

# Line Segment Detection Using Transformers without Edges

Yifan Xu\*, Weijian Xu\*, David Cheung, Zhuowen Tu  
 University of California San Diego  
 {yix081, wex041, d6cheung, ztu}@ucsd.edu

## Abstract

In this paper, we present a joint end-to-end line segment detection algorithm using Transformers that is post-processing and heuristics-guided intermediate processing (edge/junction/region detection) free. Our method, named LinE segment TRAnsformers (LETR), takes advantages of having integrated tokenized queries, a self-attention mechanism, and encoding-decoding strategy within Transformers by skipping standard heuristic designs for the edge element detection and perceptual grouping processes. We equip Transformers with a multi-scale encoder/decoder strategy to perform fine-grained line segment detection under a direct endpoint distance loss. This loss term is particularly suitable for detecting geometric structures such as line segments that are not conveniently represented by the standard bounding box representations. The Transformers learn to gradually refine line segments through layers of self-attention. In our experiments, we show state-of-the-art results on Wireframe and YorkUrban benchmarks.

## 1. Introduction

Line segment detection is an important mid-level visual process [22] useful for solving various downstream computer vision tasks, including segmentation, 3D reconstruction, image matching and registration, depth estimation, scene understanding, object detection, image editing, and shape analysis. Despite its practical and scientific importance, line segment detection remains an unsolved problem in computer vision.

Although dense pixel-wise edge detection has achieved an impressive performance [32], reliably extracting line segments of semantic and perceptual significance remains a further challenge. In natural scenes, line segments of interest often have heterogeneous structures within the cluttered background that are locally ambiguous or partially occluded. Morphological operators [27] operated on detected edges [3] often give sub-optimal results. Mid-level representations such as Gestalt laws [10] and contextual information [28] can

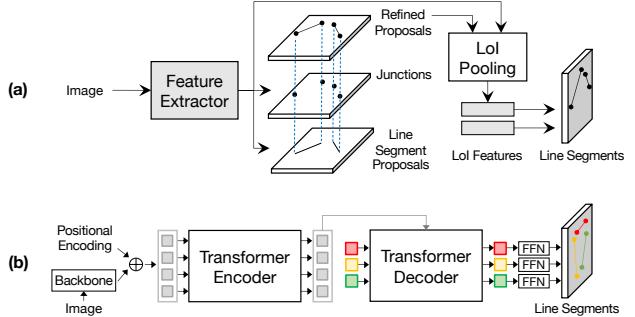


Figure 1. **Pipeline comparison** between: (a) holistically-attracted wireframe parsing (HAWP) [34] and (b) our proposed LinE segment TRAnsformers (LETR). LETR is based on a general-purpose pipeline without heuristics-driven intermediate stages for detecting junctions and generating line segment proposals.

play an important role in the perceptual grouping, but they are often hard to be seamlessly integrated into an end-to-end line segment detection pipeline. Deep learning techniques [16, 20, 14, 32] have provided greatly enhanced feature representation power, and algorithms such as [37, 33, 34] become increasingly feasible in real-world applications. However, systems like [37, 33, 34] still consist of heuristics-guided modules [27] such as edge/junction/region detection, line grouping, and post-processing, limiting the scope of their performance enhancement and further development.

In this paper, we skip the traditional edge/junction/region detection + proposals + perceptual grouping pipeline by designing a Transformer-based [29, 4] joint end-to-end line segment detection algorithm. We are motivated by the following observations for the Transformer frameworks [29, 4]: tokenized queries with an integrated encoding and decoding strategy, self-attention mechanism, and bipartite (Hungarian) matching step, capable of addressing the challenges in line segment detection for edge element detection, perceptual grouping, and set prediction; general-purpose pipelines for Transformers that are heuristics free. Our system, named LinE segment TRsformer (LETR), enjoys the modeling power of a general-purpose Transformer architecture while having its own enhanced property for detecting fine-grained geometric structures like line segments. LETR is built on top of a seminal work, DETR Transformer (DETR) [4]. However, as shown in Section 4.4 for ablation stud-

\* indicates equal contribution.

Code: <https://github.com/mlpc-ucsd/LETR>.

ies, directly applying the DETR object detector [4] for line segment detection does not yield satisfactory results since line segments are elongated geometric structures that are not feasible for the bounding box representations.

Our contributions are summarized as follows.

- We cast the line segment detection problem in a joint end-to-end fashion without explicit edge/junction/region detection and heuristics-guided perceptual grouping processes, which is in distinction to the existing literature in this domain. We achieve state-of-the-art results on the Wireframe [15] and YorkUrban benchmarks [5].
- We perform line segment detection using Transformers, based specifically on DETR [4], to realize tokenized entity modeling, perceptual grouping, and joint detection via an integrated encoder-decoder, a self-attention mechanism, and joint query inference within Transformers.
- We introduce two new algorithmic aspects to DETR [4]: first, a multi-scale encoder/decoder strategy as shown in Figure 2; second, a direct endpoint distance loss term in training, allowing geometric structures like line segments to be directly learned and detected — something not feasible in the standard DETR bounding box representations.

## 2. Related Works

### 2.1. Line Segment Detection

**Traditional Approaches.** Line detection has a long history in computer vision. Early pioneering works rely on low-level cues from pre-defined features (e.g. image gradients). Typically, line (segment) detection performs edge detection [3, 23, 7, 8, 32], followed by a perceptual grouping [13, 27, 10] process. Classic *perceptual grouping* frameworks [2, 1, 25, 21, 30] aggregate the low-level cues to form line segments in a bottom-up fashion: an image is partitioned into line-support regions by grouping similar pixel-wise features. Line segments are then approximated from line-support regions and filtered by a validation step to remove false positives. Another popular series of line segment detection approaches are based on *Hough transform* [9, 13, 24, 12] by gathering votes in the parameter space: the pixel-wise edge map of an image is converted into a parameter space representation, in which each point corresponds to a unique parameterized line. The points in the parameter space that accumulate sufficient votes from the candidate edge pixels are identified as line predictions. However, due to the limitations in the modeling/inference processes, these traditional approaches often produce sub-optimal results.

