

Constrained-Context Conditional Diffusion Models for Imitation Learning

Vaibhav Saxena¹, Yotto Koga² and Danfei Xu¹

Abstract—Offline Imitation Learning (IL) is a powerful paradigm to learn visuomotor skills, especially for high-precision manipulation tasks. However, IL methods are prone to *spurious correlation*—expressive models may focus on distractors that are irrelevant to action prediction—and are thus fragile in real-world deployment. Prior methods have addressed this challenge by exploring different model architectures and action representations. However, none were able to balance between sample efficiency, robustness against distractors, and solving high-precision manipulation tasks with complex action space. To this end, we present Constrained-Context Conditional Diffusion Model (C3DM), a diffusion model policy for solving 6-DoF robotic manipulation tasks with high precision and ability to ignore distractors. A key component of C3DM is a fixation step that helps the action denoiser to focus on task-relevant regions around the predicted action while ignoring distractors in the context. We empirically show that C3DM is able to consistently achieve high success rate on a wide array of tasks, ranging from table top manipulation to industrial kitting, that require varying levels of precision and robustness to distractors. For details, see sites.google.com/view/c3dm-imitation-learning

I. INTRODUCTION

Behavior cloning (BC) is an effective framework for learning visuomotor robotic skills from a fixed set of offline demonstrations. This is especially advantageous for tasks that require high-precision manipulation steps such as aligning and insertion, which present a hard exploration challenge for interactive methods such as Reinforcement Learning (RL). However, BC suffers from an important limitation — with limited data diversity, an end-to-end policy trained to map high-dimensional input such as images to actions often degenerates to focusing on spurious features in the environment instead of the task-relevant ones [1]. This leads to poor precision and fragile execution in real-world applications such as kitting and assembly, where acting in scenes with distractors is inevitable. This challenge is especially prominent for continuous-control problems where instead of committing to interacting with one object, the policy often collapses to an “average” of a mix of correct and incorrect solutions. Towards making robots more generally useful for real-world tasks, the question naturally arises — given limited offline demonstrations, how can we enable robots to ignore distractors in the environment and facilitate high-precision manipulation?

Prior work attempts to address this challenge through different model architectures and action representations. A

¹Vaibhav Saxena (vsaxena33@gatech.edu) and Danfei Xu (danfei@gatech.edu) are with the School of Interactive Computing, Georgia Institute of Technology, Atlanta, GA, USA

²Yotto Koga (yotto.koga@autodesk.com) is with Autodesk Research, San Francisco, CA, USA

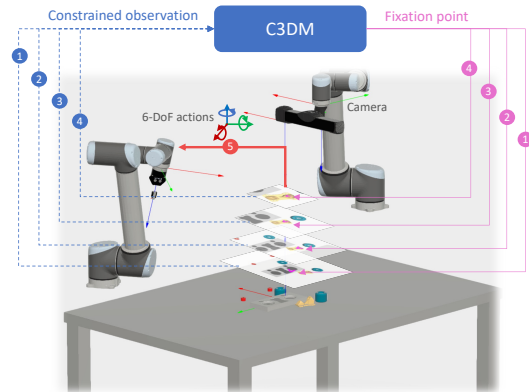


Fig. 1: Illustration showing action prediction for table-top manipulation using our conditional diffusion model (C3DM), which learns to fixate on relevant parts of the input and iteratively refine its prediction using more details about the observation.

prominent line of research [2–4] uses fully convolutional networks (FCN) that *share parameters spatially*, and by virtue of this architecture learn mappings from visual input to a spatially-quantized action space. In these models, each possible action gets assigned with a probability mass that avoids the action averaging effect in turn making the model less prone to bad precision due to distractors. While effective for 2D tasks such as pushing piles and transporting objects, this strategy prohibits tasks with complex 6-DoF action spaces as the quantization space grows exponentially with action dimensions. For modeling actions in an arbitrary continuous space, generative policy models [5–7] have recently proven to be effective methods for sample-efficient learning from limited offline demonstrations. The probabilistic nature of these models, that represent actions as either part of an energy function [5] or score function for diffusion models [7], along with the arbitrary action distribution allows these models to capture complex solutions with limited supervision. However, in the presence of distractors, these models are not immune from learning spurious correlations between images and actions resulting in inaccurate predictions. Hence, towards high-precision manipulation, we build upon one such generative model whose intermediate action representations can be used to “fixate” on task-relevant parts of the scene, and imbue in it the ability to share parameters spatially – a virtue that allows for distractor-invariant action prediction.

We present **Constrained-Context Conditional Diffusion Model (C3DM)**, a diffusion model policy that solves 6-DoF

robotic manipulation tasks with high precision and robustness to distractions. Essentially, we learn a conditional generative model on actions given input observations, where the inference distribution is fixed to a Gaussian noising process and the learned generative distribution utilizes an iteratively refining denoising process on the action variable. Key to our method is a novel denoising diffusion process that *zooms into* a part of the input image around a predicted “fixation point” at each denoising iteration, facilitating invariance to distractions during action denoising. We illustrate this process in Fig. 1. Effectively, the model fixates on the context around iteratively-refined points of interest and removes all irrelevant context by the end of the denoising process, allowing the accuracy of action prediction to reach arbitrary precision.

