

# Gen2Sim: Scaling up Simulation with Generative Models for Robotic Skill Learning

Anonymous Author(s)

Affiliation

Address

email

1       **Abstract:** We propose *Generation to Simulation (Gen2Sim)*, a method for scaling  
2 up robot skill learning in simulation by automatically generating simulation 3D  
3 assets, scenes, task definitions, task decompositions and reward functions, cap-  
4 italizing over large pre-trained generative models of language and images. We  
5 propose methods for 3D simulation asset generation from lifting open-world 2D  
6 object images using image diffusion models and LLM queries for plausible ranges  
7 of physical parameters. We then chain-of-thought prompt LLMs to parse URDF  
8 files of generated and human-developed assets to generate task descriptions, task  
9 decomposition, and corresponding reward functions, based on the assets and scene  
10 affordances. We train reinforcement learning policies in the simulation environ-  
11 ments using our generated tasks supervised by the generated reward functions. We  
12 demonstrate successful policy learning for a number of long horizon tasks using  
13 *Gen2Sim*, without any human involvement. Our work contributes hundreds of  
14 simulated assets and tasks for articulated and novel 3D object assets, taking a step  
15 towards fully autonomous robotic manipulation skill acquisition in simulation.

16       **Keywords:** Policy Learning in Simulation, Manipulation, Generative Models

## 17 1 Introduction

18 Scaling up training data has been a driving force behind the recent revolutions in language modeling  
19 [1], visual understanding [2], speech recognition [3], image generation [4], to name a few. This begs  
20 the question: can we scale up robot data to enable a similar revolution in robotic skill learning? One  
21 way to scale robot data is in the real world, by having multiple robots self-explore [5] or by collecting  
22 kinesthetic demonstrations at scale, with proper instrumentation or crowd-sourcing [6]. This is a  
23 promising direction, but safety concerns and wear and tear of the robots might be an obstacle towards  
24 autonomous real-world exploration. Another way to scale robot data is in simulation, by scaling up  
25 simulated environments and tasks, training robot policies in simulation with reinforcement learning  
26 or trajectory optimization, and then transferring them to the real world [7]. Such sim2real paradigm  
27 has seen recent successes in robot locomotion [8, 9, 10], in hand manipulation [11, 12], acrobatic  
28 flight [13, 14], and deformable object manipulation [15, 16, 17]. However, these examples, though  
29 very important and exciting, are still fairly isolated.

30 A central bottleneck towards scaling up simulation environments and tasks is the laborious manual  
31 effort needed for developing the visuals and physics of assets, their arrangement and configurations,  
32 the development of task curricula, and reward functions or programmatic demonstrations. In industry,  
33 tremendous resources have been invested in developing simulators for autonomous vehicles [18],  
34 warehouse robots, articulated objects [19], home environments [20, 21, 22], etc., many of which are  
35 proprietary and not open-sourced. Given these considerations, an important question naturally arises:  
36 How can we minimize manual effort in simulation development for diverse robotic skill learning?

37 In this paper, we explore automating the creation of simulation environments, and the development of  
38 manipulation tasks and rewards, by exploiting the latest progress in large-scale pre-trained generative

Skills: ■ Pick ■ Place ■ Turn ■ Press ■ Open ■ Close ■ Long Horizon Skills

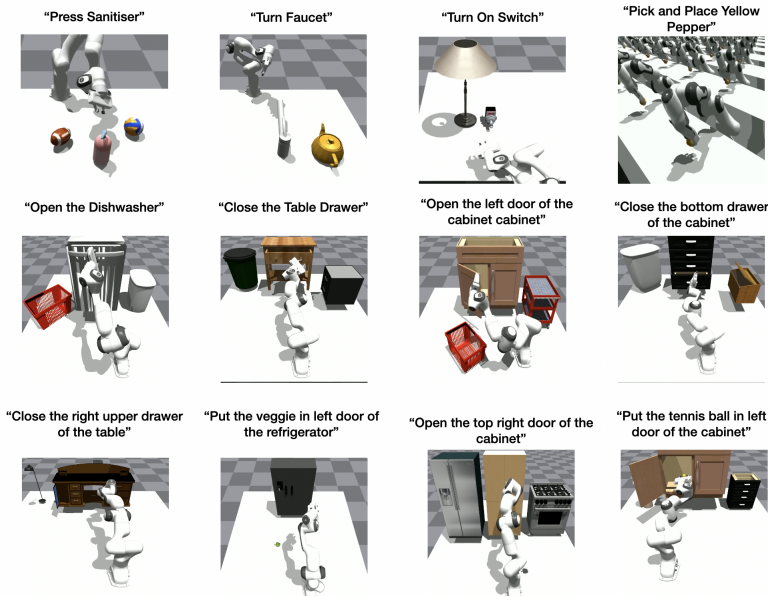


Figure 1: *Gen2Sim* is an automated generative pipeline for assets, textures, physical properties, tasks, task decompositions and corresponding rewards functions, aiming for autonomous robotic skill learning in simulation. Here we show 12 generated tasks, concerning the semantic affordances of the diverse types of assets in the scene.

39 models, aiming towards automated robotic skill learning. Our system, *Gen2Sim*, strives to automate  
 40 all stages involved in such development: from low-level 3D assets, textures, and physics properties,  
 41 to high-level task descriptions and reward functions, leading to automated skill learning in diverse  
 42 scenarios (See Figure 1). We automate 3D object asset generation by combining image diffusion  
 43 models for 3D mesh and texture generation, and LLMs for querying physical parameters information.  
 44 We showcase how LLMs and image generative models can diversify the appearances and behaviors of  
 45 assets by producing plausible ranges of textures, sizes and physical parameters, achieving “intelligent”  
 46 domain diversification. We automate task description, task decomposition and reward function  
 47 generation by few-shot prompting LLMs to predict language descriptions for semantically meaningful  
 48 tasks, concerning affordances of existing and generated 3D assets, articulated or not, alongside their  
 49 reward functions. *Gen2Sim* generates thousands of object assets and task variations without any  
 50 human involvement beyond several LLM prompt designs. We successfully train RL policies using  
 51 our auto-generated tasks and reward functions. Last, we demonstrate the usefulness of our simulation-  
 52 trained policies, by constructing a digital-twin environment of a given real scene, allowing a robot to  
 53 practice skills in the twin simulator and deploying it back to the real world to execute the task.