**Deep Learning Based Approaches.** The recent surge of deep learning based approaches has achieved much-improved performance on the line segment detection problem [15, 33, 37, 36, 34] with the use of learnable features to capture extensive context information. One typical family

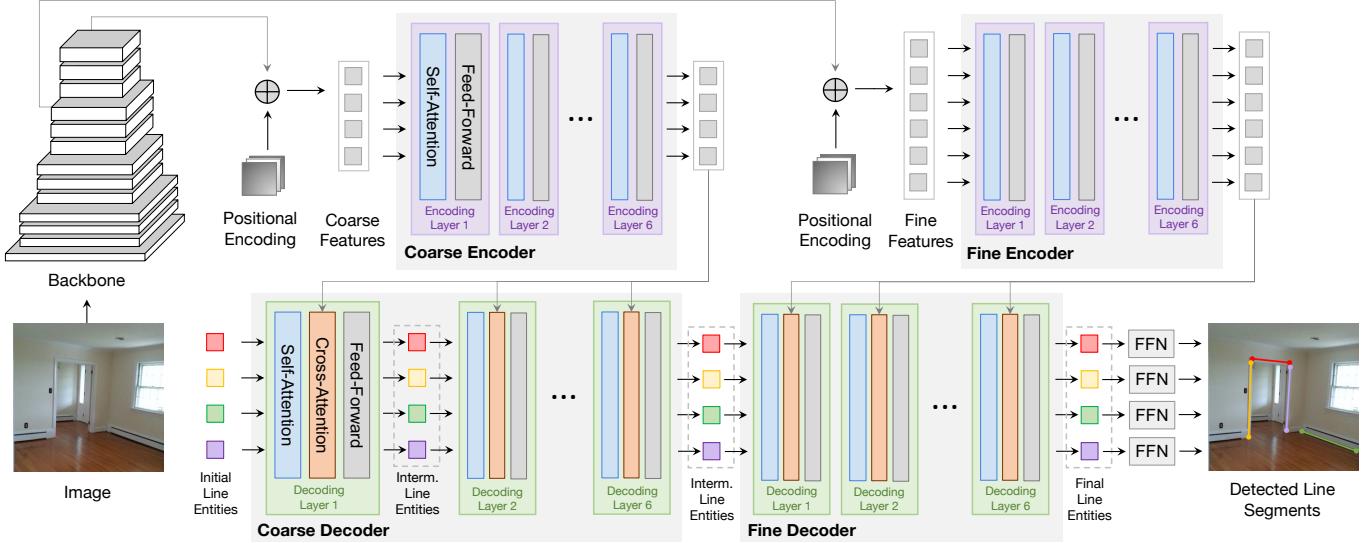
of methods is *junction-based pipelines*: Deep Wireframe Parser (DWP) [15] creates two parallel branches to predict the junction heatmap and the line heatmap, followed by a merging procedure. Motivated by [26], L-CNN [37] simplifies [15] into a unified network. First, a junction proposal module produces the junction heatmap and then converts detected junctions into line proposals. Second, a line verification module classifies proposals and removes unwanted false-positive lines. Methods like [37] are end-to-end, but they are at the instance-level (for detecting the individual line segments). Our LETR, like DETR [4], has a general-purpose architecture that is trained in a holistically end-to-end fashion. PPGNet [36] proposes to create a point-set graph with junctions as vertices and model line segments as edges. However, the aforementioned approaches are heavily dependent on high-quality junction detection, which is error-prone to various imaging conditions and complex scenarios.

Another line of approaches employs *dense prediction* to obtain a surrogate representation map and applies a post-process procedure to extract line segments: AFM [33] proposes an attraction field map as an intermediate representation that contains 2-D projection vectors pointing to associated lines. A squeeze module then recovers vectorized line segments from the attraction field map. Despite a relatively simpler design, [33] demonstrates its inferior performance compared with junction-based approaches. Recently, HAWP [34] builds a hybrid model of AFM [33], and L-CNN [37] by computing line segment proposals from the attraction field map and then refining proposals with junctions before further line verification.

In contrast, as shown in Figure 1, our approach differs from previous methods by removing heuristics-driven intermediate stages for detecting edge/junction/region proposals and surrogate prediction maps. Our approach is able to directly predict vectorized line segments while keeping competitive performances under a general-purpose framework.

### 2.2. Transformer Architecture

Transformers [29] have achieved great success in the natural language processing field and become *de facto* standard backbone architecture for many language models [29, 6]. It introduces self-attention and cross-attention modules as basic building blocks, modeling dense relations among elements of the input sequence. These attention-based mechanisms also benefit many vision tasks such as video classification [31], semantic segmentation [11], image generation [35], etc. Recently, end-to-end object detection with Transformers (DETR) [4] reformulates the object detection pipeline with Transformers by eliminating the need for hand-crafted anchor boxes and non-maximum suppression steps. Instead, [4] proposes to feed a set of object queries into the encoder-decoder architecture with interactions from the image feature sequence and generate a final set of predictions. A bipartite matching objective is then optimized to force unique



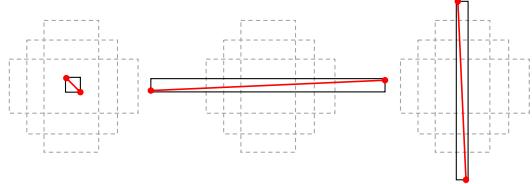
**Figure 2. Schematic illustration of our LETR pipeline:** An image is fed into a backbone network and generates two feature maps, which are then used by the coarse and the fine encoder respectively. Initial line entities are then first refined by the coarse decoder with the interaction of the coarse encoder output, and then the intermediate line entities from the coarse decoder are further refined by the fine decoder attending to the fine encoder. Finally, line segments are detected by feed-forward networks (FFNs) on top of line entities.

assignments between predictions and targets.

We introduce two new aspects to DETR [4] when realizing our LETR: 1) multi-scale encoder and decoder; 2) direct distance loss for the line segments.

### 3. Line Segment Detection with Transformers

#### 3.1. Motivation



**Figure 3. Bounding box representation.** Three difficult cases to represent line segments using bounding box diagonals. Red lines, black boxes, and gray dotted boxes refer to as line segments, the corresponding bounding boxes, and anchors respectively.

