

DeepLSD: Line Segment Detection and Refinement with Deep Image Gradients

Rémi Pautrat¹ Daniel Barath¹ Viktor Larsson² Martin R. Oswald^{1,3} Marc Pollefeys^{1,4}
¹ Department of Computer Science, ETH Zurich ² Lund University ³ University of Amsterdam
⁴ Microsoft Mixed Reality and AI Zurich Lab

Abstract

Line segments are ubiquitous in our human-made world and are increasingly used in vision tasks. They are complementary to feature points thanks to their spatial extent and the structural information they provide. Traditional line detectors based on the image gradient are extremely fast and accurate, but lack robustness in noisy images and challenging conditions. Their learned counterparts are more repeatable and can handle challenging images, but at the cost of a lower accuracy and a bias towards wireframe lines. We propose to combine traditional and learned approaches to get the best of both worlds: an accurate and robust line detector that can be trained in the wild without ground truth lines. Our new line segment detector, DeepLSD, processes images with a deep network to generate a line attraction field, before converting it to a surrogate image gradient magnitude and angle, which is then fed to any existing handcrafted line detector. Additionally, we propose a new optimization tool to refine line segments based on the attraction field and vanishing points. This refinement improves the accuracy of current deep detectors by a large margin. We demonstrate the performance of our method on low-level line detection metrics, as well as on several downstream tasks using multiple challenging datasets. The source code and models are available at <https://github.com/cvg/DeepLSD>.

1. Introduction

Line segments are ubiquitous in human-made environments and encode the underlying scene structure in a compact way. As such, line features have been used in multiple vision tasks: 3D reconstruction and Structure-from-Motion (SfM) [17, 32, 34], Simultaneous Localization and Mapping [13, 15, 25, 37, 62], visual localization [14], tracking [38], vanishing point estimation [49], etc. Thanks to their spatial extent and presence even in textureless areas, they offer a good complement to feature points [14, 15, 37].

All these applications require a robust and accurate detector to extract line features from images. Traditionally, line segments are extracted from the image gradient using

LSD [51]

HAWP [54]

Line distance field

Ours

Figure 1. Line detection in the wild. Top row: on challenging images, handcrafted methods such as LSD [51] suffer from noisy image gradients, while current learned methods like HAWP [54] were trained on wireframe images and generalize poorly. Bottom row: we combine deep learning to regress a line attraction field and a handcrafted detector to get both accurate and robust lines.

handcrafted heuristics, such as in the Line Segment Detector (LSD) [51]. These methods are fast and very accurate since they rely on low-level details of the image. However, they can suffer from a lack of robustness in challenging conditions such as in low illumination, where the image gradient is noisy. They also miss global knowledge from the scene and will detect any set of pixels with the same gradient orientation, including uninteresting and noisy lines.

Recently, deep networks offer new possibilities to tackle these drawbacks. This resurgence of line detection methods was initiated by the deep wireframe methods aiming at inferring the line structure of indoor scenes [19, 30, 53, 54, 61]. Since then, more generic deep line segment detectors have been proposed [10, 16, 20, 28, 50], including joint line detectors and descriptors [1, 36, 56]. These methods can, in theory, be trained on challenging images and, thus, gain robustness to be able to handle the extent of line segments in

an image, they can also encode some image context and distinguish between noisy and relevant lines. On the other

hand, most of these methods are fully supervised and there exists currently only a single dataset with ground truth lines, the Wireframe dataset [19]. Initially designed for wireframe parsing, this dataset is biased towards structural lines and is limited to indoor scenes. Therefore, it is not a suitable training set for generic line detectors, as illustrated in Figure 1. Additionally, similarly as with feature points [31, 45], current deep detectors are lacking accuracy and are still outperformed by handcrafted methods on easy images. The exact localization of line endpoints is often hard to obtain, as lines can be fragmented and suffer from partial occlusion. Many applications using lines consequently consider in noise lines and ignore the endpoints [34].

Based on this assessment, we propose in this work to keep the best of both worlds: use deep learning to process the image and discard unnecessary details, then use handcrafted methods to detect the line segments. We thus retain the benefits of deep learning, namely, to abstract the image and gain more robustness to illumination and noise, while at the same time retaining the accuracy of classical methods. We achieve this goal by following the tracks of two previous methods that used a dual representation of line segments, with attraction fields [53, 54]. The latter are continuous representations that are well-suited for deep learning, and we show how to leverage them as input to the traditional line detectors. Contrary to these two previous methods, we do not rely on ground truth lines to train our line attraction field, but propose instead to bootstrap existing methods to create a high-quality pseudo ground truth. Thus, our network can be trained on any dataset and be specialized towards specific applications, which we show in our experiments.

We additionally propose a novel optimization procedure to refine the detected line segments. This refinement is based on the attraction field output by the proposed network, as well as on vanishing points, optimized together with the segments. Not only can this optimization be used to effectively improve the accuracy of our prediction, but it can also be applied to other deep line detectors.

In summary, we propose the following contributions:

- We propose a method bootstrapping current detectors to create ground truth line attraction fields on any image.
- We introduce an optimization procedure that can simultaneously refine line segments and vanishing points. This optimization can be used as a stand-alone refinement to improve the accuracy of any existing deep line detector.
- We set a new record in several downstream tasks requiring line segments by combining the robustness of deep learning approaches with the precision of handcrafted methods in a single pipeline.

2. Related Work

Handcrafted Line Detectors. Detecting line segments in images is traditionally performed based on the image gradient. Early methods threshold the gradient magnitude to keep only strong edges and search for aligned sets of pixels sharing the same gradient angle. LSD [51] grows line regions, fits a rectangle to the resulting set of pixels, and finally extract a line segment. EDLines [4] grows the line regions in one direction only, orthogonal to the image gradient. Several extensions of these methods have been proposed, such as the multi-scale version of LSD, MLSLSD [41], and ELSLSD [48], a faster version of EDLines which avoids breaking lines in case of small discontinuities. AG3Line [58] proposes to actively group the seed points and adds line geometry constraints. Another approach consists in detecting full lines with the Hough transform [18] in a first step, then finding segments within these lines [12]. Since all these methods rely on low-level details of the image, they are highly accurate and fast, but lack robustness to noise and low illumination.

