# A Versatile and Differentiable Hand-Object Interaction Representation

Théo Morales[1], Omid Taheri[2], and Gerard Lacey[3]

[1]Trinity College Dublin
[2]Max Planck Institute for Intelligent Systems
[3]Maynooth University
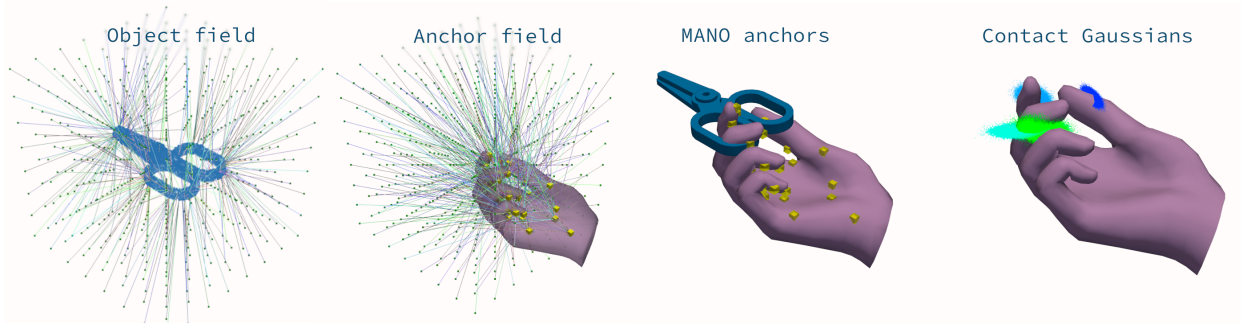tmorales@tcd.ie, omid.taheri@tuebingen.mpg.de, gerard.lacey@mu.ie

Figure 1. Decomposition of our Coarse Hand-Object Interaction Representation (CHOIR). From left to right, our representation encodes the object geometry with point-wise distances in a regular grid (coloured rays), the hand shape and pose as point-wise distances to 32 MANO anchors on the mesh surface (yellow cubes), and the hand contact points as probability densities from 3D Gaussian distributions (coloured point clouds). CHOIR is a fully-differentiable, versatile representation of the hand-object pair in object frame.

## Abstract

*Synthesizing accurate hands-object interactions (HOI) is critical for applications in Computer Vision, Augmented Reality (AR), and Mixed Reality (MR). Despite recent advances, the accuracy of reconstructed or generated HOI leaves room for refinement. Some techniques have improved the accuracy of dense correspondences by shifting focus from generating explicit contacts to using rich HOI fields. Still, they lack full differentiability or continuity and are tailored to specific tasks. In contrast, we present a Coarse Hand-Object Interaction Representation (CHOIR), a novel, versatile and fully differentiable field for HOI modelling. CHOIR leverages discrete unsigned distances for continuous shape and pose encoding, alongside multivariate Gaussian distributions to represent dense contact maps with few parameters. To demonstrate the versatility of CHOIR we design JointDiffusion, a diffusion model to learn a grasp distribution conditioned on noisy hand-object interactions or only object geometries, for both refinement and synthesis applications. We demonstrate JointDiffusion's improvements over the SOTA in both applications: it increases the contact F1 score by 5% for refinement and decreases the sim. displacement by 46% for synthesis. Our experiments show that JointDiffusion with CHOIR yield superior contact accuracy and physical realism compared to SOTA methods designed for specific tasks. Project page: https://theomorales.com/CHOIR*

## 1. Introduction

Numerous computer vision applications could benefit from highly accurate hand pose prediction in object manipulation scenarios, such as Augmented Reality (AR) or Mixed Reality (MR), human-robot collaboration, etc. However, SOTA models still struggle to generalize to novel grasps on unknown objects [12, 19], in both synthesis and reconstruction. The problem is challenging because hands are small, dexterous, with many degrees of freedom, making it hard to be accurately tracked or reconstructed. Additionally, interactions naturally come with occlusions or noisy observations, making it harder to estimate accurate

hand-object interactions. Such inaccuracies, like subtle hand-object penetrations or slightly off-positioned fingers, can significantly affect the realism of the hand-object interactions. A common approach is to train a model to reason in 3D space and predict coarse hand-object poses from images, and then refine them with a model trained on hand-object contacts [1, 16, 21, 45, 53]. This coarse-to-fine approach generates an estimate of how an unknown object is being grasped, while dense hand-object interactions (learned or simulated) serve as a test-time optimization (TTO) objective to refine the estimate.

Recently, there has been progress in dense contact map prediction, either directly from images or meshes [10, 16, 20, 21, 50]. However, they still have some limitations such as: compute intensity (typically involving point cloud processing), the need for feature engineering [16], and being uninformative in cases where the hand is approaching but has not yet touched the object. To address these issues, recent methods have proposed to define the hand pose and shape in an object-centric space using ray casting or spring systems [47, 53]. While these hand-object representations directly improve the refinement capabilities of TTO, they still require engineering, are compute-intensive, or are not fully differentiable. To address the gaps, we propose to use a *Coarse Hand-Object Interaction Representation*, named CHOIR, a novel field leveraging unsigned distances and multivariate Gaussian distribution to represent shape, pose, and contact maps for hand-object interactions. CHOIR encodes the object geometry as distances to the fixed Basis Point Set representation (BPS) [37], and the hand pose and shape as distances from the same basis points to the fixed MANO anchors proposed by [47]. In addition, CHOIR encodes coarse contact maps represented as 3D Gaussian distributions around the MANO anchors, such that dense contact maps can be inferred from probability densities. As such, it is scalable, fully differentiable, and efficient on GPUs. To demonstrate its effectiveness, we train a conditional Denoising Diffusion Probabilistic Model (DDPM), named *JointDiffusion*, to learn the distribution of hand-object interactions in CHOIR representation. We demonstrate plausible grasp synthesis alongside noisy grasp refinement through the same model architecture trained on different condition variables.

Overall, experiments demonstrate that our method outperforms baselines on denoising and generating static hand interactions and that our approach offers superior contact-based metrics. Our models and code will be available for research purposes.

In the direction of solving hand-object interaction challenges, this work makes the following key contributions:

- We propose CHOIR, a versatile and differentiable representation that encodes hand-object interactions, enhancing accuracy in contact modelling.

- Our method introduces a novel way to represent dense contact maps using Gaussian distributions, leading to more accurate hand-object contacts. In addition, we propose a novel and simple way to compute contact weights for all hand vertices.

- We employ a multimodal conditional diffusion model tailored to our CHOIR framework, which works for both synthesizing plausible grasps and refining noisy ones.

## 2. Related works

Despite many advances in hand motion tracking or reconstruction, estimating accurate hand-object interaction poses is still a challenging and unsolved problem. Recently, there has been a push towards the coarse-to-fine paradigm for hand-object interaction, where a coarse hand pose is first generated or reconstructed, and then is refined via optimization [16, 18, 19, 53] with pseudo-ground-truth or using learning-based methods [43, 44]. In this section, we review the most relevant works and their limitations.

**Hand-Object Interaction Reconstruction:** With the growing availability of rich annotated datasets for hand-object interaction [3, 32, 41], many recent works focus on simultaneously reconstructing the hands and objects from images [7, 15–17, 19, 28, 55]. Many of these works leverage deep learning techniques to estimate the hand and object poses [18, 19]. However, the initial hand and object pose from these methods are often approximate and require further refining. To achieve this, some work optimizes the results further using contact constraints or interaction constraints [42, 51]. Zhou et al. [53] proposed to use a spatiotemporal field for hand-object interaction and train a network to refine this instead. They then use the refined field in a two-step optimization process to get the refined hand poses. This is however slow due to the field not being fully differentiable and requiring a search algorithm. This pre-optimization imposes a lower bound on the optimization time. Here we propose a lightweight and fully differentiable field on which we can optimize hand meshes solely based on the L2 norm.

**Grasp Synthesis:** Grasp synthesis, split into static and dynamic domains, has received much attention recently. In the static domain, many classic methods use physical constraints to satisfy realistic grasps [5, 14, 27, 29, 36]. Newer methods take a learning-based approach and use big datasets of hand-object interactions to learn grasps [6, 7, 13, 22–24, 43, 54]. These often generate the pose parameters of a model directly [6, 7], or estimate an implicit representation for the grasp [24, 53].

Another body of work focuses on generating dynamic grasps. Similar to static grasps, some methods define contact constraints and use optimization to satisfy them

[30, 31, 42, 48, 49, 52]. For better motion realism, recent methods use reinforcement learning (RL) for hand grasp generation [2, 4, 11, 33–35, 38]. However, in both static and dynamic grasp generations, the grasps are mostly inaccurate and require further refinement. To further refine grasps there exist optimization-based methods [16, 22, 51, 53] or learning-based ones [43, 44]. Some methods like [43, 51] directly refine the hand pose. Instead of directly operating on the poses, [53] proposes to refine an implicit interaction field for the hand motions and then use it in an optimization process to refine the hand poses. This, however, is very slow due to the complicated nature of the proposed interaction field. Here we train a diffusion model on our novel CHOIR representation, where we can both generate and refine hand-object interaction by only conditioning the model on different observations.

**Hand-Object Interaction Representation:** Implicit representations are increasingly gaining traction in the field, especially for hand-object interaction representation. The Grasping Field, proposed by Karunratanakul et al. [25], introduces an SDF with hand-parts labels. While it has the advantages of our proposed interaction field, namely being coarse and distance-based, it primarily utilizes whole hand and object point clouds as inputs, leading to a high-dimensional model. Their method, however, does not emphasize grasp refinement or denoising. In contrast, ContactOpt [16] advocates for a dense contact map based on hand and object meshes. Although high-dimensional, it offers significant value in grasp refinement and denoising.

Yang et al. [47] propose CPF with coarse anchors and an innovative spring system. Despite its promising direction, it involves minimizing relatively intricate energy functions, proving time-consuming at TTO. Furthermore, its inability to handle non-full or dynamic grasps presents limitations. Jiang et al. [21] introduced a distinct approach with hand-object contact consistency reasoning. By employing a CVAE for initial predictions followed by a contact network, they present dense contact maps enriched with prior contact regions. Diverging from the above methods, Yu et al. [50] used a UV-Based 3D hand-object reconstruction for grasp optimization. While valuable, its image-centric nature is not well-suited for hands and fingers that are not in direct contact with the object. Additionally, the proposed gSDF by Chen et al. [9], despite its precise functionalities remain to be explored further.

