

Efficient and Scalable Diffusion Transformer Policies with Mixture of Expert Denoisers for Multitask Learning

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Diffusion Policies have become widely used in Goal-Conditioned Imitation Learning, offering several appealing properties, such as generating multimodal
2 and discontinuous behavior. As models are becoming larger to capture more
3 complex capabilities, their computational demands increase, as shown by recent
4 scaling laws. Therefore, continuing with the current architectures will present
5 a computational roadblock. To address this, we propose Mixture-of-Denoising
6 Experts (MoDE) as a novel policy for guided behavior generation. MoDE is able
7 to achieve competitive performance to current state-of-the-art dense transformer-
8 based Diffusion Policies while requiring fewer active parameters, reducing the
9 inference cost significantly. To achieve this, MoDE introduces a novel routing strategy that conditions the expert selection on the current noise level of the diffusion denoising process. MoDE achieves competitive or state of the art performance on
10 four established imitation learning benchmarks, including CALVIN and LIBERO.
11 In addition, we perform thorough ablations on the various components in MoDE.
12
13
14

15 1 Introduction

16 Diffusion models [1, 2], which learn to reverse a noise-adding process, have gained popularity as
17 policies for [Imitation Learning \(IL\)](#) [3, 4, 5] due to their ability to generate diverse behaviors [6]
18 and handle complex action spaces. However, their high computational cost, resulting from large
19 architectures and multiple denoising steps, limits their use in real-time robotics applications, especially
20 those with limited on-board computing resources. To address this challenge, we explore Mixture of
21 Experts (MoE) models that can scale model capacity while reducing computational requirements by
22 utilizing only a subset of parameters during each forward pass.

23 We introduce [Mixture-of-Denoising Experts Policy \(MoDE\)](#), a scalable and efficient [Mixture-of-Experts \(MoE\)](#) Diffusion Policy. Our work is inspired by prior results showcasing the multitask nature of the denoising process [7], where there is little transfer between the different phases in the denoising process. We present a novel noise-conditioned routing mechanism, that distributes tokens to our experts based on the current noise level. [MoDE](#) leverages noise-conditioned self-attention combined with a noise input token for enhanced noise-injection. Our proposed Policy surpasses previous Diffusion Policies with higher efficiency and demonstrates improved performance across
29 134 diverse tasks in challenging goal-conditioned imitation learning benchmarks: CALVIN [8] and
30 LIBERO [9]. Through comprehensive ablation studies, we investigate the impact of various design
31 decisions, including token routing strategies, noise-injection techniques, and expert distribution.
32

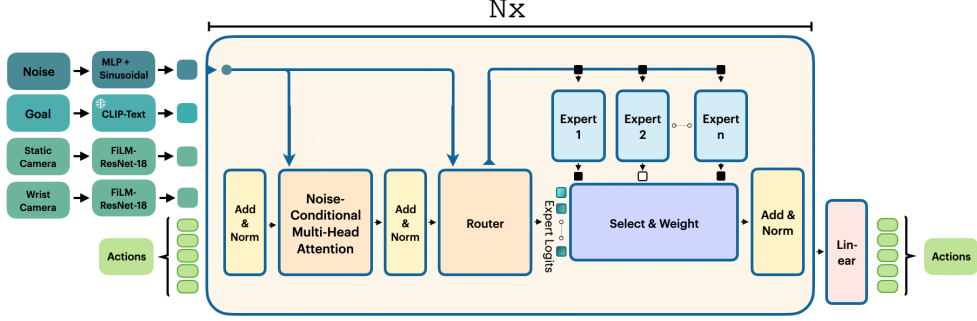


Figure 1: Overview of the proposed MoDE architecture. The model uses a transformer with causal masking from top to bottom. Each transformer block uses noise-conditional self-attention and is followed by a noise-conditioned router, that distributes tokens to specialized expert models conditioned on the current noise level. Each expert is a simple MLP with Swish-GLU activation.

33 2 Method

34 2.1 Problem Formulation

35 We consider the problem of learning a language-conditioned policy $\pi_\theta(\bar{\mathbf{a}}|\bar{\mathbf{s}}, \mathbf{g})$ given a dataset of
 36 robot demonstrations \mathcal{T} . The policy predicts a sequence of future actions $\bar{\mathbf{a}} = (\mathbf{a}, \dots, \mathbf{a}_{i+j-1})$ of
 37 length j , conditioned on the history of state embeddings $\bar{\mathbf{s}} = (\mathbf{s}_{i-h+1}, \dots, \mathbf{s}_i)$ of length h and a
 38 desired goal \mathbf{g} . The dataset contains $\tau \in \mathcal{T}$ trajectories, where trajectory consists of a sequence of
 39 triplets of state, actions, and goal $(\bar{\mathbf{s}}_i, \mathbf{a}_i, \mathbf{g})$. \mathbf{g} is a language instruction. Our policy is trained to
 40 maximize the log-likelihood of the action sequence given the context of state history and goal:

$$\mathcal{L}_{\text{IL}} = \mathbb{E} \left[\sum_{(\bar{\mathbf{s}}, \bar{\mathbf{a}}, \mathbf{g}) \in \mathcal{T}} \log \pi_\theta(\bar{\mathbf{a}}|\bar{\mathbf{s}}, \mathbf{g}) \right]. \quad (1)$$

41 2.2 Diffusion Policy

42 MoDE uses the continuous-time diffusion model of EDM [10] as a policy representation. Diffusion
 43 models are a type of generative model for generating data by initially adding noise through Gaussian
 44 perturbations and then reversing this process. MoDE applies the score-based diffusion model to
 45 represent the policy $\pi_\theta(\bar{\mathbf{a}}|\bar{\mathbf{s}}, \mathbf{g})$. The perturbation and inverse process can be described using a
 46 stochastic differential equation:

$$d\bar{\mathbf{a}} = (\beta_t \sigma_t - \dot{\sigma}_t) \sigma_t \nabla_{\mathbf{a}} \log p_t(\bar{\mathbf{a}}|\bar{\mathbf{s}}, \mathbf{g}) dt + \sqrt{2\beta_t} \sigma_t d\omega_t, \quad (2)$$

47 where β_t controls the noise injection, $d\omega_t$ refers to infinitesimal Gaussian noise, and $p_t(\bar{\mathbf{a}}|\bar{\mathbf{s}}, \mathbf{g})$ is the
 48 score function of the diffusion process, that moves samples away from regions of high-data density in
 49 the forward process. To generate new samples from noise a neural network is trained to approximate
 50 the score function $\nabla_{\bar{\mathbf{a}}} \log p_t(\bar{\mathbf{a}}|\bar{\mathbf{s}}, \mathbf{g})$ via Score matching (SM) [11]

$$\mathcal{L}_{\text{SM}} = \mathbb{E}_{\sigma, \bar{\mathbf{a}}, \epsilon} [\alpha(\sigma_t) \| D_\theta(\bar{\mathbf{a}} + \epsilon, \bar{\mathbf{s}}, \mathbf{g}, \sigma_t) - \bar{\mathbf{a}} \|_2^2], \quad (3)$$

51 where $D_\theta(\bar{\mathbf{a}} + \epsilon, \bar{\mathbf{s}}, \mathbf{g}, \sigma_t)$ is the trainable neural network. During training, we sample noise from a
 52 training distribution and add it to an action sequence. The network predicts the denoised actions and
 53 computes the SM loss.

54 After training, we can generate a new action sequence starting from random noise by approximating
 55 the reverse SDE or related ODE in discrete steps using numerical ODE integrators. Therefore, we
 56 sample noise from the prior $\mathbf{a}_T \sim \mathcal{N}(\mathbf{0}, \sigma_T^2 \mathbf{I})$ and iteratively denoise it. MoDE uses the DDIM-solver,
 57 a numerical ODE-solver designed for diffusion models [12], that allows fast denoising of actions in a
 58 few steps. MoDE uses 10 denoising steps in all our experiments.