Learned Line Detectors. Deep line detection was first introduced through the task of wireframe parsing, i.e. estimating the structural lines of a scene [19]. Several approaches have been proposed to parameterize and represent the line segments, e.g., with two endpoints [61], attraction fields [53, 54], center and offset to the endpoints [20], graphs [33, 59], and transformers [52]. Wireframes can be further improved through a Deep Hough transform [30]. All these methods are trained on a single dataset, the Wireframe dataset [19], and they are not necessarily suitable for other tasks such as visual localization and SfM. Generic deep line segment detectors have also been proposed, with a focus on efficiency [10, 16], and can improve visual localization with points and lines [14]. However, these methods are again trained solely on the Wireframe dataset and their predicted lines are biased towards structural lines and indoor scenes.

Some works also perform a joint detection and description of line segments. SOLD2 [36] introduced a self-supervised training, using the homography adaptation technique initially described in SuperPoint [11]. ELSD [56] and L2D2 [1] both propose similar networks, but ELSD is again trained on the Wireframe dataset, while L2D2 uses a novel process to extract a line ground truth from LiDAR scans. Though these approaches are a first step towards unsupervised line detection, they still lack accuracy.

Attraction Fields. This work proposes to combine deep learning methods with classical line extractors. The key component for this is to use a dual representation of lines through an attraction field. This representation was first introduced by Xue et al. [53] for the wireframe task, and later improved with HAWP [54, 55]. They represent the set of discrete lines of an image with a continuous 2D vector field, suitable for deep networks. We adopt a similar approach, with small

Figure 2. Overview of the method. (1) We generate ground truth line distance and angle fields (DF/AF) by bootstrapping LSD [51]. (2) A deep network is trained to predict the DF/AF, which is then converted to a surrogate image gradient. (3) Line segments are extracted with LSD and (4) refined based on the DF/AF.

modifications to make the prediction more accurate. While adding two angles pointing at the endpoints of the closest not exactly an attraction field, Teplyakov et al. [50] also line. Recovering the original segments from the attraction proposed to predict a line mask and line angle field with field is then straightforward.

a network, then used LSD [51] to get line segments. Our However, this representation is not optimal to obtain accurate line segments, as illustrated in Figure 3. Directly pre- instead of a simple binary mask. Attraction fields have also dicting the position of the endpoints as done in HAWP [54] been leveraged for keypoint detection [21], where 2D vec- requires a larger receptive field to be able to get information tors are voting for the closest keypoint in the image. These from far-away endpoints, so that the network will focus on detections-by-voting offer a convenient way to represent dis- higher-level details instead of low-level ones. Additionally, crete quantities through continuous ones, and are also a key- deep networks are still struggling to yield accurate keypoint aspect of our approach when it comes to generating a reliable- detections [31, 45], which holds even more for line endpoints, ground truth for line detection.

3. Hybrid Line Detector

We demonstrate how to combine the robustness of deep- we propose to restrict our network to a smaller receptive field and to let the traditional heuristics determine the endpoints. networks together with the accuracy of handcrafted line and detectors. We train a deep network to predict a line attraction We adopt a similar attraction field representation as field, convert it to a surrogate image gradient, and feed it to HAWP [54] but without the additional two angles point- a handcrafted line detector to obtain the segments. Finally, ing at endpoints, yielding only line distance field (DF) an optimization based on the attraction field is used to refine and a line angle field (AF). For every pixel in these two the lines, as depicted in Figure 2. images, the line distance field gives the distance from the current pixel to the closest point on a line, and the line angle field A returns the orientation of the closest line. These two quantities can be easily obtained from the 2D offset field $(x; y) \in \mathbb{R}^H \times \mathbb{R}^W \rightarrow \mathbb{R}^H \times \mathbb{R}^W$ pointing at the closest point on a line, where $(H; W)$ are the dimension of the image:

3.1. Line Attraction Field

Representing line segments through an attraction field was first proposed by Xue et al. [53]. They initially proposed to regress a 2D vector field for each pixel of an image, indicating the relative position of the closest point on a line. This approach allows to represent discrete quantities (the line segments) as a smooth 2-channel image well suited for deep- We define here the line angle modulo so that a pixel above or below a line would have the same angle. Adopting this

$$D = \sqrt{x^2 + y^2}; \quad A = \arctan \frac{y}{x} + \pi \mod 2\pi \quad (1)$$

(a) AFM [53] (b) HAWP [54] (c) Distance eld (d) Angle eld

Figure 3. Attraction eld parametrizations. (a) Parametrizing with 2D vectors may produce noisy angles for small vector norms. (b) Adding offsets to the endpoints requires long-range information and is not robust to noisy endpoints. We propose to decouple the distance eld (c) and line orientation eld (d).

parametrization has the advantage of separating the norm from the angle of the 2D offset. Traditional detectors are outputting the distance eld $\hat{D} \in \mathbb{R}^{H \times W}$ and the other one leveraging the image gradient magnitude and angle, so we adopt a similar representation. Furthermore, both quantities are continuous close to line segments, and the line angle is even constant close enough to a line.

3.2. Ground Truth Generation

To learn the attraction eld, a ground truth is needed. Both AFM [53] and HAWP [54] are supervised with the ground truth lines of the Wireframe dataset [19]. We explore a novel method to acquire our ground truth, by bootstrapping with a ReLU activation and outputs a normalized distance eld $\hat{D}_n \in \mathbb{R}^{H \times W}$. Inspired by SuperPoint [11] and SOLD2 [36], we propose to generate the ground truth attraction eld through homography adaptation. Given a single input image I , we warp it with N random homographies H_i , detect line segments in all the warped images using any existing line detector, and then warp back the segments into I to get a set L_i of lines. We use LSD [51] to extract lines as it is currently among the most accurate existing line detector. The next step is to aggregate all the detections together, however, aggregating discrete quantities such as lines is non-trivial. SOLD2 [36] proposed to aggregate the endpoints and line heatmaps, and recover the segments afterwards. Instead, we propose to convert the sets of lines into a distance eld D_i and angle eld A_i , and to aggregate them by taking the median value of each pixel (u, v) across all images:

$$\begin{aligned} D(u; v) &= \text{median}_{i \in [1; N]} D_i(u; v) \\ A(u; v) &= \text{median}_{i \in [1; N]} A_i(u; v) \end{aligned} \quad (2)$$

By taking the median, we remove the noisy lines that were detected in only a few images, as shown in Figure 4.

