

Dex1B: Learning with 1B Demonstrations for Dexterous Manipulation

Jianglong Ye*, Keyi Wang*, Chengjing Yuan, Ruihan Yang, Yiquan Li, Jiyue Zhu
Yuzhe Qin, Xueyan Zou, Xiaolong Wang

UC San Diego

<https://jianglongye.com/dex1b>

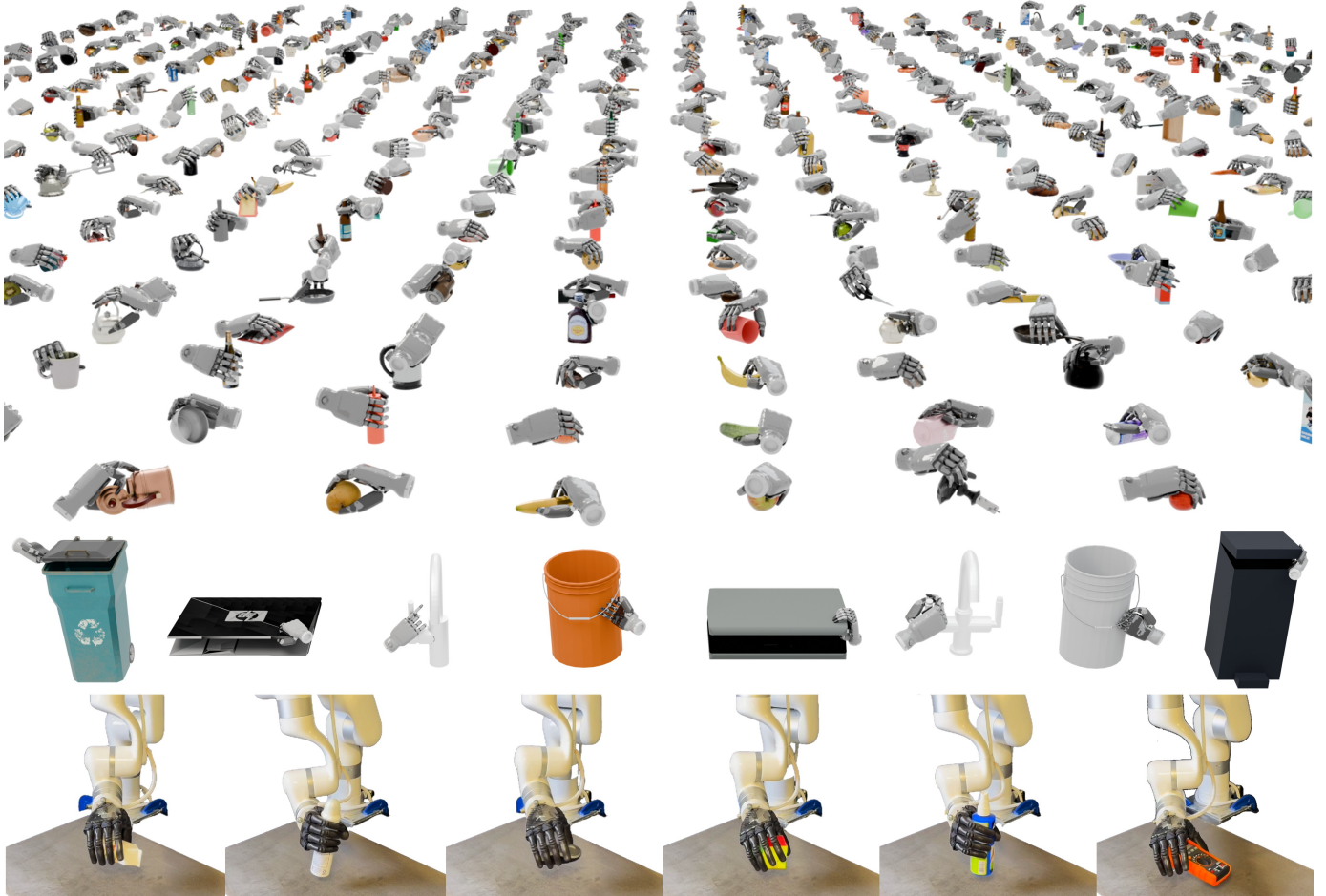


Fig. 1: The **Dex1B** benchmark consists of 1B generated high-quality demonstrations for grasping (top) and articulation (middle) tasks. At the bottom, we show the **direct sim-to-real** transfer results of our method **DexSimple** trained on Dex1B. This demonstrates that Dex1B is scalable and generalizable to real environments.

Abstract—Generating large-scale demonstrations for dexterous hand manipulation remains challenging, and several approaches have been proposed in recent years to address this. Among them, generative models have emerged as a promising paradigm, enabling the efficient creation of diverse and physically plausible demonstrations. In this paper, we introduce Dex1B, a large-scale, diverse, and high-quality demonstration dataset produced with generative models. The dataset contains one billion demonstrations for two fundamental tasks: grasping and articulation. To construct it, we propose a generative model that integrates geometric constraints to improve feasibility and applies addi-

tional conditions to enhance diversity. We validate the model on both established and newly introduced simulation benchmarks, where it significantly outperforms prior state-of-the-art methods. Furthermore, we demonstrate its effectiveness and robustness through real-world robot experiments.

I. INTRODUCTION

Dexterous manipulation with hand has been a long-standing topic in robotics. While its highly flexible and dynamic nature allows for more complex and robust manipulation skills, the high degrees of freedom (DoF) of a hand makes it very challenging to achieve its ideal function. In fact,

* Equal contribution.

with recent advancements in applications using parallel-jaw grippers [55, 18, 11, 56, 3], researchers in the community have started questioning the necessity of dexterous hands and having doubts about whether hands are only making problems harder.

We argue that dexterous hand is indeed valuable, but we just did not have enough data to capture the diverse and complex distributions required for effective dexterous manipulation. To address this data scarcity, previous works have explored various approaches, including human demonstrations [26, 38, 10, 9], optimization-based methods [47, 8, 46], reinforcement learning (RL)-based techniques [12, 57]. While these methods help generate demonstrations at a certain scale, they each have limitations: human annotation is costly and imprecise, optimization-based methods are slow and sensitive to initialization, and RL-based techniques lack data diversity.

Meanwhile, a large body of generative models [21, 27, 53, 24, 48] has been proposed in recent years to model the distribution of demonstration datasets, motivating us to explore how generative models can be leveraged for data generation. And we identify two key issues when applying generative models on data generation: i). **Feasibility**: The success rate of generative models is often lower than that of deterministic models. ii). **Diversity**: While generative models can produce more diverse actions than deterministic models, they still tend to interpolate between the seen demonstrations, which may maintain or even reduce the original level of diversity of whole dataset rather than expanding it.

To address the feasibility issue, we propose incorporating geometric constraints into the generative model, which significantly improves its performance. We also integrate optimization techniques with generative models, leveraging the strengths of both approaches: optimization ensures physically plausible results, while generative models enable efficient large-scale generation. To improve diversity, we introduce additional conditions to the generative model and prioritize sampling actions for less frequent condition values, encouraging the model to generate actions that differ more from existing ones in the dataset.

Our approach begins with an optimization-based method to construct a small yet high-quality seed dataset of dexterous manipulation demonstrations. We then train a generative model on this dataset and use it to scale up data generation efficiently. To mitigate biases introduced by optimization, we propose an debias mechanism, which systematically improves the diversity of generated data. This framework results in Dex1B, a dataset comprising one billion dexterous hand demonstrations, representing a substantial advancement in scale, diversity and quality over existing datasets. Compared to DexGraspNet [47], which operates on the object set of similar scale, our dataset offers $700\times$ more demonstrations, significantly enriching the available training data for learning-based models. Unlike previous approaches that rely solely on human annotation or optimization, our method combines optimization and neural networks, achieving a superior balance between cost, efficiency, and data quality.

Dataset	Task	Num Objects	Num Demonstrations	Method
DDG [23]	Grasping	565	6.9K	GraspIt
DexYCB [7]	Grasping	20	1K	Annot.
ContactPose [6]	Grasping	25	2306	Capture
RealDex [26]	Grasping	52	2.6K	Capture
DexGraspNet [47]	Grasping	5K	1.32M	Optim.
SynH2R [12]	Handover	1174	6K	Optim. + RL
RP1M [57]	Piano	-	1M	RL
Dex1B (Ours)	Grasping + Arti.	6K	1B	Optim. + Gen.

TABLE I: **Dataset comparison** on dexterous robotic manipulation. Our proposed Dex1B is a multi-task trajectory-level dataset with 1B demonstrations. Optim., Arti., Gen. are short for optimization, articulation and generative models.

To effectively leverage the scale and diversity of Dex1B, we introduce DexSimple, a new baseline that extends prior work [21] by incorporating conditional generation and enhanced loss functions. Despite its simplicity, DexSimple benefits from the scale and diversity of Dex1B, achieving state-of-the-art (SoTA) performance across dexterous manipulation tasks.

