SceneScript: Reconstructing Scenes With An Autoregressive Structured Language Model

Armen Avetisyan¹, Christopher Xie¹, Henry Howard-Jenkins¹, Tsun-Yi Yang¹, Samir Aroudj¹, Suvam Patra¹, Fuyang Zhang^{2†}, Duncan Frost¹, Luke Holland¹, Campbell Orme¹, Jakob Engel¹, Edward Miller¹, Richard Newcombe¹, and Vasileios Balntas¹

> ¹ Meta Reality Labs ² Simon Fraser University https://projectaria.com/scenescript



Fig. 1: (top) Given an egocentric video of an environment, SceneScript directly predicts a 3D scene representation consisting of structured scene language commands. (bottom) Our method generalizes on diverse real scenes while being solely trained on synthetic indoor environments. (last column, bottom) A notable advantage of our method is its capacity to easily adapt the structured language to represent novel scene entities. For example, by introducing a single new command, SceneScript can directly predict object parts jointly with the layout and bounding boxes.

Abstract. We introduce SceneScript, a method that directly produces full scene models as a sequence of structured language commands using an autoregressive, token-based approach. Our proposed scene representation is inspired by recent successes in transformers & LLMs, and departs from more traditional methods which commonly describe scenes as meshes, voxel grids, point clouds or radiance fields. Our method infers the set of structured language commands directly from encoded visual data using a scene language encoder-decoder architecture. To train SceneScript, we generate and release a large-scale synthetic dataset called Aria Synthetic Environments consisting of 100k high-quality indoor scenes, with photorealistic and ground-truth annotated renders of egocentric scene walkthroughs. Our method gives state-of-the art results in architectural layout estimation, and competitive results in 3D object detection. Lastly, we explore an advantage for SceneScript, which is the ability to readily adapt to new commands via simple additions to the structured language, which we illustrate for tasks such as coarse 3D object part reconstruction.

[†]Work done while the author was an intern at Meta.

1 Introduction

Scene representations play a crucial role in machine learning and computer vision applications, enabling accurate understanding of the environment. Over the years, researchers have explored various options such as meshes, voxel grids, point clouds, and implicit representations, aiming to represent complex realworld scenes with high-fidelity. Each of these exhibits distinct advantages and limitations that impact their suitability for different tasks. Meshes offer detailed geometric information but can be expensive in both computation and memory. Voxel grids provide a volumetric representation but suffer from a trade-off between resolution and memory requirements. Point clouds are efficient in representing sparse scenes but lack semantics and explicit connectivity information. Implicit representations, such as DeepSDF [22] and NeRF [15], can be infinitely precise but lack interpretability and editability. The selection of an appropriate representation directly impacts the performance and efficacy of various tasks including object recognition, scene understanding, and 3D reconstruction. In this paper, we propose a novel scene representation based on structured language commands as a more efficient and versatile solution.

Our motivation stems from the recent advancements in the field of Large Language Models (LLMs) and "next token prediction" autoregressive methods [20], coupled with recent works on exploring generation of sequences to represent geometric structures. For example, PolyGen [18] demonstrated the ability to describe 3D meshes as a sequence of vertices and faces generated using transformers [30]. Similarly, CAD-as-Language [11] showcased the effectiveness of generating Computer-Aided Design (CAD) primitives to represent 2D CAD sketches. Our main goal is to **directly infer a metrically accurate representation of a full scene** as a text-based sequence of specialized structured language commands.

Our method, denoted SceneScript, autoregressively predicts a language of hand-designed structured language commands in pure text form. This language offers several distinct advantages: 1) As pure text, it is compact and reduces memory requirements of a large scene to only a few bytes. 2) It is crisp and complete since the commands are designed to result in sharp and well-defined geometry (similar to scalable vector graphics). 3) It is interpretable, editable and semantically rich by design via the use of *high-level parametric* commands such as make_door(*door_parameters). 4) It can seamlessly integrate novel geometric entities by simply adding new structured commands to the language. 5) The fact that the scene representation is a series of language tokens similar to [20] opens up a plethora of potential new applications in the future such as ways to edit the scene, query it, or spin up chat interactions.

We mainly focus on the problems of architectural layout estimation and object detection as proxy tasks for the efficacy of our SceneScript language as a scene representation. Architectural entities such as walls, doors, and windows are highly structured entities, making them an ideal test-bed. However, one notable drawback of language models is that they require vast amounts of data for training. Since there is no existing dataset of scene walkthroughs and their corresponding structured language commands, we publicly released Aria Synthetic Environments (ASE), a synthetically generated dataset of 100k unique interior scenes. For each scene, we simulate egocentric trajectories with an entire suite of sensor data from Project Aria [14]. We also release multiple sources of ground truth including depth and instance segmentations. Importantly, for each rendered egocentric sequence, the architectural layout ground truth is given in our proposed SceneScript language.