Despite the exceptional performance achieved by the recent deep learning based approaches [37, 33, 34] on line segment detection, their pipelines still involve heuristics-driven intermediate representations such as junctions and attraction field maps, raising an interesting question: *Can we directly model all the vectorized line segments with a neural network?* A naive solution could be simply regarding the line segments as *objects* and building a pipeline following the standard object detection approaches [26]. Since the location of 2-D objects is typically parameterized as a bounding box, the vectorized line segment can be directly read from a diagonal of the bounding box associated with



**Figure 4. Line entity representation.** For each row, we show how a same line entity predicts line segments with same property in three different indoor/outdoor scenes. The top line entity is specialized for horizontal line segments in the middle of the figure, and the bottom one prefers to predict vertical line segments with a various range of lengths.

the line segment object. However, the limited choices of anchors make it difficult for standard two-stage object detectors to predict very short line segments or line segments nearly parallel to the axes (see Figure 3). The recently appeared DETR [4] eliminates the anchors and the non-maximum suppression, perfectly meets the need of line segment detection. However, the vanilla DETR still focuses on bounding box representation with a GIoU loss. We further convert the box predictor in DETR into a vectorized line segment predictor by adapting the losses and enhancing the use of multi-scale features in our designed model.

### 3.2. Overview

In a line segment detection task, a detector aims to predict a set of line segments from given images. Performing line segment detection with Transformers removes the need of explicit edge/junction/region detection [37, 34] (see Figure 1). Our LETR is built purely based on the Transformer encoder-decoder structure. The proposed line segment detection process consists of four stages:

(1) *Image Feature Extraction*: Given an image input, we obtain the image feature map  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$  from a CNN backbone with reduced dimension. The image feature is concatenated with positional embeddings to obtain spatial relations. (2) *Image Feature Encoding*: The flattened feature map  $\mathbf{x} \in \mathbb{R}^{HW \times C}$  is then encoded to  $\mathbf{x}' \in \mathbb{R}^{HW \times C}$  by a multi-head self-attention module and a feed forward network module following the standard Transformer encoding architecture. (3) *Line Segment Detection*: In the Transformer decoder networks,  $N$  learnable line entities  $\mathbf{l} \in \mathbb{R}^{N \times C}$  interact with the encoder output via the cross-attention module. (4) *Line Segment Prediction*: Line entities make line segment predictions with two prediction heads built on top of the Transformer decoder. The line coordinates are predicted by a multi-layer perceptron (MLP), and the prediction confidences are scored by a linear layer.

**Self-Attention and Cross-Attention.** We first visit the scaled dot-product attention popularized by Transformer architectures [29]. The basic scaled dot-product attention consists of a set of  $m$  queries  $Q \in \mathbb{R}^{m \times d}$ , a set of  $n$  key-value pairs notated as a key matrix  $K \in \mathbb{R}^{n \times d}$  and a value matrix  $V \in \mathbb{R}^{n \times d}$ . Here we set  $Q, K, V$  to have same feature dimension  $d$ . The attention operation  $F$  is defined as:

$$F = \text{Att}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (1)$$

In our encoder-decoder Transformer architecture, we adopt two attention modules based on the multi-head attention, namely the self-attention module (SA) and cross-attention (CA) module (see Figure 2). The SA module takes in a set of input embeddings notated as  $\mathbf{x} = [x_1, \dots, x_i] \in \mathbb{R}^{i \times d}$ , and outputs a weighted summation  $\mathbf{x}' = [x'_1, \dots, x'_i] \in \mathbb{R}^{i \times d}$  of input embeddings within  $\mathbf{x}$  following Eq.1 where  $F = \text{Att}(Q = \mathbf{x}, K = \mathbf{x}, V = \mathbf{x})$ . The CA module takes in two sets of input embeddings notated as  $\mathbf{x} = [x_1, \dots, x_i] \in \mathbb{R}^{i \times d}$ ,  $\mathbf{z} = [x_1, \dots, x_j] \in \mathbb{R}^{j \times d}$  following Eq.1 where  $F = \text{Att}(Q = \mathbf{z}, K = \mathbf{x}, V = \mathbf{x})$ .

**Transformer Encoder in LETR** is stacked with multiple encoder layers. Each encoder layer takes in image features  $\mathbf{x} \in \mathbb{R}^{HW \times c}$  from its predecessor encoder layer and processes it with a SA module to learn the pairwise relation. The output features from SA module are passed into a point-wise fully-connected layer (FC) with activation and dropout layer followed by another point-wise fully-connected (FC) layer. Layer norm is applied between SA module and first

FC layer and after second FC layer. Residual connection is added before the first FC layer and after the second FC layer to facilitate optimization of deep layers.

**Transformer Decoder in LETR** is stacked with multiple decoder layers. Each decoder layer takes in a set of image features  $\mathbf{x}' \in \mathbb{R}^{HW \times C}$  from the last encoder layer and a set of line entities  $\mathbf{l} \in \mathbb{R}^{N \times C}$  from its predecessor decoder layer. The line entities are first processed with a SA module, each line entity  $l \in \mathbb{R}^C$  in  $\mathbf{l}$  attends to different regions of image feature embeddings  $\mathbf{x}'$  via the CA module. FC layers and other modules are added into the pipeline similar to the Encoder setting above.

**Line Entity Interpretation.** The *line entities* are analogous with the *object queries* in DETR [4]. We found each line entity has its own preferred existing region, length, and orientation of potential line segment after the training process (shown in Figure 4). We discuss line entities together make better predictions through self-attention and cross-attention refinement when encountering heterogeneous line segment structures in Section 4.4 and Figure 5.

### 3.3. Coarse-to-Fine Strategy

Different from object detection, line segment detection requires the detector to consider the local fine-grained details of line segments with the global indoor/outdoor structures together. In our LETR architecture, we propose a coarse-to-fine strategy to predict line segments in a refinement process. The process allows line entities to make precise predictions with the interaction of multi-scale encoded features while having an awareness of the holistic architecture with the communication to other line entities. During the coarse decoding stage, our line entities attend to potential line segment regions, often unevenly distributed, with a low resolution. During the fine decoding stage, our line entities produce detailed line segment predictions with a high resolution (see Figure 2). After each decoding layer at both coarse and fine decoding stage, we require line entities to make predictions through two shared prediction heads to make more precise predictions gradually.

