
Short-to-Long Distillation: Learning Long-Context Policies from Short-Context Supervision

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Consistency and reactivity are two essential properties for robotic policies. Yet,
2 recent methods often trade one for the other: using long action chunks and nu-
3 merous denoising steps to improve intra-chunk consistency, at the cost of lower
4 inference frequency and slower inference speed. In this paper, we first revisit the
5 necessity of these design choices through the lens of data scaling. We find that
6 with sufficient training data, extending history action contexts can substitute for
7 future action chunks, without compromising performance; moreover, conditioning
8 on longer contexts reduces action ambiguity, lessening the need for iterative de-
9 noising. Motivated by these observations, we introduce Short-to-Long Distillation,
10 a policy distillation approach that learns a long-context few-step student policy
11 from synthetic data generated by a short-context many-step teacher policy. Central
12 to our approach are two data curation strategies: (i) on-policy noise injection to
13 broaden the coverage of action contexts, and (ii) mode-seeking chunk optimization
14 to sharpen the distribution of action labels. Empirically, our method achieves
15 strong results on diffusion policies across Push-T and RoboMimic tasks. Notably,
16 using only 1k distilled sequences, the student policies match their teachers in static
17 settings and surpass them by up to 40% in stochastic environments. Our results
18 suggest the promise of synthetic data as a scalable alternative to inductive biases
19 for robot learning.

20 1 Introduction

21 Humans solve physical tasks with both consistency and reactivity—producing coherent long-horizon
22 behaviors while adapting rapidly to unexpected changes. Robots should possess the same capabilities,
23 yet learning such policies from human demonstrations remains challenging. Often, collected demon-
24 stration datasets exhibit three properties: long temporal dependencies, large behavioral variations,
25 and limited state coverage, resulting in a highly multimodal yet sparsely covered distribution over
26 valid action sequences [7, 10]. To address these challenges, recent policies are typically built with: (i)
27 action chunking - predict and execute a chunk of future actions over multiple steps [6, 20], (ii) action
28 denoising – refine an action chunk from random noise through multiple steps [3, 6]. While these
29 designs improve consistency within a chunk, they inherently compromise reactivity: larger chunks
30 lower inference frequency and multi-step denoising slows inference speed.

31 To enhance reactivity, recent studies have explored three main approaches. First, guided infer-
32 ence improves cross-chunk consistency during inference [1, 4, 11], reducing the reliance on large
33 chunks. Second, hierarchical policies introduce low-level controllers beneath high-level chunk
34 predictors [2, 12, 19], enabling rapid local adaptation. Third, accelerated denoising reduces the
35 number of refinement steps [9, 15, 17]. Despite these advances, existing approaches remain tied to
36 multi-step chunking and denoising, imposing an inherent trade-off between long-term consistency
37 and short-term reactivity.

38 In this work, we first revisit the necessity of action chunking and action denoising. We hypothesize
39 that these design choices are not intrinsic requirements of robotic policies, but rather artifacts of

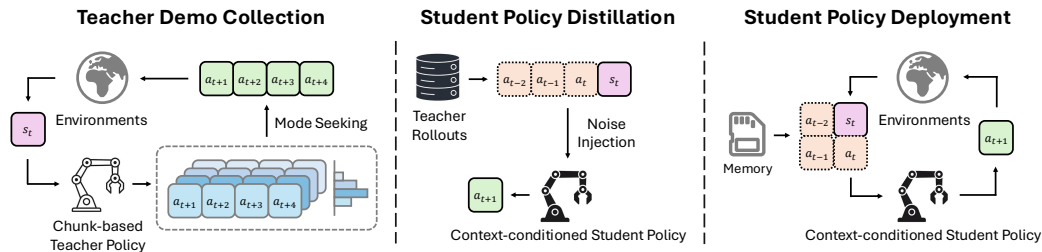


Figure 1: Illustration of Short-to-Long Distillation: A chunk-based teacher policy first interacts with the environment to generate demonstrations using mode seeking. The context-conditioned student policy is then trained on stored teacher rollouts, augmented with context noise injection to expand action coverage. At deployment, the student conditions on longer context actions while predicting shorter chunks, achieving higher reactivity without compromising performance or consistency.