While architectural layout serves as a test-bed, we demonstrate that our method SceneScript can easily be extended to new tasks via simple extensions to SceneScript language while keeping both the visual input and network architecture fixed. We illustrate this on the problem of 3D object detection, which results in a method that jointly infers architectural layout and 3D oriented bounding boxes. Additionally, we demonstrate more proof-of-concept experiments that show that our method results in a significantly lower barrier to entry for new tasks including representing coarse 3D object reconstruction, curved entities, composition of entities, and entity states.

Our core contributions are:

- We propose SceneScript, a method that jointly predicts architectural layout and object bounding boxes of a scene in the form of structured language commands given a video stream.
- We demonstrate that SceneScript can easily be extended to completely new tasks with simple additions of commands to our structured language, significantly lowering the barrier to entry for new tasks.
- We release a large-scale synthetic dataset, named Aria Synthetic Environments, comprized of 100k unique high-quality 3D indoor scenes with GT, which will enable large scale ML training of scene understanding methods.
- We show that training SceneScript on Aria Synthetic Environments leads to generalization on real scenes (see videos/demos on the project page).

2 Related Works

2.1 Layout Estimation

Layout estimation is an active research area, aiming to infer architectural elements. Scan2BIM [17] proposes heuristics for wall detection to produce 2D floorplans. Ochmann et al. [19] formulate layout inference as an integer linear program using constraints on detected walls. Shortest path algorithms around birds-eye view (BEV) free space [3] and wall plus room instance segmentation [5] have also been explored.

Furukawa et al. [10] utilize a *Manhattan world*-based multi-view stereo algorithm [9] to merge axis-aligned depth maps into a full 3D mesh of building interiors. RoomNet [13] predicts layout keypoints while assuming that a fixed set of Manhattan layouts can occur in a single image of a room. LayoutNet [38] improves on this by predicting keypoints and optimising the Manhattan room layout inferred from them. Similarly, AtlantaNet [23] predicts a BEV floor or

Table 1: Complete set of structured language commands designed for detailing architectural layouts and object bounding boxes. Supported data types can include int, float, bool. It is important to note that the language's extensibility allows for easy augmentation by introducing new commands like make_prim, make_pillar, or enhancing existing commands, such as incorporating is_double_door (bool).

idintidintidinta_xfloatwall0_idintwall0_idintclassinta_yfloatwall1_idintwall1_idintposition_xfloatb_xfloatposition_xfloatposition_zfloatposition_zfloatb_yfloatposition_zfloatposition_zfloatangle_zfloatb_zfloatwidthfloatwidthfloatscale_xfloat	<pre>make_wall (int)</pre>	make_door ((int)	make_window	(int)	make_bbox	(int)
height float height float height float scale z float	id int a_x float a_y float a_z float b_x float b_y float b_z float height float	id wall0_id position_x i position_y i position_z i width i height i	int int float float float float float	id wall0_id position_x position_y position_z width height	int int float float float float	id class position_x position_y position_z angle_z scale_z scale_x scale_z	int int float float float float float float

ceiling shape and approximates the shape contour with a polygon resulting in an *Atlanta world* prior. SceneCAD [1] uses a graph neural network to predict object-object and object-layout relationships to refine its layout prediction.

Our approach stands out by requiring neither heuristics nor explicitly defined prior knowledge about architectural scene structure. In fact, our method demonstrates geometric understanding of the scene which emerges despite learning to predict the GT scene language as a sequence of text tokens.

2.2 Geometric Sequence Modelling

Recent works have explored transformers for generating objects as text-based sequences. PolyGen [18] models 3D meshes as a sequence of vertices and faces. CAD-as-Language [11] represents 2D CAD sketches as a sequence of triplets in protobul representation, followed by a sequence of constraints. Both Sketch-Gen [21] and SkexGen [33] use transformers to generate sketches. DeepSVG [4] learns a transformer-based variational autoencoder (VAE) that is capable of generating and interpolating through 2D vector graphics images. DeepCAD [32] proposes a low-level language and architecture similarly to DeepSVG, but applies it to 3D CAD models instead of 2D vector graphics. Our approach stands out by utilising *high-level* commands, offering interpretability and semantic richness. Additionally, while low-level commands can represent arbitrarily complex geometries, they lead to prohibitively longer sequences when representing a full scene.