**Coarse Decoding.** During the coarse decoding stage, we pass image features and line entities into an encoder-decoder Transformer architecture. The encoder receives coarse features from the output of Conv5 (C5) from ResNet with  $\frac{1}{32}$  original resolution. Then, line entity embeddings attend to coarse features from the output of the encoder in the cross-attention module at each layer. The coarse decoding stage is necessary for success at fine decoding stage and its high efficiency with less memory and computation cost.

**Fine Decoding.** The fine decoder inherits line entities from the coarse decoder and high-resolution features from the fine encoder. The features to the fine encoder come from the output of Conv4 (C4) from ResNet with  $\frac{1}{16}$  original resolution. The line entity embeddings decode feature information in the same manner as the coarse decoding stage.

### 3.4. Line Segment Prediction

In the previous decoding procedure, our multi-scale decoders progressively refine  $N$  initial line entities to produce same amount final line entities. In the prediction stage. Each final entity  $l$  will be fed into a feed-forward network (FFN), which consists of a classifier module to predict the confidence  $p$  of being a line segment, and a regression module to predict the coordinates of two end points  $\hat{\mathbf{p}}_1 = (\hat{x}_1, \hat{y}_1)$ ,  $\hat{\mathbf{p}}_2 = (\hat{x}_2, \hat{y}_2)$  that parameterizes the associated line segment  $\hat{\mathbf{L}} = (\hat{\mathbf{p}}_1, \hat{\mathbf{p}}_2)$ .

**Bipartite Matching.** Generally, there are many more line entities provided than actual line segments in the image. Thus, during the *training* stage, we conduct a set-based bipartite matching between line segment predictions and ground-truth targets to determine whether the prediction is associated with an existing line segment or not: Assume there are  $N$  line segment predictions  $\{(p^{(i)}, \hat{\mathbf{L}}^{(i)}); i = 1, \dots, N\}$  and  $M$  targets  $\{\mathbf{L}^{(j)}; j = 1, \dots, M\}$ , we optimize a bipartite matching objective on a permutation function  $\sigma(\cdot) : \mathbb{Z}_+ \rightarrow \mathbb{Z}_+$  which maps prediction indices  $\{1, \dots, N\}$  to potential target indices  $\{1, \dots, N\}$  (including  $\{1, \dots, M\}$  for ground-truth targets and  $\{M + 1, \dots, N\}$  for unmatched predictions):

$$\mathcal{L}_{\text{match}} = \sum_{i=1}^N \mathbb{1}_{\{\sigma(i) \leq M\}} [\lambda_1 d(\hat{\mathbf{L}}^{(i)}, \mathbf{L}^{(\sigma(i))}) - \lambda_2 p^{(i)}] \quad (2)$$

$$\sigma^* = \arg \min_{\sigma} \mathcal{L}_{\text{match}} \quad (3)$$

where  $d(\cdot, \cdot)$  represents  $L_1$  distance between coordinates and  $\mathbb{1}_{\{\cdot\}}$  is an indicator function.  $\mathcal{L}_{\text{match}}$  takes both distance and confidence into account with balancing coefficients  $\lambda_1, \lambda_2$ . The optimal permutation  $\sigma^*$  is computed using a Hungarian algorithm, mapping  $M$  positive prediction indices to target indices  $\{1, \dots, M\}$ . During the *inference* stage, we filter the  $N$  line segment predictions by setting a fixed threshold on the confidence  $p^{(i)}$  if needed due to no ground-truth provided.

### 3.5. Line Segment Losses

We compute line segment losses based on the optimal permutation  $\sigma^*$  from the bipartite matching procedure, in which  $\{i; \sigma^*(i) \leq M\}$  represents indices of positive predictions.

**Classification Loss.** Based on a binary cross-entropy loss, we observe that hard examples are less optimized after learning rate decay and decide to apply adaptive coefficients inspired by focal loss [18] to the classification loss term  $\mathcal{L}_{\text{cls}}$ :

$$\mathcal{L}_{\text{cls}}^{(i)} = -\mathbb{1}_{\{\sigma^*(i) \leq M\}} \alpha_1 (1 - p^{(i)})^\gamma \log p^{(i)} \quad (4)$$

$$-\mathbb{1}_{\{\sigma^*(i) > M\}} \alpha_2 p^{(i)^\gamma} \log(1 - p^{(i)}) \quad (5)$$

**Distance Loss.** We compute a simple  $L_1$ -based distance loss for line segment endpoint regression:

$$\mathcal{L}_{\text{dist}}^{(i)} = \mathbb{1}_{\{\sigma^*(i) \leq M\}} d(\hat{\mathbf{L}}^{(i)}, \mathbf{L}^{(\sigma^*(i))}) \quad (6)$$

where  $d(\cdot, \cdot)$  represents the sum of  $L_1$  distances between prediction and target coordinates. The distance loss is only applied to the positive predictions. Note that we remove the GIoU loss from [4] since GIoU is mainly designed for the similarity between bounding boxes instead of line segments. Thus, the final loss  $\mathcal{L}$  of our model is formulated as:

$$\mathcal{L} = \sum_{i=1}^N \lambda_{\text{cls}} \mathcal{L}_{\text{cls}}^{(i)} + \lambda_{\text{dist}} \mathcal{L}_{\text{dist}}^{(i)} \quad (7)$$