The SOTA in grasp refinement is represented by TOCH [53], recognized for its features like accommodating approaching hands and fingers not in direct contact. While it excels in dynamic grasps, it has not been evaluated for static ones. In this work, we focus on designing a contact-dense, expressive interaction field for multiple applications.
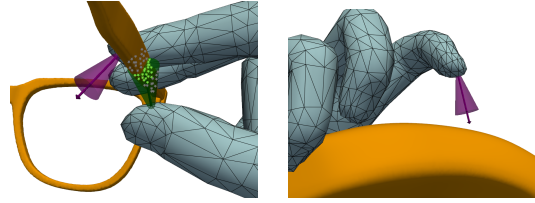


Figure 2. Illustration of the cone of tolerance used to determine the raw hand contact weights. For each hand vertex, the weights are the count of object points inside the vertex's cone defined along its normal vector. (Left) The green points on the object's surface are inside the cone, hence contributing to the hand vertex's weight while the grey points do not. (Left & Right) No object points are inside the purple cone: its vertex has a contact weight of 0.

## 3. Method

At the heart of our approach is the development of a novel representation for hand-object interactions, coined CHOIR, which addresses the limitations inherent in state-of-the-art techniques, particularly those relying on ray-based dense correspondence fields such as TOCH [53]. Our representation is designed to encode shape, pose, and contacts while remaining fully differentiable and continuous. We exploit it in two settings: (1) to refine grasps from noisy predictions of an off-the-shelf method for hand-object reconstruction, and (2) to synthesize realistic grasps given an object shape. To do so, we design a DDPM backbone based on the U-Net architecture which jointly decodes the contact parameters and unsigned distances from a shared latent space for efficient learning.

In this section, we first go through the details of the proposed representation and its implementation. We then describe the DDPM and context encoder.

### 3.1. Coarse hand-object interaction representation (CHOIR)

**Notation:** In the following, we denote a dataset sample as $X$ and a 3D vector as $\boldsymbol{x}$. Superscripts indicate a sample index while subscripts indicate a point index. We differentiate vectors from scalars by bold symbols for the former. Note that $x$ on its own denotes a generic data sample of any form.

The proposed Coarse Hand-Object Interaction Representation called *CHOIR*, is a novel field for representing hand-object interaction using unsigned distances and multivariate Gaussian distributions. The object geometry and relative hand pose are represented with unsigned distances from a common set of points, following the Basis Point Set (BPS) representation [37]. It is a lightweight 3D point cloud representation with fixed dimensionality that enables the use of convolutions with a regular point grid. Define a dataset $\mathcal{D} = \{X^1, \ldots, X^N\}$ consisting of $N$ point clouds where

each point cloud $X^n$ is composed of points $\{\boldsymbol{x}_1^n, \ldots, \boldsymbol{x}_{K_n}^n\}$. A basis point set $\mathcal{B} = \{\boldsymbol{b}^1, \ldots, \boldsymbol{b}^M\}$ is defined as a regular grid in $\mathbb{R}^3$. Then, the dataset is normalized such that $X^n$ fits in the grid and its centroid is at the origin. Finally, the BPS representation of each point cloud $X^n$ is computed as

$$
X_{\mathrm{BPS}}^n = \begin{bmatrix} \min_k \|\boldsymbol{b}^1 - \boldsymbol{x}_k^n\|_2^2 \\ \min_k \|\boldsymbol{b}^2 - \boldsymbol{x}_k^n\|_2^2 \\ \ldots \\ \min_k \|\boldsymbol{b}^M - \boldsymbol{x}_k^n\|_2^2 \end{bmatrix}. \tag{1}
$$

CHOIR represents the hand-object interaction as (a) a coarse hand pose in a canonical object frame using the MANO parametric hand mesh [40], and (b) probabilistic hand contact points.

### 3.1.a Shape and pose representation

The pose part of CHOIR is defined as a concatenation of the BPS representation of the object mesh and the distances from the BPS to the 32 pre-assigned MANO anchors proposed by [47], *i.e.*, a CHOIR specifies an object point cloud $X^n$ together with a hand mesh $H$. The anchor distances $\boldsymbol{d}_H = [d_H^1, \ldots, d_H^M]^T$ are given by

$$
d_H^j = \|\boldsymbol{b}^j - \delta_H(j)\|_2^2 \tag{2}
$$

where the function $\delta_H(j)$ returns the anchor for point $\boldsymbol{b}^j$ and hand mesh $H$. Note that the same MANO anchor can be assigned to multiple basis points. We propose two assignment schemes $\delta_H$: (1) a repeating pattern of the 32 ordered indices and (2) a shuffled version of the latter. We did not find any difference in accuracy between the two, which we show quantitatively in Appendix B.1, thus we use assignment (1).

### 3.1.b Probabilistic contact representation

Instead of representing hand contact points as a discrete vector mapping each MANO vertex to a contact weight or binary class, we opt for a lightweight and continuous representation based on 3D multivariate Gaussian distributions. Given a hand-object pair $(H, X^n)$, we first compute the contact weights $w_i$ for each MANO vertex $\boldsymbol{v}_i$. We define it as the count of all object surface points within a cone of tolerance defined at the root of $\boldsymbol{v}_i$, in the direction of the vertex normal $\boldsymbol{n}_i$ (see Fig. 2). This approach is inspired by ContactOpt's contact capsules [16], which includes vertices inside the mesh. With a cone, we exclude these vertices such that the contact map does not encode penetration patches in favour of a simpler optimization objective. Effectively, $\boldsymbol{w}_i = |\mathcal{S}|$ where $\mathcal{S}$ is the set of points $\boldsymbol{x}$ obeying the



(a) Raw hand contact weights in red.

(b) 3D Gaussian distributions as coloured point clouds.

(c) Recovered contact weights (left) *vs*. raw contact weights (right).

Figure 3. Visualization of our probabilistic contact maps (best seen in colour). (a) The raw hand contact weights are computed with our cone of tolerance method. (b) 32 3D Gaussian distributions are fitted – one for each MANO anchor – on the weights to obtain contact probability densities. (c) Comparison of the recovered probabilistic dense contact map and of the raw contact weights. Our method leaves gaps in the contact map to allow for a $2mm$ penetration and improve contact fitting.

following conditions (where we set $\lambda = 4$mm and $\kappa = \frac{4}{3}\pi$):

$$
\boldsymbol{x} \in X^n, \tag{3}
$$

$$
\|\boldsymbol{x} - \boldsymbol{v}_i\| \leq \lambda, \tag{4}
$$

$$
\arccos \frac{\boldsymbol{n}_i \cdot (\boldsymbol{x} - \boldsymbol{v}_i)^T}{\boldsymbol{n}_i \|\boldsymbol{x} - \boldsymbol{v}_i\|} \leq \kappa. \tag{5}
$$

From this discrete contact map, we fit one 3D multivariate Gaussian distribution per MANO anchor. This is done using the weighted mesh vertices such that the probability densities match the location of the vertices with the most contact weight. In effect, given the set of vertices $\mathcal{V}$ and the set of associated weights $\mathcal{W}$, we define the multiset $\mathcal{V}_w$ composed of each element $\boldsymbol{v}_i \in \mathcal{V}$ repeated $w_i \in \mathcal{W}$ times:

$$
\mathcal{V}_w = \{\underbrace{\boldsymbol{v}_i, \boldsymbol{v}_i, \ldots, \boldsymbol{v}_i}_{w_i \text{ times}} \text{ for all } \boldsymbol{v}_i \in \mathcal{V}, w_i \in \mathcal{W}\}. \tag{6}
$$

We then maximize the likelihood of the Gaussian parameters given the multiset $\mathcal{V}_w$. This allows us to encode a probabilistic dense contact map for the hand mesh as a set of 32 multivariate normal distributions (MVN). Each hand vertex then gets a contact probability by querying the probability density function of the nearest anchor's Gaussian (see Fig. 3). For anchor $j$, the MVN is parameterized by a mean vector $\boldsymbol{\mu}^j \in \mathbb{R}^3$ and a covariance matrix $\Sigma^j \in \mathbb{R}^{3 \times 3}$. Since the latter must be positive semi-definite, it can be challenging to predict it with a neural network. One approach is

4

to represent it as a lower triangular matrix $L^j$ obtained via Cholesky decomposition, such that $\Sigma^j = L^j L^{jT}$, and to enforce the diagonal entries to be positive. The final form of our probabilistic contact representation for a given hand mesh $H$ is thus:

$$c_H = \begin{bmatrix} \boldsymbol{\mu}^0 & \boldsymbol{l}^0 \\ \boldsymbol{\mu}^1 & \boldsymbol{l}^1 \\ \cdots & \\ \boldsymbol{\mu}^{31} & \boldsymbol{l}^{31} \end{bmatrix} \in \mathbb{R}^{32 \times 9} \tag{7}$$

where $\boldsymbol{l}^j \in \mathbb{R}^6$ is the vector containing the elements of and below the diagonal of the matrix $L^j$. In summary, a CHOIR (see Fig. 1) is defined as

$$x_{\text{CHOIR}} = [X_{\text{BPS}}^n \in \mathbb{R}^M, \boldsymbol{d}_H \in \mathbb{R}^M, \boldsymbol{c}_H \in \mathbb{R}^{32 \times 9}]. \tag{8}$$

By encoding coarse hand-object correspondences in this way, we can fit hand meshes onto ground-truth CHOIRs with less than 1mm absolute mean per-joint pose error. Thus, generating valid CHOIRs is the accuracy bottleneck; in the next subsection, we present our learning method.

## 3.2. Learning conditional distributions of CHOIR

Denoising Diffusion Probabilistic Models (DDPM) have recently made their prowess in distribution learning for high dimensional problems [39]. The combination of recent improvements to the U-Net architecture and the DDPM framework enables the modelling of complex relationships between context and target information. We propose to exploit these advances to model complex conditional CHOIR distributions with multiple modalities of context, such as noisy hand-object pairs or object shapes.

