

HGDiffuser: Efficient Task-Oriented Grasp Generation via Human-Guided Grasp Diffusion Models

Dehao Huang^{1,2}, Wenlong Dong^{1,2}, Chao Tang^{1,2}, Hong Zhang^{1,2} *Life Fellow, IEEE*

Abstract—Task-oriented grasping (TOG) is essential for robots to perform manipulation tasks, requiring grasps that are both stable and compliant with task-specific constraints. Humans naturally grasp objects in a task-oriented manner to facilitate subsequent manipulation tasks. By leveraging human grasp demonstrations, current methods can generate high-quality robotic parallel-jaw task-oriented grasps for diverse objects and tasks. However, they still encounter challenges in maintaining grasp stability and sampling efficiency. These methods typically rely on a two-stage process: first performing exhaustive task-agnostic grasp sampling in the 6-DoF space, then applying demonstration-induced constraints (e.g., contact regions and wrist orientations) to filter candidates. This leads to inefficiency and potential failure due to the vast sampling space. To address this, we propose the Human-guided Grasp Diffuser (HGDiffuser), a diffusion-based framework that integrates these constraints into a guided sampling process. Through this approach, HGDiffuser directly generates 6-DoF task-oriented grasps in a single stage, eliminating exhaustive task-agnostic sampling. Furthermore, by incorporating Diffusion Transformer (DiT) blocks as the feature backbone, HGDiffuser improves grasp generation quality compared to MLP-based methods. Experimental results demonstrate that our approach significantly improves the efficiency of task-oriented grasp generation, enabling more effective transfer of human grasping strategies to robotic systems. To access the source code and supplementary videos, visit <https://sites.google.com/view/hgdiffuser>.

I. INTRODUCTION

Task-oriented grasping (TOG) refers to grasping objects in a manner that is aligned with the intended task, which is the first and crucial step for robots to perform manipulation tasks [1]. For instance, when handing over a kitchen knife to a human, the blade should be grasped perpendicular to the handle to ensure a safe and efficient handover. Many objects in daily life are designed with human convenience in mind, so human task-oriented grasping inherently includes the skills necessary for manipulating the objects, such as maintaining stability and avoiding collisions with the environment during manipulation. Consequently, previous methods have proposed to utilize human grasp demonstrations as training data or reference templates for robotic task-oriented grasp generation. In this work, we specify an important form of this problem and focus on transferring human grasps from demonstrations to robotic 6-DoF parallel-jaw task-oriented

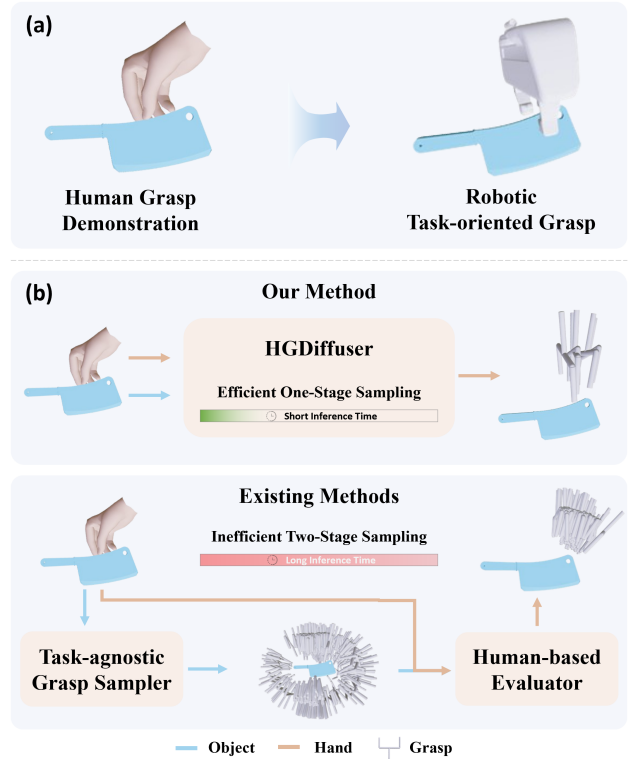


Fig. 1. (a) Demonstration-based methods, which generate robotic 6-DoF parallel-jaw task-oriented grasps by leveraging human demonstrations. (b) Comparison of existing two-stage methods and our single-stage method. Unlike two-stage methods, which require extensive sampling followed by filtering to generate grasps, our method directly generates grasps with minimal sampling, making it more efficient.

grasps due to the popularity of this type of robot end-effector, as illustrated in Figure 1(a).

To transfer human grasps to robotic 6-DoF parallel-jaw task-oriented grasps, recent works [2]–[5] have proposed a two-stage approach for task-oriented grasping. First, task-agnostic grasp sampling generates stable parallel-jaw grasp candidates. Subsequently, explicit task-oriented constraints derived from human demonstrations are applied to filter these candidates, ensuring both stability and task-specific requirements are met. This two-stage approach, leveraging a diverse set of stable grasp candidates, often succeeds in identifying high-quality task-oriented grasps. However, in the vast 6-DoF grasp sampling space, task-agnostic samplers require extensive sampling to generate sufficiently diverse candidates that satisfy both stability and task-oriented requirements, leading to inefficiency and a potential risk of failure.

To address this challenge, we propose the Human-guided Grasp Diffuser (HGDiffuser), a diffusion-based framework that leverages human task-oriented grasp demonstrations to

¹Shenzhen Key Laboratory of Robotics and Computer Vision, Southern University of Science and Technology, Shenzhen, China.

²Department of Electronic and Electrical Engineering, Southern University of Science and Technology, Shenzhen, China.

This work was supported in part by the Shenzhen Key Laboratory of Robotics and Computer Vision (ZDSYS20220330160557001).

directly generate 6-DoF parallel-jaw task-oriented grasps in a single sampling stage. Unlike conventional two-stage methods requiring extensive task-agnostic sampling and subsequent constraint-based filtering, HGDiffuser integrates explicit task-oriented constraints directly into the sampling process, significantly enhancing efficiency, as shown in Figure 1(b). Specifically, HGDiffuser utilizes the inherent guided sampling mechanism of diffusion models [6]–[8], to incorporate explicit task-oriented constraints extracted from human demonstrations—including hand-object contact regions and wrist orientations—as an additional score function to guide grasp sampling toward stable, task-compliant solutions within the 6-DoF parallel-jaw grasp manifold. By eliminating the need for exhaustive random sampling, HGDiffuser enhances efficiency and enables the direct and effective transfer of human task-oriented grasping strategies to robotic parallel-jaw grippers. Furthermore, building upon recent advancements in diffusion models, we implement Diffusion Transformer (DiT) blocks [9] as the feature backbone of HGDiffuser. This architectural choice leverages the attention mechanism inherent in transformers, enabling more effective feature fusion compared to traditional UNet-based or MLP-based backbones. The incorporation of DiT blocks improves the feature fusion of HGDiffuser. The experimental results demonstrate that HGDiffuser achieves a remarkable 81.26% reduction (from 1.019s to 0.191s) in inference time compared to the state-of-the-art (SOTA) two-stage method while maintaining competitive task-oriented grasp generation quality.