54 In summary, we make the following contributions:

- 55 • We show how pre-trained generative models of images and language can help automate 3D  
 56 asset generation and diversification, task description generation, task decomposition and  
 57 reward function generation that supports reinforcement learning of long horizon tasks in  
 58 simulation with minimal human involvement.
- 59 • We deploy our method to generate hundreds of assets, and hundreds of manipulation  
 60 tasks, their decompositions and their reward functions, for both human-developed and  
 61 automatically generated object assets.

62 Our code will be made publicly available upon publication. For videos and more qualitative results,  
 63 see our project site: <https://gen2sim.github.io/>.

## 64 2 Related Work

65 **Large Language Models for task and motion planning in robotics** Large language models (LLMs)  
66 map instructions to language subgoals [23, 24, 25, 26] or action programs [27] with appropriate  
67 plan-like or program-like prompts. LLMs trained from Internet-scale text have shown impressive  
68 zero-shot reasoning capabilities for a variety of downstream language tasks [1] when prompted  
69 appropriately, without any weight fine-tuning [28, 29, 30, 31]. LLMs were used to generate task  
70 curricula and predict skills to execute in Minecraft worlds [32, 33, 34] Following the seminal work of  
71 Code as Policies, many works map language to programs over given skills [35] or hand-designed  
72 motion planners [36]. Our work instead maps task descriptions into task decompositions and reward  
73 functions, to guide reinforcement learning in simulation, to discover skills that would achieve the  
74 generated tasks. Work of [37] also uses language for predicting reward function for robot locomotion,  
75 but does not consider task generation and decomposition or interaction with objects. Our work is the  
76 first to use LLMs for task decomposition and reward generation, as well as asset generation.

77 **Automating 3D asset creation with generative models** The traditional process of creating 3D assets  
78 typically involves multiple labor-intensive stages, including geometry modeling, shape baking, UV  
79 mapping, material creation, texturing and physics parameter estimation, where different software  
80 tools and the expertise of skilled artists are often required. It is thus desirable to automate 3D asset  
81 generation to automatically generate high-quality assets that support realistic rendering under arbitrary  
82 views and have plausible physical behaviours during force application and contacts. The lack of  
83 available 3D data and the abundance of 2D image data have stimulated interest in learning 3D models  
84 from 2D image generators [38, 39]. The availability of strong 2D image generative models based on  
85 diffusion led to high-quality 3D models from text descriptions [40, 41, 42] or single 2D images using  
86 the diffusion model as a 2D prior [43, 44, 45]. In this work, instead of a text-conditioned model,  
87 we use a view and relative pose conditioned image generative model, which we found to provide  
88 better prior for score distillation. Some methods attempt to use videos of assets and differentiable  
89 simulations to estimate their physics parameters and/or adapt the simulation environment, in an  
90 attempt to close the simulation to reality gap [46, 47, 48]. Our effort is complementary to these  
91 works.

92 **Simulation environments for robotic skill learning** In recent years, improving simulators for robot  
93 manipulation has attracted increasingly more attention. Many robotic manipulation environments and  
94 benchmarks [49, 50, 19] are built on top of either PyBullet [51] or MuJoCo [52] as their underlying  
95 physics engines, which mainly support rigid-body simulation. Recently, environments supporting  
96 soft-body manipulation, such as FleX [53], SAPIEN [19], TDW [54], SoftGym [55] and FluidLab  
97 [17] provide capabilities for simulating deformable objects and fluids. Our automated asset and task  
98 generation are not tied to any specific simulation platforms and can be used with any of them. We  
99 unleash the common sense knowledge and reasoning capabilities provided by LLMs and use them to  
100 suggest task descriptions, task decompositions, and reward functions. We then use reinforcement  
101 learning to discover solution trajectories instead of TAMP-based search.

## 102 3 Gen2Sim

103 Gen2Sim generates 3D assets from object images using image diffusion models and predicts physical  
104 parameters for them using LLMs (Section 3.1). It then prompts LLMs to generate language task  
105 descriptions and corresponding reward functions for each generated or human-developed asset,  
106 suitable to their affordances (Section 3.2). Finally, we train RL policies in the generated environments  
107 using the generated reward functions, allowing robots to acquire manipulation skills in diverse scenes  
108 and tasks. We additionally show the applicability of the simulation trained policy by constructing  
109 digital twin environment in simulation, and deploy the trained trajectory in the real world (Section  
110 3.3). See Figure 2 for our method overview.

### 111 3.1 3D Asset Generation

112 Gen2Sim automates 3D asset generation by transforming 2D images of objects to textured 3D meshes  
113 with plausible physics parameters. The images can be 1) real images taken in the robot’s environment,