Our goal is to determine the conditional distribution of hand poses $p(\boldsymbol{d}_H, \boldsymbol{c}_H | y)$ based on an observation $y$, where $\boldsymbol{d}_H$ and $\boldsymbol{c}_H$ are parts of CHOIR (see Eq. (8)). The context $y$ is either (1) a noisy hand-object pair, encoded as a CHOIR with missing contacts $\boldsymbol{c}_H$, or (2) an object point cloud encoded as $X_{\text{BPS}}^n \in \mathbb{R}^M$, *i.e.* a CHOIR with missing contacts and hand pose. To learn this distribution we separately predict the noise samples for the distance field $\boldsymbol{d}_h$ and the contact Gaussians $\boldsymbol{c}_H$, denoted $\boldsymbol{\epsilon}_d$ and $\boldsymbol{\epsilon}_c$ respectively. This is motivated by the structure of $\boldsymbol{d}_H$ which allows the use of convolutions, while $\boldsymbol{c}_H$ is a vector in $\mathbb{R}^{32 \times 9}$. Thus, our DDPM backbone (see Fig. 4) is composed of a 3D U-Net for the prediction of $\boldsymbol{\epsilon}_d$, and of a second decoder for the prediction of $\boldsymbol{\epsilon}_c$. This contact decoder is a fully connected residual network whose input is a concatenation of the latent variable $\boldsymbol{z}_t$ and the latent features from the bottleneck layer of the U-Net. This encourages the model to learn pose and contact features in a shared space, such that the latent codes are relevant to both $\boldsymbol{d}_H$ and $\boldsymbol{c}_H$.

Our U-Net implementation uses Multi-Head Self-Attention (MHSA) to encourage relevant feature extrac-

tion and Multi-Head Cross-Attention (MHCA) to condition the network on the context. The latter is embedded with an encoder identical to the U-Net encoder, with an additional spatial pooling mechanism via a shallow fully-connected network. We train one model per modality of context and include experiments on a multi-modal model in Appendix C.5.

Ultimately, we propose to learn the conditional distribution $p(\boldsymbol{d}_H, \boldsymbol{c}_H | y)$ in two settings:

1. Where $y$ is a noisy observation of a CHOIR, with missing contacts $\boldsymbol{c}_H$, such that
   $x_{\text{CHOIR}} = [X_{\text{BPS}}^n \in \mathbb{R}^M, \boldsymbol{d}_H \in \mathbb{R}^M]$.

2. Where $y$ is an object point cloud in BPS representation $X_{\text{BPS}}^n \in \mathbb{R}^M$.

We then sample from this distribution to obtain a full CHOIR. In the next subsection, we show how to obtain MANO parameters in the object coordinate system.

## 3.3. Test-Time Optimization (TTO)

While the state-of-the-art HOI fields are either not fully differentiable [53] or rely on random restarts [16, 53], fitting a hand mesh to CHOIR is done by gradient descent in two stages. Firstly, we fit a MANO mesh to the unsigned distance field of CHOIR. Secondly, we adjust the hand contacts to the nearest object points using the contact Gaussians of CHOIR. These two stages rely on distance minimization with continuous losses, giving a smoother loss landscape than methods using contact agreement objectives [16, 22].

### 3.3.a Coarse pose and shape fitting stage

The first stage's objective $\mathcal{L}_{\text{PoseShape}}$ is composed of a reconstruction loss $\mathcal{L}_{\text{rec}}$, a shape regularizer $\mathcal{L}_{\text{shape}}$ and a pose regularizer $\mathcal{L}_{\text{pose}}$:

$$\begin{aligned} \mathcal{L}_{\text{PoseShape}} &= \lambda_1 \cdot \mathcal{L}_{\text{rec}} + \lambda_2 \cdot \mathcal{L}_{\text{shape}} + \lambda_3 \cdot \mathcal{L}_{\text{pose}} \\ \mathcal{L}_{\text{rec}} &= \|\boldsymbol{d}_H - \hat{\boldsymbol{d}_H}\|_2^2 \\ \mathcal{L}_{\text{shape}} &= \|\boldsymbol{\beta}_{\text{MANO}}\|_2 \\ \mathcal{L}_{\text{pose}} &= \|\boldsymbol{\theta}_{\text{MANO}} - \boldsymbol{\theta}_{\text{MANO}}^{\text{init}}\|_2 \end{aligned} \tag{9}$$

where $\boldsymbol{d}_H$ and $\hat{\boldsymbol{d}_H}$ are the respective ground-truth and predicted anchor distances (see Eq. (2)), $\boldsymbol{\beta}_{\text{MANO}}$ and $\boldsymbol{\theta}_{\text{MANO}}$ are the MANO shape and pose parameters, respectively. The shape regularizer prevents the hand mesh from over-deforming to satisfy the reconstruction loss, while the pose regularizer prevents strong deviation from the initial MANO pose estimate $\boldsymbol{\theta}_{\text{MANO}}^{\text{init}}$. Note that in the grasp synthesis case, we remove the pose regularizer.

We minimize $\mathcal{L}_{\text{PoseShape}}$ w.r.t. the MANO parameters alongside rotation and translation of the wrist joint with the Adam optimizer [26]. We set $\lambda_1$ to 1000 to bring the loss
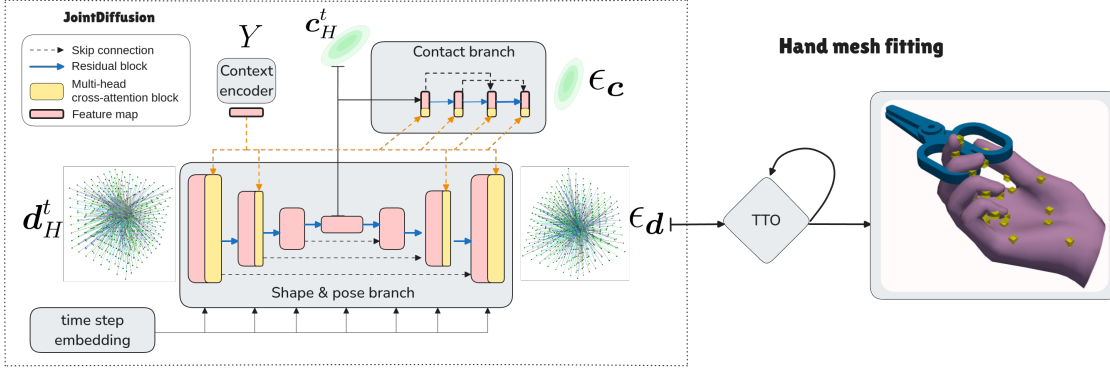
Figure 4. Architecture of *JointDiffusion*. The 3D U-Net predicts the noise sample $\epsilon_d$ for the hand distance field $d_H$. The contact prediction branch predicts the noise sample $\epsilon_c$ for the contact Gaussian parameters $c_H$ from the features of the U-Net's bottleneck. This joint learning encourages the U-Net to extract features relevant to both tasks, enhancing the accuracy of the learned CHOIR distribution.

into the millimetre scale and found $\lambda_2 = 1 \times 10^{-4}$ and $\lambda_3 = 1 \times 10^{-8}$ to work well in practice. In the grasp refinement setting, all parameters are initialized with the noisy inputs from which CHOIR observations are built, leading to fast convergence ($\sim 150$ iterations). Since the pose and shape encodings are unsigned distances, the fitting loss is obtained in very few lines of Python code (see supplementary material). This TTO stage fits a hand mesh to the predicted distance field but does not account for contacts and penetration. For this, we introduce stage two.

### 3.3.b  Dense contact fitting stage

In the second stage, we refine the obtained hand grasp by minimizing the weighted distances from the MANO vertices to their nearest object points under some constraints. The weights of vertices $v$ are obtained from the probability distribution function (PDF) $\Phi_j(v_i)$ of the nearest anchor's contact Gaussian (see Eq. (7)). For each MANO vertex $v_i$, nearest anchor $j$ and nearest set of $K$ (such as 5) object points $\mathcal{K}^n$, the reconstruction loss is:

$$\mathcal{L}_{\text{rec}} = \sum_{i=1}^{N} \sum_{k=1}^{K} \Phi_j(v_i) \cdot \|v_i - p_k\|_2^2, \quad p \in \mathcal{K}^n. \quad (10)$$

This objective is minimized in conjunction with a penetration regularizer following [16], the shape regularizer defined in Eq. (9), and a pose regularizer to avoid deviating from the initial solution, defined as

$$\mathcal{L}_{\text{pose}} = \eta_1 \cdot \|R_{\text{MANO}} - R_{\text{MANO}}^{\text{stage1}}\|_2 + \eta_2 \cdot \|t_{\text{MANO}} - t_{\text{MANO}}^{\text{stage1}}\|_2 \quad (11)$$

where $R_{\text{MANO}}$ and $t_{\text{MANO}}$ are the respective rotation matrix and translation vector of the MANO mesh, $\eta_1 = 1 \times 10^{-2}$ and $\eta_2 = 1 \times 10^{-1}$. The final objective is

$$\mathcal{L}_{\text{Contacts}} = \lambda_4 \cdot \mathcal{L}_{\text{rec}} + \lambda_5 \cdot \mathcal{L}_{\text{penetration}} + \lambda_6 \cdot \mathcal{L}_{\text{pose}} + \lambda_2 \cdot \mathcal{L}_{\text{shape}} \quad (12)$$

where $\lambda_4 = 10$, $\lambda_5 = 1000$, and $\lambda_6 = 0.5$. We optimize this loss w.r.t. the same parameters and with the same method as in the previous stage. In the next section, we evaluate the combined CHOIR + *JointDiffusion* + TTO solution on grasp refinement and synthesis benchmarks, and show how much each component contributes to the accuracy and plausibility of our grasps.

## 4. Evaluation

Our solution consists of (a) a representation (CHOIR), (b) a learning method (*JointDiffusion*) for denoising and synthesizing interactions, and (c) a hand-mesh fitting algorithm (TTO). We evaluate this solution in two settings and in Appendix C.5, we evaluate a multi-modal variant trained in both settings.