The closest work to ours is Pix2Seq [6]. Pix2Seq proposes a similar architecture to ours but experiments only with 2D object detection, thus requiring domain-specific augmentation strategies. Another closely related work is Point2Seq [34] that trains a recurrent network for autoregressively regressing continuous 3D bounding box parameters. Interestingly, they find the autoregressive ordering of parameters outperforms current standards for object detection architectures, including anchors [12] and centers [36].



Fig. 2: Aria Synthetic Environments: (top) Random samples of generated scenes showing diversity of layouts, lights and object placements. (bottom - left to right) A top down view of a scene filled with objects, a simulated trajectory (blue path), renderings of depth, RGB, and object instances, and lastly a scene pointcloud.

3 SceneScript Structured Language Commands

We first describe our structured language commands that define a full scene representation including both layout and objects. After this, we introduce our corresponding large scale training dataset: *Aria Synthetic Environments*.

3.1 Commands and Parameters

We begin with a parameterization that captures the most common layout elements. For this purpose we use three commands: make_wall, make_door, and make_window. Each command comes with a set of parameters that results in well-defined geometry. For example, the full set of parameters for make_wall specifies a gravity-aligned 2D plane, while make_door and make_window specify box-based cutouts from walls. It is worth noting that this parametrization is arbitrary, and is only made in the context of presenting a proof-of-concept SceneScript system. There are infinitely many parametrization schemes, in this work we opt for one that prioritizes ease of use and research iteration speed.

In addition to representing these three major layout entities, we aim to jointly infer objects as oriented bounding boxes. Thus, we introduce a fourth command:

make_bbox: id, class, position_x, position_y,
 position_z, angle_z, scale_x, scale_y, scale_z

This simple parametrization represents an oriented 3D bounding box that is assumed to be aligned with gravity (assuming it points in the -z direction). A summary of these commands and their respective parameters are shown in Table 1.

While we have described just four commands to capture structure and objects in an indoor environment, importantly, this text-based parametrization can readily be extended geometrically and/or even semantically to include states or other functional aspects. For example, changing the mentioned make_door command to include parameters such as open_degree allows the language to represent door states. In Section 6, we demonstrate how such extensions to the language allows for SceneScript to readily adapt to new tasks including coarse 3D object reconstruction.

3.2 Scene Definition

A single scene can be described as a sequence of our proposed structured language commands. The sequence requires no specific ordering, and can be of arbitrary length. From this definition, the 3D scene can easily be obtained by parsing the commands through a simple custom interpreter.

3.3 Training Dataset

To enable practical indoor scene reconstruction based on our structured language commands (Section 3.1), we publicly released a new dataset called *Aria Synthetic Environments*. It consists of a large number of training pairs of egocentric scene walkthroughs linked with corresponding ground truth command sequences.

Since transformers require vast amounts of data, we generated 100k synthetic scenes, which is in comparison infeasible for real-world data. Each synthetic scene comes with a floor plan model, a corresponding complete 3D scene as well as a simulated agent trajectory and a photo-realistic rendering of this trajectory. Fig. 2 illustrates the basics of *Aria Synthetic Environments*. For brevity, we refer the reader to the supplemental material for further details.

4 SceneScript Network Architecture

Our pipeline is a simple encoder-decoder architecture that consumes a video sequence and returns **SceneScript** language in a tokenized format. Figure 3 illustrates a high-level overview of our method.

We examine three encoder variants: a pointcloud encoder, a posed image set encoder, and a combined encoder. The decoder remains the same in all cases.

4.1 Input Modalities and Encoders

The encoder computes a latent scene code in the form of a 1D sequence from video walkthrough of the scene. The decoder is designed to consume these 1D sequences as input. This enables the integration of various input modalities within a unified framework. As a preliminary, for each scene we assume access to a set of M posed camera images $\{\mathbf{I}_1, ..., \mathbf{I}_M\}$, e.g., SLAM output.



Fig. 3: SceneScript core pipeline overview. Raw images & pointcloud data are encoded into a latent code, which is then autoregressively decoded into a sequence of commands that describe the scene. Visualizations are shown using a customly built interpreter. Note that for the results in this paper, the the point clouds are computed from the images using Aria MPS [25] – i.e. are not using a dedicated RGB-D / Lidar sensor.

Point Clouds. A point cloud $\mathbf{P} = {\mathbf{p_1}, ..., \mathbf{p_N}}$ consists of N points, where $\mathbf{p_i}$ is a 3D point coordinate. It can come from passive images using *SLAM* or *SfM*, or RGB-D / Lidar sensors.

