

# URDFormer: Constructing interactive Realistic Scenes from Real Images via Simulation and Generative Modeling

Anonymous Author(s)  
Affiliation  
Address  
email

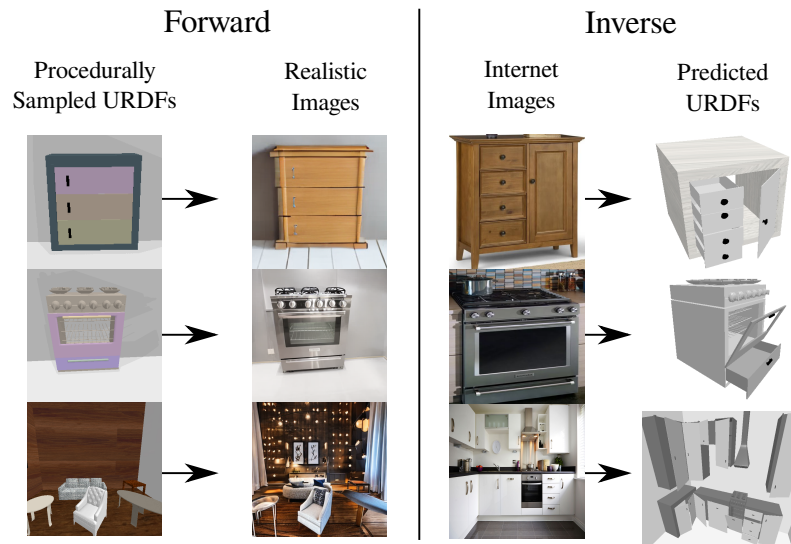


Figure 1: Our method uses generative models in a “Forward” process to produce structurally consistent realistic images from procedurally generated simulation content. We then use these generated simulation/image pairs to train an “Inverse” process that is able to estimate the underlying structure of diverse real-world images.

1       **Abstract:** Constructing accurate and targeted simulation scenes that are both vi-  
2       sually and physically realistic is a significant practical interest in domains rang-  
3       ing from robotics to computer vision. However, this process is typically done  
4       largely by hand - a graphic designer and a simulation engineer work together with  
5       predefined assets to construct rich scenes with realistic dynamic and kinematic  
6       properties. While this may scale to small numbers of scenes, to achieve the gener-  
7       alization properties that are requisite of data-driven machine learning algorithms,  
8       we require a pipeline that is able to synthesize large numbers of realistic scenes,  
9       complete with “natural” kinematic and dynamic structure. To do so, we develop  
10      models for inferring structure and generating simulation scenes from natural im-  
11      ages, allowing for scalable scene generation from web-scale datasets. To train  
12      these image-to-simulation models, we show how effective generative models can  
13      be used in generating training data, the network can be *inverted* to map from real-  
14      istic images back to complete scene models. We show how this paradigm allows  
15      us to build large datasets of scenes with semantic and physical realism, enabling a  
16      variety of downstream applications in robotics and computer vision. More visual-  
17      izations are available at: <https://sites.google.com/view/urdformer/home>

18       **Keywords:** Generative Modeling, Scene Generation, Simulation

## 19 1 Introduction

20 Simulation offers the dual advantages of scalable and cheap data collection and an easy way to  
21 encode domain-specific prior knowledge into end-to-end machine-learning problems [1, 2, 3]. This  
22 is particularly important for data-scarce problems such as robotics, where collecting real data can  
23 lead to costly and unsafe failures or may require expensive human supervision. Critical to each of  
24 these endeavors is a rich and accurate simulation environment, complete with complex scene layouts  
25 and kinematic structure. The de-facto process for generating simulation content is either manual  
26 [4] or procedural [5]. The manual process for creating simulation scenes involves the algorithm  
27 designer working to characterize, identify, and model a particular real-world scene, a painstaking  
28 and impractical process. This leads to content that is not very diverse due to the onerous human  
29 effort required. On the other hand, rule-based procedural generation methods [5, 6] have seen  
30 success in particular machine learning applications such as embodied navigation, but often struggle  
31 to capture the natural complexity of the real world. Moreover, the procedural generation process  
32 is not controllable, making it hard to generate simulation content corresponding to a *particular*  
33 real-world environment. The inability of the current status quo in the generation of simulation  
34 content - both procedural generation and manual creation, makes apparent the necessity of a targeted  
35 technique for scalable content creation in simulation that is able to retain realistic kinematic and  
36 semantic structure.

37 To enable a variety of downstream use cases, scalable content creation in simulation must be (1)  
38 realistic enough such that inferences made in simulation transfer back to the real world (2) diverse  
39 in a way that captures natural statistics so as to enable learning generalizable models and policies  
40 (3) controllable in a way that allows for targeted generation of particular scenes of interest. While  
41 a variety of methods for scene generation and inverse graphics [7, 8, 9] satisfy one or more of  
42 these criteria, to the best of our knowledge, it has proven challenging to develop content creation  
43 methods that satisfy them all. We develop methods that map directly from isolated real-world images  
44 to corresponding simulation content (expressed as a Unified Robot Description File (URDF)) that  
45 could plausibly represent the semantics, kinematics, and structure of the scene. This is an inverse  
46 mapping problem going from real-world images to kinematically accurate, interactive simulation.  
47 While inverse modeling problems in the literature have been tackled with data-driven techniques  
48 such as supervised learning, in this case, a large-scale paired dataset of realistic images and their  
49 corresponding simulation environments does not readily exist in the literature.

50 Our key idea is that we can generate a suitable dataset for inverse modeling from images to plausible  
51 simulations by leveraging controllable text-to-image generative models [10]. From a set of proce-  
52 durally or manually constructed scenes, we can generate realistic images that are representative of  
53 that particular simulation scene. This paired dataset of simulation scenes and corresponding realistic  
54 images can then be *inverted* via supervised learning to learn a model that maps from realistic images  
55 directly to plausible simulation environments. This learned model can generate realistic and diverse  
56 content directly from real-world images mined from the web without any additional annotation. The  
57 resulting models can be used in several use cases - (1) diverse generation: generating a large and  
58 diverse set of realistic simulation environments that correspond directly to real-world images, or (2)  
59 targeted generation: generating a simulation environment (or narrow distribution of environments)  
60 corresponding to a particular set of desired images.

## 61 2 Related Work

62 **Inverse-Graphics:** A variety of work in inverse-graphics focuses on inferring scene properties  
63 such as geometry, lighting, and other geometric properties from single images [11]. This work has  
64 both been optimization-based [12] and learning-based [13]. In a similar vein, a rich body of work  
65 [14] focuses on mesh reconstruction and novel view synthesis using a variety of techniques such as  
66 implicit neural fields [15, 16, 17], Gaussian splatting [18, 19], differentiable rendering [20, 21, 22]  
67 amongst many other techniques. Importantly, the focus of many of these works on inverse graphics  
68 has been on geometric reconstruction rather than our focus on scene-level simulation construction



69 complete with kinematic and semantic structure like object relationships and articulation. There  
 70 have been a number of efforts in inferring physical properties such as articulation [23, 24, 25],  
 71 friction and surface properties [26, 27, 28, 29], although these typically require either interaction  
 72 or video access. In contrast, our work focuses less on exact geometry reconstruction but rather on  
 73 generating correct scene statistics at the articulation/kinematics/positioning level for entire scenes  
 74 or complex objects from single RGB images. Our goal is a fast generation process that can scale  
 75 to generate hundreds of scenes with natural statistics, without requiring interaction or targeted data  
 76 collection per domain.