In summary, the contributions of this paper are outlined as follows.

- We present HGDiffuser, a novel diffusion-based framework that leverages the guided sampling mechanism to incorporate explicit task-oriented grasp constraints derived from human demonstrations. This approach eliminates the exhaustive task-agnostic sampling of conventional methods, significantly improving efficiency while maintaining grasp quality.
- We propose the integration of Diffusion Transformer (DiT) blocks as the core architectural component of HGDiffuser. The attention mechanism in DiT blocks enables superior feature fusion, enhancing the generated grasps’ quality.

II. RELATED WORK

Recent research on task-oriented grasp generation for parallel-jaw grippers can be categorized into three types based on the nature of the required data. Below, we provide a detailed discussion of each category.

Methods based on human demonstration data. This approach [2]–[5], [10], [11] leverages inexpensive human grasp demonstrations to generate task-oriented grasps. Human demonstrations may include static images [10], videos of identical objects [4], [12], or videos of similar objects [2], [3] for novel object-task pairs.

A core component of these methods is the transformation module that converts human grasps into parallel-jaw grasps.

Early end-to-end solutions [10], [12] directly map human grasps using manual rules or MLP networks trained on small datasets. For instance, DemoGrasp [12] assumes that humans grasp objects in a fixed manner, using the midpoint between the thumb and the index finger as the parallel-jaw grasp point and the wrist orientation as the grasp direction. Patten et al. [10] annotated a small-scale dataset of human grasps and corresponding parallel-jaw grasps, to learn a mapping using an MLP network. However, these struggle with human grasp diversity [13] and complex mapping relationships, often yielding unstable grasps.

Recent two-stage methods [2]–[5] initially generate grasp candidates through a task-agnostic parallel-jaw grasp sampler [14], then apply task-oriented constraints derived from human grasp demonstrations to select optimal grasps. While Robo-ABC [3] and FUNCTO [15] use region constraints, RTAGrasp [2] and DITTO [4] combine region and orientation constraints. This two-stage approach, leveraging a diverse set of stable grasp candidates, often identifies high-quality task-oriented grasps that satisfy both stability and task-oriented requirements. However, in the vast 6-DoF parallel-jaw grasp space, task-agnostic samplers require extensive sampling to generate sufficiently diverse candidates that satisfy both stability and task-oriented requirements, leading to inefficiency. In contrast, Our HGDiffuser addresses this by leveraging the inherent guided sampling mechanism of diffusion models to incorporate task-oriented constraints into a diffusion-based parallel-jaw task-agnostic sampler directly. This innovative single-stage approach significantly improves efficiency while reducing the risk of sampling failure.

Methods based on manually annotated task-oriented grasp data. The high cost of manual annotation and generalization limitations remain key challenges for these approaches. Murali et al. [16] introduce TaskGrasp, a manually curated large-scale TOG dataset, alongside GC-NGrasp, which improves novel object generalization using semantic knowledge from pre-built knowledge graphs. Later works like GraspGPT [17], [18] extend capabilities by integrating open-ended semantic/geometric knowledge through LLM/VLM interactions, enabling generalization to novel object-task categories beyond training data. However, these methods remain constrained to objects/tasks resembling those in costly manual annotated datasets [18].

Methods based on large-scale internet data. Benefiting from Vision-Language Models (VLMs) pre-trained on massive web data, recent studies [19]–[21] have explored zero-shot task-oriented grasp generation capabilities. However, these methods face two fundamental limitations: (1) the absence of fine-grained object understanding data essential for task-oriented grasping in pre-training datasets, and (2) their capability being limited to predicting task-related grasp regions. These constraints jointly lead to suboptimal performance.

III. PROBLEM FORMULATION

Given a single-view RGB-D image of a human demonstration I_{demo} , which shows a person naturally grasping a target

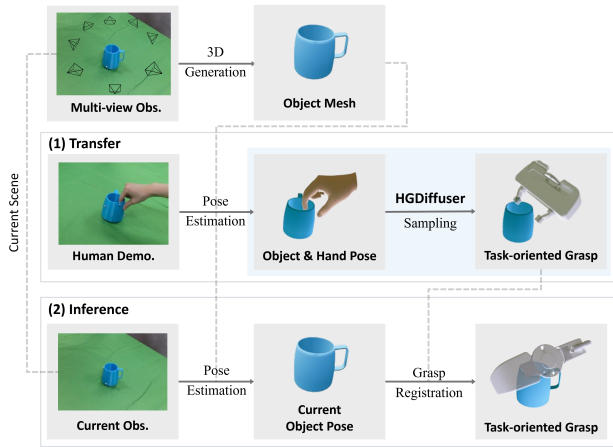


Fig. 2. Overview of our task-oriented grasping system. The task demonstrated is to handover a cup.

object o in a task-oriented manner, our goal is to reproduce the same task-oriented grasp on the same target object o using a parallel-jaw robotic gripper. The demonstration image \mathbf{I}_{demo} can be captured by the robot’s camera.

Figure 2 illustrates an overview of our task-oriented grasping system. The system comprises two phases: a transfer phase and an inference phase. Before these phases, the robot performs multi-view observations of the target object o and employs a 3D generation foundation model [22], [23] to reconstruct the object’s mesh \mathbf{M}_o , as shown in the upper part of Figure 2.

In the transfer phase, we use \mathbf{I}_{demo} and \mathbf{M}_o as inputs, employing vision foundation models for object pose estimation [24] and hand pose estimation [25]. This process yields the object point cloud $\mathbf{X}_o \in \mathbb{R}^{N \times 3}$ under the current object pose and the MANO parametric model [26] representing the human grasp $\mathbf{X}_h = \{\theta, \beta\}$, where $\theta \in \mathbb{R}^{48}$ and $\beta \in \mathbb{R}^{10}$. The task-oriented human grasp is then transformed to generate the corresponding parallel-jaw gripper grasp \mathbf{H} . During the inference phase, we estimate the current pose of target object o and register the transformed grasp \mathbf{H} .

It is crucial to emphasize that our core contribution is the development of HGDiffuser for the transformation module, while other system components are implemented using established techniques.