Specifically, we use the Semi-dense Pointclouds from Project Aria's Machine Perception Services [25], that are obtained from a visual-inertial SLAM system using Aria's monochrome cameras and IMUs. We discretize the point cloud to 5cm resolution, then employ a sparse 3D convolution library [26,27] to generate pooled features. The encoder \mathcal{E}_{geo} applies a series of down convolutions, resulting in a reduction of the number of points in the lowest level.

$$\mathbf{F}_{aeo} = \mathcal{E}_{aeo}(\mathbf{P}), \quad \mathbf{P} \in \mathbb{R}^{N \times 3}, \mathbf{F}_{aeo} \in \mathbb{R}^{K \times 512}$$
(1)

where $K \ll N$. \mathbf{F}_{geo} is a condensed latent representation of the point cloud that contains the necessary scene context. For later use in the transformer decoder, we treat \mathbf{F}_{geo} as a sequence of feature vectors where the entries \mathbf{f}_i , $i \in 1...K$ are sorted lexicographically according to the coordinate of the active site \mathbf{c}_i , $i \in$ 1...K. To incorporate positional encoding, we append the coordinates of the active sites to their respective feature vectors $\mathbf{f}_i \leftarrow \mathsf{cat}(\mathbf{f}_i, \mathbf{c}_i)$.

Point Clouds with Lifted Features. We additionally explore augmenting the point cloud with image features. From the original egocentric sequence and associated trajectory, we sample a set of M keyframes, \mathbf{I}_i where $i \in 1...M$, and compute a set of image features \mathbf{F}_i for each. We then project each point into the set of keyframe cameras and retrieve the feature vector (output by a CNN) at the pixel location:

$$\mathbf{f}_{ip} = F_i(\pi(\mathbf{p})) \qquad \mathbf{p} \in \mathbf{P}, i \in 1...M, \tag{2}$$

where $\pi(\cdot)$ represents the projection function of a 3D point into the camera. If $\pi(\mathbf{p})$ falls outside the image bounds, no feature is retrieved. We combine the set

of lifted features for a point through an average, resulting in a single feature vector for each point: $\mathbf{f}_p = 1/M \sum_{i=1}^M \mathbf{f}_{ip}$.

We form our lifted-feature point cloud, $\mathbf{P}' = {\mathbf{p}'_1, ..., \mathbf{p}'_M}$, by concatenating each point's lifted feature with the original XYZ location: $\mathbf{p}' = \mathsf{cat}(\mathbf{f}_p, \mathbf{p})$. \mathbf{P}' is then encoded into a context sequence using sparse 3D convolutions, with only the input feature channels adjusted to match the new point feature dimension.

End-to-end Encoding of Posed Views. In order to encode the egocentric sequence more directly without a pre-computed point cloud, we adopt a $2D \leftrightarrow 3D$ bidirectional transformer encoder following the form defined in RayTran [29].

In this formulation, we initialize a volume representation of the scene as a dense voxel grid of features, \mathbf{V} , that coincides with the scene geometry. In turn, we sample a subset of M keyframes, \mathbf{I}_i where $i \in 1...M$, from the full stream of posed images. And for each of these keyframes we compute image features from a CNN, \mathbf{F}_i . Repeated layers of bidirectional attention enable the image and voxel grid features to be refined iteratively in successive transformer blocks through the aggregation of view-point and global-scene information. As in RayTran [29], the interaction between the two representations is guided by the image-formation process by leveraging known camera parameters and poses. Attention in these transformer blocks is restricted by patch-voxel ray intersections, where each image patch attends to the voxels it observes and each voxel location attends to all the patches that observe it. The resulting voxel grid of features is flattened, concatenated with an encoded represention of its XYZ location, and passed to the decoder.

4.2 Language Decoder

We utilize a transformer decoder [30] to decode the scene latent code into a sequence of structured language commands. The sequence of tokens passes through an embedding layer, followed by a positional encoding layer. Together with the encoded scene code (Section 4.1), the embedded tokens are passed into the several transformer decoder layers where a causal attention mask is used to ensure autoregressive generation. More implementation details can be found in Appendix ??.

4.3 Language Tokenization

We refer to the serialization the structured language into a sequence of tokens as **tokenization**. The goal is to construct a bijective mapping between a sequence of structured language commands (Section 3) and a sequence of integer tokens that can be predicted by the transformer decoder architecture. We utilize the following schema:

[START, PART, CMD, PARAM_1, PARAM_2, ..., PARAM_N, PART, ..., STOP]

For example, a sample sequence for a make_door command may look like:

[START, PART, MAKE_DOOR, POSITION_X, POSITION_Y, POSITION_Z, WALLO_IDX, WALL1_IDX, WIDTH, HEIGHT, PART, ..., STOP]

This schema enables 1D packing of tokens without requiring fixed-size slots or padding like other sequence modelling methods such as [32]. Additionally, it does not impose any limitations on the number or the hierarchy of sub-sequences, as they are flexibly separated by a PART token. This allows for arbitrarily complex scene representations.