3.3. Learning the Line Attraction Field

To regress our line distance and angle elds, we leverage a UNet-like neural network architecture [40]. The input image of size $(H; W)$ is processed by several convolutional layers and gradually downsampled up to a factor of 8 through 3 successive average pooling operations. The features are then upsampled back to the original resolution through another series of convolutional layers and bilinear interpolation. The resulting deep features are then split into two branches, one of the angles into account for the given predicted and ground

Input image Distance eld Angle eld

Figure 4. Pseudo GT visualization. Given an input image, we generate a line distance and angle elds (color coded [5]) and use them to supervise a deep network. Noisy lines, such as the ones in the bush at the bottom, are averaged out and ignored.

While all convolutions are followed by ReLU [2] and Batch Normalization [22], the last two outputs have different activations. The angle eld is obtained through a sigmoid activation and is multiplied by π to get an angle within $[0; \pi]$.

Since the distance eld can get very small values close to lines, where we also want the highest accuracy, we adopt a special normalization. The distance eld branch ends with a ReLU activation and outputs a normalized distance eld $\hat{D}_n \in \mathbb{R}^{H \times W}$. The final distance eld is obtained through the following denormalization:

$$\hat{D} = r \cdot e^{\hat{D}_n}; \quad (3)$$

where r is a parameter in pixels that defines a region around each line. Since handcrafted methods mainly need gradient information close to line segments, we supervise our network only on pixels at a distance of less than r pixels from a line. By selecting a small value for r , large portions of the image may not have any supervision, including areas where the pseudo ground truth was not able to detect real lines, e.g. lines with small contrast. Enforcing these lines to be in the background during training, i.e. with high distance eld, provides a detrimental training signal and decreases the recall of the prediction. On the contrary, with our loose supervision, these low contrast lines are not penalized during training and our trained model can detect them, thus yielding a more complete prediction than the ground truth.

We compute the training loss by comparing with a normalized version of the ground truth $\hat{D}_n = \log \frac{\hat{D}}{r}$. Note that since we only supervise pixels with a distance eld below r , $\frac{\hat{D}}{r} \in [0; 1]$ and so $\hat{D}_n \in \mathbb{R}^+$. We compute the total loss as the sum of the losses for the distance eld and the angular eld:

$$L = L_D + L_A; \quad (4)$$

where L_D is an L1 loss between the normalized distance elds and L_A is an L2 angular loss that takes the circularity of the angles into account for the given predicted and ground

(a) A double edge

(b) HAWP [54]

(c) Ours

Figure 5. Distinguishing double edges. (a) An example of a bright-dark-bright edge and the oriented angle field. (b) Wireframe methods treat it as a single line. (c) We detect it as two lines for better accuracy.

truth angle fields $\hat{A}; A \in [0; \pi]^H \times W$:

$$\begin{aligned} L_D &= \sum_{j,j'} \hat{D}_n(D_{n,j,j'})^2; \\ L_A &= \min(\sum_{j,j'} \hat{A}_{jj'} - A_{jj'}^2; \sum_{j,j'} |j - j'| \hat{A}_{jj'} - A_{jj'}^2); \end{aligned} \quad (5)$$

3.4. Extracting Line Segments

Since handcrafted detectors are based on the image gradient, we propose to convert our distance and angle fields into a surrogate image gradient magnitude and angle:

$$\begin{aligned} M &= r \cdot \hat{D} \\ &= \hat{A} \cdot \bar{z} \end{aligned} \quad (6)$$

Our predicted angle follows the directions of the lines and is perpendicular to the image gradient, so we rotate it by \bar{z} . The maximal magnitude of a pixel on a line is

An important difference between the approaches of AFM and LSD is the gradient orientation. For an edge separating a dark from a bright area, LSD keeps track of the dark-to-bright gradient direction, while AFM does not. This becomes important when several parallel lines occur next to each other in a dark-bright-dark or bright-dark-bright pattern, as illustrated in Figure 5. For better accuracy and scale-invariance, we advocate to detect these double edges and make our predicted angle oriented based on the sign of the image gradient angle:

$$\theta = \begin{cases} d(\theta; \theta_1) < d(\theta; \theta_2) \\ \text{otherwise} \end{cases}; \quad (7)$$

where $d(\cdot; \cdot)$ is a circular distance between two angles. Now equipped with an oriented angle θ and magnitude M , we can directly apply any existing classical line segment detector. Unless stated otherwise, we always use the LSD [51]

approach in the following, due to its high accuracy. In summary, the purpose of the deep net is to suppress image noise and detect low-contrast lines, while the line segments are accurately extracted by LSD afterwards.

We also add a filtering step, leveraging the DF and AF. We sample n_f points along each line, and compute the fraction of samples whose distance function is below τ_d and angle is close enough to the line orientation with tolerance τ_a . Only the segments with enough inliers are kept.

3.5. Line Segment Refinement with Optimization

To make lines even more accurate, we propose an optimization step to refine them by leveraging the predicted DF and AF. This refinement can also be used to enhance the lines of any other detector, and we show in Section 4.4 how it can make current deep detectors much more precise.

While lines are detected independently, they usually appear in highly structured configurations in the image. In particular, lines that are parallel in 3D will share vanishing points. We propose to integrate this as soft constraints into our refinement, effectively reducing the degrees of freedom.