40 limited data. We thus conduct a set of controlled experiments by scaling demonstration data up to
 41 $10\times$ larger per task than typical academic datasets. Interestingly, we find that with sufficient training
 42 data, extending history action context can replace long action chunks, without sacrificing performance.
 43 Moreover, conditioning on longer contexts reduces action ambiguity, substantially decreasing the
 44 need for iterative denoising. Despite these promising outcomes, acquiring human demonstrations at
 45 such a scale is costly and even impractical. This raises a natural question: Can we leverage existing
 46 policies, rather than human supervision, as data generators to reach this critical scale?

47 To this end, we introduce Short-to-Long Distillation (S2L), a policy distillation framework for learning
 48 performant long-context few-step policies. S2L proceeds in two stages: we first train a short-context,
 49 long-chunk, many-step teacher policy on human demonstrations to capture consistency, and then
 50 distill its knowledge into a long-context, short-chunk, few-step student policy using synthetic rollouts
 51 generated by the teacher. Unlike conventional distillation methods [8, 16, 18], S2L operates in a
 52 heterogeneous setting where teacher and student differ in both input and output spaces, precluding
 53 direct one-to-one supervision. To bridge this structural mismatch, we propose two data curation
 54 strategies: (i) on-policy noise injection to expand the coverage of action contexts, and (ii) mode-
 55 seeking chunk optimization to sharpen the distribution of action labels. Together, these strategies
 56 enable the learned student to inherit the consistency strength of the teacher policy while discarding
 57 its reliance on chunking and denoising, thereby achieving substantially higher reactivity.

58 The main contributions of this paper are twofold: an analysis revealing how data scale influences
 59 policy design (§2) and a distillation method for learning long-context, few-step policies from existing
 60 short-context multi-step ones (§3). Empirically, we evaluate our method on diffusion policies across
 61 Push-T and RoboMimic tasks. Using only 1k distilled sequences, the student policies achieve
 62 performance comparable to their teachers in static environments and exceed them by up to 40% in
 63 stochastic environments. Our results suggest the promise of synthetic data curation as a scalable
 64 alternative to inductive biases for robot learning.

65 2 Analysis: Effect of Data Scaling on Policy Design

66 **Problem Setup.** We study imitation learning from a finite set of demonstrations $\mathcal{D} = \{(s_t, a_t)\}_{t=1}^T$.
 67 At each time step t , the policy conditions on the current observation s_t and the previous c actions to
 68 predict a chunk of h future actions: $\pi_{\theta}(a_{t:t+h-1} \mid s_t, a_{t-c:t-1})$. We assume the chunk distribution is
 69 modeled via an *iterative refinement* process (e.g., diffusion or flow-matching), requiring k denoising
 70 steps per chunk prediction.

71 Recent policies often favor *short* contexts (c), *long* chunks (h), and *many* refinement steps (k) to
 72 improve intra-chunk consistency. However, this design reduces reactivity in two ways: (i) the decision
 73 frequency decreases as h grows, and (ii) inference latency increases with k . In principle, an ideal
 74 closed-loop regime would invert this design—*long* context, *short* chunks (up to token-by-token),
 75 and *few* refinement steps—as commonly seen in language modeling. Yet, applying this strategy in
 76 robotics often hurts performance, arguably due to the comparatively limited data scale.