**Grasp refinement:** We replicate the benchmark for refining noisy grasps, proposed by Grady et al. [16], which consists of a perturbed version of the ContactPose dataset [8]. ContactPose comprises highly accurate hand poses for static grasps of 25 objects performed by 50 participants. Grady et al. [16] define large perturbations on the hand poses as 3 additive and i.i.d. noise components: (1) translation noise $\epsilon_t \sim \mathcal{N}(0, 5)$ in cm, (2) pose noise $\epsilon_\theta \sim \mathcal{N}(0, 0.05)$ in PCA space, and (3) rotation noise $\epsilon_R \sim \mathcal{N}(0, 15)$ in radians. However, they omit a validation split to have more training data. We instead split the dataset with 70% data for training, 10% for validation and the last 20% for testing. For training, we use 16 perturbed versions of each sample and 4 for validation and testing. We retrain ContactOpt [16] on these new splits and also train all methods on object splits to evaluate generalizability in Appendix C.3. Additionally, we retrain the recent SOTA method TOCH [53], specialized in denoising dynamic grasps on the GRAB [43] benchmark, which has not been evaluated on static grasps with this amount of noise.

**Grasp synthesis:** In addition to refining noisy interac-

Table 1. Evaluation of our approach on static grasp refinement against SOTA methods on the Perturbed ContactPose benchmark. * means reported figures. All methods are evaluated with one non-cherry-picked output per sample. *JointDiffusion* shows greater contact accuracy and outperforms all methods on most metrics, especially contact metrics (F1, Precision, Recall) and intersection volume, showing greater contact fidelity on the hand locations. Best results are in bold and second best are underlined.

| Method | MPJPE (mm) ↓ | R-MPJPE (mm) ↓ | IV ($cm^3$) ↓ | F1 (%) ↑ | Precision (%) ↑ | Recall (%) ↑ |
|---|---|---|---|---|---|---|
| Perturbed data | 83.02 | 21.55 | 6.99 | 1.55 | 1.88 | 2.74 |
| ContactOpt [16] | 32.88 | <u>28.17</u> | 12.83* | 17.27 | 13.24 | **34.30** |
| TOCH [53] | **26.96** | 29.24 | <u>10.14</u> | <u>22.23</u> | <u>21.46</u> | 25.09 |
| *JointDiffusion* (ours) | <u>27.69</u> | **23.54** | **6.04** | **27.20** | **25.21** | <u>32.80</u> |

| Ground truth | Observation | ContactOpt [16] | TOCH [53] | *JointDiffusion* |
|---|---|---|---|---|



Figure 5. Qualitative comparison of grasp denoising on one challenging case of the Peturbed ContactPose benchmark. Our method produces less penetration than TOCH [53], and substantially better output than ContactOpt [16] which maximizes hand-object contact.

tions, we demonstrate that our *JointDiffusion* can be used to synthesize novel interactions for unseen objects. To do this, we train *JointDiffusion* on ContactPose with the object mesh encoded in BPS representation as input to the context encoder. For quantitative comparison, we retrain a recent SOTA method in grasp synthesis: GraspTTA [22]. This method also uses test-time adaptation, which makes it a good baseline to compare against. We use the same training, validation and test splits as for Sec. 4. This benchmark evaluates the capabilities of our model to learn the complex interaction between hands and objects. In grasp refinement, partial noisy information is given to the model during inference, but in grasp synthesis, the model must generate plausible grasps without prior information other than the training data. In Appendix C.5, we evaluate our multimodal variant

on the OakInk benchmark [46] against GrabNet [43].

### 4.1. Qualitative & quantitative results

We use several key metrics to quantify hand pose and contact error, as detailed in Appendix C.1. In particular, we employ the (Root-aligned) Mean Per-Joint Pose Error (R-MPJPE and MPJPE) for either world space error (MPJPE) or object space error (R-MPJPE). However, a low pose error is not always indicative of a realistic or well-refined grasp. For denoising, recovering intended contacts is more important: the Precision score captures intended grasp locations, while a high Recall score implies fewer false negatives. The latter can be misleadingly increased by maximizing hand contact, leading to less dexterous grasps (as shown by ContactOpt [16] on Fig. 5). The F1 sore (harmonic mean of Pre-

| Input object | Sample 1 | Sample 2 | Sample 3 | Sample 4 |



Figure 6. Generated grasps obtained with *JointDiffusion* on the ContactPose benchmark. The synthesized grasps show plausible grasps with good finger-object contact and minimal penetration, showing the expressive contact modelling of CHOIR.

Table 2. Evaluation of our approach against GraspTTA [22] on static grasp generation for the ContactPose benchmark. †: contact fitting enabled. Best results are in bold, second best are underlined.

| Method | IV (cm$^3$) ↓ | SD (cm) ↓ |
|---|---|---|
| GraspTTA [22] | 5.17 | 3.81 |
| *JointDiffusion* | 8.13 | 2.07 |
| *JointDiffusion* † | **4.51** | **2.05** |

cision and Recall) is the most meaningful metric for hand contact fidelity. For synthesis, we employ the simulation displacement (SD) metric with IV to evaluate the feasibility and stability of grasps.

For grasp refinement, Tab. 1 shows that our method *JointDiffusion* outperforms two SOTA methods on most metrics, and comes second best in the remaining metrics. In particular, our method brings a 5% improvement over TOCH [53] and 10% over ContactOpt [16] in contact F1 score. Our method demonstrates the lowest R-MPJPE (−4.6mm over ContactOpt) and highest contact precision (+3.8% over TOCH), indicating a higher contact and grasp fidelity than both methods. With the highest contact precision, the lowest intersection volume (−4mm over TOCH), and a high contact recall, *JointDiffusion* demonstrates the most accurate contact inference, as seen with a challenging case on Fig. 5. This comes at the cost of a negligible penalty in absolute pose, where TOCH outperforms *Joint-Diffusion* by less than 1mm. More qualitative comparisons between *JointDiffusion* and ContactOpt are available in Appendix C.2. ContactOpt remains 2% better in hand contact Recall since it aims to maximize hand-object contact and thus reduces false negatives. However, with a 12% worse contact Precision than our method, it cannot yield the intended grasp with high accuracy, as reflected by a 10% lower F1 score and a failure on a challenging case in Fig. 5.

Table 3. Ablation study of CHOIR on grasp refinement (best seen in colour). Best metrics are in bold and second best are underlined, with improvement or degradation w.r.t. the previous row in parenthesis. A keypoint diffusion model is used as a baseline (see Appendix B.3). The BPS representation for shape and pose encoding improves all contact metrics at the cost of a slightly higher pose error. Adding the probabilistic contacts substantially improves the contact metrics with a small cost in pose accuracy.

| Method | MPJPE (mm) ↓ | IV (cm$^3$) ↓ | F1 (%) ↑ | Precision (%) ↑ | Recall (%) ↑ |
|---|---|---|---|---|---|
| KP. Baseline | **22.11** (-0.00) | 9.18 (-0.00) | 19.45 (+0.00) | 20.22 (+0.00) | 21.18 (+0.00) |
| *+BPS* | 24.69 (+2.58) | 9.12 (-0.06) | 20.77 (+1.32) | 20.86 (+0.64) | 23.38 (+2.20) |
| *+BPS +Contacts* | 27.69 (+4.00) | **6.04** (-3.14) | **27.20** (+6.43) | **25.21** (+4.35) | **32.80** (+9.42) |

For grasp synthesis, Fig. 6 shows *JointDiffusion*'s ability to generate plausible and realistic grasps. It shows minimal penetration and consistent contact between the used fingers and the object, owing to rich contact modelling through CHOIR. These results are validated quantitatively on Tab. 2, where our solution outperforms the SOTA method, GraspTTA [22], while being more versatile in applications. *JointDiffusion* reduces the intersection volume by 13% and the simulation displacement by 46%, resulting in more stable and feasible grasps. More qualitative results can be found in Appendix C.4.

## 4.2. Ablation study

To validate the design choices of CHOIR, we conduct an ablation study of its components by starting from a keypoint baseline and adding CHOIR components. The baseline uses a model similar to *JointDiffusion* to learn hand joint keypoints (see Appendix B.3). In Tab. 3, *+BPS* corresponds to *JointDiffusion* without contact representation, and *+BPS +Contacts* matches *JointDiffusion* with the full CHOIR. We evaluate each row on the Perturbed Contact-Pose benchmark. Tab. 3 shows that while the baseline yields a lower MPJPE, it also gives the highest IV and lowest con-

tact scores. The first two rows yield better pose accuracy by ignoring physical plausibility and contact fidelity, simplifying the learning objective. The last row improves on all contact and penetration metrics. Hence, CHOIR offers the best compromise for pose accuracy and contact fidelity in the denoising setting, while yielding plausible and stable grasps in the synthesis setting. These findings corroborate the previous experimental results.

## 5. Conclusion

In this work, we introduced the novel Coarse Hand-Object Interaction Representation (CHOIR), a versatile and fully-differentiable hand-object interaction field. CHOIR represents hand and object shape and pose as unsigned distance, and dense contacts using probability distributions, leading to more accurate hand-object interactions. Additionally, leveraging CHIOR we train a diffusion model, *JointDiffusion*, to both refine or generate hand-object interactions. CHOIR demonstrates improvements in pose and contact accuracy over existing representations, providing a compact representation for refining or synthesizing hand-object interaction poses.

**Limitations & Future Work:** Despite its advancements, our method is not without limitations. The reliance on the BPS representation may limit the ability to capture detailed interactions. Furthermore, the model's focus on static interactions might restrict its application in real-world scenarios. To address these limitations and extend the utility of our framework, future work will focus on: (a) investigating learnable geometry representations to capture more detailed interactions, (b) handling dynamic interactions by incorporating temporal information into CHOIR.