We first compute a set of vanishing points (VPs) associated with the predicted line segments, using the multi-model fitting algorithm Progressive-X [7]. We use a strict inlier threshold to be sure to associate only relevant lines to a VP. The optimization is then performed independently for each line and is a weighted unconstrained least square minimization of three different costs:

$$C = \alpha C_A + \beta C_D + \gamma C_V; \quad (8)$$

Given a set P of n_{opt} points uniformly sampled along a line segment, we denote each point p_i , the orientation angle of the line as θ_l , and the VP associated with the line as v . We use the following three costs:

$$\begin{aligned} C_A &= \frac{1}{n_{opt}} \sum_{p_i \in P} 1 - \cos(\hat{A}(p_i) - \theta_l); \\ C_D &= \frac{1}{n_{opt}} \sum_{p_i \in P} \hat{D}(p_i); \quad C_V = d_{VP}(l; v_l); \end{aligned} \quad (9)$$

where d_{VP} is a distance measure between a line and a VP. We adopt the perpendicular distance of the line endpoints, projected onto the infinite line passing through the center of the line and the VP, as in [49]. These objectives are thus minimizing the difference between the sampled angle and the line orientation angle, minimizing the average distance field value over the line, and minimizing the distance between the line and its VP. In case the closest VP is farther away from the line than a threshold τ_p , we drop the VP constraint as it would push the line towards a wrong VP. To avoid lines drifting or collapsing to a single point, we keep the length of the line fixed, and we only optimize the lines

Figure 6. Line detection examples. Wireframe methods [20, 54] only detect structural lines, while DeepLSD offers more generic detections.

4.1. Evaluation on Low-Level Metrics

, and a translation of the middle point in the perpendicular direction of the line.

Since the VPs are already computed, we can even optimize the VPs as well, as a by-product of our approach. Jointly optimizing lines and VPs empirically led to inferior results, mainly because some lines require more refinement than others, so that a global refinement performs worse than independently optimizing the lines. We alternate, instead, between refining the lines and refining the VPs, for a fixed number of iterations. The VP refinement is performed through a least square minimization of the distance between the VP and all associated lines, and the line-VP association is recomputed after each iteration.

3.6. Implementation Details

We train two versions of our network, one indoors on the Wireframe dataset [19], but without using the ground truth lines, and one outdoors on MegaDepth [29]. Given the large size of MegaDepth, we keep 150 scenes for training and 17 for validation, and only sample 50 images from each scene. We use the Adam optimizer [23] and an initial learning rate of $1e^{-3}$, which is divided by 10 each time the validation loss reaches a plateau. The training takes roughly 12 hours on a single NVIDIA RTX 2080 GPU. For the line detection, we set the line region to 5 pixels and ignore magnitudes below 3 when applying LSD. We use $n_s = 50$ samples in the filtering step, $D_F = 1:5$, $\sigma = \frac{1}{9}$ and accept lines with more than 50% inliers. The parameters for VP estimation are tuned for each method on a validation set, but the usual threshold t_{VP} ranges from 1 to 2 pixels. The optimization weights are empirically chosen as $\alpha_s = 1$, $\alpha_A = 1$, and $\alpha_V = 0.2$. We adopt $n_{opt} = 10$ samples, perform a fixed set of $k = 5$ alternating iterations, and optimize with Ceres [3].

4. Experiments

To evaluate the performance of our method, we cannot use labeled lines as the existing ones are usually biased towards wireframes. We are more interested in evaluating the potential to use these lines for downstream applications, such as homography estimation, 3D line reconstruction, and visual localization. We also provide a visual comparison of various line detectors in Figure 6.

We first evaluate our line detection on two challenging datasets to test the robustness of the methods. First, the HPatches dataset [6], consisting of 580 pairs of images with ground truth homographies relating them and varying illumination and viewpoint changes. Second, the RDNIM dataset [35], also with image pairs related by a homography and with challenging day-night variations. We use the night reference in our experiments to get more challenging pairs.

Similarly as in [36], we assess the repeatability and localization error metrics. For both metrics, we compute a one-to-one matching of the detected line segments between the two images of a pair using the ground truth homography. For each match, one can then compute the distance between the line in the reference image and the line of the warped image reprojected into the reference frame. We consider two line distance measures: the structural distance evaluating the average distance between the endpoints, and the orthogonal distance measuring the average distance of each endpoint of one line to their orthogonal projection to the other line. Repeatability (Rep) measures the ratio of lines whose match has an error below 3 pixels, and the localization error (LE) returns the average distance of the 50 most accurate matches.

We also compute homography estimation scores similarly as in [11]. We first match line segments between the two images, using the Line Band Descriptor (LBD) [57]. To estimate the homography, we sample minimal sets of 4 line matches and run LO-RANSAC [27] for up to 1M iterations, using the orthogonal line distance as reprojection error.

We compare in Table 1 our method to two classical detectors: LSD [51] and ELSed [48]; the best methods using attraction fields: HAWP [54], its recent update HAWPv3 trained in a self-supervised way [55], and LSDNet [50]: a similar approach as ours combining LSD and a deep network; and two generic deep line detectors: TP-LSD [20] and SOLD2 [36]. We use the implementation of the authors with the biggest model available and default parameters, except for HAWP where we use a threshold of 0.9, as it was not detecting enough lines otherwise. HAWPv3 was trained on ImageNet. For LSD, we use the implementation of Rafael Grompone¹ instead of the OpenCV one as it gets much better results. Our method is given without the final optimization in the following, unless otherwise specified.