77 *How does the optimal design of robot policies evolve as demonstration data increases?*

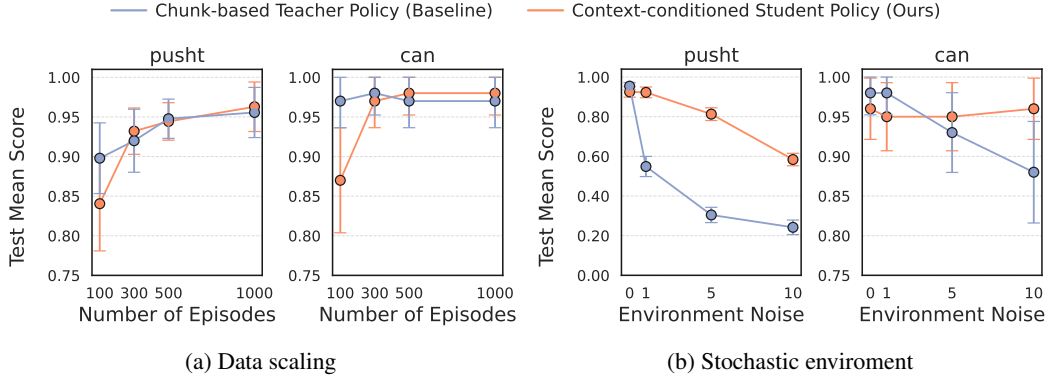


Figure 2: Effect of dataset size and context conditioning. (a) Increasing the demonstration dataset shifts the optimal design toward longer context lengths ($c \uparrow$) and shorter prediction horizons ($h \downarrow$). (b) Context-conditioned student policies surpass chunk-based teacher policies as the environment noise level increases.

78 **Scaling Analysis.** Existing robot learning benchmarks typically provide only a few hundred demon-
 79 strations per task. We hypothesize that this data scale is limiting, constraining the optimal policy
 80 design and biasing current methods toward short-context, chunk-heavy, many-step designs. To
 81 understand how the data scale influences the design space, we conduct a controlled scaling analysis
 82 on two widely used tasks, *Push-T* and *Can*, varying the amount of demonstration data up to $10\times$ the
 83 original scale. For each data scale, we keep the total modeling length $l = c + h$ fixed, but vary the
 84 context length $c \in \{0, 2, 4, 6, 8\}$ and the number of refinement steps $k \in \{1, 3, 8, 15, 100\}$.

85 Interestingly, we observe that both *action chunking* and *iterative denoising*, which dominate recent
 86 policy designs [4, 5, 14, 17], become less critical as data increases. As shown in Fig. 2, policies with
 87 *longer context* and *fewer refinement steps* gradually match—and eventually surpass—the performance
 88 of chunk-based, many-step counterparts. This result suggests that *learning long-context, few-step*
 89 *policies is not infeasible—the bottleneck lies in data.*

90 3 Method: Short-to-Long Distillation for Policy Learning

91 Motivated by the analysis in §2, we next introduce a policy distillation method for learning performant
 92 policies conditioned on longer action contexts and operating with few denoising steps. While learning
 93 such a policy directly from limited human demonstrations is generally challenging, we seek to achieve
 94 this goal by progressively expanding the data in two stages: first, we train a teacher policy π_{chunk} on
 95 human demonstrations, using a long chunk horizon and many denoising steps to prioritize intra-chunk
 96 consistency; subsequently, we train a student policy π_{context} on synthetic rollouts generated by π_{chunk} ,
 97 conditioning on longer contexts and using fewer refinement steps to improve reactivity.

98 The core challenge lies in enabling the student not only to inherit the teacher’s consistency but
 99 also to *outperform* it in reactivity. This may seem out of reach at first glance. However, we notice
 100 two structural differences between the teacher and student policies that offer valuable guidance for
 101 effective distillation:

- 102 • *Smooth chunks vs. noisy contexts.* The teacher generates temporally consistent action chunks,
 103 whereas the student conditions on histories of individual actions, which can be much noisier.
- 104 • *Ambiguous samples vs. specified labels.* The teacher may yield diverse chunk samples, whereas
 105 the student operates under extended histories that make the next actions more specified.

106 These differences imply that, even if not all raw teacher rollouts are directly useful, they can
 107 potentially be turned into effective training signals for the student. Specifically, we consider two
 108 curation strategies tailored to history action contexts and future action labels, respectively.

