

MINERVAS: Massive INterior EnviRonments VirtuAl Synthesis

Supplementary Material

Haocheng Ren^{1,*}, Hao Zhang^{1,*}, Jia Zheng², Jiaxiang Zheng², Rui Tang², Rui Wang¹, Hujun Bao¹

¹State Key Lab of CAD&CG, Zhejiang University ²Manycore Tech (Kujiale)

<https://coohom.github.io/MINERVAS>

In this supplementary material, we first show the all experimental details, including the DSL code for generating imagery data, the experimental setups, and more qualitative results. Then, we provide additional examples of random samplers used in our system.

1. Experiment on Room Layout Estimation

DSL code. We first filter the room under Manhattan-world assumption using DSL in the Scene Process Stage. We set the camera model as “panorama” and the camera resolution as 1024×512 in the Entity Process Stage. We also use the sampler of the transformation component to randomly move cameras in each room, and export corner and camera positions. Figure 1 shows the DSL code for this task.

Experimental setup. In this experiment, we use MatterportLayout [2, 5] as the real dataset. The dataset consists of 1,647 images for training, 190 images for validation, and 458 images for testing. We use our system to synthesize 120K panorama images from 80K scenes. Each panorama image corresponds to one room in scenes.

Following [5], we adopt four standard metrics: 3D IoU, 2D IoU, RMSE and the accuracy under the threshold (δ_1). We adopt HorizonNet [4] as the baseline approach. We use an Adam optimizer with an initial learning rate of 3×10^{-4} with a polynomial decay policy. We set the mini-batch size to 24. We also use two training strategies in this experiment, *i.e.*, “r” and “s + r”. In “s + r”, each batch contains 16 images from the real dataset and 8 from the synthetic dataset. For each strategy, we train the network for 30K iterations.

Qualitative results. We show more qualitative results of room layout estimation in Figure ??.

2. Experiment on Semantic Segmentation

DSL code. We set the camera model as “panorama” and the camera resolution as 1024×512 in the Entity Process

Stage, and output the customized semantic label in the Pixel Process Stage. Figure 2 shows the DSL for this task.

Experimental setup. In this experiment, we use 2D-3D-S [1] as the real data. We split the dataset into 955 for training, 84 for validation, and 373 for testing. Then, we synthesize 12k panoramic images using our system. Each panorama image corresponds to one room in scenes.

We use an SGD optimizer with an initial learning rate of 2×10^{-2} with a polynomial decay policy, momentum 0.9, and weight decay of 10^{-4} . We set the mini-batch size to 8. In “s + r”, each batch contains 4 images from the real dataset and 4 from the synthetic dataset. For each strategy, we train the whole network for 10k iterations.

Qualitative results. We show more qualitative results of semantic segmentation in Figure 3. As can be seen, training on the synthetic and real dataset achieves the best result. The boundary is more clear after using synthetic dataset. It demonstrates that our synthetic data could be used to improve the performance of the network.

3. Experiment on Depth Estimation

DSL code. We set the resolution of the image to 640×480 and horizontal field-of-view (FoV) to 53° in the Entity Process Stage, which are the same as Microsoft Kinect used in the NYUv2 dataset. Figure 4 shows the DSL for this task.

Experimental setup. In this experiment, we use NYUv2 as the real data. We split the images into 795 for training and 654 for testing. Then, we synthesize 14k images using our system.

We use an SGD optimizer with an initial learning rate 1×10^{-4} with polynomial decay policy, momentum 0.9, and weight decay 5×10^{-4} . We set the mini-batch size to 8. In “s + r”, each batch contains 4 images from the real dataset and 4 from the synthetic dataset. For each strategy, we train the whole network for 15k iterations.

Qualitative results. We show more qualitative results in Figure 5. As one can see, training on the synthetic and real dataset, the network generates more accurate estimations

*: Equal contribution.

than that only using real images for training. The noise has been significantly reduced in the depth estimation.

4. Additional Results of Random Samplers

Furniture Layout sampler. With the layout sampler, users can change the furniture arrangement in the scene process stage easily, which could further increase the diversity of the scene. We show the code example using the furniture layout sampler and then give more results in Figure 6.

