# Self-Supervised Real-to-Sim Scene Generation

Aayush Prakash        Shoubhik Debnath        Jean-Francois Lafleche

Eric Cameracci        Gavriel State        Stan Birchfield        Marc T. Law

NVIDIA

## Abstract

*Synthetic data is emerging as a promising solution to the scalability issue of supervised deep learning, especially when real data are difficult to acquire or hard to annotate. Synthetic data generation, however, can itself be prohibitively expensive when domain experts have to manually and painstakingly oversee the process. Moreover, neural networks trained on synthetic data often do not perform well on real data because of the domain gap. To solve these challenges, we propose Sim2SG, a self-supervised automatic scene generation technique for matching the distribution of real data. Importantly, Sim2SG does not require supervision from the real-world dataset, thus making it applicable in situations for which such annotations are difficult to obtain. Sim2SG is designed to bridge both the content and appearance gaps, by matching the content of real data, and by matching the features in the source and target domains. We select scene graph (SG) generation as the downstream task, due to the limited availability of labeled datasets. Experiments demonstrate significant improvements over leading baselines in reducing the domain gap both qualitatively and quantitatively, on several synthetic datasets as well as the real-world KITTI dataset.*

## 1. Introduction

Synthetic data, for which annotations can be automatically generated, is a promising solution to overcome the well-known supervised learning bottleneck of labeling real data. For this approach to succeed, such synthetic data must look like real data, in both *appearance* and *content*. Differences in appearance and content together comprise the so-called "domain gap" between synthetic and real data [55, 27]. *Appearance* refers to aspects like color, texture, shape, and lighting, whereas
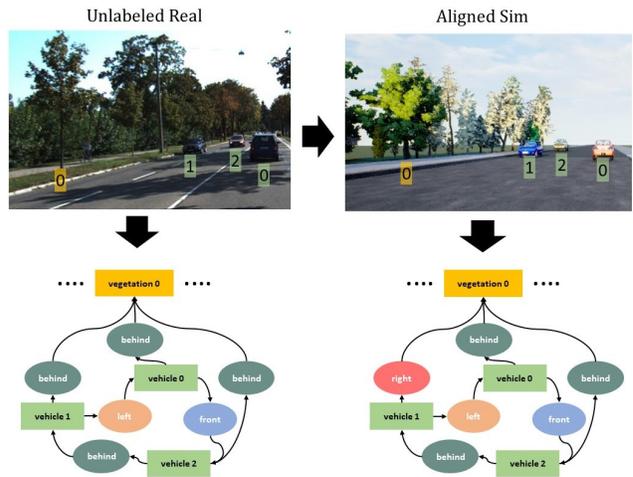


Figure 1. We present Sim2SG, a self-supervised real-to-sim scene generation technique that matches the distribution of real data, for the purpose of training a network to infer scene graphs. Sim2SG does not require costly supervision from the real-world dataset. (Only one tree is shown in the graph to avoid clutter.)

*content* refers to the number, position, and orientation of objects in the scene, as well as their relationships to one another.

At one extreme, synthetic scenes can be created manually by domain experts to reduce these gaps, but such solutions are expensive and therefore do not scale well [61, 47, 45, 18]. At the opposite extreme, domain randomization intentionally leverages these gaps to facilitate sim-to-real transfer [54, 41, 56]; these approaches, however, fail to capture the complexity and distribution of real scenes, thus fundamentally limiting their performance.

Recent approaches such as Meta-Sim [27], Meta-Sim2 [14], and SceneGen [52] attempt to reduce the content gap by automatically learning to generate synthetic data that matches the distribution of real data. However, Meta-Sim can only learn the position and orientation of objects, and therefore cannot align the

structure of the synthetic scenes (e.g., number and types of objects in a scene) to real data. Similarly, Meta-Sim2 only learns the distribution of one type of object (cars), and thus cannot match the surrounding context. Both Meta-Sim and Meta-Sim2 rely on a complex, realistic simulator designed on a set of hand-crafted heuristic rules to ensure that data are generated properly. Scene-Gen addresses these limitations but requires access to a large amount of labeled real data, which undermines the underlying purpose of synthetic data. None of these approaches addresses the appearance gap.

In this paper, we address both the content and appearance gaps, and we do so in a self-supervised manner that requires no real-world labels. Given an unlabeled real dataset, our method aims to automatically generate synthetic data that matches the distribution of the real data. See Figure 1. We propose Sim2SG (Simulation to Scene Graph), a *synthesis-by-analysis* framework, that generates scenes via self-learning [72] loop comprising of two alternating stages: *synthesis* and *analysis*. During the synthesis stage we leverage a synthetic data generator to create scenes. By ensuring that the number of objects, as well as their type and placement, are similar, the synthesized data resembles the real data, and therefore the content gap is reduced. During the analysis stage, we use the generated synthetic data for training. To further reduce both the content and appearance gaps, the corresponding latent and output distributions are aligned via Gradient Reversal Layers [19].

We show the effectiveness of our method in the scene graph generation task [11, 65, 66]. A scene graph (SG) summarizes entities in the scene and plausible relationships among them. The difficulty of hand-labeling scene graphs has limited the community to a small number of datasets [36, 29]. We experimentally demonstrate our method in three distinct environments: synthetic CLEVR [25], synthetic Dining-Sim, and real KITTI [20]. We nearly close the domain gap in the CLEVR environment and show significant improvements over respective baselines in Dining-Sim and KITTI. Through ablations, we validate our contributions regarding appearance and content gaps.

**Contributions:** Our contributions are three-fold: (1) To the best of our knowledge, we are first to do self-supervised aligned scene generation. (2) We propose a novel synthesis-by-analysis framework that addresses both the content and appearance gaps without using any real labels. (3) Experimentally, we show that Sim2SG obtains significant improvements on downstream tasks over baselines in all three scenarios: CLEVR, Dining-Sim, and KITTI. We also present ablations to illustrate the effectiveness of our technique.

## 2. Related Work

**Synthetic Data** has been used for many tasks including, but not limited to, object detection [27, 41, 56], semantic segmentation [45, 47, 58], optical flow modeling [4, 16], scene flow [38], classification [1, 3], stereo [43, 67], 3D keypoint extraction [50], object pose estimation [40, 57] and 3D reconstruction [39]. There are also several simulators [15, 12, 28, 53, 10, 13, 63] available for generating synthetic data. However, to the best of our knowledge, synthetic data has not been applied to scene graph generation.

**Domain Gap** is the performance gap when the network is trained on a synthetic source domain and evaluated on real target data. Most prior work addresses the appearance gap by image translations [8, 17, 23, 24, 32, 71], clever feature alignment [7, 34, 37, 48, 64, 30] and domain randomization [41, 54, 56]. There are few works which handle the content gap [2, 35, 51, 68] by addressing the *label shift* between the two domains. However, they do not exploit the unlabeled images from the target domain. We, on the other hand, leverage the images from the target domain to reduce the domain gap further. To this end, we exploit a scene graph generation task which is more complex than classification and relies on self-training [72]. The idea of self-training with *pseudo labels* is used in [34, 51] to learn models from the target distribution. However, the labels predicted by the model on the target are often inaccurate because of the domain gap [69]. We instead propose to rely on a synthetic data generator to produce accurate labels. Some approaches [6, 49] similar to ours also train their task model on top of domain invariant features for image classification and image segmentation. However, they do not address the content gap.

**Scene Generation** has been used in several machine learning driven approaches. For instance, some methods propose ways to learn how to synthesize indoor scenes [59, 70, 42, 31]. LayoutVAE [26] generates scenes conditioned on label sets. Deep priors and models [60, 46] synthesize scenes by sequentially placing objects. A few techniques [70, 59] also represent scenes as scene graphs with relationships among objects in the same way as our approach. SceneGen [52] learns to generate traffic scenes by modeling the attributes of all objects. However, all these methods either do not address the content gap with real data or use labeled real data for supervision. Meta-Sim [27] learns the distribution of position and rotation of objects in the scene and also requires small labeled real data. Meta-Sim2 [14] additionally learns the distribution of number and type of objects in the scene without needing labeled real data. However, it does not match the distribution of context. Both Meta-Sim and Meta-Sim2 rely on a

simulator designed on handcrafted rules. Unlike Meta-Sim and Meta-Sim2, our method captures relationships among objects and therefore generates more accurate scenes.

## 3. Reducing the domain gap

Let $\langle x_r, y_r \rangle \sim q(x,y)$ be the real data and labels (where the labels $y_r$ are not known). Our goal is to generate synthetic data $\langle x_s, y_s \rangle \sim p(x,y)$ such that the distributions $p$ and $q$ match (*i.e.*, $p \approx q$). We assume that both real and synthetic domains share the same categories of objects as well as scenario (*e.g.*, both are driving scenes). The difference between $p$ and $q$ is the domain gap.

We now study this domain gap between synthetic and real domains. Let $\phi$ be a network that encodes an input image $x$ into a latent representation $z \in \mathcal{Z}$, and let $h$ be a network that estimates some desired quantity from $z$. (In our case, $h$ infers a scene graph.) We consider as in [62] that the task error in the real domain, $\epsilon_r(\phi, h)$, is a function of three terms:

$$\epsilon_r(\phi, h) = \int q(z) e_r \, dz$$
$$= \underbrace{\int p(z) e_s dz}_{\epsilon_s(\phi,h)} + \underbrace{\int q(z)(e_r - e_s)dz}_{\epsilon^c(\phi,h)} + \underbrace{\int (q(z) - p(z)) e_s dz}_{\epsilon^a(\phi,h)}$$
(1)

where $e_r \equiv |h(\phi(x_r)) - y_r|$ is the real risk, and similarly $e_s \equiv |h(\phi(x_s)) - y_s|$ is the synthetic risk; and where the distributions of shared features for synthetic and real domains are denoted by $p(z)$ and $q(z)$, respectively. In this equation, $\epsilon_s(\phi, h)$ is the training error on the synthetic domain, $\epsilon^c(\phi, h)$ is the risk gap between the domains, and $\epsilon^a(\phi, h)$ is the feature gap between the two domains. In the following, we drop the terms $\phi, h$ for notational simplicity.

If the error $\epsilon_r$ reduces to zero on the target domain, we have closed the domain gap. To reduce this error, we must reduce both $\epsilon^c$ and $\epsilon^a$. We call the former the *content gap* because it refers to the difference in the distributions of the two domains regarding the number of objects and their class, placement, pose and scale. The difficulty of minimizing the content gap is that computing it is intractable since we do not have access to labels $y_r$ from the real domain, and hence $e_r$ cannot be computed. We call the gap $\epsilon^a$ the *appearance gap* because it refers to differences in texture, color, lighting, and reflectance of the scene. Recent work [21] also shows that latent representations $z$ frequently account for appearance. In the following, we describe our approach to reduce both the content and appearance gaps.

---

**Algorithm 1** Pseudocode for Sim2SG training
1: **Given: generator**, $X_r$     ▷ Data generator, real images
2: **while** not converged **do**
3:     ▷ Synthesis
4:     $G \leftarrow \{h(\phi(x_r)) : x_r \in X_r\}$     ▷ Scene graphs
5:     $X_s, Y_s \leftarrow \texttt{generator}(G)$ ▷ Generate aligned synthetic data
6:     ▷ Analysis
7:     $loss \leftarrow 0$
8:     **for** $(x_s, y_s; x_r) \in (X_s, Y_s; X_r)$ **do**
9:         $loss \mathrel{+}= \mathcal{L}^{\texttt{a}}(\phi(x_s), \phi(x_r))$     ▷ Appearance gap
10:       $loss \mathrel{+}= \mathcal{L}^{\texttt{c}}(h(\phi(x_s)), h(\phi(x_r)))$     ▷ Prediction gap
11:       $loss \mathrel{+}= \texttt{task}(x_s, y_s)$     ▷ SG prediction task loss
12:       $\phi, h \leftarrow \texttt{optimize}(\phi, h, loss)$     ▷ SGD step
13:     **end for**
14: **end while**

---

## 4. Simulation to Scene Graph (Sim2SG)

Figure 2 illustrates our proposed Sim2SG pipeline, which comprises two alternating steps. During synthesis (real-to-sim), synthetic data is generated to match the distribution of the unlabeled real data. During analysis (sim-to-real), a scene graph (SG) prediction network [65] is trained using the ground truth labels of the synthetic data. This is analogous to Expectation-Maximization, where synthesis is like the E-step, and analysis is like the M-step. Results from different iterations of this self-learning loop, which is initialized using synthetic data generated by structured domain randomization (SDR) [41], are shown in Figure 3.

We now describe the synthesis and analysis steps. Algorithm 1 illustrates the pseudocode of both steps.