59 **2.3 Mixture-of-Experts Denoising**

60 We now present the details of MoDE’s noise-conditioned expert routing. An overview of MoDE is
 61 shown in Figure 1. For language conditioning, MoDE uses a frozen CLIP language encoder model to
 62 generate a single latent goal vector. To encode images, MoDE uses FiLM-conditioned ResNets-18.
 63 Let $\mathbf{X} \in \mathbb{R}^{\text{tokens} \times \mathbf{D}}$ be a sequence of input tokens of dimension D , and σ_t be the noise level. Let
 64 $\phi(\sigma_t)$ encode the noise level into a token using a sinusoidal embedding followed by a small MLP,
 65 and let $\hat{\mathbf{X}}$ also contain $\phi(\sigma_t)$ as a token. MoE MoDE($\mathbf{X}, \phi(\sigma_t)$) is a composition of L transformer
 66 blocks which we will now describe in detail. Each transformer block f^i takes $\phi(\sigma_t)$ as input as well.
 67 We now define each block f^i as a composition of a self-attention (SA) layer and an MoE layer,

$$f^i(\mathbf{X}, \phi(\sigma_t)) = \text{MoE}(\text{SA}(\hat{\mathbf{X}}) + \mathbf{X}, \phi(\sigma_t)) + \mathbf{X}. \quad (4)$$

68 We add the noise token, $\phi(\sigma_t)$ to all the tokens in \mathbf{X} before computing the self-attention as done in
 69 [13].

$$\hat{\mathbf{X}} = \phi(\sigma_t) + \mathbf{X}, \quad (5)$$

70 and following [14],

$$\text{SA}(\hat{\mathbf{X}}) = \text{softmax}\left(\frac{1}{\sqrt{D}}[\hat{\mathbf{X}}W_Q][\hat{\mathbf{X}}W_K]^T\right)[\hat{\mathbf{X}}W_V]. \quad (6)$$

71 Given N experts $\{E_i\}_{i=1}^N$, we define the sparse MoE layer as

$$\text{MoE}(\mathbf{X}, \phi(\sigma_t)) = \sum_{i=1}^N \mathbf{R}(\phi(\sigma_t)) \mathbf{E}_i(\mathbf{X}), \quad (7)$$

72 where, the routing function $\mathbf{R}(\cdot) : \mathbb{R}^{\text{tokens} \times \mathbf{D}} \rightarrow \mathbb{R}^{\text{tokens} \times \mathbf{N}}$. Our noise-only conditioned routing
 73 function is defined as

$$\mathbf{R}(\phi(\sigma_t)) = \text{topk}(\text{softmax}(\phi(\sigma_t)\mathbf{W}_R), k) \quad (8)$$

74 While topk typically varies across different MoE implementations, we use multinomial sam-
 75 pling. Where we sample without replacement k elements according to their probabilities given
 76 by $\text{softmax}(\phi(\sigma_t)\mathbf{W}_R)$. We set all non-chosen elements to 0. Since sampling is a non-differentiable
 77 process, scaling the expert outputs by the routing probability, $\mathbf{E}_i(\mathbf{X})$ is needed to pass gradients to
 78 the routing function. In addition, we optimally re-normalize the chosen k elements probabilities.
 79 After the final layer, we use a linear projection layer to get the denoised action sequence. The above
 80 formulation is general enough to explore routing variants; we tested what the router was conditioned
 81 on, such as noise only or token only. We report the results in section 3.

82 In addition, to mitigate expert collapse, we adopt an additional loss function (LB) that regularizes
 83 the router called load balancing [15]. Here for a given expert E_i , we compute the fraction of the
 84 tokens that were routed to it, under topk being an arg max function, and we scale that by the average
 85 probability of routing to E_i across all tokens in a batch \mathcal{B} .

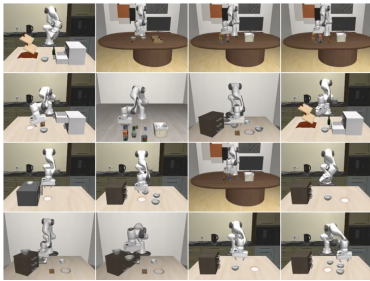
$$LB(\sigma_t) = N \sum_{n=1}^N \frac{1}{|\mathcal{B}|} \left(\sum_{i=1}^{|\mathcal{B}|} \mathbb{1}\{\mathbf{R}(\phi(\sigma_{t_i}))_n > \mathbf{0}\} \right) \frac{1}{|\mathcal{B}|} \left(\sum_{i=1}^{|\mathcal{B}|} \text{softmax}(\phi(\sigma_{t_i})\mathbf{W}_R)_n \right) \quad (9)$$

86 We use $\gamma = 0.01$ for the load-balancing loss.

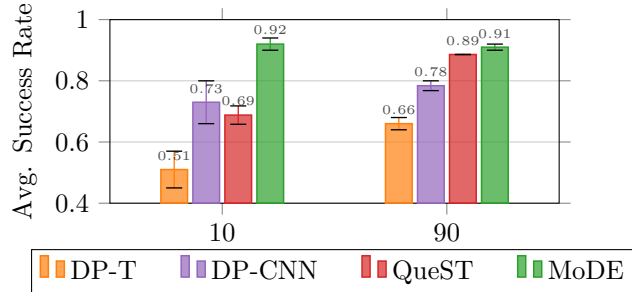
87 **3 Evaluation**

88 Our experiments aim to answer four key questions: (I) How does MoDE compare to other policies
 89 and prior diffusion transformer architectures in terms of performance?

90 We compare MoDE against prior diffusion transformer architecture [5], ensuring fair comparisons
 91 by using a similar number of active parameters. MoDE uses 8 layers with 4 experts and a latent
 92 dimension of 1024 in all experiments. Detailed hyperparameters are provided in the Appendix (
 93 Table 1).



(a) LIBERO-90 Tasks



(b) Results for LIBERO-10 and LIBERO-90

Figure 2: Visualization and Results for LIBERO environment. (a) Few example environments and tasks of the LIBERO-90 task suite. (b) Average results for both LIBERO challenges averaged over 3 seeds with 20 rollouts for each task.

94 3.1 Long-Horizon Multi-Task Experiments

95 We first evaluate MoDE on the LONG-challenge and LIBERO-90 challenge of the LIBERO bench-
 96 mark [9]. The LONG challenge requires a model to learn 10 tasks in different settings. It demands
 97 long-horizon behavior generation with several hundreds of steps for completion. The 90 variant
 98 tests policies on 90 diverse short-horizon tasks in different environments. Figure 2a visualizes a few
 99 examples of these tasks. All environments feature two cameras: a static one and a wrist-mounted
 100 camera, used to encode the current observation using FiLM-ResNets-18. We test each policy 20
 101 times on each task and report the average results over 3 seeds. We use an action chunking length
 102 of 10 and a history length of 1. MoDE and all other diffusion architectures use FiLM-conditioned
 103 ResNets-18 with a CLIP sentence embedding to encode the goal and the images.

104 **Baselines.** We compare MoDE against three state-of-the-art baselines: 1) The Diffusion Transformer
 105 (DP-T) architecture [5], which conditions on noise and observations using a cross-attention module.
 106 2) The standard Diffusion Policy CNN-based architecture (DP-CNN). 3) QueST [16], a transformer-
 107 based policy that learns discrete action representations using vector-quantized embeddings of action
 108 sequences. We tested all baselines ourselves, except for QueST, whose results are taken directly from
 109 their paper.

110 **Results.** The performance of all models on the benchmark is summarized in Figure 2b. Overall,
 111 MoDE achieves the highest average performance in both benchmarks, while the QueST baseline is
 112 the second best in the LIBERO-90 setting and the CNN-architecture is second best in the long horizon
 113 setting. These results demonstrate MoDE’s ability to learn long-horizon tasks with high accuracy.
 114 The performance gap is more pronounced in the challenging LIBERO-10 experiment, where MoDE
 115 is the first policy to achieve an over 90% success rate. Furthermore, MoDE surpasses prior best
 116 Diffusion baselines by an average of 16% in both settings, all while maintaining its computational
 117 advantage. This showcases MoDE’s ability to achieve state-of-the-art performance with a more
 118 efficient use of computational resources.

119 4 Conclusion

120 In this work, we introduced Mixture-of-Denoising Experts (MoDE), a novel Diffusion Policy that
 121 leverages a mixture of experts Transformer to enhance the performance and efficiency of diffusion
 122 policies. We also proposed a noise-conditioned routing strategy for learning specialized experts
 123 within our model. In our extensive experiments and ablation studies across diverse benchmarks, we
 124 demonstrated the advantages of MoDE to outperform prior Diffusion Policies with a lower number
 125 of parameters and 40% less FLOPS during inference. In future work, we want to experiment with
 126 more routing strategies, such as expert-choice routing [17].