Light sampler. With the light sampler, users could randomly change the light in the scene to simulate various lighting conditions. We show the code example using the light sampler and then give more results in Figure 7.

Model sampler. With the model sampler, users could replace the original model with a new model with an appropriate size and similar category. We show the code example using the model sampler and then give more results in Figure 8.

Material sampler. With the material sampler, users could sample the appropriate new material based on the original material type, according to the material information stored in the database. We show the code example using the material sampler and then give more results in Figure 9.

Custom sampler. The flexibility of our DSL allows users to create their customized sampling strategy. Figure 10 shows an example of customized trajectory sampler DSL. We can further use this customized sampler for vision tasks, *e.g.*, SLAM as in Figure 11.

References

- [1] Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017. [1](#)
- [2] Angel X. Chang, Angela Dai, Thomas A. Funkhouser, Maciej Halber, Matthias Nießner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from RGB-D data in indoor environments. In *3DV*, pages 667–676, 2017. [1](#)
- [3] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Transactions on Graphics (TOG)*, 36(4):1, 2017. [9](#)
- [4] Cheng Sun, Chi-Wei Hsiao, Min Sun, and Hwann-Tzong Chen. Horizonnet: Learning room layout with 1d representation and pano stretch data augmentation. In *CVPR*, pages 1047–1056, 2019. [1](#)
- [5] Chuhang Zou, Jheng-Wei Su, Chi-Han Peng, Alex Colburn, Qi Shan, Peter Wonka, Hung-Kuo Chu, and Derek Hoiem. Manhattan room layout reconstruction from a single 360 image: A comparative study of state-of-the-art methods. *International Journal of Computer Vision (IJCV)*, 2021. [1](#)

```

1 class ManhattanSceneFilter(SceneProcessor):
2     def filter_scene(self):
3         Valid = False
4         for room in self.shader.world.rooms:
5             corners = room.boundary
6             shift_corners = copy.deepcopy(corners)
7             shift_corners.append(shift_corners.pop(0))
8             for i in range(len(shift_corners)):
9                 direction = glm.vec2(shift_corners[i]).xy - glm.vec2(corners[i]).xy
10                direction = glm.normalize(direction)
11                shift_corners[i] = direction
12            EPSILON = 0.001
13            MANHATTAN = True
14            for i in range(len(shift_corners)):
15                cos = glm.dot(shift_corners[i], shift_corners[(i + 1) % len(shift_corners)])
16                if cos > EPSILON:
17                    MANHATTAN = False
18                    break
19            if MANHATTAN:
20                Valid = True
21        return Valid
22
23     def process(self):
24         if self.filter_scene():
25             print("get the expected scene")
26         else:
27             sys.exit(7) # Don't process this scene anymore.
28

1 class CameraSetting(EntityProcessor):
2     def process(self):
3         for camera in self.shader.world.cameras:
4             camera.set_attr("imageWidth", 1024)
5             camera.set_attr("imageHeight", 512)
6             camera.set_attr("cameraType", "PANORAMA")
7             camera.set_attr("position", x=1000, y=1000, z=800)

1 class UserOutput(StructureProcessor):
2     def process(self):
3         # write out the corners of the rooms in the scene
4         for room in self.shader.world.rooms:
5             for plane, height in zip(["floor", "ceiling"], [0, self.shader.world.levels[0].height]):
6                 corners = []
7                 for corner in room.boundary:
8                     corners.append({'x': corner[0], 'y': corner[1], 'z': height})
9                 self.shader.world.pick(
10                     corners=corners,
11                     catName=plane,
12                     type='corners',
13                     id=f'{room.roomId}_{plane}')
14
15         # write out cameras in the scene
16         for camera in self.shader.world.cameras:
17             self.shader.world.pick(
18                 type="camera",
19                 position=camera.position,
20                 id=camera.id
21             )

```

Figure 1. DSL for Manhattan room layout estimation task.

```

1 class CameraSetting(EntityProcessor):
2     def process(self):
3         for camera in self.shader.world.cameras:
4             camera.set_attr("imageWidth", 1024)
5             camera.set_attr("imageHeight", 512)
6             camera.set_attr("cameraType", "PANORAMA")

1 class LabelMapping(PixelProcessor):
2     def process(self, **kwargs):
3         self.gen_semantic(label_arch=NYU40_Mapping)

```

Figure 2. DSL for semantic segmentation task.

