



HAL
open science

Polygonal Building Segmentation by Frame Field Learning

Nicolas Girard, Dmitriy Smirnov, Justin Solomon, Yuliya Tarabalka

► **To cite this version:**

Nicolas Girard, Dmitriy Smirnov, Justin Solomon, Yuliya Tarabalka. Polygonal Building Segmentation by Frame Field Learning. CVPR 2021 - IEEE Conference on Computer Vision and Pattern Recognition, Jun 2021, Pittsburg / Virtual, United States. hal-02548545v2

HAL Id: hal-02548545

<https://inria.hal.science/hal-02548545v2>

Submitted on 31 Mar 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Polygonal Building Extraction by Frame Field Learning

Nicolas Girard¹ Dmitriy Smirnov² Justin Solomon² Yuliya Tarabalka³

¹Université Côte d’Azur, Inria ²Massachusetts Institute of Technology ³LuxCarta Technology

Abstract

While state of the art image segmentation models typically output segmentations in raster format, applications in geographic information systems often require vector polygons. To help bridge the gap between deep network output and the format used in downstream tasks, we add a frame field output to a deep segmentation model for extracting buildings from remote sensing images. We train a deep neural network that aligns a predicted frame field to ground truth contours. This additional objective improves segmentation quality by leveraging multi-task learning and provides structural information that later facilitates polygonization; we also introduce a polygonization algorithm that utilizes the frame field along with the raster segmentation. Our code is available at <https://github.com/Lydorn/Polygonization-by-Frame-Field-Learning>.

1. Introduction

Due to their success in processing large collections of noisy images, deep convolutional neural networks (CNNs) have achieved state-of-the-art in remote sensing segmentation. Geographic information systems like Open Street Map (OSM) [30], however, require segmentation data in *vector* format (e.g., polygons and curves) rather than raster format, which is generated by segmentation networks. Additionally, methods that extract objects from remote sensing images require especially high throughput to handle the volume of high-resolution aerial images captured daily over large territories of land. Thus, modifications to the conventional CNN pipeline are necessary.

Existing work on deep building segmentation generally falls into one of two general categories. The first vectorizes the probability map produced by a network

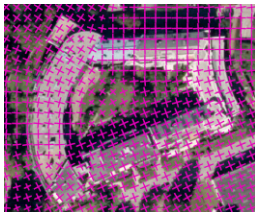


Figure 1: A frame field output by our network.

a posteriori, e.g., by using contour detection (marching squares [25]) followed by polygon simplification (Ramer–Douglas–Peucker [32, 13]). Such approaches suffer when the classification maps contain imperfections such as smoothed out corners, a common artifact of conventional deep segmentation methods. Moreover, as we show in Fig. 2, even perfect probability maps are challenging to polygonize due to shape information being lost from the discretization of the raster output. To improve the final polygons, these methods employ expensive and complex post-processing procedures. ASIP polygonization [20] uses polygonal partition refinement to approximate shapes from the output probability map based on a tunable parameter controlling the trade-off between complexity and fidelity. In [45], a decoder and a discriminator regularize output probability maps adversarially. This requires computing large matrices of pairwise discontinuity costs between pixels and involves adversarial training, which is less stable than conventional supervised learning.

Another category of deep segmentation methods learns a vector representation directly. For example, Curve-GCN [23] trains a graph convolutional network (GCN) to deform polygons iteratively, and PolyMapper [21] uses a recurrent neural network (RNN) to predict vertices one at a time. While these approaches directly predict polygon parameters, GCNs and RNNs suffer from several disadvantages. Not only are they more difficult to train than CNNs, but also their output topology is restricted to simple polygons without holes—a serious limitation in segmenting complex buildings. Additionally, adjoining buildings with common walls are common, especially in city centers. Curve-GCN and PolyMapper are unable to reuse the same polyline in adjoining buildings, yielding overlaps and gaps.

We introduce a building extraction algorithm that avoids the challenges above by adding a frame field output to a fully-convolutional network (see Fig. 1). While this has imperceptible effect on training or inference time, the frame field not only increases segmentation performance, e.g., yielding sharper corners, but also provides useful information for vectorization. Additional losses learn a valid frame field that is consistent with the segmentation. These losses regularize the segmentation, similar to [39], which includes

MRF/CRF regularization terms in the loss function to avoid extra MRF/CRF inference steps.

The frame field allows us to devise a straightforward polygonization method extending the Active Contours Model (ACM, or “snakes”) [19], which we call the Active Skeleton Model (ASM). Rather than fitting contours to image data, ASM fits a skeleton graph, where each edge connects two junction nodes with a chain of vertices (i.e., a polyline). This allows us to reuse shared walls between adjoining buildings. To our knowledge, no existing method handles this case ([41] shows results with common walls but does not provide details). Our method naturally handles large buildings and buildings with inner holes, unlike end-to-end learning methods like PolyMapper [21]. Lastly, our polygon extraction pipeline is highly GPU-parallelizable, making it faster than more complex methods.