References

- 127
- 128 [1] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in Neural*
129 *Information Processing Systems*, 33:6840–6851, 2020.
- 130 [2] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based
131 generative modeling through stochastic differential equations. In *International Conference on*
132 *Learning Representations*, 2020.
- 133 [3] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna,
134 C. Xu, J. Luo, T. Kreiman, Y. Tan, D. Sadigh, C. Finn, and S. Levine. Octo: An open-source
135 generalist robot policy. <https://octo-models.github.io>, 2023.
- 136 [4] M. Reuss, M. Li, X. Jia, and R. Lioutikov. Goal conditioned imitation learning using score-based
137 diffusion policies. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- 138 [5] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion policy:
139 Visuomotor policy learning via action diffusion. In *Proceedings of Robotics: Science and*
140 *Systems (RSS)*, 2023.
- 141 [6] X. Jia, D. Blessing, X. Jiang, M. Reuss, A. Donat, R. Lioutikov, and G. Neumann. To-
142 wards diverse behaviors: A benchmark for imitation learning with human demonstrations.
143 In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=6pPYRXKppw>.
144
- 145 [7] T. Hang, S. Gu, C. Li, J. Bao, D. Chen, H. Hu, X. Geng, and B. Guo. Efficient diffusion training
146 via min-snr weighting strategy, 2024.
- 147 [8] O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard. Calvin: A benchmark for language-
148 conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and*
149 *Automation Letters*, 2022.
- 150 [9] B. Liu, Y. Zhu, C. Gao, Y. Feng, Q. Liu, Y. Zhu, and P. Stone. Libero: Benchmarking knowledge
151 transfer for lifelong robot learning. *arXiv preprint arXiv:2306.03310*, 2023.
- 152 [10] T. Karras, M. Aittala, T. Aila, and S. Laine. Elucidating the design space of diffusion-based
153 generative models. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in*
154 *Neural Information Processing Systems*, 2022.
- 155 [11] P. Vincent. A connection between score matching and denoising autoencoders. *Neural Compu-*
156 *tation*, 23(7):1661–1674, 2011. doi:10.1162/NECO_a_00142.
- 157 [12] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. In *ICLR*, 2021.
- 158 [13] A. Hatamizadeh, J. Song, G. Liu, J. Kautz, and A. Vahdat. Diffit: Diffusion vision transformers
159 for image generation. *arXiv preprint arXiv:2312.02139*, 2023.
- 160 [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and
161 I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*,
162 30, 2017.
- 163 [15] W. Fedus, B. Zoph, and N. Shazeer. Switch transformers: Scaling to trillion parameter models
164 with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- 165 [16] A. Mete, H. Xue, A. Wilcox, Y. Chen, and A. Garg. Quest: Self-supervised skill abstractions
166 for learning continuous control. *arXiv preprint arXiv:2407.15840*, 2024.
- 167 [17] Y. Zhou, T. Lei, H. Liu, N. Du, Y. Huang, V. Zhao, A. Dai, Z. Chen, Q. Le, and J. Laudon.
168 Mixture-of-experts with expert choice routing, 2022.

- 169 [18] X. Li, M. Liu, H. Zhang, C. Yu, J. Xu, H. Wu, C. Cheang, Y. Jing, W. Zhang, H. Liu, et al.
170 Vision-language foundation models as effective robot imitators. In *International Conference on*
171 *Learning Representations*, 2024.
- 172 [19] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican,
173 M. Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in*
174 *Neural Information Processing Systems*, 35:23716–23736, 2022.
- 175 [20] T. Brooks, A. Holynski, and A. A. Efros. Instructpix2pix: Learning to follow image editing
176 instructions. In *CVPR*, 2023.
- 177 [21] K. Black, M. Nakamoto, P. Atreya, H. Walke, C. Finn, A. Kumar, and S. Levine. Zero-
178 shot robotic manipulation with pretrained image-editing diffusion models. *arXiv preprint*
179 *arXiv:2310.10639*, 2023.
- 180 [22] H. Wu, Y. Jing, C. Cheang, G. Chen, J. Xu, X. Li, M. Liu, H. Li, and T. Kong. Unleashing large-
181 scale video generative pre-training for visual robot manipulation. In *International Conference*
182 *on Learning Representations*, 2024.
- 183 [23] O. Mees, L. Hermann, and W. Burgard. What matters in language conditioned robotic imitation
184 learning over unstructured data. *IEEE Robotics and Automation Letters (RA-L)*, 7(4):11205–
185 11212, 2022.
- 186 [24] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. Film: Visual reasoning with a
187 general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*,
188 volume 32, 2018.
- 189 [25] W. Peebles and S. Xie. Scalable diffusion models with transformers. In *Proceedings of the*
190 *IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- 191 [26] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet. Learning
192 latent plans from play, 2019.
- 193 [27] N. M. M. Shafiullah, Z. J. Cui, A. Altanzaya, and L. Pinto. Behavior transformers: Cloning k
194 modes with one stone. In *Thirty-Sixth Conference on Neural Information Processing Systems*,
195 2022. URL <https://openreview.net/forum?id=agTr-vRQsa>.
- 196 [28] Z. J. Cui, Y. Wang, N. M. M. Shafiullah, and L. Pinto. From play to policy: Conditional behavior
197 generation from uncurated robot data. In *International Conference on Learning Representations*,
198 2023. URL <https://openreview.net/forum?id=c7rM7F7jQjN>.
- 199 [29] S. Lee, Y. Wang, H. Etukuru, H. J. Kim, N. M. M. Shafiullah, and L. Pinto. Behavior generation
200 with latent actions. *arXiv preprint arXiv:2403.03181*, 2024.
- 201 [30] Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution.
202 *Advances in Neural Information Processing Systems*, 32, 2019.
- 203 [31] Z. Xian, N. Gkanatsios, T. Gervet, T.-W. Ke, and K. Fragkiadaki. Chaineddiffuser: Unifying
204 trajectory diffusion and keypose prediction for robotic manipulation. In *7th Annual Conference*
205 *on Robot Learning*, 2023. URL <https://openreview.net/forum?id=W0zgY2mBTAB>.
- 206 [32] T.-W. Ke, N. Gkanatsios, and K. Fragkiadaki. 3d diffuser actor: Policy diffusion with 3d scene
207 representations. *arXiv preprint arXiv:2402.10885*, 2024.
- 208 [33] X. Li, V. Belagali, J. Shang, and M. S. Ryou. Crossway diffusion: Improving diffusion-based
209 visuomotor policy via self-supervised learning. *arXiv preprint arXiv:2307.01849*, 2023.
- 210 [34] P. M. Scheickl, N. Schreiber, C. Haas, N. Freymuth, G. Neumann, R. Lioutikov, and F. Mathis-
211 Ullrich. Movement primitive diffusion: Learning gentle robotic manipulation of deformable
212 objects. *arXiv preprint arXiv:2312.10008*, 2023.

- 213 [35] M. Reuss, Ö. E. Yağmurlu, F. Wenzel, and R. Lioutikov. Multimodal diffusion transformer:
214 Learning versatile behavior from multimodal goals. In *Robotics: Science and Systems*, 2024.
- 215 [36] A. Ajay, Y. Du, A. Gupta, J. B. Tenenbaum, T. S. Jaakkola, and P. Agrawal. Is conditional
216 generative modeling all you need for decision making? In *International Conference on Learning*
217 *Representations*, 2023. URL <https://openreview.net/forum?id=sP1fo2K9DFG>.
- 218 [37] M. Janner, Y. Du, J. Tenenbaum, and S. Levine. Planning with diffusion for flexible behavior
219 synthesis. In *International Conference on Machine Learning*, pages 9902–9915. PMLR, 2022.
- 220 [38] J. Pari, N. Shafiqullah, S. Arunachalam, and L. Pinto. The Surprising Effectiveness of Representa-
221 tion Learning for Visual Imitation. In *Proceedings of Robotics: Science and Systems*, New
222 York City, NY, USA, June 2022. doi:10.15607/RSS.2022.XVIII.010.
- 223 [39] T.-H. Wang, J. Zheng, P. Ma, Y. Du, B. Kim, A. E. Spielberg, J. B. Tenenbaum, C. Gan,
224 and D. Rus. Diffusebot: Breeding soft robots with physics-augmented generative diffusion
225 models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL
226 <https://openreview.net/forum?id=1zo4iioUEs>.
- 227 [40] Y. Du, M. Yang, B. Dai, H. Dai, O. Nachum, J. B. Tenenbaum, D. Schuurmans, and P. Abbeel.
228 Learning universal policies via text-guided video generation. *arXiv preprint arXiv:2302.00111*,
229 2023.
- 230 [41] P.-C. Ko, J. Mao, Y. Du, S.-H. Sun, and J. B. Tenenbaum. Learning to Act from Actionless
231 Video through Dense Correspondences. *arXiv:2310.08576*, 2023.
- 232 [42] A. Ajay, S. Han, Y. Du, S. Li, A. Gupta, T. Jaakkola, J. Tenenbaum, L. Kaelbling, A. Srivastava,
233 and P. Agrawal. Compositional foundation models for hierarchical planning. *arXiv preprint*
234 *arXiv:2309.08587*, 2023.
- 235 [43] J. Carvalho, A. T. Le, M. Baierl, D. Koert, and J. Peters. Motion planning diffusion: Learning and
236 planning of robot motions with diffusion models. In *2023 IEEE/RSJ International Conference*
237 *on Intelligent Robots and Systems (IROS)*, pages 1916–1923. IEEE, 2023.
- 238 [44] J. Urain, N. Funk, J. Peters, and G. Chalvatzaki. Se (3)-diffusionfields: Learning smooth cost
239 functions for joint grasp and motion optimization through diffusion. In *2023 IEEE International*
240 *Conference on Robotics and Automation (ICRA)*, pages 5923–5930. IEEE, 2023.
- 241 [45] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously
242 large neural networks: The sparsely-gated mixture-of-experts layer, 2017.
- 243 [46] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. d. l.
244 Casas, E. B. Hanna, F. Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*,
245 2024.
- 246 [47] N. Du, Y. Huang, A. M. Dai, S. Tong, D. Lepikhin, Y. Xu, M. Krikun, Y. Zhou, A. W. Yu,
247 O. Firat, B. Zoph, L. Fedus, M. Bosma, Z. Zhou, T. Wang, Y. E. Wang, K. Webster, M. Pellat,
248 K. Robinson, K. Meier-Hellstern, T. Duke, L. Dixon, K. Zhang, Q. V. Le, Y. Wu, Z. Chen, and
249 C. Cui. Glam: Efficient scaling of language models with mixture-of-experts, 2022.
- 250 [48] Z. Chi, L. Dong, S. Huang, D. Dai, S. Ma, B. Patra, S. Singhal, P. Bajaj, X. Song, X.-L. Mao,
251 H. Huang, and F. Wei. On the representation collapse of sparse mixture of experts. In A. H.
252 Oh, A. Agarwal, D. Belgrave, and K. Cho, editors, *Advances in Neural Information Processing*
253 *Systems*, 2022. URL <https://openreview.net/forum?id=mWaYC6CZf5>.
- 254 [49] H. Hazimeh, Z. Zhao, A. Chowdhery, M. Sathiamoorthy, Y. Chen, R. Mazumder, L. Hong, and
255 E. H. Chi. Dselect-k: Differentiable selection in the mixture of experts with applications to multi-
256 task learning. *CoRR*, abs/2106.03760, 2021. URL <https://arxiv.org/abs/2106.03760>.