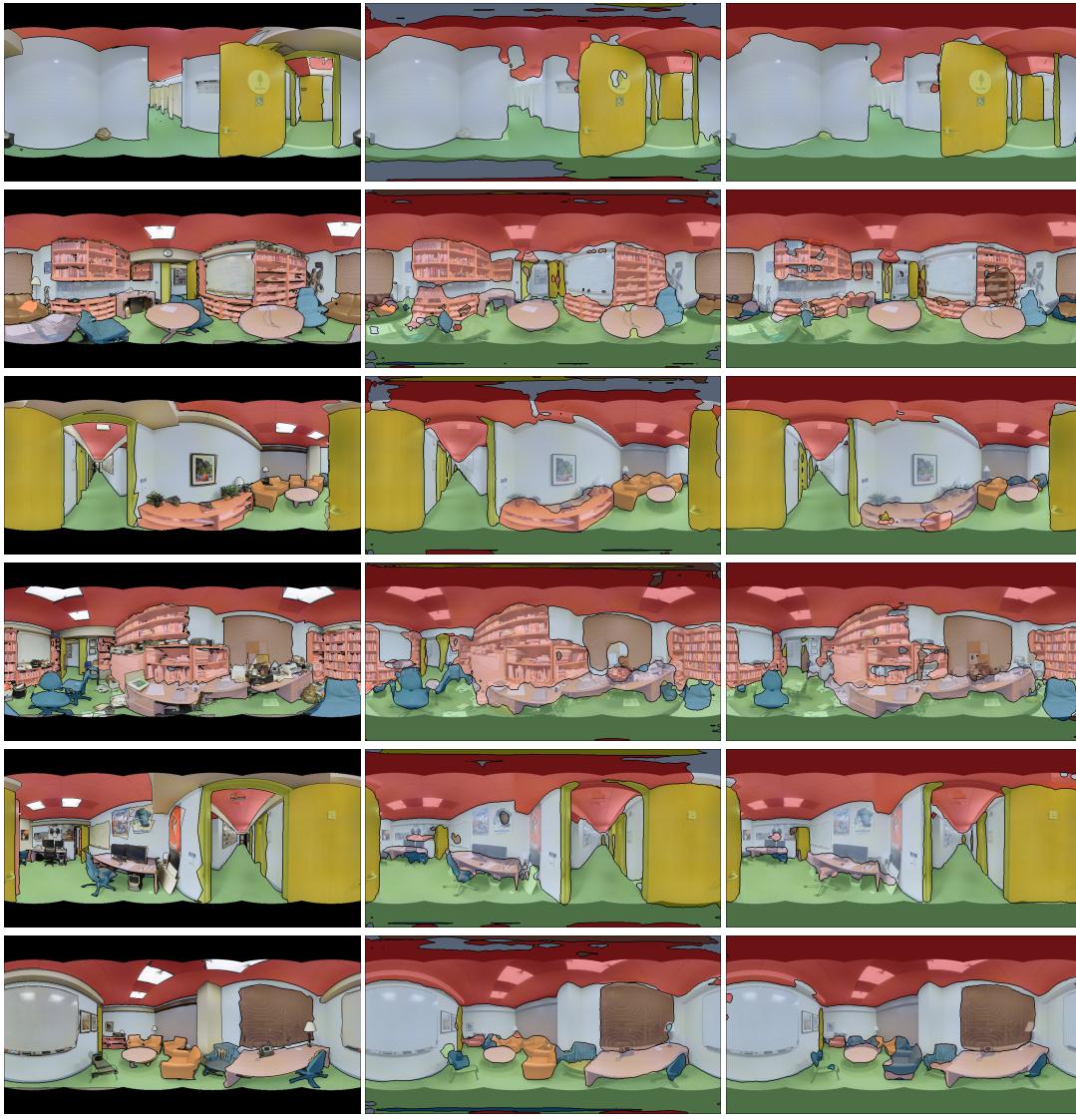


Figure 3. More qualitative results of the semantic segmentation on 2D-3D-S dataset. Different colors denote different semantic categories.

```

1 class CameraSetting(EntityProcessor):
2     def process(self):
3         for camera in self.shader.world.cameras:
4             camera.set_attr("imageWidth", 640)
5             camera.set_attr("imageHeight", 480)
6             camera.set_attr("cameraType", "PERSPECTIVE")
7             camera.set_attr("hfov", 53)

```

Figure 4. DSL for depth estimation.