¹<http://www.ipol.im/pub/art/2012/gjmr-lsd/>

		Traditional		Learned				Hybrid		Traditional		Learned				Hybrid	
		LSD	ELSED	HAWP	HAWPv3	TP-LSD	SOLD2	LSDNet	DeepLSD	LSD	ELSED	HAWP	HAWPv3	TP-LSD	SOLD2	LSDNet	DeepLSD
		[51]	[48]	[54]	[55]	[20]	[36]	[50]	(Ours)	[51]	[48]	[54]	[55]	[20]	[36]	[50]	(Ours)
HPatches [6]	Struct Rep"	0.314	0.240	0.330	0.272	0.413	0.308	0.108	0.367	0.283	0.209	0.284	0.320	0.344	0.307	0.047	0.285
	LE #	<u>1.309</u>	1.551	2.019	2.132	1.500	1.741	2.860	1.235	2.039	2.303	2.206	1.939	<u>1.779</u>	1.879	3.331	1.733
	Orth Rep"	<u>0.468</u>	0.465	0.337	0.309	0.444	0.395	0.200	0.485	0.403	0.392	0.284	0.354	0.377	0.386	0.130	<u>0.394</u>
	LE #	0.793	0.845	1.905	1.937	1.305	1.362	2.285	<u>0.818</u>	1.369	<u>1.248</u>	2.215	1.704	1.625	1.449	2.752	1.098
H estimation"		<u>0.697</u>	0.617	0.260	0.231	0.388	0.421	0.316	0.705	<u>0.468</u>	0.200	0.006	0.026	0.030	0.182	0.027	0.591
# lines / img		492.6	425.4	53.6	82.0	88.6	122.9	172.1	486.2	191.4	112.0	31.6	23.8	24.1	138.2	109.1	400.0
Time [ms]#		104	10	61	51	179	334	48	271	34	3	42	47	75	199	44	96

Table 1. Line detection evaluation on the HPatches [6] and RDNIM [35] datasets. We compare repeatability (Rep) and localization error (LE) in structural and orthogonal distances, together with homography estimation. We get the best score on homography estimation and a good trade-off between classical and learned methods for the all metrics. The best score is in bold and the second best is underlined.

From the results, the learned methods, led by TP-LSD [20], offer good repeatability, but suffer from a low localization error and inaccurate homography estimation. Handcrafted methods and our method are much more accurate, due to the fact that they do not directly regress the endpoints, but gradually grow the line segments using very low-level details. DeepLSD displays the best improvement over LSD when the changes become the most challenging, i.e. on RDNIM with strong day-night changes. It can significantly improve the localization error and homography estimation score. In spite of having a similar approach as ours, LSDNet [50] performs poorly for multiple reasons: they lose accuracy by rescaling images to a fixed low resolution, their line mask is less precise than our distance field, and their training is limited to the Wireframe dataset, while ours can be trained on more diverse images. Overall, our method offers the best trade-off between handcrafted and learned methods and consistently ranks first in the downstream task of homography estimation.

4.2. 3D Line Reconstruction

The aim of this work is to provide general-purpose lines and as such, the lines generated by DeepLSD should be suitable for 3D reconstruction. We leverage Line3D++ [17] that takes a collection of images with known poses and the associated 2D line segments, and outputs a 3D reconstruction of lines. We propose to compare our method with a few baselines on the first 4 scenes of the Hypersim dataset [39]. This synthetic - but highly realistic - dataset has the advantage of offering a ground truth mesh and 3D model, making it suitable for a quantitative evaluation. Given the ground truth mesh of the scene, we can compute the recall and precision of the 3D lines. Recall is the length in meters of all the portions of lines that are within 5 millimeters from the mesh. High values mean that many lines have been reconstructed. Precision is the percentage of predicted lines that are within 5 millimeters from the mesh. High values indicate that most of the predicted lines are on a real 3D surface.

The results can be seen in Table 2. DeepLSD obtains the best recall overall, and second best precision. While the query camera poses. We report the median translation

	ai.001.001		ai.001.002		ai.001.003		ai.001.004		Average	
	R	P	R	P	R	P	R	P	R	P
LSD [51]	183.6	95.8	61.8	95.3	85.0	88.9	225.3	91.5	213.9	92.9
SOLD2 [36]	109.9	94.7	89.3	92.8	62.0	89.0	58.6	89.1	80.0	91.4
HAWPv3 [55]	15.8	79.9	15.6	81.0	24.4	68.4	18.5	77.3	18.6	76.7
TP-LSD [20]	68.8	95.3	38.9	94.7	50.7	98.2	102.7	94.3	65.3	95.6
DeepLSD	204.8	96.5	89.5	98.1	137.8	88.0	231.1	91.9	226.1	93.6

Table 2. Line 3D reconstruction evaluation. We reconstruct lines in 3D with Line3D++ [17] and evaluate the line length recall in m (R) and precision (P) on the first 4 scenes of Hypersim [39].

TP-LSD [20] ranks first in precision, it is able to recover very few lines, as shows its average recall, which is 71% smaller than the one of DeepLSD. We provide qualitative examples of the reconstructions in the supp. material. Note that DeepLSD is able to reconstruct more lines and with a higher precision than LSD [51], the detector that is the most commonly used for line reconstruction [17].

4.3. Visual Localization

The 7Scenes dataset [47] is a well-known RGB-D dataset for visual localization, displaying 7 indoor scenes with GT poses and depth. While most scenes are already saturated for point-based localization, the Stairs scene remains very challenging for feature points. Due to the lack of texture and repeated patterns of the stairs, current point-based methods are still struggling on this scene [8]. We thus propose to evaluate our method and previous works on this particular scene, by following the pipeline of hloc [42, 43], enriched with line features. As points remain important features, we still detect SuperPoint features [11] and match them with SuperGlue [44]. We detect lines with different detectors, and match them between database and query images with the SOLD2 descriptor [36]. Since depth is available on 7Scenes, we can directly back-project lines in 3D and do not rely on line mapping. In practice, we sample points along each line, un-project them to 3D, and re-fit a line in 3D to these unprojected points. We use the solvers of [24, 26, 60] to generate poses from a minimal set of 3 features (3 points, 2 points and 1 line, 1 point and 2 lines, or 3 lines), then combine them in a

	T / R err#	Acc "
Point-only	4.7 / 1.25	53.4
LSD [51]	3.4 / 0.94	73.2
SOLD2 [36]	3.5 / 0.96	71.5
HAWPv3 [55]	3.4 / 0.93	72.1
TP-LSD [20]	3.4 / 0.98	74.2
DeepLSD	3.1 / 0.85	76.6

Figure 7. Visual localization on 7Scenes stairs [47]. We evaluate the median translation and rotation errors (cm / deg), the pose accuracy at a 5 cm / 5 deg threshold, and plot the pose accuracy curve for various thresholds.