109 **Noise-Corrupted Context Expansion.** Given a teacher rollout, we corrupt the action context with
 110 additive Gaussian noise and train the student on the perturbed history:

$$\tilde{a}_{t-c:t-1} = a_{t-c:t-1} + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I).$$

111 Random noise is sampled at every use so each iteration presents a distinct history; the student predicts
 112 $a_{t:t+h-1}$ from $(s_t, \tilde{a}_{t-c:t-1})$ with the teacher chunk as the target. This stochastic corruption expands

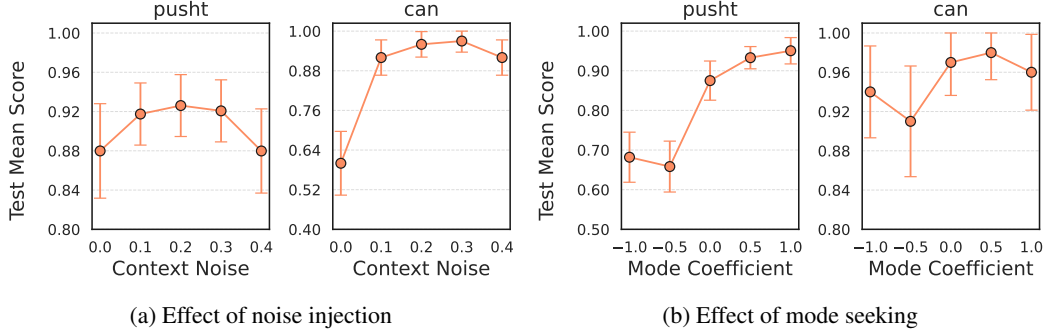


Figure 3: Effect of chunk curation. (a) Adding Gaussian noise to context actions improves student policy performance by expanding action coverage, though excessive noise eventually degrades results. (b) Increasing the degree of mode seeking during teacher demonstration enhances student performance, as clustering on denser modes provides sharper supervision.

113 context coverage and reduces overfitting to deterministic teacher behavior. The noise level σ is a
 114 hyperparameter and can be tuned based on a teacher–student discrepancy.

115 **Mode-Seeking Chunk Optimization.** Given the state and context at time t , let the teacher produce
 116 $N = 16$ candidate chunks $\mathcal{C}_t = \{\hat{a}_{t:t+h-1}^{(i)}\}_{i=1}^N$. We assign each candidate a density score $z_i =$
 117 $\phi(\hat{a}_{t:t+h-1}^{(i)})$ based on average distance to 5 nearest neighbors. We then use $d \in [-1, 1]$ controls *mode*
 118 *seeking* vs. *mode covering* with coupled truncation: for $d > 0$ keep only the densest top fraction; for
 119 $d < 0$ keep only the least-dense bottom fraction; for $d = 0$ keep all. We then select the supervision
 120 chunk from the kept set, favoring higher z_i when $d > 0$ and thereby sharpening the label distribution
 121 in the mode-seeking regime.

122 4 Experiments

123 In this section, we evaluate our method on two robotic manipulation benchmarks: Push-T [6] and
 124 RoboMimic [13]. Specifically, we use the public checkpoints from prior work [6, 11] as the teacher
 125 policies, and assess the effectiveness of the proposed curation strategies on the student performance.

126 4.1 Effect of Noise Injection

127 **Setup.** We study the effect of Gaussian noise injection on student policy performance by varying
 128 the context noise level $\sigma \in \{0.0, 0.1, 0.2, 0.3, 0.4\}$ in *Push-T* and *RoboMimic (Can)*. Student policies
 129 are trained on 500 teacher demonstrations collected with mode coefficient $d = 0.5$, using context
 130 length $c = 8$ and prediction horizon $h = 1$. All other factors (architecture, optimizer, learning rate,
 131 training steps) are held fixed.

132 **Results.** As shown in Fig. 3a, performance follows a non-monotonic trend with respect to σ .
 133 Moderate noise ($\sigma \leq 0.2$) improves generalization by expanding action coverage, while larger values
 134 degrade performance as injected noise dominates the context signal.