IV. HGDIFFUSER

An overview of the proposed HGDiffuser framework is shown in Figure 3. Given an object point cloud \mathbf{X}_o and a human grasp \mathbf{X}_h , HGDiffuser infers the corresponding task-oriented grasp \mathbf{H} for a parallel-jaw gripper. Formally, HGDiffuser learns a conditional distribution $\rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h)$, where $\mathbf{H} \in SE(3)$ represents a valid task-oriented grasp that aligns with the human demonstration.

While prior research [27]–[29] has explored VAE and diffusion models for task-agnostic grasp generation (i.e., learning $\rho(\mathbf{H} | \mathbf{X}_o)$), they rely on extensive simulated datasets for generalization. However, existing task-oriented grasping datasets lack the scale needed for similar end-to-end training. To address this problem, our approach leverages the guidance mechanism of diffusion models [8]. We train

a diffusion-based generative model on task-agnostic data to learn $\rho(\mathbf{H} | \mathbf{X}_o)$. During inference, we employ the guided sampling mechanism to incorporate the explicit task-oriented constraints $\rho(\mathbf{X}_h | \mathbf{H})$ derived from human grasp demonstrations, effectively sampling from $\rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h)$. This design allows us to generate task-oriented grasps that are stable and aligned with human demonstrations without requiring large-scale task-oriented training data.

A. DiT-based Diffusion Model for $\rho(\mathbf{H} | \mathbf{X}_o)$

For the first part of HGDiffuser, we build on the prior diffusion-based task-agnostic grasp sampler [27] and introduce the DiT blocks as the feature backbone.

Training procedure. The Denoising Score Matching (DSM) [30] is employed as the training procedure. Given a grasp dataset with $\{\mathbf{H}, \mathbf{X}_o\}$, the goal is to use a Noise Conditional Score Network (NCSN) s_θ to learn $\rho(\mathbf{H} | \mathbf{X}_o)$. The NCSN $s_\theta(\mathbf{H}, k, \mathbf{X}_o)$ is trained to estimate $\nabla_{\mathbf{H}} \log \rho_{\sigma_k}(\mathbf{H} | \mathbf{X}_o)$, where $k \in \{0, \dots, L-1\}$ denotes a noise scale among levels σ_k . The training objective minimizes the following loss function \mathcal{L}_{dsm} , using a score matching method [31]:

$$\mathcal{L}_{dsm} = \sum_{k=1}^L \mathbb{E}_{\rho_{\sigma_k}(\mathbf{H} | \mathbf{X}_o)} [\|\nabla_{\mathbf{H}} \log \rho_{\sigma_k}(\mathbf{H} | \mathbf{X}_o) - s_\theta(\mathbf{H}, k, \mathbf{X}_o)\|_2^2]$$

Inference procedure. The trained NCSN generates samples based on annealed Langevin Markov chain Monte Carlo (MCMC) [31]. A sample \mathbf{H}_L is drawn from $\mathcal{N}(\mathbf{0}, \mathbf{I})$, followed by L -step iterations (from $k = L-1$ to $k = 0$):

$$\begin{aligned} \mathbf{H}_{k-1} &= \mathbf{H}_k + \epsilon_k \nabla_{\mathbf{H}} \log \rho_{\sigma_k}(\mathbf{H} | \mathbf{X}_o) + \sqrt{2\epsilon_k} \mathbf{z} \\ &= \mathbf{H}_k + \epsilon_k s_\theta(\mathbf{H}, k, \mathbf{X}_o) + \sqrt{2\epsilon_k} \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \end{aligned}$$

where ϵ_k is a step-dependent coefficient that decreases as k decreases. The inference process repeats the central part of Figure 3 for L steps.

Network architecture. The NCSN $s_\theta(\mathbf{H}, k, \mathbf{X}_o)$ takes as input the object point cloud \mathbf{X}_o , the current noise step k , and the current grasp \mathbf{H}_k , and outputs the score $\mathbf{H}_s \in SE(3)$. The object point cloud \mathbf{X}_o is encoded using VN-PointNet [32], a $SO(3)$ -equivariant point cloud feature encoder, producing the feature vector $\mathbf{f}^o \in \mathbb{R}^d$, where d denotes the descriptor dimension. For the current grasp $\mathbf{H}_k \in SE(3)$, a predefined gripper points mapper first transforms it into gripper points $\mathbf{X}_g \in \mathbb{R}^{g \times 3}$, as illustrated in Figure 3, with g being the predefined number of points. Encoding the grasp via gripper points rather than directly from $SE(3)$ allows for more effective integration with the object point cloud features, as both are represented in Cartesian space. The gripper points \mathbf{X}_g are then processed by an MLP to generate the feature vector $\mathbf{f}^g \in \mathbb{R}^{g \times d}$. Finally, the noise step k is encoded using transformer sinusoidal position embedding [9], resulting in the feature vector $\mathbf{f}^t \in \mathbb{R}^d$.

To effectively fuse the features $\{\mathbf{f}^g, \mathbf{f}^o, \mathbf{f}^t\}$, we introduce a DiT-based feature backbone. At its core lies a transformer block with an attention mechanism, which has been extensively validated to achieve superior feature fusion capabilities across most tasks compared to other alternatives. Our