		Struct		Orth		H estim	# lines / img	Time [ms]
		Rep"	LE #	Rep"	LE #			
HAWP [54]	Baseline	0.253	1.34	0.253	1.43	0.701	40	
	Opt w/o VP	0.300	1.293	0.399	1.067	0.864	95.2	142
	Opt w/ VP	0.318	1.245	0.431	0.967	0.892		300
TP-LSD [20]	Baseline	0.273	1.379	0.342	1.269	0.658	46	
	Opt w/o VP	0.314	1.326	0.470	0.949	0.898	90.8	145
	Opt w/ VP	0.331	1.277	0.512	0.861	0.913		297
SOLD2 [36]	Baseline	0.197	1.277	0.333	0.894	0.848	297	
	Opt w/o VP	0.172	1.388	0.339	0.814	0.935	166.7	426
	Opt w/ VP	0.185	1.330	0.368	0.753	0.920		697
DeepLSD (Ours)	Baseline	0.318	0.941	0.489	0.574	0.991	68	
	Opt w/o VP	0.314	0.938	0.482	0.575	0.994	168.8	154
	Opt w/ VP	0.319	0.927	0.501	0.544	0.981		542

Table 3. Line re-nement on the Wireframe dataset [19]. We use an error threshold of 1 pixel for the repeatability metrics. The re-nement can significantly improve the localization error and homography score of inaccurate methods.

and rotation error, as well as the percentage of successfully recovered poses under various thresholds.

Figure 7 shows that DeepLSD obtains the best performance on this challenging dataset. One can highlight the large boost of performance brought by line features compared to using points only. Lines are indeed still present and well localized in indoor environments such as in this scene, and can be matched even when in low-textured scenes.

4.4. Impact of the Line Re-nement

We evaluate applying our proposed line re-nement as a post-processing step for several learned detection methods. Classical detectors are usually already accurate enough, so that our re-nement would not enhance them much. For each method, we compare the raw lines with the lines and VPs optimized by our line optimization. Table 3 shows results of line detectors on the 462 images of the test set of the Wireframe dataset [19]. The second image is obtained using a synthetic homographic warp of the first image. We use the Wireframe dataset as it has a lot of well-defined vanishing points, which can be leveraged during the optimization. We include results for our proposed optimization with and without the VP constraint to show the increased accuracy with VPs. As we want to highlight the gain in accuracy, we compute repeatability with an error threshold of only 1 pixel.

	Struct		Orth		H estim	# lines / img
	Rep"	LE #	Rep"	LE #		
Single edge	0.241	2.121	0.328	1.686	0.434	130.8
No DF normalization	0.344	1.343	0.475	0.879	0.674	439.6
HAWP with our lines	0.209	2.138	0.239	1.840	0.245	98.0
DeepLSD (Ours)	0.367	1.235	0.485	0.818	0.705	486.2

Table 4. Ablation study on the HPatches dataset [6]. We compare DeepLSD to alternatives detecting single edges, without DF normalization and with HAWP re-trained on our line GT.

Results show that the re-nement can significantly improve all metrics evaluating the accuracy of the lines, i.e. the localization error and homography estimation. This is particularly true for HAWP [54] and TP-LSD [20], with a decrease in localization error with orthogonal distance of up to 32% for both, and an improvement of homography score of 27% and 39%. The benefits brought by the re-nement are lower for our method, as its raw predicted lines are already sub-pixel accurate and the optimization is limited by the resolution of the DF and AF. Nonetheless, it can slightly improve most metrics. A limitation of this re-nement is the execution time, which grows linearly with the number of lines, and requires running two networks.

4.5. Ablation Study

We validate our design choices on the HPatches dataset [6] with low-level detector metrics. We compare our proposed approach with the same model detecting single edges instead of double ones, our network trained without the DF normalization, and a version of the HAWP [54] backbone re-trained on our line GT on the MegaDepth dataset [29]. The results of Table 4 emphasize the importance of each component. Note that re-training HAWP [54] on our lines yields poor results due to the higher number of lines compared to wireframe lines, and the fact that generic lines have often noisy endpoints, so that predicting an angle to the two endpoints is noisy as well.

5. Conclusion

We presented a hybrid line segment detector combining the robustness of deep learning and the accuracy of handcrafted detectors, using a learned surrogate image gradient as intermediate representation. Without the requirement of ground truth lines, our method can be trained on any dataset and is suitable for most tasks including line segments. Finally, we proposed a line re-nement able to improve the accuracy of our method and to bridge the gap in line localization between deep line detectors and handcrafted ones. We believe that our general-purpose lines will open new possibilities to use line segments in the wild.

Acknowledgments. We would like to warmly thank Iago Suarez for reviewing this paper and for the insightful discussions, as well as Yifan Yu for sharing his code for visual localization.