The tokenized sequence is discretized into integers at a 5cm resolution, then translated into a sequence of embeddings via learnable lookup table. Note that by designing the SceneScript language, we also design the tokenization. This tokenization scheme is notably different from standard NLP tokenization, which involves Byte-Pair Encodings (BPE) [20].

5 Results

In this section, we introduce the metrics we define to measure performance, and we discuss some qualitative and quantitative results that give insights to our proposed SceneScript method.

5.1 Metrics

To evaluate accuracy of the layout estimation, we make use of geometric metrics applied between the ground-truth room layout and our predicted **SceneScript** language. To do so, we define an *entity distance*, d_E , between a pair of entities of the same class. Each entity, E, is represented as a 3D plane segment comprising of 4 corners $\{c_1, c_2, c_3, c_4\}$. The distance between two entities, E and E' is computed as the maximum Euclidean distance between each corner and its counterpart assigned via Hungarian mathcing, i.e.: $d_E(E, E') = \max\{||c_i - c'_{\pi(i)}|| : i = 1, ..., 4\}$, where $\pi(i)$ is the permutation also found via Hungarian matching.

We threshold d_E to define the success criteria for the predicted entities. We compute the following metrics:

- F1 Score @ threshold the F1 score of the set of predictions is computed at a single d_E threshold.
- Average F1 Score the F1 score is computed across a range of entity distance thresholds and averaged.

The scores are computed for each class independently and averaged to overall score. In addition, scores are computed for each scene and averaged across the dataset. We use the following range of thresholds (cm) for the average F1 scores: $T = \{1, 2, ..., 9, 10, 15, 25, 30, 50, 75, 100\}.$



Fig. 4: Qualitative samples between our model and SOTA methods on *Aria Synthetic Environments*'s test set. Hierarchical methods like SceneCAD suffer from error cascading which leads to missing elements in the edge prediction module. RoomFormer (a 2D method extruded to 3D) primarily suffers from lightly captured scene regions which leave a unnoticeable signal in the density map.

 Table 2: Layout Estimation on Aria Synthetic Environments
 Quantitative

 comparison between our method and related recent work.
 \$\$

	F1 @5cm				Avg F1			
Method	mean	wall	door	window	mean	wall	door	window
SceneCAD '20 [1]	-	0.048	-	-	-	0.275	-	-
RoomFormer '23 [37]	0.139	0.159	0.148	0.110	0.464	0.505	0.481	0.407
Ours (Point cloud)	0.848	0.930	0.922	0.692	0.784	0.816	0.811	0.724
Ours (Lifted features)	0.903	0.943	0.959	0.806	0.801	0.818	0.822	0.764
Ours (Image-only)	0.661	0.687	0.798	0.497	0.719	0.727	0.772	0.658

5.2 Layout Estimation

We perform scene layout estimation with the three encoder variants of **SceneScript**: a baseline model with sparse3d convolution pointcloud encoder, an RGB RayTran-based feature volume encoder [29], and our proposed lifted feature point encoder. The same transformer decoder is being used in all three scenarios.

To provide comparison to existing works, we include results from two baseline methods, namely SceneCAD [1] and the recent RoomFormer [37]. For these experiments SceneCAD and RoomFormer were both trained on *Aria Synthetic Environments*. Note that SceneCAD only predicts walls.

Table 2 shows the main results for our F1-based metrics on *Aria Synthetic Environments*. SceneScript exhibits a substantial performance advantage over SOTA layout estimation baselines across multiple metrics.

Table 3: 3D Object Detection Performance comparison against state-of-the-art methods on an 3D object detection task trained and evaluated by F1-score at 0.25 and 0.5 IoU thresholds. By simply adding a make_bbox command SceneScript can achieve competitive object detection results.

(a) Aria Synthetic Environments				(b) ScanNet [8]				
Method	Input	F @.25 IoU	1 @.50 IoU	Method	Input	F @.25 IoU	1 @.50 IoU	
3DETR '21 [16] Cube R-CNN '23 [2] ImVoxelNet '22 [24] Ours	Points RGB RGB Points	0.201 0.394 0.584 0.620	0.078 0.228 0.516 0.577	3DETR '21 [16] 3DETR-m '21 [16] SoftGroup '22 [31] Ours	Points Points RGB Points RGB Points	0.480 0.536 0.622 0.506	0.349 0.407 0.573 0.406	

Both baseline methods encounter a significant decline in accuracy when dealing with finer details. See Figure 4 for qualitative comparisons between our method and baseline methods.