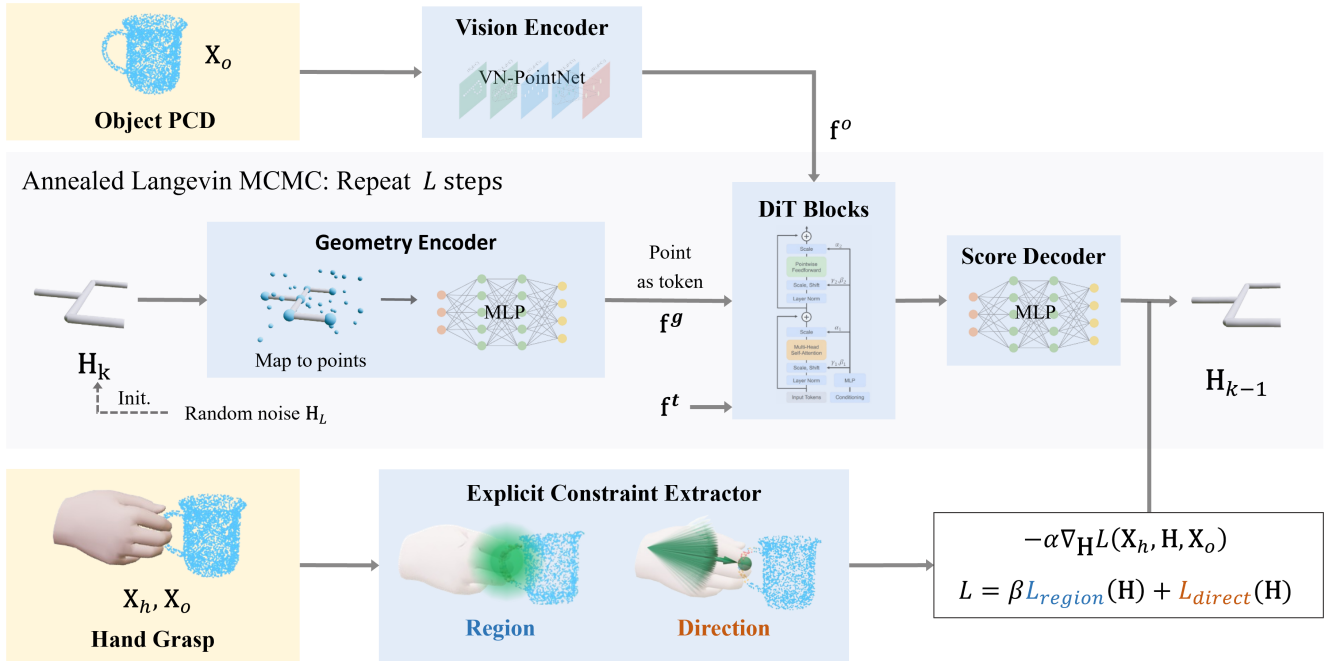


Fig. 3. An overview of HGDiffuser. The grasp generation employs annealed Langevin MCMC sampling with T steps. The input object point cloud \mathbf{X}_o is encoded into feature \mathbf{f}^o via vision encoder, while current grasp \mathbf{H}_k is processed into \mathbf{f}^g via geometry encoder. These features, along with step feature \mathbf{f}^t from sinusoidal encoding, serve as inputs to the DiT-based backbone. The fused features are decoded to produce a noise conditional score. For the input human grasp \mathbf{X}_h , explicit task-oriented constraints are extracted to construct a loss function guiding the sampling process. The noise conditional score, combined with the loss function, updates grasp \mathbf{H}_k to \mathbf{H}_{k-1} , iterating L times to output final grasp \mathbf{H}_0 .

DiT-based feature backbone comprises D DiT blocks [9] connected in series, where D is a configurable parameter adjustable based on data volume and model size. The DiT block, a variant of standard transformer block, takes input tokens and condition features as inputs and outputs fused feature tokens. It incorporates the condition feature input into the block’s adaptive layer normalization, a method demonstrated to be more efficient and effective than direct fusion using cross-attention.

Since our diffusion model learns the distribution of grasps \mathbf{H} , denoted as $\rho(\mathbf{H} | \mathbf{X}_o)$, the gripper points feature \mathbf{f}^g serves as the input tokens, while the fused object point cloud feature \mathbf{f}^o and noise step feature \mathbf{f}^t are utilized as the condition feature input. For the gripper points feature $\mathbf{f}^g \in \mathbb{R}^{g \times d}$, we propose a method inspired by image transformers, where each image patch is treated as a token. Existing methods for handling point cloud features as transformer input tokens [33] typically rely on advanced point serialization strategies and serialized attention mechanisms to address the unordered nature and large quantity of point cloud data. In our approach, however, the predefined gripper points mapper directly converts grasp \mathbf{H} into a serialized set of gripper points, analogous to the serialization of pixels in images. Consequently, we treat each gripper point feature as an individual input token. Following established practices for processing noise step features, we directly sum the object point cloud feature \mathbf{f}^o and the noise step feature \mathbf{f}^t to form the condition feature input $\mathbf{f}^c = \mathbf{f}^t + \mathbf{f}^o$.

B. Guidance-based Inference for $\rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h)$

After obtaining the diffusion model that has learned $\rho(\mathbf{H} | \mathbf{X}_o)$, representing task-agnostic grasp generation, we leverage the guided sampling mechanism inherent in

diffusion models to incorporate explicit task-oriented constraints derived from human grasp demonstration \mathbf{X}_h . This enables the diffusion model to sample from the distribution $\rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h)$, corresponding to the desired task-oriented grasp generation.

Inference procedure. To sample from the distribution $\rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h)$, we follow the approach of the aforementioned diffusion model by transforming the problem into estimating the score function $\nabla_{\mathbf{H}} \log \rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h)$. Using Bayes’ theorem, the score function of the conditional distribution decomposes into the sum of the score functions of the prior distribution and the likelihood distribution:

$$\begin{aligned} \nabla_{\mathbf{H}} \log \rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h) = \\ \nabla_{\mathbf{H}} \log \rho(\mathbf{H} | \mathbf{X}_o) + \nabla_{\mathbf{H}} \log \rho(\mathbf{X}_h | \mathbf{H}, \mathbf{X}_o) \end{aligned}$$

Here, the prior distribution component $\nabla_{\mathbf{H}} \log \rho(\mathbf{H} | \mathbf{X}_o)$ corresponds to our trained NSCN $s_{\theta}(\mathbf{H}, k, \mathbf{X}_o)$. The likelihood distribution component $\nabla_{\mathbf{H}} \log \rho(\mathbf{X}_h | \mathbf{H}, \mathbf{X}_o)$ represents guidance derived from human grasp demonstration. In the context of class-conditional image generation tasks, this term is known as Classifier Guidance [8], typically estimated using a pre-trained image classifier. The image classifier’s understanding of image categories guides unconditionally generative models in generating images of specific categories. In our task, we leverage the task-relevant information in human grasps to guide the generation of task-oriented parallel-jaw grasps. Specifically, we introduce a non-network-based, differentiable loss function $L(\mathbf{X}_h, \mathbf{H}, \mathbf{X}_o)$, which evaluates the probability of a parallel-jaw grasp \mathbf{H} meeting task-oriented requirements based on the human grasp demonstration \mathbf{X}_h . The lower the loss, the higher the corresponding probability. The estimated score