Our key contributions are as follows:

- We introduce a novel iterative data generation pipeline that combines optimization and generative models to generate large-scale dexterous demonstrations for grasping and articulation tasks. Using this approach, we construct Dex1B, the largest and most diverse dexterous demonstration dataset to date.
- We propose a simple yet effective baseline method that incorporates enhanced loss functions while supporting conditional generation, making it particularly well-suited for our iterative pipeline and policy deployment.
- Leveraging Dex1B and DexSimple, we achieve a 22% performance improvement over previous best methods on grasp synthesis tasks, setting a new benchmark in dexterous manipulation.

II. RELATED WORK

Dexterous Hand Manipulation. Dexterous hand manipulation is a pivotal area in robotics, concentrating on precise, multi-fingered grasping and manipulation. Early research in this field primarily addressed control-based methods, laying the foundation for dexterous manipulation through studies in caging to grasping techniques [39], grasp synthesis [40], and manipulability in underactuated hands [36]. Further developments in control-based approaches simplified the search space for multi-fingered grasps [34, 35], as well as optimization processes [58, 15], leading to highly precise, agile, and safe manipulation. However, these methods generally lack generalization across diverse environments and use cases.

Subsequent research shifted towards learning-based approaches to enhance flexibility and scalability [1, 33]. This includes generating pose vectors directly [22, 14, 52], utilizing intermediate representations [41, 49], and leveraging contact maps [5, 45]. While learning-based methods offer increased adaptability, they are sensitive to data quality and scope, a

limitation addressed in this work. Recent works [29, 42] employ RL-based methods to transfer robust grasping capabilities of dexterous robotic hands to the real world, using point cloud and RGB inputs, respectively.

Manipulation Dataset. Existing manipulation datasets can be broadly categorized into synthetic and real-world collections. Synthetic datasets like ObMan [20] and DDGdata [30] utilize the GraspIt planner to generate grasping poses but suffer from limited diversity due to naive search strategies. Real-world datasets with human hand poses offer more natural interactions, such as HO3D [19] which leverages 2D keypoint annotations and physics constraints, and DexYCB [7] which captures multi-view RGBD recordings. ContactDB [4] and ContactPose [6] further enhance grasp understanding by incorporating thermal imaging and detailed contact information, respectively. However, these real-world datasets are inherently restricted to human hand structures and common daily hand poses. In contrast, our approach leverages optimization and neural networks to generate diverse manipulation trajectories that transcend these limitations. Our work is related to prior work [28], which has a similar goal of expanding the distribution of generative models. Our approach differs in two ways: (1) Scale: our pipeline aims to generate 1B demonstrations rather than thousands. (2) Diversity: we increase diversity by conditioning on object geometry instead of using a classifier. The recent work DexGraspNet2.0 [54] proposes learning a diffusion model over large-scale optimized demonstration datasets. In contrast, our work integrates the generative model directly into the data generation pipeline instead of relying solely on optimization. Additionally, we unify grasping and articulation tasks in both our model design and benchmark, while they focus on grasping only. We presents the differences of several representative manipulation datasets in Tab. I.

III. DEX1B BENCHMARK

We introduce a comprehensive benchmark for two fundamental dexterous manipulation tasks: **grasping** and **articulation**. In the grasping task, the robot hand must reach for and lift an object, whereas the articulation task requires the hand to manipulate an articulated object to achieve a specific degree of opening. Our benchmark consists of over 6,000 diverse objects and provides one billion demonstrations across three dexterous hands: the Shadow Hand, the Inspire Hand, and the Ability Hand. Each demonstration consists of a complete action sequence, from initial reaching to object manipulation.

To generate these demonstrations, we synthesize key hand poses at critical interaction points with the object, while the remaining action sequences—such as reaching, lifting, and opening—are generated using motion planning. The evaluation of our benchmark is conducted with ManiSkill [44, 50].

Overview of Data Generation. Broadly, hand pose generation for dexterous manipulation can be approached through optimization-based methods or generative models. While optimization methods can be effective, they are often computationally expensive, especially for large-scale generation, and tend to bias the dataset toward simpler cases. On the other hand,

generative models rely on an initial dataset to learn meaningful data distributions. In this work, we integrate both approaches to leverage their strengths.

As illustrated in Figure 2, we begin by constructing a small-scale seed dataset using optimization. This seed dataset serves as the foundation for training a generative model to learn its underlying data distribution. The trained generative model is then used to produce additional demonstrations. However, since the generative model inherently inherits the biases of the seed dataset, we introduce a debiasing strategy to enhance diversity. Specifically, we condition the generative model on targeted factors to generate hand poses under less frequently observed conditions, thereby expanding the dataset beyond the initial distribution. By iteratively refining the generative model through repeated training and debiasing operations, we construct our final dataset, Dex1B, which achieves both diversity and robustness in dexterous manipulation demonstrations.

Optimization for Seed Dataset. To generate the seed dataset, we implement an efficient optimization method for hand pose synthesis based on previous work [47, 8], while including new features like scene-level collision avoidance and support for various hands. Although the optimization process is well-engineered (1,000 grasps per minutes on a single GPU), generating one billion demonstrations remains computationally expensive. Therefore, we only use optimization to create a small-scale seed dataset (around 5 million poses). A dexterous hand pose is parameterized by a tuple $g = (T, R, \theta)$, where $T \in \mathbb{R}^3$ for global translation, $R \in SO(3)$ for global rotation, and $\theta \in \mathbb{R}^d$ for robot hand joint angles ($d = 22$ for Shadow hand, $d = 6$ for Inspire and Ability hand). The object geometry is represented by mesh O . Unlike previous methods [47, 8], which use link meshes to model hand geometry, we approximate the hand using manually defined spheres (around 10 spheres for each link). This approximation significantly accelerates the optimization process.

Our optimization energy function for the grasping task E_{grasp} is given by:

$$E_{\text{grasp}} = E_{\text{fc}} + w_{\text{dis}}E_{\text{dis}} + w_{\text{sdf}}E_{\text{sdf}} + w_jE_j + w_sE_s,$$

where E_{fc} is the force closure energy term, E_{dis} minimizes contact distance, E_{sdf} prevents penetration, E_j enforces joint limits, and E_s avoids self-collisions. w_x represents the weight for the corresponding terms. The penetration term E_{sdf} is formulated based on the sphere-based hand representation. A sphere is defined by (c, r) , where the center $c \in \mathbb{R}^3$ and the radius $r \in \mathbb{R}^+$. The center c is transformed with the link pose using forward kinematics. The SDF penetration term is given by:

$$E_{\text{sdf}} = \sum_i r_i - SDF(c_i, O),$$

where $SDF(\cdot)$ is the signed distance function (SDF) query between a 3D point and a mesh. To prevent scene collisions, we incorporate SDF queries for other meshes (e.g., the table), and the final SDF for a sphere is computed as the maximum value across different meshes. Detailed formulation of other energy terms can be found in previous works [47, 25].

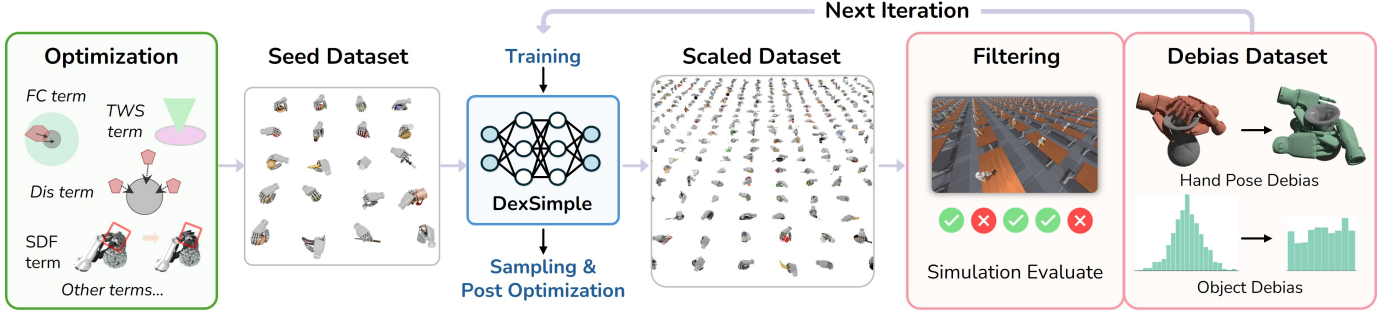


Fig. 2: **DexIB** demonstration collection. The engine takes object assets and hand pose initialization as input, using a control-based optimization algorithm to generate the Seed dataset. Then the Seed dataset is used as the training data for DexSimple, else for Dex1Bfor the last iteration. Then DexSimple will generate a scaled proposal dataset with π as the scaling ratio. For the proposal dataset, we then use the simulation critic and debiased algorithm to create the debiased dataset for optimization refinement.