## 4. Experiments

### 4.1. Datasets

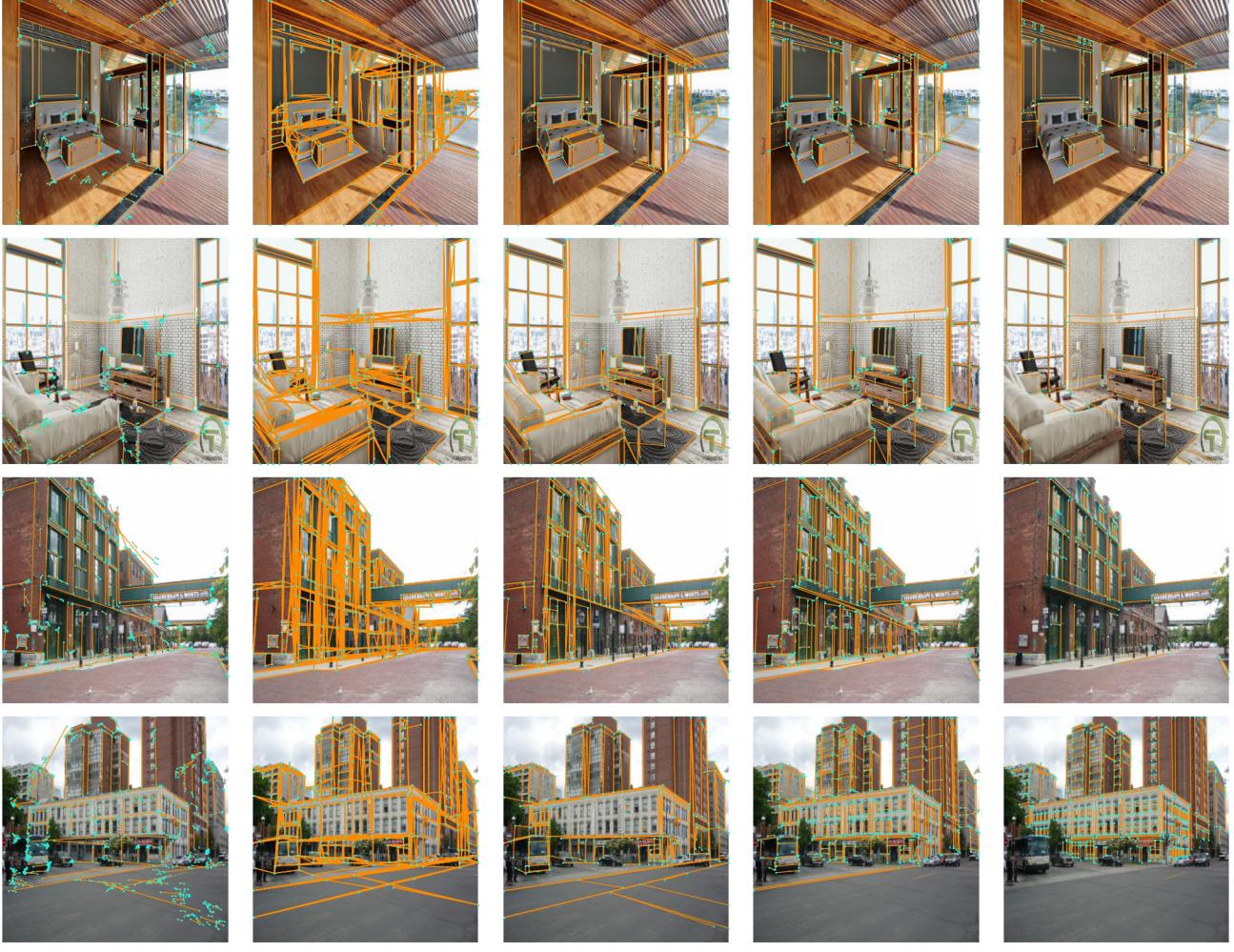
We train and evaluate our model on the ShanghaiTech *Wireframe* dataset [15], which consists of 5000 training images and 462 testing images. We also evaluate our model on the *YorkUrban* dataset [5] with 102 testing images from both indoor scenes and outdoor scenes.

Through all experiments, we conduct data augmentations for the training set, including random horizontal/vertical flip, random resize, random crop, and image color jittering. At the training stage, we resize the image to ensure the shortest size is at least 480 and at most 800 pixels while the longest size is at most 1333. At the evaluation stage, we resize the image with the shortest side at least 1100 pixels.

### 4.2. Implementation

**Networks.** We adopt both ResNet-50 and ResNet-101 as our feature backbone. For an input image  $X \in \mathbb{R}^{H_0 \times W_0 \times 3}$ , the coarse encoder takes in the feature map from the Conv5 (C5) layer of ResNet backbone with resolution  $x \in \mathbb{R}^{H \times W \times C}$  where  $H = \frac{H_0}{32}, W = \frac{W_0}{32}, C = 2048$ . The fine encoder takes in a higher resolution feature map ( $H = \frac{H_0}{16}, W = \frac{W_0}{16}, C = 1024$ ) from the Conv4 (C4) layer of ResNet. Feature maps are reduced to 256 channels by a 1x1 convolution and are fed into the Transformer along with the *sine/cosine* positional encoding. Our coarse-to-fine strategy consists of two independent encoder-decoder structures processing multi-scale image features. Each encoder-decoder structure is constructed with 6 encoder and 6 decoder layers with 256 channels and 8 attention heads.

**Optimization.** We train our model using 4 Titan RTX GPUs through all our experiments. Model weights from DETR [4] with ResNet-50 and ResNet-101 backbone are loaded as pre-training, and we discuss the effectiveness of pre-training in Section 5. We first train the coarse encoder-decoder for 500 epochs until optimal. Then, we freeze the weights in the coarse Transformer and train the fine Transformer initialized by coarse Transformer weights for 325 epochs (including a 25-epoch focal-loss fine-tuning). We adopt deep supervision [17, 32] for all decoder layers following DETR [4]. FFN prediction head weights are shared through all decoder layers. We use AdamW as the model



(a) AFM [33]

(b) LCNN [37]

(c) HAWP [34]

(d) LETR (ours)

(e) Ground-Truth

Figure 5. **Qualitative evaluation of line detection methods.** From left to right: the columns are the results from AFM [33], LCNN [37], HAWP [34], LETR (ours) and the ground-truth. From top to bottom: the top two rows are the results from the Wireframe test set, and the bottom two rows are the results from the YorkUrban test set.

optimizer and set weight decay as  $10^{-4}$ . The initial learning rate is set to  $10^{-4}$  and is reduced by a factor of 10 every 200 epochs for the coarse decoding stage and every 120 epochs for the fine prediction stage. We use 1000 line entities in all reported benchmarks unless specified elsewhere. To mitigate the class imbalance issue, we also reduce the classification weight for background/no-object instances by a factor of 10.

### 4.3. Evaluation Metric

We evaluate our results based on two heatmap-based metrics,  $AP^H$  and  $F^H$ , which are widely used in previous LSD task[37, 15], and Structural Average Precision (sAP) which is proposed in L-CNN [37]. On top of that, we evaluate the result with a new metric, Structural F-score (sF), for a more comprehensive comparison.