## References

[1] Ahmed Tawfik Aboukhadra, Jameel Nawaz Malik, Ahmed Elhayek, Nadia Robertini, and Didier Stricker. Thor-net: End-to-end graformer-based realistic two hands and object reconstruction with self-supervision. *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 1001–1010, 2022. 2

[2] Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Józefowicz, Bob McGrew, Jakub W. Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, Jonas Schneider, Szymon Sidor, Josh Tobin, Peter Welinder, Lilian Weng, and Wojciech Zaremba. Learning dexterous in-hand manipulation. *Int. J. Robotics Res.*, 39(1), 2020. 3

[3] L. Ballan, A. Taneja, J. Gall, L. Van Gool, and M. Pollefeys. Motion Capture of Hands in Action Using Discriminative Salient Points. In *European Conference on Computer Vision (ECCV)*, pages 640–653, 2012. 2

[4] Kevin Bergamin, Simon Clavet, Daniel Holden, and James Richard Forbes. Drecon: Data-driven responsive control of physics-based characters. *ACM Trans. Graph.*, 38(6), nov 2019. 3

[5] Jeannette Bohg, Antonio Morales, Tamim Asfour, and Danica Kragic. Data-driven grasp synthesis—a survey. *IEEE Transactions on Robotics*, 30:289–309, 2014. 2

[6] Samarth Brahmbhatt, Cusuh Ham, Charles C. Kemp, and James Hays. ContactDB: Analyzing and predicting grasp contact via thermal imaging. In *Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[7] 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 *European Conference on Computer Vision (ECCV)*, volume 12358, pages 361–378, 2020. 2

[8] Samarth Brahmbhatt, Chengcheng Tang, Christopher D. Twigg, Charles C. Kemp, and James Hays. Contactpose: A dataset of grasps with object contact and hand pose. *ArXiv*, abs/2007.09545, 2020. 6, 12, 21

[9] Zerui Chen, Shizhe Chen, Cordelia Schmid, and Ivan Laptev. gsdf: Geometry-driven signed distance functions for 3d hand-object reconstruction. *ArXiv*, abs/2304.11970, 2023. 3

[10] Hoseong Cho, Chanwoo Kim, Jihyeon Kim, Seongyeong Lee, Elkhan Ismayilzada, and Seungryul Baek. Transformer-based unified recognition of two hands manipulating objects. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4769–4778, 2023. 2

[11] Sammy Christen, Muhammed Kocabas, Emre Aksan, Jemin Hwangbo, Jie Song, and Otmar Hilliges. D-Grasp: Physically plausible dynamic grasp synthesis for hand-object interactions. In *Computer Vision and Pattern Recognition (CVPR)*, 2022. 3

[12] Enric Corona, Albert Pumarola, Guillem Alenyà, Francesc Moreno-Noguer, and Grégory Rogez. GanHand: Predicting human grasp affordances in multi-object scenes. In *Computer Vision and Pattern Recognition (CVPR)*, pages 5030–5040, 2020. 1

[13] Enric Corona, Albert Pumarola, Guillem Alenyà, Francesc Moreno-Noguer, and Gregory Rogez. GanHand: Predicting human grasp affordances in multi-object scenes. In *Computer Vision and Pattern Recognition (CVPR)*, pages 5030–5040, 2020. 2

[14] Sahar El-Khoury, Anis Sahbani, and Philippe Bidaud. 3D objects grasps synthesis: A survey. In *IFToMM World Congress on Mechanism and Machine Science*, 2011. 2

[15] Guillermo Garcia-Hernando, Shanxin Yuan, Seungryul Baek, and Tae-Kyun Kim. First-person hand action benchmark with RGB-D videos and 3D hand pose annotations. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[16] Patrick Grady, Chengcheng Tang, Christopher D. Twigg, Minh Vo, Samarth Brahmbhatt, and Charles C. Kemp.

Contactopt: Optimizing contact to improve grasps. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1471–1481, 2021. 2, 3, 4, 5, 6, 7, 8, 13, 18, 19, 22

[17] Shreyas Hampali, Mahdi Rad, Markus Oberweger, and Vincent Lepetit. Honnotate: A method for 3d annotation of hand and object poses. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 3193–3203. Computer Vision Foundation / IEEE, 2020. 2

[18] Yana Hasson, Bugra Tekin, Federica Bogo, Ivan Laptev, Marc Pollefeys, and Cordelia Schmid. Leveraging photometric consistency over time for sparsely supervised hand-object reconstruction. *CoRR*, abs/2004.13449, 2020. 2

[19] Yana Hasson, Gül Varol, Dimitris Tzionas, Igor Kalevatykh, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *Computer Vision and Pattern Recognition (CVPR)*, pages 11807–11816, 2019. 1, 2

[20] Junxing Hu, Hongwen Zhang, Zerui Chen, Mengcheng Li, Yunlong Wang, Yebin Liu, and Zhen Sun. Learning explicit contact for implicit reconstruction of hand-held objects from monocular images. *ArXiv*, abs/2305.20089, 2023. 2

[21] Hanwen Jiang, Shaowei Liu, Jiashun Wang, and Xiaolong Wang. Hand-object contact consistency reasoning for human grasps generation. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11087–11096, 2021. 2, 3

[22] Hanwen Jiang, Shaowei Liu, Jiashun Wang, and Xiaolong Wang. Hand-object contact consistency reasoning for human grasps generation. In *Proceedings of the International Conference on Computer Vision*, 2021. 2, 3, 5, 7, 8, 21

[23] Korrawe Karunratanakul, Adrian Spurr, Zicong Fan, Otmar Hilliges, and Siyu Tang. A skeleton-driven neural occupancy representation for articulated hands. In *2021 International Conference on 3D Vision (3DV)*, pages 11–21, 2021. 2

[24] Korrawe Karunratanakul, Jinlong Yang, Yan Zhang, Michael J. Black, Krikamol Muandet, and Siyu Tang. Grasping Field: Learning implicit representations for human grasps. In *International Conference on 3D Vision (3DV)*, pages 333–344, 2020. 2

[25] Korrawe Karunratanakul, Jinlong Yang, Yan Zhang, Michael J. Black, Krikamol Muandet, and Siyu Tang. Grasping field: Learning implicit representations for human grasps. *2020 International Conference on 3D Vision (3DV)*, pages 333–344, 2020. 3

[26] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. 5

[27] Paul G. Kry and Dinesh K. Pai. Interaction capture and synthesis. *Transactions on Graphics (TOG)*, 25(3):872–880, 2006. 2

[28] T. Kwon, B. Tekin, J. Stühmer, F. Bogo, and M. Pollefeys. H2O: Two Hands Manipulating Objects for First Person Interaction Recognition. In *International Conference on Computer Vision (ICCV)*, pages 10138–10148, 2021. 2

[29] Ying Li, Jiaxin L. Fu, and Nancy S. Pollard. Data-driven grasp synthesis using shape matching and task-based pruning. *Transactions on Visualization and Computer Graphics (TVCG)*, 13(4):732–747, 2007. 2

[30] Karen C. Liu. Dextrous manipulation from a grasping pose. *Transactions on Graphics (TOG)*, 28(3):59, 2009. 3

[31] Igor Mordatch, Zoran Popovic, and Emanuel Todorov. Contact-invariant optimization for hand manipulation. In *Symposium on Computer Animation (SCA)*, pages 137–144, 2012. 3

[32] Iason Oikonomidis, Nikolaos Kyriazis, and Antonis A. Argyros. Full DOF tracking of a hand interacting with an object by modeling occlusions and physical constraints. In *International Conference on Computer Vision (ICCV)*, 2011. 2

[33] Soohwan Park, Hoseok Ryu, Seyoung Lee, Sunmin Lee, and Jehee Lee. Learning predict-and-simulate policies from unorganized human motion data. *ACM Trans. Graph.*, 38(6), nov 2019. 3

[34] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. DeepMimic: Example-guided deep reinforcement learning of physics-based character skills. *Transactions on Graphics (TOG)*, 37(4):143:1–143:14, 2018. 3

[35] Xue Bin Peng, Glen Berseth, Kangkang Yin, and Michiel Van De Panne. Deeploco: Dynamic locomotion skills using hierarchical deep reinforcement learning. *ACM Trans. Graph.*, 36(4), jul 2017. 3

[36] Nancy S. Pollard and Victor Brian Zordan. Physically based grasping control from example. In *International Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*, pages 311–318, 2005. 2

[37] Sergey Prokudin, Christoph Lassner, and Javier Romero. Efficient learning on point clouds with basis point sets. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4331–4340, 2019. 2, 3

[38] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. In *Proceedings of Robotics: Science and Systems*, Pittsburgh, Pennsylvania, June 2018. 3

[39] Robin Rombach, A. Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10674–10685, 2021. 5

[40] Javier Romero, Dimitrios Tzionas, and Michael J. Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, Nov. 2017. 4

[41] Srinath Sridhar, Franziska Mueller, Michael Zollhoefer, Dan Casas, Antti Oulasvirta, and Christian Theobalt. Real-time joint tracking of a hand manipulating an object from RGB-D input. In *European Conference on Computer Vision (ECCV)*, 2016. 2

[42] Omid Taheri, Vasileios Choutas, Michael J. Black, and Dimitrios Tzionas. GOAL: Generating 4D whole-body motion for hand-object grasping. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3

[43] Omid Taheri, Nima Ghorbani, Michael J. Black, and Dimitrios Tzionas. GRAB: A dataset of whole-body human

grasping of objects. In *European Conference on Computer Vision (ECCV)*, volume 12349, pages 581–600, 2020. 2, 3, 6, 7, 21

[44] Omid Taheri, Yi Zhou, Dimitrios Tzionas, Yang Zhou, Duygu Ceylan, Soren Pirk, and Michael J. Black. Grip: Generating interaction poses using latent consistency and spatial cues. 2023. 2, 3

[45] Rong Wang, Wei Mao, and Hongdong Li. Interacting hand-object pose estimation via dense mutual attention. *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 5724–5734, 2022. 2

[46] Lixin Yang, Kailin Li, Xinyu Zhan, Fei Wu, Anran Xu, Liu Liu, and Cewu Lu. OakInk: A large-scale knowledge repository for understanding hand-object interaction. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 7, 21

[47] 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. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11077–11086, 2020. 2, 3, 4

[48] Yuting Ye and C. Karen Liu. Synthesis of detailed hand manipulations using contact sampling. *Transactions on Graphics (TOG)*, 31(4):41:1–41:10, 2012. 3

[49] Yuting Ye and Karen C. Liu. Synthesis of detailed hand manipulations using contact sampling. *Transactions on Graphics (TOG)*, 31(4):41:1–41:10, 2012. 3

[50] Ziwei Yu, Linlin Yang, You Xie, Ping Chen, and Angela Yao. Uv-based 3d hand-object reconstruction with grasp optimization. In *British Machine Vision Conference*, 2022. 2, 3

[51] He Zhang, Yuting Ye, Takaaki Shiratori, and Taku Komura. ManipNet: Neural manipulation synthesis with a hand-object spatial representation. *Transactions on Graphics (TOG)*, 40(4):121:1–121:14, 2021. 2, 3

[52] Wenping Zhao, Jianjie Zhang, Jianyuan Min, and Jinxiang Chai. Robust realtime physics-based motion control for human grasping. *Transactions on Graphics (TOG)*, 32(6):207:1–207:12, 2013. 3

[53] Keyang Zhou, Bharat Lal Bhatnagar, Jan Eric Lenssen, and Gerard Pons-Moll. Toch: Spatio-temporal object-to-hand correspondence for motion refinement. In *European Conference on Computer Vision*, 2022. 2, 3, 5, 6, 7, 8, 12, 18, 19, 22

[54] Tianqiang Zhu, Rina Wu, Xiangbo Lin, and Yi Sun. Toward human-like grasp: Dexterous grasping via semantic representation of object-hand. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15721–15731, 2021. 2

[55] C. Zimmermann, D. Ceylan, J. Yang, B. Russell, M. Argus, and T. Brox. FreiHAND: A Dataset for Markerless Capture of Hand Pose and Shape From Single RGB Images. In *International Conference on Computer Vision (ICCV)*, pages 813–822, 2019. 2

## A. Supplementary material

## B. Method details

In this section, we include additional information regarding our representation and learning method.

