

Learning Factorized Diffusion Policies for Conditional Action Diffusion

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Diffusion models have emerged as a promising choice for learning
2 robot skills from demonstrations. However, diffusion models are neither robust
3 to visual distribution shifts nor sample-efficient for policy learning. In this work,
4 we present ‘Factorized Diffusion Policies’ abbreviated as FDP, a novel theoretical
5 framework to learn action diffusion models without the need to jointly con-
6 dition on all observational modalities such as proprioception and vision. Using
7 our factored approach leads to 10% absolute performance improvement for ten
8 RLBench and four Adroit tasks when compared to a standard diffusion policy
9 which jointly conditions on all modalities. Moreover, FDP results in 25% higher
10 absolute performance across five RL Bench tasks with distribution shifts such as
11 visual changes or distractors, where existing diffusion policies fail catastrophically.
12 Our real-world experiments show that FDP is safe and relatively robust to
13 deploy against visual distractors and appearance changes when compared to stan-
14 dard diffusion policies. Videos are available at <https://fdp-policy.github.io/fdp-policy/>.
15

16 1 Introduction

17 Diffusion models have emerged as a promising choice for learning robot skills from demonstrations
18 [1]. Following various diffusion models, several generative models originally proposed in the vision
19 literature have been used for robot learning, exploiting properties such as one-step inference [2, 3]
20 and multimodal priors [4]. However, unlike computer vision, conditioning is critical in robotics due
21 to the numerous observational modalities that influence the robot’s action choices. Humans prioritize
22 different sensory modalities according to the specific requirements of the task [5]. Humans have
23 also been shown to prioritize the more reliable modality between vision and haptics [6]. Naturally,
24 based on the task, robot skills should also depend more strongly on certain observational modes
25 than others. For instance, repetitive motions like dance are more likely to depend on the robot’s
26 proprioception, while search and rescue is conditioned strongly on its vision.

27 However, the current method of training diffusion policies jointly conditions the action diffusion
28 process on all the observational modalities for every task [1]. This is a monolithic joint condition-
29 ing approach – “when all you have is a hammer, everything looks like a nail”. Learning the full
30 conditional action distribution makes Diffusion Policies sensitive to distribution shifts in any of the
31 modalities. We show that learning the full conditional results in low sample efficiency, brittleness
32 to distribution shifts. In this work, we propose a novel *theoretical framework* ‘Factorized Diffusion
33 Policies’ FDP for learning action diffusion models that decouples observational modalities for pri-
34 oritization. At its core, FDP learns a *residual model* using some input modalities that have been
35 omitted while training a base model with *prioritized inputs*. The base and residual model outputs
36 are then composed to obtain samples from the full conditional action distribution. In addition, we
37 present an architecture that enables efficient learning of the residual model in the FDP framework.
38 We demonstrate that prioritization of modalities may yield significant gains in sample efficiency
39 and naturally improves policy robustness to distribution shifts in the residual observations. Our
40 contributions are as follows.

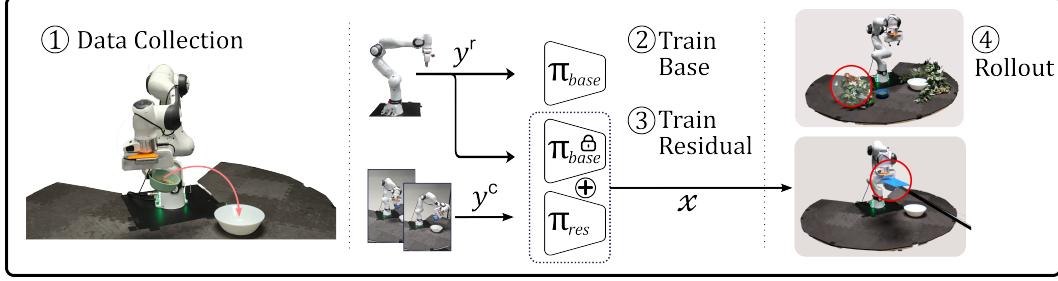


Figure 1: Training and inference for learning visuomotor policies using FDP with vision as a residual over proprioception. FDP is robust to deployment with distractors and camera occlusions.

- 41 1. We introduce Factorized Diffusion Policies (FDP), a novel theoretical framework for training dif-
42 fusion models on robot demonstration data that decouples observation modalities for prioritiza-
43 tion. We derive a novel loss function for learning a residual model on top of a policy trained with
44 prioritized modalities, and propose an efficient architectural implementation to ease its learning.
- 45 2. Our experiments show that prioritization of observational modalities produces significant sample
46 efficiency gains in several RLBench [7] and Adroit hand manipulation [8] environments. We
47 show through several distractor experiments on RLBench that learning a visual residual model
48 using FDP results in policies that are 25% more performant over standard diffusion policies.
- 49 3. We collect demonstrations across several task environments on a real robot and evaluate both
50 FDP and standard diffusion policy in the original environments as well as in modified versions
51 with visual distractors and appearance changes. In our real-world experiments, FDP outperforms
52 diffusion policies by over 40% in the presence of distractors, occlusions and appearance changes.

53 2 Background and Related Work

54 **Diffusion Models.** Gaussian diffusion models [9] learn the reverse diffusion kernel $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$ for
55 a fixed forward kernel that adds Gaussian noise at each step $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 -
56 \alpha_t)\mathcal{I})$, such that $q(\mathbf{x}_T) \approx \mathcal{N}(0, \mathcal{I})$. Here, $t <= T$ is the diffusion time step and α_t is the noise
57 schedule. For training the model, maximization of the evidence lower bound on the log-likelihood of
58 the data distribution $\log q(\mathbf{x}_0)$ yields the commonly used loss function in Equation 1 [10, 11].

$$\mathcal{L}_t(\theta) = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \epsilon_0 \sim \mathcal{N}(0, \mathcal{I})} [\lambda_t [\|\epsilon_0 - \hat{\epsilon}_{\theta}(\mathbf{x}_t, t)\|_2^2]] \quad (1)$$

59 Here, λ_t , a function of α_t is the weighting parameter for different time steps, usually taken as 1
60 [10]. The model is trained to predict the noise ϵ_0 added to the data sample \mathbf{x}_0 to generate the noisy
61 sample \mathbf{x}_t taken as input to the network.