Our main contributions are:

- (i) a learned frame field aligned to object tangents, which improves segmentation via multi-task learning;
- (ii) coupling losses between outputs for self-consistency, further leveraging multi-task learning; and
- (iii) a fast polygonization method leveraging the frame field, allowing complexity tuning of a corner-aware simplification step and handling non-trivial topology.

2. Related work

ASIP polygonization [20] inputs an RGB image and a probability map of objects (e.g., buildings) detected in the image (e.g., by a neural network). Then, starting from a polygonal partition that oversegments the image into convex cells, the algorithm refines the partition while labeling its cells by semantic class. The refinement process is an optimization with terms that balance fidelity to the input against complexity of the output polygons. The configuration space is explored by splitting and merging the polygonal cells. As the fidelity and complexity terms can be balanced with a coefficient, the fidelity-to-complexity ratio can be tuned. However, there does not exist a systematic approach for interpreting or determining this coefficient. While ASIP post-processes the output of a deep learning method, recent approaches aim for an end-to-end pipeline.

CNNs are successful at converting grid-based input to grid-based output for tasks where each output pixel depends on its local neighborhood in the input. In this setting, it is straightforward and efficient to train a network for supervised prediction of segmentation probability maps. The paragraphs below, however, detail major challenges when using such an approach to extract polygonal buildings.

First, the model needs to produce variable-sized outputs to capture varying numbers of objects, contours, and vertices. This requires complex architectures like recurrent neural networks (RNNs) [18], which are not as efficiently

trained as CNNs and need multiple iterations at inference time. Such is the case for PolyMapper [21], Polygon-RNN [5], and Polygon-RNN++ [1]. Curve-GCN [23] predicts a fixed number of vertices simultaneously.

A second challenge is that the model must make discrete decisions of whether to add a contour, whether to add a hole to an object, and with how many vertices to describe a contour. Adding a contour is solved by object detection: a contour is predicted for each detected object. Adding holes to an object is more challenging, but a few methods detect holes and predict their contours. One model, BSP-Net [8], circumvents this issue by combining predicted convex shapes for the final output, producing shapes in a compact format, with potential holes inside. To our knowledge, the number of vertices is not a variable that current deep learning models can optimize for; discrete decisions are difficult to pose differentially without training techniques such as the straight-through estimator [3] or reinforcement learning [37, 28, 27].

A third challenge is that, unlike probability maps, the output structure of polygonal building extraction is not grid-like. Within the network, the grid-like structure of the image input has to be transformed to a more general planar graph structure representing building outlines. City centers have the additional problem of adjoining buildings that share a wall. Ideally, the output geometry for such a case would be a collection of polygons, one for each individual building, which share polylines corresponding to common walls. Currently, no existing deep learning method tackles this case. Our method solves it but is not end-to-end. PolyMapper [21] tackles the individual building and road network extraction tasks. As road networks are graphs, they propose a novel sequentialization method to reformulate graph structures as closed polygons. Their approach might work in the case of adjoining buildings with common walls. Their output structure, however, is less adapted to GPU computation, making it less efficient. RNNs such as PolyMapper [21], Polygon-RNN [5], and Polygon-RNN++ [1] perform beam search at inference to prune off improbable sequences, which requires more vertex predictions than are used in the final output and is inefficient. The DefGrid [14] module is a non-RNN approach where the network processes polygonal superpixels. It is more complex than our simple fully-convolutional network and is still subject to the rounded corner problem.

The challenges above demand a middle ground between learning a bitmap segmentation followed by a hand-crafted polygonization method and end-to-end methods, aiming to be easily-deployable, topologically flexible w.r.t. holes and common walls, and efficient. A step in this direction is the machine-learned building polygonization [44] that predicts building segmentations using a CNN, uses a generative adversarial network to regularize building boundaries, and

learns a building corner probability map, from which vertices are extracted. In contrast, our model predicts a frame field both as additional geometric information (instead of a building corner probability map) and as a way to regularize building boundaries (instead of adversarial training). The addition of this frame field output is similar in spirit to DiResNet [12], a road extraction neural network that outputs road direction in addition to road segmentation, first introduced in [2]. The orientation is learned for each road pixel by a cross-entropy classification loss whose labels are orientation bins. This additional geometric feature learned by the network improves the overall geometric integrity of the extracted objects (in their case road connectivity). The differences to our method include the following: (1) our frame fields encode two orientations instead of one (needed for corners), (2) we use a regression loss instead of a classification loss, and (3) we use coupling losses to promote coherence between segmentation and frame field.

3. Method

Our key idea is to help the polygonization method solve ambiguous cases caused by discrete probability maps by asking the neural network to output missing shape information in the form of a frame field (see Fig. 2). This practically does not increase training and inference time, allows for simpler and faster polygonization, and regularizes the segmentation—solving the problem of small misalignments of the ground truth annotations that yield rounded corners if no regularization is used.