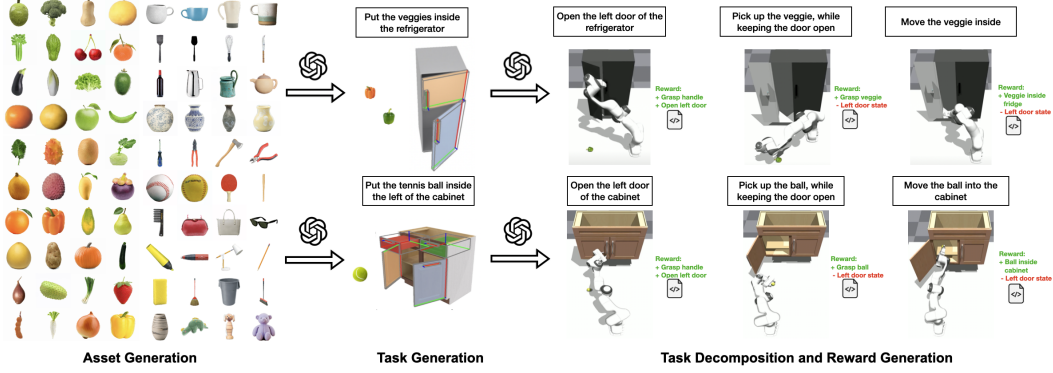


Figure 2: **The Gen2Sim pipeline:** Gen2Sim first generates 3d assets by lifting (generated) 2D images to 3D, and then use both generated assets and assets obtained from other publicly available datasets to populate environments. Afterwards, it queries LLMs to generate meaningful tasks given the scene description, performs task decomposition, generates policy training supervision (reward functions), and yields automated skill learning.

124 2) real images provided by Google search under relevant category names, e.g., “avocado”, or 3)  
 125 images generated by pre-trained text-conditioned diffusion models, such as stable diffusion [56],  
 126 prompted appropriately to generate uncluttered images of the relevant objects, e.g., “an image of an  
 127 individual avocado”. We query GPT-4 [57] for a list of object categories relevant for manipulation  
 128 tasks to search online for or to generate, instead of manually designing it. (Check out our project site  
 129 for a detailed list of the objects we generated.) Given a real or generated 2D image of an object, we  
 130 lift it to a 3D model using Score Distillation Sampling [40, 41], initially developed in [40, 58] for  
 131 text-to-3D lifting. We provide background on image diffusion models below, before we describe our  
 132 3D model fitting approach.

### 123 3.1.1 Image diffusion models

134 A diffusion model learns to model a probability distribution  $p(x)$  by inverting a process that gradually  
 135 adds noise to the image  $x$ . The diffusion process is associated with a variance schedule  $\{\beta_t \in$   
 136  $(0, 1)\}_{t=1}^T$ , which defines how much noise is added at each time step. The noisy version of sample  
 137  $x$  at time  $t$  can then be written  $x_t = \sqrt{\alpha_t}x + \sqrt{1 - \alpha_t}\epsilon$  where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ , is a sample from a  
 138 Gaussian distribution (with the same dimensionality as  $x$ ),  $\alpha_t = 1 - \beta_t$ , and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ . One  
 139 then learns a denoising neural network  $\hat{\epsilon} = \epsilon_\phi(x_t; t)$  that takes as input the noisy image  $x_t$  and the  
 140 noise level  $t$  and tries to predict the noise component  $\epsilon$ . Diffusion models can be easily extended to  
 141 draw samples from a distribution  $p(x|\mathbf{c})$  conditioned on a prompt  $\mathbf{c}$ , where  $\mathbf{c}$  can be a text description,  
 142 a camera pose, and image semantic map, *etc* [4, 59, 60]. Conditioning on the prompt can be done  
 143 by adding  $\mathbf{c}$  as an additional input of the network  $\epsilon_\phi$ . For 3D lifting, we build on Zero-1-to-3 [61], a  
 diffusion model for novel object view synthesis that conditions on an image view of an object and a  
 relative camera rotation around the object to generate plausible images for the target object viewpoint,  
 $\mathbf{c} = [I_1, \pi]$ . It is trained on a large collection  $\mathcal{D}' = \{(x^i, \mathbf{c}^i)\}_{i=1}^N$  of images paired with views and  
 relative camera orientations as conditioning prompt by minimizing the loss:

$$\mathcal{L}_{\text{diff}}(\phi; \mathcal{D}') = \frac{1}{|\mathcal{D}'|} \sum_{x^i, \mathbf{c}^i \in \mathcal{D}'} \|\epsilon_\phi(\sqrt{\alpha_t}x^i + \sqrt{1 - \alpha_t}\epsilon, \mathbf{c}^i, t) - \epsilon\|^2.$$

### 138 3.1.2 Image-to-3D Mesh using Score Distillation Sampling

139 Given an image and relative camera pose 2D diffusion model  $p(I|[I_0, \pi])$ , we extract from it a 3D  
 140 rendition of the input image  $I_0$ , represented by a differential 3D representation using Score Distillation  
 141 Sampling (SDS). We do so by randomly sampling a camera pose  $\pi$ , rendering a corresponding view  $I_\pi$ ,  
 142 assessing the likelihood of the view based on the model  $p(I_\pi|[I_0, \pi])$ , and updating the differentiable  
 143 3D representation to increase the likelihood of the generated view based on the model. Specifically,

144 the denoiser network is frozen and the 3D model is updated as:

$$\nabla(\theta)\mathcal{L}_{SDS}(\theta; \pi, \mathbf{c}, t) = \mathbb{E}_{t, \epsilon}[w(t)(\epsilon_{\phi}(a_t I + \sigma_t \epsilon; t, \mathbf{c}) - \epsilon) \cdot \nabla_{\theta} I],$$

145 where  $I = R(\theta, \pi)$  is the image rendered from a given viewpoint  $\pi$ . We use a two-stage fitting,  
146 wherein the first stage an instantNGP NeRF representation [62] is used, similar to RealFusion [43],  
147 and in the second stage a mesh-based representation is initialized from the NeRF and finetuned  
148 differentially, similar to Fantasia3D [41]. More information of our score distillation sampling can be  
149 found in our website.

150 SDS was initially developed in [40, 58] for text-to-3D lifting. Gen2Sim considers an image as  
151 input for 3D lifting instead. In Section 4, we compare against RealFusion [43] and Fantasia3D  
152 [41] that also consider image-to-3D lifting by textual inversion for diffusion adaptation and by an  
153 image re-projection loss, respectively. We show our proposed pipeline generates more faithful 3D  
154 models from images because the image likelihood provided by the view and pose conditioned image  
155 generative model [44] is more informative than a generic or personalized text-conditioned one.