### 4.1. Synthesis Step

In this step, we aim to generate synthetic data that matches the distribution of unlabeled real data, thereby reducing the content gap. As mentioned previously, minimizing $\epsilon^c$ in the current form is not tractable. Nevertheless, a sufficient condition for $\epsilon^c$ to be zero is for $y_s = y_r$ and $h(\phi(x_r)) = h(\phi(x_s))$. This leads to splitting the content gap into two terms:

$$\epsilon^c = \int q(z)(e_r - e_s)dz$$
$$\simeq \underbrace{\int q(z)(y_s - y_r)dz}_{\epsilon^{c,label}} + \underbrace{\int q(z)(h(\phi(x_r)) - h(\phi(x_s))dz}_{\epsilon^{c,pred}}$$
(2)

where $\epsilon^{c,label}$ refers to the gap between the ground truth labels of synthetic and real data, and where $\epsilon^{c,pred}$ refers to the gap between the output of the network $h$ trained on synthetic and real data.

Since we do not have access to the labels $y_r$ of the target (real) domain, we propose to estimate $y_r$ through
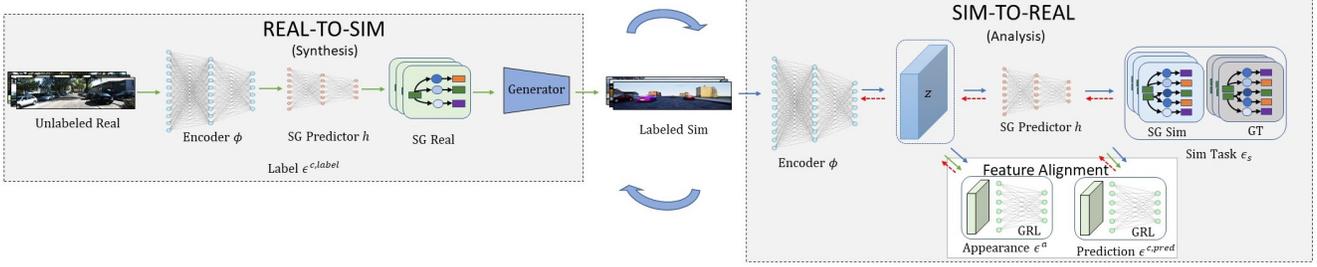
Figure 2. Overview of Sim2SG, a self-supervised synthesis-by-analysis framework that generates synthetic data in a loop comprising of two alternating stages: synthesis and analysis. During synthesis (real-to-sim), we infer scene graphs from real data to generate synthetic data that matches the distribution of real data, thus bridging the content label gap, $\epsilon^{c,label}$. During analysis (sim-to-real), we train a scene graph prediction network ($\phi$, $h$) using the synthetic data. Additionally, Gradient Reversal Layers (GRLs) are used to align the latent features $z$ and output features $h(z)$ to bridge the appearance $\epsilon^a$ and content prediction $\epsilon^{c,pred}$ gap, respectively.


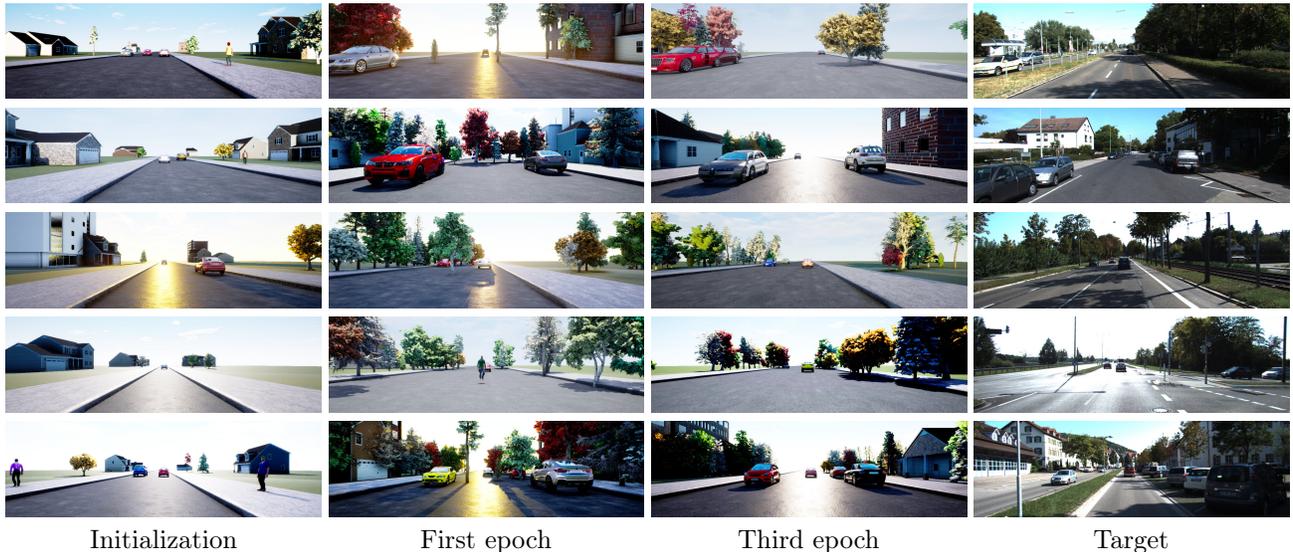
| Initialization | First epoch | Third epoch | Target |

Figure 3. As the self-learning loop progresses, synthetic data increasingly matches the content of the target image. From left to right: synthetic data initialized by SDR [41], after first and third epochs, and corresponding real KITTI target images.

a self-supervised method that infers scene graphs from the target data. Our scene graphs consist of nodes and edges. Each node $o_i = \langle b_i, c_i \rangle$ consists of a bounding box $b_i = \{u_i, v_i, w_i, h_i\}$ (center, width, and height of the 2D box) and category $c_i$ (such as car, pedestrian, building, etc.). Each edge captures a relationship $r_i = \langle o_i, p, o_j \rangle$ where $p$ is a predicate (such as behind, left, on, etc.). Using known camera parameters and a flat ground plane assumption, the scene graph can be mapped to a full 3D scene. Some parameters (*e.g.*, texture or pose) and some context (ground, sky, lighting), that are not part of the scene graph are randomized. One advantage of generating synthetic data is that we can augment the relationships of the inferred scene graph by adding new relationships, simply by reasoning in the 3D representation. A synthetic data generator is then used to render the 3D scenes. Note that the weights of our networks ($\phi$, $h$) remain unchanged during this step.

## 4.2. Analysis Step

In this step, we train the encoder $\phi$ and predictor $h$ of our scene graph (SG) generation model on the synthetic data generated by the preceding step. The task loss function is from Yang *et al.* [65], which includes cross entropy loss for object classification and relationship classification, and $\ell_1$ loss for bounding boxes.

Although using synthetic data aligns the content label gap $\epsilon^{c,label}$, it does not align the appearance gap $\epsilon^a$ nor content prediction gap $\epsilon^{c,pred}$. To address the appearance gap, we align the feature distributions $p(z)$ and $q(z)$, since $p(z) = q(z)$ is a sufficient condition for $\epsilon^a$ to be zero [62]. To do this, during training $z$ is passed through a Gradient Reversal Layer (GRL) [19, 7], followed by a domain classifier $D^a$ that learns to classify the input as synthetic or real. The GRL acts as an identity function during forward propagation, but flips the sign of the gradients during back propagation. The

classifier $D^a$ provides an additional loss term for training $\phi$. We minimize $D^a$'s loss w.r.t. its own parameters while maximizing w.r.t. the network parameters of $\phi$. The loss function for training is as follows:

$$\mathcal{L}^a = -\sum_x \Big[ d_i \log D^a(\phi(x)) + (1-d_i)\log(1-D^a(\phi(x))) \Big] \tag{3}$$

where $x \in x_s, x_r$, $d_i = 0$ for synthetic images $x_s$ and $d_i = 1$ for real images $x_r$.

Similarly, we seek to minimize the content prediction gap $\epsilon^{c,pred}$ by matching the distributions of the outputs $h(\phi(x_s))$ and $h(\phi(x_r))$ using the same GRL-based technique. This is based on the observation that the output of the scene graph generation model should be similar for corresponding inputs from the synthetic and real domains. During training, the output $h(z)$ is passed through a GRL, followed by a domain classifier $D^c$. The classifier provides an additional loss term for training both $\phi$ and $h$. The loss function is computed as:

$$\mathcal{L}^c = -\sum_z \Big[ d_i \log D^c(h(z)) + (1-d_i)\log(1-D^c(h(z))) \Big] \tag{4}$$

where $z \in \phi(x_s), \phi(x_r)$. Fully matching the source and target appearance distributions without regard to content may cause detrimental results [48, 62]. We introduce a warm-up period during which content is aligned without regard to appearance or content prediction.

## 5. Experiments

We evaluate Sim2SG in three different environments with increasing complexity: CLEVR [25], our own Dining-Sim using ShapeNet [5], and KITTI [20]. In each environment we have a fully labeled synthetic domain and unlabeled target domain with labeled test data. We use scene graph (SG) generation as the downstream task and implement the encoder $\phi$ using ResNet-101 [22] and the predictor $h$ using Graph R-CNN [65].

Quantitative evaluation metrics include detection mAP (mean average precision) @0.5 IoU (Intersection over Union in 2D) and relationship triplet recall @20 or @50 [29]. Note that relationship triplet recall implicitly includes object detection recall. All the mean and standard deviations are based on five runs, and all results are reported after saturation. Details of the environments, training parameters, and hyperparameters are in the appendix.

For notational simplicity, we use $\sigma^{c,label}$ to refer to our synthesis step (whose purpose is to reduce the content label gap $\epsilon^{c,label}$), $\sigma^a$ to refer to the first GRL (whose purpose is to reduce the appearance gap $\epsilon^a$), and $\sigma^{c,pred}$ to refer to the second GRL (whose purpose is to

| Method | mAP@0.5 IoU | Recall@20 |
|---|---|---|
| SDR [41] | 0.723 ±0.053 | 0.356 ±0.047 |
| Ours ($\sigma^{c,label}$) | 0.832 ±0.046 | 0.493 ±0.064 |
| Ours ($\sigma^a$) | 0.821 ±0.048 | 0.815 ±0.026 |
| Ours ($\sigma^a$, $\sigma^{c,label}$) | **0.892** ±0.024 | **0.888** ±0.018 |

Table 1. Results of Sim2SG on the CLEVR target domain. Aligning both appearance and content yields the best results.

reduce the content prediction gap $\epsilon^{c,pred}$). Similarly, we use $\sigma^c$ to refer to the combination of $\sigma^{c,label}$ and $\sigma^{c,pred}$, whose collective purpose is to reduce the content gap.

### 5.1. Experiments on the CLEVR dataset

The goal of the experiments on the CLEVR environment [25] is to show that Sim2SG can address the domain gap in a simple controlled environment. There are three classes of objects: cube, sphere and cylinder; and four relationships: front, behind, left and right. The source and target domains were generated using disjoint object colors, object materials, margin between objects, and number of objects, to ensure a significant domain gap. Furthermore, the target images have more complex texture via application of a style transfer network. We use 1000 labeled images for source domain, 1000 unlabeled images for target domain for training, and 200 labeled source and target images for evaluation.

**Results:** Quantitative evaluation of Sim2SG is reported in Table 1. Compared with the baseline using SDR [41], our techniques for label alignment $\sigma^{c,label}$ and appearance alignment $\sigma^a$ *drastically reduce* the domain gap. The best results are achieved by combining these two ideas. With such a simple environment, however, we found $\sigma^{c,pred}$ to have negligible effect, and therefore it is not reported. Qualitative results of scene graph recall are shown in first row of Figure 4.