62 **Connection to Score-based Models.** Song et al. [12] presented a unified framework showing that
63 both diffusion models [9, 10] and score-based models [13] can be interpreted as discretizations of
64 different forward stochastic differential equations (SDEs). Denoising score matching (DSM) [14]
65 is used to learn the score $\nabla_{\tilde{\mathbf{x}}} \log q_{\sigma}(\tilde{\mathbf{x}})$ at different noise scales σ required for sampling from the
66 data distribution via the corresponding reverse-time SDEs [15]. Explicit Score Matching (ESM)
67 [16, 14] was proposed to estimate this score by minimizing the Fisher divergence with the Gaussian-
68 smoothed data distribution $q_{\sigma}(\tilde{\mathbf{x}}) = \int q(\mathbf{x}) \mathcal{N}(\tilde{\mathbf{x}}; \mathbf{x}, \sigma^2 I) d\mathbf{x}$. DSM alleviates the computational dif-
69 ficulties of ESM [14, 17, 18], and is shown in Equation 2, where $s_{\theta}(\tilde{\mathbf{x}})$ represents the learned score
70 model.

$$\mathcal{J}_{\sigma_t}(\theta) \stackrel{ESM}{=} \mathbb{E}_{q_{\sigma_t}(\tilde{\mathbf{x}})} \left[\frac{1}{2} \|\nabla_{\tilde{\mathbf{x}}} \log q_{\sigma_t}(\tilde{\mathbf{x}}) - s_{\theta}(\tilde{\mathbf{x}})\|_2^2 \right] \stackrel{DSM}{=} \mathbb{E}_{q_{\sigma_t}(\mathbf{x}, \tilde{\mathbf{x}})} \left[\frac{1}{2} \|\nabla_{\tilde{\mathbf{x}}} \log q_{\sigma_t}(\tilde{\mathbf{x}}|\mathbf{x}) - s_{\theta}(\tilde{\mathbf{x}})\|_2^2 \right] + C \quad (2)$$

71 Diffusion models use a forward transition kernel $q(\tilde{\mathbf{x}}|\mathbf{x}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_{t-1}, (1 - \bar{\alpha}_t)\mathcal{I})$ with dis-
72 crete time and $\bar{\alpha}_i = \prod_{j=1}^i \alpha_j$, yielding the loss shown in Equation 1, while score-based model typ-
73 ically use $\mathcal{N}(\mathbf{x}_t; \mathbf{x}_{t-1}, \sigma_t^2 \mathcal{I})$, where α_t and σ_t are respective noise scales. Based on the equivalence
74 of Equations 1 and 2, an optimal diffusion model learned using Equation 1, is related to the score of

75 the diffused data distribution by $\epsilon_\theta^*(\mathbf{x}_t, t) / \sqrt{1 - \bar{\alpha}_t} = \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$ [12, 11]. Typically, diffusion
 76 models generate samples via progressive denoising through the reverse diffusion process [10], while
 77 score-matching models sample from the data distribution using Langevin dynamics [19, 20].

78 **Relevant work in Robotics.** Sample efficiency and generalization are of primary importance in
 79 robotics, as scaling the collection of multimodal data is difficult and the number of variations of
 80 tasks is unbounded. While generative model families such as diffusion [1], score-based models
 81 [21], stochastic interpolants [4], and flows [2] have been applied in robotics, they do not address
 82 these limitations. Prior compositional works have tried to address these problems by composing
 83 learned constraints to generalize to new task combinations in manipulation [22] and planning [23], or
 84 composing distributions across heterogeneous modalities for tool use [24]. However, all the previous
 85 works compose learned or analytical distributions, limiting their application to combinations of
 86 existing solutions. Instead, in our FDP framework, we learn a residual over a base policy that, when
 87 composed with the base policy, provides samples corresponding to the data distribution. Recent
 88 augmentation-based methods [25, 26] improve generalization but add a substantial computational
 89 overhead and remain vulnerable to visual failures like temporary camera occlusions or dynamic
 90 scene changes. In contrast, FDP is an algorithmic improvement that achieves robustness to such
 91 perturbations without data augmentations as demonstrated in our real-world experiments.

92 3 Methodology

93 Assume that we have robot demonstrations $D = \{(\mathbf{x}, \mathbf{y})_i\}$ where $i = 1..N$, consisting of actions \mathbf{x}
 94 and different observational modalities $\mathbf{y}^{1:M}$, such as images or point clouds from different cameras
 95 and proprioception data. We are interested in learning $p(\mathbf{x}|\mathbf{y})$ from the data such that given a task
 96 description, current camera images, state of the robot, and other observations, we can sample an
 97 action \mathbf{x} with a high likelihood in the data distribution. Most treatments of diffusion models have
 98 been studied primarily in the context of single-modality distributions, such as those over image
 99 pixels [10, 27, 13]. This formulation has been directly adopted by the robotics community [1, 21,
 100 28], leading to the popular optimization objective shown in Equation 3.

$$\mathcal{L}_t(\theta) = \mathbb{E}_{(\mathbf{x}_0, \mathbf{y}) \sim q(\mathbf{x}_0, \mathbf{y}), \epsilon_0 \sim \mathcal{N}(0, I)} [\|\epsilon_0 - \hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)\|_2^2] \quad (3)$$

101 Here, the network ϵ_θ is conditioned on the observation \mathbf{y} , and is trained to predict the noise added
 102 to the action sample \mathbf{x} . Although prior work in robotics adopts this conditional formulation [1], it
 103 is assumed without formal justification that the trained network maximizes the log-likelihood of the
 104 conditional distribution $p(\mathbf{x}|\mathbf{y})$. Hence, we present our first result as follows.