We empirically validate the capabilities of our model on a wide array of simulated and real-world tasks that require varying levels of precision, going to as low as 0.5 cm in the tolerance for position prediction. In all experiments, we compare our model against its variant that masks away distractions without any zooming (C3DM-Mask) and show that this, too, beats baselines but is less effective than our method (C3DM) that can predict actions using higher levels of detail. Through such training on a variety of input contexts, our method not only becomes precise but also very sample efficient, which we empirically demonstrate when we deploy our model on a real robot after training on just 20 human demonstrations.

II. RELATED WORKS

A. Visual Imitation Learning

Imitation learning (IL) has been proven effective for robotic manipulation tasks [8–10]. Recent works [11–15] use deep neural networks to map directly from image observations to action, and demonstrated visuomotor learning for complex and diverse manipulation tasks. However, these policies tend to generalize poorly to new situations due to spurious connections between pixels and actions [1,16] and requires extensive training data to become robust. To address the challenge, recent works [16–18] have incorporated visual attention [19] as inductive biases in the policy. For example, VIOLA [18] uses object bounding priors to select salient regions as policy input. Our constrained context formulation can be similarly viewed as a visual attention mechanism that iteratively refines with the diffusion process. Another line of work exploits the spatial invariance of fully convolution model and discretized action space to improve learning sample efficiency [2–4]. Notably, PerAct [4] show strong performance on 6DoF manipulation tasks. However, to implement discretized 3D action space, PerAct requires a large voxel-wise transformer model that is expensive to train and evaluate, and it also requires 3D point cloud input data. Our method instead adopts an implicit model that can model arbitrary action space and only requires 2D image input.

B. Diffusion Policy Models

Diffusion models have shown remarkable performance in modeling complex data distribution such as high-resolution

images [20,21]. More recently, diffusion models have been applied to decision making [7,22–24] and show promising results in learning complex human action [7,25,26] and conditionally generating new behaviors [22,23]. Closely related to our work is DiffusionPolicy [7] that uses a diffusion model to learn visuomotor policies. However, despite the demonstrated capability in modeling multimodal actions, DiffusionPolicy still learns an end-to-end mapping between the input and the score function and is thus prone to spurious correlation, as we will demonstrate empirically. Our method introduces a constrained context formulation that enables the denoising diffusion process to fixate on relevant parts of the input and iteratively refine its predictions.

C. Implicit Policy Models

Closely related to diffusion policy models are implicit-model policies [5,27], which represent distributions over actions using Energy-Based Models (EBMs) [28]. Finding optimal actions with an energy function-based policies can be done through sampling or gradient-based procedures such as Stochastic gradient Langevin dynamics. Alternatively, actions prediction can be implicitly represented as a distance field to the optimal action [6]. Similar to diffusion policies, implicit models can represent complex action distributions but are also prone to spurious correlations in imitation, as we will demonstrate empirically.

III. BACKGROUND

Diffusion models are a class of generative models that introduce a hierarchy of latent variables and are trained to maximize the variational lower-bound on the log-likelihood of observations (ELBO). To compute the ELBO, this class of models fixes the inference distribution on the introduced latent variables to a *diffusion process*, that can be aggregated to directly infer each latent variable from the observation. Say \mathbf{a}_0 represents the random variable that we want to model, we introduce T latent variables $\{\mathbf{a}_t\}_{t=1}^T$ each of which can be inferred from the observation \mathbf{a}_0 using

$$q(\mathbf{a}_t|\mathbf{a}_0) := \sqrt{1 - \beta(t)} \cdot \mathbf{a}_0 + \sqrt{\beta(t)} \cdot \boldsymbol{\varepsilon}_t, \quad (1)$$

where $\boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, I)$ is the standard normal Gaussian noise, and $\beta(t)$ is a fixed “noise schedule” that determines the variance of the noising process and satisfies $\beta(t) < 1 \forall t$ and $\beta(t) \rightarrow 1$ as $t \rightarrow T$. We call the index t “time,” and in principle the cardinality of the set $\{\mathbf{a}_t\}_{t=1}^T$ should tend to infinity. The choice of $q(\mathbf{a}_t|\mathbf{a}_0)$ can play a significant role in the learning process of the generative distribution, and we will see in this paper how it affects the performance of our action prediction model during evaluation.

Each of the latent variables can be interpreted as noisy variants of \mathbf{a}_0 , and hence the generative distribution over \mathbf{a}_0 can be modeled using time-conditioned distributions $p(\mathbf{a}_0|\mathbf{a}_t;t)$ where p is modeled using a neural network f_θ , also known as the “score function.” f_θ is trained to predict the standard normal noise that was used to infer \mathbf{a}_t , that is,

$$\theta = \min_{\theta} \text{MSE}(f_\theta(\mathbf{a}_t,t), \boldsymbol{\varepsilon}_t) \quad (2)$$