**Ablations:** We conduct two ablation experiments on the CLEVR dataset to confirm that appearance alignment $\sigma^a$ and label alignment $\sigma^{c,label}$ work as intended. In the first experiment, the source and target domains are created with the same color and texture but different number and placement of objects. Thus, there is a content gap, but no appearance gap. We observe that label alignment $\sigma^{c,label}$ closes the domain gap from 0.76 to 0.996 for Recall@20, whereas $\sigma^a$ leads to performance degradation. In the second experiment, the source and target domains have the same number of objects and placement but use different color and texture. Thus, there is an appearance gap, but no content gap. We observe that appearance alignment $\sigma^a$ reduces the domain gap, achieving 0.938 versus 0.339 for Recall@20 for the baseline, whereas label alignment $\sigma^{c,label}$ fails to have significant improvement on rela-
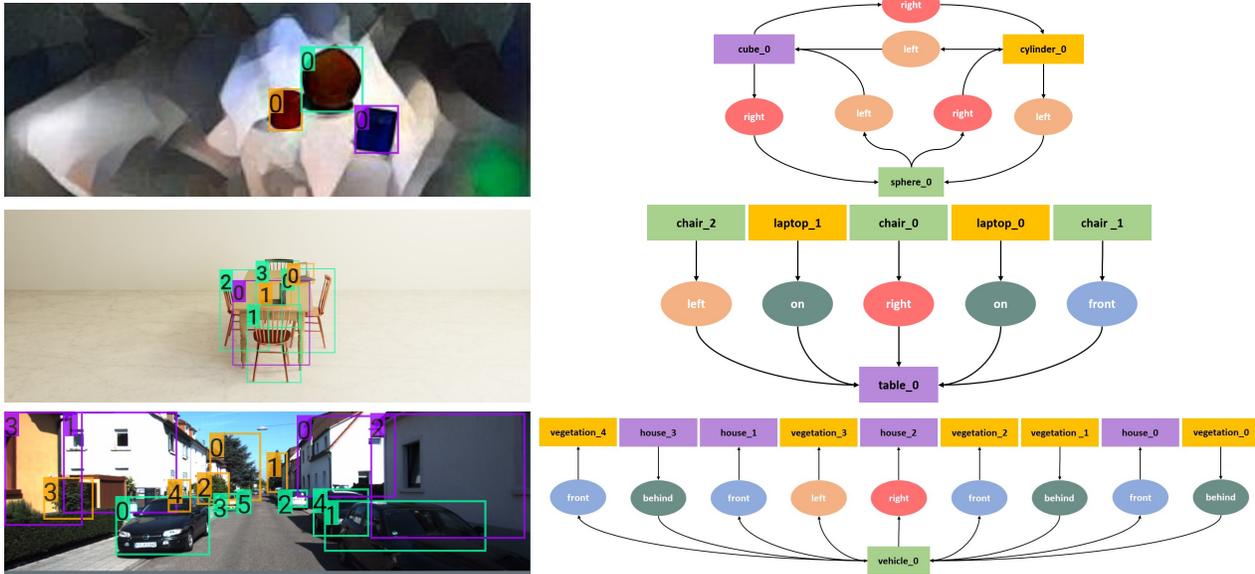
Figure 4. Qualitative results of Sim2SG on the target domain for CLEVR (top row), Dining Sim (middle row) and KITTI environments (bottom row). Objects are color coded. For better visibility, we only show partial scene graph for KITTI.

| Method | mAP@0.5 IoU | Recall@50 |
|---|---|---|
| SDR [41] | 0.584 $\pm$0.049 | 0.331 $\pm$0.064 |
| Ours ($\sigma^{c,label}$) | 0.713 $\pm$0.038 | 0.501 $\pm$0.044 |
| Ours ($\sigma^c$, $\sigma^a$) | **0.729** $\pm$0.015 | **0.547** $\pm$0.015 |

Table 2. Results of Sim2SG on the Dining-Sim target domain. In the last row, $\sigma^c$ includes both $\sigma^{c,label}$ and $\sigma^{c,pred}$.

tionship triplet recall. These experiments show that our label alignment procedure $\sigma^{c,label}$ reduces the content gap, while appearance alignment $\sigma^a$ addresses the appearance gap.

### 5.2. Experiments on Dining-Sim

We created a dataset that we call Dining-Sim by placing objects created from ShapeNet objects [5] in realistic arrangements, with more complex textures and realistic lighting. This dataset has three classes of objects: chair, table and laptop; and there are five relationships: front, behind, left, right, and on. Quantitative results, shown in Table 2, agree with the findings of the previous section. Since this data is more complex, the warm-up period mentioned earlier is necessary, and therefore the label alignment $\sigma^{c,label}$ must be performed first. This label alignment drastically improves performance on the target domain, and the combination of label alignment $\sigma^{c,label}$, appearance alignment $\sigma^a$, and prediction alignment $\sigma^{c,pred}$ achieve the best results. For reference, the oracle performance on the target domain is 0.904 mAP@0.5 IoU and 0.846 Recall@50. Qualitative results are illustrated in the second row of Figure 4.

### 5.3. Real-world experiments on KITTI

In this section, we validate our approach on the real-world KITTI dataset [20]. For the synthetic domain, we implemented from scratch a simplified version of SDR [41], with a fixed camera, a subset of object classes, no post processing, and without curved road splines, all for faster data generation. The number of lanes, sidewalks, and various objects, along with their positions, pose, color, texture, as well as lighting settings are randomly selected. Note that, unlike SDR, heuristic rules are not used to avoid collisions or ensure realistic placement, which further demonstrates the power of our proposed method in automatically generating useful training data. We use four object classes: car, pedestrian, vegetation, house, and four relationships: front, left, right, behind. Relationships are constrained so that at least one node is a car, e.g., *car behind car*, *vegetation left car*, and so on. Although we have shown the 'on' relationship to work in the Dining-Sim environment, such relationships are trivial to predict in this driving environment because they are always true: cars are always *on* the road, pedestrians are *on* the sidewalk, and so forth. Therefore, we did not include them in the experiments.

A subset of the real KITTI data [20] from the 2D object detection suite are used as the target domain. We use 6000 unlabeled KITTI images for training, and 550 labeled KITTI images for evaluation. The latter images include not only the provided annotations for cars and pedestrians, but also annotations that we added for

**Figure 5.** Qualitative results of objects detected on three different KITTI images. Left: SDR fails to detect many objects and yields a large number of false positives (mislabels), leading to poor scene graphs (not shown). Middle: Meta-Sim improves on false-positives, but still fails to detect some objects. Right: Our method detects objects correctly with fewer false positives, thus generating more accurate scene graphs. (Cars in green, vegetation in yellow, buildings in purple.)

vegetation and houses, along with relationships.[1] For training, we used 6000 labeled synthetic images, and another 1000 labeled synthetic images for validation.

**Baselines**: We compare against methods for which code is publicly available. In particular, we compare against unsupervised scene-generation methods, including structured domain randomization (SDR) [41] and Meta-Sim [27]. We also compare against a popular self-learning method based on pseudo labels extracted by training directly on the real KITTI data after pre-training on SDR synthetic data [72]. Furthermore, we compare against unsupervised domain adaptation methods for object detection, which are based on aligning features from the source and target domains: DA Faster R-CNN [7], GPA [64], and SAPNet [30]. According to a recent survey paper [33], the latter two are state-of-the-art leading techniques. All baselines are trained on 6000 images from SDR, using the hyperparameters provided by the original authors.

**Results:** We evaluate these baselines, along with our method, on three KITTI evaluation modes for 2D object detection: easy, moderate and hard, which are based on object size, occlusion and truncation. For all three modes, our method yields significantly improved results over the baselines. Table 3 shows the object detection and scene graph generation results for KITTI hard; other results are in the appendix. The last three rows of the table show that most improvements come from label alignment $\sigma^{c,label}$ and appearance alignment $\sigma^a$. The combination of all three, $\sigma^{c,label}$, $\sigma^a$ and $\sigma^{c,pred}$, achieves the best results overall. Note that $\sigma^{c,label}$ alone is able to beat most baselines, thus demonstrating the

efficacy of our synthesis step at automatically generating realistic training data without heuristics. Qualitative results are illustrated in the third row of Figure 4.

Insight as to why our method outperforms SDR [41] and Meta-Sim [27] can be gained by viewing the images generated by the various methods, shown in Figure 6. Our method generates synthetic data that better matches the distribution of the real data, because SDR lacks access to the object distributions in KITTI, and because Meta-Sim cannot align the structure of the scenes (e.g., the number of objects) and lacks the notion of relationships in its scene graphs. (Note that Meta-Sim2 [14] also lacks relationships, but code is not publicly available.) Furthermore, our method scales better with scene complexity compared with Meta-Sim, which requires passing expensive numerical gradients through a renderer. As a result, our training time is 12 hours on a single NVIDIA V100 GPU, compared with 72 hours for Meta-Sim. Figure 5 compares our method on the downstream task with SDR and Meta-Sim.

The reason our method outperforms domain adaptation techniques [7, 64, 30] is because feature alignment without content alignment is not effective, as we discussed briefly at the end of Section 3. Nevertheless, it may be possible to combine our label alignment technique $\sigma^{c,label}$ with domain adaptation methods, which we leave for future work. Our method always produces proper ground truth for the scene via synthetic data generator, even if the scene is not generated exactly according to the corresponding real image. Self-learning [72] methods based on pseudolabels from KITTI, on the other hand, can have errors in the ground truth due to incorrectly classified objects or imprecise bounding boxes, as explained in [69].

---

[1]Our annotations are publicly available at https://research.nvidia.com/publication/2021-08_Sim2SG

| Method | Car | Pedestrian | House | Vegetation | mAP@0.5IoU | Recall@50 |
|---|---|---|---|---|---|---|
| SDR [41] | 0.382 ±0.029 | 0.168 ±0.017 | 0.211 ±0.023 | 0.174 ±0.010 | 0.234 ±0.006 | 0.070 ±0.007 |
| Meta-Sim [27] | 0.413 ±0.009 | 0.197 ±0.027 | 0.236 ±0.009 | 0.164 ±0.023 | 0.253 ±0.003 | 0.075 ±0.005 |
| Self-learning [72] | 0.312 ±0.006 | 0.167 ±0.015 | 0.191 ±0.003 | 0.263 ±0.006 | 0.233 ±0.004 | 0.062 ±0.003 |
| DA Faster R-CNN [7] | 0.424 ±0.028 | 0.170 ±0.029 | 0.200 ±0.041 | 0.169 ±0.014 | 0.241 ±0.014 | 0.074 ±0.015 |
| GPA [64] | 0.174 ±0.040 | 0.011 ±0.016 | 0.106 ±0.031 | 0.059 ±0.027 | 0.087 ±0.020 | 0.015 ±0.005 |
| SAPNet [30] | 0.362 ±0.054 | 0.085 ±0.051 | 0.116 ±0.021 | 0.067 ±0.022 | 0.157 ±0.024 | – |
| Ours ($\sigma^{c,label}$) | 0.410 ±0.009 | **0.262** ±0.025 | 0.240 ±0.010 | 0.229 ±0.036 | 0.285 ±0.003 | 0.104 ±0.006 |
| Ours ($\sigma^{c,label}$, $\sigma^a$) | 0.493 ±0.004 | 0.252 ±0.014 | 0.247 ±0.012 | 0.253 ±0.020 | 0.311 ±0.311 | 0.127 ±0.004 |
| Ours ($\sigma^{c,label}$, $\sigma^a$, $\sigma^{c,pred}$) | **0.501** ±0.006 | 0.241 ±0.018 | **0.254** ±0.010 | **0.269** ±0.014 | **0.316** ±0.004 | **0.139** ±0.004 |

Table 3. Results on KITTI hard after training on labeled synthetic data and unlabeled real data. The class specific AP values for 2D object detection are reported at 0.5 IoU. The last column shows relationship triplet recall for scene graph generation.
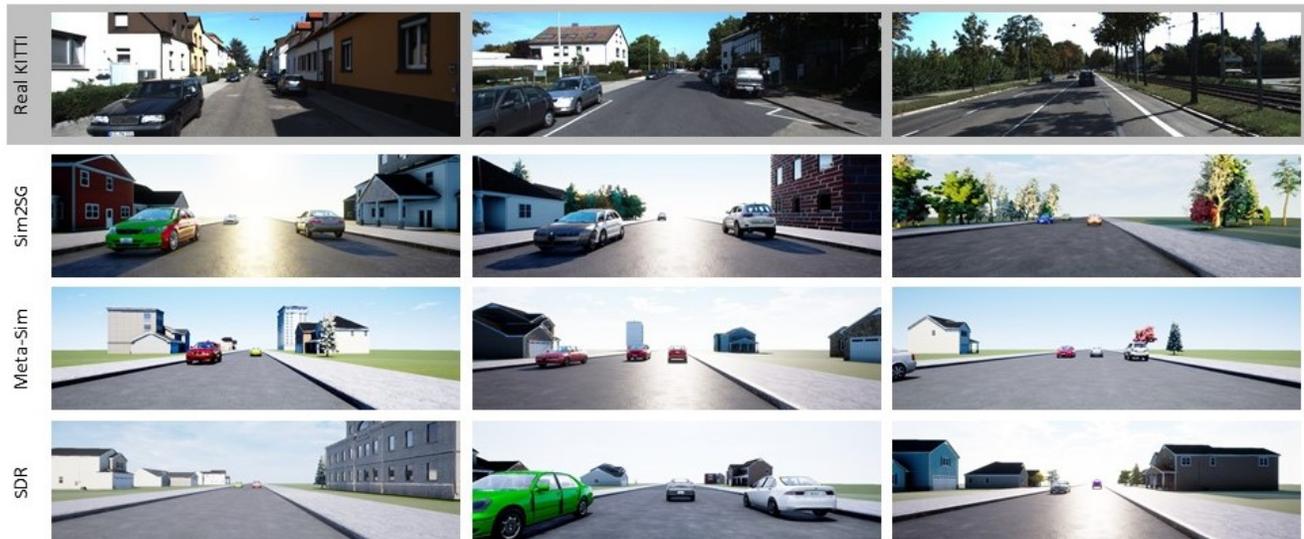


Figure 6. Synthetic images generated by Sim2SG (ours), Meta-Sim [27], and SDR [41] with the corresponding KITTI samples. Our method aligns the number and placement of both cars and context (vegetation, houses) better than other methods.

