

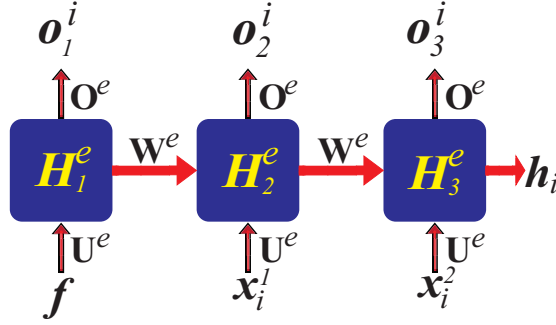
# Supplemental material for “SeqXY2SeqZ: Structure Learning for 3D Shapes by Sequentially Predicting 1D Occupancy Segments From 2D Coordinates”

Zhizhong Han<sup>1,2</sup>, Guanhui Qiao<sup>1</sup>, Yu-Shen Liu<sup>1\*</sup>, and Matthias Zwicker<sup>2</sup>

<sup>1</sup> School of Software, BNRist, Tsinghua University, Beijing, P. R. China

<sup>2</sup> Department of Computer Science, University of Maryland, College Park, USA  
h312h@umd.edu, qiaogh18@mails.tsinghua.edu.cn, liuyushen@tsinghua.edu.cn,  
zwicker@cs.umd.edu

## 1 The encoding process



**Fig. 1.** The encoding process.

In the encoder RNN shown in Fig. 1, we use the following equations to calculate the hidden states at each step,

$$H_a^e = \text{ReLU}(\mathbf{U}^e \mathbf{f}), a = 1, \quad (1)$$

$$H_a^e = \text{ReLU}(\mathbf{U}^e \mathbf{x}_i^{a-1} + \mathbf{W}^e \mathbf{H}_{a-1}^e), a > 1, \quad (2)$$

where  $\mathbf{U}^e$  and  $\mathbf{W}^e$  are learnable parameters.  $\mathbf{f}$  is the shape condition which can be either learnable parameter in autoencoding experiment or feature vector

\* Corresponding Author. This work was supported by National Key R&D Program of China (2018YFB0505400) and NSF (award 1813583).

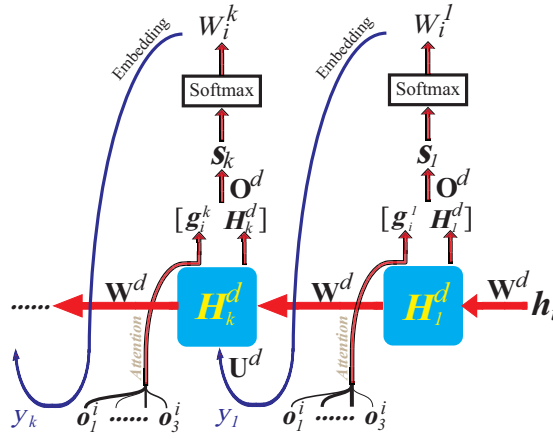
extracted from images in single image reconstruction experiment, and  $\{x_i^{a-1}\}$  are location embeddings which are also learnable parameter, and  $a \in [1, 2, 3]$ .

Then, we can calculate the output  $\mathbf{o}_a^i$  at each step as defined below,

$$\mathbf{o}_a^i = \mathbf{O}^e \mathbf{H}_a^e + \mathbf{b}^e, \quad (3)$$

where  $\mathbf{O}^e$  and  $\mathbf{b}^e$  are learnable parameters.

## 2 The decoding process



**Fig. 2.** The decoding process.

In the decoder shown in Fig. 2, we use attention weights  $d_a^k$  learned by attention mechanism [1] at each step  $k$  to aggregate the output  $\mathbf{o}_a^i$  at all steps of encoder below,

$$\mathbf{g}_i^k = \sum_{a=1}^3 d_a^k \times \mathbf{o}_a^i. \quad (4)$$

The hidden state at each step is calculated as follows,

$$\mathbf{H}_k^d = \text{ReLU}(\mathbf{W}^d \mathbf{h}_i), k = 1, \quad (5)$$

$$\mathbf{H}_k^d = \text{ReLU}(\mathbf{U}^d \mathbf{y}_{k-1} + \mathbf{W}^d \mathbf{H}_{k-1}^d), k > 1, \quad (6)$$

where  $\mathbf{W}^d$  and  $\mathbf{U}^d$  are learnable parameters,  $\{\mathbf{y}_{k-1}\}$  are location embeddings which are learnable parameters.

Then, we combine the hidden state and the aggregated output from encoder to learn a feature for start and end location prediction as follows,

$$\mathbf{s}_k = \mathbf{O}^d[\mathbf{g}_i^k \mathbf{H}_k^d] + \mathbf{b}^d, \quad (7)$$

where  $\mathbf{O}^d$  and  $\mathbf{b}^d$  are learnable parameters.