105 **Lemma 3.1.** *The diffusion loss function $\mathcal{L}_t(\theta)$ as defined in Equation 3, in expectation over the time-
 106 steps $1 \leq t \leq T$, maximizes the variational lower bound on the log-likelihood of the conditional
 107 data distribution $\log q(\mathbf{x}|\mathbf{y})$, under a Markovian noising process $\hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1})$ and the conditional
 108 reverse transition kernel as $\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})$.*

109 The proof for Lemma 3.1 is presented in Appendix C.2. In Equation 3, $\epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t)$ arises from
 110 the reparametrization of the reverse transition kernel $q_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}^{1:M})$, and from a score-based
 111 perspective, it learns the score of the full action conditional $\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{y}^{1:M})$ times a constant.
 112 We argue that learning the full conditional directly is restrictive in several aspects of robot learning.
 113 Firstly, it necessitates the joint collection of the robot action and all observational modalities. This
 114 restriction makes it impossible to add observational modalities later. Secondly, the model is vulnera-
 115 ble to even small distribution shifts in *any* modality. These shifts require a prohibitively large amount
 116 of data to address when the observation modalities are high-dimensional. Finally, among the mul-
 117 tiple observation modalities it is hard to pinpoint the level of each mode’s task dependent influence
 118 with limited data. Hence, we present FDP, a method to add structure and decouple observational
 119 modalities in the score of the full action conditional $\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{y}^{1:M})$. By factorizing modalities
 120 using Bayes’ theorem and learning residuals for subsequent terms, FDP effectively encodes task
 121 requirements and learns policies robust to distribution shifts in residual modalities.

122 **3.1 Factorized Diffusion Policies**

123 Let $\mathbf{y}^{1:k}$ be the prioritized observational modalities of the M total modalities, where $\mathbf{y}^{1:k} \equiv \mathbf{y}^1, \dots, \mathbf{y}^k$
 124 and $1 \leq k < M$. To decouple the observational modalities, we utilize Bayes' theorem on the score of
 125 the full action conditional to obtain the following.

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}^{1:M}; \theta, \phi) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}^{1:k}; \theta) + \nabla_{\mathbf{x}_t} \log p(\mathbf{y}^{k+1:M} | \mathbf{x}_t, \mathbf{y}^{1:k}; \phi) \quad (4)$$

126 To prioritize modalities $\mathbf{y}^{1:k}$, we propose to first learn a diffusion policy π_{base} : $\epsilon_\theta(\mathbf{x}_t, \mathbf{y}^{1:k}, t)$ that
 127 corresponds to the first score term on the right-hand side of Equation 4. To learn the second score
 128 term, explicitly training a classifier $p(\mathbf{y}^{k+1:M} | \mathbf{x}_t, \mathbf{y}^{1:k})$ [29] is impractical due to the high dimen-
 129 sionality and continuity of observational modalities $\mathbf{y}^{1:M}$, such as images. Hence, we employ ex-
 130 plicit score matching [16, 14] as shown in Equation 5.

$$D_F^t = \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:M})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1:M} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1:M} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] \quad (5)$$

131 Here, observational modalities $\mathbf{y}^{1:M}$ can be noised with a Gaussian kernel $\mathcal{N}(\tilde{\mathbf{y}}; \mathbf{y}, \tau^2 I)$ of variance
 132 τ^2 that is small enough such that $p_\tau(\tilde{\mathbf{y}}^i) \approx p(\mathbf{y}^i)$. Chao et al. [18] show that the empirical score is
 133 difficult to estimate for large datasets and derive the denoising likelihood score matching (DLSM)
 134 objective for conditional distributions, which forms the basis for our next result.

135 **Theorem 3.2.** *Explicit score matching for $\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1:M} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})$, as expressed in Equation
 136 5 with \mathbf{x} is noised with the diffusion transition kernel $\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}, (1 - \alpha_t) \mathcal{I})$, is equal up to a
 137 constant to the following loss:*

$$L_{\text{res}}^t(\phi) = \mathbb{E}_{p_\tau(\mathbf{x}_0, \mathbf{y}^{1:M}, \tilde{\mathbf{y}}^{1:M})} \mathbb{E}_{\epsilon_0 \sim \mathcal{N}(0, \mathcal{I})} \left[\frac{1}{2} \|\epsilon_0 - \epsilon_\theta(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}, t) - \hat{\epsilon}_\phi(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:M}, t)\|_2^2 \right] \quad (6)$$

138 The proof for Theorem 3.2 is presented in Appendix C.3. Equation 6 allows us to train the score
 139 of the classifier π_{res} : $\hat{\epsilon}_\phi(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:M}, t)$ as a residual over frozen π_{base} to predict noise ϵ_0 added to the
 140 action \mathbf{x}_0 . Learning π_{res} as a residual of π_{base} ensures that the model does not overfit modalities
 141 $y_{k+1:M}$, but only learns correlations to bridge the gap between the expected score and the predicted
 142 score of the model π_{base} trained on the prioritized modalities $\mathbf{y}^{1:k}$. Hence, policies learned in this
 143 factorized way are naturally robust to distribution shifts in the residual modalities. Moreover, ex-
 144 plicit prioritization of $\mathbf{y}^{1:k}$ by training π_{base} prior to learning the residual leads to sample efficiency,
 145 as the model learns correlations with the stronger modality without having to attend to other modal-
 146 ities. Since diffusion models are trained on discrete time steps, the residual is learned on the same
 147 time discretization as used for π_{base} . Once trained, actions can be sampled from the conditional dis-
 148 tribution $p(\mathbf{x} | \mathbf{y}^{1:M})$ using reverse diffusion [10] on the composition [30] of π_{base} and π_{res} :

$$\mathbf{x}_{t-1} \sim \mathcal{N}\left(\mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon(\mathbf{x}_t, \mathbf{y}^{1:M}, t) \right), \sqrt{1 - \alpha_t} \mathcal{I}\right) \quad (7)$$

$$\epsilon(\mathbf{x}_t, \mathbf{y}^{1:M}, t) = \epsilon_\theta(\mathbf{x}_t, \mathbf{y}^{1:k}, t) + \hat{\epsilon}_\phi(\mathbf{x}_t, \mathbf{y}^{1:M}, t) \quad (8)$$

149 The specific instantiations of FDP for combinations of modalities are presented in Appendix C.1.
 150 Note that the base mode π_{base} : $\hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}^{1:k}, t)$ can be further decomposed with respect to observa-
 151 tional modalities. In this work, the residual model π_{res} is learned in expectation over data collected
 152 jointly for all modalities $\mathbf{y}^{1:M}$, of which $\mathbf{y}^{1:k}$ are used for training π_{base} . However, we emphasize
 153 that an important feature of our learning formulation is that it enables decoupled data collection for
 154 additional conditionals that could then be used to learn the residual. This potentially may alleviate
 155 some difficulties encountered for scaling coupled data in robotics and is left for future work.