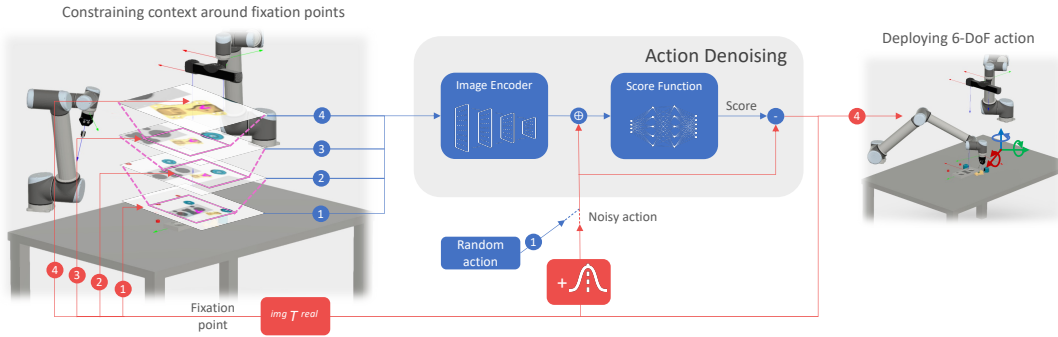


Fig. 2: **Constrained-Context Conditional Diffusion Model (C3DM)** for visuomotor policy learning. Here we illustrate our iterative refinement procedure (4 timesteps) wherein we constrain our input context around a “fixation point” predicted by the model (\wedge) at each refinement step. Subsequently, we refine the predicted action by fixating only on the useful part of the context, hence removing distractions and making use of higher levels of detail in the input.

Finally, to generate samples of \mathbf{a}_0 , we begin with samples of \mathbf{a}_T , which for the chosen noising process in Equation (1) should result in a sample from the standard normal Gaussian. Then, we “denoise” \mathbf{a}_T using the learned score function $f_\theta(\mathbf{a}_T, T)$ and further noise the resulting sample to obtain the next latent \mathbf{a}_{T-1} . This is one step of the “iterative refinement” process. We continue to refine the sample for T steps after which we return the last denoised sample, which is a sample from the distribution $p(\mathbf{a}_0)$. In this work, we focus on modeling the generative distribution over 6-dimensional action variables conditioned on image observations.

IV. METHOD

We present a **Constrained-Context Conditional Diffusion Model (C3DM)** that predicts continuous actions conditioned on input image observations. Our model has the ability to identify “fixation points” in the input, use them to ignore distractions, and process different levels of detail to iteratively refine predicted actions. As the model queries for more levels of detail, it constrains its observation context so as to focus on the useful parts of the input and ignore distractions. We assume access to a transform between the observation and action spaces which allows us to determine fixation points in the context and constrain our observations throughout the action generation process.

Environment Our setup consists of a robot that receives image input from a 2D camera and outputs a 6-DoF gripper pose in the camera frame (3D position and Euler rotation angles around z - x - y axes) as action. Our action space is modeled implicitly in that it can be queried using an input image and a candidate action to eventually infer the target action. We call this candidate action a “noisy action,” and C3DM models a score function that it uses to de-noise these actions.

Training data Our training data consists of N observations-action pairs, $D = \{\mathbf{o}^{(i)}, \mathbf{a}^{(i)}\}_{i=1}^N$ from an expert demonstrator, which we use to train our conditional generative model. Our observations are RGB images that can come from a camera in simulation or in the real world.

A. Action noising

For each observation $\mathbf{o}^{(i)}$ and target action $\mathbf{a}^{(i)}$, we sample K “noisy” actions $\{\tilde{\mathbf{a}}_k^{(i)}\}_{k=1}^K$, by adding noise vectors from a fixed, time-conditioned, noising process, using timesteps $\{t_k\}_{k=1}^K$, $t_k \sim \text{Unif}(0, 1)$ (i.e. $T = 1$ in our setup). While a standard noising process as in Equation (1) would theoretically suffice, in practice around $t = 1$ the model is supervised to learn an identity function [29] which could hamper training and eventually evaluation performance. Building on this motivation we ablate over two noising processes for generating noisy actions, that are

$$\tilde{\mathbf{a}}_k^{(i)} = \mathbf{a}^{(i)} + \sqrt{\beta(t_k)} \cdot \boldsymbol{\varepsilon}_k^{(i)}, \text{ and} \quad (3)$$

$$\tilde{\mathbf{a}}_k^{(i)} = \sqrt{1 - \beta(t_k)} \cdot \mathbf{a}^{(i)} + \sqrt{\beta(t_k)} \cdot \boldsymbol{\varepsilon}_k^{(i)}. \quad (4)$$

In either, $\boldsymbol{\varepsilon}_k^{(i)} \sim \mathcal{N}(\mathbf{0}, I)$ with \mathcal{N} being a normal distribution and I the identity matrix, and $\beta(t_k)$ is a tunable noise schedule which we keep as a linear mapping in our setup. As we describe in Section V-C.3, we find the noising process in Equation (3) (no drift) to perform better in practice than the latter (drift). We also note that using Equation (3) as the noising process would result in \mathbf{a}_k being sampled from a uniform random distribution when $t_k \sim 1$ instead of a standard normal distribution.

B. Context constraining

To imbue our model with the ability to identify and ignore distractions in the input, we determine a “fixation point” in the observation space around which we can constrain the context while iteratively refining the action. Moreover, at each constraining step our method also attains a higher level of detail by “zooming” into the context. Obtaining higher levels of input detail provides additional information for the model to make accurate prediction, especially for small objects. In our experiments, we ablate using the fixation point for simply masking distractions in the input (C3DM-Mask), against zooming into the context that additionally obtains higher levels of detail, which is our final method.