While the force closure energy term E_{fc} is suitable for the grasping task, achieving force closure in the articulation task is usually difficult and unnecessary. Instead, the articulation task requires the hand pose to generate specific forces and torques, such as rotating the top part of a laptop or pulling a drawer. Therefore, we replace the force closure energy term E_{fc} with the task wrench space term E_{tws} introduced in [8], which approximates the difference between a target wrench space and the current wrench space of a given hand pose. The optimization energy function for the articulation task E_{arti} is defined as:

$$E_{arti} = E_{tws} + w_{dis}E_{dis} + w_{sdf}E_{sdf} + w_jE_j + w_sE_s.$$

The target wrench space for articulated objects with revolute joints consists of a torque aligned with the joint axis and arbitrary forces. The wrench space for objects with prismatic joints consists of forces sampled from a 30° cone aligned with the joint axis and zero torque.

The optimization process is implemented using WarpLang [31] and accelerated with its BVH mesh structure. The optimized hand poses are evaluated using the simulator, and the successful ones are retained as a seed dataset.

Generative Models for Scaling-up Demonstrations. Generative models are widely adopted for capturing the distribution of action demonstrations. However, applying these models for data generation still presents several challenges: i). **Feasibility:** The success rate of generative models is often lower than that of deterministic models, leading to a higher proportion of infeasible samples. ii). **Limited Diversity:** While generative models can produce more diverse actions than deterministic models, they still tend to interpolate between the demonstrations, which may maintain or even reduce the original level of diversity of whole dataset rather than expanding it.

To address the feasibility issue, we first incorporate geometric constraints during the generation process, enabling our model to outperform state-of-the-art generative models (see Sec. IV for model details and Sec. V for experiments). In addition, we apply a post-optimization step to the sampled hand poses to prevent penetration and ensure that the fingers

closely cover the object. The energy function for the post-optimization stage of both tasks, E_{post} , is defined as:

$$E_{post} = w_{dis}E_{dis} + w_{sdf}E_{sdf} + w_jE_j + w_sE_s.$$

We exclude the task-specific terms E_{fc} and E_{tws} here since we only aim to make slight adjustments to the finger positions. Although we use optimization in this stage, the overall data generation, combined with generative models, remains significantly more efficient than pure optimization. The sampling process is approximately 100 times faster, and the number of iterations in post-optimization is substantially lower than in the pure optimization (100 vs. 6000). The refined hand poses are then evaluated using the simulator.

To improve diversity, we encourage the generative model to sample actions that differ more from existing actions in the dataset while maintaining success rate. To achieve this, we introduce an additional condition to the generative model and prioritize sampling actions for less frequent conditions. Specifically, we associate each hand pose with a single 3D point on the object. We first define the heading direction $v \in \mathbb{R}^3$ of a hand pose as the vector from the palm center to the midpoint between the thumb tip and the middle finger tip. The closest point along this direction is then assigned as the associated point of the hand pose. We adapt our generative model to take the feature vector of a 3D point as a condition for generating corresponding actions. During data generation, we first statistically compute the probability of each point associated with existing actions on the object and then sample new actions inversely proportional to this probability. Additionally, we statistically count the number of existing actions for each object and sample more actions for the more challenging ones.

After increasing the dataset size and diversity, retraining the model on the expanded dataset can further improve its performance. This *iterative data generation* process can be repeated multiple times to progressively refine both the model and the dataset.

Motion Planning. With the key-frame action, motion planning is still required to complete both tasks. For the reaching

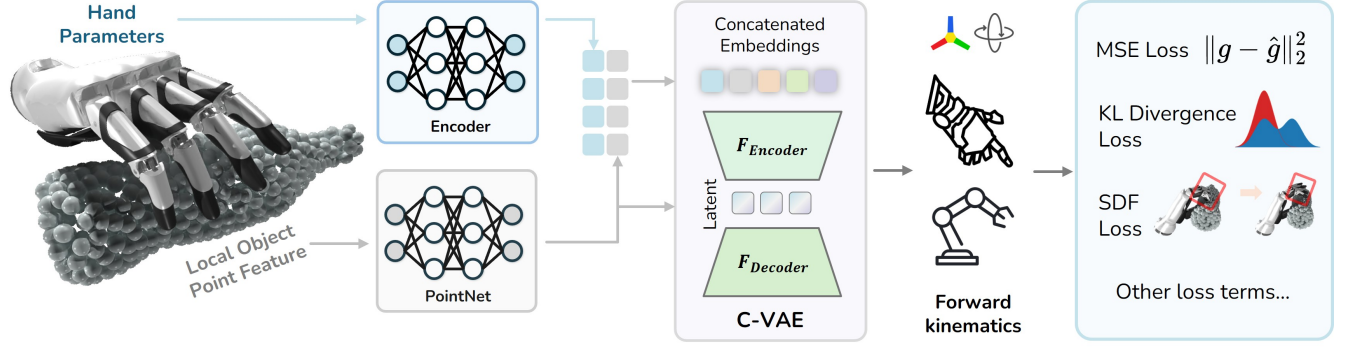


Fig. 3: **DexSimple Pipeline.** Our model takes in hand parameters and object point clouds as fixed input for CVAE, while root rotation, translation, and joint value as optional conditions. Those are combined as the input embeddings for CVAE, while point cloud embeddings are re-emphasized at the latent space. The output of CVAE is the forward kinematics define the hand pose trajectory optimized by the effective loss function.

stage in both tasks, motion planning can be formulated as an optimization problem that maximizes smoothness while avoiding collisions. Given a manually defined starting action and a goal action, we first linearly interpolate between them, and then optimize the intermediate actions using the following energy function:

$$E_{\text{reach}} = w_{\text{smooth}} \sum_{i=1}^N \|g_i - g_{i-1}\|^2 + w_{\text{sdf}} \sum_{i=0}^N E_{\text{sdf}}(g_i),$$

where $\{g_0, g_1, \dots, g_N\}$ denotes the sequence of hand poses along the trajectory, w_{smooth} and w_{sdf} are weights for smoothness and collision avoidance respectively, and $E_{\text{sdf}}(g_i)$ is computed based on the sphere-based SDF errors with respect to nearby scene meshes. Minimizing E_{reach} produces a trajectory that is both smooth and collision-free. Compared to other motion planning libraries, this simple optimization is naturally suitable for large-scale parallel data generation.

After reaching, for the grasping task, we will execute an over-shoot action to grasp the object and increase the height of the target action to lift it. For the articulation task, we will follow the trajectory based on the given joint axis (rotate along the revolute joint axis or translate along the prismatic joint axis). The complete planned trajectory is executed in the simulator for evaluation.

IV. DEXSIMPLE MODEL

While a large body of generative models [21, 27, 53, 24, 48] have been proposed for dexterous hand manipulation in recent years, their use for data generation or policy deployment remains limited. In this work, we revisit the simple CVAE model and demonstrate that incorporating an SDF-based geometric constraint during training enables it to outperform state-of-the-art methods by a large margin. Furthermore, we integrate additional condition over the base model to support diverse data generation.

Vision Encoder and CVAE. We employ a point cloud $P \in \mathbb{R}^{N \times 3}$ as the visual input, using a full point cloud sampled

from the object mesh for data generation and a single-view depth map for policy deployment. We utilize PointNet [37] to encode the point cloud into a global feature vector $f_{\text{obj}} \in \mathbb{R}^d$ and local feature vectors $f_p \in \mathbb{R}^d$ for each point $p \in \mathbb{R}^3$:

$$f_{\text{obj}}, \{f_p\}_{p \in P} = \text{PointNet}(P).$$

The VAE model uses a multi-layer perceptron (MLP) to encode the hand pose g into the mean and standard deviation vectors of a latent distribution. A sample is drawn from this distribution and passed to the MLP decoder to reconstruct the original hand pose. After concatenating conditional vectors (e.g., the global point cloud feature vector f_{obj}) to both the inputs of the VAE encoder and decoder, the CVAE model can generate samples under a given condition:

$$\begin{aligned} \mu, \sigma &= \text{Enc}(g, f_{\text{obj}}), \\ z &= \mu + \sigma \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \\ \hat{g} &= \text{Dec}(z, f_{\text{obj}}). \end{aligned} \quad (1)$$

In our work, we simply concatenate additional vectors to incorporate more conditions.

Objectives. The CVAE training is supervised by the standard reconstruction loss $\mathcal{L}_R = \|g - \hat{g}\|_2^2$ and the KL divergence loss $\mathcal{L}_{\text{KL}} = D_{\text{KL}}(\mathcal{N}(\mu, \sigma^2) \parallel \mathcal{N}(0, I))$. To enforce geometric constraints, we introduce an SDF-based loss. Different from the optimization stage, building a BVH structure for each object at each iteration during training is time-consuming. We therefore use a sampled point cloud to represent the object geometry while still using spheres to represent the hand geometry. The SDF loss is then given by:

$$\mathcal{L}_{\text{SDF}} = \sum_{c \in \mathcal{C}} \max\left(0, r_c - \min_{p \in P} \|c - p\|_2\right)^2,$$

where \mathcal{C} is the set of spheres representing the hand geometry and r_c is the radius of sphere c . This loss is not only straightforward to implement in PyTorch, but we also empirically find that training with the point-sphere SDF loss is more

stable compared to the mesh-sphere SDF loss used in the optimization stage. With this SDF loss term, our proposed DexSimple outperforms SOTA methods by a large margin (see Sec.V).