Note in Table 3 that the detection accuracy of the pedestrian category does not improve with $\sigma^a$ and $\sigma^{c,pred}$. The reason for this limitation is that pedestrians are an under-represented class in KITTI, not to mention small and hard to detect. While our Sim2SG method can align the label distribution, it cannot address class imbalance in the target domain. Nevertheless, our method is not restricted to the types or number of relations. Given a simulator to generate scenarios such that objects and their relations are detectable, our method should extend to handle them.

**Ablations:** As briefly discussed at the end of Section 3, we run the label alignment $\sigma^{c,label}$ before appearance alignment $\sigma^a$ and prediction alignment $\sigma^{c,pred}$ to address the fact that feature alignment can be detrimental if the content of both domains are not already aligned. We indeed found that performance drops significantly when we train Sim2SG without $\sigma^{c,label}$ and evaluate in the same setting as Table 3. Sim2SG with $\sigma^a$ and $\sigma^{c,pred}$ gives a 0.246 mAP@0.5 IoU for detection and

0.076 Recall@50 for relationship triplets, while adding $\sigma^{c,label}$ yields a significant boost of 0.316 mAP@0.5 IoU for detection and 0.139 Recall@50 for relationship triplets on KITTI Hard. These results confirm the effectiveness of $\sigma^{c,label}$ and therefore the importance of the entire Sim2SG framework.

## 6. Conclusion

In this work, we have proposed Sim2SG, a self-supervised real-to-sim automatic scene generation technique that matches the distribution of real data, for the purpose of training a network to infer scene graphs. The method bridges both the content and appearance gap with real data without requiring any costly supervision on the real-world dataset. The approach achieves significant improvements over baselines in all three environments for which we tested: CLEVR, a Dining-Sim dataset, and real KITTI data. The latter of these demonstrates the ability of our method to perform sim-to-real transfer.

# References

[1] David Acuna, Guojun Zhang, Marc T. Law, and Sanja Fidler. f-domain adversarial learning: Theory and algorithms. In *International Conference on Machine Learning (ICML)*, 2021. 2

[2] Kamyar Azizzadenesheli, Anqi Liu, Fanny Yang, and Animashree Anandkumar. Regularized learning for domain adaptation under label shifts. In *International Conference on Learning Representations (ICLR)*, 2019. 2

[3] João Borrego, Atabak Dehban, Rui Figueiredo, Plinio Moreno, Alexandre Bernardino, and José Santos-Victor. Applying domain randomization to synthetic data for object category detection. *arXiv:1807.09834*, 2018. 2

[4] Daniel J. Butler, Jonas Wulff, Garrett B. Stanley, and Michael J. Black. A naturalistic open source movie for optical flow evaluation. In *ECCV*, 2012. 2

[5] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An information-rich 3D model repository. *arXiv:1512.03012*, 2015. 5, 6, 12

[6] Wei-Lun Chang, Hui-Po Wang, Wen-Hsiao Peng, and Wei-Chen Chiu. All about structure: Adapting structural information across domains for boosting semantic segmentation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[7] Yuhua Chen, Wen Li, Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Domain adaptive faster r-cnn for object detection in the wild. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2, 4, 7, 8, 18, 19

[8] Yun-Chun Chen, Yen-Yu Lin, Ming-Hsuan Yang, and Jia-Bin Huang. CrDoCo: Pixel-level domain transfer with cross-domain consistency. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[9] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, 2018. 12

[10] Adam Crespi, Cesar Romero, Srinivas Annambhotla, Jonathan Hogins, and Alex Thaman. Unity perception, 2020. https://blogs.unity3d.com/2020/06/10/. 2

[11] Bo Dai, Yuqi Zhang, and Dahua Lin. Detecting visual relationships with deep relational networks. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2

[12] Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha Kembhavi, Eric Kolve, Roozbeh Mottaghi, Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, Luca Weihs, Mark Yatskar, and Ali Farhadi. RoboTHOR: An Open Simulation-to-Real Embodied AI Platform. In *CVPR*, 2020. 2

[13] Maximilian Denninger, Martin Sundermeyer, Dominik Winkelbauer, Youssef Zidan, Dmitry Olefir, Mohamad Elbadrawy, Ahsan Lodhi, and Harinandan Katam. Blenderproc. *arXiv:1911.01911*, 2019. 2

[14] Jeevan Devaranjan, Amlan Kar, and Sanja Fidler. Meta-Sim2: Unsupervised learning of scene structure for synthetic data generation. In *ECCV*, 2020. 1, 2, 7

[15] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *CORL*, 2017. 2

[16] Philipp Fischer, Alexey Dosovitskiy, Eddy Ilg, Philip Häusser, Caner Hazırbaş, Vladimir Golkov, Patrick van der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *IEEE International Conference on Computer Vision (ICCV)*, 2015. 2

[17] Geoffrey French, Michal Mackiewicz, and Mark Fisher. Self-ensembling for visual domain adaptation. In *ICLR*, 2018. 2

[18] Adrien Gaidon, Qiao Wang, Yohann Cabon, and Eleonora Vig. Virtual worlds as proxy for multi-object tracking analysis. In *CVPR*, 2016. 1

[19] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(1):1–35, 2016. 2, 4

[20] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In *CVPR*, 2012. 2, 5, 6, 15, 19

[21] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations (ICLR)*, 2019. 3

[22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2016. 5, 12

[23] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isol, Kate Saenko Alexei A. Efros, and Trevor Darrell. CyCADA: Cycle-consistent adversarial domain adaptation. In *ICML*, 2018. 2

[24] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *ECCV*, 2018. 2

[25] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, 2017. 2, 5, 12

[26] Akash Abdu Jyothi, Thibaut Durand, Jiawei He, Leonid Sigal, and Greg Mori. LayoutVAE: Stochastic scene layout generation from a label set. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 2

[27] Amlan Kar, Aayush Prakash, Ming-Yu Liu, Eric Cameracci, Justin Yuan, Matt Rusiniak, David Acuna, Antonio Torralba, and Sanja Fidler. Meta-Sim: Learning to generate synthetic datasets. *IEEE/CVF Interna-*

*tional Conference on Computer Vision (ICCV)*, 2019. 1, 2, 7, 8, 18, 19

[28] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. AI2-THOR: An interactive 3D environment for visual AI. *arXiv:1712.05474*, 2017. 2

[29] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, and et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73, Feb 2017. 2, 5

[30] Congcong Li, Dawei Du, Libo Zhang, Longyin Wen, Tiejian Luo, Yanjun Wu, and Pengfei Zhu. Spatial attention pyramid network for unsupervised domain adaptation. In *ECCV*, 2020. 2, 7, 8, 18, 19

[31] Manyi Li, Akshay Gadi Patil, Kai Xu, Siddhartha Chaudhuri, Owais Khan, Ariel Shamir, Changhe Tu, Baoquan Chen, Daniel Cohen-Or, and Hao Zhang. Grains: Generative recursive autoencoders for indoor scenes. *ACM Transactions on Graphics (TOG)*, 38(2):1–16, 2019. 2

[32] Peilun Li, Xiaodan Liang, Daoyuan Jia, and Eric P. Xing. Semantic-aware Grad-GAN for virtual-to-real urban scene adaption. In *BMVC*, 2018. 2

[33] Wanyi Li, Fuyu Li, Yongkang Luo, Peng Wang, et al. Deep domain adaptive object detection: a survey. In *IEEE Symposium Series on Computational Intelligence (SSCI)*, 2020. 7

[34] Yunsheng Li, Lu Yuan, and Nuno Vasconcelos. Bidirectional learning for domain adaptation of semantic segmentation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[35] Zachary C. Lipton, Yu-Xiang Wang, and Alex Smola. Detecting and correcting for label shift with black box predictors. In *ICML*, 2018. 2

[36] Cewu Lu, Ranjay Krishna, Michael Bernstein, and Li Fei-Fei. Visual relationship detection with language priors. In *European Conference on Computer Vision (ECCV)*, 2016. 2

[37] Yawei Luo, Liang Zheng, Tao Guan, Junqing Yu, and Yi Yang. Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[38] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *CVPR*, 2016. 2

[39] John McCormac, Ankur Handa, Stefan Leutenegger, and Andrew J Davison. SceneNet RGB-D: 5M photorealistic images of synthetic indoor trajectories with ground truth. In *ICCV*, 2017. 2

[40] Matthias Mueller, Vincent Casser, Jean Lahoud, Neil Smith, and Bernard Ghanem. Sim4CV: A photo-realistic simulator for computer vision applications. In *arXiv:1708.05869*, 2017. 2

[41] Aayush Prakash, Shaad Boochoon, Mark Brophy, David Acuna, Eric Cameracci, Gavriel State, Omer Shapira, and Stan Birchfield. Structured domain randomization: Bridging the reality gap by context-aware synthetic data. In *International Conference on Robotics and Automation (ICRA)*, 2019. 1, 2, 3, 4, 5, 6, 7, 8, 14, 15, 16, 17, 18, 19

[42] Siyuan Qi, Yixin Zhu, Siyuan Huang, Chenfanfu Jiang, and Song-Chun Zhu. Human-centric indoor scene synthesis using stochastic grammar. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[43] Weichao Qiu and Alan Yuille. UnrealCV: Connecting computer vision to Unreal Engine. In *European Conference on Computer Vision (ECCV)*, 2016. 2

[44] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015. 12

[45] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *ECCV*, 2016. 1, 2

[46] Daniel Ritchie, Kai Wang, and Yu-an Lin. Fast and flexible indoor scene synthesis via deep convolutional generative models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[47] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio Lopez. The SYNTHIA Dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *CVPR*, 2016. 1, 2

[48] Kuniaki Saito, Yoshitaka Ushiku, Tatsuya Harada, and Kate Saenko. Strong-weak distribution alignment for adaptive object detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 5

[49] Yu Sun, Eric Tzeng, Trevor Darrell, and Alexei A. Efros. Unsupervised domain adaptation through self-supervision. *arXiv:1909.11825*, 2019. 2

[50] Supasorn Suwajanakorn, Noah Snavely, Jonathan Tompson, and Mohammad Norouzi. Discovery of latent 3D keypoints via end-to-end geometric reasoning. *NeurIPS*, 2018. 2

[51] Shuhan Tan, Xingchao Peng, and Kate Saenko. Class-imbalanced domain adaptation: An empirical odyssey, 2020. 2

[52] Shuhan Tan, Kelvin Wong, Shenlong Wang, Sivabalan Manivasagam, Mengye Ren, and Raquel Urtasun. SceneGen: Learning to generate realistic traffic scenes. *arXiv:2101.06541*, 2021. 1, 2

[53] Thang To, Jonathan Tremblay, Duncan McKay, Yukie Yamaguchi, Kirby Leung, Adrian Balanon, Jia Cheng,

and Stan Birchfield. NDDS: NVIDIA deep learning dataset synthesizer, 2018. https://github.com/NVIDIA/Dataset_Synthesizer. 2

[54] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *IROS*, 2017. 1, 2

[55] Antonio Torralba and Alexei A. Efros. Unbiased look at dataset bias. In *CVPR 2011*, 2011. 1

[56] Jonathan Tremblay, Aayush Prakash, David Acuna, Mark Brophy, Varun Jampani, Cem Anil, Thang To, Eric Cameracci, Shaad Boochoon, and Stan Birchfield. Training deep networks with synthetic data: Bridging the reality gap by domain randomization. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2018. 1, 2

[57] Jonathan Tremblay, Thang To, Balakumar Sundaralingam, Yu Xiang, Dieter Fox, and Stan Birchfield. Deep object pose estimation for semantic robotic grasping of household objects. In *CoRL*, 2018. 2

[58] Apostolia Tsirikoglou, Joel Kronander, Magnus Wrenninge, and Jonas Unger. Procedural modeling and physically based rendering for synthetic data generation in automotive applications. *arXiv:1710.06270*, 2017. 2

[59] Kai Wang, Yu-An Lin, Ben Weissmann, Manolis Savva, Angel X Chang, and Daniel Ritchie. PlanIT: Planning and instantiating indoor scenes with relation graph and spatial prior networks. *ACM Transactions on Graphics (TOG)*, 38(4):1–15, 2019. 2

[60] Kai Wang, Manolis Savva, Angel X Chang, and Daniel Ritchie. Deep convolutional priors for indoor scene synthesis. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. 2

[61] Magnus Wrenninge and Jonas Unger. Synscapes: A photorealistic synthetic dataset for street scene parsing. In *arXiv:1810.08705*, 2018. 1

[62] Yifan Wu, Ezra Winston, Divyansh Kaushik, and Zachary Lipton. Domain adaptation with asymmetrically-relaxed distribution alignment, 2019. 3, 4, 5

[63] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, et al. SAPIEN: A simulated part-based interactive environment. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[64] Minghao Xu, Hang Wang, Bingbing Ni, Qi Tian, and Wenjun Zhang. Cross-domain detection via graph-induced prototype alignment. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 7, 8, 18, 19

[65] Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, and Devi Parikh. Graph R-CNN for scene graph generation. In *ECCV*, 2018. 2, 3, 4, 5, 12, 18

[66] Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. Neural motifs: Scene graph parsing with global context. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[67] Yi Zhang, Weichao Qiu, Qi Chen, Xiaolin Hu, and Alan Yuille. UnrealStereo: Controlling hazardous factors to analyze stereo vision. In *arXiv:1612.04647*, 2016. 2

[68] Han Zhao, Remi Tachet des Combes, Kun Zhang, and Geoffrey J. Gordon. On learning invariant representation for domain adaptation. In *ICML*, 2019. 2

[69] Zhedong Zheng and Yi Yang. Rectifying pseudo label learning via uncertainty estimation for domain adaptive semantic segmentation. *International Journal of Computer Vision*, 2021. 2, 7

[70] Yang Zhou, Zachary While, and Evangelos Kalogerakis. SceneGraphNet: Neural message passing for 3D indoor scene augmentation. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 2

[71] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networkss. In *ICCV*, 2017. 2

[72] Yang Zou, Zhiding Yu, B. V. K. Vijaya Kumar, and Jinsong Wang. Domain adaptation for semantic segmentation via class-balanced self-training. In *ECCV*, 2018. 2, 7, 8, 18, 19

# A. Appendix

## A.1. Architecture details

**Encoder $\phi$ and SG Predictor $h$** We use ResNet 101 [22] pretrained on ImageNet as the encoder neural network. We use the Faster R-CNN [44] and Graph Convolution Network based architecture from Graph R-CNN [65] to implement the SG Predictor $h$.