135 4.2 Effect of Mode Seeking

136 **Setup.** We then examine how the mode-seeking coefficient $d \in \{-0.5, 0.0, 0.5, 1.0\}$ in teacher
 137 demonstration collection affects student policy performance on *Push-T* and *RoboMimic (Can)*. Each
 138 teacher chunk is sampled $N = 16$ times, with $d < 0$ promoting diversity, $d = 0$ sampling uniformly,
 139 and $d > 0$ emphasizing high-density modes. Student policies are trained under the same conditions
 140 as above (500 demos, $c = 8$, $h = 1$), with all other variables controlled.

141 **Results.** As shown in Fig. 3b, *Push-T* shows consistent gains as d increases, highlighting the
 142 benefits of clustering on dense modes. *Can* exhibits a similar upward trend but with greater variance
 143 across runs. We attribute this noisiness to the binary nature of the task’s binary success metric, which
 144 may obscure finer-grained differences in policy quality.

145 **References**

- 146 [1] C. Agia, R. Sinha, J. Yang, Z.-a. Cao, R. Antonova, M. Pavone, and J. Bohg. Unpacking Failure
147 Modes of Generative Policies: Runtime Monitoring of Consistency and Progress. *arXiv preprint*
148 *arXiv:2410.04640*, Oct. 2024.
- 149 [2] L. Ankile, A. Simeonov, I. Shenfeld, M. Torne, and P. Agrawal. From Imitation to Refinement –
150 Residual RL for Precise Assembly. *arXiv preprint arXiv:2407.16677*, Dec. 2024.
- 151 [3] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman,
152 B. Ichter, S. Jakubczak, T. Jones, L. Ke, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch,
153 L. X. Shi, J. Tanner, Q. Vuong, A. Walling, H. Wang, and U. Zhilinsky. $\$ \pi_0 \$$: A Vision-
154 Language-Action Flow Model for General Robot Control. *arXiv preprint arXiv:2410.24164*,
155 Nov. 2024.
- 156 [4] K. Black, M. Y. Galliker, and S. Levine. Real-Time Execution of Action Chunking Flow
157 Policies, June 2025. URL <http://arxiv.org/abs/2506.07339>. arXiv:2506.07339 [cs].
- 158 [5] C. Cheang, S. Chen, Z. Cui, Y. Hu, L. Huang, T. Kong, H. Li, Y. Li, Y. Liu, X. Ma, H. Niu,
159 W. Ou, W. Peng, Z. Ren, H. Shi, J. Tian, H. Wu, X. Xiao, Y. Xiao, J. Xu, and Y. Yang. GR-3
160 Technical Report. *arXiv preprint arXiv:2507.15493*, July 2025.
- 161 [6] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion Policy:
162 Visuomotor Policy Learning via Action Diffusion. In *Robotics: Science and Systems XIX*.
163 Robotics: Science and Systems Foundation, July 2023. ISBN 978-0-9923747-9-2.
- 164 [7] O. X.-E. Collaboration, A. Padalkar, A. Pooley, A. Mandlekar, A. Jain, A. Tung, A. Bewley,
165 A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Singh, A. Garg, A. Brohan, A. Raffin, A. Wahid,
166 B. Burgess-Limerick, B. Kim, B. Schölkopf, B. Ichter, C. Lu, C. Xu, C. Finn, C. Xu, C. Chi,