77 Generating indoor scenes is a  
 78 long-standing problem in com-  
 79 puter vision and machine learn-  
 80 ing. This has been approached  
 81 by building learned generative  
 82 models of indoor scenes [30,  
 83 31, 32, 33] and floorplans [34,  
 84 35, 36], while others have pro-  
 85 duced text-to-scene models [37,  
 86 38]. While generating scenes  
 87 this way can be promising, these

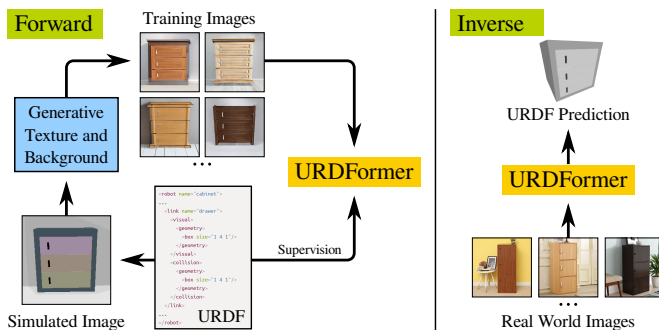


Figure 2: An overview of the training and application of URDFormer . During the forward process, existing simulation assets are first used to generate a large paired dataset of simulation assets and realistic rendered images. This paired dataset is used to train the URDFormer inverse model that can predict URDFs from RGB images. This model can then be used with real-world images to generate novel simulations.

88 methods either fail to achieve  
 89 the targeted generation of com-  
 90 plex scenes with articulation and  
 91 complex kinematic structure in-  
 92 tact or require extremely expen-  
 93 sive inference processes to do  
 94 so. On the other hand, procedural generation techniques have been popular in generating grid-world  
 95 environments [39, 40, 41, 42] and in generating home environments at scale [5]. These scenes are  
 96 diverse but are not controllable to particular target scenes or complete with physical properties and  
 97 articulation. Other techniques such as [43, 44] are able to generate large datasets of interactive  
 98 scenes but require interactive scanning with either a phone or other hardware for dataset generation  
 99 specific to indoor scenes. URDFormer is able to generate realistic, diverse, and controllable scenes  
 100 while retaining rich kinematic and semantic structure from internet images alone.

101 **Data Augmentation with Generative Models** Our work is certainly not the first [45] to use syn-  
 102 thetic data generated by generative models for training networks that can then be deployed on real  
 103 data. These models have been used the context of data augmentation [46, 47, 48], representation  
 104 learning via self supervised learning [49, 50, 51], model selection [52] and even applications like  
 105 healthcare [53]. In contrast, our work shows that controllable generative modeling can be used to  
 106 generate datasets that are suitable for inverse modeling for creating simulation assets at scale.

### 107 3 URDFormer : Generating Interactive Simulation Environments by 108 Learning Inverse Models from Generated Datasets

#### 109 3.1 Controlled Generation of Paired Datasets with Generative Models

110 Given a simulated scene  $z$  (drawn from a dataset such as [54], or procedurally generated), we use the  
 111 fact that controllable generative models are both diverse and realistic enough to take an unrealistic  
 112 rendering of a scene in simulation and generate a distribution of corresponding *realistic* images. This  
 113 allows the scene in simulation with unrealistic appearance and texture to be translated into a diverse  
 114 set of visually realistic images that plausibly match the same underlying environment. To ensure  
 115 piecewise consistency and realism of the generated images, we use two different dataset generation  
 116 techniques for the scene structure and object structure respectively. These share the same conceptual  
 117 ideas but differ to account for consistency properties in each case.

118 **Scene-Level Dataset Generation:** To generate training data for the scene model, we feed the ren-  
 119 dered image from simulation along with a templated text prompt to an image-and-text guided dif-  
 120 fusion model [10]. This model generates a new image that attempts to simultaneously match the  
 121 content described in the text prompt while retaining the global scene layout from the provided im-  
 122 age. We found that this model is able to reliably maintain the scene layout, but it may change some  
 123 individual components of the scene, for example replacing objects with a different but plausible  
 124 category, or changing the number of components under an object such as the drawers or handles.  
 125 Despite these failures, the large-scale structural consistency still provides a useful source of training  
 126 data. After running our simulated image through the generated model, we have realistic images  
 127 that contain known high-level object positions and spatial relationships, but unknown category and  
 128 low-level part structures. This means that the scene model dataset contains complete images, but  
 129 incomplete labels. Rather than complete  $(x, z)$  pairs, we have a dataset  $\mathcal{D}_{\text{scene}} = \{(x, \tilde{z})\}$  of  $(x, \tilde{z})$   
 130 pairs where  $\tilde{z}$  only contains the bounding boxes, transforms and parents of the high-level (non-part)  
 131 objects  $\tilde{z} = \{(b_1, T_1, p_1) \dots (b_n, T_n, p_n)\}$ .

132 **Object-Level Dataset Generation:** The process for generating object-level training data is similar,  
 133 but requires more care due to the tendency of generative models to modify low-level details. For  
 134 objects with complex kinematic structure, such as cabinets, we procedurally generate a large number  
 135 of examples of these objects and render them in isolation from different angles. Rather than using  
 136 the generative model to construct entirely new images, we use it to produce diverse texture images,  
 137 which are overlaid in the appropriate locations on the image using perspective warping. We then  
 138 change the background of the image using the generative model with appropriate masking derived  
 139 from the original render. For less complex objects that do not have important part-wise structure, we  
 140 simply replace the rendered image with a new sample from the image-and-text guided generative  
 141 model. Unlike the scene dataset which contains complete images but partial labels, the object dataset  
 142 contains partial images in the sense that they contain only a single object, but complete labels for  
 143 the object and its kinematic parts. We can say that this dataset  $\mathcal{D}_{\text{object}}$  contains  $(\tilde{x}, z)$  pairs where  $\tilde{x}$   
 144 is an image of a single object rather than a full scene (hence the partial  $x$ ), and  $z$  is complete for the  
 145 single object and its parts. The result of these two data generation processes is a high-level scene  
 146 structure dataset  $\mathcal{D}_{\text{scene}}$ , and a low-level object dataset  $\mathcal{D}_{\text{object}}$ .

### 147 3.2 URDFormer : Learning Inverse Generative Models for Scene Synthesis

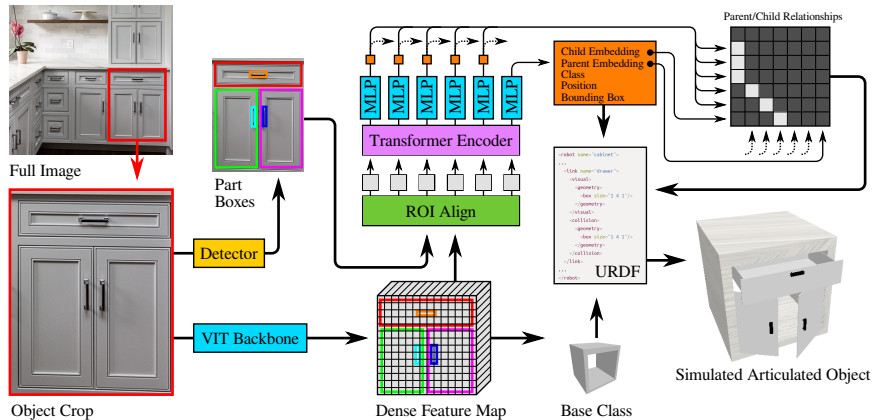


Figure 3: Architecture of URDFormer : an inverse model (URDFormer) that predicts simulation parameters from RGB images. URDFormer can translate web-scraped real-world RGB images of scenes into complete simulation assets. The model shown here is used to estimate the part structure of an individual object. When estimating the scene structure, the Object Crop image would be replaced by an image of the entire scene.