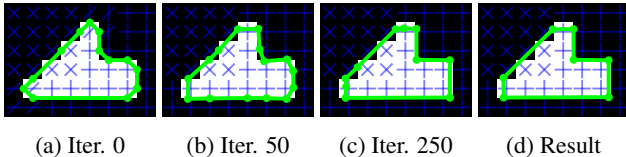


Figure 2: Even a perfect classification map can yield incorrect polygonization due to a locally ambiguous probability map, as shown in (a), the output of marching squares. Our polygonization method iteratively optimizes the contour (b-d) to align to a frame field, yielding better results as our frame field (blue) disambiguates between slanted walls and corners, preventing corners from being cut off.

3.1. Frame fields

We provide the necessary background on frame fields, a key part of our method. Following [42, 11], a frame field is a *4-PolyVector field*, which assigns four vectors to each point of the plane. In the case of a frame field, however, the first two vectors are constrained to be opposite to the other two, i.e., each point is assigned a set of vectors $\{u, -u, v, -v\}$. At each point in the image, we consider the two directions

that define the frame as two complex numbers $u, v \in \mathbb{C}$. We need two directions (rather than only one) because buildings, unlike organic shapes, are regular structures with sharp corners, and capturing directionality at these sharp corners requires two directions. To encode the directions in a way that is agnostic to relabeling and sign change, we represent them as coefficients of the following polynomial:

$$f(z) = (z^2 - u^2)(z^2 - v^2) = z^4 + c_2 z^2 + c_0. \quad (1)$$

We denote (1) above by $f(z; c_0, c_2)$. Given a (c_0, c_2) pair, we can easily recover one pair of directions defining the corresponding frame:

$$\begin{cases} c_0 = u^2 v^2 \\ c_2 = -(u^2 + v^2) \end{cases} \iff \begin{cases} u^2 = -\frac{1}{2} \left(c_2 + \sqrt{c_2^2 - 4c_0} \right) \\ v^2 = -\frac{1}{2} \left(c_2 - \sqrt{c_2^2 - 4c_0} \right). \end{cases} \quad (2)$$

In our approach, inspired by [4], we learn a smooth frame field with the property that, along building edges, at least one field direction is aligned to the polygon tangent direction. At polygon corners, the field aligns to *both* tangent directions, motivating our use of PolyVector fields rather than vector fields. Away from polygon boundaries, the frame field does not have any alignment constraints but is encouraged to be smooth and not collapse to a line field. Like [4], we formulate the field computation variationally, but, unlike their approach, we use a neural network to learn the field at each pixel, which is also explored in [38]. To avoid sign and ordering ambiguity, we learn a (c_0, c_2) pair per pixel rather than (u, v) .

3.2. Frame field learning

We describe our method, illustrated in Fig. 3. Our network takes a $3 \times H \times W$ image I as input and outputs a pixel-wise classification map and a frame field. The classification map contains two channels, \hat{y}_{int} corresponding to building interiors and \hat{y}_{edge} to building boundaries. The frame field contains four channels corresponding to the two complex coefficients $\hat{c}_0, \hat{c}_2 \in \mathbb{C}$, as in §3.1 above.

Segmentation losses. Our method can be used with any deep segmentation model as a backbone; in our experiments, we use the U-Net [33] and DeepLabV3 [7] architectures. The backbone outputs an F -dimensional feature map $\hat{y}_{backbone} \in \mathbb{R}^{F \times H \times W}$. For the segmentation task, we append to the backbone a fully-convolutional block (taking $\hat{y}_{backbone}$ as input) consisting of a 3×3 convolutional layer, a batch normalization layer, an ELU nonlinearity, another 3×3 convolution, and a sigmoid nonlinearity. This segmentation head outputs a segmentation map $\hat{y}_{seg} \in \mathbb{R}^{2 \times H \times W}$. The first channel contains the object interior segmentation map \hat{y}_{int} and the second contains the contour segmentation map \hat{y}_{edge} . Our training is supervised—each input image

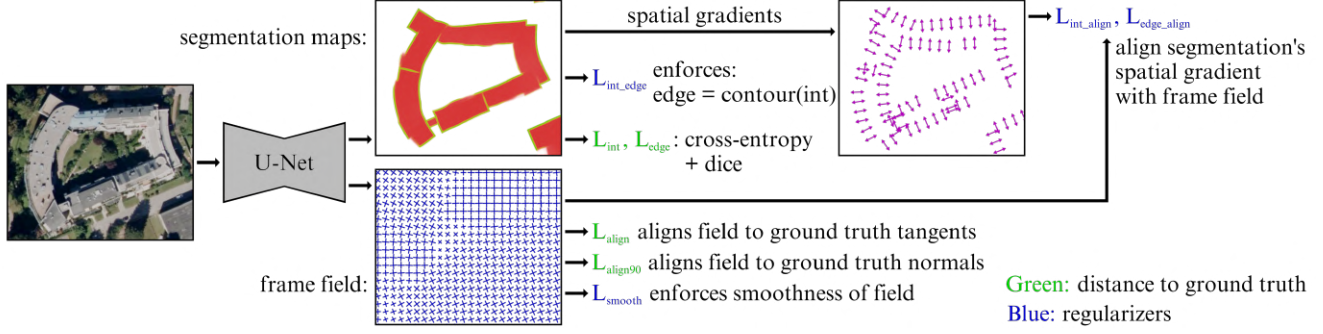


Figure 3: Given an overhead image, our model outputs an edge mask, interior mask, and frame field. The loss aligns the masks and field to ground truth data, enforces smoothness of the frame field, and ensures consistency between the outputs.