### 156 3.1.3 Generating plausible physical properties

157 The visual and collision parameters of an asset are generated from the Image-to-Mesh pipeline  
158 discussed above. To define 3D sizes and physics parameters for the generated 3D meshes, we query  
159 GPT-4 regarding the range of plausible width, height, and length for each object, and the range of  
160 mass given the object category. We then scale the generated 3D mesh based on the produced size  
161 range. We feed the mass and 3D mesh information to MeshLab [63] to get the inertia matrix for the  
162 asset. Our prompts for querying GPT for mass and 3D object size can be found on our website. We  
163 wrap the generated mesh information, its semantic name, as well as the physical parameters into  
164 URDF files to be loaded into our simulator.

## 165 3.2 Task Generation, Decomposition and Reward Function

166 Given either generated assets or assets obtained from publically available datasets, we prompt LLMs  
167 to suggest meaningful manipulation tasks considering their affordances, to decompose these tasks into  
168 subtasks when possible, and to generate reward functions for each subtask. We train reinforcement  
169 learning policies for each (sub)task using the generated reward functions, and then chain them together  
170 to solve long horizon tasks, which would have been impossible without LLMs’ decomposition.

171 Prompts to generate task descriptions, task decompositions and rewards functions contain three  
172 elements:

173 **1. Asset descriptions** We use combinations of assets we generate using the method of Section 3.1,  
174 as well as articulated assets from PartNet Mobility [19] and GPartNet dataset [64]. We populate  
175 our simulation environment with randomly sampled assets. Then, we extract information from the  
176 URDF files including link names, joints with their types, and limits, using automated scripts. For  
177 example, an asset `microwave` has parts [`door`, `handle`, and `body`], and joint [`door-joint`]  
178 of type `revolute` with a joint position range  $[0, 1]$ . We then feed the extracted configurations of  
179 the assets to the LLM, with one example shown below:

```
180 The environment contains the following assets:  
181 1. asset_name: "microwave"  
182   part_configuration:  
183     Part 1: "body"  
184     Part 2: "door"  
185       - link_name: "link_0"  
186       - joint_name: "joint_0"  
187       - joint_type: "revolute"  
188       - limit: [0, 1]  
189     Part 3: "handle"  
190       - link_name: "handle_0"  
191       - joint_name: "handlejoint_0"  
192       - joint_type: "fixed"  
193 2. asset_name: "cup"  
194   part_configuration:
```

```

195 Part 1: cup
196     - link_name: "base"
197     - joint_name: "base_joint"
198     - joint_type: "fixed"

```

199 **2. Instructions** These are instructions that regulate the response from the LLM. It includes function  
200 APIs that can be used by the LLM to query the pose of the robot end-effector, as well as different  
201 assets in the given environment:

```

202 List meaningful manipulation tasks that can be performed
203 in this environment. Give subtask decomposition and the
204 order of execution to solve the task. Also, provide the
205 reward function for each subtask.
206
207 The following tasks can be performed in this environment:
208 1. Open the Microwave Door
209 2. Close the Microwave Door
210 3. Pick Cup
211 4. Place Cup
212 5. Put the Cup in the Microwave
213     This task needs to be decomposed into sub-tasks:
214         - Open the Microwave
215         - Pick Cup
216         - Place the Cup in the Microwave

```

217 **3. Task and Decomposition Examples** are question-to-language pairs that present few-shot in-  
218 context demonstrations of how tasks can be decomposed into subtasks.

```

219 List meaningful manipulation tasks that can be performed
220 in this environment. Give subtask decomposition and the
221 order of execution to solve the task. Also, provide the
222 reward function for each subtask.
223
224 The following tasks can be performed in this environment:
225 1. Open the Microwave Door
226 2. Close the Microwave Door
227 3. Pick Cup
228 4. Place Cup
229 5. Put the Cup in the Microwave
230     This task needs to be decomposed into sub-tasks:
231         - Open the Microwave
232         - Pick Cup
233         - Place the Cup in the Microwave

```

234 **4. Examples of reward functions** are task to reward function pairs that present few-shot demonstra-  
235 tions of how tasks can be translated to reward functions. For the following example task, we provide  
236 example reward functions composed of 1) distance reward: distance between the end-effector and the  
237 target part, and 2) state reward: distance between the current and the target pose of an articulated  
238 asset, link, or joint. Note that the following is just an example for the LLM to use as a reference.

```

239 Task: OpenMicrowaveDoor
240 Task Description: open the door of the microwave
241 ```
242 def compute_reward(env):
243     # reward function
244     door_handle_pose = env.get_pose_by_link_name("microwave", "handle_0")
245     gripper_pose = env.get_robot_gripper_pose()
246     distance_gripper_to_handle = torch.norm(door_handle_pose - gripper_pose, dim=-1)
247     door_state = env.get_state_by_joint_name("microwave", "joint_0")
248     cost = distance_gripper_to_handle - door_state
249     reward = - cost
250
251     # success condition
252     target_door_state = env.get_limits_by_joint_name("microwave", "joint_0")["upper"]
253     success = torch.abs(door_state - target_door_state) < 0.1
254
255     return reward, success
256 ```

```

257 We provide a comprehensive list of our prompts on our website. We show in Section 4 that our method  
258 can generalize across assets, suggest diverse and plausible tasks, and reward functions automatically,  
259 without any additional human involvement.