148 Given the datasets  $\mathcal{D}_{\text{object}} = (\tilde{x}, z)$  and  $\mathcal{D}_{\text{scene}} = (x, \tilde{z})$  constructed as described above, we can use  
 149 supervised learning methods to learn an *inverse model* that maps images of a complex object or  
 150 scene to the corresponding simulation asset. In order to take advantage of these partially complete  
 151 datasets, we must add some structure to our prediction model. We do this by splitting our learned

152 inverse model in correspondence with the split in our forward model: we train one network  $f_{\theta}^{-1}$  to  
153 predict the high-level scene structure using dataset  $\mathcal{D}_{\text{scene}}$  and another network  $g_{\phi}^{-1}$  to predict the  
154 low-level part structure of objects using  $\mathcal{D}_{\text{object}}$ .

155 To model both the scene-level prediction model ( $f_{\theta}^{-1}$ ) and the low-level part prediction model ( $g_{\phi}^{-1}$ ),  
156 we propose a novel network architecture - URDFFormer, that takes an RGB image and predicts URDF  
157 primitives as shown in Figure 3. Note that both the scene-level prediction and the low-level part  
158 prediction use the same network architecture, the scene-level simply operates on full images with  
159 object components segmented, while the part-level operates on crops of particular objects with parts  
160 segmented. In the URDFFormer architecture, the image is first fed into a ViT visual backbone[55]  
161 to extract global features. We then obtain bounding boxes of the objects in the image using the  
162 masks rendered from the original procedurally generated scene in simulation (these are known at  
163 training time, and can be easily extracted using segmentation models at test time). We then use  
164 ROI alignment [56] to extract features for each of these bounding boxes. These feature maps are  
165 combined with an embedding of the bounding box coordinates and then fed through a transformer  
166 [57] to produce a feature for each object in the scene. An MLP then decodes these features into an  
167 optional class label (used only when training the object-level model), and a discretized 3D position  
168 and bounding box. In addition, it also produces a child embedding and a parent embedding that  
169 are used to predict the hierarchical relationships in the scene (object to its parent and so on). To  
170 construct these relationships, the network uses a technique from scene graph generation [58] that  
171 produces an  $n \times n$  relationship score matrix by computing the dot product of every possible parent  
172 with every possible child. The scene model also has learned embeddings for six different root objects  
173 corresponding to the four walls, the floor, and the ceiling so that large objects like countertops and  
174 sinks can be attached to the room.

175 Due to the unpredictable nature of the generative transforms that are used to make the scene im-  
176 age realistic, which may change class identities, only the position, bounding box, and relationship  
177 information are used when computing the high-level scene structure. To compute the class labels  
178 for the top-level objects, we use max-pooling of the dense ViT features along with an MLP in the  
179 part-prediction model  $g_{\phi}^{-1}$ . To generate a full estimate of the scene description from a natural image  
180 at test time, the image and a list of high-level bounding boxes are first fed to the scene prediction  
181 model  $f_{\theta}^{-1}$ , which predicts the location and parent for each object. The image regions correspond-  
182 ing to these boxes are then extracted and further segmented to produce part-level bounding boxes.  
183 Each of these image regions and the corresponding part boxes are then fed into the part prediction  
184 model to compute the kinematic structure of the low-level parts. This nested prediction structure  
185 can be used to generate entire scenes from web-scraped RGB images drawn from any image dataset  
186 to generate novel simulation content both at the scene level and at the object level.

## 187 4 Experiments

### 188 4.1 Phase 1: (Forward) Paired Dataset Generation

189 To synthesize the initial paired dataset, we first procedurally generate a set of URDF representations  
190 of scenes in simulation both for global scenes like kitchens and for single objects like fridges, cab-  
191 inets, and drawers. These initially generated simulation scenes are shown in Fig4 (Left). We can  
192 then follow the procedure outlined in Section 3.1 for the controlled generation of paired datasets  
193 to generate a large dataset of simulation scenes and paired realistic RGB images as shown in Fig4  
194 (Right) (More visualizations and videos are available on the website). For objects with diverse  
195 parts, we observe that depth-guided stable diffusion [10] often ignores the semantic details of local  
196 parts, leading to inconsistency issues shown as Fig 7 in Appendix A.1. To overcome this issue, we  
197 use images of texture to guide diffusion models to generate large and diverse texture templates and  
198 randomly choose one template and warp it back to the original part region using perspective trans-  
199 formation. We apply in-painting models for smoothing the boundary of the parts and generating  
200 content for the background. We visualize this process in A.1 Fig 6. In total, we generated 260K



Figure 4: Qualitative results on (forward) paired dataset generation. Left: Original simulation images. Right: Generated realistic images that match the URDF descriptions of the scene on the left.

201 images for global scenes of kitchens and living rooms, and 235K images of 14 types of objects such  
 202 as cabinets, ovens, and fridges. Details of the dataset can be found in Appendix B.1.

#### 203 4.2 Phase 2: (Inverse) Real-World URDF Prediction

204 Given the generated paired dataset shown in Fig 4, we next evaluate how successful a trained inverse  
 205 model is at generating simulation scenes representing unseen real-world test images.

206 **Real World Dataset:** We create two types of datasets for evaluation: (a) Obj300 includes URDFs  
 207 of 300 internet images of individual objects from 5 categories including 100 cabinets, 50 ovens, 50  
 208 dishwashers, 50 fridges and 50 washing machines. (b) Global scenes include URDFs of 80 internet  
 209 images including 54 kitchens and 26 living rooms. For each scene, we manually label the bounding  
 210 box for each object and its parts, as well as the URDF primitives including mesh types, parent  
 211 bounding box ID, positions, and scales relative to its parent. We use the mesh types such as "left  
 212 door", and "right door" to infer link axis and joint types. All the position values and scale values are  
 213 discretized into 12 bins.

214 **Evaluation Metrics:** Evaluating entire scenes is challenging given the mixed structure and sub-  
 215 jective nature of human labelling. We adopt an edit-distance based metric for structural comparison,  
 216 and use a small dataset of manually labelled examples for evaluation.

217 (1) Edit Distance with Bounding Box Offset: We evaluate our predicted scene structure using a tree  
 218 edit-distance metric. This method requires access to a predicted and ground-truth kinematic tree.  
 219 We start at the root of the kinematic tree and use the Hungarian method to compute the lowest-cost  
 220 assignment between the children of the predicted root and the children of the ground truth root where  
 221 the cost is based on their spatial coordinates. If there are more predicted children than ground truth,  
 222 the unassigned predicted children and all of their descendants are marked as **False Positive** edits.  
 223 Conversely, if there are more ground truth children than predicted children, the unmatched ground  
 224 truth children and all of their descendants are marked as **False Negative** edits. We then compare the  
 225 spatial parameters of the matched predicted and ground truth children. If they are not close enough



Table 1: Comparison with baseline methods: trained with random colors, selected textures, and random textures, as well as prompt guided BLIP2. URDFormer with generated realistic textures predicts more accurate simulation content from unseen real-world images.

	Obj300 ( $\downarrow$ )	Global (Obj) ( $\downarrow$ )	Global (Parts) ( $\downarrow$ )
URDFormer (Random Colors)	1.08	10.81	19.62
URDFormer (Selected Textures)	0.63	9.87	19.11
URDFormer (Random Textures)	1.22	11.85	18.67
Guided BLIP2	4.27	14.64	24.58
<b>URDFormer (Generated Textures (ours))</b>	<b>0.42</b>	<b>9.51</b>	<b>18.21</b>

226 to each other according to a fixed threshold, the predicted child and its descendants are marked as  
 227 **False Positives**, and the ground truth child and its descendants are marked as **False Negatives**. If  
 228 the two are close enough, the class label of the predicted child is compared against the class label of  
 229 the ground truth child. If they do not match, we add a **Class Incorrect** edit. Regardless of whether  
 230 the classes match, this process is recursively applied to the matching children. To compute a single  
 231 score, we assign weights to these edits based on their position in the hierarchy and sum them. For  
 232 the experiments in this paper, we assigned a weight of 1.0 to all edits at the top level corresponding  
 233 to objects, a weight of 0.5 to the parts such as cabinet doors, and a weight of 0.1 to all objects further  
 234 down the hierarchy such as handles and knobs attached to doors.