*Heatmap-based metrics,  $AP^H$ ,  $F^H$ :* Prediction and ground

truth lines are first converted to heatmaps by rasterizing the lines, and we generate the precision-recall curve comparing each pixel along their confidence. Then we can use the curve to calculate  $F^H$  and  $AP^H$ .

*Structural-based metrics, sAP[37], sF:* Given a set of ground truth line and a set of predicted lines, for each ground-truth line  $L$ , we define a predicted line  $\hat{L}$  to be a match of  $L$  if their  $L_2$  distance is smaller than the pre-defined threshold  $\vartheta \in \{10, 15\}$ . Over the set of lines matched to  $L$ , we select the line with the highest confidence as a true positive and treat the rest as candidates for false positives. If the set of matching lines is empty, we would regard this ground-truth line as false negative. Each predicted line would be matched to at most one ground truth line, and if a line isn't matched to any ground-truth line, then it is considered as a false positive. The matching is recomputed at each confidence

Table 1. **Comparison to prior work on Wireframe and YorkUrban benchmarks.** Our proposed LETR reaches state-of-the-art performance except sAP<sup>10</sup> and sAP<sup>15</sup> slightly worse than HAWP [34] in Wireframe. FPS Results for LETRs are tested on a single Tesla V100. Results for other prior works are adopted from HAWP paper.

Method	Wireframe Dataset						YorkUrban Dataset						FPS
	sAP <sup>10</sup>	sAP <sup>15</sup>	sF <sup>10</sup>	sF <sup>15</sup>	AP <sup>H</sup>	F <sup>H</sup>	sAP <sup>10</sup>	sAP <sup>15</sup>	sF <sup>10</sup>	sF <sup>15</sup>	AP <sup>H</sup>	F <sup>H</sup>	
LSD [30]	/	/	/	/	55.2	62.5	/	/	/	/	50.9	60.1	<b>49.6</b>
DWP [15]	5.1	5.9	/	/	67.8	72.2	2.1	2.6	/	/	51.0	61.6	2.24
AFM [33]	24.4	27.5	/	/	69.2	77.2	9.4	11.1	/	/	48.2	63.3	13.5
L-CNN [37]	62.9	64.9	61.3	62.4	82.8	81.3	26.4	27.5	36.9	37.8	59.6	65.3	15.6
HAWP [34]	66.5	68.2	64.9	65.9	86.1	83.1	28.5	29.7	39.7	40.5	61.2	66.3	29.5
<b>LETR (ours)</b>	<b>65.2</b>	<b>67.7</b>	<b>65.8</b>	<b>67.1</b>	<b>86.3</b>	<b>83.3</b>	<b>29.4</b>	<b>31.7</b>	<b>40.1</b>	<b>41.8</b>	<b>62.7</b>	<b>66.9</b>	5.04

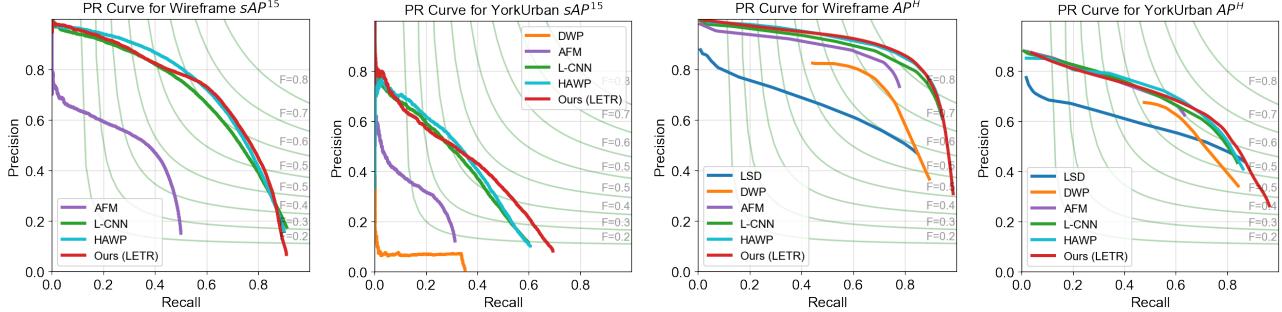


Figure 6. **Precision-Recall (PR) curves.** PR curves of sAP<sup>15</sup> and AP<sup>H</sup> for DWP[15], AFM[33], L-CNN[37], HAWP[34] and LETR (ours) on Wireframe and YorkUrban benchmarks.

level to produce the precision-recall curve, and we consider sAP as the area under this curve. Considering F<sup>H</sup> as the complementary F-score measurement for AP<sup>H</sup>, we evaluate the F-score measurement for sAP, denoted as sF, to be the best balanced performance measurement.

#### 4.4. Results and Comparisons

We summarize quantitative comparison results between LETR and previous line segment detection methods in Table 1. We report results for LETR with ResNet-101 backbone for Wireframe dataset and results with ResNet-50 backbone for York dataset. Our LETR achieves new state-of-the-art for all evaluation metrics on YorkUrban Dataset [5]. In terms of heatmap-based evaluation metrics, our LETR is consistently better than other models for both benchmarks and outperforms HAWP [34] by 1.5 for AP<sup>H</sup> on YorkUrban Dataset. We show PR curve comparison in Figure 6 on sAP<sup>15</sup> and AP<sup>H</sup> for both Wireframe [15] and YorkUrban benchmarks. In Figure 6, we notice the current limitation of LETR comes from lower precision prediction when we include fewer predictions compare to HAWP. When we include all sets of predictions, LETR predicts slightly better than HAWP and other leading methods, which matches our hypothesis that holistic prediction fashion can guide line entities to refine low confident predictions (usually due to local ambiguity and occlusion) with high confident predictions.

We also show both Wireframe and YorkUrban line segment detection qualitative results from LETR and other competing methods in Figure 5. The top two rows are indoor scene detection results from the Wireframe dataset, while the

bottom two rows are outdoor scene detection results from the YorkUrban dataset.

#### 5. Ablation Study

**Compare with Object Detection Baselines.** We compare LETR results with two object detection baseline where the line segments are treated as 2-D objects within this context in Table 2. We see clear limitations for using bounding box diagonal for both Faster R-CNN and DETR responding to our motivation in Section 3.1.