is labeled with ground truth y_{int} and y_{edge} , corresponding to rasterized polygon interiors and edges, respectively. We then use a linear combination of the cross-entropy loss and Dice loss [36] for loss L_{int} applied on the interior output as well as loss L_{edge} applied on the contour (edge) output.

Frame field losses. In addition to the segmentation masks, our network outputs a frame field. We append another head to the backbone via a fully-convolutional block consisting of a 3×3 convolutional layer, a batch normalization layer, an ELU nonlinearity, another 3×3 convolution, and a tanh nonlinearity. This frame field block inputs the concatenation of the output features of the backbone and the segmentation output: $[\hat{y}_{backbone}, \hat{y}_{seg}] \in \mathbb{R}^{(F+2) \times H \times W}$. It outputs the frame field with $\hat{c}_0, \hat{c}_2 \in \mathbb{C}^{H \times W}$. The corresponding ground truth label is an angle $\theta_\tau \in [0, \pi)$ of the unsigned tangent vector of the polygon contour. We use three losses to train the frame field:

$$L_{align} = \frac{1}{HW} \sum_{x \in I} y_{edge}(x) |f(e^{i\theta_\tau}; \hat{c}_0(x), \hat{c}_2(x))|^2, \quad (3)$$

$$L_{align90} = \frac{1}{HW} \sum_{x \in I} y_{edge}(x) |f(e^{i\theta_\tau^\perp}; \hat{c}_0(x), \hat{c}_2(x))|^2, \quad (4)$$

$$L_{smooth} = \frac{1}{HW} \sum_{x \in I} (\|\nabla \hat{c}_0(x)\|^2 + \|\nabla \hat{c}_2(x)\|^2), \quad (5)$$

where θ_w is the direction of w ($w = \|w\|_2 e^{i\theta_w}$), and $\tau^\perp = \tau - \frac{\pi}{2}$. Each loss measures a different property of the field:

- L_{align} enforces alignment of the frame field to the tangent directions. This term is small when the polynomial $f(\cdot; \hat{c}_0, \hat{c}_2)$ has a root near $e^{i\theta_\tau}$, implicitly implying that one of the field directions $\{\pm u, \pm v\}$ is aligned with the tangent direction τ . Since (1) has no odd-degree terms, this term has no dependence on the sign of τ , as desired.
- $L_{align90}$ prevents the frame field from collapsing to a line field by encouraging it to also align with τ^\perp .
- L_{smooth} is a Dirichlet energy measuring the smoothness of $\hat{c}_0(x)$ and $\hat{c}_2(x)$ as functions of location x in the image. Smoothly-varying \hat{c}_0 and \hat{c}_2 yield a smooth frame field.

Output coupling losses. We add coupling losses to ensure mutual consistency between our network outputs:

$$L_{int_align} = \frac{1}{HW} \sum_{x \in I} f(\nabla \hat{y}_{int}(x); \hat{c}_0(x), \hat{c}_2(x))^2, \quad (6)$$

$$L_{edge_align} = \frac{1}{HW} \sum_{x \in I} f(\nabla \hat{y}_{edge}(x); \hat{c}_0(x), \hat{c}_2(x))^2, \quad (7)$$

$$L_{int_edge} = \frac{1}{HW} \sum_{x \in I} \max(1 - \hat{y}_{int}(x), \|\nabla \hat{y}_{int}(x)\|_2) \cdot \|\nabla \hat{y}_{int}(x)\|_2 - \hat{y}_{edge}(x). \quad (8)$$

- L_{int_align} aligns the spatial gradient of the predicted interior map \hat{y}_{int} with the frame field (analogous to (3)).
- L_{edge_align} aligns the spatial gradient of the predicted edge map \hat{y}_{edge} with the frame field (analogous to (3)).
- L_{int_edge} makes the predicted edge map be equal to the norm of the spatial gradient of the predicted interior map. This loss is applied outside of buildings (hence the $1 - \hat{y}_{int}(x)$ term) and along building contours (hence the $\|\nabla \hat{y}_{int}(x)\|_2$ term) and is not applied inside buildings, so that common walls between adjoining buildings can still be detected by the edge map.