References

- [1] Hichem Abdellali, Robert Frohlich, Viktor Vilagos, and Zoltan Kato. L2D2: Learnable line detector and descriptor. In International Conference on 3D Vision (3DV), 2021. 1, 2
- [2] Abien Fred Agarap. Deep learning using rectified linear units (ReLU). In arXiv, 2018. 4
- [3] Sameer Agarwal and Keir Mierle. Ceres solver. <http://ceres-solver.org>. 6
- [4] C. Akinlar and C. Topal. EDLines: Real-time line segment detection by edge drawing. International Conference on Image Processing (ICIP), 2011. 2
- [5] Simon Baker, Daniel Scharstein, J. P. Lewis, Stefan Roth, Michael J. Black, and Richard Szeliski. A database and evaluation methodology for optical flow. International Conference on Computer Vision (ICCV), 2007. 4
- [6] Vassileios Balntas, Karel Lenc, Andrea Vedaldi, and Krystian Mikolajczyk. Hpatches: A benchmark and evaluation of handcrafted and learned local descriptors. Computer Vision and Pattern Recognition (CVPR), 2017. 6, 7, 8
- [7] Daniel Barath and Jiri Matas. Progressive-X: Efficient, any-time, multi-model fitting algorithm. In International Conference on Computer Vision (ICCV), 2019. 5
- [8] Eric Brachmann and Carsten Rother. Visual camera re-localization from rgb and rgb-d images using deep learning. IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI), 44, 2022. 7
- [9] Federico Camposeco, Andrea Cohen, Marc Pollefeys, and Torsten Sattler. Hybrid Camera Pose Estimation. Computer Vision and Pattern Recognition (CVPR), 2018. 7
- [10] Xili Dai, Xiaojun Yuan, Haigang Gong, and Yi Ma. Fully convolutional line parsing. In arXiv, 2021. 1, 2
- [11] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperPoint: Self-supervised interest point detection and description. In Computer Vision and Pattern Recognition Workshops (CVPRW), 2018. 2, 4, 6, 7
- [12] James H. Elder, Emilio J. Aldama, Yiming Qian, and Ron Tal. MCMLSD: A probabilistic algorithm and evaluation framework for line segment detection. In arXiv, 2020. 2
- [13] Qiang Fu, Jialong Wang, Hongshan Yu, Islam Ali, Feng Guo, Yijia He, and Hong Zhang. PL-VINS: Real-time monocular visual-inertial SLAM with point and line features. In arXiv, 2020. 1
- [14] Shuang Gao, Jixiang Wan, Yishan Ping, Xudong Zhang, Shuzhou Dong, Jijun Li, and Yandong Guo. Pose refinement with joint optimization of visual points and lines. In arXiv, 2021. 1, 2
- [15] Ruben Gomez-Ojeda, Francisco-Angel Moreno, David Zuñiga-Noël, Davide Scaramuzza, and Javier Gonzalez-Jimenez. PL-SLAM: A stereo SLAM system through the combination of points and line segments. IEEE Transactions on Robotics, 35, 2019. 1
- [16] Geonmo Gu, Byungsoo Ko, SeoungHyun Go, Sung-Hyun Lee, Jingeun Lee, and Minchul Shin. Towards real-time and light-weight line segment detection. Conference on Artificial Intelligence (AAAI), 2022. 1, 2
- [17] Manuel Hofer, Michael Maurer, and Horst Bischof. Efficient 3d scene abstraction using line segment detection. Computer Vision and Image Understanding (CVIU), 157, 2017. 1, 7
- [18] Paul VC Hough. Method and means for recognizing complex patterns, 1962. US Patent 3,069,654. 2
- [19] Kun Huang, Yifan Wang, Zihan Zhou, Tianjiao Ding, Shenghua Gao, and Yi Ma. Learning to parse wireframes in images of man-made environments. Computer Vision and Pattern Recognition (CVPR), 2018. 1, 2, 4, 6, 8
- [20] Siyu Huang, Fangbo Qin, Pengfei Xiong, Ning Ding, Yijia He, and Xiao Liu. TP-LSD: Tri-points based line segment detector. In European Conference on Computer Vision (ECCV), 2020. 1, 2, 6, 7, 8
- [21] Zhaoyang Huang, Han Zhou, Yijin Li, Bangbang Yang, Yan Xu, Xiaowei Zhou, Hujun Bao, Guofeng Zhang, and Hongsheng Li. VS-Net: Voting with segmentation for visual localization. In Computer Vision and Pattern Recognition (CVPR), 2021. 3
- [22] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariance shift. In International Conference on Machine Learning (ICML), 2015. 4
- [23] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. International Conference on Learning Representations (ICLR), 2014. 6
- [24] Zuzana Kukelova, Jan Heller, and Andrew Fitzgibbon. Efficient intersection of three quadrics and applications in computer vision. In Computer Vision and Pattern Recognition (CVPR), 2016. 7
- [25] Manuel Lange, Claudio Raisch, and Andreas Schilling. LVO: Line only stereo visual odometry. In 2019 International Conference on Indoor Positioning and Indoor Navigation (IPIN), 2019. 1
- [26] Viktor Larsson. PoseLib - Minimal Solvers for Camera Pose Estimation, 2020. 7
- [27] Karel Lebeda, Jiri Matas, and Ondrej Chum. Fixing the Locally Optimized RANSAC. In British Machine Vision Conference (BMVC), 2012. 6
- [28] Hao Li, Huai Yu, Jinwang Wang, Wen Yang, Lei Yu, and Sebastian Scherer. ULSD: Unified line segment detection across pinhole, fisheye, and spherical cameras. International Journal of Photogrammetry and Remote Sensing (ISPRS), 2021. 1
- [29] Zhengqi Li and Noah Snavely. MegaDepth: Learning single-view depth prediction from internet photos. Computer Vision and Pattern Recognition (CVPR), 2018. 6, 8
- [30] Yancong Lin, Silvia L Pintea, and Jan C van Gemert. Deep hough-transform line priors. In European Conference on Computer Vision (ECCV), 2020. 1, 2
- [31] Philipp Lindenberger, Paul-Edouard Sarlin, Viktor Larsson, and Marc Pollefeys. Pixel-Perfect Structure-from-Motion with Featuremetric Refinement. International Conference on Computer Vision (ICCV), 2021. 2, 3
- [32] André Mateus, Omar Tahri, A. Pedro Aguiar, Pedro U. Lima, and Pedro Miraldo. On incremental structure from motion using lines. IEEE Transactions on Robotics, 38, 2022. 1
- [33] Quan Meng, Jiakai Zhang, Qiang Hu, Xuming He, and Jingyi Yu. LGNN: A context-aware line segment detector. ACM International Conference on Multimedia (MM), 2020. 2

