

# ULSD: Unified Line Segment Detection across Pinhole, Fisheye, and Spherical Cameras

Hao Li<sup>\*1</sup>, Huai Yu<sup>\*1,2</sup>, Wen Yang<sup>1</sup>, Lei Yu<sup>1</sup> and Sebastian Scherer<sup>2</sup>

**Abstract**—Line segment detection is essential for high-level tasks in computer vision and robotics. Currently, most state-of-the-art (SOTA) methods are dedicated to detecting straight line segments in undistorted pinhole images, thus distortions on fisheye or spherical images may largely degenerate their performance. Targeting at the unified line segment detection (ULSD) for both distorted and undistorted images, we propose to represent line segments with the Bezier curve model. Then the line segment detection is tackled by the Bezier curve regression with an end-to-end network, which is model-free and without any undistortion preprocessing. Experimental results on the pinhole, fisheye, and spherical image datasets validate the superiority of the proposed ULSD to the SOTA methods both in accuracy and efficiency (40.6fps for pinhole images). The source code is available at <https://github.com/lh9171338/Unified-Line-Segment-Detection>.

## I. INTRODUCTION

Line segment detection is one of the most fundamental problems in computer vision and robotics, which can facilitate many high-level vision tasks such as image matching [1], camera calibration [2], [3], structure from motion [4], [5], and visual SLAM [6], [7], [8], [9]. However, most current line segment detection methods model line segments as straight lines, thus cannot be directly applied to the distorted images from fisheye or spherical cameras, which are widely used in indoor camera localization [10], [11], room layout estimation [12], [13], and other tasks.

The existing methods for distorted line segment detection are almost model-based, i.e., dependent on camera distortion parameters. One category of methods among them requires rectification of lens distortion before applying straight line segments detector. The other methods such as extended Hough transform [17] and RANSAC-based methods [18], utilize camera distortion parameters to model the distorted line segments and can be directly applied to the distorted images. However, the performance of model-based methods is largely dependent on the accuracy of camera distortion parameters, which might be even unavailable in hand.

Therefore, the model-free approach for distorted line segment detection, i.e., independent of camera distortion parameters, is applaudable in practice. As a kind of model-free representation, two endpoints model has been commonly used in straight line segments detectors [?], [29], [30], but is not enough to represent curved line segments in distorted

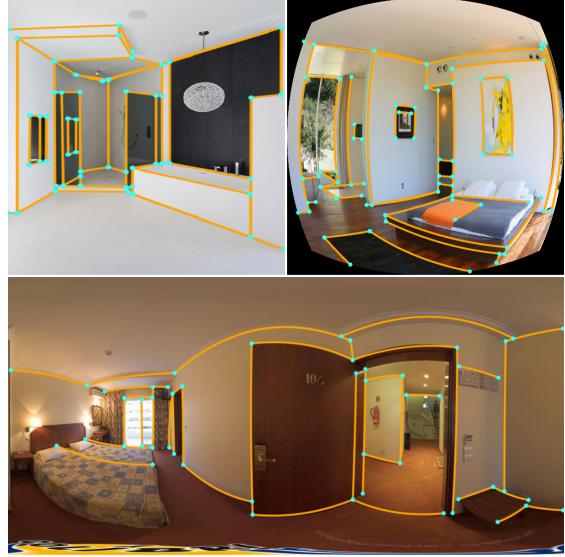


Fig. 1: Demonstration of the line segment detection on a pinhole (top-left), fisheye (top-right), and spherical image (bottom) with the proposed ULSD method.

images. Considering the strong fitting ability of the Bezier curve which has been successfully applied to arbitrarily-shaped text detection [19], we adopt the Bezier curve as a unified representation for line segments in both distorted and undistorted images. Similar to the two endpoints model, the Bezier curve can be represented as a vector parameterized by its equipartition points. Analogically, the line segment detection based on the Bezier curve model can be efficiently tackled with the coordinate regression for equipartition points, and the line classification which can be implemented by the center point detection of the target line segment [20].

Exploiting the Bezier curve model and detection method based on the center point, we design an end-to-end network for unified line segment detection (ULSD). The proposed ULSD is a model-free approach that can take both distorted and undistorted images from the pinhole, fisheye or spherical cameras as input, and directly output vectorized line segments, as shown in Fig. 1. The performance of ULSD is evaluated on various datasets which demonstrates the superiority to state-of-the-art (SOTA). As far as we know, the proposed ULSD is the first deep learning-based method to unify the line segment detection for both distorted and undistorted images. The main contribution of this work is three folded:

- We propose a model-free line segment representation

<sup>\*</sup>Equal contribution.

<sup>1</sup>Hao Li, Huai Yu, Wen Yang and Lei Yu are with the Electronic Information School, Wuhan University, Wuhan 430072, China {lihao2015, yuhuai, yangwen, ly.wd}@whu.edu.cn

<sup>2</sup>Sebastian Scherer is with the Robotics Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA basti@andrew.cmu.edu

for unified line segment detection based on the Bezier curve model, which is independent of camera distortion parameters.