235 (2) Edit Distance with IoU: Similar to bounding box offset, we simply replace the spatial coordinate  
 236 cost with IoU between two bounding boxes. We define levels of threshold based on overlapping  
 237 areas: ED IoU<sub>0.25</sub>, ED IoU<sub>0.5</sub>, ED IoU<sub>0.75</sub>. We show evaluation using both metrics in ablations,  
 238 but in general, we found the two metrics yield the same performance, thus we only use edit distance  
 239 with a bounding box for baseline evaluation.

240 **Baselines** We compare URDFormer against several other baselines in Table 1. In particular, to show  
 241 the importance of pixel realism, we compare with training on (1) Random Colors (2) Selected Re-  
 242 alistic Textures (3) Random Textures (Visualizations of their differences are in Appendix A.3). In  
 243 addition, we also compare our method against recent Vision-Language Models with guided prompts:  
 244 Guided BLIP2. In particular, (1) Random Colors randomly selects RGB values for each part inside  
 245 the scene and (2) Selected Realistic Textures manually selects texture images for corresponding ob-  
 246 jects. (3) Random Textures selects random images. (4) Guided BLIP2 takes a sequence of question  
 247 prompts and guides pretrained BLIP2 models [59] to output the URDF primitives in the valid format  
 248 (Please check Appendix C.3 for prompt details). We observe that training with generated realistic  
 249 visual features improves the generalization to real-world images. Although trained on large real-  
 250 world datasets, BLIP2 fails to reason about the 3D structure of the scene as well as the kinematics  
 251 structure of individual objects, showing using a more structured and targeted dataset is important  
 252 during training. Here Global (Obj) represents the evaluation of high-level position/parent reasoning,  
 253 while Global (Parts) represents the evaluation of the full scene including the high-level and detailed  
 254 kinematic structure of each object.

255 **Ablations** To study how different components of URDFormer impact the performance, we perform  
 256 an ablation study on (1) Do backbones pretrained on real-world images help with generalization?  
 257 (2) What are the important features of learning 3D kinematic structures, as shown in Table 2. In par-  
 258 ticular, we train URDFormer with three types of backbones: (1) vit-small-patch16-224 trained from  
 259 scratch (2) finetune vit-small-patch16-224 pretrained on ImageNet (3) finetune vit-small-patch16-  
 260 224 trained in [60] on 197K kitchen scenes and evaluate on 54 real-world kitchen images. We  
 261 observe that finetuning the vision backbone that is pretrained on real images performs better than  
 262 training from scratch, and pretrained in [60] achieves the best performance, which is likely due to  
 263 the fact that it was trained on more diverse datasets than ImageNet. We observe that both training  
 264 with only image features and training with only bounding box features decrease the performance,  
 265 indicating the importance of both spatial and visual features.

266 **Qualitative Results:** We show the qualitative results of our URDF predictions in Fig 5. We use  
 267 the same color to represent the same mesh type for better visualization. We observe that training

Table 2: Ablation study on training with different vision backbones and input features, showing training using both visual/spatial features, with a backbone pretrained on diverse real images achieves higher performance.

	ED Box ( $\downarrow$ )	ED IoU <sub>0.25</sub> ( $\downarrow$ )	ED IoU <sub>0.5</sub> ( $\downarrow$ )	ED IoU <sub>0.75</sub> ( $\downarrow$ )
Scratch	7.00	6.15	8.37	14.48
Pretrained on ImageNet	6.33	5.48	7.74	13.85
<b>Pretrained MAE</b>	<b>5.70</b>	<b>5.11</b>	<b>7.07</b>	<b>13.41</b>
Pretrained MAE (No bbox)	6.19	5.26	7.63	14.11
only with bbox	7.04	6.52	8.26	14.26

268 with data generated using the method described in section 3.1 provides diverse visual information  
 269 compared to baseline methods such as random colors or random textures. This is important for  
 270 distinguishing mesh types such as stove and fridge, and reasoning about structure relations such as  
 271 "cabinet on the right" and "cabinet in the front".

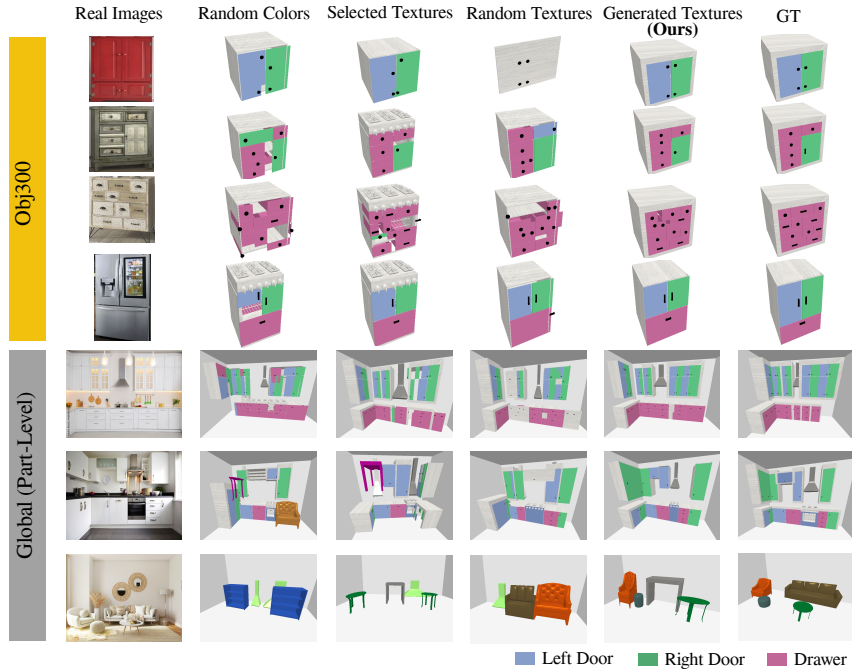


Figure 5: Evaluations of generated simulations on unseen real-world images. The left-most column indicates the real-world image input and each column indicates the performance of an inverse URDF prediction model trained with different training sets. We evaluate training datasets generated using random colors, selected textures, random textures, and textures generated with pre-trained generative models (ours), and compare these with ground truth URDF labels.

## 272 5 Conclusion

273 In this work, we presented URDFormer, a general-purpose, scalable technique for generating simulation  
 274 content from real-world RGB images. We first generate a large-scale paired dataset of procedurally  
 275 generated simulation content and a corresponding realistic RGB image using pre-trained  
 276 controllable generative models. We then use our generated paired dataset to train an inverse model  
 277 that maps directly from single RGB images to corresponding representations of scenes or complex  
 278 objects in simulation. This inverse model can then be used with large image datasets of real-world  
 279 RGB images to scalably generate simulation data complete with kinematic and semantic structure,  
 280 without requiring any hand-crafting or hand-designing of these simulation assets. We show in our  
 281 experimental results the efficacy of this scheme in generating assets at scale from real-world datasets  
 282 of RGB images.