The overall loss for the base model is defined as:

$$\mathcal{L} = \lambda_R \mathcal{L}_R + \lambda_{KL} \mathcal{L}_{KL} + \lambda_{SDF} \mathcal{L}_{SDF},$$

where λ_R , λ_{KL} , and λ_{SDF} are weights that balance different loss terms.

Conditions for Data Generation. As mentioned in Sec. III, each hand pose is associated with a single object 3D point p by finding the closest point along its heading direction v . To achieve this, we condition the CVAE on the corresponding local object feature vector f_p .

V. EXPERIMENTS

We firstly evaluate the effectiveness of the proposed generative model, DexSimple, for grasp synthesis on the DexGraspNet [47] benchmark. Then, we provide details on the synthesized Dex1B demonstration dataset and compare it with human-annotated demonstration datasets on both lifting and articulation tasks. Additionally, we downscale Dex1B and evaluate the performance of methods trained on it. Finally, ablation studies are conducted to validate our design choices.

A. Grasping Synthesis Evaluation

Grasping is essential in most manipulation tasks, we firstly evaluate the proposed method’s effectiveness in grasp synthesis using the DexGraspNet [47] benchmark. We train DexSimple solely with the benchmark’s provided training data, reducing the output to a single frame and omitting conditioning during training. The validation is conducted in the Isaac Gym simulator [32] using ShadowHand [13]. We adhere to the metrics established in the benchmark to ensure fair comparisons with baseline methods, which are divided into two categories: quality (Success Rate, Q_1 -score, Penetration) and diversity (H mean and H std). We follow the implementations¹ in [47, 27]. We compare with DDG [23], GraspTTA [21], the generation module in UniDexGrasp [51] (abbreviated as UDG), and UGG [27].

We present quantitative results in Table II. Many grasp generation methods, such as UGG [27], commonly employ post-optimization to enhance performance. To ensure a fair comparison, we indicate the use of post-optimization (abbreviated as “Opt”) in the table. The results show that the proposed generative model, DexSimple, outperforms all baseline methods by a large margin. In terms of quality, DexSimple (with post-optimization) achieves the highest success rate (86.0%), the highest Q_1 score (0.125), and the lowest penetration (0.13). For diversity, DexSimple outperforms baseline with a higher entropy mean of 8.56.

UGG [27] proposes a learning-based discriminator to filter grasping, which can be applied to our method. With this filtering, the success rate increases to 92.6%. It is worth noting

Method	Setting		Quality			Diversity	
	Opt	Filter	SR \uparrow	Q_1 \uparrow	Pen \downarrow	H mean \uparrow	H std \downarrow
DDG [23]			67.5	0.058	0.17	5.68	1.99
UGG [27]			43.6	0.026	0.43	8.33	0.30
DexSimple			63.7	0.075	0.29	8.53	0.25
UDG [51]	✓		23.3	0.056	0.15	6.89	0.08
GraspTTA [21]	✓		24.5	0.027	0.68	6.11	0.56
UGG [27]	✓		64.1	0.036	0.17	8.31	0.28
DexSimple	✓		86.0	0.125	0.13	8.56	0.15
UGG [27]	✓	✓	72.7	0.063	0.14	7.17	0.07
DexSimple	✓	✓	92.6	0.132	0.12	8.56	0.16

TABLE II: **Grasping synthesis** results on the DexGraspNet [47] benchmark. The proposed generative model, DexSimple, significantly outperforms all baseline methods. Some evaluation results are taken from UGG [27]. Opt, SR, and Pen are short for Optimization, Success Rate, and Penetration, respectively.

the success rate of DexSimple without post-optimization and filtering is slightly lower than that of DDG [23]; this is expected as our method is a generative model while DDG [23] is a regression model, and our method achieves much higher diversity (8.53 vs. 5.68).

B. Dataset Analysis

Tasks Definition. While the proposed iterative data generation pipeline can be applied to multiple hands, we take ShadowHand as an example to detail our data curation process. Beyond the grasping synthesis task in Sec. V-A, we focus on two tasks: **Grasping** and **Articulation**. The goal of the grasping task is to reach an object placed on a table, grasp it, and lift it to a specified height (0.4 m) while maintaining contact between the hand and the object. We show example trajectories from Dex1B for the grasping task in Fig. 5. The articulation task requires reaching an articulated structure (e.g., a laptop, box, or faucet) and opening it to increase its joint angle by 0.5 while maintaining contact between the hand and the object.

For the grasping task, we utilize all 5751 object assets collected by DexGraspNet [47] and exclude all objects that cannot stand stably on the table. Note that the settings between DexGraspNet and our dataset are different. Our focus is on table-top tasks, while DexGraspNet focuses on grasping objects in free space. Approximately 90% of grasping demonstration in DexGraspNet would collide with the table. For the articulation task, we utilize 650+ articulation assets collected by DexArt [2] and ManiSkill [44]. We adopt the official training/testing splits provided by previous works [47, 16, 2]. **Dataset Curation.** Our dataset curation starts by optimizing a seed dataset with 5 million demonstrations. We first train DexSimple on this seed dataset to create a 50 million proposal dataset. After optimization refinement, we use the ManiSkill [44]/SAPIEN [50] simulation to filter out unsuccessful trajectories and rebalance the data as described in Sec. III. We then retrain DexSimple on the 50 million dataset to produce a 500 million proposal dataset, and repeat this

¹Please refer to supp. material for details.

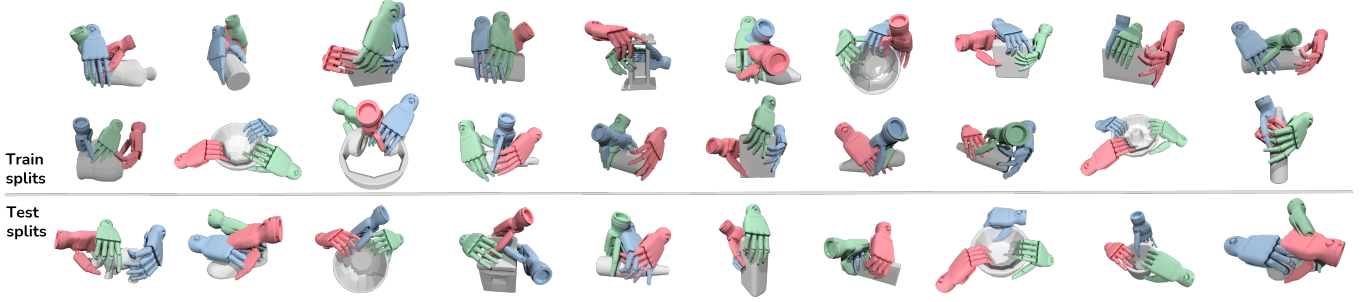


Fig. 4: **Diverse demonstrations** for objects from train/test splits. We show only the contact frame for clarity.

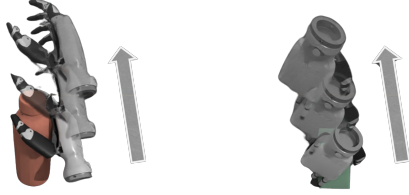


Fig. 5: **Lifting trajectory** from Dex1B dataset.

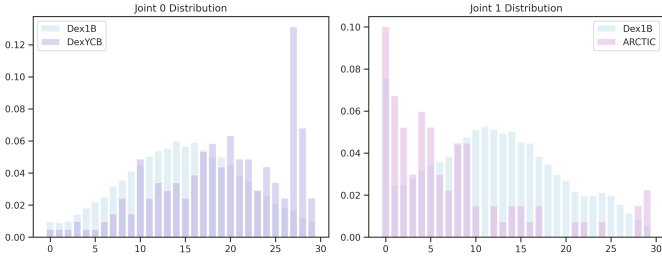


Fig. 6: **Probability distribution of joint values** from Dex1B and DexYCB/ARCTIC. The distribution of Dex1B is more evenly spread, centering around the mean joint values.

process. Finally, we collect a dataset with 950 million (around 1 billion) successful trajectories.

Implementation Details. Compared to DexGraspNet [47], our implementation of pure optimization-based grasp generation is 30 times faster, requiring only 2 minutes to generate 2000 grasps for 6000 steps on a single RTX-3090, while also achieving a higher success rate (27% vs. 20% for Shadow-Hand). When initialized with the network, our method is even 700 times faster, including CVAE sampling and 100 post-optimization steps. All CUDA and geometry-related operations are implemented using warp-lang [31].