- We design a unified end-to-end line segment detection network, which can be directly applied to both distorted and undistorted images from the pinhole, fisheye, and spherical cameras.
- We construct fisheye and spherical image datasets for line segment detection tasks in distorted images. With these datasets, the performance of the proposed ULSD is evaluated.

## II. RELATED WORK

Line segment detection is an attractive research topic during the last two decades in computer vision and robotics. Most of the related methods work well for undistorted pinhole images. For distorted images from fisheye or spherical cameras, the common practice is to undistort images and then deploy the straight-line models. Thus we will mainly review the straight line detection and distorted line detection methods.

### A. Straight Line Segment Detection

For pinhole images, a lot of line segment detection methods have been proposed. Among them, traditional approaches such as [24], [25], [26] detect lines based on the edge or gradient. The main drawback is that they are sensitive to noise and the detected lines are often fragmented. Recently, deep learning-based works such as [27], [28] significantly improve the performance by leveraging the deep features. Compared to the local edge or gradient features, the learning-based features are more robust to noise. Huang et al. [21] propose the wireframe representation and provide the first high-quality wireframe dataset. Compared with traditional line segment representation, wireframe parsing leverages the constraint of endpoint junctions, thus the output line segments are of higher quality in terms of line completeness and robustness to noise. Their method detects wireframe by two deep neural networks combined with a heuristic wireframe fusion algorithm. Zhou et al. propose the first end-to-end trainable neural network named L-CNN [29] for wireframe parsing. Based on L-CNN, Xue et al. introduce a holistic attraction field map (HAFM) [30] to represent line segments and achieve SOTA performance in accuracy and efficiency. Although HAFM brings significant improvement of performance, it needs big efforts to expand from the straight line to the distorted line.

### B. Distorted Line Segment Detection

As far as we know, there is no specially-designed method to detect distorted line segments for fisheye or spherical images. But there are some domains closely related. In the straight line-based fisheye image rectification field, various distorted lines detection methods have been used. These approaches can be divided into geometric-based methods [17], [18] and deep learning-based methods [2]. In [17], an extended Hough transform is utilized to detect the straight

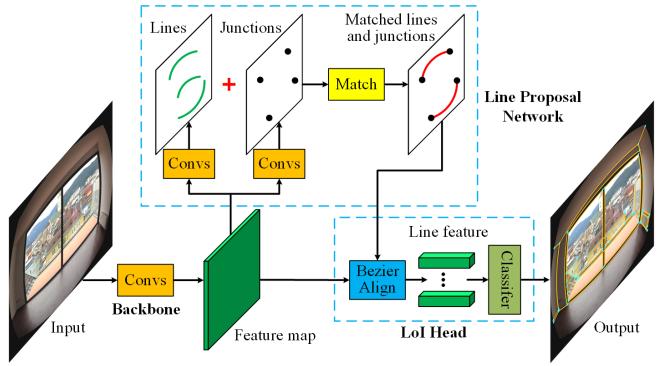


Fig. 2: An overview of our network architecture.

lines. The work [18] proposes a 5-points RANSAC method for robust line extraction. By leveraging the strong capability of networks, LaRecNet [2] can obtain more accurate results of distorted lines extraction. As for spherical images, there are some related works about line-based spherical camera localization [10], [11]. In [11], the authors first utilize Canny edge detection in the 2D equirectangular image, and then deploy spherical Hough transform to extract lines. In general, most of the above works are based on the traditional Hough transform or RANSAC algorithm. These methods are generally sensitive to noise and their detection accuracies are far from satisfactory. Additionally, there is no existing learning-based method modeling the distorted line segments in fisheye or spherical images.

## III. OVERVIEW OF THE PROPOSED METHOD

### A. Bezier Curve Representation

Since the most general line representation with the straight connection of two endpoints cannot fit lines in arbitrarily distorted images, we introduce the Bezier curve as a unified parameterized representation. The Bezier curve uses the Bernstein Polynomials as its basis to represent a parametric curve. Its definition is shown in Eq. 1.

$$B(t) = \sum_{i=0}^n b_i B_{i,n}(t), 0 \leq t \leq 1 \quad (1)$$

where  $t$  is the proportional coefficient of a point on the curve,  $n$  represents the order of the Bezier curve,  $b_i$  represents the  $i$ -th control point, and  $B_{i,n}(t) = C_n^i t^i (1-t)^{n-i}$  represents the Bernstein basis Polynomial.

According to Eq. 1, the interpolation formula for a Bezier curve can be obtained, as shown in Eq. 2:

$$\begin{bmatrix} B_{0,n}(t_0) & \cdots & B_{n,n}(t_0) \\ B_{0,n}(t_1) & \cdots & B_{n,n}(t_1) \\ \vdots & \ddots & \vdots \\ B_{0,n}(t_{m-1}) & \cdots & B_{n,n}(t_{m-1}) \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix} = \begin{bmatrix} p_0 \\ p_1 \\ \vdots \\ p_{m-1} \end{bmatrix} \quad (2)$$

where  $m$  is the number of the interpolation points,  $p_i$  represents the  $i$ -th interpolation point.