156 **3.2 Architectural Implementations of FDP**

157 The base and residual models in FDP, denoted by π_{base} and π_{res} , can be instantiated using standard
 158 architectures such as UNet [31] or DiT [32]. FDP involves the additional step of learning π_{res}
 159 as a residual over a frozen π_{base} . During inference we compose the outputs of these as shown in
 160 Figure 2 [b]. However, we find this late-stage residual learning to be inefficient in practice and
 161 propose a more integrated way to compose π_{base} and π_{res} , as shown in Figure 2. This architecture
 162 enables a simplified training objective for the residual model, equivalent to Equation 3. Instead of

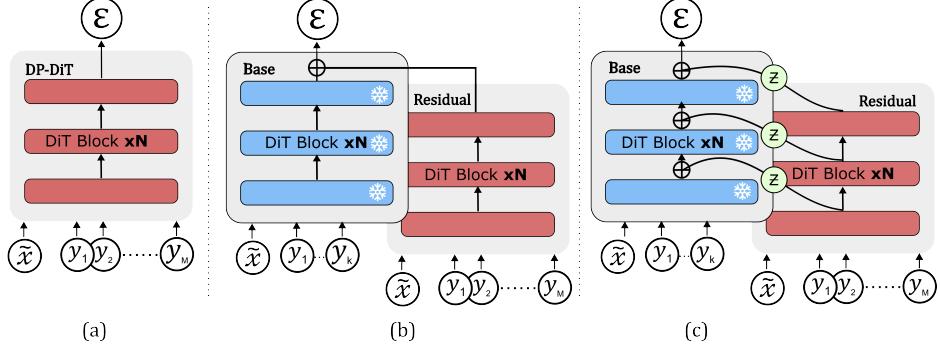


Figure 2: Architectural representations for [a] diffusion policy that jointly conditions on all observational modalities, [b] simple FDP architecture that composes the score outputs from π_{base} and π_{res} and [c] FDP architecture with block-wise composition with a zero layer Z applied on π_{res} .

learning a residual for the final score output, π_{res} learns the blockwise residual over the intermediate outputs of the frozen π_{base} . Specifically, let $\mathcal{F}_{\text{base}}^i$ and $\mathcal{F}_{\text{res}}^i$ denote the i -th DiT block outputs of the base and residual models, respectively. Then the composed output at level i can be written as $\mathcal{F}_{\text{base}}^i(x', y'^{1:k}) + \mathcal{Z}(\mathcal{F}_{\text{res}}^i(x', y'^{1:M}))$, where x' and $y'^{1:M}$ are layer inputs. Similar to Zhang et al. [33], \mathcal{Z} is a zero-initialized layer to avoid harmful updates at the start of the training and to ensure that gradient updates to the residual model improve the predictions of the composed model over π_{base} . Crucially, we find that preserving the diversity of π_{base} is essential: overfitting the base model leaves little residual signal to learn, reducing generalization. Our experiments show that selecting the π_{base} checkpoint with the lowest validation loss provides a good foundation for residual learning. Our residual model is structured following the Vision Transformers architecture [34]. In π_{res} , all observational modalities are passed through self-attention layers after encoding. Our visual residual model encodes camera images into a single patch to reduce computational overhead. Complete implementation details and architectural ablations are provided in Appendix D and H respectively.

176

177 4 Simulation Experiments

178 We train and evaluate FDP and related baselines in ten tasks of RL Bench [7] and four tasks of Adroit
 179 [35] and Robomimic [36] each. RL Bench is a large-scale simulation benchmark for robotic manipulation,
 180 offering multiple sensory modalities. For our experiments, we use several predefined single-task environments,
 181 and train policies on data collected using an in-built-planner with joint states as the action space. To study the impact of distribution shifts in observational modalities, we also evaluate
 182 the trained policies on six modified RL Bench environments that introduce appearance changes and visual distractors.
 183 Adroit [35, 8] is a simulated 24-DoF anthropomorphic hand platform widely used for dexterous manipulation research.
 184 Each task provides a low-dimensional hand-crafted environment state alongside rich proprioceptive feedback, which is used to train policies for dynamic manipulation.
 185 RoboMimic [36] provides human demonstration data collected across four RoboSuite [37] environments,
 186 featuring tasks with varying levels of precision and horizon length. More details in the Appendix E and our webpage <https://fdp-policy.github.io/fdp-policy/>.
 187

188 **Baselines.** For evaluation of sample efficiency in visuomotor tasks, we compare against several
 189 approaches that differ in the way in which they probabilistically model generative policy learning.
 190 However, for all approaches, we choose DiT-small ($\sim 90M$) [32] as our model architecture. We
 191 implement Diffusion Policy [1] using DiT, referred to as DP-DiT in our results. For comparison, we
 192 also include UNet [31] implemented by Chi et al. [1] in our baselines as DP-UNet. We reformulate
 193 POCO [24] to compose the modalities of proprioception y^r and vision y^c . We train the motion π_{base}
 194 and vision π_{res} models independently, prior to sampling from the composed distribution [30] using
 195 $\epsilon(x_t, y^r, y^c, t) = \hat{\epsilon}_\phi(x_t, y^r, y^c, t) + \lambda * \hat{\epsilon}_\theta(x_t, y^r, t)$. Here, $\lambda = 0.1$ based on POCO’s ablations [24].
 196 We also report results for classifier-free guidance [38] as CFG, where we train a single model and
 197 198

199 switch out the vision modality with a probability of 0.2. We then sample using $\epsilon(x_t, y^r, y^c, t) =$
200 $\lambda_1 * \hat{\epsilon}_\theta(x_t, y^r, y^c, t) + \lambda_2 * \hat{\epsilon}_\theta(x_t, y^r, \phi, t)$, where we set $\lambda_1 = 1.1$ and $\lambda_2 = 0.1$, as suggested by
[38]. For real-world and distractor experiments in simulation, we compare against DP-DiT.