During training, we use the same set of sampled timesteps $\{t_k\}_{k=1}^K$ that we use to noise actions to determine the size of

our constraint, which in turn is fixated around the location of the target action $\mathbf{a}^{(i)}$ in the observation. Hence our fixation point is a point in the image which we determine as

$$\mathbf{p}^{(i)} := \text{img}T^{\text{real}} \text{pos}(\mathbf{a}^{(i)}), \quad (5)$$

where $\text{img}T^{\text{real}}$ is a matrix that transforms points in the camera (real) frame to image frame, and $\text{pos}(\mathbf{a}^{(i)})$ extracts the 2-D position of actuation on the x-y plane from the action. Finally, the constrained context can be formulated as

$$\mathbf{o}_k^{(i)} := C(\mathbf{o}^{(i)}; \mathbf{p}^{(i)}, t_k), \quad (6)$$

where C either masks parts of the observation far away from $\mathbf{p}^{(i)}$ (C3DM-Mask) or zooms into the context fixated at $\mathbf{p}^{(i)}$ (C3DM). We contrast these behaviors in Figure 4.

C. Denoising network

We build a model f_θ parameterized by a neural network that is supervised to predict the standard noise vector $\boldsymbol{\varepsilon}_k^{(i)}$ given the constrained input observation $\mathbf{o}_k^{(i)}$ and noisy action $\tilde{\mathbf{a}}_k^{(i)}$. That is, we train our model using the following MSE loss:

$$\mathcal{L}(D) = \frac{1}{N} \frac{1}{K} \sum_{i=1}^N \sum_{k=1}^K \|f_\theta(\mathbf{o}_k^{(i)}, \tilde{\mathbf{a}}_k^{(i)}) - \boldsymbol{\varepsilon}_k^{(i)}\|^2. \quad (7)$$

We summarize the training procedure in Algorithm 1.

D. Action inference

While testing, we use our model’s output to reverse the noising process and obtain de-noised actions initialized from uniform random actions. We break down each step of the iterative refinement process further into two steps, where we first use the intermediate de-noised action to obtain the fixation point in the observation using Equation (5), followed by the constrained context using Equation (6). We then use this constrained context to condition the denoising function in the next iteration of the refinement process. We illustrate this action refinement process in Figure 2.

Algorithm 1 C3DM - Training

Require: $D = \{\mathbf{o}^{(i)}, \mathbf{a}^{(i)}\}_{i=1}^N, K, \text{max_iters}, \text{img}T^{\text{cam}}, \theta, f_\theta$
for all $n_iter \in \{1, \dots, \text{max_iters}\}$ **do**
 $L \leftarrow 0$
 for all $k \in \{1, \dots, K\}$ **do**
 $t_k \sim \text{Unif}(0, 1)$ ▷ sampled timestep
 for all $i \in \{1, \dots, N\}$ **do**
 $\boldsymbol{\varepsilon}_k^{(i)} \sim \mathcal{N}(\mathbf{0}, I)$ ▷ sampled noise
 $\tilde{\mathbf{a}}_k^{(i)} \leftarrow \mathbf{a}^{(i)} + \sqrt{\beta(t_k)} \cdot \boldsymbol{\varepsilon}_k^{(i)}$ ▷ noisy action
 $\mathbf{p}^{(i)} \leftarrow \text{img}T^{\text{real}} \text{pos}(\mathbf{a}^{(i)})$ ▷ fixation point
 $\mathbf{o}_k^{(i)} \leftarrow C(\mathbf{o}^{(i)}; \mathbf{p}^{(i)}, t_k)$ ▷ constrained context
 $L \leftarrow L + \|f_\theta(\mathbf{o}_k^{(i)}, \tilde{\mathbf{a}}_k^{(i)}) - \boldsymbol{\varepsilon}_k^{(i)}\|^2$
 end for
 end for
 $\theta \leftarrow \theta - \frac{1}{NK} \nabla_\theta L$
end for

TABLE I: Simulation (\star) and real-robot (\dagger) tasks and their desiderata.

Task	Ignore distractions	Precision	Position tolerance
sweeping-piles \star	✗	✗	n/a
place-red-in-green \star	✓	✓	4.0 cm
kitting-part \star	✓	✓	2.0 cm
hang-cup \star	✗	✓	1.0 cm
two-part-assembly \star	✗	✓	0.5 cm
block-in-bowl \dagger			
- big blocks	✓	✗	4.5 cm
- small blocks	✓	✓	2.5 cm
screw-in-hole \dagger	✗	✓	1.0 cm

V. EXPERIMENTS

We evaluate C3DM on 5 tasks in simulation with varying tolerance for success. We compare our method with two state-of-the-art implicit methods for behavior cloning (BC), an explicit BC method, as well as two ablations that highlight the key contributions of our method. We show that C3DM can predict actions with high precision, and identify and ignore distractions in table-top manipulation tasks. Finally, we employ our model for solving tasks using a real robot showcasing sample efficiency when learning from just 20 human demonstrations.