Final loss. Because the losses (L_{int} , L_{edge} , L_{align} , $L_{align90}$, L_{smooth} , L_{int_align} , L_{edge_align} , and L_{int_edge}) have distinct units, we compute a normalization coefficient for each loss by averaging its value over a random subset of the training dataset using a randomly-initialized network. Losses are then normalized by this coefficient before being linearly combined. This normalization aims to rescale losses such that they are easier to balance. More details are in the supplementary materials.

3.3. Frame field polygonization

The main steps of our polygonization method are shown in Fig. 4. It is inspired by the Active Contour Model (ACM) [19]. ACM is initialized with a given contour and

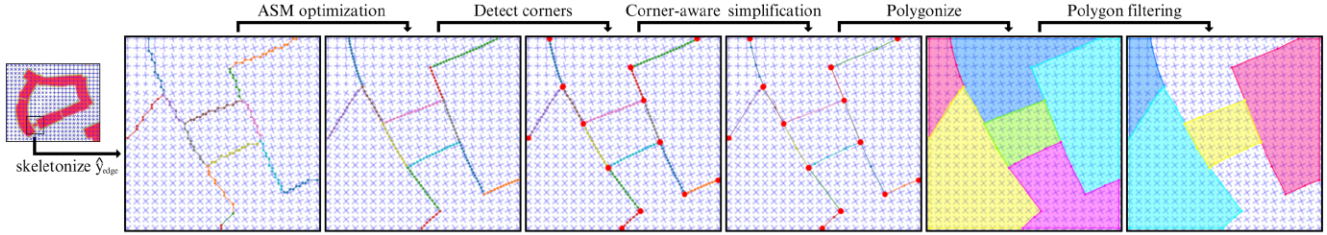


Figure 4: Overview of our post-processing polygonization algorithm. Given an interior classification map and frame field (Fig. 3) as input, we optimize the contour to align to the frame field using an Active Skeleton Model (ASM) and detect corners using the frame field, simplifying non-corner vertices.

minimizes an energy function $E_{contour}^*$, which moves the contour points toward an optimal position. Usually this energy is composed of a term to fit the contour to the image and additional terms to limit the amount of stretch and/or curvature. The optimization is performed by gradient descent. Overall the ACM lends itself perfectly for parallelized execution on the GPU, and the optimization can be performed using an automatic differentiation module included in deep learning frameworks. We adapt ACM so that the optimization is performed on a skeleton graph instead of contours, giving us the Active Skeleton Model (ASM). We call the skeleton graph the graph of connected pixels of the skeleton image obtained by the thinning method [43] applied on the building wall probability map y_{edge} . The following energy terms are used:

- $E_{probability}$ fits the skeleton paths to the contour of the building interior probability map $y_{int}(v)$ at a certain probability threshold ℓ (set to 0.5 in practice).
- $E_{frame\ field\ align}$ aligns each edge of the skeleton graph to the frame field.
- E_{length} ensures that the node distribution along paths remains homogeneous as well as tight.

Details about our data structure (designed for GPU computation), definition and computation of our energy terms, and explanation of our corner-aware simplification step can be found in the supplementary materials.

4. Experimental setup

4.1. Datasets

Our method requires ground truth polygonal building annotations (rather than raster binary masks) so that the ground truth angle for the frame field can be computed by rasterizing separately each polygon edge and taking the edge’s angle. Thus, for each pixel we get a θ_τ value, which is used in L_{align} .

We perform experiments on these datasets (more details in the supplementary material):

- CrowdAI Mapping Challenge dataset [34] (*CrowdAI dataset*): 341438 aerial images of size 300×300 pixels

with associated ground truth polygonal annotations.

- Inria Aerial Image Labeling dataset [26] (*Inria dataset*): 360 aerial images of size 5000×5000 pixels. Ten cities are represented, making it more varied than the *CrowdAI dataset*. However, the ground truth is in the form of raster binary masks. We thus create the *Inria OSM dataset* by taking OSM polygon annotations and correcting their misalignment using [15]. We also create the *Inria Polygonized dataset* by converting the original ground truth binary masks to polygon annotations with our polygonization method (see supplementary materials).
- *Private dataset*: 57 satellite images for training with sizes varying from 2000×2000 pixels to 20000×20000 pixels, captured over 30 different cities from all continents with three different types of satellites. This is our most varied and challenging dataset. However, the building outline polygons were manually labeled precisely by an expert, ensuring the best possible ground truth. Results for this private dataset are in the supplementary material.

4.2. Backbones

The first backbone we use is *U-Net16*, a small U-Net [33] with 16 starting hidden features (instead of 64 in the original). We also use *DeepLab101*, a DeepLabV3 [7] model that utilizes a ResNet-101 [17] encoder. Our best performing model is *UResNet101*—a U-Net with a ResNet-101 [17] encoder (pre-trained on ImageNet [10]). We observed that the pre-trained ResNet-101 encoder achieves better final performance than random initialization. For the *UResNet101*, we additionally use distance weighting for the cross-entropy loss, as done for the original U-Net [33].