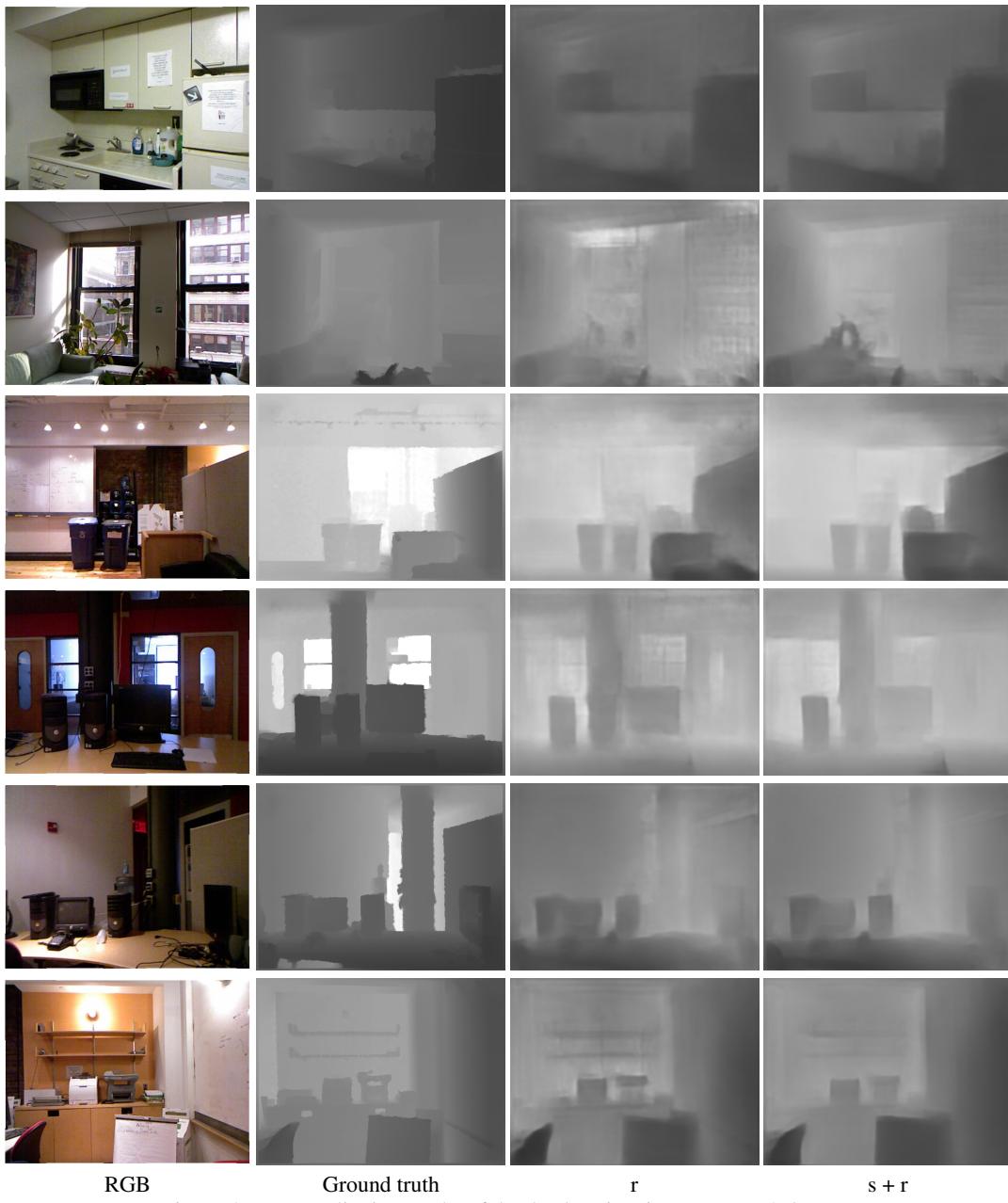


Figure 5. More qualitative results of the depth estimation on NYUv2 dataset.

```

1 class FurnitureLayoutSampler(SceneProcessor):
2     def process(self):
3         for room in self.shader.world.rooms:
4             room.randomize_layout(self.shader.world)

```



Figure 6. More results generated by the scene-level sampler.

```

1 class FurnitureLayoutSampler(SceneProcessor):