Encoder Ablation. The results demonstrate that **SceneScript** is robust to the encoding strategy chosen to encode the egocentric capture. It is able to infer highly accurate scene layouts in all configurations tested, and in each case SceneScript outperforms the included baselines by a significant margin.

Relative comparison of the encoder strategies reveals that leveraging the pointclouds from a highly specialized and well-tuned system is still able to offer advantages of an, almost, entirely learned approach such as RayTran [29]. A light extension in the form of lifted image features can widen this gap even further. In particular, we observe that the discrepancy between the encoders becomes particularly apparent as the complexity of the scene increases in the form of increased room count. In the Appendix, we also show a quantitative evaluation of per-entity error distances, which aids in further attribution of relative performance gains between the encoding methods.

Object Detection 5.3

In this section, we perform evaluation of SceneScript for object detection on both Aria Synthetic Environments and ScanNet [8]. For comparison, we include recent and state-of-the-art baseline methods, namely Cube-RCNN [2], ImVoxel-Net [24], 3DETR [16], and SoftGroup [31].

Worth noting is that SceneScript does not predict a confidence score for each make_bbox command. This means that fair computation of the conventional mAP metric for this task is not possible, as detections cannot be ranked across scenes. Among other issues, this results in a metric that varies with the order in which scenes are evaluated. We therefore report F1-score-based metrics, which do not exhibit this order variance. Further discussion of this, and mAP numbers for the baselines for reference, can be found in the Appendix.



Fig. 5: Example scene reconstructions on scenes from *Aria Synthetic Environments*. (left) Visualisation of the decomposed meshes used to create make_prim training pairs. (right) Views of full scene predictions, as well as close ups highlighting the fidelity of object reconstruction through the prediction volumetric primitives enabled by make_prim.

In Table 3a, all methods were trained on Aria Synthetic Environments. 3DETR performs poorly due to its encoder being designed for relatively clean datasets such as ScanNet [8], while semi-dense point clouds from Aria Synthetic Environments exhibit more noise and less uniformity (see Figure 3 for an example). Cube R-CNN and ImVoxelNet both operate on RGB images, and detections are tracked for the entire sequence via a tracker [2] to provide competitive performance. In Table 3b, our method provides similar performance to both 3DETR and SoftGroup.

Through the addition of the make_bbox command, SceneScript demonstrates object detection performance on par with SOTA baselines on multiple object detection datasets. This result illustrates the extensibility of a token-based approach, and that our proposed SceneScript language representation does not suffer compared to specialised object detection network architectures.

6 Extending the SceneScript Structured Language

A key advantage offered by SceneScript's structured language prediction paradigm is that the expressiveness of its reconstruction can be tailored without requiring a change to the method. Up to now, we have focussed on showcasing the efficacy of SceneScript for representing simple layout elements and objects as bounding boxes. In this section, we showcase this characteristic by increasing the fidelity of our scene representation by introducing coarse 3D object reconstruction.

6.1 Objects as Volumetric Primitives

We turn to a language based on volumetric primitives, motivated by works such as [28,35]. Using simple primitives such as cuboids and extruded cylinders, enables us to coarsely represent arbitrary object categories while maintaining object semantics (e.g. tabletops can be represented by a single cuboid). Thus, this language can describe many object categories simultaneously.



Fig. 6: Example scene reconstructions on real scenes with the addition of the make_prim command. Note that SceneScript was trained only on synthetic data.

This representation requires only one additional command over the layout and box commands already discussed previously, namely:

make_prim: bbox_id, prim_num, class, center_x, center_y
, center_z, angle_x, angle_y, angle_z, scale_x,
scale_y, scale_z

This command and its parameters describe a volumetric primitive (cuboid or extruded cylinder) via its 3D center, 3D rotation, and 3D scale. The prim_num parameter can be associated with semantics, e.g. tabletops of different tables typically have the same prim_num.

Dataset. To obtain ground truth make_prim commands that align with the objects in Aria Synthetic Environments, we first run an extension of Yang et al. [35] to obtain cuboid and extruded cylinder primitives of a database of 3D CAD models (ABO [7], which was used to populate Aria Synthetic Environments). See Figure 5 (left) for example decompositions. For this proof-of-concept experiment, we use three categories: table, chair, and sofa. We then convert these decomposed primitives into make_prim commands that are aligned with the objects in the dataset, which results in training pairs.