### 260 3.3 Sequential Reinforcement Learning for Long Horizon Tasks

261 We train reinforcement learning policies using Proximal Policy Optimization (PPO) [65] maximizing  
262 the generated reward functions for each subtask. We train RL for each generated subtask in temporal  
263 order. Once training for a subtask converges, we proceed to the next subtask. The initial state  
264 of the gripper and the environment are sampled from the resulting states of the previous subtask  
265 execution. This ensures policies can be temporally chained upon training. Note that while training  
266 till convergence doesn't guarantee successful policy training, since we decompose the high-level task  
267 into very fine-grained simplistic subtasks, such a heuristic works practically well in our experiments.  
268 Our RL policies are trained per environment using privileged information of the simulation state  
269 to facilitate learning. Such learned policies can be used as demonstration data and distilled into  
270 vision-language transformer policies ([6, 66, 67]), and we leave this to our future work.

## 271 4 Experiments

272 Our experiments aim to answer the following questions:

- 273 1. Can Gen2Sim generate plausible geometry, appearance, and physics for diverse types of objects  
274 and parts, without human expertise and with minimal human involvement?
- 275 2. Can Gen2Sim generate task language goals and reward functions for novel object categories, novel  
276 assets with different part configurations, and a combination of multiple assets in an environment?
- 277 3. Can the generated environments and reward function lead to successful learning of RL policies?

### 278 4.1 Asset Generation

279 We compare our image-to-3D lifting with two baselines:

280 1. *RealFusion* [43], which uses textual inversion of [68] to  
281 learn a text embedding for the depicted object concept in  
282 an image, and uses text-conditioned diffusion model that  
283 uses this text embedding in the text prompt during score  
284 distillation.

285 2. *Make-It-3D* [44], which uses the same NeRF and tex-  
286 tured mesh two-stage fitting with SDS as ours do, but does  
287 not use a view and pose conditioned generative model,  
288 rather a text-based diffusion model, similar to [40].

289 We show comparisons on several example objects in Fig-  
290 ure 3, with images rendered from 4 different views. Our  
291 generates more plausible 3D assets as our diffusion prior  
292 comes from an image and pose-conditioned model in  
293 comparison to approaches like Fantasia3D or RealFusion  
294 which uses text conditioning. For more visualizations and  
295 details of our generated assets, with diversified textures  
296 and their behaviors under gravity and collisions, please  
297 refer to our website.

### 298 4.2 Automated Skill Learning

299 We make use of GPU-parallel data sampling in IsaacGym [69] which enables fast and stable conver-  
300 gence of our policies. Our robotic setup uses a Franka Panda arm with a mobile base. It is equipped  
301 with either a parallel gripper or a suction cup, depending on the task needs. The suction gripper is  
302 only used in pick-and-place tasks where grasping varied geometric objects was intrinsically hard with  
303 RL. In all other task categories, we use the parallel jaw gripper. Our state representation for PPO  
304 includes the robot's joint position  $q \in \mathbb{R}^{11}$ , velocity  $\dot{q} \in \mathbb{R}^{11}$  (7-DoF arm,  $x$  and  $y$  for the mobile  
305 base and 2 extra DoFs from the gripper jaws), orientation of the gripper  $r \in SO(3)$ , and poses and

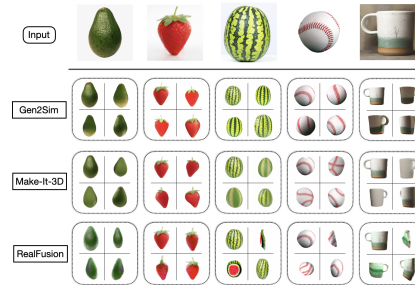


Figure 3: **Left: 3D asset generation** from Gen2Sim, RealFusion [43] and Make-It-3D [44]. Gen2Sim uses a view and camera pose conditioned image generative model during score distillation, which helps generate more accurate 3D geometry in comparison to the baselines.

306 joint configurations of the assets present in the scene. We use position control and at each timestep  $t$ ,  
307 our policy produces target gripper pose and configurations, which is then converted to target robot  
308 configurations by computing inverse kinematics. A low-level PID torque controller provided by  
309 IsaacGym is then used to produce low-level joint torque commands.

310 Gen2Sim generates diverse tasks, plausible natural language task descriptions, task decompositions  
311 and reward functions automatically for hundreds of assets, with different category labels and number  
312 of joints, based on the examples provided by the prompt. We show some examples of such generated  
313 task descriptions and their corresponding language descriptions in Figure 1 and more on our website.  
314 We show example task decompositions in Figure 2. At the time of submission, our pipeline has  
315 generated hundreds of tasks, which we will release upon publication. Note that our method can be  
316 queried endlessly to generate more tasks and provide task-specific policy demonstrations, which  
317 could be used for policy distillation in the future.

318 We provide all prompts in our website, alongside examples of GPT’s responses. Only one example is  
319 included in our prompt for task decomposition and reward generation; it concerns the task of “putting  
320 a cup in a Microwave”. We show then the prediction of GPT4 regarding task and reward function  
321 for instances of Door, DishWasher, fruits, veggies and others. Note that the articulation structure  
322 structure across all of the assets differ significantly, but our method can effectively generalize. Also,  
323 Gen2Sim capitalizes on the common sense knowledge of LLMs regarding object affordances, and  
324 thus can produce meaningful ways for interacting with the assets, such as “press the sanitizer” and  
325 “turn the faucet”, as shown in Figure 1. The rewards generated by GPT can be well optimized with  
326 off-the-shelf RL algorithms [65] to learn useful manipulation policies, and the policies are able to  
327 solve the tasks upon convergence.

## 328 5 Limitations

329 There still remain two limitations that need to be addressed for the proposed system to materialize  
330 into a platform for large-scale robot skill learning that are deployable in real-world, as identified  
331 below:

332 **1. Beyond rigid asset generation:** The assets we can currently generate are rigid or mostly rigid  
333 objects, which do not deform significantly under external forces. For articulated assets, we are  
334 using existing manually designed and labelled datasets ([19, 64]). To generate articulated objects,  
335 deformable objects and liquids, accurate fine-grained video perception is required in combination  
336 with generative priors to model the temporal dynamics of their geometry and appearance. This is an  
337 exciting and challenging direction for future work.