Table 2. **Comparison with object detection baselines** on Wireframe [15].

Method	sAP <sup>10</sup>	sAP <sup>15</sup>	sF <sup>10</sup>	sF <sup>15</sup>
Faster R-CNN	38.4	40.7	51.5	53.0
Vanilla DETR	53.8	57.2	57.2	59.0
LETR (ours)	65.2	67.7	65.8	67.1

**Effectiveness of Multi-Stage Training.** We compare the effectiveness of different modules in LETR in Table 3. During the coarse decoding stage, LETR reaches 62.3 and 65.2 for sAP<sup>10</sup> and sAP<sup>15</sup> with encoding features from the C5 layer of ResNet backbone, and 63.8 and 66.5 with the one from C4 of ResNet backbone. The fine decoder reaches 64.7 and 67.4 for sAP<sup>10</sup> and sAP<sup>15</sup> by improving the coarse prediction with fine-grained details from high-resolution features. We then adjust the data imbalance problem with focal loss to reach 65.2 and 67.7 for sAP<sup>10</sup> and sAP<sup>15</sup>.

As shown in Figure 7 (a), we found it is necessary to train the fine decoding stage after the coarse decoding stage converges. Training both stages together as a one-stage model

results a significant worse performance after 400 epochs.

**Effect of Number of Queries.** We found a large number of line entities is essential to the line segment detection task by experimenting on a wide range of the number of line entities (See Figure 7 (c), and using 1000 line entities is optimal for the Wireframe benchmark which contains 74 line segments in average).

Table 3. **Effectiveness of modules.** Ablation study of the architecture design and learning aspects in the proposed LETR on Wireframe dataset. (C) indicates the indexed feature used for coarse decoder; (F) indicates the indexed feature used for fine decoder.

Coarse Decoding	Fine Decoding	Focal Loss	Feature Index	sAP <sup>10</sup>	sAP <sup>15</sup>
✓			C5(C)	62.3	65.2
✓			C4(C)	63.8	66.5
✓	✓		C5(C), C4(F)	64.7	67.4
✓	✓	✓	C5(C), C4(F)	65.2	67.7

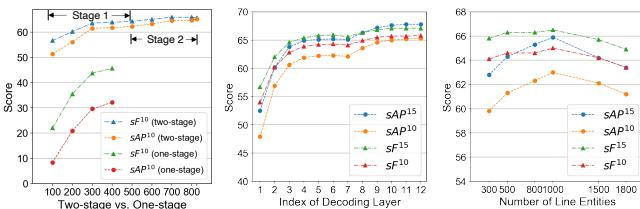


Figure 7. (a) **Multi-stage vs. single-stage training.** We compare results training coarse and fine layers in single stages and multi-stages (b) **Number of decoding layers.** We evaluate the performance of outputs from each decoding layer. The 1-6 layers are coarse decoder layers and 7-12 layers fine decoder layers. (c) **Number of line entities.** We test LETR (coarse decoding stage only) with different numbers of line entities on Wireframe.

**Effect of Image Upsampling.** All algorithms see the same input image resolution ( $640 \times 480$  typically). However, some algorithms try more precise predictions by upsampling images. To understand the impact of upsampling, we train and test HAWP and LETR under multiple upsampling scales. In Table 4 below, higher training upsampling resolution improves both methods. LETR obtains additional gains with higher test upsampling resolution.

Table 4. **Effectiveness of upsampling** with Wireframe dataset. LETR uses ResNet-101 backbone. \* Our LETR-512 resizes original image with the shortest size in a range between 288 and 512 † Our LETR-800 resizes original image with the shortest size in a range between 480 and 800.

	Train Size	Test Size	sAP <sup>10</sup>	sAP <sup>15</sup>	sF <sup>10</sup>	sF <sup>15</sup>
HAWP	512	512	65.7	67.4	64.7	65.8
HAWP	832	832	<b>67.7</b>	<b>69.1</b>	65.5	66.4
HAWP	832	1088	65.7	67.1	64.3	65.1
LETR	512*	512	61.1	64.1	63.1	64.8
LETR	800†	800	64.3	67.0	65.5	66.9
LETR	800†	1100	65.2	67.7	<b>65.8</b>	<b>67.1</b>

**Effectiveness of Pretraining.** We found model pretraining is essential for LETR to obtain state-of-the-art results. With DETR pretrained weights for COCO object detection [19], our coarse-stage-only model converges at 500 epochs. With CNN backbone pretrained weights for ImageNet classification, our coarse-stage-only model converges to a lower

score at 900 epochs. Without pretraining, LETR is difficult to train due to the limited amount of data in the Wireframe benchmark.

Table 5. **Effectiveness of pretraining.** We train LETR (coarse decoding stage only) with two variants. ImageNet represents LETR with ImageNet pretrained ResNet backbone. COCO represents LETR with COCO pretrained DETR weights.

Method	Epochs	sAP <sup>10</sup>	sAP <sup>15</sup>	sF <sup>10</sup>	sF <sup>15</sup>
ImageNet	900	58.4	62.0	62.4	64.6
COCO	500	62.3	65.2	64.3	65.9

## 6. Visualization

We demonstrate LETR’s coarse-to-fine decoding process in Figure 8. The first two columns are results from the coarse decoder receiving decoded features from the C5 ResNet layer. While the global structure of the scene is well-captured efficiently, the low-resolution features prevent it from making predictions precisely. The last two columns are results from the fine decoder receiving decoded features from the C4 ResNet layer and line entities from the coarse decoder. The overlay of attention heatmaps depicts more detailed relations in the image space, which is the key to the detector performance. This finding is also shown in Figure 7(b), where the decoded output after each layer has consistent improvement with the multi-scale encoder-decoder strategy.

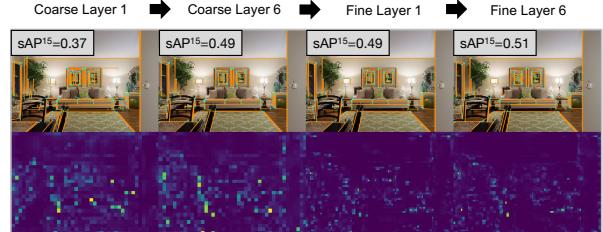


Figure 8. **Visualization of LETR coarse-to-fine decoding process.** From top to bottom: The 1<sup>st</sup> row shows line segment detection results based on line entities after different layers and the 2<sup>nd</sup> row shows its corresponding overlay of attention heatmaps. From left to right: The 1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup>, 4<sup>th</sup> columns are coarse decoder layer 1, coarse decoder layer 6, fine decoder layer 1, fine decoder layer 6, respectively.