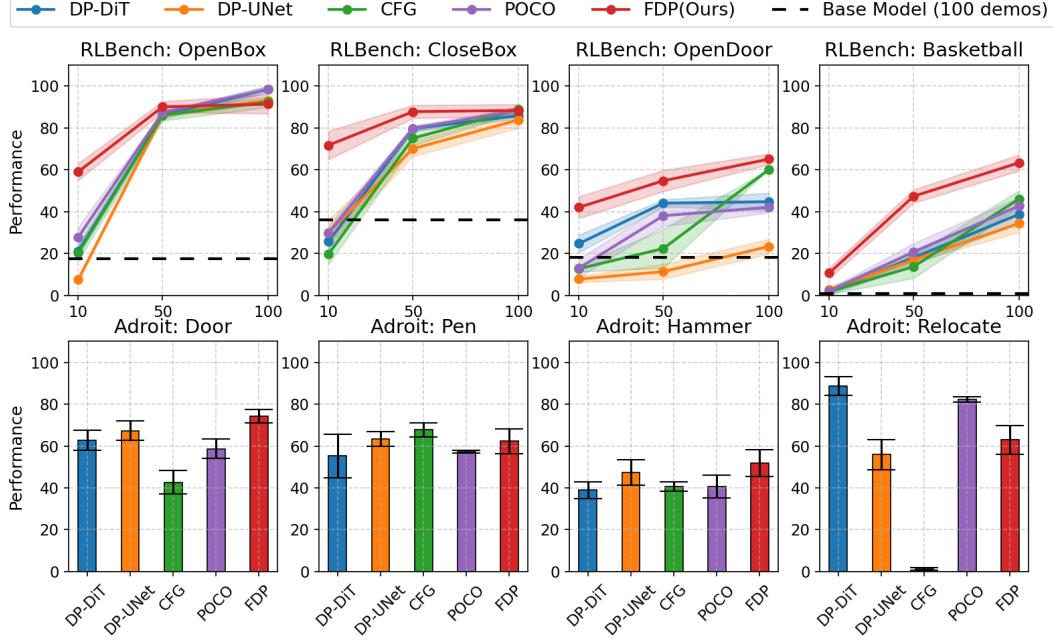


Figure 3: Evaluation for FDP and baselines on four RLBench and four Adroit tasks. FDP results in sample efficient policies at low number of demonstrations. More plots in Appendix F.

201
202 **Research Question 1: Can task-specific prioritization of modalities using FDP lead to sam-**
203 **ple efficiency gains in learning visuomotor tasks?** Prioritization of proprioception using the FDP
204 framework outperforms all baselines in four tasks each of RLBench and Adroit, across different
205 number of demonstrations as shown in Figure 3. In RLBench, FDP achieves 20% higher per-
206 formance on average with 10 demonstrations and 10% higher performance on average with 100
207 demonstrations over the strongest baseline. FDP results in sample-efficient policies, especially with
208 low number of demonstrations as the model is able to attend strongly to proprioception, only learn-
209 ing a residual for the visual observations. Note that a fully trained π_{base} motion model on 100
210 demonstrations fails in these tasks in isolation, implying that vision is required to complete these
211 tasks successfully. In the Adroit environments, proprioception prioritized FDP performs better by
212 $\sim 10\%$ on average over DP-DiT for Door, Pen and Hammer tasks. For the Relocate task, which
213 involves grasping a ball placed randomly on the table and relocating it to a random goal location, the
214 robot action is strongly influenced by the environment state specifying the ball and the goal location.
215 Hence, learning a motion π_{base} does not work in favor of improving policy performance, as the effect
216 of the state of the environment is learned as a residual over π_{base} . Prioritization of proprioception
217 will lead to sample-efficiency in repetitive tasks. Tasks that correlate heavily with robot propriocep-
218 tion are not uncommon as the robot is solving them in the first person view, and can move close to
219 the object if required. More results on these and six other RLBench tasks in Appendix F.

220 **Research Question 2: Does learning the visual modality as a residual in the FDP framework**
221 **result in robustness to distractors and appearance changes?** We present the results of the pol-
222 icy evaluations in the distractor environments in Table 1.1. Both DP-DiT and FDP are trained
223 on 100 demonstrations collected in the original environment and evaluated in three settings: the
224 original environment, an environment augmented with distractors, and an environment with visual
225 modifications to the manipulated objects. FDP significantly outperforms DP-DiT in both distractor-
226 augmented and visually modified environments by more than 40%. Additionally, we further collect

227 five demonstrations in each modified environment to investigate the benefits of few-shot adaptation
 228 to out-of-distribution data. Notably, FDP responds more effectively to additional demonstrations in
 229 the modified settings, improving its performance by 15% on average over DP’s 10%. In particu-
 230 lar, FDP updates only the residual model π_{res} with new demonstrations, adjusting the conditional
 231 distribution on visual modalities $p(y^c|x, y^r)$ without modifying the full conditional action distribu-
 232 tion $p(x|y^r, y^c)$. We extend this setting to point clouds and learn a visual residual on DP3 [39], as
 233 compared to DP3 with RGB inputs. Point clouds are sample-efficient for policy learning as they
 234 effectively encode the geometric structure of the scene in a single modality [39, 40, 41]. However,
 235 our distractor experiments show that FDP with a visual residual learned over DP3 is $\sim 20\%$ more
 236 performant than DP3 with RGB inputs. Further experimental details are in Appendix F.

Table 1.1: Robustness to Visual Distractors (100 demos)
(FDP significantly improves generalization to visual changes.)

Task	Environment	DP-DiT	FDP
OpenBox	Original	98.3 \pm 1.5	91.3 \pm 4.5
	Zero-shot color	43.3 \pm 2.5	46.7 \pm 1.5
	5 demos color	34.7 \pm 3.5	76.7 \pm 0.6
	Zero-shot distractors	1.7 \pm 2.1	16.7 \pm 2.1
	5 demos distractors	42.3 \pm 4.0	53.3 \pm 2.3
Basketball in Hoop	Original	38.7 \pm 4.2	63.3 \pm 3.8
	Zero-shot color	29.0 \pm 8.7	63.0 \pm 1.0
	5 demos color	13.0 \pm 0.0	45.0 \pm 2.6
	Zero-shot distractors	0.7 \pm 1.2	56.3 \pm 3.2
	5 demos distractors	2.7 \pm 1.2	39.7 \pm 4.9
Open Door	Original	44.7 \pm 4.2	65.0 \pm 2.6
	Zero-shot color1	0.0 \pm 0.0	14.3 \pm 3.1
	5 demos color1	17.7 \pm 3.1	52.0 \pm 7.2
	Zero-shot color2	0.3 \pm 0.6	30.7 \pm 2.5
	5 demos color2	20.0 \pm 5.2	53.7 \pm 3.8

Table 1.2: Block Pick Success Rates

(FDP performs better in tasks with less variation.)