**GRL** For appearance alignment $\sigma^a$, we use a 2-layer 2D convolution neural network based discriminator with ReLU activation. For prediction alignment $\sigma^{c,pred}$ we use 2-layer MLP based discriminator. We also scale the gradients to the encoder network $\phi$ from the discriminator by a factor of 4 in above cases.

## A.2. Experiments

### A.2.1 CLEVR

**Setup** The source and target domains of the CLEVR [25] environment leverage Blender [9] to render 320 x 240 images and corresponding ground truth scene graphs. The source and target domains were generated using disjoint object colors, object materials, margin between objects, and number of objects, to ensure a significant domain gap. We use 4 objects of colors (blue, green, magenta, yellow) and material (metal) for source domain. We use 2–3 objects of colors (pink, brown, white) and material (rubber) for target domain. Additionally, we transform the target by using a style transfer network [2]. For both domains, we sample each class and their size (small, medium & large) with equal probability. The environment has three lights and a fixed camera. We add a small random jitter to their initial positions during the rendering process. Some samples of source and target domain are shown in Figure 7.

**Training Details** We train our model for $10^5$ iterations with label alignment $\sigma^{c,label}$ and appearance alignments $\sigma^a$. We optimize the model using a SGD optimizer with learning rate of $10^{-4}$ and momentum of 0.9. We train our model using a batch size 4 on NVIDIA DGX workstations. We report saturation peak performance in all our tables. We give equal regularization weights to source task loss $\sigma_s$, appearance alignment $\sigma^a$ and label alignment $\sigma^{c,label}$.

**Results** More qualitative results of Sim2SG evaluated on the target domain for CLEVR are shown in Figure 8. We see better recall and fewer false positive

object detections leading to more accurate scene graphs. Label alignment $\sigma^{c,label}$ improves object recall, but occasionally introduces some false positive detections. Our appearance alignment $\sigma^a$ helps in reducing such false positives as shown in Figure 9. The full quantitative results of ablation are present in Table 4.

### A.2.2 Dining-Sim

**Setup** The Dining-Sim environment is written using Pixar's USD API[3] and rendered with a proprietary renderer. The source domain is rendered with 2 spp (samples per pixels) followed by denoiser. We select 1 chair (cantilever chair), 1 table (workshop table) and 1 laptop (PC). We randomly place chair and table on the floor and laptop on the floor as well as on the table with a random orientation. The asset for each subcategory is randomly chosen from a list of subcategory specific ShapeNet [5] assets. We also ensure that objects do not overlap by applying collision avoidance with simple box collision volumes. A subset of 4 to 5 simple materials that vary only in diffuse colour is created for each of the walls, floor, chair and table. Laptops use the original asset texture.

The target domain is rendered using path tracing with 20 spp (samples per pixels) followed by denoiser. We use 4 chairs (Windsor chair), 1 table (kitchen table) and 2 laptops (MacBook). We first place the table with a random orientation and position on the floor. We then place the four chairs at each side of the table, oriented towards the table centre. Two laptops are then placed randomly on the table surface with a random rotation. The asset for each subcategory is randomly chosen from a list of subcategory specific ShapeNet [5] assets. For materials, we use a subset of 4 to 6 physically based, highly detailed materials for each of the walls, floor, chair and table. As with the source domain, laptops use the original asset texture.

Both domains share room parameters: a fixed camera (60 degree field of view, positioned at far side of the room) and 3 fixed spherical lights. Samples from the source and the target domains are shown in Figure 10. There are five kinds of relationships - front, behind, left, right, and on with table as subject. We use 5000 labeled images from source, 5000 unlabeled images from target for training and 1000 labeled images from both source & target domains for evaluation. We use 1024 x 768 image resolution for training and evaluation.

**Details of Synthesis Step** We describe how we generate synthetic data by inferring scene graphs from target domain in detail. We filter the objects and relation-

---

[2] https://github.com/pytorch/examples/tree/master/fast_neural_style
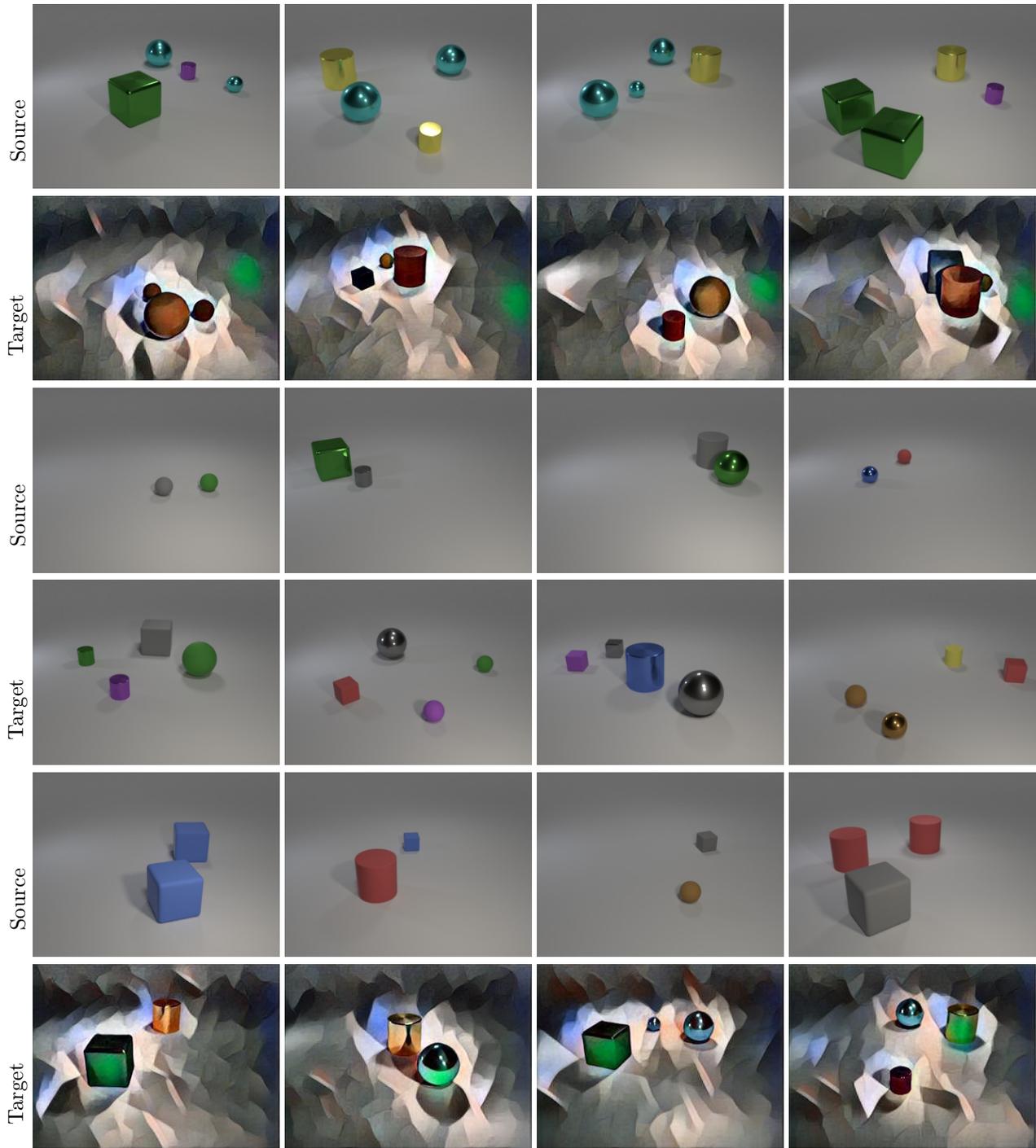
[3] https://graphics.pixar.com/usd/docs/index.html

Figure 7. Samples from source and target distributions from Clevr environment. Row 1-2: Source and Target differ in both appearance and content. Row 3-4: Source and Target differ in content but have same appearance. Row 5-6: Source and Target differ in appearance but have same content.

ships among them using an adaptive threshold (details in the next paragraph) for the generation. We assume access to camera parameters (intrinsic and extrinsic both). We project a ray from the camera through the pixel corresponding to the bounding box bottom-centre. The 3D coordinate of the object is then the intersection of the ray and the ground plane, which we assume to be flat at elevation 0. We place each object in the 3D

| Method | mAP@0.5 IoU | Recall@20 | Method | mAP@0.5 IoU | Recall@20 |
|---|---|---|---|---|---|
| SDR [41] | 0.675 | 0.339 | SDR [41] | **1.000** | 0.760 |
| Ours ($\sigma^{c,label}$) | 0.923 | 0.646 | Ours ($\sigma^a$) | 0.970 | 0.722 |
| Ours ($\sigma^a$) | **0.938** | **0.938** | Ours ($\sigma^{c,label}$) | **1.000** | **0.996** |

Table 4. Left (resp. right): Source and target domains have different (resp. similar) appearance but similar (resp. different) content distribution. All the evaluations are on the target domain.
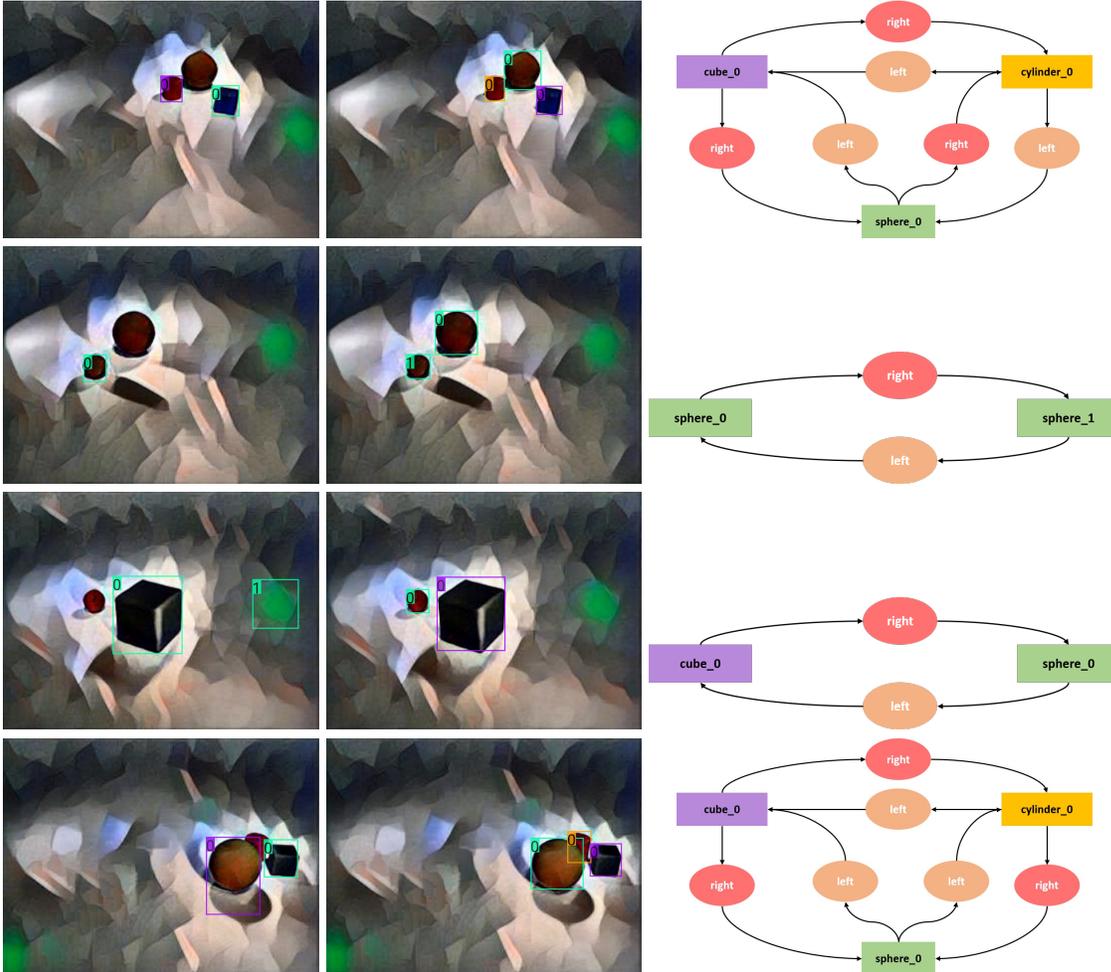