- [34] Branislav Micusik and Horst Wildenauer. Structure from motion with line segments under relaxed endpoint constraints. *International Journal of Computer Vision (IJCV)* 24, 2017. 1, 2
- [35] Rémi Pautrat, Viktor Larsson, Martin R. Oswald, and Marc Pollefeys. Online invariance selection for local feature descriptors. In *European Conference on Computer Vision (ECCV)* 2020. 6, 7
- [36] Rémi Pautrat, Juan-Ting Lin, Viktor Larsson, Martin R. Oswald, and Marc Pollefeys. SOLD2: Self-supervised occlusion-aware line description and detection. *Computer Vision and Pattern Recognition (CVPR)* 2021. 1, 2, 4, 6, 7, 8
- [37] Albert Pumarola, Alexander Vakhitov, Antonio Agudo, Alberto Sanfeliu, and Francese Moreno-Noguer. PL-SLAM: Real-time monocular visual SLAM with points and lines. In *International Conference on Robotics and Automation (ICRA)* 2017. 1
- [38] Meixiang Quan, Zheng Chai, and Xiao Liu. LOF: Structure-aware line tracking based on optical flow. *arXiv*, 2021. 1
- [39] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M. Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. *International Conference on Computer Vision (ICCV)* 2021. 7
- [40] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)* 2015. 4
- [41] Yohann Salin, Renaud Marlet, and Pascal Monasse. Multiscale line segment detector for robust and accurate SfM. In *International Conference on Pattern Recognition (ICPR)* 2016. 2
- [42] Paul-Edouard Sarlin. Visual localization made easy with hloc. <https://github.com/cvg/Hierarchical-Localization/>. 7
- [43] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. *Computer Vision and Pattern Recognition (CVPR)* 2019. 7
- [44] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperGlue: Learning feature matching with graph neural networks. *Computer Vision and Pattern Recognition (CVPR)* June 2020. 7
- [45] Paul-Edouard Sarlin, Ajaykumar Unagar, Viktor Larsson, Hugo Germain, Carl Toft, Victor Larsson, Marc Pollefeys, Vincent Lepetit, Lars Hammarstrand, Fredrik Kahl, and Torsten Sattler. Back to the Feature: Learning Robust Camera Localization from Pixels to Pose. *Computer Vision and Pattern Recognition (CVPR)* 2021. 2, 3
- [46] Torsten Sattler et al. RansacLib - A Template-based *SAC Implementation, 2019. 7
- [47] Jamie Shotton, Ben Glocker, Christopher Zach, Shahram Izadi, Antonio Criminisi, and Andrew Fitzgibbon. Scene coordinate regression forests for camera relocalization in rgb-d images. In *Computer Vision and Pattern Recognition (CVPR)* 2013. 7, 8
- [48] Iago Suárez, Joé M. Buenaposada, and Luis Baumela. ELSEd: Enhanced line segment drawing. *Pattern Recognition*, 2022. 2, 6, 7
- [49] Jean-Philippe Tardif. Non-iterative approach for fast and accurate vanishing point detection. *International Conference on Computer Vision (ICCV)* 2009. 1, 5
- [50] Lev Teplyakov, Leonid Erlygin, and Evgeny Shvets. Lsdnet: Trainable modification of Lsd algorithm for real-time line segment detection. *IEEE Access* 10, 2022. 1, 3, 6, 7
- [51] Rafael Grompone Von Gioi, Jeremie Jakubowicz, Jean-Michel Morel, and Gregory Randall. LSD: A fast line segment detector with a false detection control. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)* 32(4):722–732, 2008. 1, 2, 3, 4, 5, 6, 7, 8
- [52] Yifan Xu, Weijian Xu, David Cheung, and Zhuowen Tu. Line segment detection using transformers without edge maps. *Computer Vision and Pattern Recognition (CVPR)* 2021. 2
- [53] Nan Xue, Song Bai, Fudong Wang, Gui-Song Xia, Tianfu Wu, and Liangpei Zhang. Learning attraction field representation for robust line segment detection. *Computer Vision and Pattern Recognition (CVPR)* 2019. 1, 2, 3, 4
- [54] Nan Xue, Tianfu Wu, Song Bai, Fudong Wang, Gui-Song Xia, Liangpei Zhang, and Philip HS Torr. Holistically-attracted wireframe parsing. In *Computer Vision and Pattern Recognition (CVPR)* 2020. 1, 2, 3, 4, 5, 6, 7, 8
- [55] Nan Xue, Tianfu Wu, Song Bai, Fu-Dong Wang, Gui-Song Xia, Liangpei Zhang, and Philip H.S. Torr. Holistically-attracted wireframe parsing: From supervised to self-supervised learning. *arXiv*, 2022. 2, 6, 7, 8
- [56] Haotian Zhang, Yicheng Luo, Fangbo Qin, Yijia He, and Xiao Liu. Elsd: Efficient line segment detector and descriptor. In *International Conference on Computer Vision (ICCV)* 2021. 1, 2
- [57] Lilian Zhang and Reinhard Koch. An efficient and robust line segment matching approach based on lbd descriptor and pairwise geometric consistency. *Journal of Visual Communication and Image Representation* 24, 2013. 6
- [58] Yongjun Zhang, Dong Wei, and Yansheng Li. AG3line: Active grouping and geometry-gradient combined validation for fast line segment extraction. *Pattern Recognition* 113, 2021. 2
- [59] 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 *Computer Vision and Pattern Recognition (CVPR)* 2019. 2
- [60] Lipu Zhou, Jiamin Ye, and Michael Kaess. A stable algebraic camera pose estimation for minimal configurations of 2d/3d point and line correspondences. *Asian Conference on Computer Vision (ACCV)* 2018. 7
- [61] Yichao Zhou, Haozhi Qi, and Yi Ma. End-to-end wireframe parsing. In *International Conference on Computer Vision (ICCV)*, 2019. 1, 2
- [62] Xingxing Zuo, Xiaojia Xie, Yong Liu, and Guoquan Huang. Robust visual SLAM with point and line features. *International Conference on Intelligent Robots and Systems (IROS)* 2017. 1