Variations	Model	10 demos	50 demos	Distractors
Small	FDP	73.7 \pm 3.8	98.7 \pm 1.5	99.3 \pm 0.6
	DP-DiT	29.7 \pm 3.1	95.3 \pm 3.2	0.0 \pm 0.0
Medium	FDP	21.3 \pm 3.5	55.0 \pm 2.6	58.3 \pm 3.1
	DP-DiT	12.0 \pm 1.0	69.0 \pm 7.0	2.0 \pm 1.0
Large	FDP	6.3 \pm 3.1	20.3 \pm 3.5	0.7 \pm 1.2
	DP-DiT	3.3 \pm 0.6	45.7 \pm 7.1	0.0 \pm 0.0

Table 1.3: Robomimic Lowdim Task Success Rates (100 demos)

(Evaluating FDP at precise and long-horizon manipulation.)

Task	DP-DiT	CFG	POCO	FDP
Lift	99.0 \pm 1.7	98.7 \pm 0.6	98.7 \pm 1.5	99.7 \pm 0.6
Can	99.0 \pm 1.0	98.7 \pm 1.5	98.7 \pm 1.5	99.7 \pm 0.6
Square	80.3 \pm 4.6	80.0 \pm 3.0	76.3 \pm 6.0	58.0 \pm 6.6
Toolhang	60.0 \pm 7.5	60.7 \pm 3.5	55.7 \pm 7.5	45.7 \pm 3.8

Table 1: Tests for robustness and the effects of factorization across domains.

237 **Research Question 3: How sensitive is the task performance to prioritization perception for vi-
 238 suomotor tasks?** In visuomotor tasks, prioritizing proprioception over learning the full conditional
 239 distribution can be advantageous when the object placement diversity is low or when robustness to
 240 visual distribution shifts is critical. To demonstrate this, we construct three RLBench environments
 241 for a block-picking task, each featuring an increasingly larger object placement area. The results
 242 are presented in Table 1.2. As expected, FDP is sample-efficient and outperforms DP-DiT across all
 243 variation scales with only 10 demonstrations. However, with increasing number of demonstrations,
 244 DP-DiT eventually surpasses FDP at larger variation scales. Notably, DP-DiT still fails when visual
 245 distractors are introduced, whereas FDP remains robust and outperforms DP-DiT even at higher
 246 scales of task-variation in distractor-augmented environments. We also evaluated the performance
 247 of DP-DiT and FDP for fine-manipulation tasks on the Robomimic dataset. Although FDP is more
 248 sample efficient for the tasks of Lift and Can, it achieves a lower success rate than DP-DiT for
 249 Square and Toolhang, as shown in Table 1.3. This is to be expected, as fine-manipulation tasks
 250 present a bottleneck in the joint state-action distribution, and FDP factorizes the distribution into
 251 components where some modalities are learned as residuals over the others.

252 5 Real-world Experiments

253 We evaluate FDP and the DP-DiT baseline across four real-world domains and report their task suc-
 254 cess rates. The domains are – *Close Drawer* as a simple task where the robot has to push the drawer;
 255 *Put Block in Bowl* that assesses the policy’s ability to perform precise pick-and-place actions; *Pour*
 256 in *Bowl* to evaluate the policy’s dexterity in operating near joint limits and *Fold Towel* to assess
 257 effectiveness in manipulating deformable objects.

258 We collect 50 demonstrations per domain on a Franka FR3 robot using a 6D space mouse, recording
 259 both proprioceptive and visual observations from two cameras—one mounted on the gripper and a

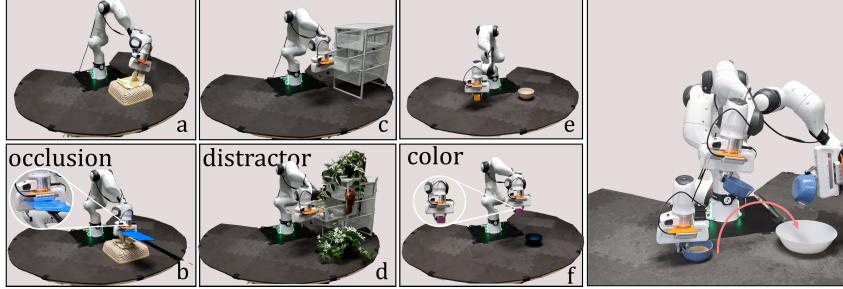


Figure 4: Task domains and their variations. In **occlusion**, the visual input is blocked using a board; **distractor**, flower pots and toys are introduced into the scene; and in **color**, the color of the manipulated object is altered during evaluation.

260 static camera covering the workspace. The trained policies are evaluated on four task variations in
 261 each domain: (a) **default**: an in-distribution setup matching the conditions used during demon-
 262 stration collection; (b) **color**: the object’s color is altered to test robustness to visual appearance
 263 changes; (c) **distractor**: novel, unseen objects such as vegetation props and soft toys are added
 264 to the scene to introduce clutter; and (d) **occlusion**: visual input is intermittently blocked during
 265 policy rollout to simulate partial observability. Figure 4 shows different task domains and their vari-
 266 ations used in our experiments. More details on the robot system setup can be found in Appendix G.
 267 We use 10 rollouts in each experiment and report the task success rate as shown in Table 2.

Task Domain	default		color		distractor		occlusion	
	DP	FDP	DP	FDP	DP	FDP	DP	FDP
Close Drawer	90	90	90	90	10	80	0	80
Put Block in Bowl	80	80	0	60	0	60	10	60
Pour in Bowl	70	80	40	80	20	60	10	50
Fold Towel	40	60	40	70	30	70	10	50

Table 2: Success rates (%) of Diffusion Policies (DP) and Factorized Diffusion Policies (FDP) across real-world tasks with 10 rollouts for each condition.

268 **Result Analysis.** We find that FDP is robust to distribution shifts in the environment. DP regularly
 269 produces unachievable robot actions under **distractor** and **occlusion** settings, often triggering
 270 safety stops, resulting in task failure. In contrast, FDP guided by its motion prior, consistently gen-
 271 erates stable actions even under severe occlusions and cluttered scenes, yielding an average absolute
 272 performance improvement of 40% over DP. In the **default** experiment we observe that the FDP
 273 policy outperforms DP in the pouring and towel-folding tasks, which require precise object manip-
 274 ulation. With just 50 demonstrations, DP overfits in these fine-grained tasks due to limited motion
 275 diversity, whereas FDP, leverages its residual guidance and effectively learns robust policies.