A. Baselines

We consider two state-of-the-art implicit learning methods, Implicit Behavioral Cloning (IBC) [5] and Neural Grasp Distance Fields (NGDF) [6] for comparison, as well as an explicit behavior cloning model (Conv. MLP). We also implemented a baseline diffusion model without context-constraining, which we call Conditional Diffusion. Additionally, we ablate against a version of C3DM, called C3DM-Mask, that is trained to only remove distractions in the observation by identifying and masking irrelevant details.

B. Model Architecture and Training Details

We use a deep convolutional neural network with residual connections (ResNet-18 [30]) to process images. We flatten the output embedding and concatenate it with query actions, which are then further processed by 4 fully-connected feedforward layers with skip connections. We use ReLU activations for all intermediate layers. We implemented all baselines with the same backbone architecture, tune learning rate in the range $[10^{-4}, 10^{-3}]$, and train using the Adam optimizer [31] with a batch size of 100 demonstrations.

C. Simulation Experiments

1) Tasks

We evaluate our method on 5 different tasks in simulation, listed in Table I with their corresponding desiderata. Since IBC experimented predominantly on tasks with “push” primitives, we include the *sweeping-piles* task for a fair comparison. The remaining tasks, *place-red-in-green*, *hang-cup*, *kitting-part*, and *two-part-assembly* are based on “pick” / “grasp”

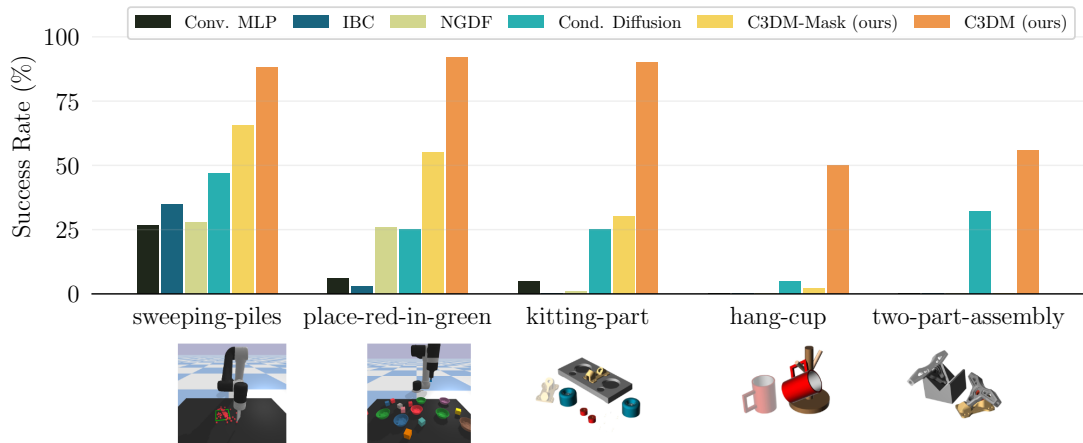


Fig. 3: (Top) Success rates for manipulation tasks in simulation (average across 100 rollouts, peak performance in 500 epochs of training). (Bottom) Illustration of the simulated evaluation tasks.

and “place” primitives. *Place-red-in-green* and *sweeping-piles* were made available by [32] as part of the Ravens simulation benchmark, and we built the remaining tasks on an Autodesk Research robotics platform [33]. The Autodesk platform connects with a variety of different robotic hardware, peripherals and assembly workcells, with a design, to simulate, to workcell deployment pipeline. For each task we assume access to 1000 demos from an oracle demonstrator, each demo being one image-action pair. Tasks in Ravens require 2-D and yaw actions only, hence we only collect 100 demos and employ rotation augmentation to obtain a total of 1000 demos. For tasks in the Autodesk platform that require full 6-DoF actions we collect 1000 demos.

Sweeping-piles. This task requires the model to output the parameters of a push primitive, that is the starting location and orientation of a pusher and ending location, to push piles of small objects into a target zone.

Place-red-in-green. This is a slightly more precision-requiring task where the robot is supposed to pick up a red 4x4x4 cm block using a suction cup and place it in a green bowl. The table is also laid with distractor blocks that can hinder model precision leading to imprecise picks.

Kitting-part. In this task, 5 parts of a skateboard truck assembly are laid out on a tabletop, and the robot is tasked to grasp the truck base on its slotted end and place correctly in a kit. To succeed in the task, our hypothesis is that the robot needs to fixate on and observe the part in detail, while ignoring distractions that can hinder precision.

Hang-cup. This high-precision task requires the model to predict the pose of the gripper for picking and placing that should result in the cup hanging by its handle on the hook of a Y-shaped hanger. In the best possible scenario where the cup is oriented such that the handle’s plane is perpendicular to the hanger’s hook, the tolerance for error in position prediction is only 1 cm.

Two-part-assembly. This is a very high precision assembly task where the robot is tasked to pick up a skateboard truck hanger and insert it in the hole of the truck base. This task is a subroutine of assembling a full skateboard truck.

2) Main results

C3DM can identify and ignore distractions in table-top manipulation. Both C3DM-Mask and C3DM show good performance on *place-red-in-green* due to their ability to ignore distractor objects leading to precise locations for picking. On the contrary, baseline methods struggle to fixate on the target object that needed to be picked. In *kitting-part*, while baseline methods are able to predict actions approximately around the target action, they were not precise enough to succeed given the low tolerance of this task. C3DM learns to ignore unnecessary objects on the table, as shown in Figure 4, leading to precise pick and place predictions. We observed failures when the distractor objects were too close to the object and C3DM could not push them out of the field of view, leading to imprecise pick predictions.