Results. We show qualitative results on *Aria Synthetic Environments* in Figure 5. In Figure 6, we show inferences in a few real-world environments despite only having trained on our simulated dataset. These results demonstrate that **SceneScript**'s general purpose architecture is capable of coarsely reconstructing objects of multiple categories through addition of a new command.

6.2 Further Extensions

In this section, we have explored just one extension of SceneScript's structured scene language in order to demonstrate its flexibility. The result was greatly increased expressiveness of the scene reconstruction. In the Appendix, we include

additional explorations that can further increase the fidelity and accuracy of the reconstructions. These explorations include reconstructing curved walls, inferring object states (e.g., door opening angles), as well as direct prediction of parametric models using object models deployed commonly by tech artists (e.g., blender geometry nodes) for object reconstruction.

7 Interactive Scene Reconstruction

Scene reconstruction often occurs offline, on pre-recorded walkthroughs of an environment, or derivatives of them, such as a point cloud. However, we take advantage of SceneScript's inference occurring at an interactive rate, which takes between 2-5s depending on the size of the environment, by implementing streaming of SceneScript's live reconstructions into a modern VR headset. This enables an interactive experience where the user can see the reconstruction overlaid onto the environment they are exploring in real-time. See the video recording on the website for examples. Visualisations of this interface are included in Figure 1.

8 Limitations and Future Work

SceneScript exhibits certain limitations that should be acknowledged. First, the structured language commands are manually defined, which requires human intervention at this stage. Secondly, due to the higher-level nature of our commands, it can be challenging to capture fine-grained geometric details with extremely high precision (i.e. mm). As a consequence, the resulting reconstructions based on this representation tend to lead to simpler and coarser geometries, potentially missing intricate nuances at the very high detail level. These limitations potentially highlight areas for future research and optimization, aiming to automate the command definition process and explore techniques to enhance the representation's ability to capture intricate geometric details accurately. However, we believe that the ability to built scene representations that are based on structured language commands will be a key part in complex and efficient methods of scene reconstruction in the future.

9 Conclusion

We introduced SceneScript, a novel reconstruction method that is based on a tokenized scene representation. SceneScript autoregressively predicts a sequence of structured scene language commands given a video stream of an indoor space. This tokenized scene representation is compact, editable and intepretable. In addition, we showed that a strength of SceneScript is the ability to extend to arbitrary elements (e.g. Bezier curves, object part decomposition) with minimal changes. This research opens up new directions in representing 3D scenes as language, bringing the 3D reconstruction community closer to the recent advances in large language models such as GPT-class models [20].

References

- Avetisyan, A., Khanova, T., Choy, C., Dash, D., Dai, A., Nießner, M.: Scenecad: Predicting object alignments and layouts in rgb-d scans. In: Computer Vision– ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16. pp. 596–612. Springer (2020)
- Brazil, G., Kumar, A., Straub, J., Ravi, N., Johnson, J., Gkioxari, G.: Omni3D: A large benchmark and model for 3D object detection in the wild. In: CVPR. IEEE, Vancouver, Canada (June 2023)
- Cabral, R., Furukawa, Y.: Piecewise planar and compact floorplan reconstruction from images. In: 2014 IEEE Conference on Computer Vision and Pattern Recognition (2014)
- Carlier, A., Danelljan, M., Alahi, A., Timofte, R.: Deepsvg: A hierarchical generative network for vector graphics animation. In: Advances in Neural Information Processing Systems (2020)
- Chen, J., Liu, C., Wu, J., Furukawa, Y.: Floor-sp: Inverse cad for floorplans by sequential room-wise shortest path. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2019)
- Chen, T., Saxena, S., Li, L., Fleet, D.J., Hinton, G.: Pix2seq: A language modeling framework for object detection. In: International Conference on Learning Representations (ICLR) (2022)
- Collins, J., Goel, S., Deng, K., Luthra, A., Xu, L., Gundogdu, E., Zhang, X., Yago Vicente, T.F., Dideriksen, T., Arora, H., Guillaumin, M., Malik, J.: Abo: Dataset and benchmarks for real-world 3d object understanding. CVPR (2022)
- Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: Proc. Computer Vision and Pattern Recognition (CVPR), IEEE (2017)
- 9. Furukawa, Y., Curless, B., Seitz, S.M., Szeliski, R.: Manhattan-world stereo. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition (2009)
- Furukawa, Y., Curless, B., Seitz, S.M., Szeliski, R.: Reconstructing building interiors from images. In: 2009 IEEE 12th international conference on computer vision (2009)
- Ganin, Y., Bartunov, S., Li, Y., Keller, E., Saliceti, S.: Computer-aided design as language. Advances in Neural Information Processing Systems 34, 5885–5897 (2021)
- Girshick, R.: Fast r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 1440–1448 (2015)
- Lee, C.Y., Badrinarayanan, V., Malisiewicz, T., Rabinovich, A.: Roomnet: End-toend room layout estimation. In: Proceedings of the IEEE international conference on computer vision (2017)
- 14. Meta: Project aria. https://projectaria.com/ (2022), accessed: 2023-08-30
- Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. Communications of the ACM 65(1), 99–106 (2021)
- Misra, I., Girdhar, R., Joulin, A.: An end-to-end transformer model for 3d object detection. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 2906–2917 (2021)
- Murali, S., Speciale, P., Oswald, M.R., Pollefeys, M.: Indoor scan2bim: Building information models of house interiors. In: 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (2017)