### B.1. CHOIR: Anchor assignment

Table 4. Average reconstruction error for MANO meshes fitted onto ground-truth CHOIRs with the *ordered* and *random* anchor assignment schemes. Mean Per-Joint Pose Error (MPJPE) and Mean Per-Vertex Pose Error (MPVPE) are averaged across the entire ContactPose [8] dataset.

|              | Ordered | Random |
|--------------|---------|--------|
| MPJPE (mm)   | 0.18    | 0.19   |
| MPVPE (mm)   | 0.22    | 0.22   |

Tab. 4 shows that both the *ordered* and *random* anchor assignment schemes produce the same reconstruction error when fitting a ground-truth CHOIR from the Contact-Pose [8] dataset. The Mean Per-Joint Pose Error (MPJPE) and Mean Per-Vertex Pose Error (MPVPE) metrics were averaged across the entire dataset. Note that with ground-truth hand-object meshes, the obtained CHOIR allows fitting a MANO mesh with less than 1mm error.

### B.2. Test-Time Optimization: Fitting loss

The Python code for the stage 1 of the TTO loss fits in a few lines of code:

```python
anchor_dist = torch.cdist(
    bps, anchors
) # Anchors predicted in TTO
distances = torch.gather(
    anchor_dist, 2, anchor_ids
)
choir_loss = F.mse_loss(
    distances, choir[..., -1]
) # Agreement of anchors and CHOIR
```

Source Code 1. Minimal Python code for the stage 1 TTO loss.

### B.3. Keypoint baseline

To evaluate the expressiveness and efficacy of each component of CHOIR, we design a diffusion model backbone that allows us to fit a simpler alternative to CHOIR. This simpler representation only encodes the hand pose and shape as 21 MANO joints $\boldsymbol{j}_H \in \mathbb{R}^{21 \times 3}$ and 32 MANO anchors $\boldsymbol{a}_H \in \mathbb{R}^{32 \times 3}$. The object is encoded as a vector of $K$ randomly sampled surface points $\boldsymbol{p}_O \in \mathbb{R}^{K \times 3}$ where we set $K = 4096$ to match CHOIR which uses a grid of

$16 \times 16 \times 16$ basis points. The final keypoint representation is defined as

$$\boldsymbol{r}_{\text{kp}} = [\boldsymbol{p}_O \in \mathbb{R}^{K \times 3}, \boldsymbol{j}_H \in \mathbb{R}^{21 \times 3}, \boldsymbol{a}_H \in \mathbb{R}^{32 \times 3}]. \quad (13)$$

However, as in *JointDiffusion*, this model learns to predict the hand part only, defined as

$$\boldsymbol{r}_{\text{kp}}^{\text{H}} = [\boldsymbol{j}_H \in \mathbb{R}^{21 \times 3}, \boldsymbol{a}_H \in \mathbb{R}^{32 \times 3}] \quad (14)$$

The backbone of this diffusion model is composed only of residual blocks made of multi-layer perceptrons (MLPs). We use 4 residual blocks with a hidden dimensionality of 512.

In effect, in this baseline, we only replace the 3D U-Net component of *JointDiffusion* with a residual MLP and remove the contact prediction branch, while keeping cross-attention and the same timestep conditioning scheme. The context encoder is also replaced with a residual MLP of hidden dimensionality 2048. We experimented with a PointNet++-based encoder but observed a degradation in performance.

### B.4. Runtime costs

To evaluate the computational costs of CHOIR, we timed its computation and that of TOCH [53] for 50 grasps on an RTX 2080Ti and Intel i9-7900X. On average, TOCH takes $\sim 8.89$s ($\pm 3.99$) while CHOIR takes $\sim 0.13$s ($\pm 0.015$), a $68\times$ reduction. When looking at the total inference time, including the model representation computation, forward pass and TTO, ours converges in $\sim 49$s ($\pm 16$) and TOCH in $\sim 23$s ($\pm 4.3$). Our diffusion model accounts for $\sim 13$s of the total (27%), hence is a major runtime bottleneck. Diffusion Models are inherently slow, but they are becoming faster, and new alternatives with similar properties can be easily integrated since our representation is agnostic to the learning method.

## C. Additional experiments and results

### C.1. Evaluation metrics

In our experiments, we use the following metrics to evaluate the fitted hand mesh to the predicted CHOIR:

- **Mean Per-Joint Pose Error (MPJPE)/(R-MPJPE) (mm)**: L2 norm between ground-truth (GT) and predicted hand joints. We compute both absolute (MPJPE) and root-aligned (R-MPJPE) metrics. The former tells us about the position of the hand around the object, and the latter tells us about the hand grasp error regardless of the spatial pose.

- **Intersection Volume (IV) (cm$^3$)**: A measure of hand-object mesh penetration. It is computed by voxelizing the hand and object meshes (1mm voxels) and computing the volume of the intersecting voxels.

- **Hand contact F1/precision/recall (%)**: The precision and recall scores are measured on binary hand contact maps obtained by upsampling the MANO mesh and computing the Chamfer distance to the object point cloud. Hand vertices within 2mm of their nearest object point are considered in contact, to emulate soft tissue deformation as in [16]. A high precision means a low false positives count, while a high recall means a low false negatives count. The F1 score is the harmonic mean of both and is a measure of predictive performance.

- **Simulation Displacement (SD) (cm)**: The distance of displacement of the object in world space when applying inward forces to the hand grasp in a physics simulation. This tells how stable the grasp is, since more hand-object contact patches result in higher friction and therefore lower displacement.

## C.2. Perturbed ContactPose

We show a qualitative comparison of our method *vs*. ContactOpt [16] on several objects. Fig. 7 shows failure cases in some challenging cases. While ContactOpt [16] fails to produce a plausible grasp for each object and noisy input, our method delivers satisfying results that still closely match the contacts of the ground-truth hand pose. Further qualitative samples are shown in Fig. 8, Fig. 9, and Fig. 10, where our method demonstrates fidelity in the reconstructed finger contacts, as opposed to ContactOpt [16].

| Method | Ground truth | Observation | Prediction |
|--------|-------------|-------------|------------|

Figure 7. Failure cases on a comparison of *JointDiffusion* and ContactOpt for the Perturbed ContactPose benchmark. While ContactOpt consistently fails at producing a plausible mesh after multiple restarts, our method results in minimal penetration and respected finger contacts with only one sample.

Figure 8. Qualitative comparison of *JointDiffusion vs*. ContactOpt on the Perturbed ContactPose benchmark. Our method, *JointDiffusion*, produces plausible grasps and maintains the fidelity of finger contacts while ContactOpt fails in challenging cases even with several random restarts. Our method only draws one sample and performs TTO without random restarts.

| Method | Ground truth | Observation | Prediction |
|--------|--------------|-------------|------------|



Figure 9. Qualitative comparison of *JointDiffusion vs*. ContactOpt on the Perturbed ContactPose benchmark. Our method, *JointDiffusion*, produces plausible grasps and maintains the fidelity of finger contacts while ContactOpt fails in challenging cases even with several random restarts. Our method only draws one sample and performs TTO without random restarts.

| Method | Ground truth | Observation | Prediction |
|---|---|---|---|

Figure 10. Qualitative comparison of *JointDiffusion vs*. ContactOpt on the Perturbed ContactPose benchmark. Our method, *JointDiffusion*, produces plausible grasps and maintains the fidelity of finger contacts while ContactOpt fails in challenging cases even with several random restarts. Our method only draws one sample and performs TTO without random restarts.

## C.3. Object splits experiment

To evaluate the generalizability of our method in the grasp refinement setting, we retrain all methods on the Perturbed ContactPose benchmark [16] with object splits instead of subject splits. We hold 2 objects out of the validation split, and reserve 5 objects for the test split, namely: *doorknob, eyeglasses, apple, bowl, toothbrush*. This increases the difficulty of the benchmark, as all test objects were unseen during training. For a method to perform well in this setting, it must learn generalizable hand-object interaction in latent space. Tab. 5 shows that our method outperforms ContactOpt [16] on most contact-based metrics, and TOCH [53] on all metrics. ContactOpt [16] retains an edge on the recall score since it maximizes the hand-object contact ratio and therefore minimizes false negatives, but at the cost of less contact fidelity since its precision score is significantly lower than *JointDiffusion*. However, TOCH [53] fails to generalize to these objects, which can be explained by the lack of object representation in the TOCH field. We consider this task to be a main challenge in hand-object interaction understanding and will focus on object generalization in future work.

## C.4. Grasp synthesis

Fig. 11 and Fig. 12 show samples of our generative model given an object mesh as input. The model is trained on the improved Perturbed ContactPose benchmark [16], *i.e.* all objects are seen during training. *JointDiffusion* generates visually plausible grasps with consistent finger contacts and minimal mesh penetration. In addition, to enhance visibility, we provide non-cherry-picked supplementary videos of generated hand grasps.