## References

- 283
- 284 [1] J. Collins, S. Chand, A. Vanderkop, and D. Howard. A review of physics simulators for robotic  
285 applications. *IEEE Access*, 9:51416–51431, 2021. doi:10.1109/ACCESS.2021.3068769. URL  
286 <https://doi.org/10.1109/ACCESS.2021.3068769>.
- 287 [2] Y. S. Narang, K. Storey, I. Akinola, M. Macklin, P. Reist, L. Wawrzyniak, Y. Guo,  
288 Á. Moravánszky, G. State, M. Lu, A. Handa, and D. Fox. Factory: Fast contact for robotic as-  
289 sembly. In K. Hauser, D. A. Shell, and S. Huang, editors, *Robotics: Science and Systems XVIII*,  
290 *New York City, NY, USA, June 27 - July 1, 2022*, 2022. doi:10.15607/RSS.2022.XVIII.035.  
291 URL <https://doi.org/10.15607/RSS.2022.XVIII.035>.
- 292 [3] M. Müller, V. Casser, J. Lahoud, N. Smith, and B. Ghanem. Sim4cv: A photo-realistic sim-  
293 ulator for computer vision applications. *Int. J. Comput. Vis.*, 126(9):902–919, 2018. doi:  
294 [10.1007/s11263-018-1073-7](https://doi.org/10.1007/s11263-018-1073-7). URL <https://doi.org/10.1007/s11263-018-1073-7>.
- 295 [4] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, M. Deitke, K. Ehsani,  
296 D. Gordon, Y. Zhu, et al. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint*  
297 *arXiv:1712.05474*, 2017.
- 298 [5] M. Deitke, E. VanderBilt, A. Herrasti, L. Weihs, K. Ehsani, J. Salvador, W. Han, E. Kolve,  
299 A. Kembhavi, and R. Mottaghi. Proctor: Large-scale embodied ai using procedural genera-  
300 tion. *Advances in Neural Information Processing Systems*, 35:5982–5994, 2022.
- 301 [6] A. Raistrick, L. Lipson, Z. Ma, L. Mei, M. Wang, Y. Zuo, K. Kayan, H. Wen, B. Han, Y. Wang,  
302 A. Newell, H. Law, A. Goyal, K. Yang, and J. Deng. Infinite photorealistic worlds using  
303 procedural generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
304 *CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 12630–12641. IEEE, 2023. doi:  
305 [10.1109/CVPR52729.2023.01215](https://doi.org/10.1109/CVPR52729.2023.01215). URL [https://doi.org/10.1109/CVPR52729.2023.](https://doi.org/10.1109/CVPR52729.2023.01215)  
306 [01215](https://doi.org/10.1109/CVPR52729.2023.01215).
- 307 [7] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum. Deep convolutional inverse graph-  
308 ics network. *Advances in neural information processing systems*, 28, 2015.
- 309 [8] S. Lunz, Y. Li, A. W. Fitzgibbon, and N. Kushman. Inverse graphics GAN: learning to generate  
310 3d shapes from unstructured 2d data. *CoRR*, abs/2002.12674, 2020. URL [https://arxiv.](https://arxiv.org/abs/2002.12674)  
311 [org/abs/2002.12674](https://arxiv.org/abs/2002.12674).
- 312 [9] M. Jaques, M. Burke, and T. M. Hospedales. Physics-as-inverse-graphics: Unsupervised phys-  
313 ical parameter estimation from video. In *8th International Conference on Learning Representations*,  
314 *ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL  
315 <https://openreview.net/forum?id=BJeKwTNFvB>.
- 316 [10] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image syn-  
317 thesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer*  
318 *vision and pattern recognition*, pages 10684–10695, 2022.
- 319 [11] P. W. Battaglia, J. B. Hamrick, and J. B. Tenenbaum. Simulation as an engine of physical scene  
320 understanding. *Proceedings of the National Academy of Sciences*, 110(45):18327–18332,  
321 2013.
- 322 [12] S. Agarwal, Y. Furukawa, N. Snavely, I. Simon, B. Curless, S. M. Seitz, and R. Szeliski.  
323 Building rome in a day. *Communications of the ACM*, 54(10):105–112, 2011.
- 324 [13] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove. DeepSDF: Learning con-  
325 tinuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF*  
326 *conference on computer vision and pattern recognition*, pages 165–174, 2019.

- 327 [14] T. Samavati and M. Soryani. Deep learning-based 3d reconstruction: a survey. *Artif. Intell.*  
328 *Rev.*, 56(9):9175–9219, 2023. doi:10.1007/s10462-023-10399-2. URL [https://doi.org/](https://doi.org/10.1007/s10462-023-10399-2)  
329 [10.1007/s10462-023-10399-2](https://doi.org/10.1007/s10462-023-10399-2).
- 330 [15] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf:  
331 Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*,  
332 65(1):99–106, 2021.
- 333 [16] K. Park, U. Sinha, J. T. Barron, S. Bouaziz, D. Goldman, S. Seitz, and R. Martin-Brualla.  
334 Deformable neural radiance fields. <https://arxiv.org/abs/2011.12948>, 2020.
- 335 [17] K. Zhang, G. Riegler, N. Snavely, and V. Koltun. NERF++: Analyzing and improving neural  
336 radiance fields. <https://arxiv.org/abs/2010.07492>, 2020.
- 337 [18] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis. 3d gaussian splatting for real-time  
338 radiance field rendering. *ACM Transactions on Graphics (TOG)*, 42(4):1–14, 2023.
- 339 [19] J. Luiten, G. Kopanas, B. Leibe, and D. Ramanan. Dynamic 3d gaussians: Tracking by persis-  
340 tent dynamic view synthesis. *arXiv preprint arXiv:2308.09713*, 2023.
- 341 [20] T. H. Nguyen-Phuoc, C. Li, S. Balaban, and Y. Yang. Rendernet: A deep convolutional net-  
342 work for differentiable rendering from 3d shapes. *Advances in neural information processing*  
343 *systems*, 31, 2018.
- 344 [21] H. Kato, D. Beker, M. Morariu, T. Ando, T. Matsuoka, W. Kehl, and A. Gaidon. Differentiable  
345 rendering: A survey. *arXiv preprint arXiv:2006.12057*, 2020.
- 346 [22] H.-T. D. Liu, M. Tao, and A. Jacobson. Paparazzi: surface editing by way of multi-view image  
347 processing. *ACM Trans. Graph.*, 37(6):221–1, 2018.
- 348 [23] Z. Xu, J. Wu, A. Zeng, J. B. Tenenbaum, and S. Song. Densephysnet: Learning dense physical  
349 object representations via multi-step dynamic interactions. *arXiv preprint arXiv:1906.03853*,  
350 2019.
- 351 [24] Z. Xu, Z. He, and S. Song. Universal manipulation policy network for articulated objects.  
352 *IEEE Robotics and Automation Letters*, 7(2):2447–2454, 2022.
- 353 [25] B. DeMoss, P. Duckworth, N. Hawes, and I. Posner. Ditto: Offline imitation learning with  
354 world models. *arXiv preprint arXiv:2302.03086*, 2023.
- 355 [26] J. Wu, E. Lu, P. Kohli, B. Freeman, and J. Tenenbaum. Learning to see physics via visual  
356 de-animation. *Advances in Neural Information Processing Systems*, 30, 2017.
- 357 [27] L. S. Piloto, A. Weinstein, P. Battaglia, and M. Botvinick. Intuitive physics learning in a  
358 deep-learning model inspired by developmental psychology. *Nature human behaviour*, 6(9):  
359 1257–1267, 2022.
- 360 [28] J. R. Kubricht, K. J. Holyoak, and H. Lu. Intuitive physics: Current research and controversies.  
361 *Trends in cognitive sciences*, 21(10):749–759, 2017.
- 362 [29] F. de Avila Belbute-Peres, K. Smith, K. Allen, J. Tenenbaum, and J. Z. Kolter. End-to-end  
363 differentiable physics for learning and control. *Advances in neural information processing*  
364 *systems*, 31, 2018.
- 365 [30] D. Ritchie, K. Wang, and Y.-a. Lin. Fast and flexible indoor scene synthesis via deep convolu-  
366 tional generative models. In *Proceedings of the IEEE/CVF conference on computer vision and*  
367 *pattern recognition*, pages 6182–6190, 2019.
- 368 [31] M. Li, A. G. Patil, K. Xu, S. Chaudhuri, O. Khan, A. Shamir, C. Tu, B. Chen, D. Cohen-Or, and  
369 H. Zhang. Grains: Generative recursive autoencoders for indoor scenes. *ACM Transactions on*  
370 *Graphics (TOG)*, 38(2):1–16, 2019.