An  $n$ -th order Bezier curve can be determined by its  $n+1$  control points. However, the control points lack geometric meaning and may be located outside the image, so it's difficult to directly learn the positions of the control points.

Therefore, our network tries to predict the positions of the equipartition points of the Bezier curve instead, and then use Eq. 2 to calculate the control points through the least square method. As shown in Fig. 3, a 3rd Bezier curve can be determined by 4 control points  $b_i$ , but the prediction of  $b_1, b_2$  is non-trivial. Therefore, we represent the Bezier curve by using its  $n + 1$  equipartition points  $p_0, p_1, p_2, p_3$ .

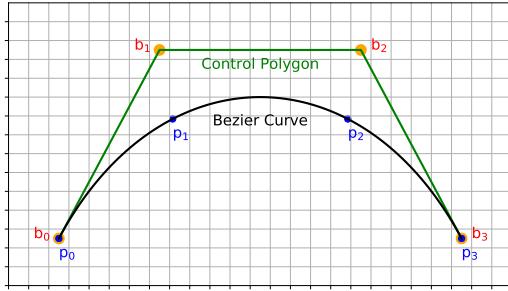


Fig. 3: A 3rd order Bezier curve.  $b_i$  represents the control points, and  $p_i$  represents the equipartition points.

### B. Overall Network Architecture

Based on the Bezier curve representation, our network is designed to detect line segments in arbitrarily distorted images (from the pinhole, fisheye, and spherical camera). Fig. 2 illustrates the network architecture. It mainly contains three modules: 1) a feature extraction backbone that takes a single image as input and outputs a shared feature map for the successive modules; 2) a Line Proposal Network (LPN) which outputs the candidate line segments; 3) an LoI (Line of Interest) head module which classifies the candidate line segments using the line features obtained through the BezierAlign module. Our pipeline is similar to the Faster R-CNN [31].

### C. Backbone

We choose the stacked hourglass network [32] as the backbone for its efficiency and effectiveness. Taking an image with size  $H \times W \times 3$  as input, the stacked hourglass network first downsamples the input image via convolution layers, then extracts features through multiple hourglass modules, and finally outputs the feature map with size  $H_b \times W_b \times C$ . The feature map is shared by the subsequent LPN and LoI head.

### D. Line Proposal Network

The Line Proposal Network contains four sub-modules: junction prediction module, line prediction module, line and junction matching module, and line sample module.

1) **Junction Prediction Module:** Junction prediction is addressed as a classification and regression problem. The input image with spatial size  $H \times W$  is divided into  $W_b \times H_b$  bins, same as the spatial size of the feature map. For each bin  $b$ , the network predicts whether there exists a junction inside it. If a junction  $p$  is inside bin  $b$ , it will also predict the offset vector from  $p$  to the center  $b$  of the bin. Therefore, the network outputs a junction confidence map  $J$  and a junction

offset map  $O$ . The ground truth of the two maps can be obtained by Eq. 3 and Eq. 4.

$$J(b) = \begin{cases} 1 & \exists p \in V : p \text{ inside } b \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

$$O(b) = \begin{cases} (b - p) & \exists p \in V : p \text{ inside } b \\ 0 & \text{otherwise.} \end{cases} \quad (4)$$

where  $V$  is the set of junctions.

$J$  and  $O$  are predicted by two decoder head consisting of two convolution layers separately. In the training phase, the binary cross-entropy loss and smooth  $l_1$  loss are used to predict  $J$  and  $O$  respectively. The total loss of junction prediction is the weighted sum of the two losses.

$$L_{junc} = \lambda_{conf}^j L_{conf}^j + \lambda_{offset}^j L_{offset}^j \quad (5)$$

Furthermore, the non-maximum suppression (NMS) is applied in  $J$  to remove duplicates, and only the top- $K$  junctions with the highest confidence are kept for the line and junction matching module.

2) **Line Prediction Module:** For an arbitrarily distorted line segment represented by an  $n$ -th order Bezier curve, the line prediction module tries to predict the location of the center point of the line segment and the offset vectors from the equipartition points to the center point. Center point prediction is the same as junction prediction. If  $n$  is even, the center point is one of the  $n + 1$  equipartition points, whose offset vector is 0, so only  $n$  offset vectors need to predict. The smooth  $l_1$  loss is used to predict the offset vectors, and the total loss is shown in Eq. 6.

$$L_{line} = \lambda_{center} L_{center} + \lambda_{offset} \sum_{i=1}^m L_{offset} \quad (6)$$

where  $m = n$  if  $n$  is even, otherwise  $m = n + 1$ .

3) **Line and Junction Matching Module:** To improve the quality of the line segment proposals, a line and junction matching module is adopted. The matching strategy is similar to HAWP [30]. A line segment proposal is kept if and only if its two endpoints can be matched with two junction proposals based on the Euclidean distance, and then the two endpoints of the line proposals are replaced by the two matched junction proposals. If there are multiple line proposals matched with the same pair of junction proposals, only the one with the shortest distance is kept.