276 6 Conclusion

277 We present Factorized Diffusion Policies (FDP), a novel theoretical framework for prioritization of
 278 observation modalities in policy learning. We provide probabilistic grounding for diffusion policy
 279 learning and reveal the pitfalls of learning a full conditional on all the observational modalities. FDP
 280 decouples the modalities and proposes a framework for their selective prioritization. We derive a
 281 novel loss function to realize the decoupling of modalities and support it with a novel architecture
 282 for efficient training. Through extensive experiments, we demonstrate several benefits of modality
 283 prioritization, including improved sample efficiency and increased robustness to visual distractors
 284 and camera occlusions when learning a residual for vision. FDP opens new avenues for future
 285 research, such as scalable integration of diversely sourced observational modalities for robot policy
 286 learning. Finally, our real-world experiments highlight that FDP maintains strong performance even
 287 under significant visual disruptions, outperforming diffusion policies by over 40%.

288 **7 Limitations**

289 FDP is a theoretical framework to decouple observational modalities for robot policy learning. Pre-
290 dominantly, we present results for visuomotor tasks, but our method is generic and can be extended
291 to other modalities. We see the following issues with our FDP framework –

292 **7.1 Prioritizing Modalities**

293 We focus on the benefits and pitfalls of prioritizing proprioception and alternatively learning a residual
294 for vision in our experiments. However, for applications in a broader scope, the choice of
295 modalities to be prioritized will need to be studied and is not answered in this work. Understanding
296 which modality to prioritize for a particular task can be a challenging question for diffusion models
297 in general and might be severely task dependent. This might indicate that there might be an inference
298 time choice of composing modality that an agent might have to make.

299 **7.2 Factorizing modalities is not a general solution for all tasks**

300 As we show in our experiments, decoupling modalities may not be the right choice for every task or
301 skill. This is because some tasks require the full joint distribution of observations. We also want to
302 point out here that it is challenging to know whether a task requires the full joint distribution or the
303 factored distribution in a specific prioritized order. Future approaches could learn to automatically
304 attend on the right modality or the joint distribution, much like humans do.

305 **7.3 Architectural Choices**

306 Moreover, the framework can also benefit from further architectural improvements that realize a
307 better trade-off between the strength of the guidance imparted by the residual model π_{res} and its
308 robustness to perturbations in its inputs. The current training setup requires a two-step process
309 for learning π_{base} and π_{res} that presents a time and computational overhead over training standard
310 diffusion policies which might be cumbersome at deployment.

311 **7.4 Baselines outside of Diffusion based policy models**

312 In this work we only compare to diffusion based policy models as we are attempting to improve
313 their robustness and extend their capabilities of factorization. However, a large scale comparison
314 against different type of policy models is left to be done. For now we do not think comparisons
315 against non-diffusion policy types is critical but it is desirable to understand when to use which type
316 of policy for a robot.

317 **7.5 Large Vision Action Models**

318 There are large scale vision action models that can perform tasks specified by language in visual
319 environments sometimes even zero shot. Here we are studying how to learn individual skills using
320 diffusion based behavior cloning approaches. These skills might then be used in a larger stack of a
321 vision-action model. The question of sample efficiency and robustness to distractors will always be
322 important independent of the scale of the models themselves.

323 **References**

- 324 [1] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion
325 policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics*
326 *Research*, page 02783649241273668, 2023.
- 327 [2] Q. Zhang, Z. Liu, H. Fan, G. Liu, B. Zeng, and S. Liu. Flowpolicy: Enabling fast and robust
328 3d flow-based policy via consistency flow matching for robot manipulation. In *Proceedings of*
329 *the AAAI Conference on Artificial Intelligence*, volume 39, pages 14754–14762, 2025.
- 330 [3] A. Prasad, K. Lin, J. Wu, L. Zhou, and J. Bohg. Consistency policy: Accelerated visuomotor
331 policies via consistency distillation. *arXiv preprint arXiv:2405.07503*, 2024.
- 332 [4] K. Chen, E. Lim, K. Lin, Y. Chen, and H. Soh. Don’t start from scratch: Behavioral refinement
333 via interpolant-based policy diffusion. *arXiv preprint arXiv:2402.16075*, 2024.
- 334 [5] B. Wahn and P. König. Is attentional resource allocation across sensory modalities task-
335 dependent? *Advances in cognitive psychology*, 13(1):83, 2017.
- 336 [6] M. O. Ernst and M. S. Banks. Humans integrate visual and haptic information in a statistically
337 optimal fashion. *Nature*, 415(6870):429–433, 2002.
- 338 [7] S. James, Z. Ma, D. R. Arrojo, and A. J. Davison. Rlbench: The robot learning benchmark &
339 learning environment. *CoRR*, abs/1909.12271, 2019. URL <http://arxiv.org/abs/1909.12271>.
- 341 [8] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. D4rl: Datasets for deep data-driven
342 reinforcement learning, 2020.
- 343 [9] J. Sohl-Dickstein, E. A. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised
344 learning using nonequilibrium thermodynamics, 2015.
- 345 [10] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models, 2020. URL <https://arxiv.org/abs/2006.11239>.
- 347 [11] C. Luo. Understanding diffusion models: A unified perspective, 2022. URL <https://arxiv.org/abs/2208.11970>.
- 349 [12] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based
350 generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*,
351 2020.
- 352 [13] Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution,
353 2020.
- 354 [14] P. Vincent. A connection between score matching and denoising autoencoders. *Neural
355 computation*, 23(7):1661–1674, 2011.
- 356 [15] B. D. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their
357 Applications*, 12(3):313–326, 1982.
- 358 [16] A. Hyvärinen and P. Dayan. Estimation of non-normalized statistical models by score matching.
359 *Journal of Machine Learning Research*, 6(4), 2005.
- 360 [17] Y. Song and D. P. Kingma. How to train your energy-based models, 2021.
- 361 [18] C.-H. Chao, W.-F. Sun, B.-W. Cheng, Y.-C. Lo, C.-C. Chang, Y.-L. Liu, Y.-L. Chang, C.-P.
362 Chen, and C.-Y. Lee. Denoising likelihood score matching for conditional score-based data
363 generation, 2022. URL <https://arxiv.org/abs/2203.14206>.