In the proposed DexSimple, we adopt PointNet [37] as the visual encoder, extracting object features in a dimension of 256. For the CVAE, both the input and output dimensions are $N_{\text{frame}} \times N_{\text{DOF}}$, and the latent vector dimension is set to 256. **Dataset Comparison.** We demonstrate the quality of Dex1B by comparing it to two large-scale, human-annotated, trajectory-level datasets: DexYCB [7] and ARCTIC [16]. DexYCB includes 20 objects, with approximately 500 trajectories for each hand. ARCTIC includes 10 articulation objects with a total of 301 trajectories. We follow [53, 10] to generate robot demonstrations from the DexYCB and ARCTIC

Method	Training Data	Eval on DexYCB		Eval on Dex1B	
		Train set	Test set	Train set	Test set
BC w. PointNet	DexYCB [7]	34.72	3.03	1.02	2.56
DexSimple	DexYCB [7]	43.49	21.21	23.68	22.80
BC w. PointNet	Dex1B (ours)	33.02	31.82	31.40	28.54
DexSimple	Dex1B (ours)	47.17	53.02	49.58	45.40

(a) Lifting task comparison on DexYCB [7] and Dex1B.

Method	Training Data	Eval on ARCTIC		Eval on Dex1B	
		Train set	Test set	Train set	Test set
BC w. PointNet	ARCTIC [16]	41.03	25.62	37.65	30.16
DexSimple	ARCTIC [16]	48.75	23.08	49.16	51.57
BC w. PointNet	Dex1B (ours)	57.50	63.67	64.74	56.88
DexSimple	Dex1B (ours)	72.00	73.49	77.05	64.79

(b) Articulation task comparison on ARCTIC [16] and Dex1B.

TABLE III: Benchmarks on (a) **lifting tasks** with DexYCB [7] and our datasets, and (b) **articulation tasks** with ARCTIC [16] and our datasets. Models trained on Dex1B consistently outperform those trained on DexYCB/ARCTIC across various tasks, baselines, and splits.

dataset, including retargeting human demonstrations to robot trajectories and adding noise to generate a larger number of physically plausible demonstrations. We collected 62% and 64% of all trajectories from the DexYCB and ARCTIC datasets, respectively, that successfully achieve task goals.

We highlight the diversity of Dex1B. Using the proposed debiasing approach in Sec.III, the diversity across object categories, joint values, and wrist poses in our dataset can be easily enhanced by generating additional samples for underrepresented data. Besides, we discretize the range of joint angles into bins and estimate a probability distribution over them. Fig. 6 shows the distribution of two joints. Unlike DexYCB and ARCTIC, which often have joint values concentrate at the limits, the distribution of Dex1B is more evenly spread, centering around the mean joint values. This is achieved by debiasing approach and including a regularization term that discourages the hand from getting too close to the joint limits during optimization. Qualitative results are presented in Fig. 4.

Benchmarks. We benchmark two methods for grasping and articulation tasks on our datasets, and compare them with the

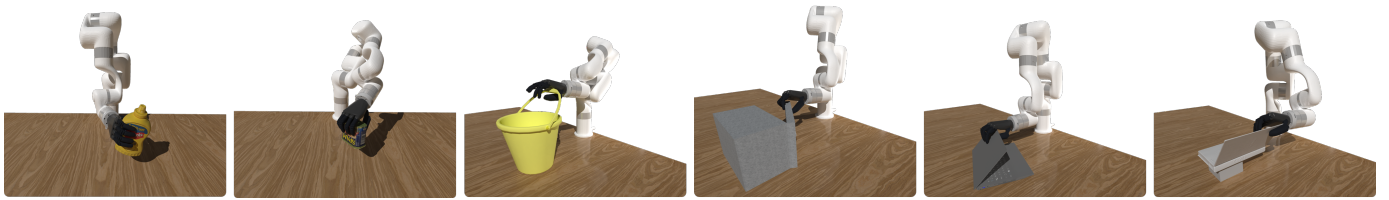


Fig. 7: **Qualitative results** for both grasping and articulation tasks. We show only the contact frame for clarity.

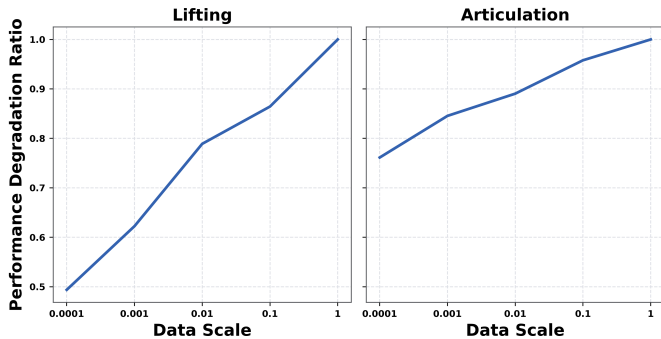


Fig. 8: **Scaling the number** of demonstrations used for training. For both tasks, our model consistently improves with more training data.

same methods trained on DexYCB [7] and ARCTIC [16]. In addition to the proposed DexSimple, we implement a vanilla behavioral cloning with PointNet [37] (referred to as BC w. PointNet). This model takes the object point cloud, current hand joint values, and poses as input to predict chunked actions for the next $n = 50$ frames. The predicted actions are then merged using a temporal weighting technique from ACT [55].

The results are reported in Tab. III. When comparing models trained on Dex1B to those trained on DexYCB/ARCTIC, we consistently find that the former outperforms the latter across tasks, baselines and splits. This suggests that supervised learning methods perform better when trained on our larger and more diverse Dex1B dataset. Tab. III also demonstrates that the proposed generative method, DexSimple, achieves better performance than the regression-based BC baselines on both the relatively small DexYCB/ARCTIC dataset and the larger-scale Dex1B. For lifting task, it also can be clearly observed that models trained on DexYCB struggle to generalize to unseen objects. Qualitative results are shown in Fig. 7.

C. Scaling the Dataset

To investigate the effect of training data size on performance, we reduce the amount of training data and analyze its impact on the success rates of both the lifting and articulation tasks. As shown in Fig. 8, the performance degradation ratio increases as data is reduced, illustrating that the success rates of the proposed DexSimple consistently improve with more training data. Notably, we observe that performance degradation is more pronounced for the lifting task than for the articulation task as training data decreases. We hypothesize this is because lifting relies heavily on stable object grasping, requiring a precise geometric understanding of individual objects, which becomes more challenging with reduced data.

Method	Quality			Diversity	
	Success Rate \uparrow	Q_1 \uparrow	Penetration \downarrow	H mean \uparrow	H std \downarrow
w/o. \mathcal{L}_{sdf}	0.7	0.001	0.92	8.58	0.16
w/o. \mathcal{L}_D	42.0	0.044	0.23	8.65	0.16
Full Model	63.7	0.075	0.29	8.53	0.25

TABLE IV: **Ablation Study** of the geometric loss terms in grasp synthesis. Both \mathcal{L}_{sdf} and \mathcal{L}_D are crucial for the grasping quality.

In contrast, the articulation task, which emphasizes trajectory execution, shows greater resilience to data reduction as it can adapt to unseen objects through a more generalized approach to motion. This suggests that while both tasks benefit from larger datasets, lifting requires a more extensive dataset to achieve stable performance, whereas articulation maintains reasonable performance even with less data.

D. Ablation Study

In this section, we ablate the geometric loss terms of the DexSimple generative model using the DexGraspNet benchmark, with results detailed in Tab. IV. The *sphere-representation SDF loss*, denoted as \mathcal{L}_{sdf} , is designed to provide fine-grained geometric guidance, which plays a crucial role in preventing the model from penetrating the object during grasp synthesis. This loss term is essential for grasping quality, as *removing \mathcal{L}_{sdf} causes a drastic drop in success rate from 63.7 to 0.7*. Without \mathcal{L}_{sdf} , the model lacks precise spatial awareness, leading to significant failures in grasp execution. On the other hand, the *distance loss \mathcal{L}_D* is responsible for encouraging the hand to make stable contact with the object surface, which enhances the grasp’s stability. This loss has a notable impact on both the success rate and the Q_1 quality metric. Although \mathcal{L}_D *slightly increases the penetration value*, it significantly contributes to an improved success rate and Q_1 score, highlighting its importance in achieving reliable grasps. The diversity metrics, represented by H mean and H std, are only minimally impacted by both loss terms, indicating that these geometric losses focus more on grasp quality than diversity. In summary, *both \mathcal{L}_{sdf} and \mathcal{L}_D are indispensable for high-quality grasp synthesis*, as they address different aspects of the grasping process—object penetration prevention and stable contact establishment, respectively.