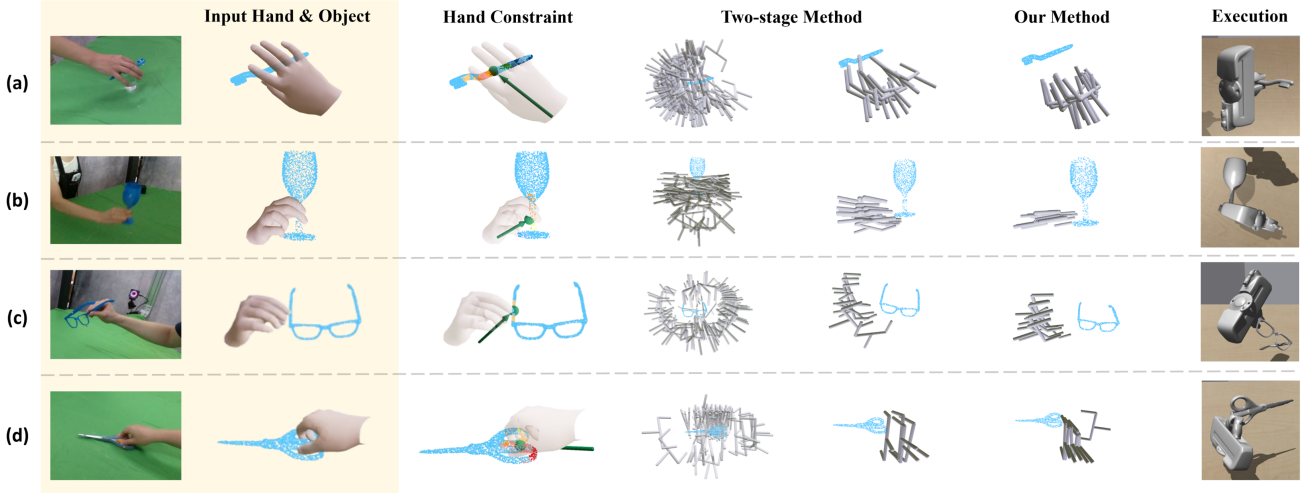


Fig. 4. Qualitative results of our method and Ours-TS method. The object categories and tasks are as follows: (a) toothbrush and brushing, (b) wine glass and pouring, (c) eyeglasses and handing over, (d) scissors and using. More results are provided in the supplementary material.

function $\nabla_{\mathbf{H}} \log \rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h)$ is expressed as:

$$\nabla_{\mathbf{H}} \log \rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h) = \mathbf{s}_{\theta}(\mathbf{H}, k, \mathbf{X}_o) - \alpha \nabla_{\mathbf{H}} L(\mathbf{X}_h, \mathbf{H}, \mathbf{X}_o)$$

where α is a scaling parameter. Sampling from the distribution $\rho(\mathbf{H} | \mathbf{X}_o, \mathbf{X}_h)$ using the annealed Langevin MCMC method follows the equation:

$$\mathbf{H}_{k-1} = \mathbf{H}_k + \epsilon_k [\mathbf{s}_{\theta}(\mathbf{H}, k, \mathbf{X}_o) - \alpha \nabla_{\mathbf{H}} L(\mathbf{X}_h, \mathbf{H}, \mathbf{X}_o)] + \sqrt{2\epsilon_k} \mathbf{z}$$

Explicit task-oriented constraint extraction. As illustrated in the lower part of Figure 3, we extract explicit task-oriented constraints from the human grasp demonstration to compute the loss function $L(\mathbf{X}_h, \mathbf{H}, \mathbf{X}_o)$, which serves as the guidance. To satisfy the task-oriented requirements of parallel-jaw grasps, two constraints are typically essential: the task-oriented grasp region constraint and the task-oriented grasp orientation constraint [2], [17]. For instance, in using a brush to clean a table, the parallel-jaw gripper should grasp the brush handle in an upright orientation away from the bristle end. This combination of grasp region and orientation ensures minimal collision with the environment (e.g., the table or debris) during the cleaning process. It reduces the force required to hold the brush.

These constraints are naturally embedded in the task-oriented human grasp, which inherently contains information about both aspects. We establish two corresponding loss functions to formalize these constraints: L_{region} and L_{direct} . Specifically, since the task-oriented grasp region of the human hand often overlaps significantly with that of the parallel-jaw gripper, we compute all contact points [34] from \mathbf{X}_h and \mathbf{X}_o , then derive the center point P_{region} . L_{region} is formulated as a RELU function based on the Euclidean distance to P_{region} , which remains zero when the distance is below a predefined threshold τ_{region} and increases linearly beyond it, as follows:

$$L_{region}(\mathbf{H}) = \text{RELU}(\|\text{Trans}(\mathbf{H}) - P_{region}\|_2 - \tau_{region})$$

Here, the Trans function extracts the translational component of \mathbf{H} , and τ_{region} is the predefined distance threshold. Regarding grasp orientation, prior research has demonstrated that human grasps often avoid collisions between the arm and the environment during task execution [35], [36]. To ensure

the parallel-jaw gripper’s relative position resembles that of the human hand, we align its orientation with the human wrist orientation. From the human hand \mathbf{X}_h , we compute the corresponding wrist orientation $D_{direct} \in SO(3)$ [26]. L_{direct} is formulated as a RELU function based on the geodesic distance to D_{direct} , as follows:

$$L_{direct}(\mathbf{H}) = \text{RELU}(\text{GeoDistance}(\text{Rot}(\mathbf{H}), D_{direct}) - \tau_{direct})$$

The Rot function extracts the rotational component of \mathbf{H} , GeoDistance computes the geodesic distance, and τ_{direct} is the predefined angular threshold. The final loss function is obtained as the weighted sum of L_{region} and L_{direct} :

$$L = \beta L_{region} + L_{direct}$$

where β is a scaling parameter.

V. EXPERIMENTAL RESULTS

In this section, we compare HGDdiffuser with two categories of existing methods for generating 6-DoF parallel-jaw task-oriented grasps from human demonstrations: (1) rule-based direct conversion methods and (2) two-stage frameworks involving candidate generation and task-constrained filtering. This comparison aims to evaluate the efficiency of HGDdiffuser. We also examine the impact of DiT blocks through an ablation study. Furthermore, we assess the practical applicability of HGDdiffuser through real-world experiments.