- 364 [19] G. O. Roberts and R. L. Tweedie. Exponential convergence of langevin distributions and their
 365 discrete approximations. 1996.
- 366 [20] G. O. Roberts and J. S. Rosenthal. Optimal scaling of discrete approximations to langevin
 367 diffusions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 60(1):
 368 255–268, 1998.
- 369 [21] M. Reuss, M. Li, X. Jia, and R. Lioutikov. Goal-conditioned imitation learning using score-
 370 based diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.
- 371 [22] W. Liu, J. Mao, J. Hsu, T. Hermans, A. Garg, and J. Wu. Composable part-based manipulation.
 372 *arXiv preprint arXiv:2405.05876*, 2024.
- 373 [23] Z. Yang, J. Mao, Y. Du, J. Wu, J. B. Tenenbaum, T. Lozano-Pérez, and L. P. Kaelbling. Com-
 374 positional diffusion-based continuous constraint solvers. *arXiv preprint arXiv:2309.00966*,
 375 2023.
- 376 [24] L. Wang, J. Zhao, Y. Du, E. H. Adelson, and R. Tedrake. Poco: Policy composition from and
 377 for heterogeneous robot learning. *arXiv preprint arXiv:2402.02511*, 2024.
- 378 [25] T. Yu, T. Xiao, A. Stone, J. Tompson, A. Brohan, S. Wang, J. Singh, C. Tan, J. Peralta,
 379 B. Ichter, et al. Scaling robot learning with semantically imagined experience. *arXiv preprint*
 380 *arXiv:2302.11550*, 2023.
- 381 [26] Z. Chen, S. Kiami, A. Gupta, and V. Kumar. Genaug: Retargeting behaviors to unseen situa-
 382 tions via generative augmentation. *arXiv preprint arXiv:2302.06671*, 2023.
- 383 [27] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models, 2022. URL <https://arxiv.org/abs/2010.02502>.
- 385 [28] X. Liu, K. Y. Ma, C. Gao, and M. Z. Shou. Diffusion models in robotics: A survey. 2025.
- 386 [29] P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis, 2021. URL <https://arxiv.org/abs/2105.05233>.
- 388 [30] Y. Du, C. Durkan, R. Strudel, J. B. Tenenbaum, S. Dieleman, R. Fergus, J. Sohl-Dickstein,
 389 A. Doucet, and W. Grathwohl. Reduce, reuse, recycle: Compositional generation with energy-
 390 based diffusion models and mcmc, 2023.
- 391 [31] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical im-
 392 age segmentation. In *Medical image computing and computer-assisted intervention–MICCAI*
 393 *2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part*
 394 *III 18*, pages 234–241. Springer, 2015.
- 395 [32] W. Peebles and S. Xie. Scalable diffusion models with transformers, 2023. URL <https://arxiv.org/abs/2212.09748>.
- 397 [33] L. Zhang, A. Rao, and M. Agrawala. Adding conditional control to text-to-image diffusion
 398 models, 2023. URL <https://arxiv.org/abs/2302.05543>.
- 399 [34] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. De-
 400 hghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transform-
 401 ers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 402 [35] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine.
 403 Learning complex dexterous manipulation with deep reinforcement learning and demonstra-
 404 tions. *arXiv preprint arXiv:1709.10087*, 2017.
- 405 [36] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese,
 406 Y. Zhu, and R. Martín-Martín. What matters in learning from offline human demonstrations
 407 for robot manipulation. *arXiv preprint arXiv:2108.03298*, 2021.

- 408 [37] Y. Zhu, J. Wong, A. Mandlekar, R. Martín-Martín, A. Joshi, S. Nasiriany, Y. Zhu, and K. Lin.
409 robosuite: A modular simulation framework and benchmark for robot learning. In *arXiv*
410 preprint arXiv:2009.12293, 2020.
- 411 [38] J. Ho and T. Salimans. Classifier-free diffusion guidance, 2022. URL <https://arxiv.org/abs/2207.12598>.
- 413 [39] Y. Ze, G. Zhang, K. Zhang, C. Hu, M. Wang, and H. Xu. 3d diffusion policy: Generalizable
414 visuomotor policy learning via simple 3d representations. *arXiv preprint arXiv:2403.03954*,
415 2024.
- 416 [40] H. Zhu, Y. Wang, D. Huang, W. Ye, W. Ouyang, and T. He. Point cloud matters: Rethinking
417 the impact of different observation spaces on robot learning. *Advances in Neural Information
418 Processing Systems*, 37:77799–77830, 2024.
- 419 [41] T.-W. Ke, N. Gkanatsios, and K. Fragkiadaki. 3d diffuser actor: Policy diffusion with 3d scene
420 representations. *arXiv preprint arXiv:2402.10885*, 2024.
- 421 [42] D. Wang, S. Hart, D. Surovik, T. Kelestemur, H. Huang, H. Zhao, M. Yeatman, J. Wang,
422 R. Walters, and R. Platt. Equivariant diffusion policy. *arXiv preprint arXiv:2407.01812*, 2024.
- 423 [43] J. Yang, Z.-a. Cao, C. Deng, R. Antonova, S. Song, and J. Bohg. Equibot: Sim (3)-equivariant
424 diffusion policy for generalizable and data efficient learning. *arXiv preprint arXiv:2407.01479*,
425 2024.
- 426 [44] S. Haldar and L. Pinto. Point policy: Unifying observations and actions with key points for
427 robot manipulation. *arXiv preprint arXiv:2502.20391*, 2025.
- 428 [45] M. Torne, A. Simeonov, Z. Li, A. Chan, T. Chen, A. Gupta, and P. Agrawal. Reconciling reality
429 through simulation: A real-to-sim-to-real approach for robust manipulation. *arXiv preprint
430 arXiv:2403.03949*, 2024.
- 431 [46] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu. Viola: Imitation learning for vision-based manipulation
432 with object proposal priors. In *Conference on Robot Learning*, pages 1199–1210. PMLR,
433 2023.
- 434 [47] C. Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- 435 [48] P. Dhariwal and A. Nichol. Diffusion models beat gans on image synthesis. *Advances in neural
436 information processing systems*, 34:8780–8794, 2021.
- 437 [49] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>.