338 **2. Simulation to reality gap.** Improving fidelity and efficiency of simulators is an active area of  
339 research that our method will dramatically benefit from, and we plan to work on. Also, combining  
340 explicit physics engines with learnt residual models from simulation and real world alignment [70]  
341 for decreasing the simulation to reality gap is an exciting research direction we plan to pursue.

## 342 6 Conclusion

343 We have presented Gen2Sim, a pipeline for automating the development of simulation environments,  
344 tasks and reward functions with pre-trained generative models of vision and language. We presented  
345 methods that create and augment geometry, textures and physics of object assets from single images,  
346 parse URDF files of assets, generate task descriptions, decompositions and reward python functions,  
347 and train reinforcement learning policies to solve the generated long horizon tasks. Addressing the  
348 limitations including generating diverse assets with more complex physical properties, and transfer  
349 trained policy to real world using realistic vision input in a closed-loop manner, are direct avenues  
350 for our future work. We believe generative models of images and language will play an important  
351 role in automating and supersizing robot training data in simulation, and in crossing the sim2real  
352 gap, necessary for delivering robot generalists in the real world. Gen2Sim takes one first step in that  
353 direction.



354 **References**

- 355 [1] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan,  
356 P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child,  
357 A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray,  
358 B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei.  
359 Language models are few-shot learners, 2020. URL [https://arxiv.org/abs/2005.](https://arxiv.org/abs/2005.14165)  
360 [14165](https://arxiv.org/abs/2005.14165).
- 361 [2] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,  
362 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from  
363 natural language supervision. *CoRR*, abs/2103.00020, 2021. URL [https://arxiv.org/](https://arxiv.org/abs/2103.00020)  
364 [abs/2103.00020](https://arxiv.org/abs/2103.00020).
- 365 [3] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever. Robust speech  
366 recognition via large-scale weak supervision, 2022.
- 367 [4] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis  
368 with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
369 *and Pattern Recognition*, pages 10684–10695, 2022.
- 370 [5] S. Levine, P. Pastor, A. Krizhevsky, and D. Quillen. Learning hand-eye coordination for robotic  
371 grasping with deep learning and large-scale data collection, 2016.
- 372 [6] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Haus-  
373 man, A. Herzog, J. Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv*  
374 *preprint arXiv:2212.06817*, 2022.
- 375 [7] S. Höfer, K. E. Bekris, A. Handa, J. C. Gamboa, M. Mozifian, F. Golemo, C. G. Atkeson,  
376 D. Fox, K. Goldberg, J. Leonard, C. K. Liu, J. Peters, S. Song, P. Welinder, and M. White.  
377 Sim2real in robotics and automation: Applications and challenges. *IEEE Trans Autom. Sci.*  
378 *Eng.*, 18(2):398–400, 2021. doi:10.1109/TASE.2021.3064065. URL [https://doi.org/](https://doi.org/10.1109/TASE.2021.3064065)  
379 [10.1109/TASE.2021.3064065](https://doi.org/10.1109/TASE.2021.3064065).
- 380 [8] Z. Fu, A. Kumar, A. Agarwal, H. Qi, J. Malik, and D. Pathak. Coupling vision and proprioception  
381 for navigation of legged robots. In *Proceedings of the IEEE/CVF Conference on Computer*  
382 *Vision and Pattern Recognition*, pages 17273–17283, 2022.
- 383 [9] A. Kumar, Z. Fu, D. Pathak, and J. Malik. Rma: Rapid motor adaptation for legged robots,  
384 2021.
- 385 [10] I. Radosavovic, T. Xiao, B. Zhang, T. Darrell, J. Malik, and K. Sreenath. Learning humanoid  
386 locomotion with transformers. *arXiv preprint arXiv:2303.03381*, 2023.
- 387 [11] T. Chen, J. Xu, and P. Agrawal. A system for general in-hand object re-orientation. In  
388 *Conference on Robot Learning*, pages 297–307. PMLR, 2022.
- 389 [12] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino,  
390 M. Plappert, G. Powell, R. Ribas, J. Schneider, N. Tezak, J. Tworek, P. Welinder, L. Weng,  
391 Q. Yuan, W. Zaremba, and L. Zhang. Solving rubik’s cube with a robot hand, 2019.
- 392 [13] E. Kaufmann, A. Loquercio, R. Ranftl, M. Müller, V. Koltun, and D. Scaramuzza. Deep drone  
393 acrobatics. *arXiv preprint arXiv:2006.05768*, 2020.
- 394 [14] A. Loquercio, E. Kaufmann, R. Ranftl, M. Müller, V. Koltun, and D. Scaramuzza. Learning  
395 high-speed flight in the wild. *Science Robotics*, 6(59):eabg5810, 2021.
- 396 [15] X. Lin, Z. Huang, Y. Li, J. B. Tenenbaum, D. Held, and C. Gan. Diffskill: Skill abstraction  
397 from differentiable physics for deformable object manipulations with tools. *arXiv preprint*  
398 *arXiv:2203.17275*, 2022.