- 257 [50] S. Roller, S. Sukhbaatar, A. Szlam, and J. Weston. Hash layers for large sparse models. *CoRR*,
258 abs/2106.04426, 2021. URL <https://arxiv.org/abs/2106.04426>.
- 259 [51] M. Lewis, S. Bhosale, T. Dettmers, N. Goyal, and L. Zettlemoyer. Base layers: Simplifying
260 training of large, sparse models, 2021.
- 261 [52] J. Obando-Ceron, G. Sokar, T. Willi, C. Lyle, J. Farebrother, J. Foerster, G. K. Dziugaite,
262 D. Precup, and P. S. Castro. Mixtures of experts unlock parameter scaling for deep rl. *arXiv*
263 *preprint arXiv:2402.08609*, 2024.
- 264 [53] O. Celik, D. Zhou, G. Li, P. Becker, and G. Neumann. Specializing versatile skill libraries
265 using local mixture of experts. In A. Faust, D. Hsu, and G. Neumann, editors, *Proceedings*
266 *of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine Learning*
267 *Research*, pages 1423–1433. PMLR, 08–11 Nov 2022. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v164/celik22a.html)
268 [press/v164/celik22a.html](https://proceedings.mlr.press/v164/celik22a.html).
- 269 [54] O. Celik, A. Taranovic, and G. Neumann. Acquiring diverse skills using curriculum reinforce-
270 ment learning with mixture of experts. *arXiv preprint arXiv:2403.06966*, 2024.
- 271 [55] K. Hansel, J. Urain, J. Peters, and G. Chalvatzaki. Hierarchical policy blending as inference for
272 reactive robot control. In *2023 IEEE International Conference on Robotics and Automation*
273 *(ICRA)*, pages 10181–10188. IEEE, 2023.
- 274 [56] A. T. Le, K. Hansel, J. Peters, and G. Chalvatzaki. Hierarchical policy blending as optimal
275 transport. In *Learning for Dynamics and Control Conference*, pages 797–812. PMLR, 2023.
- 276 [57] C. Riquelme, J. Puigcerver, B. Mustafa, M. Neumann, R. Jenatton, A. S. Pinto, D. Keysers, and
277 N. Houlsby. Scaling vision with sparse mixture of experts, 2021.
- 278 [58] O. Mees, A. Eitel, and W. Burgard. Choosing smartly: Adaptive multimodal fusion for object
279 detection in changing environments. In *2016 IEEE/RSJ International Conference on Intelligent*
280 *Robots and Systems (IROS)*, pages 151–156. IEEE, 2016.
- 281 [59] D. Blessing, O. Celik, X. Jia, M. Reuss, M. X. Li, R. Lioutikov, and G. Neumann. Informa-
282 tion maximizing curriculum: A curriculum-based approach for learning versatile skills.
283 In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL
284 <https://openreview.net/forum?id=7eW6NzSE4g>.
- 285 [60] M. X. Li, O. Celik, P. Becker, D. Blessing, R. Lioutikov, and G. Neumann. Curriculum-based
286 imitation of versatile skills. In *2023 IEEE International Conference on Robotics and Automation*
287 *(ICRA)*, pages 2951–2957, 2023. doi:10.1109/ICRA48891.2023.10160543.
- 288 [61] B. Park, S. Woo, H. Go, J.-Y. Kim, and C. Kim. Denoising task routing for diffusion models.
289 *arXiv preprint arXiv:2310.07138*, 2023.
- 290 [62] H. Go, Y. Lee, J.-Y. Kim, S. Lee, M. Jeong, H. S. Lee, and S. Choi. Towards practical plug-and-
291 play diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and*
292 *pattern recognition*, pages 1962–1971, 2023.
- 293 [63] B. Park, H. Go, J.-Y. Kim, S. Woo, S. Ham, and C. Kim. Switch diffusion transformer:
294 Synergizing denoising tasks with sparse mixture-of-experts. *arXiv preprint arXiv:2403.09176*,
295 2024.
- 296 [64] Y. Lee, J. Kim, H. Go, M. Jeong, S. Oh, and S. Choi. Multi-architecture multi-expert diffusion
297 models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages
298 13427–13436, 2024.
- 299 [65] H. Bharadhwaj, J. Vakil, M. Sharma, A. Gupta, S. Tulsiani, and V. Kumar. Roboagent:
300 Generalization and efficiency in robot manipulation via semantic augmentations and action
301 chunking, 2023.