VI. REAL-WORLD EXPERIMENTS

We demonstrate the effectiveness of the proposed method in the real world through direct sim-to-real deployment. We explore two platforms: xArm with an Ability Hand and H1

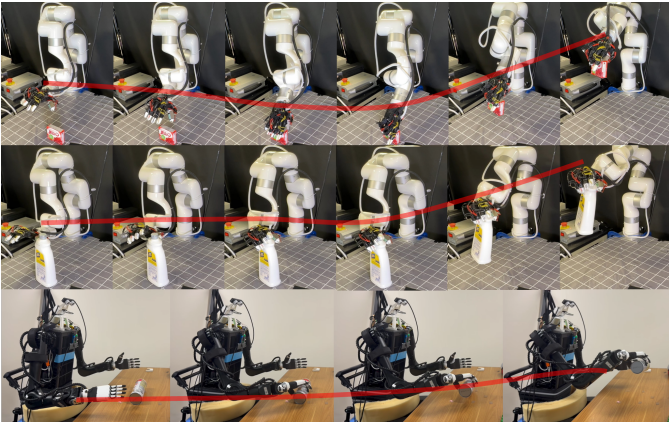


Fig. 9: We directly deploy the predicted pose to demonstrate the effectiveness of the proposed method in the real world.

Method	Obj-1	Obj-2	Obj-3	Obj-4	Obj-5	Obj-6	Obj-7	Obj-8	Obj-9	Obj-10	Mean
DexSampler	2/5	3/5	5/5	3/5	1/5	4/5	2/5	5/5	1/5	3/5	58%
Ours	4/5	5/5	5/5	5/5	5/5	5/5	5/5	5/5	5/5	4/5	96%

with an Inspire Hand. We mount a camera in a third-person view for xArm and an egocentric view for H1. The camera pose is calibrated using hand-eye calibration. We then take the partial point cloud observation from the camera as input to the model. Additionally, we sample 128 poses and select a valid inverse kinematics (IK) solution. Finally, we apply motion planning to this pose. The successful grasping trajectories are visualized in Fig. 9.

We compare with DexDiffuser [48] on *XArm+Ability Hand*. We evaluate 5 trials per object on 10 unseen objects (shown in the figure). Both methods take partial point clouds as input and use motion planning to execute the grasp poses. As DexDiffuser is originally trained for the Allegro Hand, we retrain its model using our dataset and platform. We only retrained the DexSampler module and omitted the DexEvaluator used in their paper. The results below demonstrate the better sim-to-real transfer performance of our proposed method.



VII. CONCLUSION AND LIMITATIONS

In this paper, we present **Dex1B**, a synthetic dataset for dexterous hand manipulation, containing 1 billion demonstrations. We introduce an iterative data generation pipeline that integrates optimization techniques with learning-based approaches to efficiently generate manipulation demonstrations. It begins by creating the see dataset using pure optimization methods, which is then used to train our generative model, **DexSimple**. The model accelerates the data generation loop by producing a proposal dataset that is further refined through optimization. The refined datasets are subsequently verified and debiased

for quality assurance. In our experiments, we demonstrate that DexSimple, enhanced with geometric loss, achieves a 22-point improvement over previous state-of-the-art methods on the DexGraspNet benchmark. Additionally, benchmarks for lifting and articulation tasks highlight the effectiveness of both Dex1B and DexSimple, showcasing their utility in advancing dexterous hand manipulation research.

Limitations. We state several key limitations of our method here: i) Since our method focuses on key-frame action generation, it operates in an open-loop manner when deployed in the real world, making it prone to the sim-to-real gap and control/observation errors. ii) Although the generative model is significantly more efficient than optimization, our method still relies on simulation to filter successful data. The simulation itself remains time-consuming, and reducing its runtime during data generation is a potential direction for future work. iii) Our method mainly considers single-object scene settings. For multi-object scenarios, a stronger vision backbone may be necessary.

VIII. ACKNOWLEDGMENT

This work was supported, in part, by NSF CAREER Award IIS-2240014, Qualcomm Innovation Fellowship, and gifts from Meta.

REFERENCES

- [1] OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning dexterous in-hand manipulation. *IJRR*, 2020.
- [2] Chen Bao, Helin Xu, Yuzhe Qin, and Xiaolong Wang. Dexart: Benchmarking generalizable dexterous manipulation with articulated objects. In *CVPR*, 2023.
- [3] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. $\pi 0$: A vision-language-action flow model for general robot control. *arXiv*, 2024.
- [4] Samarth Brahmbhatt, Cusuh Ham, Charles C Kemp, and James Hays. Contactdb: Analyzing and predicting grasp contact via thermal imaging. In *CVPR*, 2019.
- [5] Samarth Brahmbhatt, Ankur Handa, James Hays, and Dieter Fox. Contactgrasp: Functional multi-finger grasp synthesis from contact. In *IROS*, 2019.
- [6] Samarth Brahmbhatt, Chengcheng Tang, Christopher D Twigg, Charles C Kemp, and James Hays. Contactpose: A dataset of grasps with object contact and hand pose. In *ECCV*, 2020.
- [7] Yu-Wei Chao, Wei Yang, Yu Xiang, Pavlo Molchanov, Ankur Handa, Jonathan Tremblay, Yashraj S Narang, Karl Van Wyk, Umar Iqbal, Stan Birchfield, et al. Dexycb: A benchmark for capturing hand grasping of objects. In *CVPR*, 2021.

- [8] Jiayi Chen, Yuxing Chen, Jialiang Zhang, and He Wang. Task-oriented dexterous grasp synthesis via differentiable grasp wrench boundary estimator. *IROS*, 2024.
- [9] Yuanpei Chen, Chen Wang, Yaodong Yang, and Karen Liu. Object-centric dexterous manipulation from human motion data. In *CoRL*, 2024.
- [10] Zoey Qiuyu Chen, Karl Van Wyk, Yu-Wei Chao, Wei Yang, Arsalan Mousavian, Abhishek Gupta, and Dieter Fox. Dextranet: Real world multi-fingered dexterous grasping with minimal human demonstrations. *RSS IL workshop*, 2022.
- [11] Cheng Chi, Zhenjia Xu, Chuer Pan, Eric Cousineau, Benjamin Burchfiel, Siyuan Feng, Russ Tedrake, and Shuran Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots. In *RSS*, 2024.
- [12] Sammy Christen, Lan Feng, Wei Yang, Yu-Wei Chao, Otmar Hilliges, and Jie Song. Synh2r: Synthesizing hand-object motions for learning human-to-robot handovers. In *ICRA*, 2024.
- [13] Shadow Robot Company. Shadow hand, 2005. URL <https://www.shadowrobot.com/dexterous-hand-series/>.
- [14] Enric Corona, Albert Pumarola, Guillem Alenya, Francesc Moreno-Noguer, and Grégory Rogez. Gan-hand: Predicting human grasp affordances in multi-object scenes. In *CVPR*, 2020.
- [15] Hongkai Dai, Anirudha Majumdar, and Russ Tedrake. Synthesis and optimization of force closure grasps via sequential semidefinite programming. *ISRR*, 2018.
- [16] Zicong Fan, Omid Taheri, Dimitrios Tzionas, Muhammed Kocabas, Manuel Kaufmann, Michael J. Black, and Otmar Hilliges. ARCTIC: A dataset for dexterous bimanual hand-object manipulation. In *CVPR*, 2023.
- [17] Carlo Ferrari, John F Canny, et al. Planning optimal grasps. In *ICRA*, 1992.
- [18] Zipeng Fu, Tony Z. Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. In *CoRL*, 2024.
- [19] Shreyas Hampali, Mahdi Rad, Markus Oberweger, and Vincent Lepetit. Honnotate: A method for 3d annotation of hand and object poses. In *CVPR*, 2020.
- [20] Yana Hasson, Gul Varol, Dimitrios Tzionas, Igor Kalevatykh, Michael J Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *CVPR*, 2019.
- [21] Hanwen Jiang, Shaowei Liu, Jiashun Wang, and Xiaolong Wang. Hand-object contact consistency reasoning for human grasps generation. In *ICCV*, 2021.
- [22] Hanwen Jiang, Shaowei Liu, Jiashun Wang, and Xiaolong Wang. Hand-object contact consistency reasoning for human grasps generation. In *ICCV*, 2021.
- [23] Min Liu, Zherong Pan, Kai Xu, Kanishka Ganguly, and Dinesh Manocha. Deep differentiable grasp planner for high-dof grippers. In *RSS*, 2020.
- [24] Shaowei Liu, Yang Zhou, Jimei Yang, Saurabh Gupta, and Shenlong Wang. Contactgen: Generative contact modeling for grasp generation. In *ICCV*, 2023.
- [25] Tengyu Liu, Zeyu Liu, Ziyuan Jiao, Yixin Zhu, and Song-Chun Zhu. Synthesizing diverse and physically stable grasps with arbitrary hand structures using differentiable force closure estimator. *RA-L*, 2021.
- [26] Yumeng Liu, Yaxun Yang, Youzhuo Wang, Xiaofei Wu, Jiamin Wang, Yichen Yao, Sören Schwertfeger, Sibe Yang, Wenping Wang, Jingyi Yu, Xuming He, and Yuexin Ma. Realdex: Towards human-like grasping for robotic dexterous hand. In *IJCAI*, 2024.
- [27] Jiabin Lu, Hao Kang, Haoxiang Li, Bo Liu, Yiding Yang, Qixing Huang, and Gang Hua. Ugg: Unified generative grasping. In *ECCV*, 2024.
- [28] Qingkai Lu, Mark Van der Merwe, and Tucker Hermans. Multi-fingered active grasp learning. In *IROS*, 2020.
- [29] Tyler Ga Wei Lum, Martin Matak, Viktor Makoviychuk, Ankur Handa, Arthur Allshire, Tucker Hermans, Nathan D. Ratliff, and Karl Van Wyk. Dextrah-g: Pixels-to-action dexterous arm-hand grasping with geometric fabrics. In *CoRL*, 2024.
- [30] Jens Lundell, Francesco Verdoja, and Ville Kyrki. Ddgc: Generative deep dexterous grasping in clutter. *RA-L*, 2021.
- [31] Miles Macklin. Warp: A high-performance python framework for gpu simulation and graphics. <https://github.com/nvidia/warp>, March 2022.
- [32] Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin, David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance gpu-based physics simulation for robot learning. In *NeurIPS Datasets and Benchmarks*, 2021.
- [33] Anusha Nagabandi, Kurt Konolige, Sergey Levine, and Vikash Kumar. Deep dynamics models for learning dexterous manipulation. In *CoRL*, 2020.
- [34] Jean Ponce, Steve Sullivan, J-D Boissonnat, and J-P Merlet. On characterizing and computing three-and four-finger force-closure grasps of polyhedral objects. In *ICRA*, 1993.
- [35] Jean Ponce, Steve Sullivan, Attawith Sudsang, Jean-Daniel Boissonnat, and Jean-Pierre Merlet. On computing four-finger equilibrium and force-closure grasps of polyhedral objects. *IJRR*, 1997.
- [36] Domenico Prattichizzo, Monica Malvezzi, Marco Gabiccini, and Antonio Bicchi. On the manipulability ellipsoids of underactuated robotic hands with compliance. *RAS*, 2012.
- [37] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, 2017.
- [38] Yuzhe Qin, Yueh-Hua Wu, Shaowei Liu, Hanwen Jiang, Ruihan Yang, Yang Fu, and Xiaolong Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. In *ECCV*, 2022.
- [39] Alberto Rodriguez, Matthew T Mason, and Steve Ferry.