- 399 [16] Z. Xu, Z. Xian, X. Lin, C. Chi, Z. Huang, C. Gan, and S. Song. Roboninja: Learning an adaptive  
400 cutting policy for multi-material objects. *arXiv preprint arXiv:2302.11553*, 2023.
- 401 [17] Z. Xian, B. Zhu, Z. Xu, H.-Y. Tung, A. Torralba, K. Fragkiadaki, and C. Gan. Fluidlab: A  
402 differentiable environment for benchmarking complex fluid manipulation. In *International  
403 Conference on Learning Representations*, 2023.
- 404 [18] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun. Carla: An open urban driving  
405 simulator, 2017.
- 406 [19] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang, L. Yi,  
407 A. X. Chang, L. J. Guibas, and H. Su. Sapien: A simulated part-based interactive environment,  
408 2020.
- 409 [20] S. Srivastava, C. Li, M. Lingelbach, R. Martín-Martín, F. Xia, K. Vainio, Z. Lian, C. Gokmen,  
410 S. Buch, C. K. Liu, S. Savarese, H. Gweon, J. Wu, and L. Fei-Fei. Behavior: Benchmark for  
411 everyday household activities in virtual, interactive, and ecological environments, 2021.
- 412 [21] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun,  
413 J. Malik, D. Parikh, and D. Batra. Habitat: A platform for embodied ai research, 2019.
- 414 [22] C. Gan, J. Schwartz, S. Alter, D. Mrowca, M. Schrimpf, J. Traer, J. D. Freitas, J. Kubilius,  
415 A. Bhandwadar, N. Haber, M. Sano, K. Kim, E. Wang, M. Lingelbach, A. Curtis, K. Feigelis,  
416 D. M. Bear, D. Gutfreund, D. Cox, A. Torralba, J. J. DiCarlo, J. B. Tenenbaum, J. H. McDermott,  
417 and D. L. K. Yamins. Threedworld: A platform for interactive multi-modal physical simulation,  
418 2021.
- 419 [23] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu. Hierarchical planning for long-horizon ma-  
420 nipulation with geometric and symbolic scene graphs. *CoRR*, abs/2012.07277, 2020. URL  
421 <https://arxiv.org/abs/2012.07277>.
- 422 [24] D. Xu, R. Martín-Martín, D. Huang, Y. Zhu, S. Savarese, and L. Fei-Fei. Regression planning  
423 networks. *CoRR*, abs/1909.13072, 2019. URL <http://arxiv.org/abs/1909.13072>.
- 424 [25] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners:  
425 Extracting actionable knowledge for embodied agents. *arXiv preprint arXiv:2201.07207*, 2022.
- 426 [26] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch,  
427 Y. Chebotar, P. Sermanet, N. Brown, T. Jackson, L. Luu, S. Levine, K. Hausman, and B. Ichter.  
428 Inner monologue: Embodied reasoning through planning with language models, 2022. URL  
429 <https://arxiv.org/abs/2207.05608>.
- 430 [27] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as  
431 policies: Language model programs for embodied control, 2022. URL <https://arxiv.org/abs/2209.07753>.
- 433 [28] J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. H. Chi, Q. Le, and D. Zhou. Chain of thought  
434 prompting elicits reasoning in large language models. *CoRR*, abs/2201.11903, 2022. URL  
435 <https://arxiv.org/abs/2201.11903>.
- 436 [29] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig. Pre-train, prompt, and predict: A  
437 systematic survey of prompting methods in natural language processing. *CoRR*, abs/2107.13586,  
438 2021. URL <https://arxiv.org/abs/2107.13586>.
- 439 [30] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan,  
440 P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child,  
441 A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray,  
442 B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei.  
443 Language models are few-shot learners, 2020.

- 444 [31] S. C. Y. Chan, A. Santoro, A. K. Lampinen, J. X. Wang, A. Singh, P. H. Richemond, J. Mc-  
445 Clelland, and F. Hill. Data distributional properties drive emergent in-context learning in  
446 transformers, 2022.
- 447 [32] X. Zhu, Y. Chen, H. Tian, C. Tao, W. Su, C. Yang, G. Huang, B. Li, L. Lu, X. Wang, Y. Qiao,  
448 Z. Zhang, and J. Dai. Ghost in the minecraft: Generally capable agents for open-world  
449 environments via large language models with text-based knowledge and memory, 2023.
- 450 [33] S. Lifshitz, K. Paster, H. Chan, J. Ba, and S. McIlraith. Steve-1: A generative model for  
451 text-to-behavior in minecraft, 2023.
- 452 [34] G. Wang, Y. Xie, Y. Jiang, A. Mandlkar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar. Voyager:  
453 An open-ended embodied agent with large language models, 2023.
- 454 [35] H. Ha, P. Florence, and S. Song. Scaling up and distilling down: Language-guided robot skill  
455 acquisition. *arXiv preprint arXiv:2307.14535*, 2023.
- 456 [36] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei. Voxposer: Composable 3d value  
457 maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023.
- 458 [37] W. Yu, N. Gileadi, C. Fu, S. Kirmani, K.-H. Lee, M. G. Arenas, H.-T. L. Chiang, T. Erez,  
459 L. Hasenclever, J. Humplik, et al. Language to rewards for robotic skill synthesis. *arXiv preprint*  
460 *arXiv:2306.08647*, 2023.
- 461 [38] T. Nguyen-Phuoc, C. Li, L. Theis, C. Richardt, and Y.-L. Yang. Hologan: Unsupervised learning  
462 of 3d representations from natural images, 2019.
- 463 [39] E. R. Chan, M. Monteiro, P. Kellnhofer, J. Wu, and G. Wetzstein. pi-gan: Periodic implicit  
464 generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF*  
465 *conference on computer vision and pattern recognition*, pages 5799–5809, 2021.
- 466 [40] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion.  
467 *arXiv*, 2022.
- 468 [41] R. Chen, Y. Chen, N. Jiao, and K. Jia. Fantasia3d: Disentangling geometry and appearance for  
469 high-quality text-to-3d content creation. *arXiv preprint arXiv:2303.13873*, 2023.
- 470 [42] C.-H. Lin, J. Gao, L. Tang, T. Takikawa, X. Zeng, X. Huang, K. Kreis, S. Fidler, M.-Y. Liu,  
471 and T.-Y. Lin. Magic3d: High-resolution text-to-3d content creation. In *IEEE Conference on*  
472 *Computer Vision and Pattern Recognition (CVPR)*, 2023.
- 473 [43] L. Melas-Kyriazi, C. Rupprecht, I. Laina, and A. Vedaldi. Realfusion: 360 reconstruction of  
474 any object from a single image. In *CVPR*, 2023. URL [https://arxiv.org/abs/2302.](https://arxiv.org/abs/2302.10663)  
475 [10663](https://arxiv.org/abs/2302.10663).
- 476 [44] J. Tang, T. Wang, B. Zhang, T. Zhang, R. Yi, L. Ma, and D. Chen. Make-it-3d: High-fidelity 3d  
477 creation from a single image with diffusion prior, 2023.
- 478 [45] B. Shen, X. Yan, C. R. Qi, M. Najibi, B. Deng, L. Guibas, Y. Zhou, and D. Anguelov. Gina-3d:  
479 Learning to generate implicit neural assets in the wild, 2023.
- 480 [46] E. Heiden, C. E. Denniston, D. Millard, F. Ramos, and G. S. Sukhatme. Probabilistic inference  
481 of simulation parameters via parallel differentiable simulation, 2022.
- 482 [47] E. Heiden, M. Macklin, Y. Narang, D. Fox, A. Garg, and F. Ramos. Disect: A differentiable  
483 simulator for parameter inference and control in robotic cutting, 2022.
- 484 [48] K. Wang, W. R. J. I. au2, S. Lu, X. Huang, J. Booth, R. Kramer-Bottiglio, M. Aanjaneya, and  
485 K. Bekris. Real2sim2real transfer for control of cable-driven robots via a differentiable physics  
486 engine, 2023.