167 C. Huang, C. Chan, C. Pan, C. Fu, C. Devin, D. Driess, D. Pathak, D. Shah, D. Büchler,
168 D. Kalashnikov, D. Sadigh, E. Johns, F. Ceola, F. Xia, F. Stulp, G. Zhou, G. S. Sukhatme,
169 G. Salhotra, G. Yan, G. Schiavi, G. Kahn, H. Su, H.-S. Fang, H. Shi, H. B. Amor, H. I.
170 Christensen, H. Furuta, H. Walke, H. Fang, I. Mordatch, I. Radosavovic, I. Leal, J. Liang,
171 J. Abou-Chakra, J. Kim, J. Peters, J. Schneider, J. Hsu, J. Bohg, J. Bingham, J. Wu, J. Wu,
172 J. Luo, J. Gu, J. Tan, J. Oh, J. Malik, J. Booher, J. Tompson, J. Yang, J. J. Lim, J. Silvério, J. Han,
173 K. Rao, K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund,
174 K. Kawaharazuka, K. Zhang, K. Rana, K. Srinivasan, L. Y. Chen, L. Pinto, L. Fei-Fei, L. Tan,
175 L. Ott, L. Lee, M. Tomizuka, M. Spero, M. Du, M. Ahn, M. Zhang, M. Ding, M. K. Srirama,
176 M. Sharma, M. J. Kim, N. Kanazawa, N. Hansen, N. Heess, N. J. Joshi, N. Suenderhauf,
177 N. Di Palo, N. M. M. Shafiullah, O. Mees, O. Kroemer, P. R. Sanketi, P. Wohlhart, P. Xu,
178 P. Sermanet, P. Sundaresan, Q. Vuong, R. Rafailov, R. Tian, R. Doshi, R. Martín-Martín,
179 R. Mendonca, R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Moore,
180 S. Bahl, S. Dass, S. Sonawani, S. Song, S. Xu, S. Haldar, S. Adebola, S. Guist, S. Nasiriany,
181 S. Schaal, S. Welker, S. Tian, S. Dasari, S. Belkhale, T. Osa, T. Harada, T. Matsushima, T. Xiao,
182 T. Yu, T. Ding, T. Davchev, T. Z. Zhao, T. Armstrong, T. Darrell, V. Jain, V. Vanhoucke, W. Zhan,
183 W. Zhou, W. Burgard, X. Chen, X. Wang, X. Zhu, X. Li, Y. Lu, Y. Chebotar, Y. Zhou, Y. Zhu,
184 Y. Xu, Y. Wang, Y. Bisk, Y. Cho, Y. Lee, Y. Cui, Y.-H. Wu, Y. Tang, Y. Zhu, Y. Li, Y. Iwasawa,
185 Y. Matsuo, Z. Xu, and Z. J. Cui. Open X-Embodiment: Robotic Learning Datasets and RT-X
186 Models. *arXiv preprint arXiv:2310.08864*, Dec. 2023.
- 187 [8] G. Hinton, O. Vinyals, and J. Dean. Distilling the Knowledge in a Neural Network. *arXiv*
188 *preprint arXiv:1503.02531*, Mar. 2015.
- 189 [9] S. H. Høeg, Y. Du, and O. Egeland. Streaming Diffusion Policy: Fast Policy Synthesis with
190 Variable Noise Diffusion Models. *arXiv preprint arXiv:2406.04806*, Oct. 2024.
- 191 [10] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany,
192 M. K. Srirama, L. Y. Chen, K. Ellis, P. D. Fagan, J. Hejna, M. Itkina, M. Lepert, Y. J. Ma,
193 P. T. Miller, J. Wu, S. Belkhale, S. Dass, H. Ha, A. Jain, A. Lee, Y. Lee, M. Memmel, S. Park,
194 I. Radosavovic, K. Wang, A. Zhan, K. Black, C. Chi, K. B. Hatch, S. Lin, J. Lu, J. Mercat,
195 A. Rehman, P. R. Sanketi, A. Sharma, C. Simpson, Q. Vuong, H. R. Walke, B. Wulfe, T. Xiao,
196 J. H. Yang, A. Yavary, T. Z. Zhao, C. Agia, R. Bajjal, M. G. Castro, D. Chen, Q. Chen, T. Chung,