- 371 [32] M. Keshavarzi, A. Parikh, X. Zhai, M. Mao, L. Caldas, and A. Y. Yang. Scenegen: Generative  
372 contextual scene augmentation using scene graph priors. *arXiv preprint arXiv:2009.12395*,  
373 2020.
- 374 [33] D. A. Hudson and C. L. Zitnick. Generative adversarial transformers. *CoRR*, abs/2103.01209,  
375 2021. URL <https://arxiv.org/abs/2103.01209>.
- 376 [34] R. Hu, Z. Huang, Y. Tang, O. Van Kaick, H. Zhang, and H. Huang. Graph2plan: Learning  
377 floorplan generation from layout graphs. *ACM Transactions on Graphics (TOG)*, 39(4):118–1,  
378 2020.
- 379 [35] N. Nauata, S. Hosseini, K.-H. Chang, H. Chu, C.-Y. Cheng, and Y. Furukawa. House-gan++:  
380 Generative adversarial layout refinement network towards intelligent computational agent for  
381 professional architects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
382 *Pattern Recognition*, pages 13632–13641, 2021.
- 383 [36] K. Wang, X. Xu, L. Lei, S. Ling, N. Lindsay, A. X. Chang, M. Savva, and D. Ritchie. Roomi-  
384 noes: Generating novel 3d floor plans from existing 3d rooms. *Comput. Graph. Forum*, 40(5):  
385 57–69, 2021. doi:10.1111/cgf.14357. URL <https://doi.org/10.1111/cgf.14357>.
- 386 [37] A. Chang, M. Savva, and C. D. Manning. Learning spatial knowledge for text to 3d scene  
387 generation. In *Proceedings of the 2014 conference on empirical methods in natural language*  
388 *processing (EMNLP)*, pages 2028–2038, 2014.
- 389 [38] A. Chang, W. Monroe, M. Savva, C. Potts, and C. D. Manning. Text to 3d scene generation  
390 with rich lexical grounding. *arXiv preprint arXiv:1505.06289*, 2015.
- 391 [39] A. Khalifa, P. Bontrager, S. Earle, and J. Togelius. Pcgrl: Procedural content generation via  
392 reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence and*  
393 *Interactive Digital Entertainment*, volume 16, pages 95–101, 2020.
- 394 [40] S. Earle, M. Edwards, A. Khalifa, P. Bontrager, and J. Togelius. Learning controllable content  
395 generators. In *2021 IEEE Conference on Games (CoG)*, pages 1–9. IEEE, 2021.
- 396 [41] M. Dennis, N. Jaques, E. Vinitzky, A. Bayen, S. Russell, A. Critch, and S. Levine. Emergent  
397 complexity and zero-shot transfer via unsupervised environment design. *Advances in neural*  
398 *information processing systems*, 33:13049–13061, 2020.
- 399 [42] L. Gisslén, A. Eakins, C. Gordillo, J. Bergdahl, and K. Tollmar. Adversarial reinforcement  
400 learning for procedural content generation. In *2021 IEEE Conference on Games (CoG), Copen-*  
401 *hagen, Denmark, August 17-20, 2021*, pages 1–8. IEEE, 2021. doi:10.1109/CoG52621.2021.  
402 9619053. URL <https://doi.org/10.1109/CoG52621.2021.9619053>.
- 403 [43] Z. Li, T. Yu, S. Sang, S. Wang, M. Song, Y. Liu, Y. Yeh, R. Zhu, N. B. Gun-  
404 davarapu, J. Shi, S. Bi, H. Yu, Z. Xu, K. Sunkavalli, M. Hasan, R. Ramamoor-  
405 thi, and M. Chandraker. Openrooms: An open framework for photorealistic in-  
406 door scene datasets. In *IEEE Conference on Computer Vision and Pattern Recog-*  
407 *niton, CVPR 2021, virtual, June 19-25, 2021*, pages 7190–7199. Computer Vision  
408 Foundation / IEEE, 2021. doi:10.1109/CVPR46437.2021.00711. URL [https://openaccess.thecvf.com/content/CVPR2021/html/Li\\_OpenRooms\\_An\\_Open\\_](https://openaccess.thecvf.com/content/CVPR2021/html/Li_OpenRooms_An_Open_Framework_for_Photorealistic_Indoor_Scene_Datasets_CVPR_2021_paper.html)  
409 [Framework\\_for\\_Photorealistic\\_Indoor\\_Scene\\_Datasets\\_CVPR\\_2021\\_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Li_OpenRooms_An_Open_Framework_for_Photorealistic_Indoor_Scene_Datasets_CVPR_2021_paper.html).  
410
- 411 [44] M. Deitke, R. Hendrix, A. Farhadi, K. Ehsani, and A. Kembhavi. Phone2proc: Bringing  
412 robust robots into our chaotic world. In *IEEE/CVF Conference on Computer Vision and Pat-*  
413 *tern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 9665–9675.  
414 IEEE, 2023. doi:10.1109/CVPR52729.2023.00932. URL [https://doi.org/10.1109/](https://doi.org/10.1109/CVPR52729.2023.00932)  
415 [CVPR52729.2023.00932](https://doi.org/10.1109/CVPR52729.2023.00932).