4) **Line Sample Module:** The Line sample module is used to sample positive and negative line proposals for training the classifier of the LoI head. A line segment proposal is assigned with a positive label if there is a ground truth line segment and their distance calculated by Eq. 7 is less than a predefined threshold  $\eta$ . Otherwise, it's assigned with a negative label. Then, we can obtain a positive and a negative sample set. Finally, a certain number of positive and negative line segment proposals are randomly sampled from the two sets respectively. Apart from the above positive samples, we also sample some positive line segments from the ground truth to increase the number of positive samples and help

cold-start [29] the training at the beginning.

$$d(l, l') = \min \left( \sum_{i=0}^n \|p_i - p'_i\|^2, \sum_{i=0}^n \|p_i - p'_{n-i}\|^2 \right), \\ l = (\mathbf{p}_0, \dots, \mathbf{p}_n), l' = (\mathbf{p}'_0, \dots, \mathbf{p}'_n) \quad (7)$$

#### E. LoI Head Module

The LoI (Line of Interest) head module takes a list of candidate line segments together with the feature map  $F$  as input and predicts whether or not each candidate line segment is true. To extract the fixed-length line feature vector from the feature map, previous methods [29], [30] adopt LoI Pooling layer [29] which is based on the linear interpolation of straight line segments. However, LoI Pooling does not work for the distorted line segments. Therefore, we introduce the BezierAlign layer. Based on Eq. 2, it can uniformly sample  $N_p$  points from a distorted line segment.

The feature for each sampled point is computed from  $F$  using bi-linear interpolation. After a 1D max-pooling operator, all the features from the  $N_p$  points are concatenated as the line feature vector. After the BezierAlign operation, we feed all feature vectors into a classifier consisting of multiple fully-connected layers followed by a sigmoid layer, and finally get the confidence of each candidate line segment. The binary-cross entropy loss is used in the LoI head module. To balance the loss of positive and negative samples, the two losses are calculated and weighted separately. The total loss of the LoI head module is shown in Eq. 8.

$$L_{cls} = \lambda_{pos} L_{pos} + \lambda_{neg} L_{neg} \quad (8)$$

The loss of the whole network is the sum of all above losses.

$$L = L_{junc} + L_{line} + L_{cls} \quad (9)$$

## IV. EXPERIMENTS

In this section, we will discuss the datasets, network implementation details, and the experimental results.

#### A. Datasets

**Pinhole image.** The common Wireframe dataset [21] and YorkUrban dataset [22] are used. The former contains 5000 training images and 462 testing images. The latter contains 102 images.

**Fisheye image.** Currently, there is no available public fisheye line segment detection dataset, thus we build the F-Wireframe dataset and F-YorkUrban dataset by distorting images from the Wireframe dataset and YorkUrban dataset with the fisheye distortion model. The synthesized dataset is the same size as the original dataset.

**Spherical image.** Since no public spherical line segment dataset is available, we build a dataset by manually annotating images from the SUN360 dataset [23] and then augmenting by flip (horizontal, vertical, and horizontal-vertical) and periodical shifting (horizontal). Finally, the dataset contains 5200 training images and 68 testing images.

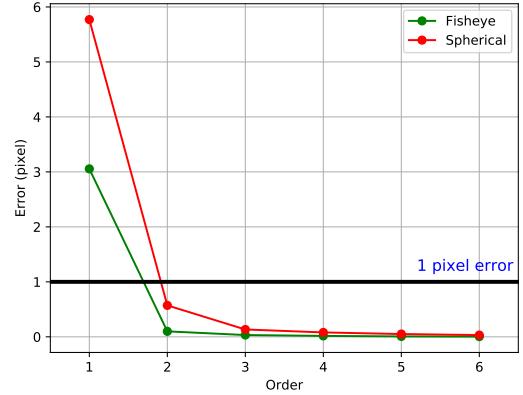


Fig. 4: Illustration of the fitting errors of Bezier curves with different orders for fisheye and spherical images.

#### B. Implementation Details

For the pinhole and fisheye networks, the input image is resized to  $(H, W) = (512, 512)$  and the feature map spatial size is  $(H_b, W_b) = (128, 128)$ . For the spherical network, the input image size and the feature map spatial size are  $(512, 1024)$  and  $(128, 256)$  respectively. The network's hyper-parameter settings are the same as L-CNN and HAWP for a fair comparison. Our network is trained using the Adam optimizer [33]. The learning rate, weight decay, and training batch size are set to  $4 \times 10^{-4}$ ,  $1 \times 10^{-4}$ , and 6 respectively. We use step decay as the learning rate scheduler. All experiments are conducted on a single NVIDIA GTX 2080Ti GPU.

As for the order setting of the Bezier curve, the 1st order is enough for pinhole images. While for fisheye and spherical images, the order is determined by experimental evaluations. We measure the fitting errors of Bezier curves with the order ranging from the 1 to 6. As shown in Fig. 4, the 2nd order Bezier curve is enough to make the fitting error less than 1 pixel for line segments in both spherical and fisheye images. Even though, the fitting error can be further decreased by increasing the order, this will largely increase the computational complexity. To balance the accuracy of line representation and efficiency of the network, we leave the setting for the order of the Bezier curve model in the following section by evaluating the performance of ULSD with the order of 2 to 4 for fisheye and spherical images.