- 197 J. Drake, E. P. Foster, J. Gao, D. A. Herrera, M. Heo, K. Hsu, J. Hu, D. Jackson, C. Le, Y. Li,
198 K. Lin, R. Lin, Z. Ma, A. Maddukuri, S. Mirchandani, D. Morton, T. Nguyen, A. O’Neill,
199 R. Scalise, D. Seale, V. Son, S. Tian, E. Tran, A. E. Wang, Y. Wu, A. Xie, J. Yang, P. Yin,
200 Y. Zhang, O. Bastani, G. Berseth, J. Bohg, K. Goldberg, A. Gupta, A. Gupta, D. Jayaraman, J. J.
201 Lim, J. Malik, R. Martín-Martín, S. Ramamoorthy, D. Sadigh, S. Song, J. Wu, M. C. Yip, Y. Zhu,
202 T. Kollar, S. Levine, and C. Finn. DROID: A Large-Scale In-The-Wild Robot Manipulation
203 Dataset. *arXiv preprint arXiv:2403.12945*, Mar. 2024.
- 204 [11] Y. Liu, J. I. Hamid, A. Xie, Y. Lee, M. Du, and C. Finn. Bidirectional Decoding: Improving
205 Action Chunking via Closed-Loop Resampling, Dec. 2024. URL [http://arxiv.org/abs/
206 2408.17355](http://arxiv.org/abs/2408.17355). arXiv:2408.17355 [cs].
- 207 [12] X. Ma, S. Patidar, I. Haughton, and S. James. Hierarchical Diffusion Policy for Kinematics-
208 Aware Multi-Task Robotic Manipulation. *arXiv preprint arXiv:2403.03890*, Mar. 2024.
- 209 [13] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese,
210 Y. Zhu, and R. Martín-Martín. What Matters in Learning from Offline Human Demonstrations
211 for Robot Manipulation. In *Proceedings of the 5th Conference on Robot Learning*, pages
212 1678–1690. PMLR, Jan. 2022.
- 213 [14] NVIDIA, J. Bjorck, F. Castañeda, N. Cherniadev, X. Da, R. Ding, L. J. Fan, Y. Fang, D. Fox,
214 F. Hu, S. Huang, J. Jang, Z. Jiang, J. Kautz, K. Kundalia, L. Lao, Z. Li, Z. Lin, K. Lin, G. Liu,
215 E. Llontop, L. Magne, A. Mandlekar, A. Narayan, S. Nasiriany, S. Reed, Y. L. Tan, G. Wang,
216 Z. Wang, J. Wang, Q. Wang, J. Xiang, Y. Xie, Y. Xu, Z. Xu, S. Ye, Z. Yu, A. Zhang, H. Zhang,
217 Y. Zhao, R. Zheng, and Y. Zhu. GR00T N1: An Open Foundation Model for Generalist
218 Humanoid Robots. *arXiv preprint arXiv:2503.14734*, Mar. 2025.
- 219 [15] A. Prasad, K. Lin, J. Wu, L. Zhou, and J. Bohg. Consistency Policy: Accelerated Visuomotor
220 Policies via Consistency Distillation. *arXiv preprint arXiv:2405.07503*, May 2024.
- 221 [16] A. A. Rusu, S. G. Colmenarejo, C. Gulcehre, G. Desjardins, J. Kirkpatrick, R. Pascanu, V. Mnih,
222 K. Kavukcuoglu, and R. Hadsell. Policy Distillation. *arXiv preprint arXiv:1511.06295*, Jan.
223 2016.
- 224 [17] M. Shukor, D. Aubakirova, F. Capuano, P. Kooijmans, S. Palma, A. Zouitine, M. Aractingi,
225 C. Pascal, M. Russi, A. Marafioti, S. Alibert, M. Cord, T. Wolf, and R. Cadene. SmolVLA:
226 A Vision-Language-Action Model for Affordable and Efficient Robotics. *arXiv preprint
227 arXiv:2506.01844*, June 2025.
- 228 [18] L. Wang and K.-J. Yoon. Knowledge Distillation and Student-Teacher Learning for Visual
229 Intelligence: A Review and New Outlooks. *IEEE Transactions on Pattern Analysis and Machine
230 Intelligence*, 44(6):3048–3068, June 2022.
- 231 [19] H. Xue, J. Ren, W. Chen, G. Zhang, Y. Fang, G. Gu, H. Xu, and C. Lu. Reactive Diffusion
232 Policy: Slow-Fast Visual-Tactile Policy Learning for Contact-Rich Manipulation. *arXiv preprint
233 arXiv:2503.02881*, Apr. 2025.
- 234 [20] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning Fine-Grained Bimanual Manipulation
235 with Low-Cost Hardware. *arXiv preprint arXiv:2304.13705*, Apr. 2023.