Table 5. Quantitative evaluation of our approach on static grasp refinement against ContactOpt [16] on the Perturbed ContactPose benchmark with object splits. * means reported figures. *JointDiffusion* is evaluated with one non-cherry-picked generated grasp per sample. *JointDiffusion* shows greater contact accuracy and outperforms ContactOpt [16] on most contact metrics, although ContactOpt [16] retains a greater recall score due to its objective which maximizes the hand-object contact ratio, hence reducing false negatives. Best results are in bold, second best are underlined.

| Method | MPJPE (mm) ↓ | R-MPJPE (mm) ↓ | IV (cm$^3$) ↓ | F1 (%) ↑ | Precision (%) ↑ | Recall (%) ↑ |
|---|---|---|---|---|---|---|
| Perturbed data | 83.02 | 21.55 | 6.99 | 1.55 | 1.88 | 2.74 |
| ContactOpt [16] | **35.05** | **29.13** | <u>12.83</u>* | <u>15.39</u> | <u>12.04</u> | **30.36** |
| TOCH [53] | 48.27 | 51.13 | 17.63 | 11.18 | 10.74 | 13.54 |
| *JointDiffusion* | 42.54 | 29.55 | **2.90** | **21.40** | **21.94** | <u>23.05</u> |



Figure 11. Qualitative evaluation of our method, *JointDiffusion*, trained on the object modality of input for grasp synthesis. Each sample is generated from the same input, the object mesh in canonical pose. *JointDiffusion* produces plausible grasps with minimal mesh penetration and consistent finger contacts.

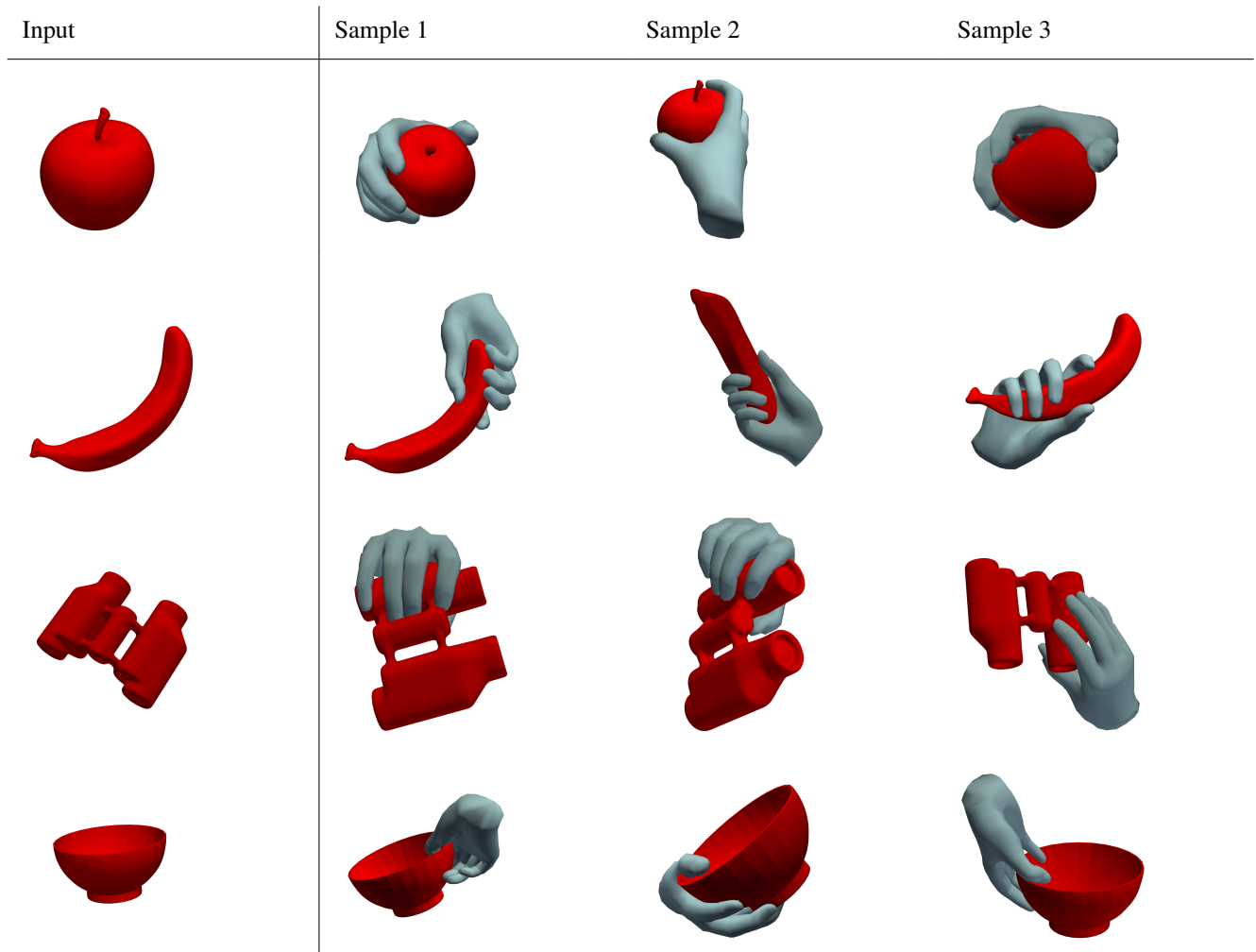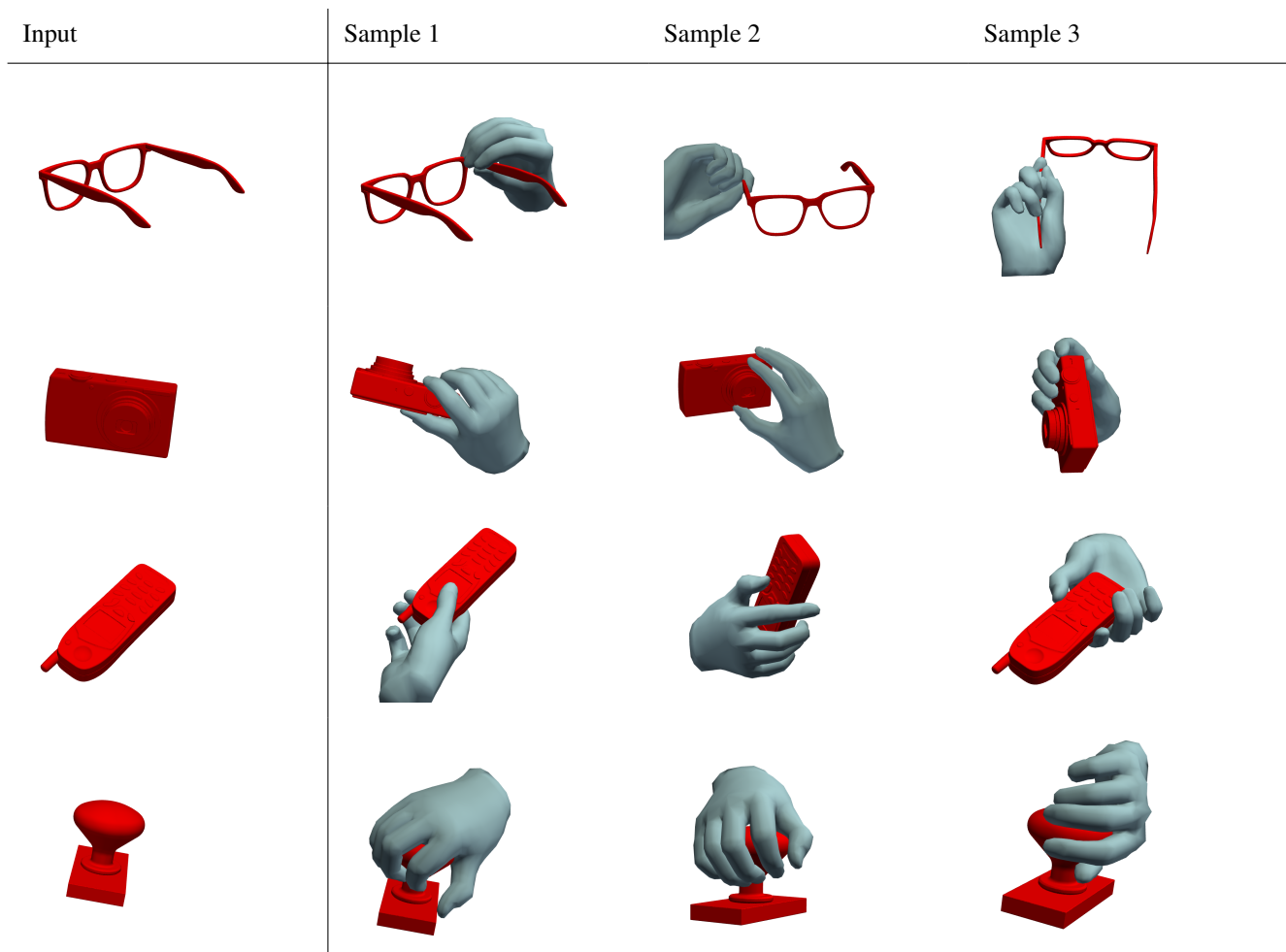| Input | Sample 1 | Sample 2 | Sample 3 |
|-------|----------|----------|----------|



Figure 12. Qualitative evaluation of our method, *JointDiffusion*, trained on the object modality of input for grasp synthesis. Each sample is generated from the same input, the object mesh in canonical pose. *JointDiffusion* produces plausible grasps with minimal mesh penetration and consistent finger contacts.

## C.5. Multimodal model: grasp refinement & synthesis

We further explore our model expressiveness by jointly training two context encoders along with the diffusion backbone of *JointDiffusion*, as opposed to separately trained models for object conditioning and noisy hand-object pair conditioning. Fig. 13 shows qualitative results of the grasp synthesis from this model, while Fig. 14 shows qualitative results of the grasp denoising task for the same model. We trained two multimodal models: one on the ContactPose [8] dataset, and one on the OakInk [46] dataset which we only evaluate on grasp synthesis.

A quantitative evaluation on the denoising task is shown on Tab. 7, and one on the generation task is shown on Tab. 6. For the latter, the increase in simulation displacement (SD) for our method with contact fitting suggests that some hand penetration is helpful to a stable grasp. Note that the synthetic nature of most OakInk samples results in incorrect vertex normals, adversely affecting our penetration regularization loss and performance. This could be solved with a different approach to penetration regularization, such as via the signed distance function.