#### C. Results and Comparisons

To evaluate the performance of the proposed ULSD, comparisons are made to the conventional methods LSD [26] and spherical Hough transform (SHT), deep learning-based methods DWP [21], AFM [27], L-CNN [29], and HAWP [30]. For pinhole images targeting to detect straight line segments, the USLD with 1st order Bezier curve (ULSD<sup>1</sup>) is trained on the Wireframe dataset and tested on both two pinhole image datasets. For fisheye and spherical images with distorted line segments, the ULSD with the 2nd, 3rd, and 4th order Bezier curve (ULSD<sup>2</sup>, ULSD<sup>3</sup>, ULSD<sup>4</sup>) are respectively evaluated by training on the F-Wireframe dataset and SUN360 dataset.

TABLE I: Quantitative results and comparisons on the Wireframe dataset and YorkUrban dataset.

Method	Wireframe Dataset			YorkUrban Dataset			FPS
	sAP <sup>10</sup>	msAP	mAP <sup>J</sup>	sAP <sup>10</sup>	msAP	mAP <sup>J</sup>	
LSD [26]	9.5	9.3	17.2	9.4	9.4	15.4	<b>50.9</b>
DWP [21]	6.8	6.6	38.6	2.7	2.7	23.4	2.3
AFM [27]	24.3	23.4	24.3	9.1	8.9	12.5	14.3
L-CNN [29]	62.9	62.1	59.3	26.4	26.1	30.4	13.7
HAWP [30]	<b>66.5</b>	<b>65.7</b>	60.2	<b>28.5</b>	<b>28.1</b>	<b>31.7</b>	30.9
ULSD <sup>1</sup> (ours)	66.4	65.6	<b>61.4</b>	27.4	27.0	31.0	<b>40.6</b>

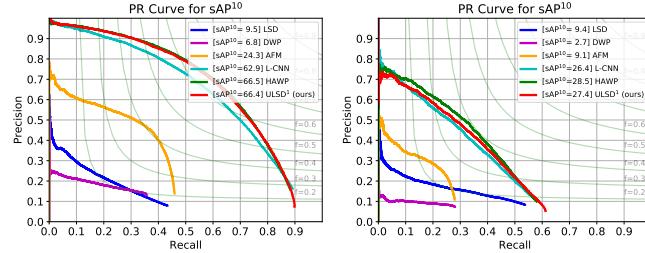


Fig. 5: Precision-Recall (PR) curves of sAP<sup>10</sup> on the Wireframe dataset (the left plot) and YorkUrban dataset (the right plot).

**Quantitative Results.** Quantitative performance in the accuracy of line segment detection is evaluated based on metrics including structural average precision (sAP) of line segments under the threshold of 5, 10, 15 pixels, mean structural average precision (msAP) over the threshold of 5, 10, and 15 pixels and vectorized junction mean AP (mAP<sup>J</sup>). And the efficiency of the aforementioned algorithms is further evaluated based on the frames per second (FPS).

The quantitative results and comparisons for pinhole images are shown in Table. I and Fig. 5. The proposed ULSD<sup>1</sup> obtains much higher detection accuracy compared with traditional LSD. Furthermore, compared to the deep learning-based methods, the proposed ULSD achieves the accuracy comparable to the SOTA, but with remarkable improvement in efficiency by at least 31%, which comes from the high efficiency of the ULSD’s line segment prediction module.

TABLE II: Quantitative results and comparisons on the F-Wireframe dataset and F-YorkUrban dataset.

Method	F-Wireframe Dataset			F-YorkUrban Dataset			FPS
	sAP <sup>10</sup>	msAP	mAP <sup>J</sup>	sAP <sup>10</sup>	msAP	mAP <sup>J</sup>	
LSD [26]	4.3	4.3	11.1	5.2	5.1	10.7	<b>47.9</b>
L-CNN [29]	43.4	42.9	44.2	19.9	19.6	26.4	14.3
HAWP [30]	46.3	45.6	43.8	21.5	21.2	26.4	31.5
HAWP <sup>+</sup>	56.4	55.4	-	25.8	25.4	-	31.5
ULSD <sup>2</sup> (ours)	<b>61.2</b>	<b>60.2</b>	<b>56.3</b>	<b>30.2</b>	<b>29.6</b>	32.6	<b>36.8</b>
ULSD <sup>3</sup> (ours)	60.3	59.3	56.1	28.6	28.0	31.5	36.5
ULSD <sup>4</sup> (ours)	59.9	59.0	<b>56.3</b>	30.1	<b>29.6</b>	<b>33.1</b>	36.3

Then by deploying higher-order Bezier curve representation, the quantitative performance for different algorithms is evaluated on fisheye and spherical datasets. For the fair comparison, the performance of SOTA straight line segment detection, HAWP, is further evaluated with an additional prepossessing by rectifying the fisheye images with randomly corrupted camera distortion parameters (10% random noise). And we denote it as HAWP<sup>+</sup>. The quantitative results are

TABLE III: Quantitative results and comparisons on the SUN360 dataset. The input image for L-CNN and HAWP is resized to 512 × 512, but ULSD remains the original size 512 × 1024, thus the speed of ULSD is a little slower than HAWP’s.

