

ScoreHOI: Physically Plausible Reconstruction of Human-Object Interaction via Score-Guided Diffusion

Supplementary Material

In this appendix, we describe the DDIM sampling loop in our methods in pseudo-code in Section A. We also provide additional detailed implementations and qualitative comparisons in Section B and Section D. Furthermore, we append rendered 3D object models and the source code along with this appendix.

A. DDIM Refinement Loop

For a deeper understanding of the DDIM sampling loop, we illustrate a pseudo-code implementation in Algorithm A. The algorithm describe the DDIM_Loop in detail.

B. Detailed Implementations

We employ the PointNeXt [42] model, pre-trained on the ModelNet40 [55] dataset, as our affordance-aware regressor. ModelNet40 is a extensively utilized benchmark dataset for 3D shape classification and retrieval tasks, comprising 12,311 CAD models across 40 distinct object categories, including airplanes, chairs, tables, and cars. PointNeXt maintains the fundamental hierarchical architecture of PointNet++ [41] while integrating contemporary deep learning techniques to augment performance and efficiency. Through the incorporation of point cloud awareness, the model is capable of comprehending object geometry across a variety of human-object interaction scenarios. For the contact predictor and mesh regressor, we leverage the contact estimation transformer and the contact-based refinement methodology proposed by [37]. The contact predictor receives human and object feature tokens, along with estimated human and object mesh vertices, as inputs to generate contact masks through two symmetrical 4-layer transformer blocks. During our contact-driven iterative refinement process, we iteratively update the human and object meshes to enhance contact interactions, thereby refining the contact prediction results with each iteration. After N iterations of optimization, we obtain the x_0^N for final refinement. The mesh regressor, also constructed with two analogous 4-layer transformer branches, takes the updated human and object feature tokens and contact masks as input to further refine the mesh results.

C. Ablation Studies

About Diffusion Priors. We demonstrate the generation results starting from different noise levels without physical guidance, as shown in the Figure A. The IG-Adapter successfully constructs plausible human-object interaction pat-

Table A. Efficiency comparison with template-free methods.

	HDM	InterTrack	ScoreHOI	ScoreHOI-F
FPS \uparrow	0.0047	0.0012	0.2895	2.0080

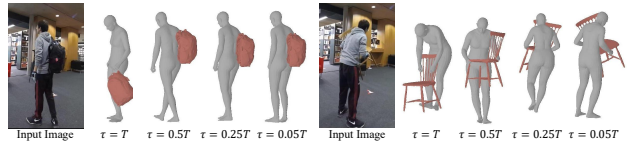


Figure A. Generation results starting from different noise levels.

terns, confirming that the diffusion model has acquired a valid prior distribution.

Comparison with template-free methods. We compared ScoreHOI with template-free methods, including HDM [57] and InterTrack [61]. As shown in the Table A, ScoreHOI surpasses these methods in efficiency. Additionally, these methods output point clouds, which are less practical for applications like robotic data collection due to the additional time and uncertainty introduced when fitting SMPL parameters.

D. More Qualitative Results

To further substantiate the performance of our ScoreHOI, we present additional qualitative results of human-object interaction reconstruction in Figure B. Our methodology exhibits superior performance across diverse interaction patterns, such as sitting, carrying, grasping, and lifting. By incorporating physical constraints, our ScoreHOI achieves a higher degree of accuracy and physical fidelity in the reconstruction of human and object meshes.

E. Geometry Model Demos

The demo 3D geometry mesh results are available under the **demo** directory. Reviewers are able to assess the performance from any perspective utilizing software such as MeshLab or Blender.

Algorithm A Score-Guided DDIM refinement loop

- 1: **Input:** latent parameters \mathbf{x}_0^n at step n , denoising model ϵ_ϕ , image features \mathbf{c}_I , geometry features \mathbf{c}_G , gradient step size ρ , noise level τ , DDIM step size Δt , estimated contact masks $\{\mathbf{M}_i\}_{i \in \{h,o,f\}}$
 - 2: **Output:** latent parameters \mathbf{x}_0^{n+1} for next sampling step $n + 1$
 - 3: $\mathbf{x}_\tau = \text{DDIMInvert}(\mathbf{x}_0^n, \mathbf{c}_I, \mathbf{c}_G)$ ▷ Run DDIM inversion until noise level τ
 - 4: **for** $t = \tau$ to Δt with step size Δt **do**
 - 5: $\tilde{\epsilon} \leftarrow \epsilon_\phi(\mathbf{x}_t^n, t, \mathbf{c}_I, \mathbf{c}_G)$ ▷ Predict noise
 - 6: Initialize computational graph for \mathbf{x}_t^n
 - 7: $\hat{\mathbf{x}}_0^n \leftarrow \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t^n - \sqrt{1 - \alpha_t})\tilde{\epsilon}$ ▷ Predict one-step denoised result
 - 8: $L_P \leftarrow \text{PhysicalGuidance}(\hat{\mathbf{x}}_0^n, \{\mathbf{M}_i\}_{i \in \{h,o,f\}})$ ▷ Compute physical guidance loss
 - 9: $\tilde{\epsilon}' \leftarrow \tilde{\epsilon} + \rho\sqrt{1 - \alpha_t}\nabla_{\mathbf{x}_t^n} L_P$ ▷ Compute modified noise after score-guidance
 - 10: $\hat{\mathbf{x}}_0^{n'} \leftarrow \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t^n - \sqrt{1 - \alpha_t})\tilde{\epsilon}'$ ▷ Predict one-step denoised result with modified noise
 - 11: $\mathbf{x}_{t-\Delta t}^n \leftarrow \sqrt{\alpha_{t-\Delta t}}\hat{\mathbf{x}}_0^{n'} + \sqrt{1 - \alpha_{t-\Delta t}}\tilde{\epsilon}'$ ▷ DDIM sampling step
 - 12: **end for**
 - 13: $\mathbf{x}_0^{n+1} \leftarrow \hat{\mathbf{x}}_0^{n'}$ ▷ Update \mathbf{x}_0^n for next generation
 - 14: **return** \mathbf{x}_0^{n+1}
-



Figure B. **Extra Qualitative comparisons.** We highlight the contact interaction within each picture.