Table 6. Evaluation of our multimodal model on static grasp generation against two state-of-the-art methods on two benchmarks. *JointDiffusion* outperforms GraspTTA [22] on the ContactPose benchmark [8] and is on par with GrabNet [43] on the OakInk benchmark [46]. We used reported metrics for GrabNet [43] from the OakInk paper [46] and sampled one grasp per dataset sample for our method on both benchmarks. Best results are in bold.

| Method | ContactPose [8] | | OakInk [46] | |
|---|---|---|---|---|
| | IV (cm$^3$) ↓ | SD (cm) ↓ | IV (cm$^3$) ↓ | SD (cm) ↓ |
| GraspTTA [22] | 5.17 | **3.81** | - | - |
| GrabNet [43] | - | - | 6.60 | **1.21** |
| *JointDiffusion* | **5.13** | 5.80 | **5.98** | 5.84 |

Table 7. Evaluation of our approach on static grasp refinement against two SOTA methods and our baseline on the Perturbed ContactPose benchmark. * means reported figures. Our multimodal model is marked with †. Both *JointDiffusion* variants were evaluated with one non-cherry-picked generated grasp per sample. While our baseline yields better reconstruction accuracy in absolute pose, our full model *JointDiffusion* shows greater contact accuracy and outperforms ContactOpt [16] and TOCH [53] on almost all metrics. The multimodal version still outperforms these baselines on contact-based metrics and IV score for grasp refinement, while also being able to do grasp synthesis. Best results are in bold, second best are underlined.

| Method | MPJPE (mm) ↓ | R-MPJPE (mm) ↓ | IV (cm³) ↓ | F1 (%) ↑ | Precision (%) ↑ | Recall (%) ↑ |
|---|---|---|---|---|---|---|
| Perturbed data | 83.02 | 21.55 | 6.99 | 1.55 | 1.88 | 2.74 |
| ContactOpt [16] | 32.88 | 28.17 | 12.83* | 17.27 | 13.24 | **34.30** |
| TOCH [53] | **26.96** | 29.24 | 10.14 | 22.23 | 21.46 | 25.09 |
| *JointDiffusion* | 27.69 | **23.54** | 6.04 | **27.20** | **25.21** | 32.80 |
| *JointDiffusion* † | 35.45 | 33.10 | **5.62** | 24.88 | 23.87 | 29.24 |

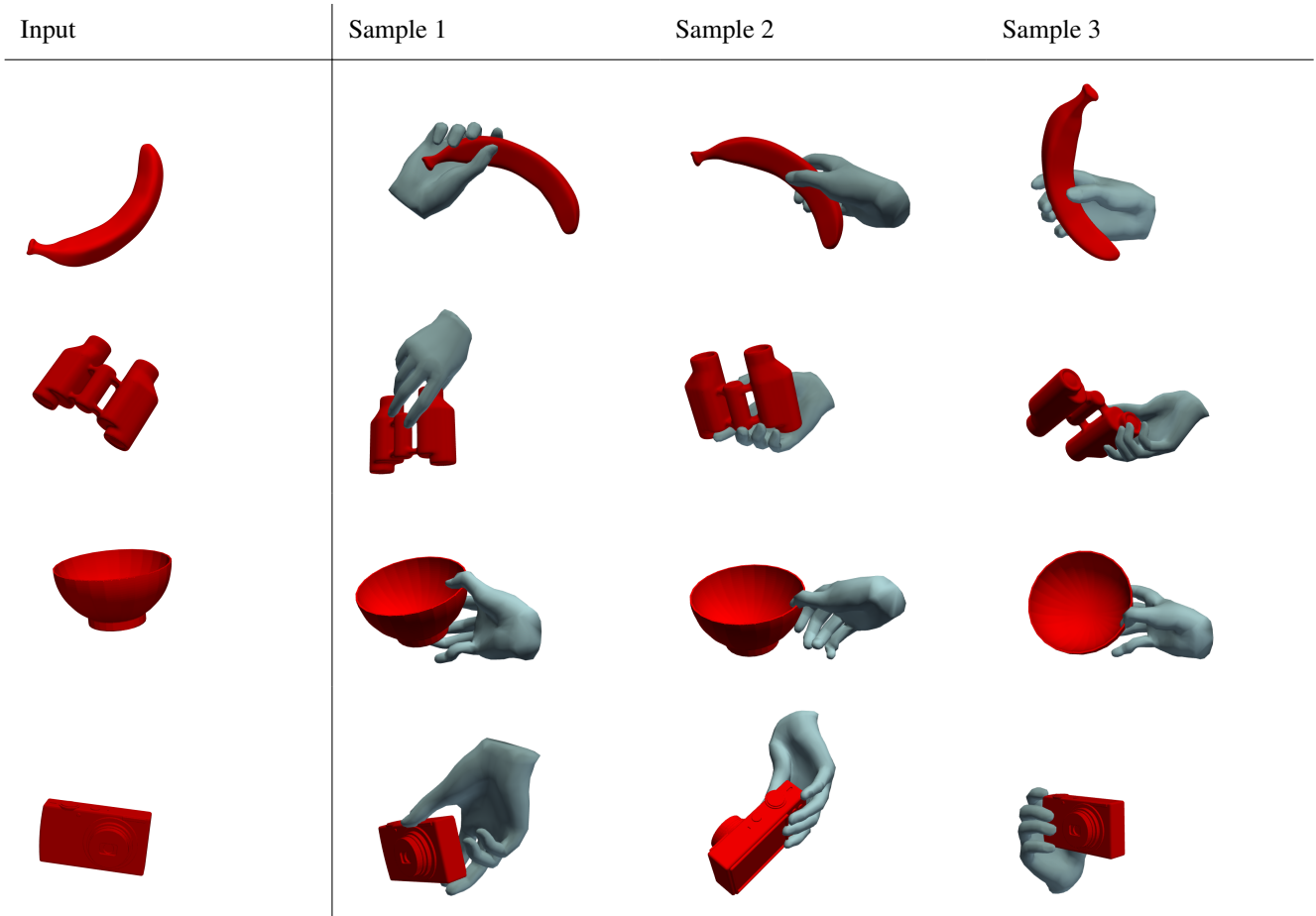| Input | Sample 1 | Sample 2 | Sample 3 |
|---|---|---|---|



Figure 13. Qualitative evaluation of our multimodal *JointDiffusion*, trained on both object and noisy hand-object pair modalities, in the grasp synthesis setting. Each sample is generated from the same input, the object mesh in canonical pose. *JointDiffusion* produces plausible grasps with minimal mesh penetration and consistent finger contacts.
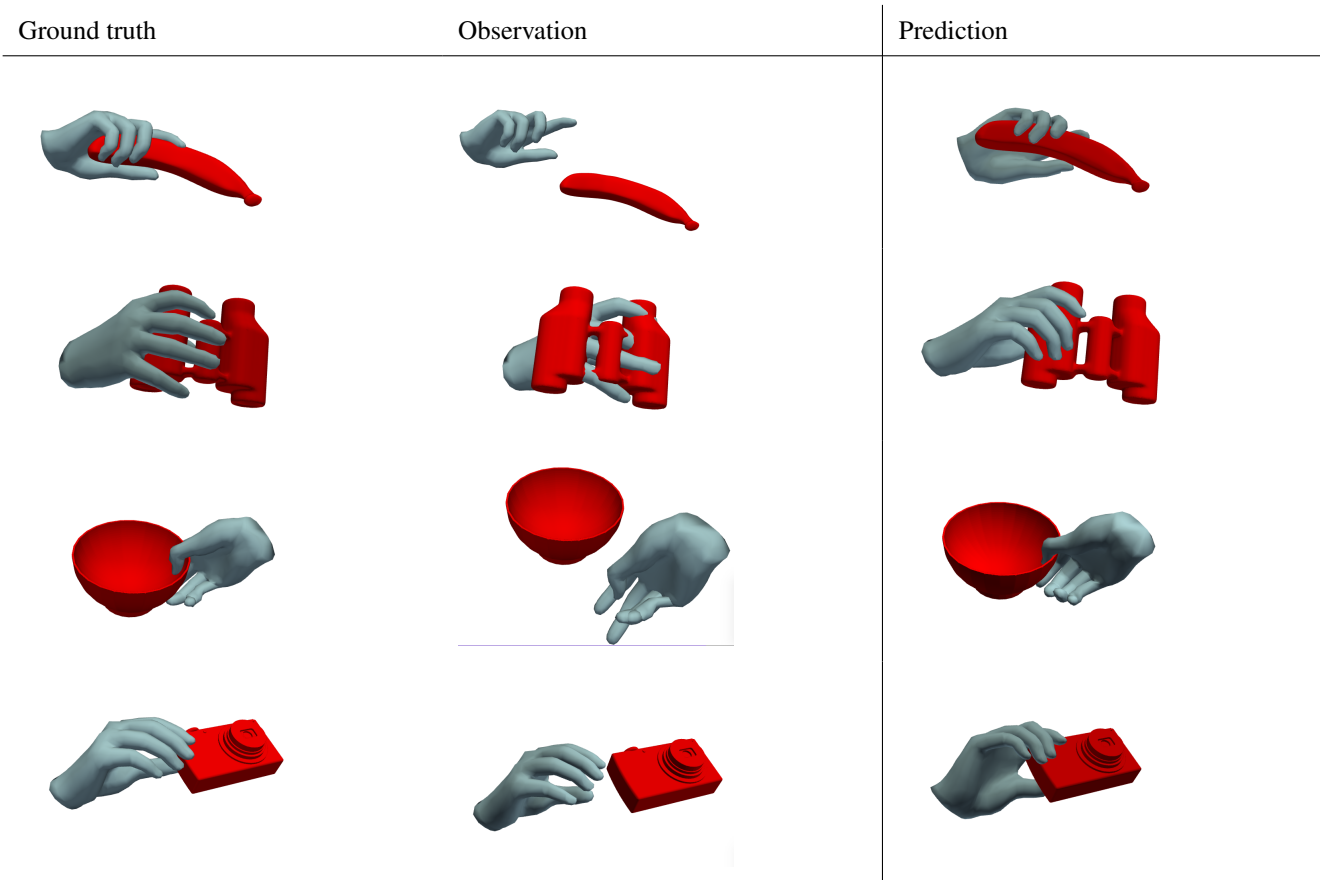
Figure 14. Qualitative evaluation of our multimodal *JointDiffusion*, trained on both object and noisy hand-object pair modalities, in the grasp refinement setting. *JointDiffusion* produces plausible grasps with minimal mesh penetration and respects finger contacts from the ground-truth mesh.