Finally, we predict the start or end location  $w_i^k$  at each step by a softmax layer defined below,

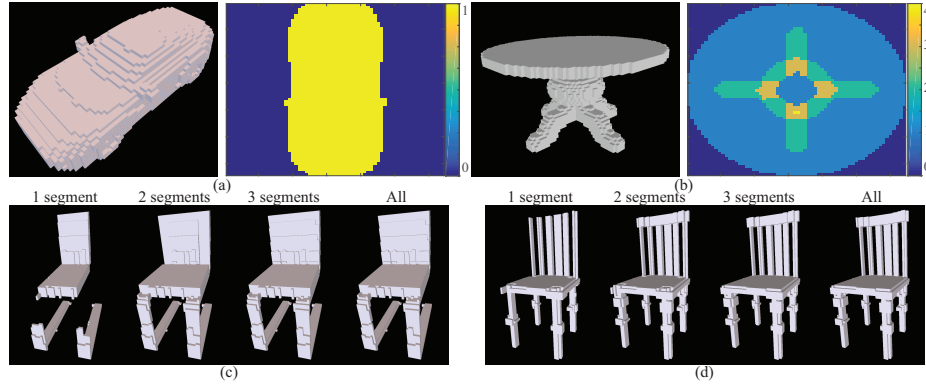
$$\mathbf{y}_k = \mathbf{W}\mathbf{s}_k + \mathbf{b}, \quad (8)$$

$$w_i^k = \arg \max \text{softmax}(\mathbf{y}_k), \quad (9)$$

where  $\mathbf{W}$  and  $\mathbf{b}$  are learnable parameters.

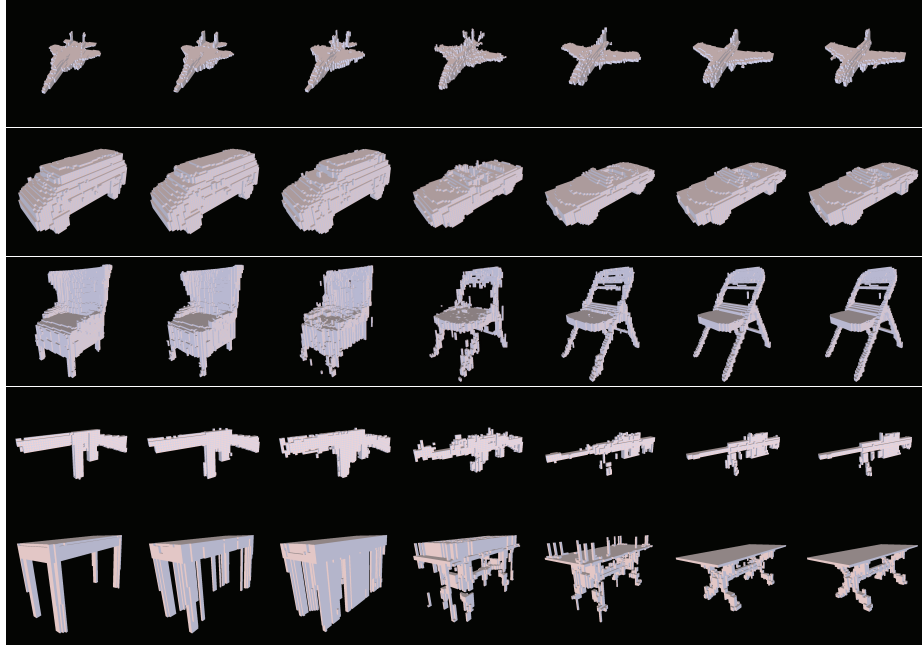
### 3 More analysis

**Occupancy segment visualization.** With the auto-encoded shapes, we justify the efficiency of our voxel tubelization by visualizing the number of predicted occupancy segments at each 2D coordinate in Fig. 3. For a simple car in Fig. 3 (a), it is enough to use only one occupancy segment to represent the geometry at each 2D coordinate. Although the table in Fig. 3 (b) is more complex, we can still leverage no more than three occupancy segments to represent the geometry at almost all 2D coordinates. In Fig. 3 (c) and (d), we visualize the reconstructed chairs with different numbers of occupancy segments. For complex structures in chairs, it is still enough to reconstruct almost a whole shape using only two occupancy segments.



**Fig. 3.** The efficiency demonstration of occupancy segments at 2D coordinates. We show the numbers of occupancy segments at all 2D coordinates as an image for a car (a) and table (b). We also visualize the shapes reconstructed by different numbers of occupancy segments at all 2D coordinates in (c) and (d). In all cases, only few occupancy segments are needed to represent a complex shape.

**Interpolation.** We visualize the shape condition space learned in auto-encoding experiment by visualizing the interpolation between two shapes. We interpolate the features between two learned shape conditions, which are further leveraged to reconstruct shapes shown in Fig. 4. The transition shows how one shape is gradually transformed to another one by growing or shrinking occupancy segments in each tube at all 2D locations.



**Fig. 4.** Interpolated shapes in the learned 3D feature space. We use the learned shape conditions of two shapes to establish a line in the feature space, and then uniformly sample five features along this line by interpolating the two shape conditions. Finally, we generate a interpolated shape using each one of the five interpolated features as shape condition.

## References

1. Bahdanau, D., Cho, K., Bengio, Y.: Neural machine translation by jointly learning to align and translate. CoRR **abs/1409.0473** (2014)