- 302 [66] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation
303 with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.
- 304 [67] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal,
305 E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv*
306 *preprint arXiv:2302.13971*, 2023.
- 307 [68] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Haus-
308 man, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. Joshi, R. Julian,
309 D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manjunath,
310 I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao,
311 M. Ryoo, G. Salazar, P. Sanketi, K. Sayed, J. Singh, S. Sontakke, A. Stone, C. Tan, H. Tran,
312 V. Vanhoucke, S. Vega, Q. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-1:
313 Robotics transformer for real-world control at scale. In *arXiv preprint arXiv:2212.06817*, 2022.
- 314 [69] B. Zitkovich, T. Yu, S. Xu, P. Xu, T. Xiao, F. Xia, J. Wu, P. Wohlhart, S. Welker, A. Wahid,
315 Q. Vuong, V. Vanhoucke, H. Tran, R. Soricut, A. Singh, J. Singh, P. Sermanet, P. R. Sanketi,
316 G. Salazar, M. S. Ryoo, K. Reymann, K. Rao, K. Pertsch, I. Mordatch, H. Michalewski, Y. Lu,
317 S. Levine, L. Lee, T.-W. E. Lee, I. Leal, Y. Kuang, D. Kalashnikov, R. Julian, N. J. Joshi,
318 A. Irpan, brian ichter, J. Hsu, A. Herzog, K. Hausman, K. Gopalakrishnan, C. Fu, P. Florence,
319 C. Finn, K. A. Dubey, D. Driess, T. Ding, K. M. Choromanski, X. Chen, Y. Chebotar, J. Carbajal,
320 N. Brown, A. Brohan, M. G. Arenas, and K. Han. RT-2: Vision-language-action models transfer
321 web knowledge to robotic control. In *7th Annual Conference on Robot Learning*, 2023. URL
322 <https://openreview.net/forum?id=XMQgwiJ7KSX>.
- 323 [70] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess,
324 A. Dubey, C. Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to
325 robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- 326 [71] J. Gu, S. Kirmani, P. Wohlhart, Y. Lu, M. G. Arenas, K. Rao, W. Yu, C. Fu, K. Gopalakrishnan,
327 Z. Xu, P. Sundaresan, P. Xu, H. Su, K. Hausman, C. Finn, Q. Vuong, and T. Xiao. Rt-trajectory:
328 Robotic task generalization via hindsight trajectory sketches. In *International Conference on*
329 *Learning Representations*, 2024.
- 330 [72] O. X.-E. Collaboration, A. Padalkar, A. Pooley, A. Jain, A. Bewley, A. Herzog, A. Irpan,
331 A. Khazatsky, A. Rai, A. Singh, A. Brohan, A. Raffin, A. Wahid, B. Burgess-Limerick, B. Kim,
332 B. Schölkopf, B. Ichter, C. Lu, C. Xu, C. Finn, C. Xu, C. Chi, C. Huang, C. Chan, C. Pan, C. Fu,
333 C. Devin, D. Driess, D. Pathak, D. Shah, D. Büchler, D. Kalashnikov, D. Sadigh, E. Johns,
334 F. Ceola, F. Xia, F. Stulp, G. Zhou, G. S. Sukhatme, G. Salhotra, G. Yan, G. Schiavi, H. Su,
335 H.-S. Fang, H. Shi, H. B. Amor, H. I. Christensen, H. Furuta, H. Walke, H. Fang, I. Mordatch,
336 I. Radosavovic, I. Leal, J. Liang, J. Kim, J. Schneider, J. Hsu, J. Bohg, J. Bingham, J. Wu,
337 J. Wu, J. Luo, J. Gu, J. Tan, J. Oh, J. Malik, J. Tompson, J. Yang, J. J. Lim, J. Silvério, J. Han,
338 K. Rao, K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund,
339 K. Kawaharazuka, K. Zhang, K. Majd, K. Rana, K. Srinivasan, L. Y. Chen, L. Pinto, L. Tan,
340 L. Ott, L. Lee, M. Tomizuka, M. Du, M. Ahn, M. Zhang, M. Ding, M. K. Srirama, M. Sharma,
341 M. J. Kim, N. Kanazawa, N. Hansen, N. Heess, N. J. Joshi, N. Suenderhauf, N. D. Palo,
342 N. M. M. Shafiullah, O. Mees, O. Kroemer, P. R. Sanketi, P. Wohlhart, P. Xu, P. Sermanet,
343 P. Sundaresan, Q. Vuong, R. Rafailov, R. Tian, R. Doshi, R. Martín-Martín, R. Mendonca,
344 R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Moore, S. Bahl, S. Dass,
345 S. Song, S. Xu, S. Haldar, S. Adebola, S. Guist, S. Nasiriany, S. Schaal, S. Welker, S. Tian,
346 S. Dasari, S. Belkhale, T. Osa, T. Harada, T. Matsushima, T. Xiao, T. Yu, T. Ding, T. Davchev,
347 T. Z. Zhao, T. Armstrong, T. Darrell, V. Jain, V. Vanhoucke, W. Zhan, W. Zhou, W. Burgard,
348 X. Chen, X. Wang, X. Zhu, X. Li, Y. Lu, Y. Chebotar, Y. Zhou, Y. Zhu, Y. Xu, Y. Wang,
349 Y. Bisk, Y. Cho, Y. Lee, Y. Cui, Y. hua Wu, Y. Tang, Y. Zhu, Y. Li, Y. Iwasawa, Y. Matsuo,
350 Z. Xu, and Z. J. Cui. Open X-Embodiment: Robotic learning datasets and RT-X models.
351 <https://arxiv.org/abs/2310.08864>, 2023.

Hyperparameter	CALVIN ABCD	CALVIN ABC	LIBERO-10	LIBERO-90
Number of Transformer Layers	8	8	8	8
Number Experts	4	4	4	4
Attention Heads	8	8	8	8
Action Chunk Size	10	10	10	10
History Length	1	1	1	1
Embedding Dimension	1024	1024	1024	1024
Image Encoder	FiLM-ResNet18	FiLM-ResNet18	FiLM-ResNet18	FiLM-ResNet18
Goal Lang Encoder	CLIP ViT-B/32	CLIP ViT-B/32	CLIP ViT-B/32	CLIP ViT-B/32
Attention Dropout	0.3	0.3	0.3	0.3
Residual Dropout	0.1	0.1	0.1	0.1
MLP Dropout	0.1	0.1	0.1	0.1
Optimizer	AdamW	AdamW	AdamW	AdamW
Betas	[0.9, 0.9]	[0.9, 0.9]	[0.9, 0.9]	[0.9, 0.9]
Learning Rate	1e-4	1e-4	1e-4	1e-4
Transformer Weight Decay	0.05	0.05	0.05	0.05
Other weight decay	0.05	0.05	0.05	0.05
Batch Size	512	512	512	512
Train Steps in Thousands	30	25	15	40
σ_{\max}	80	80	80	80
σ_{\min}	0.001	0.001	0.001	0.001
σ_t	0.5	0.5	0.5	0.5
EMA	True	True	True	True
Time steps	Exponential	Exponential	Exponential	Exponential
Sampler	DDIM	DDIM	DDIM	DDIM
Parameter Count (Millions)	460	460	460	460

Table 1: Summary of all the Hyperparameters for the **MoDE** policy used in our experiments.353 **A.1 Experiments Details**354 **A.1.1 CALVIN Benchmark**

355 The CALVIN benchmark [8] is an established **IL** benchmark for learning language-conditioned
356 behavior from human play data. In contrast to other benchmarks the data does not contain structured
357 demonstrations, where the robot completes one task, but instead, the dataset was collected by humans,
358 that randomly interact with the environment. From these long-horizon trajectories across the 4
359 different settings, the authors randomly cut out short sequences with 64 frames and labeled them
360 with the task label. While the dataset offers models to train on the unlabeled part too, we restricted
361 **MoDE** to only train on the labeled parts. The Franke Emika Panda robot is controlled using a
362 Delta-End-Effector Space with a discrete gripper. We use two cameras to encode the current scene: a
363 static camera and a wrist one and predict the next 10 actions, before receiving the next observations
364 and generating another set of 10 actions.

365 **CALVIN ABC.** We train **MoDE** and our dense transformer baseline for 25k training steps with a
366 batch size of 512 on a 4 GPU Cluster Node with 4 A6000 NVIDIA GPUs for 6.5 hours with all
367 1000 rollouts at the end of training. We report the mean results averaged over 3 seeds as done in
368 all relevant prior work. All baselines are reported from the original paper given the standardized
369 evaluation protocol of CALVIN [8].

370 **CALVIN ABCD.** We train **MoDE** and our dense transformer baseline for 30k training steps with
371 a batch size of 512 on a 4 GPU Cluster Node with 4 A6000 NVIDIA GPUs for 7.5 hours with all
372 1000 rollouts at the end of training. We report the mean results averaged over 3 seeds as done in all
373 relevant prior work.

374 **A.1.2 LIBERO Benchmark**

375 **LIBERO-10.** The LIBERO-10 benchmark consists of 50 demonstrations for 10 different tasks
376 that are all labeled with a text instruction. The Franka Emika Panda robot is controlled using an
377 end-effector controller. Similar to CALVIN all models have access to two camera inputs: a static one
378 and a wrist camera. We train **MoDE** and our dense transformer baseline for 50 epochs with a batch
379 size of 512 on a 4 GPU Cluster Node with 4 A6000 NVIDIA GPUs for 2 hours with all 200 rollouts

Model	Block Push	Relay Kitchen	CAL ABC	CAL ABCD	L-10	Average
Dense T	0.96±0.02	3.73±0.12	2.83±0.19	4.13±0.11	0.91±0.02	0.839±0.144
Token-Router	0.97±0.01	3.85±0.03	2.67±0.04	4.29±0.08	0.90±0.01	0.845±0.161
σ_t -Router	<i>0.97±0.01</i>	3.79±0.04	<i>2.79±0.16</i>	4.30±0.02	0.92±0.02	0.851±0.151

Table 2: Overview of the performance of all different token routing strategies used for MoDE across 5 benchmarks. We mark the best result for each environment in **bold** and the second best in *curly*. We use CAL to represent CALVIN. To average the results, we normalize all scores and compute the average over all environments.

380 at the end of training. The benchmark does require to test models on 10 different long-horizon tasks.
381 We test each task 20 times for each model and report the final average performance overall 10 tasks.

382 **LIBERO-90.** The LIBERO-10 benchmark consists of 50 demonstrations for 90 different tasks
383 that are all labeled with a text instruction. The Franka Emika Panda robot is controlled using an
384 end-effector controller. We train MoDE and our dense transformer baseline for 50k steps with a
385 batch size of 512 on a 4 GPU Cluster Node with 4 A6000 NVIDIA GPUs for 12 hours with all 1800
386 rollouts at the end of training. The benchmark does require to test models on 90 different tasks in
387 many different environments. We test each task 20 times for each model and report the final average
388 performance overall 90 tasks.

389 A.2 Baselines

390 Below we explain several baselines used in the experiments in detail:

391 **Diffusion Policy-CNN/T** Inspired by [5], we evaluate extension of the DDPM based Diffusion Policy
392 framework for goal-conditioned Multi-task learning. We evaluate two versions: the CNN-based
393 variant and the Diffusion-Transformer variant, that is conditioned on context and noise using cross-
394 attention. For our experiments we also use EDM-based Diffusion framework for fair comparison
395 against MoDE. We optimized the ideal number of layers and latent dimension for the Transformer
396 baseline and our final version uses 8 layers with a latent dimension of 1024. Larger or smaller variants
397 resulted in lower average performance.