Method	SUN360 Dataset					FPS
	sAP <sup>5</sup>	sAP <sup>10</sup>	sAP <sup>15</sup>	msAP	mAP <sup>J</sup>	
SHT	0.9	1.7	2.5	1.7	3.4	0.05
L-CNN [29]	39.8	42.5	43.6	42.0	34.8	12.6
HAWP [30]	41.7	44.7	45.8	44.1	33.1	<b>25.4</b>
ULSD <sup>2</sup> (ours)	<b>61.9</b>	<b>67.6</b>	<b>69.8</b>	<b>66.4</b>	<b>47.3</b>	24.8
ULSD <sup>3</sup> (ours)	60.9	66.7	68.7	65.4	47.0	24.6
ULSD <sup>4</sup> (ours)	60.3	66.1	68.0	64.8	<b>47.3</b>	24.4

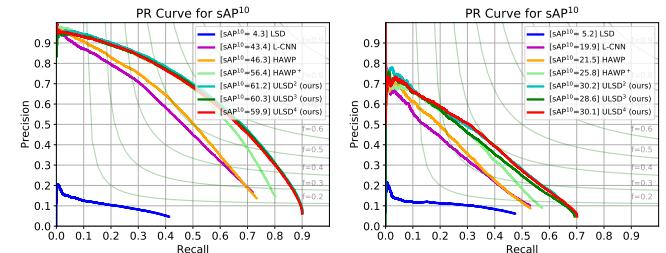


Fig. 6: PR curves of sAP<sup>10</sup> on the F-Wireframe dataset (the left plot) and F-YorkUrban dataset (the right plot).

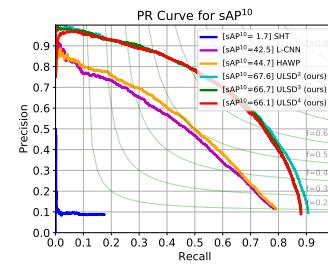


Fig. 7: PR curves of sAP<sup>10</sup> on the SUN360 dataset.

shown in Table. II, Table. III, and the corresponding PR curves in Fig. 6 and Fig. 7. For both fisheye and spherical images, the proposed ULSD outperforms the SOTA in both accuracy and efficiency. It is worth noting that although the rectification brings some performance improvement for HAWP<sup>+</sup> comparing to its origin, the proposed ULSD still outperforms HAWP<sup>+</sup>. Even more, ULSD is a model-free method that is independent of camera distortion parameters.

Comparing the results of ULSD on different order Bezier representations, we can observe that ULSD<sup>2</sup> is slightly better than higher-order ULSD<sup>3</sup> and ULSD<sup>4</sup> for both accuracy and efficiency. Although higher-order Bezier representation can improve the line fitting precision (Fig. 4), less than 1 pixel fitting error does not make big differences for the precision evaluation. Thus the higher-order Bezier curve representation does not effectively improve the line detection accuracy, instead makes the model more complex and difficult to learn, which results in lower accuracy. Thus, we evaluate the qualitative results with ULSD<sup>2</sup> for distorted images.

**Qualitative Results.** The qualitative results are shown in Fig. 8. Since LSD and SHT are based on the edge or gradient, they detect some noise edges without geometric meaning. Besides, they also produce a lot of fragmented line