Figure 8. Qualitative results of Sim2SG on the target domain for CLEVR. First column shows that the SDR [41] fails to either detect objects or have high number of false positives (mislabels) leading to poor scene graph. Our method detects objects better, has fewer false positives and ultimately generates more accurate scene graphs as shown in second and third column respectively. Objects are color coded.

scene by picking a random 3D asset according its type (class) and assigning random pose in the range 0°–360°. We assume context like ground, wall as described in the previous paragraph. We refine the 3D scene further according to the predicted relationships among objects. For example, we use "on" relationship to refine object placements by adjusting the object (laptop or chair) elevation to match the table top. We then render the 3D scene.

**Training Details** We optimize the model using a SGD optimizer with learning rate of $10^{-4}$ and momentum of 0.9. We train our model using a batch size 2 on NVIDIA DGX workstations. We report saturation peak performance in all our tables. We give equal weights to source task loss $\sigma_s$, appearance alignment $\sigma^a$, prediction alignment $\sigma^{c,pred}$ and label alignment $\sigma^{c,label}$.

We first train the model with label alignment $\sigma^{c,label}$) for 6 epochs each with $10^4$ iterations and score threshold of 0.5. Then we add appearance alignment $\sigma^a$ and
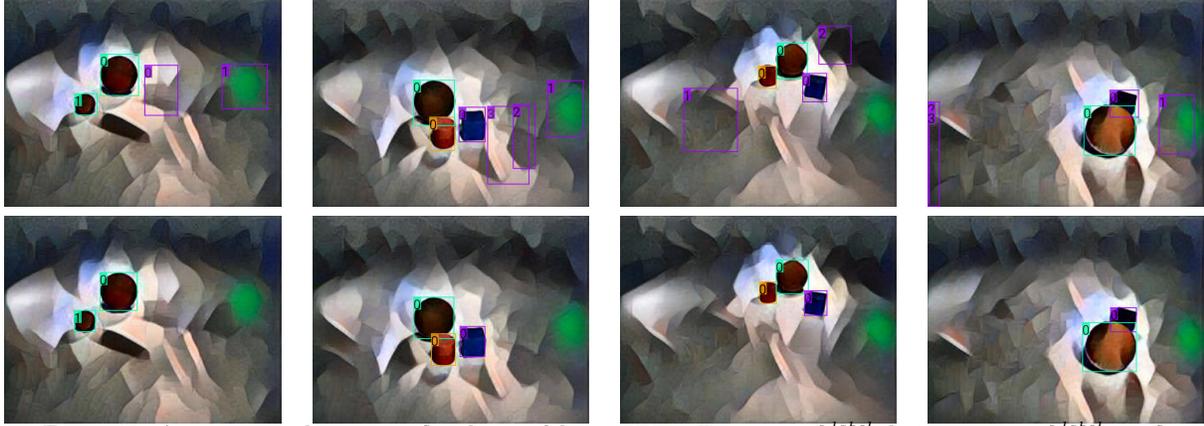
Figure 9. Appearance alignment $\sigma^a$ reducing false positive. Top row: $\sigma^{c,label}$, bottom row: $\sigma^{c,label} + \sigma^a$

prediction alignment $\sigma^{c,pred}$ and train for an additional $2 \times 10^4$ iterations. It takes 12 hours for full training including rendering time.

**Results** We present the full quantitative results in Table 5 and qualitative results in the Figure 11. We observe that the combination of all alignment terms $\sigma^{c,label}$, $\sigma^a$ and $\sigma^{c,pred}$ gives the best relationship triplet recall of 0.547@50. In order to keep our approach as general as possible, we do not enforce strict rules on object placements and prefer to randomize parameters that are not predicted such as orientation as illustrated in the qualitative results of label alignment $\sigma^{c,label}$ in Figure 12. When target domain assets are too dissimilar from the assets in the source domain, it often results in incorrect reconstructions as shown in Figure 12 (last column). We also observe that after label alignment $\sigma^{c,label}$, the model occasionally has false positive detections, particularly in areas of the floor that have intricate patterns. We qualitatively show that these false positives disappear with the addition of appearance alignment $\sigma^a$ term (Figure 13).

**Ablations** We conduct two sets of experiments on the Dining-Sim environment. The first experiment studies appearance gap: source domain has different appearance but similar content from the target domain. The source domain is generated using the target generation scheme but using source dataset materials as shown in Figure 10. The appearance gap because of photo-realistic texture in target domain is from 0.846 Recall@50 to 0.625 Recall@50. We observe that the appearance alignment $\sigma^a$ helps reduce the appearance gap, increasing relationship triplet recall from 0.625@50 to 0.821@50. For reference, the oracle performance on the target domain is 0.846 Recall@50. Similarly, the second experiment studies content gap where the source and target use the same materials but have different assets,

object positions and number of objects. We accomplish this by modifying the source generation scheme to select materials from the target dataset. Samples of source and target are shown in the third and fourth rows of Figure 10. We observe that the label alignment $\sigma^{c,label}$ term helps reduce the content gap and increase relationship triplet recall from 0.468@50 to 0.539@50. The relatively modest improvement makes sense as the two domains still differ in content (source and target domain assets differ). Full results are in Table 6.

### A.2.3 Drive-Sim

**Setup** As mentioned in Section 5.3, we use an Unreal Engine 4[4] based driving simulator akin to [41] to generate synthetic data. We have cars (1-2 per lane), trees(1-3), houses/buildings(1-3), pedestrians(0-2), sidewalk(2), roads(2-6). We do not have poles, street signs or any other objects. We have straight roads. We use realistic random placements, e.g. cars can only be placed on a lane, pedestrians on sidewalk, houses on ground and trees on both sidewalk and ground. We randomize the time of the day, cloud density and use directional light. We assume real world scale. We place our camera at a car height on a random right lane with fixed camera parameters (0 yaw, 0 pitch, 90 fov). We add realistic texture and color to each object similar to [41]. We use 1242 x 375 image resolution for training and evaluation.

**Details on Synthesis Step** We describe how we generate synthetic data by inferring scene graphs from KITTI [20] in detail. During synthesis stage, we infer the scene graphs from KITTI and further filter the objects and relationships among them using a confidence threshold of 0.2. We do not have access to KITTI camera parameters and we use the camera parameters
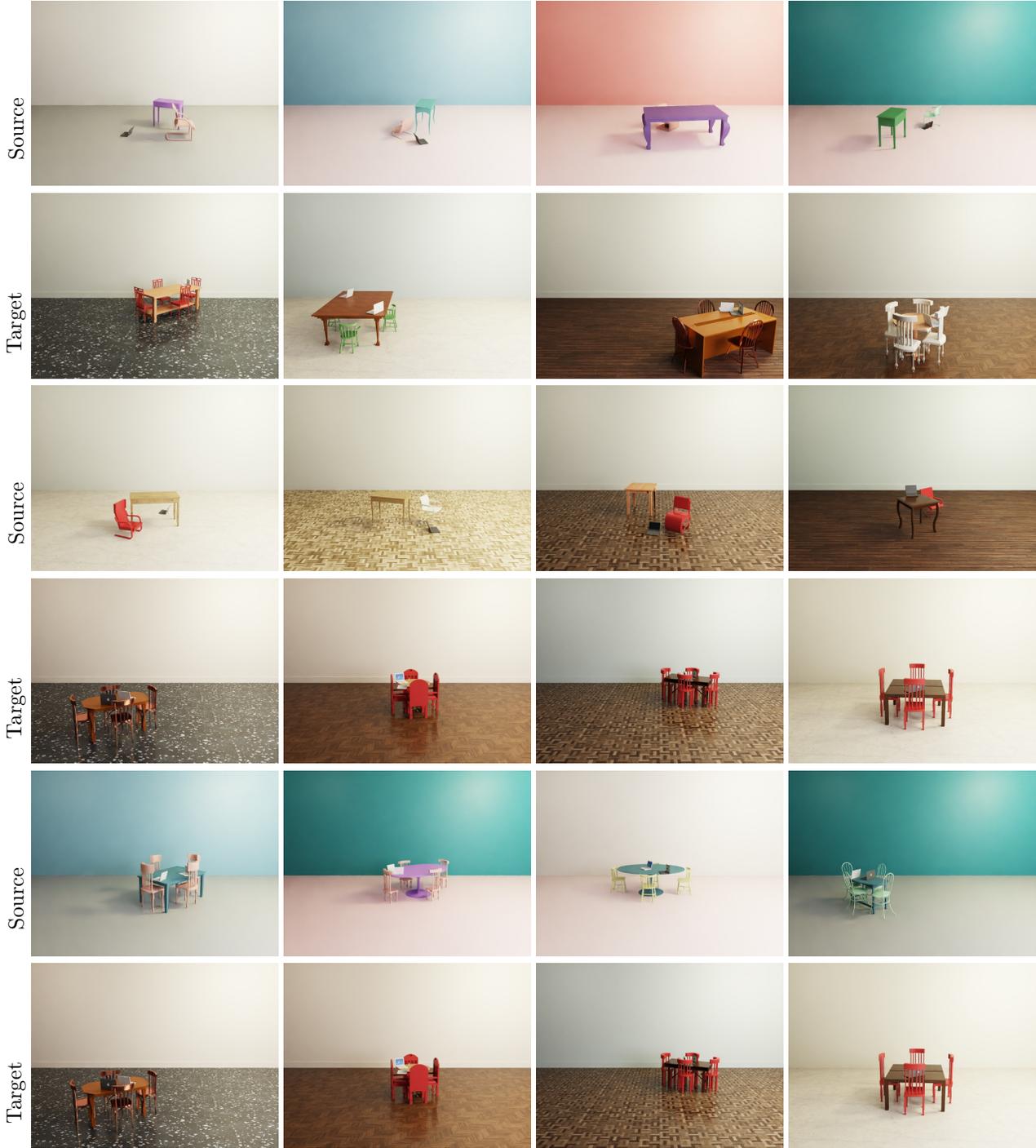
---

[4] https://www.unrealengine.com/

Figure 10. Samples from source and target distributions for Dining-Sim. Row 1-2: Source and Target domains differ in both appearance and content. Row 3-4: Source and Target differ in content but have same appearance. Row 5-6: Source and Target differ in appearance but have same content.

| Method | Chair AP | Table AP | Laptop AP | mAP @0.5 IoU | Recall@50 |
|---|---|---|---|---|---|
| SDR [41] | **0.842** ±0.038 | 0.519 ±0.088 | 0.392 ±0.051 | 0.584 ±0.049 | 0.331 ±0.064 |
| Ours ($\sigma^{c,label}$) | 0.737 ±0.043 | 0.724 ±0.030 | 0.608 ±0.047 | 0.713 ±0.038 | 0.501 ±0.044 |
| Ours ($\sigma^{c,label}$, $\sigma^a$, $\sigma^{c,pred}$) | 0.770 ±0.022 | **0.757** ±0.037 | **0.659** ±0.005 | **0.729** ±0.015 | **0.547** ±0.015 |