398 **RoboFlamingo.** RoboFlamingo [18] is a **Vision-Language-Models (VLM)** finetuned for behavior
399 generation. The authors use a 3 billion parameter Flamingo model [19] and fine-tune it on CALVIN
400 by freezing the forward blocks and only fine-tuning a new Perceiver Resampler module to extract
401 features from a frozen vision-transformer image encoder and the cross-attention layers to process
402 the image features. Finally, a new action head is learned to generate actions. Overall, the finetuning
403 requires training approx. 1 billion of the parameters. We report the reported results from the paper
404 since they use the standard CALVIN evaluation suite.

405 **SuSIE.** This model first finetunes Instruct2Pix, an image-generation diffusion model, that generates
406 images conditioned on another image and a text description [20] on the local CALVIN robotics
407 domain and uses it as a high-level goal generator. The low-level policy is a **Convolutional neural
408 network (CNN)**-based Diffusion Policy, that predicts the next 4 actions given the current state
409 embedding and desired sub-goal from the high-level policy [21].

410 **GR-1** A causal GPT-Transformer model [22], that has been pretrained on large-scale generative
411 video prediction of human videos. Later, the model is finetuned using co-training of action prediction
412 and video prediction on CALVIN. We report the results directly from their paper for the CALVIN
413 benchmark.

414 A.3 Scaling Multi-Task Experiments

415 Next, we evaluate MoDE efficacy on the demanding **CALVIN Language-Skills Benchmark** [8],
416 an established image-based benchmark for **IL**. This benchmark contains a large dataset of human-
417 recorded demonstrations. First, MoDE is tested on the **ABCD→D** challenge, which involves 22, 966

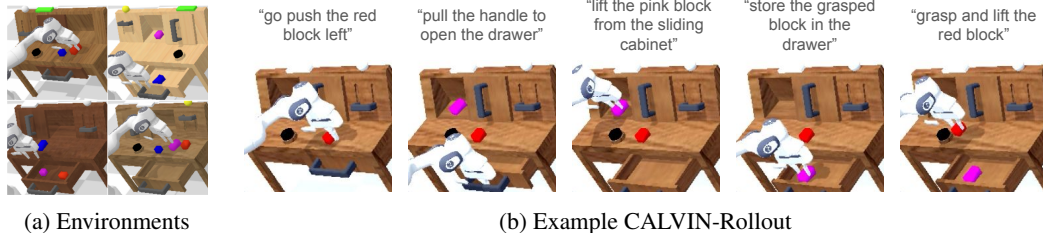


Figure 3: Overview of the CALVIN environment. (a) CALVIN contains four different settings (A,B,C,D) with different configurations of slides, drawers and textures. (b) Example rollout consisting of 5 tasks in sequence. The next goal is only given to the policy, if it manages to complete the prior.

Train→Test	Method	Active Params in Million	PrT	No. Instructions in a Row (1000 chains)					Avg. Len.
				1	2	3	4	5	
ABCD→D	Diff-P-CNN	321	×	86.3%	72.7%	60.1%	51.2%	41.7%	3.16±0.06
	Diff-P-T	194	×	78.3%	53.9%	33.8%	20.4%	11.3%	1.98±0.09
	RoboFlamingo	1000	✓	96.4%	89.6%	82.4%	74.0%	66.0%	4.09±0.00
	GR-1	130	✓	94.9%	89.6%	84.4%	78.9%	73.1%	4.21±0.00
	MoDE	277	×	96.6%	90.6%	86.6%	80.9%	75.5%	4.30±0.02
ABC→D	Diff-P-CNN	321	×	63.5%	35.3%	19.4%	10.7%	6.4%	1.35±0.05
	Diff-P-T	194	×	62.2%	30.9%	13.2%	5.0%	1.6%	1.13±0.02
	RoboFlamingo	1000	✓	82.4%	61.9%	46.6%	33.1%	23.5%	2.47±0.00
	SuSIE	860+	✓	87.0%	69.0%	49.0%	38.0%	26.0%	2.69±0.00
	GR-1	130	✓	85.4%	71.2%	59.6%	49.7%	40.1%	3.06±0.00
	MoDE	277	×	87.7%	69.8%	52.11%	40.2%	29.1%	2.79±0.18

Table 3: Performance comparison on the two CALVIN challenges. The table reports average success rates for individual tasks within instruction chains and the average rollout length (Avg. Len.) to complete 5 consecutive instructions, based on 1000 chains. Zero standard deviation indicates methods without reported average performance. 'PrT' denotes models requiring policy pretraining. Parameter counts exclude language encoders.

418 interaction sequences across four environments (A, B, C, D), with each consisting of 64 timesteps
 419 and 34 diverse tasks. These tasks require the acquisition of complex, sequential behaviors and the
 420 ability to chain together different skills. Figure 3a depicts the diverse configurations of interactive
 421 elements within these environments. This particular challenge examines the scaling abilities of
 422 policies trained on a rich variety of data and skills across multiple settings. All policies are tested on
 423 1000 instructions chains consisting of 5 tasks in sequence in environment D following the official
 424 protocol of CALVIN [8]. One example rollout with 5 different tasks is visualized in Figure 3b. In
 425 terms of scoring, the model receives 1 point for completing a task and only progresses to the next
 426 task upon completion of the prior one. We report the average sequence length over 3 seeds with 1000
 427 instruction chains each.

428 **Baselines.** We test MoDE against several methods specialized for learning language-conditioned
 429 behavior and against other baseline diffusion policy architectures. We also compare MoDE against
 430 RoboFlamingo and GR-1. RoboFlamingo is a fine-tuned Vision-Language-Action model, that
 431 contains around 3 billion parameters and has been pre-trained on diverse internet data. GR-1 is a
 432 generative decoder-only Transformer pretrained on large-scale video generation and then co-finetuned
 433 on CALVIN [22]. If available, we report the average performance of all prior work directly from their
 434 paper, given the standard evaluation protocol in CALVIN [23].

435 **Results.** Our findings, outlined in Table 3 reveal that MoDE outperforms all other policies in terms
 436 of average success rate. Moreover, MoDE outperforms well-established baselines like RoboFlamingo
 437 and GR-1, which depend on extensive internet-scale pretraining for their results. Notably, while GR-1
 438 uses fewer active parameters (130M compared to MoDE’s 277M), it operates with a history length
 439 of 10 and 14 tokens for each timestep. Despite this, MoDE proves more computationally efficient,
 440 requiring fewer FLOPs during inference (7.03 vs 7.93 GFLOPs for GR-1). The combination of sota

441 performance, lower computational demands, and no need for resource-intensive pretraining positions
 442 **MoDE** as a highly practical solution for multitask settings.

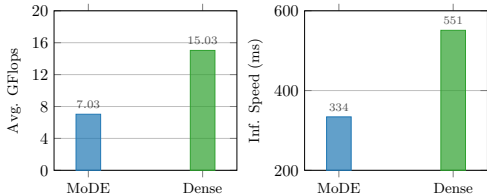
443 **A.4 Zero-shot Generalization Experiments**

444 Finally, we then extend our investigation to the **ABC**→**D** challenge in the second phase, testing
 445 **MoDE**’s zero-shot generalization abilities. In this experiment, models are only trained on data from
 446 the first three CALVIN environments A,B,C and tested on the unseen setting of environment D, which
 447 has different positions of relevant objects and texture of the table. This requires policies, that are able
 448 to generalize their learned behavior to new environment configurations and different textures, which
 449 is especially challenging.

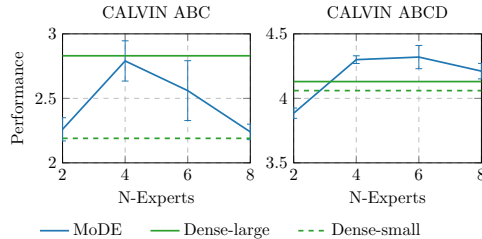
450 **Baselines.** For this experiment, we compare **MoDE** against the previous CALVIN baselines, with
 451 the addition of SuSIE [21]. A hierarchical policy utilizing a finetuned image-generation model,
 452 Instruct2Pix [20], to generate goal images, which guide a low-level diffusion policy. The high-level
 453 goal generation model requires large-scale pretraining. SuSIE’s results are based on 100 rollouts only,
 454 without standard deviation, due to the computational cost of generating subgoal images.

455 **Results.** The results of this experiment are summarized in the lower part of Table 3. **MoDE**
 456 outperforms all tested baselines except for GR-1 and surpasses all other Diffusion Policy architectures
 457 by a wide margin. While GR-1 is slightly better, it requires expensive large-scale pre-training with
 458 32 GPUS and 7 days of total training time to achieve these results. In contrast, **MoDE** trains on
 459 4 GPUs in 8 hours and achieves similar performance without additional pretraining requirements.
 460 Therefore, in response to Question (I), we affirmatively conclude that Mixture-of-Experts models not
 461 only enhance scaling performance but also demonstrate strong zero-shot generalization capabilities.