2     def process(self):
3         for room in self.shader.world.rooms:
4             room.randomize_layout(self.shader.world)

```

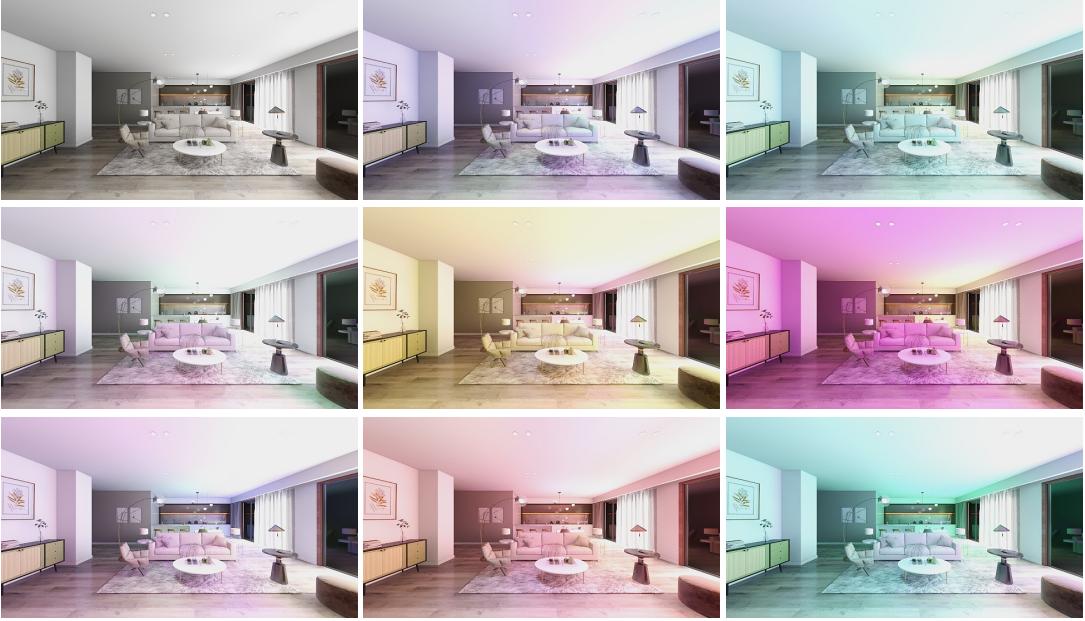


Figure 7. More results generated by the light sampler.

```

1 class ModelSampler(EntityProcessor):
2     def process(self):
3         for instance in self.shader.world.instances:
4             self.shader.world.replace_model(id=instance.id)