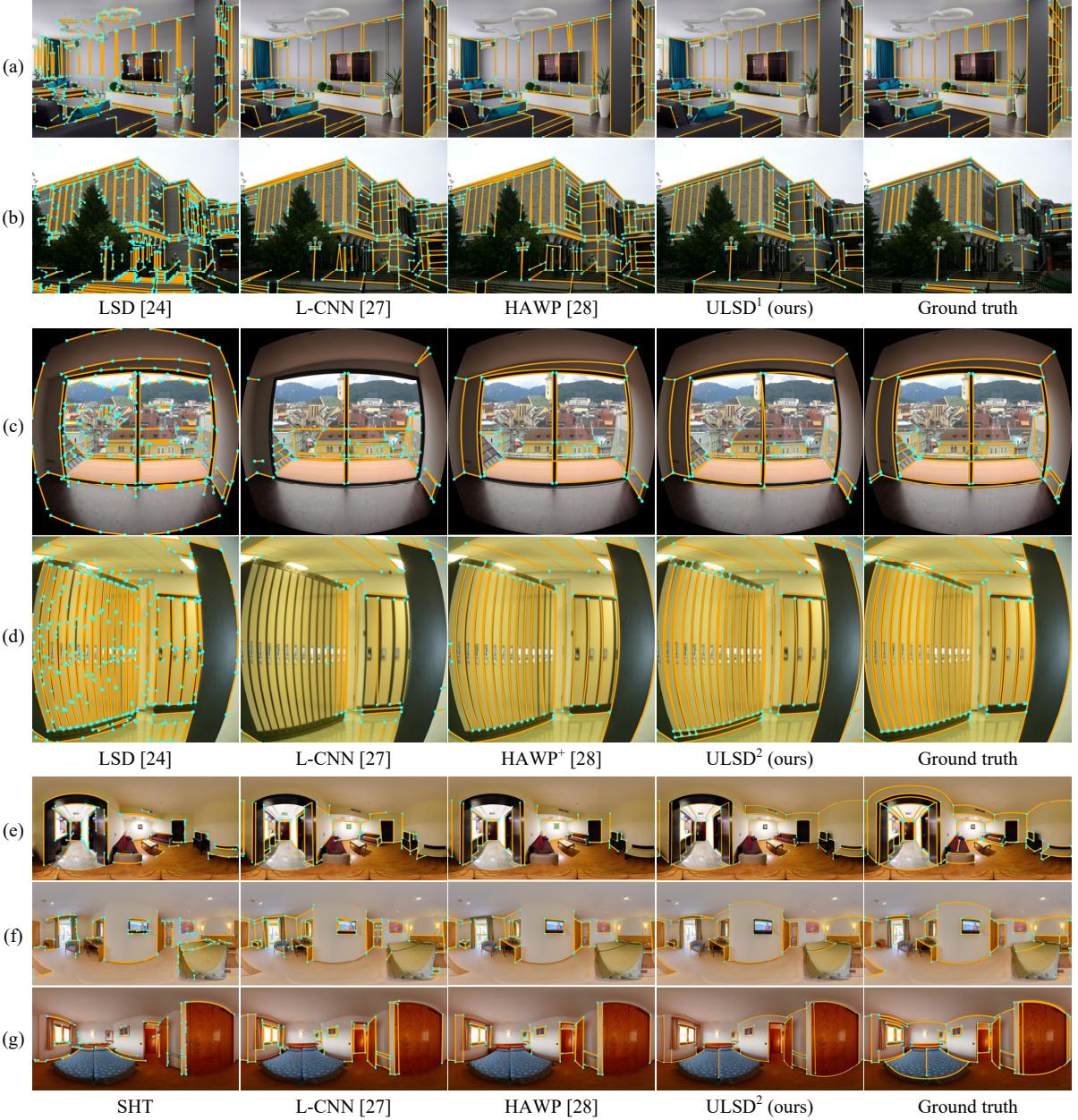


Fig. 8: Qualitative results and comparisons. Rows: (a)-(b) pinhole case, (c)-(d) fisheye case, (e)-(g) spherical case

segments due to not exploiting the constraint of junctions. By leveraging the learning-based features as well as the constraints of junctions, deep models L-CNN and HAWP can produce high-quality straight line segments detection results for pinhole images. However, restricted by the two-endpoint representation, they are incapable of detecting the distorted line segments for fisheye and spherical images. Even though, the performance of SOTA methods for straight line segments detection can be further improved by rectification preprocessing, our proposed ULSD armed with the parameterized Bezier curve model for line segments has the best performance both in accuracy and efficiency. Furthermore, ULSD can directly extract line segments in distorted or undistorted images.

## V. CONCLUSION

In this paper, we present a unified line segment detection method (ULSD) to detect arbitrarily distorted lines in pinhole, fisheye and, spherical images. With the novel Bezier curve representation, our network can formulate arbitrarily distorted line models in the three core modules: backbone, Line Proposal Network, and LoI head. Compared with the SOTA methods L-CNN [29] and HAWP [30], our ULSD achieves competitive results for pinhole images. More importantly, it obtains much higher line detection accuracy and efficiency for the fisheye and spherical images. With the real-time line detection speed using a single GPU, the proposed ULSD has great potential for online visual tasks such as SLAM and 3D reconstruction.