- From caging to grasping. *IJRR*, 2012.
- [40] Carlos Rosales, Raúl Suárez, Marco Gabiccini, and Antonio Bicchi. On the synthesis of feasible and prehensile robotic grasps. In *ICRA*, 2012.
 - [41] Lin Shao, Fabio Ferreira, Mikael Jorda, Varun Nambiar, Jianlan Luo, Eugen Solowjow, Juan Aparicio Ojea, Oussama Khatib, and Jeannette Bohg. Unigrasp: Learning a unified model to grasp with multifingered robotic hands. *RA-L*, 2020.
 - [42] Ritvik Singh, Arthur Allshire, Ankur Handa, Nathan Ratliff, and Karl Van Wyk. Dextrah-rgb: Visuomotor policies to grasp anything with dexterous hands. *arXiv preprint arXiv:2412.01791*, 2024.
 - [43] Balakumar Sundaralingam, Siva Kumar Sastry Hari, Adam Fishman, Caelan Garrett, Karl Van Wyk, Valts Blukis, Alexander Millane, Helen Oleynikova, Ankur Handa, Fabio Ramos, et al. Curobo: Parallelized collision-free robot motion generation. In *ICRA*, 2023.
 - [44] Stone Tao, Fanbo Xiang, Arth Shukla, Yuzhe Qin, Xander Hinrichsen, Xiaodi Yuan, Chen Bao, Xinsong Lin, Yulin Liu, Tse kai Chan, Yuan Gao, Xuanlin Li, Tongzhou Mu, Nan Xiao, Arnav Gurha, Zhiao Huang, Roberto Calandra, Rui Chen, Shan Luo, and Hao Su. Maniskill3: Gpu parallelized robotics simulation and rendering for generalizable embodied ai. *arXiv*, 2024.
 - [45] Dylan Turpin, Liquan Wang, Eric Heiden, Yun-Chun Chen, Miles Macklin, Stavros Tsogkas, Sven Dickinson, and Animesh Garg. Grasp’d: Differentiable contact-rich grasp synthesis for multi-fingered hands. In *ECCV*, 2022.
 - [46] Dylan Turpin, Tao Zhong, Shutong Zhang, Guanglei Zhu, Eric Heiden, Miles Macklin, Stavros Tsogkas, Sven Dickinson, and Animesh Garg. Fast-grasp’d: Dexterous multi-finger grasp generation through differentiable simulation. In *ICRA*, 2023.
 - [47] Ruicheng Wang, Jialiang Zhang, Jiayi Chen, Yinzhen Xu, Puhao Li, Tengyu Liu, and He Wang. Dexgraspnet: A large-scale robotic dexterous grasp dataset for general objects based on simulation. In *ICRA*, 2023.
 - [48] Zehang Weng, Haofei Lu, Danica Kragic, and Jens Lundell. Dexdiffuser: Generating dexterous grasps with diffusion models. *RA-L*, 2024.
 - [49] Albert Wu, Michelle Guo, and C Karen Liu. Learning diverse and physically feasible dexterous grasps with generative model and bilevel optimization. *CoRL*, 2022.
 - [50] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, et al. Sapien: A simulated part-based interactive environment. In *CVPR*, 2020.
 - [51] Yinzhen Xu, Weikang Wan, Jialiang Zhang, Haoran Liu, Zikang Shan, Hao Shen, Ruicheng Wang, Haoran Geng, Yijia Weng, Jiayi Chen, et al. Unidexgrasp: Universal robotic dexterous grasping via learning diverse proposal generation and goal-conditioned policy. In *CVPR*, 2023.
 - [52] Lixin Yang, Xinyu Zhan, Kailin Li, Wenqiang Xu, Jiefeng Li, and Cewu Lu. Cpf: Learning a contact potential field to model the hand-object interaction. In *ICCV*, 2021.
 - [53] Jianglong Ye, Jiashun Wang, Binghao Huang, Yuzhe Qin, and Xiaolong Wang. Learning continuous grasping function with a dexterous hand from human demonstrations. *RA-L*, 2023.
 - [54] Jialiang Zhang, Haoran Liu, Danshi Li, XinQiang Yu, Haoran Geng, Yufei Ding, Jiayi Chen, and He Wang. Dexgraspnet 2.0: Learning generative dexterous grasping in large-scale synthetic cluttered scenes. In *CoRL*, 2024.
 - [55] Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. In Kostas E. Bekris, Kris Hauser, Sylvia L. Herbert, and Jingjin Yu, editors, *RSS*, 2023.
 - [56] Tony Z Zhao, Jonathan Tompson, Danny Driess, Pete Florence, Kamyar Ghasemipour, Chelsea Finn, and Ayzaan Wahid. Aloha unleashed: A simple recipe for robot dexterity. In *CoRL*, 2024.
 - [57] Yi Zhao, Le Chen, Jan Schneider, Quankai Gao, Juho Kannala, Bernhard Schölkopf, Joni Pajarinen, and Dieter Buehler. Rplm: A large-scale motion dataset for piano playing with bi-manual dexterous robot hands. *arXiv preprint arXiv:2408.11048*, 2024.
 - [58] Yu Zheng and Chee-Meng Chew. Distance between a point and a convex cone in n -dimensional space: Computation and applications. *T-RO*, 2009.

Dex1B: Learning with 1B Demonstrations for Dexterous Manipulation

Appendix

I. GRASPING SYNTHESIS EVALUATION DETAILS

We evaluate the performance of the proposed generative model, DexSimple, using the DexGraspNet [47] benchmark. This benchmark includes 5355 objects, with each object associated with approximately 200 grasps across five scales: $\{0.06, 0.08, 0.1, 0.12, 0.15\}$. We follow the train/test splits provided by the benchmark, which divide it into 4229 objects for training and 1126 objects for testing. We train DexSimple only using the training data provided by the benchmark. The results of DDG [23], GraspTTA [21], UDG [51] and UGG [27] are taken from the UGG [27] paper.

In the evaluation, we adhere to the metrics established in DexGraspNet [47], which are divided into two categories: Quality (Success Rate, Q_1 -score, Penetration) and Diversity (H mean and H std). We detail each metric here.

Quality Metrics. i). *Success rate (%)*: A grasp is considered successful if it withstands at least one of the six gravity directions in the Isaac Gym simulator [32] and maintains a maximal penetration depth of less than 0.5 cm. ii). *Q_1 -score*: Q_1 [17] is the radius of the inscribed sphere of the ConvexHull ($\cup_i w_i$), indicating the norm of the smallest wrench that can destabilize the grasp. The contact threshold is set to 1 cm. iii). *Maximal Penetration Depth (cm)*: This metric measures the penetration depth from the object point cloud to the hand meshes.

Diversity Metrics. *Diversity* in the DexGraspNet benchmark is evaluated using joint angle entropy. The range of joint motion is divided into 10000 bins², and all generated samples are used to estimate a probability distribution. The entropy is calculated based on this distribution. The reported values in main text refer to the mean and standard deviation across all joints (H mean and H std).