## 7. Conclusion

In this paper, we presented LETR, a line segment detector based on a multi-scale encoder/decoder Transformer structure. By casting the line segment detection problem in a holistically end-to-end fashion, we perform set prediction without explicit edge/junction/region detection and heuristics-guided perceptual grouping processes. A direct endpoint distance loss allows geometric structures beyond bounding box representations to be modeled and predicted.

**Acknowledgment.** This work is funded by NSF IIS-1618477 and NSF IIS-1717431. We thank Justin Lazarow, Feng Han, Ido Durst, Yuezhou Sun, Haoming Zhang, and Heidi Cheng for valuable feedbacks.

## References

- [1] Michael Boldt, Richard Weiss, and Edward Riseman. Token-based extraction of straight lines. *IEEE Transactions on Systems, Man, and Cybernetics*, 19(6):1581–1594, 1989.
- [2] J. Brian Burns, Allen R. Hanson, and Edward M. Riseman. Extracting straight lines. *IEEE Trans. Pattern Anal. Mach. Intell.*, 8(4):425–455, 1986.
- [3] John Canny. A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence*, (6):679–698, 1986.
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Eur. Conf. Comput. Vis.*, 2020.
- [5] Patrick Denis, James H Elder, and Francisco J Estrada. Efficient edge-based methods for estimating manhattan frames in urban imagery. In *European conference on computer vision*, 2008.
- [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.
- [7] Piotr Dollár, Zhuowen Tu, and Serge Belongie. Supervised learning of edges and object boundaries. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2006.
- [8] Piotr Dollár and C Lawrence Zitnick. Structured forests for fast edge detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2013.
- [9] Richard O Duda and Peter E Hart. Use of the hough transformation to detect lines and curves in pictures. *Communications of the ACM*, 15(1):11–15, 1972.
- [10] James H Elder and Richard M Goldberg. Ecological statistics of gestalt laws for the perceptual organization of contours. *Journal of Vision*, 2(4):5–5, 2002.
- [11] Jun Fu, Jing Liu, Haijie Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3146–3154, 2019.
- [12] Yasutaka Furukawa and Yoshihisa Shinagawa. Accurate and robust line segment extraction by analyzing distribution around peaks in hough space. *Computer Vision and Image Understanding*, 92(1):1–25, 2003.
- [13] Nicolas Guil, Julio Villalba, and Emilio L Zapata. A fast hough transform for segment detection. *IEEE Transactions on Image Processing*, 4(11):1541–1548, 1995.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [15] Kun Huang, Yifan Wang, Zihan Zhou, Tianjiao Ding, Shenghua Gao, and Yi Ma. Learning to parse wireframes in images of man-made environments. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 626–635, 2018.
- [16] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. *Adv. Neural Inform. Process. Syst.*, 2012.
- [17] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial intelligence and statistics*, pages 562–570, 2015.
- [18] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Int. Conf. Comput. Vis.*, pages 2999–3007, 2017.
- [19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Eur. Conf. Comput. Vis.*, pages 740–755. Springer, 2014.
- [20] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. *IEEE Conf. Comput. Vis. Pattern Recog.*, 2015.
- [21] Xiaohu Lu, Jian Yao, Kai Li, and Li Li. Cannylines: A parameter-free line segment detector. In *2015 IEEE International Conference on Image Processing (ICIP)*, pages 507–511. IEEE, 2015.
- [22] David Marr. Vision: A computational investigation into the human representation and processing of visual information, henry holt and co. Inc., New York, NY, 2(4.2), 1982.
- [23] David R Martin, Charless C Fowlkes, and Jitendra Malik. Learning to detect natural image boundaries using local brightness, color, and texture cues. *IEEE transactions on pattern analysis and machine intelligence*, 26(5):530–549, 2004.
- [24] Jiri Matas, Charles Galambos, and Josef Kittler. Robust detection of lines using the progressive probabilistic hough transform. *Computer vision and image understanding*, 78(1):119–137, 2000.
- [25] Marcos Nieto, Carlos Cuevas, Luis Salgado, and Narciso García. Line segment detection using weighted mean shift procedures on a 2d slice sampling strategy. *Pattern Analysis and Applications*, 14(2):149–163, 2011.
- [26] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [27] Stephen M Smith and J Michael Brady. Susan—a new approach to low level image processing. *International journal of computer vision*, 23(1):45–78, 1997.
- [28] Zhuowen Tu. Auto-context and its application to high-level vision tasks. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2008.
- [29] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [30] R G von Gioi, J Jakubowicz, J M Morel, and G Randall. LSD: A Fast Line Segment Detector with a False Detection Control. *IEEE Trans. Pattern Anal. Mach. Intell.*, 32(4):722–732, 2010.
- [31] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018.
- [32] Saining Xie and Zhuowen Tu. Holistically-nested edge detection. In *Proceedings of the IEEE international conference on computer vision*, pages 1395–1403, 2015.
- [33] Nan Xue, Song Bai, Fudong Wang, Gui-Song Xia, Tianfu Wu, and Liangpei Zhang. Learning attraction field representation for robust line segment detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.

- [34] Nan Xue, Tianfu Wu, Song Bai, Fudong Wang, Gui-Song Xia, Liangpei Zhang, and Philip HS Torr. Holistically-attracted wireframe parsing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2788–2797, 2020.
- [35] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. In *International Conference on Machine Learning*, pages 7354–7363. PMLR, 2019.
- [36] Ziheng Zhang, Zhengxin Li, Ning Bi, Jia Zheng, Jinlei Wang, Kun Huang, Weixin Luo, Yanyu Xu, and Shenghua Gao. Ppgnet: Learning point-pair graph for line segment detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019.
- [37] Yichao Zhou, Haozhi Qi, and Yi Ma. End-to-end wireframe parsing. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 962–971, 2019.