- 16 Avetisyan et al.
- Nash, C., Ganin, Y., Eslami, S.A., Battaglia, P.: Polygen: An autoregressive generative model of 3d meshes. In: International conference on machine learning. pp. 7220–7229. PMLR (2020)
- Ochmann, S., Vock, R., Klein, R.: Automatic reconstruction of fully volumetric 3d building models from oriented point clouds. ISPRS journal of photogrammetry and remote sensing 151, 251–262 (2019)
- 20. OpenAI: Gpt-4 technical report (2023)
- Para, W., Bhat, S., Guerrero, P., Kelly, T., Mitra, N., Guibas, L.J., Wonka, P.: Sketchgen: Generating constrained cad sketches. In: Advances in Neural Information Processing Systems (2021)
- 22. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (2019)
- 23. Pintore, G., Agus, M., Gobbetti, E.: Atlantanet: inferring the 3d indoor layout from a single 360 image beyond the manhattan world assumption. In: Proceedings of the European conference on computer vision (ECCV) (2020)
- Rukhovich, D., Vorontsova, A., Konushin, A.: Imvoxelnet: Image to voxels projection for monocular and multi-view general-purpose 3d object detection. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 2397–2406 (2022)
- Somasundaram, K., Dong, J., Tang, H., Straub, J., Yan, M., Goesele, M., Engel, J.J., De Nardi, R., Newcombe, R.: Project aria: A new tool for egocentric multimodal ai research. arXiv preprint arXiv:2308.13561 (2023)
- Tang, H., Liu, Z., Li, X., Lin, Y., Han, S.: TorchSparse: Efficient Point Cloud Inference Engine. In: Conference on Machine Learning and Systems (MLSys) (2022)
- Tang, H., Liu, Z., Zhao, S., Lin, Y., Lin, J., Wang, H., Han, S.: Searching Efficient 3D Architectures with Sparse Point-Voxel Convolution. In: European Conference on Computer Vision (ECCV) (2020)
- Tulsiani, S., Su, H., Guibas, L.J., Efros, A.A., Malik, J.: Learning shape abstractions by assembling volumetric primitives. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
- Tyszkiewicz, M.J., Maninis, K.K., Popov, S., Ferrari, V.: Raytran: 3d pose estimation and shape reconstruction of multiple objects from videos with ray-traced transformers. In: Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part X. pp. 211–228. Springer (2022)
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. NeurIPS **30** (2017)
- Vu, T., Kim, K., Luu, T.M., Nguyen, X.T., Yoo, C.D.: Softgroup for 3d instance segmentation on 3d point clouds. In: CVPR (2022)
- Wu, R., Xiao, C., Zheng, C.: Deepcad: A deep generative network for computeraided design models. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (2021)
- Xu, X., Willis, K.D., Lambourne, J.G., Cheng, C.Y., Jayaraman, P.K., Furukawa, Y.: Skexgen: Autoregressive generation of cad construction sequences with disentangled codebooks. In: International Conference on Machine Learning (ICML) (2022)
- 34. Xue, Y., Mao, J., Niu, M., Xu, H., Mi, M.B., Zhang, W., Wang, X., Wang, X.: Point2seq: Detecting 3d objects as sequences. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8521–8530 (2022)

- Yang, K., Chen, X.: Unsupervised learning for cuboid shape abstraction via joint segmentation from point clouds. ACM Transactions on Graphics (TOG) 40(4), 1–11 (2021)
- Yin, T., Zhou, X., Krahenbuhl, P.: Center-based 3d object detection and tracking. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (2021)
- 37. Yue, Y., Kontogianni, T., Schindler, K., Engelmann, F.: Connecting the dots: Floorplan reconstruction using two-level queries. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2023)
- Zou, C., Colburn, A., Shan, Q., Hoiem, D.: Layoutnet: Reconstructing the 3d room layout from a single rgb image. In: Proceedings of the IEEE conference on computer vision and pattern recognition (2018)