A. Quantitative Comparison

Baselines We compare HGDdiffuser to the following methods: (1) DemoGrasp [12], which generates parallel-jaw task-oriented grasps based on the direct conversion of human task-oriented grasps using manually designed rules. Specifically, the grasp direction of the parallel-jaws corresponds to the wrist orientation, while the grasp point is aligned with the center of the contact region between the index finger, thumb, and the object. (2) RTAGrasp [2], Robo-ABC [3] and DITTO [4], all of which employ a similar two-stage approach, represent state-of-the-art methods for generating parallel-jaw task-oriented grasps from human demonstrations. The first stage utilizes Contact-GraspNet [14] as a task-agnostic grasp

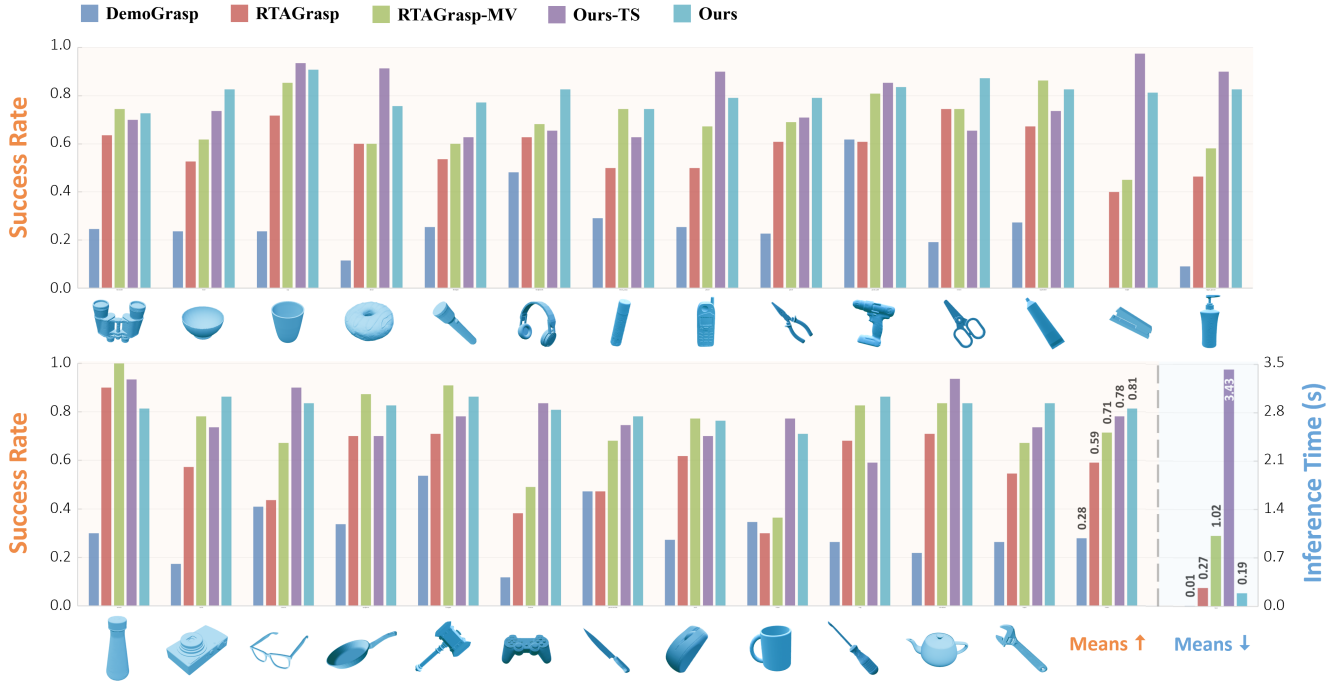


Fig. 5. Quantitative results. In the bottom-right section, we compare our method with baseline approaches in terms of average task-oriented grasping success rate and average inference time. The remaining sections present the average success rates across 24 object categories (out of 236 total object instances) from the dataset.

sampler to generate candidate grasps from the object’s point cloud. In the second stage, these candidates are filtered based on explicit task-oriented constraints derived from human demonstrations, mainly considering wrist orientation and hand-object contact points. In our results, we collectively refer to these three methods as RTAGrasp.

Since Contact-GraspNet operates on partially observed point clouds, these methods rely on partial observations as input. In contrast, HGDdiffuser utilizes the complete object point cloud. To ensure a fair comparison, we aggregate grasps generated by Contact-GraspNet from four different viewpoints, aiming for a more comprehensive coverage of the object’s surface. This enhanced version is denoted as RTAGrasp-MultiView (RTAGrasp-MV). Additionally, we introduce a two-stage variant of our approach, which follows a method similar to the two-stage baselines but replaces Contact-GraspNet with our sampler, which does not incorporate human demonstration guidance. This variant is denoted as Ours-Two-stage (Ours-TS), enabling a more equitable comparison between our single-stage and existing two-stage methods.

Dataset We evaluate HGDdiffuser and baselines on the OakInk dataset [37]. Specifically, we use a subset from OakInk-Shape, selecting 340 object instances across 33 object categories, with approximately 10 instances randomly chosen per category. Each object instance includes an object mesh and a corresponding task-oriented human grasp demonstration. The object mesh is used to generate object point clouds via sampling, while the human grasp, represented using MANO representation, is directly used as input.

Metrics Following prior works, we evaluate the task-oriented grasping success for each object instance and compute the average success rate across all instances. The stability of

task-oriented grasps is automatically evaluated using the NVIDIA Isaac Gym simulation platform [38], while the task relevance is manually assessed based on the corresponding human demonstration. Additionally, we evaluate the inference time to compare the efficiency of different methods.

Implementation Details All experiments are conducted on a desktop PC equipped with a single Nvidia RTX 3090 GPU. HGDdiffuser is optimized using the Adam optimizer [39] with a weight decay of 0.0001 and a learning rate of 0.0001. The model is trained for 500 epochs with a batch size of 32.

TABLE I. COMPARISON OF DIFFERENT GRASP SAMPLING QUANTITIES

Method	#	Success Rate (%)	Inference Time (s)
Ours-TS	100	71.18	0.671
	200	75.35	1.352
	500	78.12	3.428
	1000	77.68	6.793
Ours	1	81.21	0.191
	100	81.38	0.716

Results Figure 5 presents the comparison results against the baseline methods. DemoGrasp achieves a success rate of only 27.85%, primarily because only human demonstrations that conform to manually designed rules can be successfully converted into stable grasps. This limitation aligns with our observations in Section II. RTAGrasp, leveraging a two-stage framework, improves the success rate to 59.06%. However, since grasp candidates are generated from single-view point clouds, they often fail to cover the entire object, potentially missing the grasps corresponding to human demonstrations, which leads to unsuccessful attempts. RTAGrasp-MV aggregates grasps from multiple viewpoints, providing more comprehensive object coverage and improving the success rate to 71.47%. However, this enhancement comes at the

TABLE II. ABLATION STUDY ON FEATURE BACKBONE

Method	Success Rate (%)	Inference Time (s)
GraspLDM	74.32	0.183
Ours w/o DiT	71.35	0.163
Ours	80.65	0.167

cost of increased inference time, reaching 1.019s due to the additional viewpoint processing. Our two-stage variant, Ours-TS, achieves a success rate of 78.12% but requires the longest inference time of 3.428s. In contrast, our HGDiffuser (Ours) integrates human demonstration guidance directly into the sampling process, effectively filtering out grasp candidates that fail to satisfy task constraints in the first stage. As a result, HGDiffuser not only outperforms both RTAGrasp-MV and Ours-TS in terms of success rate but also significantly reduces inference time.