```

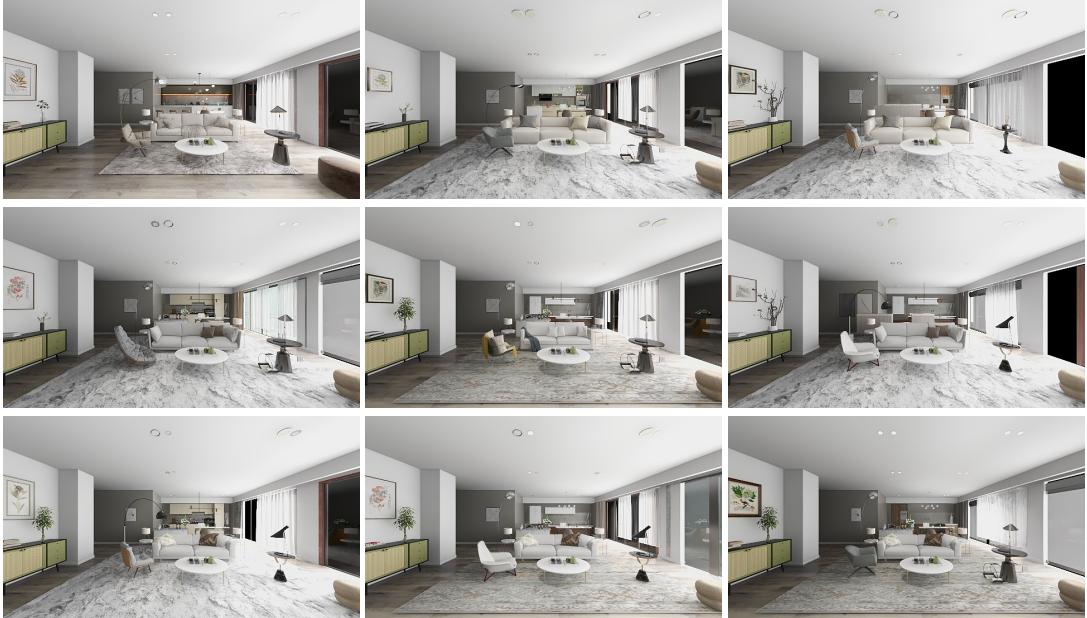


Figure 8. More results generated by the mesh sampler.

```

1 class MaterialSampler(EntityProcessor):
2     def process(self):
3         for instance in self.shader.world.instances:
4             if instance.label_name in ['sofa', 'floor', 'carpet']:
5                 self.shader.world.replace_material(id=instance.id)