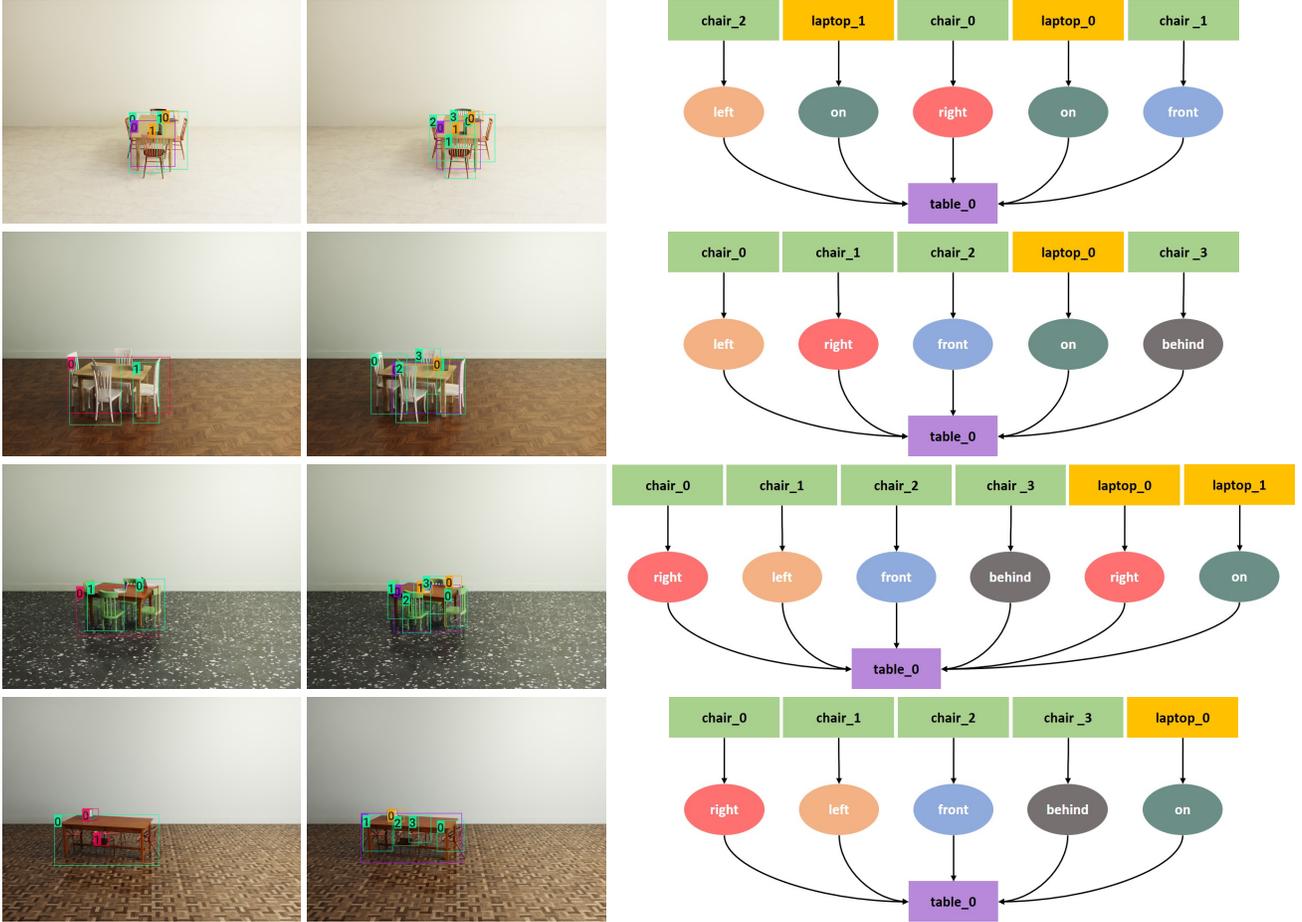Table 5. Quantitative results of Sim2SG on a target domain in Dining-Sim environment.

Figure 11. Qualitative results of Sim2SG on the target domain for Dining-Sim. First column shows that the SDR [41] fails to either detect objects or have high number of false positives (mislabels) leading to poor scene graph. Our method detects objects better, has fewer false positives and ultimately generates more accurate scene graphs as shown in second and third column respectively. Objects are color coded.



Figure 12. Scenes generated by our method (top) for target samples (bottom) in Dining-Sim environment.

described in the previous paragraph. Using the assumed camera parameters (both intrinsic and extrinsic) we project a ray from the camera through the pixel corresponding to the bounding box bottom-centre. The 3D coordinate of the object is then the intersection of the ray and the ground plane, which we assume to be
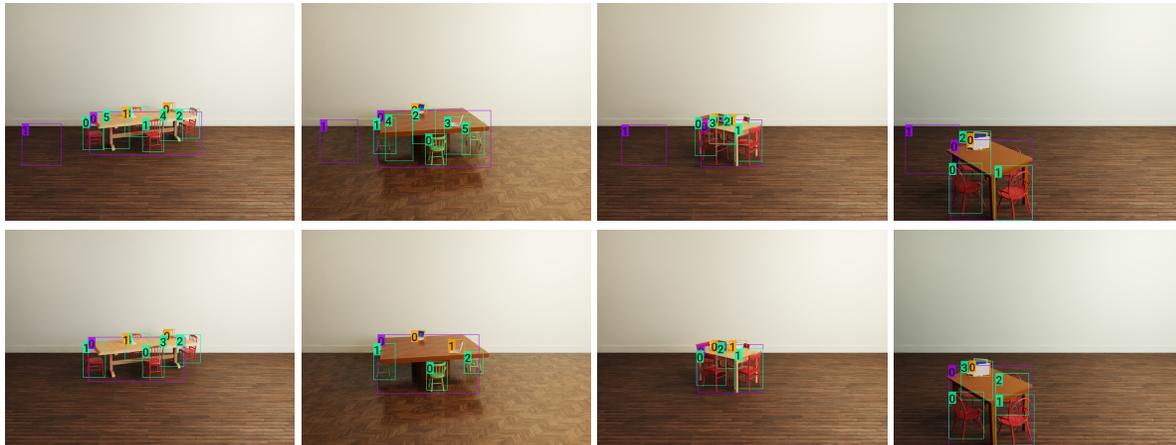
Figure 13. Appearance alignment $\sigma^a$ reducing false positive. Top row: $\sigma^{c,label}$, bottom row: $\sigma^{c,label} + \sigma^a$

| Method | mAP @0.5 IoU | Recall@50 | Method | mAP @0.5 IoU | Recall@50 |
|---|---|---|---|---|---|
| SDR [41] | 0.772 ±0.043 | 0.625 ±0.076 | SDR [41] | 0.676 ±0.011 | 0.468 ±0.006 |
| Ours ($\sigma^a$) | **0.878** ±0.001 | **0.821** ±0.006 | Ours ($\sigma^{c,label}$) | **0.737** ±0.024 | **0.539** ±0.006 |

Table 6. Dining-Sim ablations. Left (resp. right): Source and target domains have different (resp. similar) appearance but similar (resp. different) content distribution. All the evaluations are on the target domain.

flat at elevation 0. We place each object in the 3D scene by picking a random 3D asset according its type (class) and assigning random pose in the range 0°–360° (except cars that are aligned to the lane). We assume contexts like road, ground, sky, sidewalk as described in the previous paragraph. We refine the 3D scene further according to the predicted relationships among objects. We also assume a consistent lane width, and number of roads are determined by positions of the detected vehicles in the scene. We place multiple Trees (*i.e. Vegetation*) if the projected 3D volume permits. We then render the 3D scene.

**Training Details** We optimize the model using a SGD optimizer with learning rate of $10^{-4}$ and momentum of 0.9. We train our model using a batch size 2 on NVIDIA DGX workstations. We report saturation peak performance in all our tables. We give equal regularization weights to source task loss $\sigma_s$, appearance alignment $\sigma^a$, prediction alignment $\sigma^{c,pred}$ and label alignment $\sigma^{c,label}$.

We first train our model using label alignment $\sigma^{c,label}$) for 3 epochs each with $10^4$ iterations. We then add appearance alignment $\sigma^a$ and prediction alignment $\sigma^{c,pred}$ and train for an additional 60k iterations. This makes sense as $\sigma^a$ works better when content/labels are aligned between the two domains. The total training takes 12 hours including the rendering time.

**Baselines**: We adapt domain adaptation baselines [7, 64] to our framework by using the same backbone (Resnet 101) and SG Predictor (GraphRCNN [65])

network as Sim2SG, but their loss function. We do not adapt SAPNet [30]. We train these baselines on 6000 images from the source domain [41] using the same optimizer and learning rate as Sim2SG for $6 \times 10^4$ iterations. We found GPA [64] and SAPNet [30] detection performance to be lower than that reported in their work especially for pedestrian, vegetation and house classes. It is worth noting that their reported class-wise performance numbers only overlap with some of the classes in our work.

We train [27] for 40 epochs with a batch size of 16 and learning rate 0.001 as per the authors. We then obtain 6000 images and train it on Sim2SG framework (Resnet 101 backbone and GraphRCNN SG predictor) for $6 \times 10^4$ iterations using the same optimizer and learning rate as Sim2SG. For self-learning based on pseudo labels [72], we obtain the pseudo labels on KITTI images using the most confident predictions by synthetic pretrained GraphRCNN network (as per the authors). We then train these labeled KITTI images on Sim2SG framework for 60k iterations using the same optimizer and learning rate as Sim2SG.

**KITTI Annotation** We use the existing bounding box annotations of Vehicle and Pedestrians. We annotate Trees and Houses/Buildings of all sizes, occlusion and truncation in KITTI. We use the available camera parameters to project the 2D bounding box into 3D space to help us annotate spatial relationships–front, behind, left and right.

| Method | Car | Pedes. | House | Veg. | mAP | Recall@50 |
|---|---|---|---|---|---|---|
| SDR [41] | 0.488 ±0.007 | 0.214 ±0.025 | 0.223 ±0.022 | 0.177 ±0.010 | 0.276 ±0.005 | 0.112 ±0.009 |
| Meta-Sim [27] | 0.575 ±0.008 | 0.227 ±0.024 | 0.252 ±0.008 | 0.174 ±0.033 | 0.307 ±0.003 | 0.143 ±0.007 |
| Self-learning [72] | 0.466 ±0.009 | 0.215 ±0.016 | 0.189 ±0.025 | 0.265 ±0.008 | 0.284 ±0.006 | 0.129 ±0.006 |
| DA-FasterRCNN [7] | 0.523 ±0.036 | 0.209 ±0.038 | 0.203 ±0.037 | 0.171 ±0.012 | 0.277 ±0.017 | 0.119 ±0.022 |
| GPA [64] | 0.248 ±0.053 | 0.016 ±0.026 | 0.097 ±0.030 | 0.063 ±0.031 | 0.109 ±0.022 | 0.028 ±0.009 |
| SAPNet [30] | 0.420 ±0.052 | 0.124 ±0.035 | 0.018 ±0.004 | 0.042 ±0.010 | 0.151 ±0.010 | – |
| Ours ($\sigma^{c,label}$) | 0.566 ±0.033 | **0.310** ±0.029 | 0.261 ±0.009 | 0.242 ±0.040 | 0.345 ±0.002 | 0.193 ±0.010 |
| Ours ($\sigma^{c,label}$, $\sigma^a$) | 0.606 ±0.021 | 0.309 ±0.013 | 0.272 ±0.008 | 0.260 ±0.021 | 0.362 ±0.007 | 0.220 ±0.010 |
| Ours ($\sigma^{c,label}$, $\sigma^a$, $\sigma^{c,pred}$) | **0.623** ±0.033 | 0.301 ±0.018 | **0.283** ±0.007 | **0.274** ±0.015 | **0.370** ±0.005 | **0.240** ±0.003 |
| SDR | 0.412 ±0.006 | 0.174 ±0.018 | 0.215 ±0.022 | 0.177 ±0.011 | 0.245 ±0.002 | 0.085 ±0.008 |
| Meta-Sim | 0.455 ±0.040 | 0.203 ±0.026 | 0.242 ±0.009 | 0.176 ±0.029 | 0.269 ±0.007 | 0.093 ±0.005 |
| Self-learning | 0.377 ±0.007 | 0.174 ±0.014 | 0.197 ±0.002 | 0.263 ±0.006 | 0.253 ±0.004 | 0.077 ±0.005 |
| DA-FasterRCNN | 0.472 ±0.028 | 0.181 ±0.031 | 0.203 ±0.041 | 0.168 ±0.013 | 0.256 ±0.012 | 0.091 ±0.019 |
| GPA | 0.201 ±0.043 | 0.016 ±0.026 | 0.106 ±0.031 | 0.065 ±0.036 | 0.096 ±0.026 | 0.018 ±0.007 |
| SAPNet | 0.419 ±0.068 | 0.098 ±0.028 | 0.017 ±0.005 | 0.038 ±0.008 | 0.143 ±0.017 | – |
| Ours ($\sigma^{c,label}$) | 0.471 ±0.033 | **0.273** ±0.027 | 0.244 ±0.010 | 0.233 ±0.036 | 0.305 ±0.006 | 0.128 ±0.008 |
| Ours ($\sigma^{c,label}$, $\sigma^a$) | 0.511 ±0.002 | 0.266 ±0.014 | 0.251 ±0.013 | 0.256 ±0.021 | 0.321 ±0.004 | 0.155 ±0.005 |
| Ours ($\sigma^{c,label}$, $\sigma^a$, $\sigma^{c,pred}$) | **0.529** ±0.029 | 0.249 ±0.017 | **0.262** ±0.011 | **0.270** ±0.015 | **0.328** ±0.007 | **0.170** ±0.004 |
| SDR | 0.382 ±0.029 | 0.168 ±0.017 | 0.211 ±0.023 | 0.174 ±0.010 | 0.234 ±0.006 | 0.070±0.007 |
| Meta-Sim | 0.413 ±0.009 | 0.197 ±0.027 | 0.236 ±0.009 | 0.164 ±0.023 | 0.253 ±0.003 | 0.075 ±0.005 |
| Self-learning | 0.312 ±0.006 | 0.167 ±0.015 | 0.191 ±0.003 | 0.263 ±0.006 | 0.233 ±0.004 | 0.062 ±0.003 |
| DA-FasterRCNN | 0.424 ±0.028 | 0.170 ±0.029 | 0.200 ±0.041 | 0.169 ±0.014 | 0.241 ±0.014 | 0.074 ±0.015 |
| GPA | 0.174 ±0.040 | 0.011 ±0.016 | 0.106 ±0.031 | 0.059 ±0.027 | 0.087 ±0.020 | 0.015 ±0.005 |
| SAPNet | 0.362 ±0.054 | 0.085 ±0.051 | 0.116 ±0.021 | 0.067 ±0.022 | 0.157 ±0.024 | – |
| Ours ($\sigma^{c,label}$) | 0.410 ±0.009 | **0.262** ±0.025 | 0.240 ±0.010 | 0.229 ±0.036 | 0.285 ±0.003 | 0.104 ±0.006 |
| Ours ($\sigma^{c,label}$, $\sigma^a$) | 0.493 ±0.004 | 0.252 ±0.014 | 0.247 ±0.012 | 0.253 ±0.020 | 0.311 ±0.311 | 0.127 ±0.004 |
| Ours ($\sigma^{c,label}$, $\sigma^a$, $\sigma^{c,pred}$) | **0.501** ±0.006 | 0.241 ±0.018 | **0.254** ±0.010 | **0.269** ±0.014 | **0.316** ±0.004 | **0.139** ±0.004 |