II. BENCHMARK DETAILS

The proposed Dex1B dataset contains 3,491 objects for training and 933 objects for testing in the lifting task, as well as 43 articulation objects for training and 21 articulation objects for testing in the articulation task.

The lifting and articulation tasks are conducted in the ManiSkill [44]/SAPIEN [50] simulation environments. The physical parameters are listed in Tab. V.

Lifting Task Definition. The hand begins from a randomly sampled pose and joint configuration. The goal of the lifting task is to reach an object placed on a plane (table), grasp it, and lift it to a height of 0.4 m. Additionally, the task requires at least two fingers to maintain contact with the object in the final frame. We show example trajectories from Dex1B for the lifting task in Fig. 10.

²The value 10000 is taken from the UGG implementation, and the entropy results obtained using this value are consistent with those reported in [47, 27]. However, the descriptions provided in [47, 27] are inaccurate.

Parameter	Value
Object Mass	0.1 kg
Simulation Frequency	100
Control Frequency	25
Contact Offset	0.001
Solver Position Iterations	30

TABLE V: Physical parameters for the simulation.

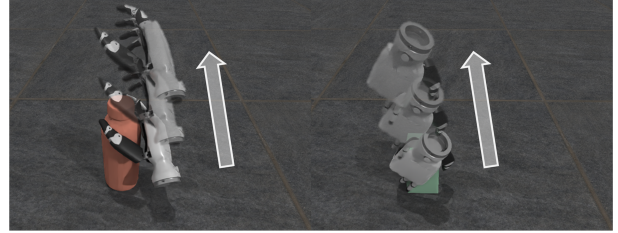


Fig. 10: Lifting trajectory from Dex1B dataset.

Articulation Task Definition. The hand begins from a randomly sampled pose and joint configuration. The goal of the articulation task is to reach the interactable link of an articulated structure (e.g., the top link of a laptop) and open it to increase its joint angle by 0.5 radians. Additionally, the task requires at least two fingers to maintain contact with the object in the final frame. We show example trajectories from Dex1B for the lifting task in Fig. 11.

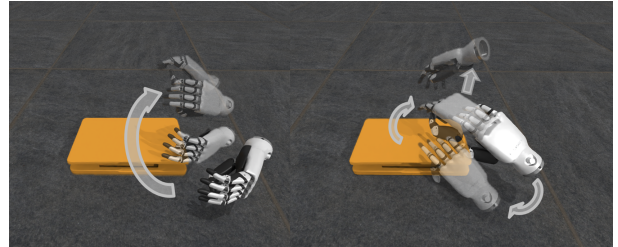


Fig. 11: Articulation trajectory from Dex1B dataset.

III. ITERATIVE DATA ENGINE DETAILS

The Iterative Data Engine begins with a controllable optimization algorithm for hand pose and trajectory optimization. Both the grasping synthesis and motion planning can be formulated as optimization problems [25, 47, 43]. We first describe the pure optimization-based data generation process for both grasping and motion planning.

A. Grasping Synthesis as Optimizaiton

Differentiable Force Closure. Following [25, 47], we adopt differentiable force closure estimator as an energy term for grasping optimization. This term encourages a set of contact points to form force closure, can be expressed as

$$E_{fc} = \|Gc\|_2 \quad (2)$$

where

$$G = \begin{bmatrix} I_3 & \cdots & I_3 \\ [x_1]_{\times} & \cdots & [x_n]_{\times} \end{bmatrix}$$

$$[x_i]_{\times} = \begin{bmatrix} 0 & -x_i^{(z)} & x_i^{(y)} \\ x_i^{(z)} & 0 & -x_i^{(x)} \\ -x_i^{(y)} & x_i^{(x)} & 0 \end{bmatrix}$$

Here x represents $n = 4$ contact points, $c \in \mathbb{R}^{n \times 3}$ represents the contact normal vectors, which can be computed given object mesh O and x . As in DexGraspNet, the contact points are randomly chosen from a predefined set of candidate contact points at each iteration.

SDF Energy. To prevent penetration, [47, 21] represent objects as point clouds and compute the penetration distance between object point clouds and hand link meshes. In contrast, we represent objects as meshes and hand links as spheres, calculate the sphere-mesh signed distance function (SDF), and leverage a BVH structure to accelerate computations.

$$E_{\text{sdf}} = \max \left(-\min_s (\text{SDF}(s_{\text{center}}) + s_{\text{radius}}), 0 \right) \quad (3)$$

Here, s denotes the sphere, $s_{\text{center}} \in \mathbb{R}^{n \times 3}$ is the 3D position of the sphere center, and $s_{\text{radius}} \in \mathbb{R}^+$ is the sphere radius. $\text{SDF}(s_{\text{center}})$ is the SDF query operation (inside is negative). We implement this term using warp-lang [31]. The resulting SDF term is far more accurate and efficient.

Other Energy Terms. We follow DexGraspNet [47] to add contact distance (E_D), self-penetration energy (E_S), and joint limit energy (E_J):

$$E_D = \sum_i d(x_i, O),$$

$$E_S = \sum_{p \neq q} \max(\delta - d(p, q), 0), \quad (4)$$

$$E_J = \sum_i (\max(\theta_i - \theta_i^{\max}, 0) + \max(\theta_i^{\min} - \theta_i, 0))$$

Here, contact distance (E_D) is the distance between predefined contact points x_i and object mesh O . Self-penetration energy (E_S) is the threshold $\delta = 0.02$ minus the distance between joints p, q . Joint limit energy (E_J) is the deviation for joint angles θ from their joint limits.

The complete energy function for grasping synthesis is as follows:

$$E = E_{\text{fc}} + w_{\text{sdf}} E_{\text{sdf}} + w_D E_D + w_S E_S + w_J E_J, \quad (5)$$

The weights are set as follows: $w_{\text{sdf}} = 100.0$, $w_D = 100.0$, $w_S = 10.0$, $w_J = 1.0$.

B. Motion Planning as Optimizaiton

We formulate motion planning as an optimization problem incorporating SDF energy (E_{sdf}), self-penetration energy (E_S), joint limit energy (E_J), and smoothness energy (E_{smooth}). The first three energy terms are extensions of single-frame grasping formulations to multi-frame settings. The last smoothness energy ensures velocity smoothness:

$$E_{\text{smooth}} = \sum_{t=2}^T \|(g_t - g_{t-1}) / \Delta t\|^2 \quad (6)$$

Here, $g = (T, R, \theta)$ represents the grasp tuple with $T \in \mathbb{R}^3$ for global translation, $R \in SO(3)$ for global rotation, and $\theta \in \mathbb{R}^d$ for joint angles.

Continuous Euler Angles Optimization. Free hand is typically implemented using 6 root joints (3 prismatic joints for translation, 3 revolute joints for rotation) in the simulation. The 3 revolute joints are equivalent to intrinsic Euler angles (x-y-z). However, converting from other rotation representations (such as 6D rotation) to intrinsic Euler angles is not necessarily continuous across time steps. To address this, we use an optimization-based approach for rotation conversion. Specifically, we use a smoothness energy term (E_{smooth}) and a rotation difference term, which indicates the rotation angle difference between current rotation and target rotation, to optimize the revolute joint angles over a trajectory.

C. Task-specific Trajectory Optimization

Both the lifting and articulation demonstrations can be divided into three stages: pre-grasping, grasping, and post-grasping. The grasping stage involves grasp synthesis, while the pre-grasping stage focuses on motion planning. In the post-grasping stage, lifting requires raising the hand's height (z-value), whereas articulation requires rotating the hand along the axis of the articulated structure. For all demonstrations, we follow this data generation process: i) Grasp synthesis; ii) Sampling the starting pose and perform motion planning; iii) Generating the task-specific post-grasping trajectory.

IV. DEXSIMPLE MODEL DETAILS

While generative models [21, 24, 27] have been extensively studied for dexterous manipulation, we introduce several differentiable loss terms inspired by optimization-based methods, SDF loss, distance loss, and smoothness loss, and demonstrate that these losses significantly enhance performance. For the SDF loss, due to the lack of cache support for BVH structures in warp-lang, we use point-sphere SDF queries instead of mesh-sphere SDF queries during large-scale training.

The proposed DexSimple is a CVAE architecture with a PointNet encoder. The architecture details and training hyperparameters for DexSimple are provided in Tab. VI.

Parameter	Value
Num. Points	1024
PointNet Layer Sizes	(3, 64, 128, 1024, 256)
CVAE In/Out Dim.	$N_{\text{frames}} \times N_{\text{DOF}}$
CVAE Layer Sizes	(256, 512, 256)
Learning rate	1e-5
MSE Loss Weight	1.0
KL Loss Weight	1e-4
SDF Loss Weight	1e-4
Distance Loss Weight	1e-4
Smoothness Loss Weight	1e-5
Optimizer	Adam

TABLE VI: Parameters for DexSimple.