- 487 [49] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba.  
488 Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 489 [50] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, D. Gordon, Y. Zhu,  
490 A. Gupta, and A. Farhadi. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint*  
491 *arXiv:1712.05474*, 2017.
- 492 [51] E. Coumans and Y. Bai. Pybullet, a python module for physics simulation for games, robotics  
493 and machine learning. <http://pybullet.org>, 2016.
- 494 [52] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012*  
495 *IEEE/RSJ international conference on intelligent robots and systems*, pages 5026–5033. IEEE,  
496 2012.
- 497 [53] M. Macklin, M. Müller, N. Chentanez, and T.-Y. Kim. Unified particle physics for real-time  
498 applications. *ACM Transactions on Graphics (TOG)*, 33(4):1–12, 2014.
- 499 [54] C. Gan, J. Schwartz, S. Alter, M. Schrimpf, J. Traer, J. De Freitas, J. Kubilius, A. Bhandwaldar,  
500 N. Haber, M. Sano, et al. Threedworld: A platform for interactive multi-modal physical  
501 simulation. *arXiv preprint arXiv:2007.04954*, 2020.
- 502 [55] X. Lin, Y. Wang, J. Olkin, and D. Held. Softgym: Benchmarking deep reinforcement learning  
503 for deformable object manipulation. *arXiv preprint arXiv:2011.07215*, 2020.
- 504 [56] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis  
505 with latent diffusion models, 2022.
- 506 [57] OpenAI. Gpt-4 technical report, 2023.
- 507 [58] H. Wang, X. Du, J. Li, R. A. Yeh, and G. Shakhnarovich. Score jacobian chaining: Lifting  
508 pretrained 2d diffusion models for 3d generation, 2022.
- 509 [59] Y. Li, H. Liu, Q. Wu, F. Mu, J. Yang, J. Gao, C. Li, and Y. J. Lee. Gligen: Open-set grounded  
510 text-to-image generation. *arXiv preprint arXiv:2301.07093*, 2023.
- 511 [60] L. Zhang and M. Agrawala. Adding conditional control to text-to-image diffusion models.  
512 *arXiv preprint arXiv:2302.05543*, 2023.
- 513 [61] R. Liu, R. Wu, B. V. Hoorick, P. Tokmakov, S. Zakharov, and C. Vondrick. Zero-1-to-3:  
514 Zero-shot one image to 3d object, 2023.
- 515 [62] T. Müller, A. Evans, C. Schied, and A. Keller. Instant neural graphics primitives with a  
516 multiresolution hash encoding. *ACM Transactions on Graphics*, 41(4):1–15, jul 2022. doi:10.  
517 1145/3528223.3530127. URL <https://doi.org/10.1145%2F3528223.3530127>.
- 518 [63] G. R. M. C. Paolo Cignoni, Alessandro Muntoni. MeshLab.
- 519 [64] H. Geng, H. Xu, C. Zhao, C. Xu, L. Yi, S. Huang, and H. Wang. Gapartnet: Cross-category  
520 domain-generalizable object perception and manipulation via generalizable and actionable parts,  
521 2023.
- 522 [65] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization  
523 algorithms, 2017.
- 524 [66] S. Christen, W. Yang, C. Pérez-D’Arpino, O. Hilliges, D. Fox, and Y.-W. Chao. Learning  
525 human-to-robot handovers from point clouds, 2023.
- 526 [67] A. Jaegle, S. Borgeaud, J.-B. Alayrac, C. Doersch, C. Ionescu, D. Ding, S. Koppula, D. Zoran,  
527 A. Brock, E. Shelhamer, et al. Perceiver io: A general architecture for structured inputs &  
528 outputs. *arXiv preprint arXiv:2107.14795*, 2021.

- 529 [68] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-Or. An  
530 image is worth one word: Personalizing text-to-image generation using textual inversion, 2022.
- 531 [69] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin,  
532 A. Allshire, A. Handa, et al. Isaac gym: High performance gpu-based physics simulation for  
533 robot learning. *arXiv preprint arXiv:2108.10470*, 2021.
- 534 [70] A. Zeng, S. Song, J. Lee, A. Rodriguez, and T. Funkhouser. Tossingbot: Learning to throw  
535 arbitrary objects with residual physics, 2020.