Table 7. Evaluation on three modes of KITTI : easy (top), moderate (middle), hard (bottom) when training on the labeled synthetic data and unlabeled real data. The class specific AP and mAP are reported at 0.5 IoU.

| Method | Recall@20 | Recall@50 | Recall@100 |
|---|---|---|---|
| SDR | 0.098 | 0.131 | 0.146 |
| Meta-Sim | 0.109 | 0.149 | 0.164 |
| Ours ($\sigma^{c,label}$, $\sigma^a$, $\sigma^{c,pred}$) | **0.184** | **0.235** | **0.252** |
| SDR | 0.067 | 0.088 | 0.099 |
| Meta-Sim | 0.071 | 0.094 | 0.104 |
| Ours ($\sigma^{c,label}$, $\sigma^a$, $\sigma^{c,pred}$) | **0.132** | **0.167** | **0.181** |
| SDR | 0.053 | 0.071 | 0.079 |
| Meta-Sim | 0.058 | 0.076 | 0.085 |
| Ours ($\sigma^{c,label}$, $\sigma^a$, $\sigma^{c,pred}$) | **0.107** | **0.137** | **0.150** |

Table 8. Recall on three modes of KITTI : easy (top), moderate (middle), hard (bottom) when training on the labeled synthetic data and unlabeled real data. The evaluation is performed once.

**Results**  Full quantitative evaluations results are in Table 7 on all KITTI [20] evaluation criteria– easy, moderate and hard. In all three criteria, Sim2SG is able to achieve significantly better results (higher detection mAP @0.5 IoU and relationship triplet recall @ 50) than all the baselines. We also report recall @20, 100 for few baselines and our approach in Table 8. More qualitative results of label alignment $\sigma^{c,label}$ is in Figure 14. We show qualitative improvements (better object recall and fewer false positive detections) over SDR [41] and Meta-Sim [27] in Figure 15 and the corresponding accurate and full scene graphs in Figure 16.
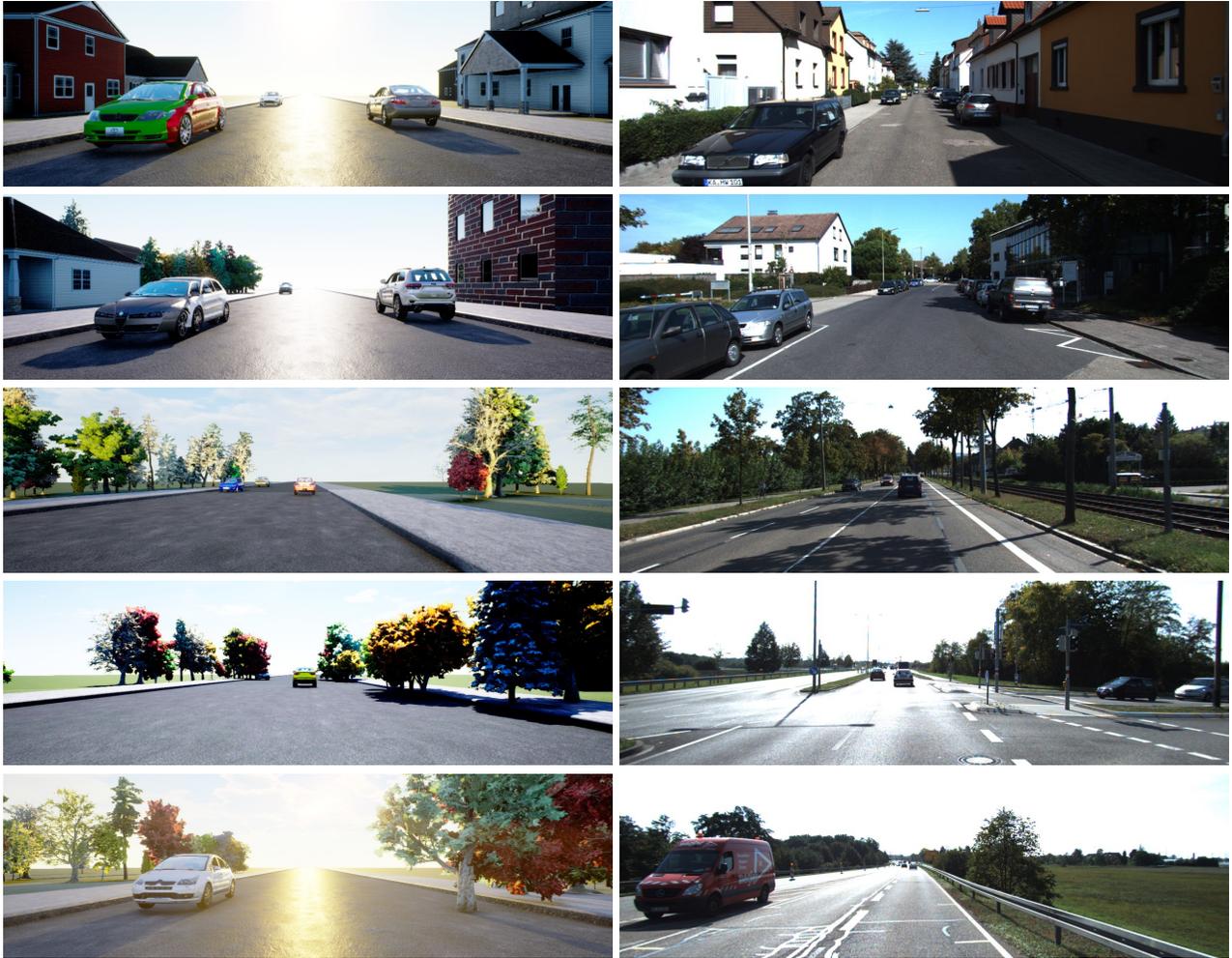
Figure 14. Scenes generated by our method (left) for target KITTI samples (right).

Figure 15. Qualitative results of objects detected on three different KITTI images. Top: SDR fails to detect many objects and yields a large number of false positives (mislabels), leading to poor scene graphs (not shown). Middle: Meta-Sim improves on false-positives, but still fails to detect some objects. Bottom: Our method detects objects correctly with fewer false positives, thus generating more accurate scene graphs. (Cars in green, vegetation in yellow, buildings in purple.)
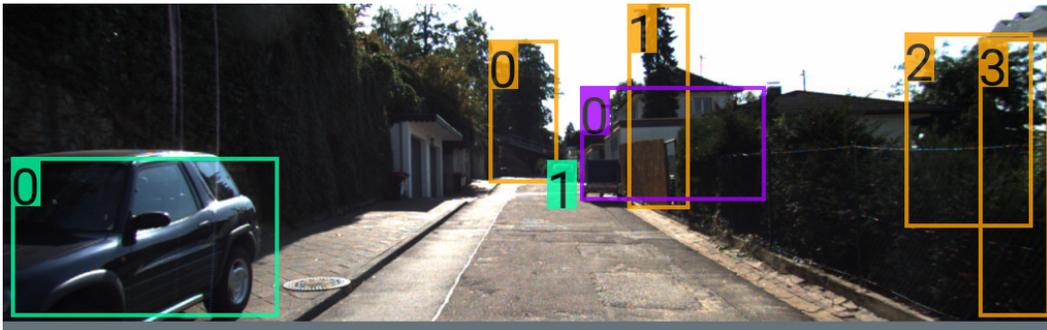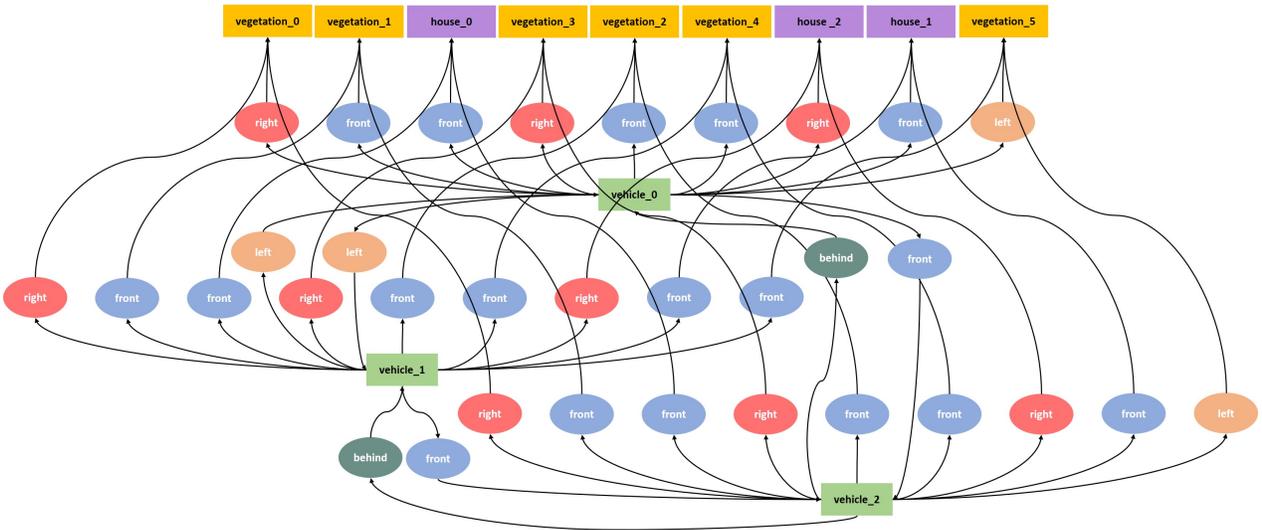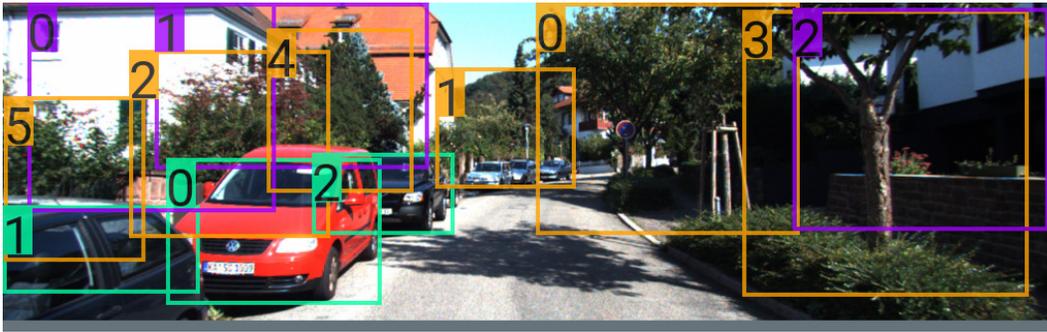
Figure 16. Qualitative results of Sim2SG on KITTI. Sim2SG generates accurate scene graphs.