4.3. Ablation study and additional experiments

We perform an ablation study to validate various components of our method (results in Tables 1, 2, and 4):

- “No field” removes the frame field output for comparison to pure segmentation. Only interior segmentation L_{int} , edge segmentation L_{edge} and interior/edge coupling $L_{int\ edge}$ losses remain.
- “Simple poly.” uses a baseline polygonization algo-

rithm (marching-squares contour detection followed by the Ramer–Douglas–Peucker simplification) on the interior classification map learned by our full method. This allows us to study the improvement of our polygonization method from leveraging the frame field.

Additional experiments in the supplementary material include: “no coupling losses” removes all coupling losses ($L_{int\ align}$, $L_{edge\ align}$, $L_{int\ edge}$) to determine whether enforcing consistency between outputs has an impact; “no $L_{align90}$,” “no $L_{int\ edge}$,” “no $L_{int\ align}$ and $L_{edge\ align}$,” and “no L_{smooth} ” all remove the specified losses; “complexity vs. fidelity” varies the simplification tolerance parameter ϵ to demonstrate the trade-off between complexity and fidelity of our corner-aware simplification procedure.

4.4. Metrics

The standard metric for image segmentation is Intersection over Union (IoU), which is then used to compute other metrics such as MS COCO [22], Average Precision (AP), and Average Recall (AR)—along with variants AP_{50} , AP_{75} , AR_{50} , AR_{75} . Since we aim to produce clean geometry, it is important to measure contour regularity, not captured by the area-based metrics IoU, AP, and AR. Moreover, as annotations are bound to have some alignment noise, only optimizing IoU will favor blurry segmentations with rounded corners over sharp segmentations, as the blurry ones correspond to the shape expectation of the noisy ground truth annotation; segmentation results with sharp corners may even yield a lower IoU than segmentations with rounded corners. We thus introduce the *max tangent angle error* metric that compares the tangent angles between predicted polygons and ground truth annotations, penalizing contours not aligned with the ground truth. It is computed by uniformly sampling points along a predicted contour, computing the angle of the tangent for each point, and comparing it to the tangent angle of the closest point on the ground truth contour. The *max tangent angle error* is the maximum tangent angle error over all sampled points. More details about the computation of these metrics can be found in the supplementary material.

5. Results and discussion

5.1. CrowdAI dataset

We visualize our polygon extraction results for the *CrowdAI dataset* and compare them to other methods in Fig. 5. The ASIP polygonization method [20] inputs the probability maps of a U-Net variant [9] that won the CrowdAI challenge. All methods perform well on common building types, e.g., houses and residential buildings, but we can see that results of ASIP are less regular than PolyMapper and ours. For more complex building shapes (e.g., not rectangular or with a hole inside), ASIP outputs reasonable re-

Method	Mean <i>max tangent angle errors</i> ↓
UResNet101 (no field), simple poly.	51.9°
UResNet101 (with field), simple poly.	45.1°
U-Net variant [9], ASIP poly. [20]	44.0°
UResNet101 (with field), ASIP poly. [20]	38.3°
U-Net variant [9], UResNet101 our poly.	36.6°
PolyMapper [21]	33.1°
UResNet101 (with field), our poly.	31.9°

Table 1: Mean *max tangent angle errors* over all the original validation polygons of the *CrowdAI dataset* [34].

sults, albeit still not very regular. However, the PolyMapper approach of object detection followed by polygonal outline regression does not work in the most difficult cases. It does not support nontrivial topology by construction, but also, it struggles with large complex buildings. We hypothesize that PolyMapper suffers from the fact that there are not many complex buildings and does not generalize as well as fully-convolutional networks.

We report results on the original validation set of the *CrowdAI dataset* for the *max tangent angle error* in Table 1 and MS COCO metrics in Table 2. “(with field)” refers to models trained with our full frame field learning method, “(no field)” refers to models trained without any frame field output, “mask” refers to the output raster segmentation mask of the network, “our poly.” refers to our frame field polygonization method, and “simple poly.” refers to the baseline polygonization of marching squares followed by Ramer-Douglas-Peucker simplification. We also applied our polygonization method to the same probability maps used by the ASIP polygonization method (U-Net variant [9]) for fair comparison of polygonization methods.

In Table 1, “simple poly.” performs better using “(with field)” segmentation compared to “(no field)” because of a regularization effect from frame field learning. PolyMapper performs significantly better than “simple poly.” even though it is not explicitly regularized. Our frame field learning and polygonization method is necessary to decrease the error further and compare favorably to PolyMapper.

In Table 2, our UResNet101 (with field) outperforms most previous works, except “U-Net variant [9], ASIP poly. [20]” due to the U-Net variant being the winning entry to the challenge. However our polygonization applied after that same U-Net variant achieves better max tangent angle error and AP than ASIP but worse AR. The same is true when applying ASIP to our UResNet101 (with field): it has slightly worse AP, AR, and max tangent angle error. However, the ASIP method also results in better max tangent angle error when using our UResNet101 (with field) compared to using the U-Net variant.