Action refinement with a fixated gaze can help predict 6-DoF gripper poses with high precision. C3DM is the only method that could succeed substantially on the *hang-cup* and *two-part-assembly* tasks due to its ability to precisely predict the full 6-DoF action. C3DM beat all other baselines and the C3DM-Mask ablation showing the importance of action refinement with a fixated gaze. We illustrate how our learned score field creates this fixation in Figure 5. The cases where our model did fail were in which the model fixated on spurious locations early in the refinement process leading to the correct pick location being outside the region of view, as well as when small errors in the action prediction led to irrecoverable scenarios.

Iterative refinement with more levels of input detail can solve visual challenges in tasks. The *sweeping-piles* task poses the challenge of inferring the location of small particles and the target zone from an image, for which C3DM outperforms all considered baselines. We note that IBC’s performance on this task is lower than that reported in their work for planar sweeping, because our version of this task is more visually complex with the target zone marked by a thin square as opposed to being marked by a color-filled region. The *kitting-part* task also presents a complex-looking part to

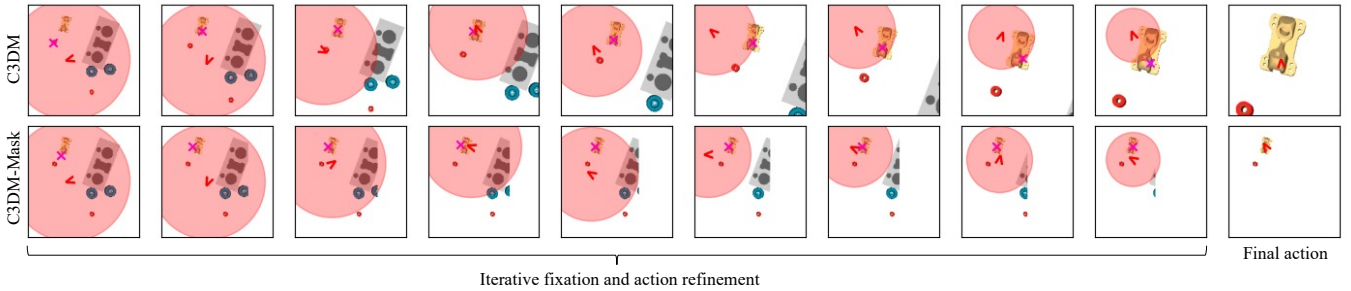


Fig. 4: Constrained-context action refinement on the *kitting-part* (pick) task using C3DM (top) and C3DM-Mask (bottom). The red region (overlaid on the top-down image capture) depicts the standard deviation of the predicted latent action, \wedge represents the latent action (orientated with grasp yaw), and \times is the fixation point. We observe how the model fixates on the target part (yellow) as it refines its action for grasping.

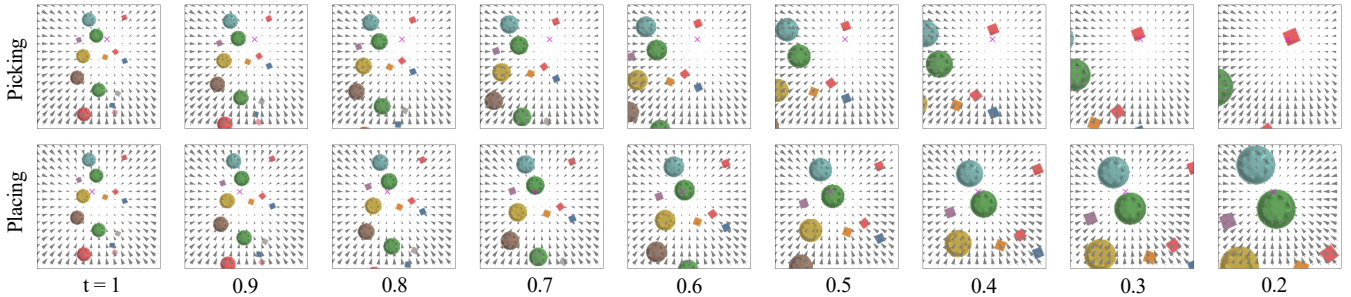


Fig. 5: Score fields and the induced fixation points (denoted by \times) when inferring the pick (top) and place (bottom) actions for the *place-red-in-green* task using C3DM. For both picking and placing, we note that the model is distracted at first with the entire table-top in view. As it zooms in, our model fixates closer and closer to the object of interest, which are a red block for picking and a green bowl for placing.

TABLE II: Comparing model performance after training for 200 epochs with and without drift in the diffusion process.

Method	C3DM	C3DM-Mask
Drift	60%	40%
No drift	79%	47%

infer the pose of in order to predict a stable grasp. C3DM can iteratively refine its pose by obtaining higher levels of detail in the input leading to a substantially low grasp error of 0.72 cm and 11.28° and eventually a higher success rate compared to all baseline methods, including the C3DM-Mask ablation, which can only ignore distractions in the input.

