974

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994 995

996

997

1003

1004

1005

1006

1007 1008

1009

1026

1027

1028

1078

1079

972 **A. Supplementary Materials** 973

A.1. Pseudocode for Line Vectorization

975 Algorithm 1 gives a more detailed description for line 976 vectorization. The algorithm takes the C-junction set V_C and T-junction set V_T as the input and outputs a vectorized wireframe (V, E). In the first stage (Lines 2-3), we find the lines among C-junctions according to the line confidence function. The procedure PRUNE-LINES greedily removes the lines with the lowest confidence that either intersect with other lines (Line 32) or are too close to other lines in term of the polar angle (Line 35). In the second stage (Lines 4-25), we add the T-junctions into the wireframe. From Lines 6-14, we find the T-junctions that are on the existing wireframe. We first adjust the positions of those T-junctions by projecting them onto the line (Line 9) and then add them to the candidate T-junction set V' (Line 10). Because the degree of a T-junction is always one, we try to find the connection with the highest confidence for those candidates T-junctions (Lines 15-25). We repeat the this process until V, V', and E remain the same in the last iteration.

A.2. Line Assignments for Vanishing Points

In Equation (4), we need to find the set of lines $A_i \subseteq E$ corresponding to the vanishing point *i*. Mathematically, we define the objective function

$$\min_{\mathsf{A}} \sum_{i}^{3} \sum_{(\boldsymbol{u},\boldsymbol{v})\in\mathsf{A}_{i}} \|(\boldsymbol{u}-\boldsymbol{V}_{i})\times(\boldsymbol{u}-\boldsymbol{v})\|_{2},$$

where $\|(\cdot) \times (\cdot)\|_2$ can be understood as the parallelogram area formed by two vectors. Since each line in this equation is mutually independent, we can solve this optimization problem by greedily assigning each line to the best vanishing point *i* to minimize the objective function.

A.3. Sampled Failure Cases and Discussions

1010 Figure 7 demonstrates some failure cases in our SceneCity dataset. We found that our pipeline might not 1011 1012 work well on the scenes in which there are many lines and junctions that are close to each other. This is because the 1013 resolution of the output heat map is 128×128 , so any de-1014 1015 tail whose size is below two or three pixels might get lost 1016 during the vectorization stage. Therefore, one of our future 1017 work is to explore the possibility of using high-resolution input and output images. There are also issues in the 3D 1018 1019 depth refinement stage. When the scene is complex, finding 1020 the assignment A_i for each line can be hard, due to the er-1021 ror in the junction position and line direction. In addition, the term contributed by erroneous lines in Equation (4) can 1022 make the depth of some junctions inaccurate. Such problem 1023 might potentially be alleviated by increasing the resolution 1024 1025 of the input and output images, using a more data-driven

	102
Algorithm 1 Edge Vectorization Algorithm	103
Require: Candidate C-junction set V_C . T-junction set V_T .	103
Require: Hyper-parameters n_c and n_0 .	103
Ensure: Wireframe (V, E).	103
1: procedure Vectorize(V_C , V_T)	103
$2: \forall \leftarrow \forall C$	103
3: $E \leftarrow Prune-Lines(\{(u, v) u, v \in V, c(u, v) > \eta_c\})$	103
4: $V' \leftarrow \emptyset$	103
5: while V, V', or E change in the last iteration do	103
6: for all $w \in V_T$ do	103
7: for all $e = (u, v) \in E$ do	104
8: if w is near e then	104
9: project w to the line e	104
10: $V' \leftarrow V' \cup \{w\}$	104
11: break	104
12: end if	104
13: end for	104
14: end for	104
15: for all $u \in V'$ do	104
16: $\boldsymbol{v} \leftarrow \operatorname{argmax}_{\boldsymbol{v} \in V \cup V'} c(\boldsymbol{u}, \boldsymbol{v})$	104
17: if $c(\boldsymbol{u}, \boldsymbol{v}) \geq \eta_c$ then	105
18: $V' \leftarrow V' \setminus \{u\}$	105
19: $V \leftarrow V \cup \{u\}$	105
20: $E \leftarrow E \cup \{(u, v)\}$	105
21: end if	105
22: end for	105
23: $V_T \leftarrow V_T \setminus (V \cup V')$	105
24: end while	105
25: $E \leftarrow Prune-Lines(E)$	105
26: return (V, E)	106
27: end procedure	106
28: procedure Prune-Lines(E)	106
29: sort E w.r.t confidence values in descending order	106
30: $E' \leftarrow \emptyset$	106
31: for all $e \in E$ do	106
32: if $\exists e' \in E' : e$ intersects with e' then	106
33: continue	106
34: end if	106
35: if $\exists e' \in E' : e' \cap e \neq \emptyset$ and $\angle(e, e') < \eta_\circ$ then	106
36: continue	107
$37: \qquad \text{end if} \qquad \qquad$	107
$38: \qquad E' \leftarrow E' \cup \{e\}$	107
39: end for	107
40: return E'	107
41: ena procedure	107
	107
	107





Figure 7: Failure cases on the SceneCity dataset.

method, designing a better objective function, or employing a RANSIC approach in those two stages.

A.4. Network Comparison

In this section, we discuss the difference between our network and the one in [12]. Our network is similar as their line detection network with the difference in the following perspectives:

1. In each hourglass module, they use two consecutive residual modules (RM) at each spatial resolution while we only use one RM, resulting less parameters in each hourglass module. Note that our design is the same as the original hourglass paper [23], which enables us to use more RM in each hourglass module and reduce the computational complexity since more computation is allocated in lower resolution stages. We adopt such design since we find using two RM gives negligible gains to the performance compared with only one.

2. We apply the intermediate supervision to the stacked hourglass network. For each hourglass modules, the loss term associated with the predicted heat maps is added to the final loss. [12] does not such intermediate supervision in their method. We find such intermediate supervision vital in both our synthetic dataset as well as their 2D dataset in terms of both accuracy and robustness.

3. We observe that using 2 stacked hourglass modules gives a similar performance as the 5 stacked hourglass modules in [12]. On the other hand, using 2 stacks consumes less memory, which gives us more design flexibility. For example, we are able to utilize a larger batch size to make the gradients more stable.