Runtimes. We compare runtimes in Table S4. ASIP does not have a GPU implementation. In their paper they give

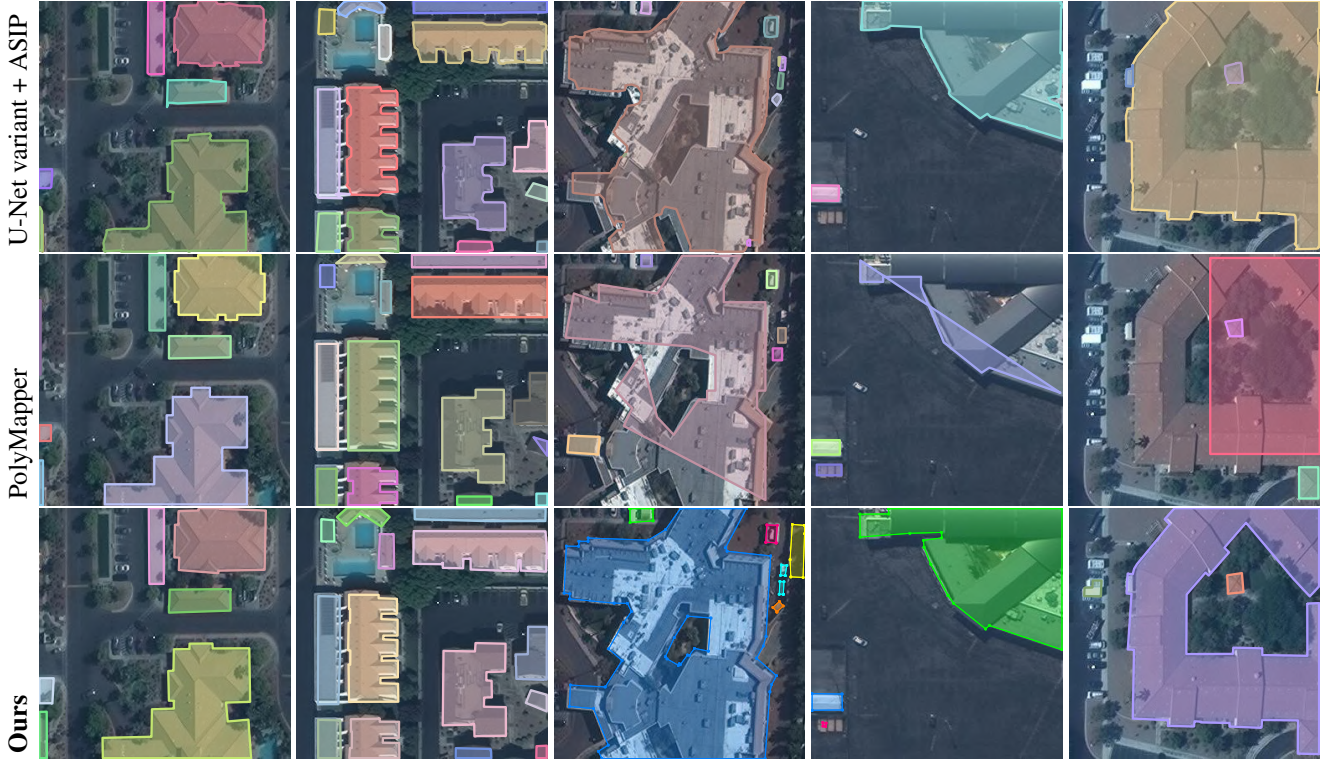


Figure 5: Example building extraction results on CrowdAI test images. Buildings become more complex from left to right. (top) U-Net variant [9] + ASIP [20], (middle) PolyMapper [21], and (bottom) **ours**: UResNet101 (full), frame field polygonization.

Method	AP ↑	AP ₅₀ ↑	AP ₇₅ ↑	AR ↑	AR ₅₀ ↑	AR ₇₅ ↑
UResNet101 (no field), mask	62.4	86.7	72.7	67.5	90.5	77.4
UResNet101 (no field), simple poly.	61.1	87.4	71.2	64.7	89.4	74.1
UResNet101 (with field), mask	64.5	89.3	74.6	68.1	91.0	77.7
UResNet101 (with field), simple poly.	61.7	87.7	71.5	65.4	89.9	74.6
UResNet101 (with field), our poly.	61.3	87.5	70.6	65.0	89.4	73.9
UResNet101 (with field), ASIP poly. [20]	60.0	86.3	69.9	64.0	88.8	73.4
U-Net variant [9], UResNet101 our poly.	67.0	92.1	75.6	73.2	93.5	81.1
Mask R-CNN [16] [35]	41.9	67.5	48.8	47.6	70.8	55.5
PANet [24]	50.7	73.9	62.6	54.4	74.5	65.2
PolyMapper [21]	55.7	86.0	65.1	62.1	88.6	71.4
U-Net variant [9], ASIP poly. [20]	65.8	87.6	73.4	78.7	94.3	86.1

Table 2: AP and AR results on the *CrowdAI dataset* [34] for all polygonization experiments.