462 **A.5 Computational Efficiency of MoDE**



(a) Computational efficiency comparison between **MoDE** and Dense-Transformer model with the same number of parameters. Left: FLOP count for both model variants. Right: Average inference speed over 10 forward passes. **MoDE** demonstrates superior efficiency with lower FLOP count and faster inference.



(b) Scaling performance of **MoDE** and Dense-**MoDE** on CALVIN ABC and ABCD environments, showing average performance for 2 to 8 experts using best-performing variants for each environment.

463 We compute the average time for **MoDE** and a dense transformer baseline with a similar number of
 464 parameters for 10 forward passes to assess the computational efficiency of both. The results of this
 465 ablation are summarized in Figure 4a. **MoDE** is around 40% faster than the dense transformer model
 466 with lower FLOPS and fewer parameters than the dense model. Given the prior experimental results,
 467 **MoDE** does increase over dense models in terms of performance, efficiency and inference speed.

468 **A.6 Ablation Studies**

469 To thoroughly evaluate **MoDE**’s performance and design choices, we conducted a series of ablation
 470 studies. These experiments address our research questions: the computational efficiency of **MoDE**
 471 (Question II), the impact of different routing strategies (Question III), and the distribution of tokens
 472 to experts (Question IV).

473 **A.6.1 What design decisions affect MoDE’s performance?**

474 First, we assess the impact of various design decisions
475 on MoDE’s performance. We ablate the choice of noise-
476 conditioning and various MoE strategies on the LIBERO-10
477 benchmark. The results are summarized in Table 4.

478 **Noise-Injection Ablations.** Our experiments reveal the im-
479 portance of proper noise conditioning in MoDE. The full
480 MoDE model, which uses both input noise tokens and noise-
481 conditioned self-attention, achieves the best performance with
482 an average success rate of 0.92. Removing the input noise token
483 slightly decreases performance to 0.90, highlighting the comple-
484 mentary nature of both conditioning methods. Using only
485 the noise token for conditioning, without noise-conditioned self-
486 attention, further reduces performance to 0.85. Interestingly,
487 using FiLM noise conditioning [24], a common approach in
488 image-diffusion models [25], yields the lowest performance in
489 this group at 0.81. These results underscore the effectiveness of our proposed noise conditioning
490 strategy in MoDE, demonstrating a clear performance advantage of 0.08 over the FiLM approach.

491 **MoE Ablations.** Next, we ablate several design decisions regarding Mixture-of-Experts. First, we test
492 the topk number of used experts. Setting topk to 1 only marginally lowers the average performance
493 from 0.92 to 0.91. MoDE maintains robust performance even with a single expert. We also examine
494 the effect of using a shared expert, where the model consistently employs the same expert in all
495 cases. This approach achieves a comparable average success rate of 0.90. Different choices for the
496 token-distribution loss are also tested. While MoDE uses $\gamma = 0.01$ as a default value, we experiment
497 with γ values of 0.1 and 0.001, which result in average success rates of 0.90 and 0.86, respectively.
498 These results indicate that a γ value of 0.01 strikes the best performance.

499 **Latent Dimension.** We investigate the impact of varying the latent dimension in MoDE, testing
500 dimensions of 256, 512, and 1024 (our default). The results show that increasing the latent dimension
501 from 256 to 512 yields a modest performance improvement from 0.86 to 0.87, while further increasing
502 to 1024 provides a more substantial boost to 0.92. This suggests that a larger latent dimension allows
503 MoDE to capture more complex representations, leading to improved performance.

	Avg. Success.
MoDE	0.92
- Input Noise Token	0.90
- Noise-cond Attention	0.85
FiLM Noise Conditioning	0.81
TopK=1	0.91
Shared Expert	0.90
$\gamma = 0.1$	0.90
$\gamma = 0.001$	0.86
256 Embed Dim	0.86
512 Embed Dim	0.87

Table 4: Ablation Studies for MoDE on LIBERO-10. All results are averaged over 3 seeds with 20 rollouts each.

504 **A.7 Detailed Experimental Results**

505 We summarize the ablations regarding the choice of routing in Table 2. Therefore, we test two 2
506 different routing strategies across 5 benchmarks.

507 **A.8 Additional Ablation Studies**

508 **A.8.1 Optimal Routing Strategy for Diffusion Transformers**

509 Next, we answer Question (II) by testing different routing strategies for our diffusion-transformer
510 policy across several environments. We test two different token routing strategies: 1) Token-only
511 conditioned Routing and 2) Noise-only Token Routing. (1) is commonly used in LLMs, where the
512 routing is decided based on the tokens only. We test these strategies in five experiments and report the
513 average performance over 3 seeds: Noise-only Routing achieves an average normalized performance
514 of 0.851, slightly outperforming Token-only Routing, which achieves 0.845. Detailed results are
515 summarized in Table 2 in the Appendix. The results demonstrate the effectiveness of our proposed
516 routing strategy. While the performance difference is small, Noise-only Routing offers an additional
517 advantage: the ability to predict all used experts based on noise levels once before roll-outs, enabling
518 faster inference. This is particularly beneficial for robotics applications.

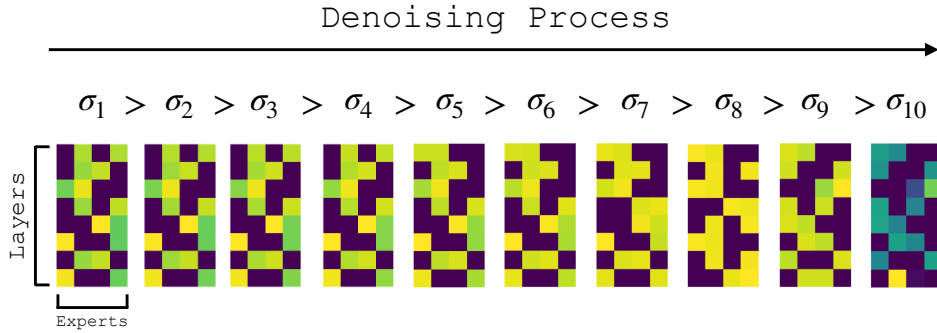


Figure 5: Visualized Expert Utilization. The average usage of all experts in MoDE across all layers is shown. Purple color corresponds to low usage and yellow color to high one, and each image is separately normalized. The average activation shows that MoDE learns to utilize different experts for different noise levels.

519 A.8.2 How does the model distribute the tokens to different experts?

520 To address Question IV, we analyzed how MoDE distributes tokens to different experts using a
 521 pre-trained model. Figure 5 visualizes the average usage of each expert in each model layer during
 522 inference across various noise levels, using 10 denoising steps for clarity. Our analysis reveals
 523 that MoDE learns to utilize different experts for various noise levels, suggesting that the router has
 524 specialized for different noise regimes. A transition in expert utilization occurs around σ_8 . In the first
 525 layer, the model learns an expert specialized for low-noise levels, primarily used in the last denoising
 526 step at σ_{\min} . These findings affirmatively answer Question IV, demonstrating that MoDE effectively
 527 distributes tokens across experts based on noise levels.

528 A.8.3 How does the model scale with more experts?

529 Finally, we analyze the effect of increasing the number of experts in MoDE. The results are presented
 530 in Figure 4b, where we evaluate MoDE on the CALVIN ABCD and CALVIN ABC benchmarks using
 531 2, 4, 6, and 8 experts. For comparison, we include two dense MoDE baselines: Dense-small and
 532 Dense-large. Dense-small shares the same latent dimensionality as MoDE, while Dense-large is scaled
 533 up to 2024 dimensions, matching MoDE’s overall parameter count. Our analysis focuses on how
 534 scaling affects both general performance (CALVIN ABCD) and zero-shot generalization (CALVIN
 535 ABC). In the ABCD environment, MoDE with 4 experts achieves the best performance. Interestingly,
 536 increasing beyond 4 experts degrades performance, possibly due to overfitting or increased routing
 537 complexity. In the zero-shot generalization (ABC), MoDE with 4 experts still performs well. Notably,
 538 the Dense-small variant consistently underperforms across both tasks, underscoring the efficiency of
 539 the MoE architecture in utilizing parameters more effectively. The Dense-small variant consistently
 540 underperforms. Overall, MoDE demonstrates that it can achieve comparable or superior performance
 541 to dense transformer models while requiring fewer computational resources.

542 A.9 State-based Experiments

543 We conduct additional experiments with MoDE on two established multi-task state-based environ-
 544 ments:

545 **Relay Kitchen.** We utilize the Franka Kitchen environment from [26] to evaluate models. This virtual
 546 kitchen environment allows human participants to manipulate seven objects using a VR interface:
 547 a kettle, a microwave, a sliding door, a hinged door, a light switch, and two burners. The resulting
 548 dataset consists of 566 demonstrations collected by the original researchers, where each participant
 549 performed four predetermined manipulation tasks per episode. The Franka Emika Panda robot is
 550 controlled via a 9-dimensional action space representing the robot’s joint and end-effector positions.