- 416 [45] P. Eigenschink, T. Reutterer, S. Vamosi, R. Vamosi, C. Sun, and K. Kalcher. Deep generative  
417 models for synthetic data: A survey. *IEEE Access*, 11:47304–47320, 2023. doi:10.1109/  
418 ACCESS.2023.3275134. URL <https://doi.org/10.1109/ACCESS.2023.3275134>.
- 419 [46] Z. Q. Chen, S. C. Kiami, A. Gupta, and V. Kumar. Genau: Retargeting behaviors to unseen  
420 situations via generative augmentation. In K. E. Bekris, K. Hauser, S. L. Herbert, and J. Yu,  
421 editors, *Robotics: Science and Systems XIX, Daegu, Republic of Korea, July 10-14, 2023*,  
422 2023. doi:10.15607/RSS.2023.XIX.010. URL <https://doi.org/10.15607/RSS.2023.XIX.010>.
- 424 [47] T. Yu, T. Xiao, J. Tompson, A. Stone, S. Wang, A. Brohan, J. Singh, C. Tan, D. M. J. Per-  
425 alta, K. Hausman, B. Ichter, and F. Xia. Scaling robot learning with semantically imagined  
426 experience. In K. E. Bekris, K. Hauser, S. L. Herbert, and J. Yu, editors, *Robotics: Science  
427 and Systems XIX, Daegu, Republic of Korea, July 10-14, 2023*, 2023. doi:10.15607/RSS.2023.  
428 XIX.027. URL <https://doi.org/10.15607/RSS.2023.XIX.027>.
- 429 [48] B. Trabucco, K. Doherty, M. Gurinas, and R. Salakhutdinov. Effective data augmentation  
430 with diffusion models. *CoRR*, abs/2302.07944, 2023. doi:10.48550/arXiv.2302.07944. URL  
431 <https://doi.org/10.48550/arXiv.2302.07944>.
- 432 [49] S. Fu, N. Tamir, S. Sundaram, L. Chai, R. Zhang, T. Dekel, and P. Isola. Dreamsim: Learning  
433 new dimensions of human visual similarity using synthetic data. *CoRR*, abs/2306.09344, 2023.  
434 doi:10.48550/arXiv.2306.09344. URL <https://doi.org/10.48550/arXiv.2306.09344>.
- 435 [50] Y. Tian, L. Fan, P. Isola, H. Chang, and D. Krishnan. Stablerep: Synthetic images from text-  
436 to-image models make strong visual representation learners. *CoRR*, abs/2306.00984, 2023.  
437 doi:10.48550/arXiv.2306.00984. URL <https://doi.org/10.48550/arXiv.2306.00984>.
- 438 [51] A. Jahanian, X. Puig, Y. Tian, and P. Isola. Generative models as a data source for mul-  
439 terview representation learning. In *The Tenth International Conference on Learning Rep-  
440 resentations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL  
441 <https://openreview.net/forum?id=qhAeZjs7dCL>.
- 442 [52] A. Shoshan, N. Bhonker, I. Kviatkovsky, M. Fintz, and G. G. Medioni. Synthetic data for  
443 model selection. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett,  
444 editors, *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Hon-  
445 olulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 31633–  
446 31656. PMLR, 2023. URL <https://proceedings.mlr.press/v202/shoshan23a.html>.
- 447 [53] E. Choi, S. Biswal, B. A. Malin, J. Duke, W. F. Stewart, and J. Sun. Generating multi-label  
448 discrete patient records using generative adversarial networks. In F. Doshi-Velez, J. Fackler,  
449 D. C. Kale, R. Ranganath, B. C. Wallace, and J. Wiens, editors, *Proceedings of the Machine  
450 Learning for Health Care Conference, MLHC 2017, Boston, Massachusetts, USA, 18-19 Au-  
451 gust 2017*, volume 68 of *Proceedings of Machine Learning Research*, pages 286–305. PMLR,  
452 2017. URL <http://proceedings.mlr.press/v68/choi17a.html>.
- 453 [54] K. Mo, S. Zhu, A. X. Chang, L. Yi, S. Tripathi, L. J. Guibas, and H. Su. Partnet: A  
454 large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding.  
455 In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach,  
456 CA, USA, June 16-20, 2019*, pages 909–918. Computer Vision Foundation / IEEE, 2019.  
457 doi:10.1109/CVPR.2019.00100. URL [http://openaccess.thecvf.com/content\\_CVPR\\_2019/html/Mo\\_PartNet\\_A\\_Large-Scale\\_Benchmark\\_for\\_Fine-Grained\\_and\\_Hierarchical\\_Part-Level\\_3D\\_CVPR\\_2019\\_paper.html](http://openaccess.thecvf.com/content_CVPR_2019/html/Mo_PartNet_A_Large-Scale_Benchmark_for_Fine-Grained_and_Hierarchical_Part-Level_3D_CVPR_2019_paper.html).
- 460 [55] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. De-  
461 ghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transform-  
462 ers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

- 463 [56] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE*  
464 *international conference on computer vision*, pages 2961–2969, 2017.
- 465 [57] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polo-  
466 sukhin. Attention is all you need. *Advances in neural information processing systems*, 30,  
467 2017.
- 468 [58] J. Yang, W. Peng, X. Li, Z. Guo, L. Chen, B. Li, Z. Ma, K. Zhou, W. Zhang, C. C. Loy, and  
469 Z. Liu. Panoptic video scene graph generation. In *CVPR*, 2023.
- 470 [59] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with  
471 frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- 472 [60] I. Radosavovic, T. Xiao, S. James, P. Abbeel, J. Malik, and T. Darrell. Real-world robot  
473 learning with masked visual pre-training. In *Conference on Robot Learning*, pages 416–426.  
474 PMLR, 2023.

475 **Appendix**

476 **A (Forward) Data Generation**

477 **A.1 Part-Consistency**

478 We compare our part-wise generation method with other approaches qualitatively in Fig 7. In particular, we observe that depth-guided or in-painting stable diffusion models [10] often ignore local consistency, making it difficult to render high-quality images that are paired with the simulation content. Instead, we only use stable diffusion to change the style of texture and warp the original texture to each part of the object, as shown in Fig6 and the masks of each part can be obtained directly from the simulator.

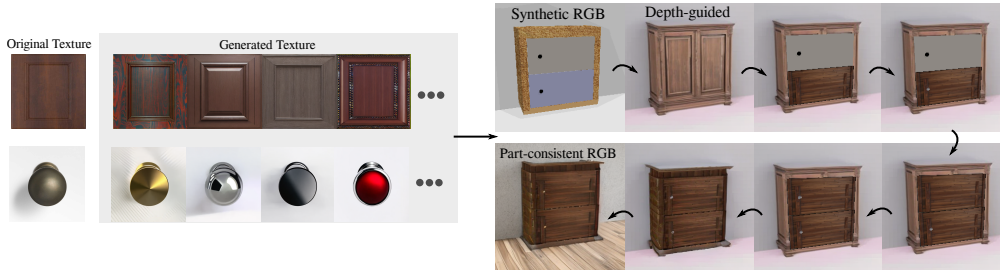


Figure 6: Paired dataset generation using texture and prompt templates to guide Stable Diffusion [10] and create a diverse texture dataset, which can be then warped on the targeted individual part of the object, as described in Section 3.1

483



Figure 7: Qualitative comparison among different rendering methods: depth-guided diffusion models, inpainting stable diffusion and part-wise generation

484 **A.2 Dataset Generation**

485 Qualitative details of our dataset of articulated, rigid objects as well as the full scenes are shown in Fig8, the top row in each section represents the original synthetic images from the simulation, the bottom images are the generated pair images that match the original kinematic structures.

488 **A.3 Baseline Data**

489 We visualize the different training data for baseline methods shown in Table 1: URDFormer with random colors, selected textures, random textures, and generated textures. All baseline inputs are