3) Additional Observations

Table II shows comparison between two diffusion processes for training our model, one that drifts the latent action to the origin (Equation (4)), and the other a pure diffusion (Equation (3)). We observed the latter to perform better in practice and is the default choice for all our main results. We also show the success rates when varying the number of refinement steps during action refinement in Figure 6, and as expected we observe a rising trend in success as number of refinement steps increase.

D. Real-robot Experiments

We demonstrate the efficacy of our model for ignoring distractions and precise action prediction on a real robot for

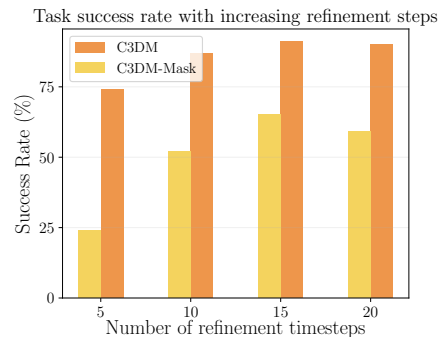


Fig. 6: Success rate on the *place-red-in-green* task with increasing number of refinement timesteps.

two tasks, *block-in-bowl* and *screw-in-hole*. The setup consists of a Franka Emika Panda arm and an Intel RealSense D435 camera mounted viewing the workspace top-down. Success rates averaged over 20 trials are summarized in Table III.

1) Data Collection, Training and Evaluation

We collected 20 human demonstrations using a space-mouse teleoperator for each task which required picking target objects and placing them into a goal location. Each demo consisted of an initial image observation and the locations for pick and place actions in the robot’s coordinate frame. We trained a baseline Conditional Diffusion model, C3DM-Mask, and C3DM for 200 epochs. During testing, we evaluated all

TABLE III: Success rates for tasks on a real robot.

Method	block-in-bowl		screw-in-hole	
	big blocks	small blocks	pick	pick+place
Cond. Diffusion	45%	0%	5%	0%
C3DM-Mask	60%	40%	0%	0%
C3DM	55%	65%	80%	60%



Fig. 7: Real-world experiment setup and tasks.

models using a linear schedule of 10 timesteps.

2) Model Performance on Real Robots

C3DM can ignore distractions in real images. We performed experiments on two variants of the *block-in-bowl* task, one with big blocks with edge size 4.5 cm and other with small blocks of edge size 2.5 cm. Example image observations are shown in Figure 7. We observed that C3DM was able to effectively identify and ignore distractions in the input leading to precise actions for both variants. While C3DM-Mask was able to ignore distractions well, it lacked precision when evaluated with smaller blocks.

Iterative context constraining leads to precise actions and high sample efficiency. We performed experiments on a *screw-in-hole* task where the robot is supposed to pick up a hex-head screw of edge size 1 cm, base diameter 1 cm, and place it in a receptacle with a hole of diameter 1.5 cm. Successfully completing this task requires high precision as a pick prediction that is slightly off from the center of the screw would result in the screw falling on the table instead of being clamped by the parallel-jaw gripper of the robot. We observed that C3DM was significantly more successful in completing this task. While other methods failed due to both imprecise action prediction as well as the lack of generalization capability given a small number of demonstrations, C3DM was able to generalize to all test locations as well as be precise within the 1 cm tolerance needed for this task.

VI. CONCLUSION

We presented a Constrained-Context Conditional Diffusion Model (C3DM) for visuomotor imitation learning in high-precision manipulation tasks. We demonstrated that the fixation-based context constraining allows our diffusion model to iteratively refine sampled actions while removing distractors from the input, achieving high success rate in a wide range of tasks requiring varying levels of precision and robustness against distraction. We also identified failure modes in our experiments, and believe that future work with goal-conditioning policies and training on action trajectories (rather than sub-task goals) will alleviate failures and further render our model applicable for longer horizon manipulation.