551 The 30-dimensional observation space contains information about the current state of the relevant
 552 objects in the environment. As a desired goal state, we randomly sample future states as a desired
 553 goal to reach.

554 For this experiment, we train all models for 40k training steps with a batch size of 1024 and evaluate
 555 them 100 times as done in prior work [27, 28, 4] to guarantee a fair evaluation. All reported results
 556 are averaged over 4 seeds. We train our models on a local PC RTX with an RTX 3070 GPU for
 557 approx. 2 hours for each run with the additional experimental rollouts.

558 **Block Push.** The PyBullet environment features an XArm robot tasked with pushing two blocks
 559 into two square targets within a plane. The desired order of pushing the blocks and the specific
 560 block-target combinations are sampled from the set of 1000 demonstrations as a desired goal state.
 561 The demonstrations used for training our models were collected using a hard-coded controller that
 562 selects a block to push first and independently chooses a target for that block. After pushing the first
 563 block to a target, the controller pushes the second block to the remaining target. This approach results
 564 in four possible modes of behavior, with additional stochasticity arising from the various ways of
 565 pushing a block into a target. The models only get a credit, if the blocks have been pushed in the
 566 correct target position and order. We consider a block successfully pushed if its center is within 0.05
 567 units of a target square.

568 All models were trained on a dataset of 1000 controller-generated demonstrations under these
 569 randomized conditions. All models have been trained for 60k steps with a batch size
 570 of 1024. To evaluate them we follow prior work [27, 28, 4] and test them on 100 dif-
 571 ferent instructions and report the average result over 4 seeds. We train our models on
 572 a local PC RTX 3070 GPU for approx. 3 hours for each run with a final evaluation.
 573 Demonstrations are sourced from a scripted oracle,
 574 which first pushes a randomly chosen block to a
 575 selected square, followed by the other block to a
 576 different square. The policies are conditioned to
 577 push the blocks in the desired configuration using
 578 a goal state-vector. We chose an action sequence
 579 length of 1 given a history length of 4 for these
 580 experiments, which are inspired by our dense dif-
 581 fusion transformer baseline BESO [4].

582 **Baselines.** In this setting, we compare **MoDE**
 583 against several SOTA goal-conditioned policies.
 584 We test two transformer architectures, C-BeT [28]
 585 and VQ-BeT [29], that predict discretized actions
 586 with an offset. C-BeT uses k-means clustering
 587 together with an offset vector while VQ-BeT leverages residual Vector Quantization to embed actions
 588 into a hierarchical latent space. Further, we test against a dense diffusion policy transformer model
 589 BESO [4]. BESO uses the same continuous-time diffusion policy combined with a dense transformer
 590 to predict a single action given a sequence of prior states. To enable a fair comparison, we chose the
 591 same hyperparameters for BESO and **MoDE** in both settings. We test all models averaged over 4
 592 seeds and report the mean values directly from prior work [29].

593 **Results.** The results of both experiments are summarized in Table 5. **MoDE** achieves a new SOTA
 594 performance on both benchmarks and outperforms the dense transformer variant of BESO in both
 595 settings. Further, **MoDE** achieves higher performance compared to other policy representation
 596 methods such as VQ-BeT and C-BeT.

	Block Push	Relay Kitchen
C-BeT	0.87±(0.07)	3.09±(0.12)
VQ-BeT	0.87±(0.02)	3.78±(0.04)
BESO	0.96±(0.02)	3.73±(0.05)
MoDE	0.97±(0.01)	3.79±(0.02)

Table 5: Comparison of the performance of different policies on the state-based goal-conditioned relay-kitchen and block-push environment averaged over 4 seeds. **MoDE** outperforms the dense transformer variant BESO and other policy representations on all baselines.

597 A.10 Related Work

598 B Related Work

599 **Diffusion in Robotics.** In recent years, Diffusion Models [30, 1, 10] have gained widespread
600 adoption in the context of robotics. They are used as a policy representation for Imitation Learning [5,
601 4, 31, 32, 33, 34, 35] and in Offline Reinforcement Learning [36, 37, 38]. Other applications of
602 diffusion models in robotics include robot design generation [39], video-generation [40, 41, 42] and
603 motion planning [43, 44]. The most common architecture for using diffusion models as a policy in
604 robotics is a CNN with additional FiLM conditioning [24] to guide the generation based on context
605 information. Recently, the transformer architecture has been adopted as a strong alternative backbone
606 for diffusion policies, specifically in IL. Examples include Octo [3], BESO [4] and 3D-Diffusion-
607 Actor [32]. However, no prior work considers using a Mixture of Experts architecture for improving
608 the computational efficiency and inference time and solely relies on dense transformer architectures.

609 **Mixture-of-Experts.** MoE are a class of models where information is selectively routed through
610 the model. The modern version of MoE was introduced in [45], where a routing or gating network
611 conditionally chooses a subset of experts to send an input to. After Transformers [14] proved to be an
612 effective model that scales well with data, they were modified to have expert feed-forward networks at
613 each block of the model in [15] which presented Switch Transformers. Switch Transformers laid the
614 groundwork that is still widely adopted in different Large-Language-Models (LLM) [46, 47]. This
615 allowed for more total parameters while keeping the forward and backward FLOPs smaller than their
616 dense counterpart, thus yielding significant performance gains. However, training both the router and
617 experts in parallel is a non-trivial optimization problem, and it can yield suboptimal solutions such as
618 expert collapse where experts learn similar functions instead of specializing [48]. In addition, router
619 collapse occurs when the router selects a small subset of the experts and doesn't utilize all the experts.
620 This is mitigated with load balancing losses [45, 15] which encourage the router to distribute inputs
621 more evenly across experts. Multiple works have explored different methods to perform routing, such
622 as expert choice routing [17], differential k-selection [49], frozen hashing functions [50], and linear
623 assignment [51].

624 In the context of robotics, MoE models are used in many settings without being combined with a
625 transformer architecture. Several works use a mixture of small MLP policies, that focus on different
626 skills in Reinforcement Learning [52, 53, 54] or for robot motion generation [55, 56], another body of
627 work utilizes combinations of small CNNs robot perception [57, 58]. Further applications include
628 learning multimodal behavior using a mixture of Gaussian policies [59, 60]. Despite the extensive
629 usage of MoE in many domains, no prior work has tried to utilize MoE together with Diffusion
630 Policies for scalable and more efficient Diffusion Policies.

631 **Multi Task Learning in Diffusion Models.** It has been shown that the denoising process is multi-task
632 [7]. Leveraging this idea, works have taken architectures that are suited for multi-task learning. Some
633 works have explicitly scheduled which parameters are specialized to which stage in the denoising
634 process [61, 62]. In extension to this [63] uses the scheduling as guidance during training but also
635 learns how to modulate representations based on the denoising stage. Finally, some works have
636 employed different architectures based on the denoising stage [64].

637 **Transformers for Robot Learning.** Transformer models have become the standard network archi-
638 tecture for many end-to-end robot learning policies in the last few years. They have been combined
639 with different policy representations in the context of IL. One area of research focuses on generating
640 sequences of actions with Variational Autoencoder (VAE) models [65, 66]. These action-chunking
641 transformer models typically use an encoder-decoder transformer as a policy architecture. Several Dif-
642 fusion Policies, such as Octo [3], BESO [4], ChainedDiffuser [31] and 3D-Diffusion-Actor leverage
643 a transformer model as a policy backbone. Another direction of research treats behavior generation
644 as discrete next-token predictions similar to auto-regressive language generation [67]. C-Bet, RT-1,
645 and RT-2 use discretized action binning to divide seen actions into k -classes [28, 27, 68, 69], while
646 VQ-BeT [29] learns latent actions with residual Vector Quantization. Several works have shown the

647 advantages of using pre-trained LLM or VLM as a policy backbone, which are then finetuned for
648 action generation [70, 71, 72, 18]. None of the recent work considers using any Mixture-of-Expert
649 architecture for policy learning. MoDE is the first architecture to leverage MoE architecture combined
650 with diffusion for behavior generation.

651 **B.1 Limitations**

652 MoDE still has certain limitations. In our experiments, we find that MoDE exhibits a slightly higher
653 standard deviation compared to the baselines. We hypothesize that the router’s initialization has
654 significant impact on overall optimization, requiring future work on stabilizing routing models. In
655 addition, when visualizing expert utilization, in some of our experiments we noticed that only a
656 subset of the total experts were being utilized - a phenomenon known as expert collapse [48]. In
657 addition to load balancing regularization, having more inductive biases that encourage the router to
658 fully utilize all experts are needed.