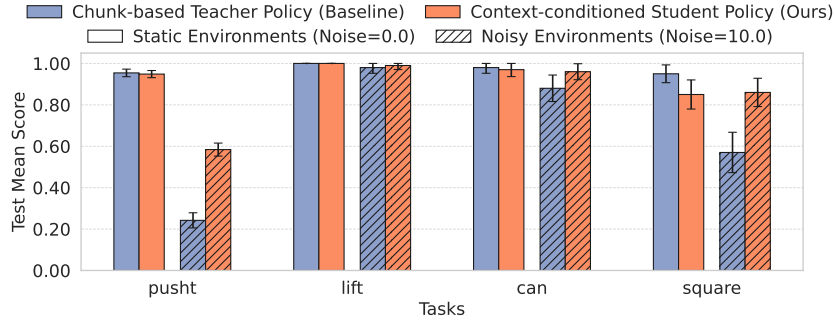


Figure 4: Context-conditioned student policies match the performance of chunk-based teacher policies in static environments, and surpass them in stochastic environments across four Push-T and RoboMimic tasks.

236 A Implementation Details

237 We parameterize the student policy with a U-Net backbone and augment it with FiLM layers for
 238 conditioning on observations, diffusion timesteps, and context actions following standard diffusion
 239 policy[6] architectures.

240 **Teacher demo collection.** During chunk-based teacher demonstration collection, we employ early
 241 stopping to mitigate mode collapse, which can otherwise reduce performance in mode-seeking
 242 settings. The chunk-based teacher policy has context length $c = 0$ and action horizon $h = 8$. We use
 243 16 samples and 3 modes, which is sufficient to capture the multimodal behavior of teacher policies in
 244 both Push-T and Robomimic tasks.

245 **Student policy distillation.** For context-conditioned student policy distillation, we set context
 246 length $c = 8$ and action horizon $h = 1$. Training uses a batch size of 256 with 100 diffusion steps for
 247 both training and inference. Rollouts are generated with a DDPM sampler.

248 Unless otherwise noted, all experiments share the same architecture, optimizer, diffusion scheduler,
 249 and saved teacher demo data. To ensure convergence and fair comparison, we train each model for up
 250 to 500 epochs and report results from the checkpoint achieving the best validation performance.

251 **Environment Setup.** We evaluate student policy performance on four robot manipulation bench-
 252 marks:

253 *Push-T:* We use the Push-T environment from Chi et al. [6], where the goal is to push a T-shaped
 254 block to a designated target location. The action space consists of two-dimensional end-effector
 255 velocity commands, and policies are conditioned on keypoint-based observations.

256 *Robomimic:* We adopt three tasks from the Robomimic suite [13] — Lift, Can, and Square. Policies
 257 are conditioned on state-based observations, with continuous action spaces following the suite’s
 258 default setup.

Name	Value
Context length	0
Action horizon	8
Number of samples	16
Number of modes	3
Mode coefficient	0.5

Table 1: Default hyperparameters for teacher demonstration collection.

Name	Value
Context length	8
Action horizon	1
Context noise	0.2
Number of teacher episodes	500
Batch size	256
Epochs	300
Learning rate	1×10^{-4}
Weight decay	1×10^{-6}
Optimizer	AdamW
Training diffusion steps	100
Inference diffusion steps	100
Diffusion scheduler	DDPM
Number of test environments	100
Environment noise	0.0

Table 2: Default hyperparameters for student policy distillation.