The number of grasp samples significantly impacts both success rate and inference time. For two-stage methods, increasing the number of first-stage grasp samples enhances the likelihood of filtering more stable, task-oriented grasps during the second stage. Table I presents the performance of Ours-TS and Ours under different grasp sampling quantities. As the number of sampled grasps increases, Our-TS success rate improves from 71.18% to 78.12%, but at the expense of significantly longer inference time, ranging from 0.671s to 3.428s. In contrast, Ours is almost unaffected by the number of grasp samples and can achieve the highest success rate by sampling just a single grasp. Ours-TS employs our grasp sampler without human guidance, which maintains consistent performance in 6-DoF grasp sampling compared to Ours. This result demonstrates the significant efficiency advantage of our single-stage sampling over the two-stage approach.

Figure 4 presents the qualitative experimental results of Ours-TS and Ours. Combined with the inference time evaluation in Table I, it can be observed that our method achieves comparable grasp generation quality to two-stage methods while significantly reducing inference time.

B. Ablation Study

We conduct our evaluation on the same OakInk dataset as in previous experiments. Since the introduced DiT blocks are expected to enhance the overall performance of the sampler and to ensure comparability with other existing samplers, we evaluate the performance of the sampler in generating task-agnostic grasps, considering only whether the grasp succeeds in the simulation platform without imposing task constraints. To establish a strong baseline, we include the state-of-the-art diffusion-based grasp sampler GraspLDM [29] in our evaluation. Additionally, we assess a variant of our method without DiT blocks (Ours w/o DiT), which corresponds to the existing approach GraspDiff [27], utilizing an MLP-based feature backbone.

Table II presents the ablation study results. The Ours w/o DiT variant achieves an average grasp success rate of 71.35%, which is 2.97% lower than the state-of-the-art GraspLDM. By incorporating DiT blocks as feature backbones and designing corresponding tokenized inputs, our full model Ours attains a success rate of 80.65%, surpassing

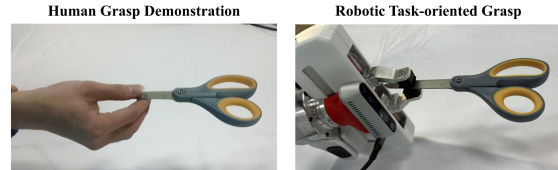


Fig. 6. An example of real-world experiments.

TABLE III. QUANTITATIVE RESULTS OF REAL-WORLD EXPERIMENTS

Stage	Perception	Planning	Action
Success Rate	26 / 30	22 / 30	20 / 30

GraspLDM by 6.33%, while maintaining a nearly unchanged inference time. These results demonstrate that the integration of DiT blocks effectively leverages attention mechanisms to extract more informative features, thereby improving the performance of HGDiffuser.

C. Real-world Experiments

To assess the practicality and applicability of our proposed method, we perform a quantitative evaluation through physical experiments, with the results presented in Table III. The videos of a subset of these experiments are provided in the Supplementary Materials. Our experimental setup utilizes a table-mounted 7-DoF Franka Research 3 arm with a Franka hand, enhanced by a wrist-mounted RealSense D435i camera. We conducted experiments on 30 object-task pairs (including 10 objects and eight tasks), with one human demonstration collected for each pair. Following the evaluation framework of [17], we measure success rates across three stages: perception, planning, and action. Our system achieves an 86.67% success rate in the perception stage, validating the practical applicability of our method.

Physical experiments identify two perception challenges hindering task-oriented grasping in our method. First, the task-agnostic grasp sampler, trained on a large-scale dataset [40], often fails to generate stable grasps for specific object regions, such as the headband of headphones or the handles of scissors and teapots. Even with human-guided sampling, the resulting grasps remain unstable. Second, discrepancies between the object’s pose during grasping and demonstration necessitate pose estimation [24]. However, inaccurate reconstructed meshes or partial occlusion lead to pose estimation errors, causing imprecise grasps—a common issue in demonstration-based methods.

VI. CONCLUSION

In this work, we propose HGDiffuser, a diffusion-based framework that leverages human grasp demonstrations to generate robotic 6-DoF parallel-jaw task-oriented grasps. The human grasp demonstrations are processed to create explicit task-oriented constraints, which are then used to guide the sampling of a pre-trained task-agnostic diffusion model. Compared to existing two-stage methods, HGDiffuser eliminates the need for extensive sampling in the vast task-agnostic grasp space, resulting in significantly higher efficiency and comparable or higher accuracy than the approaches. Evaluation on the OakInk dataset demonstrates the superiority of HGDiffuser over existing methods on generation efficiency.