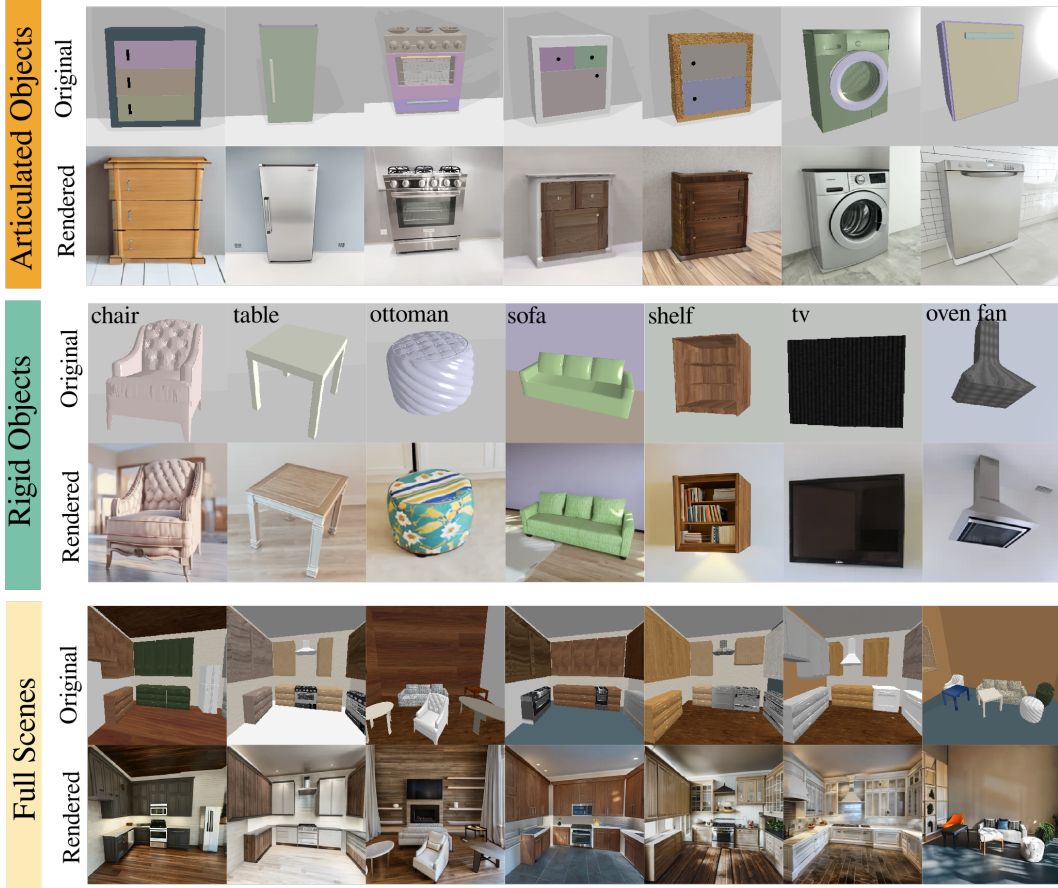


Figure 8: Data Generation for articulated objects, rigid objects, and full scenes. The top row in each section represents the original synthetic images from the simulation, the bottom images are the generated pair images that match the original kinematic structures.

491 captured from the same camera angles. As shown in Fig 9, the generated texture shows high pixel  
 492 realism that is closer to the distribution of the real world. As shown in Table 1, training on such data  
 493 improves performance in predicting URDF structures from real-world images during the test time.

## 494 B Training Details of URDFormer

### 495 B.1 Dataset

496 Our training dataset includes 267K global scene labels (197K kitchen scenes and 70K living room  
 497 scenes) and 235K objects, which include 14 types of objects including cabinet, oven, dishwasher,  
 498 washer, fridge, oven fan, shelf, tv, sofa, chair, square table, ottoman, coffee table and stuffed toy.  
 499 Among these objects, 5 categories are articulated: cabinet, oven, dishwasher, washer, and fridge.  
 500 These articulated objects include part meshes in from 8 types: drawer, left door, right door, oven  
 501 door, down door, circle door, handle and knob.

### 502 B.2 Training Details

503 We use a pretrained vit-small-patch16-224 trained in [60] as the vision backbone, which outputs  
 504 the global image features dimensions of  $14 \times 14 \times 384$ . To predict the base type, the global features  
 505 are first max-pooled followed by a MLP to predict a class type over 14 object types. We then  
 506 perform ROI alignment on cropped features with bounding boxes of the objects or parts. In ROI  
 507 Alignment, we set the spatial scale=1 / 16 and the sampling ratio=2. The ROI size is set to 14.



Figure 9: Comparison among baseline methods with different training input: Random Colors, selected textures, random textures and generated textures. Generated textures shows photo-realism that closer to the real-world distribution.

508 The roi-aligned features are then fed into a 3-layer MLP followed by a norm layer. To compute  
 509 the positional encoder, we feed the bounding box coordinates into a 3-layer MLP as well as a norm  
 510 layer. These normalized roi features together with the normalized spatial features are summed as  
 511 the token features and feed into the transformer, which are then fed into MLPs to compute URDF  
 512 primitives: position start (relative to parent), position end (relative to parent), mesh type, and parent-  
 513 child relation matrix. Here instead of regressing to a position value, we treat it as a classification  
 514 problem, where we discretize the x,y, and z axis of the parent mesh into 12 bins. During training, the  
 515 maximum sequence length is set to 32, which means the maximum number of bounding boxes per  
 516 image is 32. All baseline methods (URDFFormer with random colors, selected textures, and random  
 517 textures) are trained on one A40 GPU with a batch size of 256. All baselines are trained with an  
 518 equal number of epochs and evaluated using the last checkpoint.

## 519 C Experiment Details

### 520 C.1 Details for Dataset Assets

521 We procedurally generate scenes using both rigid and articulated objects. In particular, we collected  
 522 9 categories of common rigid objects in the kitchen and living room and 5 categories of common  
 523 articulated objects for kitchens, and randomly rescale them during data generation.

### 524 C.2 Prompts for Data Generation

#### 525 Textures:

526 (1) material prompt: 'bright', 'colorful', 'modern', 'multicolor', 'fancy color', 'accent', 'glass',  
 527 'chestnut', 'Oakwood', 'Maplewood', 'Cherrywood', 'Birchwood', 'Walnut', 'Mahogany', 'Pine',  
 528 'Beech', 'Ash', 'Hickory', 'Teak', 'Rosewood', 'Alder', 'Cedar', 'Bamboo', 'Plywood', 'Acacia',  
 529 'Poplar', 'fir'.

530 (2) full texture prompt: "a material wooden panel texture, high resolution, 4k, photorealistic".

531 **Objects:** "A object name, nice detailed, fancy, photorealistic, inside a home, 4k, natural light"

#### 532 Full Scenes:

533 (1) style prompt: "bright", "warm", "modern", "mediterranean", "vintage", "contemporary", "tran-  
534 sitional".

535 (2) kitchen: "a high-resolution picture of a bright style kitchen, very pretty, very natural lighting,  
536 ultra-high resolution, 8k, 16k, natural light, photorealistic, realism."

537 (3) living room: "a high-resolution picture of a bright style living room, with sofa, chairs, tv, ot-  
538 toman, floor lamps, etc, very pretty, very natural lighting, ultra-high resolution, 8k, 16k, natural  
539 light, photorealistic, realism".

### 540 C.3 Prompts for BLIP2

541 In this section, we show examples of how we guide Vision-Language Models such as BLIP2 [59] to  
542 produce answers that can be converted into comparison results with ours.

543 **Global Parent Prompt:** "which of the wall is this object most likely on? choose one from 'floor',  
544 'ceiling', 'front wall', 'left wall' and 'right wall'"

545 **Object Base Prompt:** "what's the name of the object. choose one word from cabinet, oven, dish-  
546 washer, washer, fridge, oven fan, shelf, tv, sofa, chair, square table, ottoman, coffee table and stuffed  
547 toy."

548 **Object Position Prompt:** "This image has a width of 512 and height of 512, the object box coordi-  
549 nate x is at 215, if the object scale is from 0 to 12, where do you imagine putting this bounding box  
550 relative to the object along the length in the 3D space. Choose from an integer from 0 to 12"

### 551 C.4 Compare with Other Scene Generation Methods

552 We compare our pipeline with other methods of scene generation 10. In particular, we evaluate on  
553 (1) If the generated content follows the real-world structure (2) If the method works only on RGB  
554 images (3) if the method is fully automatic without human interaction with the scene (4) If it is  
555 scalable (5) if it can be applied to global scenes and (6) if the generated scenes are fully articulated.

	Real World Distribution	RGB	Fully Automatic	Scalable	Scene Layout	Articulated Objects
Ditto	✓	✗	✗	✗	✗	✓
Ditto in the house	✓	✗	✗	✗	✓	✓
ProThor	✗	N/A	✓	✓	✓	✗
Phone2Proc	✗	✗	✗	✗	✓	✗
Ours	✓	✓	✓	✓	✓	✓

Figure 10: Comparison against different approaches in scene generation: Ditto, Ditto in the house, ProcThor, Phone2Proc.