```

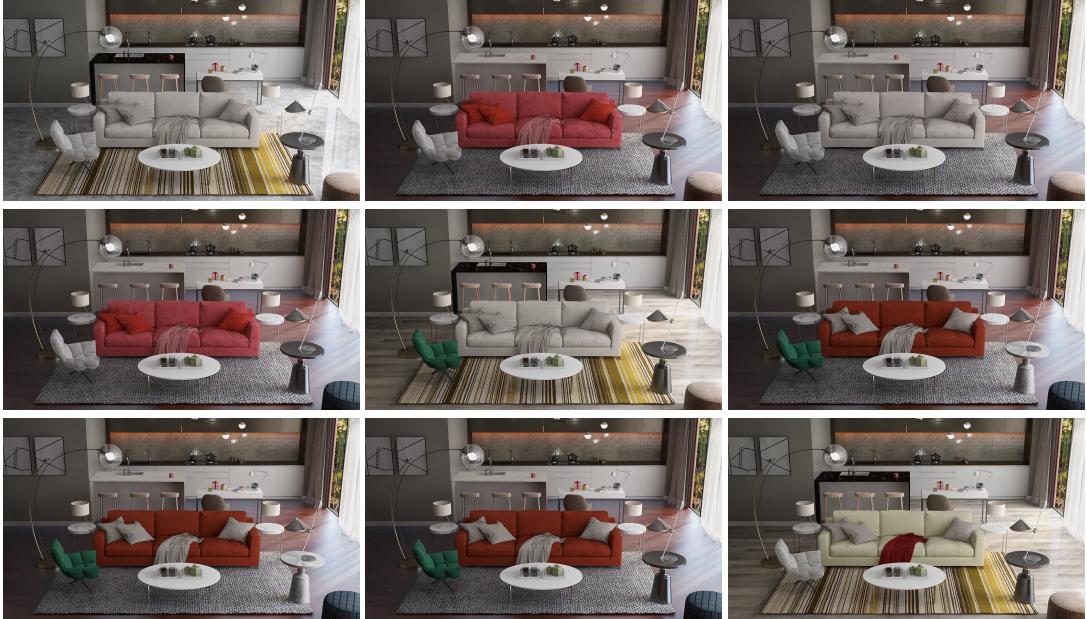


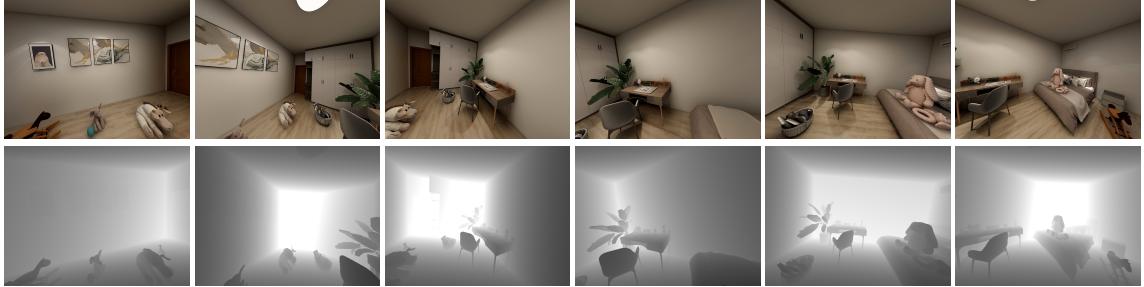
Figure 9. More results generated by the material sampler.

```

1 class CustomTrajectorySampler(EntityProcessor):
2
3     def calculate_vel(self, velocity, step_time):
4         mu = glm.normalize(glm.vec2(np.random.normal(0, 1, size=2)))
5         FORCE = 100
6         f = FORCE * mu
7         PHI, A, C_D, M = 1.204, 0.09, 0.1, 1.0
8         d = -0.5 * glm.normalize(velocity) * PHI * A * C_D * glm.dot(velocity, velocity)
9         velocity = velocity + (d + f) * M * step_time
10        S_MAX = 10
11        if S_MAX < sqrt(glm.dot(velocity, velocity)):
12            velocity = S_MAX * glm.normalize(velocity)
13        return velocity
14
15    def process(self):
16        for camera in self.shader.world.cameras:
17            self.shader.world.delete_entity(camera)
18
19        key_points = []
20        for room in self.shader.world.rooms:
21            # init
22            room_points = []
23            room_polygon = room.gen_polygon()
24            camera_vel = glm.normalize(glm.vec2(np.random.normal(0, 1, size=2)))
25            camera_pos = glm.vec2(room.position)
26            # calculate road points
27            length_of_trajectory, delta_time, scale = 100, 0.03, 1000
28            for i in range(length_of_trajectory):
29                # Camera
30                camera_vel = self.calculate_vel(camera_vel, delta_time)
31                new_position = camera_pos + scale * camera_vel * delta_time
32                next_camera_point = Point(tuple(new_position.xy))
33                if not next_camera_point.within(room_polygon):
34                    p1, p2 = nearest_points(room_polygon, next_camera_point)
35                    normal = glm.normalize(glm.vec2(p1.x - p2.x, p1.y - p2.y))
36                    camera_vel = glm.reflect(camera_vel, normal)
37                    camera_pos = camera_pos + scale * camera_vel * delta_time
38                    room_points.append(list(camera_pos))
39
40        key_points.append(room_points)
41
42        self.make_traj(
43            imageHeight=480,
44            imageWidth=640,
45            keyPoints=key_points,
46            speed=1200,
47            fps=3,
48            speedMode=1,
49            pitchMode=1,
50            pitch=[-10, 10],
51            hfov=70,
52            vfov=55,
53            height=1400,
54            heightMode=1,
55            cameraType="PERSPECTIVE"
56        )

```

Figure 10. DSL for custom trajectory sampler.



(a)



(b)

Figure 11. Results of custom trajectory sampler. (a) Sampled trajectory. (b) 3D reconstruction result by BundleFusion [3].