REFERENCES

- [1] R. Detry, J. Papon, and L. Matthies, "Task-oriented grasping with semantic and geometric scene understanding," in *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2017, pp. 3266–3273.
- [2] W. Dong, D. Huang, J. Liu, C. Tang, and H. Zhang, "RTAGrasp: Learning task-oriented grasping from human videos via retrieval, transfer, and alignment," 2025.
- [3] Y. Ju, K. Hu, G. Zhang, G. Zhang, M. Jiang, and H. Xu, "Robo-abc: Affordance generalization beyond categories via semantic correspondence for robot manipulation," in *European conference on computer vision*. Springer, 2024, pp. 222–239.
- [4] T. W. Nick Heppert, Max Argus, "DITTO: Demonstration imitation by trajectory transformation," in *2024 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2024.
- [5] Y. Cai, J. Gao, C. Pohl, and T. Asfour, "Visual imitation learning of task-oriented object grasping and rearrangement," *arXiv preprint arXiv:2403.14000*, 2024.
- [6] S. Wu, Y. Zhu, Y. Huang, K. Zhu, J. Gu, J. Yu, Y. Shi, and J. Wang, "AffordDP: Generalizable diffusion policy with transferable affordance," *arXiv preprint arXiv:2412.03142*, 2024.
- [7] S. Yan, Z. Zhang, M. Han, Z. Wang, Q. Xie, Z. Li, Z. Li, H. Liu, X. Wang, and S.-C. Zhu, "M2diffuser: Diffusion-based trajectory optimization for mobile manipulation in 3d scenes," *arXiv preprint arXiv:2410.11402*, 2024.
- [8] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in neural information processing systems*, vol. 34, pp. 8780–8794, 2021.
- [9] W. Peebles and S. Xie, "Scalable diffusion models with transformers," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2023, pp. 4195–4205.
- [10] T. Patten and M. Vinze, "Imitating task-oriented grasps from human demonstrations with a low-DoF gripper," in *Proceedings of the sixteenth international conference on autonomic and autonomous systems*, 2020, p. 7.
- [11] M. Kokic, D. Kragic, and J. Bohg, "Learning task-oriented grasping from human activity datasets," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3352–3359, 2020.
- [12] P. Wang, F. Manhardt, L. Minciullo, L. Garattoni, S. Meier, N. Navab, and B. Busam, "DemoGrasp: Few-shot learning for robotic grasping with human demonstration," in *2021 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2021, pp. 5733–5740.
- [13] T. Feix, J. Romero, H.-B. Schmiedmayer, A. M. Dollar, and D. Kragic, "The grasp taxonomy of human grasp types," *IEEE Transactions on human-machine systems*, vol. 46, no. 1, pp. 66–77, 2015.
- [14] M. Sundermeyer, A. Mousavian, R. Triebel, and D. Fox, "Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes," in *2021 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2021, pp. 13 438–13 444.
- [15] C. Tang, A. Xiao, Y. Deng, T. Hu, W. Dong, H. Zhang, D. Hsu, and H. Zhang, "Functo: Function-centric one-shot imitation learning for tool manipulation," *arXiv preprint arXiv:2502.11744*, 2025.
- [16] A. Murali, W. Liu, K. Marino, S. Chernova, and A. Gupta, "Same object, different grasps: Data and semantic knowledge for task-oriented grasping," in *Conference on robot learning*. PMLR, 2021, pp. 1540–1557.
- [17] C. Tang, D. Huang, W. Ge, W. Liu, and H. Zhang, "Graspnet: Leveraging semantic knowledge from a large language model for task-oriented grasping," *IEEE Robotics and Automation Letters*, 2023.
- [18] C. Tang, D. Huang, W. Dong, R. Xu, and H. Zhang, "Foundationgrasp: Generalizable task-oriented grasping with foundation models," *IEEE Transactions on Automation Science and Engineering*, 2025.
- [19] R. Mirjalili, M. Krawez, S. Silenzi, Y. Blei, and W. Burgard, "Langrasp: Using large language models for semantic object grasping," *arXiv preprint arXiv:2310.05239*, 2023.
- [20] H. Li, W. Mao, W. Deng, C. Meng, R. Zhang, F. Jia, T. Wang, H. Fan, H. Wang, and X. Deng, "SegGrasp: Zero-shot task-oriented grasping via semantic and geometric guided segmentation," *arXiv preprint arXiv:2410.08901*, 2024.
- [21] S. Jin, J. Xu, Y. Lei, and L. Zhang, "Reasoning grasping via multi-modal large language model," *arXiv preprint arXiv:2402.06798*, 2024.
- [22] Y. Mu, T. Chen, S. Peng, Z. Chen, Z. Gao, Y. Zou, L. Lin, Z. Xie, and P. Luo, "Robotwin: Dual-arm robot benchmark with generative digital twins (early version)," *arXiv preprint arXiv:2409.02920*, 2024.
- [23] T. Chen, Y. Mu, Z. Liang, Z. Chen, S. Peng, Q. Chen, M. Xu, R. Hu, H. Zhang, X. Li, and others, "G3Flow: Generative 3D Semantic Flow for Pose-aware and Generalizable Object Manipulation," *arXiv preprint arXiv:2411.18369*, 2024.
- [24] B. Wen, W. Yang, J. Kautz, and S. Birchfield, "Foundationpose: Unified 6d pose estimation and tracking of novel objects," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 17 868–17 879.
- [25] G. Pavlakos, D. Shan, I. Radosavovic, A. Kanazawa, D. Fouhey, and J. Malik, "Reconstructing hands in 3d with transformers," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 9826–9836.
- [26] J. Romero, D. Tzionas, and M. J. Black, "Embodied hands: Modeling and capturing hands and bodies together," *arXiv preprint arXiv:2201.02610*, 2022.
- [27] J. Urain, N. Funk, J. Peters, and G. Chalvatzaki, "Se (3)-diffusionfields: Learning smooth cost functions for joint grasp and motion optimization through diffusion," in *2023 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2023, pp. 5923–5930.
- [28] A. Mousavian, C. Eppner, and D. Fox, "6-dof graspnet: Variational grasp generation for object manipulation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 2901–2910.
- [29] K. R. Barad, A. Orsula, A. Richard, J. Dentler, M. Olivares-Mendez, and C. Martinez, "Graspplm: Generative 6-dof grasp synthesis using latent diffusion models," *IEEE access : practical innovations, open solutions*, 2024.
- [30] A. Hyvärinen and P. Dayan, "Estimation of non-normalized statistical models by score matching," *Journal of Machine Learning Research*, vol. 6, no. 4, 2005.
- [31] Y. Song and S. Ermon, "Improved techniques for training score-based generative models," *Advances in neural information processing systems*, vol. 33, pp. 12 438–12 448, 2020.
- [32] C. Deng, O. Litany, Y. Duan, A. Poulernard, A. Tagliasacchi, and L. J. Guibas, "Vector neurons: A general framework for so (3)-equivariant networks," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 12 200–12 209.
- [33] X. Wu, L. Jiang, P.-S. Wang, Z. Liu, X. Liu, Y. Qiao, W. Ouyang, T. He, and H. Zhao, "Point transformer V3: Simpler faster stronger," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 4840–4851.
- [34] L. Yang, X. Zhan, K. Li, W. Xu, J. Li, and C. Lu, "Cpf: Learning a contact potential field to model the hand-object interaction," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 11 097–11 106.
- [35] D. Saito, K. Sasabuchi, N. Wake, J. Takamatsu, H. Koike, and K. Ikeuchi, "Task-grasping from a demonstrated human strategy," in *2022 IEEE-RAS 21st international conference on humanoid robots (humanoids)*. IEEE, 2022, pp. 880–887.
- [36] K. Wang, Y. Fan, and I. Sakuma, "Robot grasp planning: a learning from demonstration-based approach," *Sensors*, vol. 24, no. 2, p. 618, 2024.
- [37] L. Yang, K. Li, X. Zhan, F. Wu, A. Xu, L. Liu, and C. Lu, "Oakink: A large-scale knowledge repository for understanding hand-object interaction," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 20 953–20 962.
- [38] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, and others, "Isaac gym: High performance gpu-based physics simulation for robot learning," *arXiv preprint arXiv:2108.10470*, 2021.
- [39] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [40] C. Eppner, A. Mousavian, and D. Fox, "Acronym: A large-scale grasp dataset based on simulation," in *2021 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2021, pp. 6222–6227.