## REFERENCES

- [1] N. Xue, G. Xia, X. Bai, L. Zhang, and W. Shen, “Anisotropic-scale junction detection and matching for indoor images,” *IEEE Transactions on Image Processing*, vol. 27, no. 1, pp. 78–91, 2018.
- [2] Z. Xue, N. Xue, G. Xia, and W. Shen, “Learning to calibrate straight lines for fisheye image rectification,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1643–1651.
- [3] Mi Zhang, Jian Yao, Menghan Xia, Kai Li, Yi Zhang, and Yaping Liu, “Line-based multi-label energy optimization for fisheye image rectification and calibration,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 4137–4145.
- [4] B. Micusik and H. Widenauer, “Structure from motion with line segments under relaxed endpoint constraints,” in *2014 2nd International Conference on 3D Vision*, vol. 1, 2014, pp. 13–19.
- [5] G. Schindler, P. Krishnamurthy, and F. Dellaert, “Line-based structure from motion for urban environments,” in *Third International Symposium on 3D Data Processing, Visualization, and Transmission (3DPVT’06)*, 2006, pp. 846–853.
- [6] D. Zou, Y. Wu, L. Pei, H. Ling, and W. Yu, “Structvto: Visual-inertial odometry with structural regularity of man-made environments,” *IEEE Transactions on Robotics*, vol. 35, no. 4, pp. 999–1013, 2019.
- [7] T. Lemaire and S. Lacroix, “Monocular-vision based slam using line segments,” in *Proceedings 2007 IEEE International Conference on Robotics and Automation (ICRA)*, 2007, pp. 2791–2796.
- [8] H. Zhou, D. Zou, L. Pei, R. Ying, P. Liu, and W. Yu, “Structslam: Visual slam with building structure lines,” *IEEE Transactions on Vehicular Technology*, vol. 64, no. 4, pp. 1364–1375, 2015.
- [9] H. Yu, W. Zhen, W. Yang, J. Zhang, and S. Scherer, “Monocular camera localization in prior lidar maps with 2d-3d line correspondences,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020.
- [10] T. Goto, S. Pathak, Y. Ji, H. Fujii, A. Yamashita, and H. Asama, “Spherical camera localization in man-made environment using 3d-2d matching of line information,” in *Proceedings of the International Workshop on Advanced Image Technology 2017 (IWAIT2017)*, 2017.
- [11] T. Goto, S. Pathak, Y. Ji, H. Fujii, A. Yamashita, and H. Asama, “Line-based global localization of a spherical camera in manhattan worlds,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 2296–2303.
- [12] C. Zou, A. Colburn, Q. Shan, and D. Hoiem, “Layoutnet: Reconstructing the 3d room layout from a single rgb image,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 2051–2059.
- [13] C. Fernandez-Labrador, J. M. Facil, A. Perez-Yus, C. Demonceaux, J. Civera, and J. J. Guerrero, “Corners for layout: End-to-end layout recovery from 360 images,” *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1255–1262, 2020.
- [14] V. R. Kumar, S. A. Hiremath, M. Bach, S. Milz, C. Witt, C. Pinard, S. Yogamani, and P. Mäder, “Fisheyedistancenet: Self-supervised scale-aware distance estimation using monocular fisheye camera for autonomous driving,” in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 574–581.
- [15] M. Toromanoff, E. Wirbel, F. Wilhelm, C. Vejarano, X. Perrotton, and F. Moutarde, “End to end vehicle lateral control using a single fisheye camera,” in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 3613–3619.
- [16] Z. Cui, L. Heng, Y. C. Yeo, A. Geiger, M. Pollefeys, and T. Sattler, “Real-time dense mapping for self-driving vehicles using fisheye cameras,” in *2019 International Conference on Robotics and Automation (ICRA)*, 2019, pp. 6087–6093.
- [17] Q. Zhang and S. Kamata, “Fisheye image correction based on straight-line detection and preservation,” in *2015 IEEE International Conference on Systems, Man, and Cybernetics*, 2015, pp. 1793–1797.
- [18] Z. Mi, X. Hu, J. Yao, L. Zhao, J. Li, and J. Gong, “Line-based geometric consensus rectification and calibration from single distorted manhattan image,” *IEEE Access*, vol. PP, pp. 1–1, 10 2019.
- [19] Y. Liu, H. Chen, C. Shen, T. He, L. Jin, and L. Wang, “Abcnet: Real-time scene text spotting with adaptive bezier-curve network,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 9806–9815.
- [20] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, “Centernet: Keypoint triplets for object detection,” in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 6568–6577.
- [21] K. Huang, Y. Wang, Z. Zhou, T. Ding, S. Gao, and Y. Ma, “Learning to parse wireframes in images of man-made environments,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 626–635.
- [22] P. Denis, J. H. Elder, and F. J. Estrada, “Efficient edge-based methods for estimating manhattan frames in urban imagery,” in *Computer Vision – ECCV 2008*, 2008, pp. 197–210.
- [23] J. Xiao, K. A. Ehinger, A. Oliva, and A. Torralba, “Recognizing scene viewpoint using panoramic place representation,” in *2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 2695–2702.
- [24] D. Ballard, “Generalizing the hough transform to detect arbitrary shapes,” *Pattern Recognition*, vol. 13, no. 2, pp. 111 – 122, 1981.
- [25] J. Matas, C. Galambos, and J. Kittler, “Robust detection of lines using the progressive probabilistic hough transform,” *Computer Vision and Image Understanding*, vol. 78, no. 1, pp. 119 – 137, 2000.
- [26] R. Grompone von Gioi, J. Jakubowicz, J. Morel, and G. Randall, “Lsd: A fast line segment detector with a false detection control,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 4, pp. 722–732, 2010.
- [27] N. Xue, S. Bai, F. Wang, G. Xia, T. Wu, and L. Zhang, “Learning attraction field representation for robust line segment detection,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1595–1603.
- [28] Z. Zhang, Z. Li, N. Bi, J. Zheng, J. Wang, K. Huang, W. Luo, Y. Xu, and S. Gao, “Ppgnet: Learning point-pair graph for line segment detection,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 7098–7107.
- [29] Y. Zhou, H. Qi, and Y. Ma, “End-to-end wireframe parsing,” in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 962–971.
- [30] N. Xue, T. Wu, S. Bai, F. Wang, G. S. Xia, L. Zhang, and P. H. S. Torr, “Holistically-attracted wireframe parsing,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 2785–2794.
- [31] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [32] A. Newell, K. Yang, and J. Deng, “Stacked hourglass networks for human pose estimation,” in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., 2016, pp. 483–499.
- [33] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *International Conference on Learning Representations*, 12 2014.