REFERENCES

- [1] P. De Haan, D. Jayaraman, and S. Levine, “Causal confusion in imitation learning,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [2] A. Zeng, P. Florence, J. Tompson, S. Welker, J. Chien, M. Attarian, T. Armstrong, I. Krasin, D. Duong, V. Sindhwani, *et al.*, “Transporter networks: Rearranging the visual world for robotic manipulation,” in *Conference on Robot Learning*. PMLR, 2021, pp. 726–747.
- [3] M. Shridhar, L. Manuelli, and D. Fox, “Cliport: What and where pathways for robotic manipulation,” in *Conference on Robot Learning*. PMLR, 2022, pp. 894–906.
- [4] —, “Perceiver-actor: A multi-task transformer for robotic manipulation,” 2022.
- [5] P. Florence, C. Lynch, A. Zeng, O. Ramirez, A. Wahid, L. Downs, A. Wong, J. Lee, I. Mordatch, and J. Tompson, “Implicit behavioral cloning,” 2021. [Online]. Available: <https://arxiv.org/abs/2109.00137>
- [6] T. Weng, D. Held, F. Meier, and M. Mukadam, “Neural grasp distance fields for robot manipulation,” 2023.
- [7] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song, “Diffusion policy: Visuomotor policy learning via action diffusion,” 2023.
- [8] S. Schaal, “Is imitation learning the route to humanoid robots?” *Trends in Cognitive Sciences*, vol. 3, pp. 233–242, 1999.
- [9] A. Billard, S. Calinon, R. Dillmann, and S. Schaal, “Robot programming by demonstration,” in *Springer Handbook of Robotics*, 2008.
- [10] P. Englert and M. Toussaint, “Learning manipulation skills from a single demonstration,” *The International Journal of Robotics Research*, vol. 37, no. 1, pp. 137–154, 2018.
- [11] T. Zhang, Z. McCarthy, O. Jow, D. Lee, X. Chen, K. Goldberg, and P. Abbeel, “Deep imitation learning for complex manipulation tasks from virtual reality teleoperation,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 5628–5635.
- [12] C. Finn, T. Yu, T. Zhang, P. Abbeel, and S. Levine, “One-shot visual imitation learning via meta-learning,” in *Conference in Robot Learning*, vol. abs/1709.04905, 2017.
- [13] A. Mandlekar, D. Xu, R. Martín-Martín, S. Savarese, and L. Fei-Fei, “Learning to generalize across long-horizon tasks from human demonstrations,” *arXiv preprint arXiv:2003.06085*, 2020.
- [14] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese, Y. Zhu, and R. Martín-Martín, “What matters in learning from offline human demonstrations for robot manipulation,” *arXiv preprint arXiv:2108.03298*, 2021.
- [15] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, *et al.*, “Rt-1: Robotics transformer for real-world control at scale,” *arXiv preprint arXiv:2212.06817*, 2022.
- [16] C. Wang, R. Wang, D. Xu, A. Mandlekar, L. Fei-Fei, and S. Savarese, “Generalization through hand-eye coordination: An action space for learning spatially-invariant visuomotor control,” *CoRR*, vol. abs/2103.00375, 2021. [Online]. Available: <https://arxiv.org/abs/2103.00375>
- [17] P. Abolghasemi, A. Mazaheri, M. Shah, and L. Bölöni, “Pay attention! - robustifying a deep visuomotor policy through task-focused attention,” *CoRR*, vol. abs/1809.10093, 2018. [Online]. Available: <http://arxiv.org/abs/1809.10093>
- [18] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu, “Viola: Imitation learning for vision-based manipulation with object proposal priors,” *arXiv preprint arXiv:2210.11339*, 2022.
- [19] V. Mnih, N. Heess, A. Graves, and K. Kavukcuoglu, “Recurrent models of visual attention,” *CoRR*, vol. abs/1406.6247, 2014. [Online]. Available: <http://arxiv.org/abs/1406.6247>
- [20] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, “Score-based generative modeling through stochastic differential equations,” *CoRR*, vol. abs/2011.13456, 2020. [Online]. Available: <https://arxiv.org/abs/2011.13456>
- [21] J. Ho, A. Jain, and P. Abbeel, “Denosing diffusion probabilistic models,” *CoRR*, vol. abs/2006.11239, 2020. [Online]. Available: <https://arxiv.org/abs/2006.11239>
- [22] A. Ajay, Y. Du, A. Gupta, J. Tenenbaum, T. Jaakkola, and P. Agrawal, “Is conditional generative modeling all you need for decision-making?” *arXiv preprint arXiv:2211.15657*, 2022.
- [23] Z. Zhong, D. Rempe, D. Xu, Y. Chen, S. Veer, T. Che, B. Ray, and M. Pavone, “Guided conditional diffusion for controllable traffic simulation,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 3560–3566.

- [24] U. A. Mishra, S. Xue, Y. Chen, and D. Xu, "Generative skill chaining: Long-horizon skill planning with diffusion models," in *7th Annual Conference on Robot Learning*, 2023.
- [25] T. Pearce, T. Rashid, A. Kanervisto, D. Bignell, M. Sun, R. Georgescu, S. V. Macua, S. Z. Tan, I. Momennejad, K. Hofmann, *et al.*, "Imitating human behaviour with diffusion models," *arXiv preprint arXiv:2301.10677*, 2023.
- [26] M. Reuss, M. Li, X. Jia, and R. Lioutikov, "Goal-conditioned imitation learning using score-based diffusion policies," *arXiv preprint arXiv:2304.02532*, 2023.
- [27] D. Jarrett, I. Bica, and M. van der Schaar, "Strictly batch imitation learning by energy-based distribution matching," *Advances in Neural Information Processing Systems*, vol. 33, pp. 7354–7365, 2020.
- [28] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang, "A tutorial on energy-based learning," *Predicting structured data*, vol. 1, no. 0, 2006.
- [29] T. Salimans and J. Ho, "Progressive distillation for fast sampling of diffusion models," *CoRR*, vol. abs/2202.00512, 2022. [Online]. Available: <https://arxiv.org/abs/2202.00512>
- [30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, vol. abs/1512.03385, 2015. [Online]. Available: <http://arxiv.org/abs/1512.03385>
- [31] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014. [Online]. Available: <http://arxiv.org/abs/1412.6980>
- [32] A. Zeng, P. Florence, J. Tompson, S. Welker, J. Chien, M. Attarian, T. Armstrong, I. Krasin, D. Duong, A. Wahid, V. Sindhwani, and J. Lee, "Transporter networks: Rearranging the visual world for robotic manipulation," 2020. [Online]. Available: <https://arxiv.org/abs/2010.14406>
- [33] Y. Koga, H. Kerrick, and S. Chitta, "On cad informed adaptive robotic assembly," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 10 207–10 214.