Method	Time (sec) ↓	Hardware
PolyMapper [21]	0.38	GTX 1080Ti
ASIP [20]	0.15	Laptop CPU
Ours	0.04	GTX 1080Ti

Table 3: Average times to extract buildings from a 300×300 pixel patch. **Ours** refers to UResNet101 (with field), our poly. ASIP’s time does not include model inference.

an average runtime of 1-3s on CPU with ~10% CPU utilization. Assuming perfect parallelization, they estimate their average runtime to be 0.15s with 100% CPU utiliza-

tion. Their method uses a priority queue for optimizing the polygonal partitioning with various geometric operators and is harder to implement on GPU. Our efficient data structure makes our building extraction competitive with prior work.

5.2. Inria OSM dataset

U-Net16 (no field), simple poly. U-Net16 (with field), **our** poly.



Figure 6: Small crop of *Inria dataset* results.

Method	mIoU \uparrow	Mean <i>max tangent angle errors</i> \downarrow
Eugene Khvedchenya ² , simple poly.	80.7%	52.2°
ICTNet [6], simple poly.	80.1%	52.1°
UResNet101 (no field), simple poly.	73.2%	52.0°
Zorzi et al. [44] poly.	74.4%	34.5°
UResNet101 (with field), our poly.	74.8%	28.1°

Table 4: IoU and mean *max tangent angle errors* for polygon extraction methods on the *Inria polygonized dataset*.

The *Inria OSM dataset* is more challenging than the *CrowdAI dataset* because it contains more varied areas (e.g., countryside, city center, residential, and commercial) with different building types. It also contains adjacent buildings with common walls, which our edge segmentation output can detect. The mean IoU on test images of the output classification maps is 78.0% for the U-Net16 trained with a frame field compared to 76.9% for the U-Net16 with no frame field. The IoU does not significantly penalize irregular contours, but, by visually inspecting segmentation outputs as in Fig. 6, we can see the effect of the regularization. Our method successfully handles complex building shapes which can be very large, with blocks of buildings featuring common walls and holes. See the supplementary materials for more results.

5.3. Inria polygonized dataset

Eugene Khvedchenya², simple poly. UResNet101 (with field), **our poly.**



Figure 7: Crop of an *Inria polygonized dataset* test image.

The *Inria polygonized dataset* with its associated challenge¹ allows us to directly compare to other methods trained on the same ground truth, even though it does not consider learning of separate buildings. In Table 4, our

¹<https://project.inria.fr/aerialimagelabeling/leaderboard/>

method matches [44] in terms of mIoU, with lower *max tangent angle error*. The two top methods on the leaderboard (ICTNet [6] and “Eugene Khvedchenya”) achieve a mIoU over 80%, but they lack contour regularity with high *max tangent angle error*; they also only output segmentation masks, needing a *posteriori* polygonization to extract polygonal buildings. Fig. 7 shows the cleaner geometry of our method. The ground truth of the *Inria polygonized dataset* has misalignment noise, yielding imprecise corners that produce rounded corners in the prediction if no regularization is applied. See the supplementary materials for more results.

6. Conclusion

We improve on the task of building extraction by learning an additional output to a standard segmentation model: a frame field. This motivates the use of a regularization loss, leading to more regular contours, e.g., with sharp corners. Our approach is efficient since the model is a single fully-convolutional network. The training is straightforward, unlike adversarial training, direct shape regression, and recurrent networks, which require significant tuning and more computational power. The frame field adds virtually no cost to inference time, and it disambiguates tough polygonization cases, making our polygonization method less complex. Our data structure for the polygonization makes it parallelizable on the GPU. We handle the case of holes in buildings as well as common walls between adjoining buildings. Because of the skeleton graph structure, common wall polylines are naturally guaranteed to be shared by the buildings on either side. As future work, we could apply our method to any image segmentation network, including multi-class segmentation, where the frame field could be shared between all classes.

7. Acknowledgements

Thanks to ANR for funding the project EPITOME ANR-17-CE23-0009 and to Inria Sophia Antipolis - Méditerranée “Nef” computation cluster for providing resources and support. The MIT Geometric Data Processing group acknowledges the generous support of Army Research Office grant W911NF2010168, of Air Force Office of Scientific Research award FA9550-19-1-031, of National Science Foundation grant IIS-1838071, from the CSAIL Systems that Learn program, from the MIT-IBM Watson AI Laboratory, from the Toyota-CSAIL Joint Research Center, from a gift from Adobe Systems, from an MIT.nano Immersion Lab/NCSoft Gaming Program seed grant, and from the Skoltech-MIT Next Generation Program. This work was also supported by the National Science Foundation Graduate Research Fellowship under Grant No. 1122374.

References

- [1] David Acuna, Huan Ling, Amlan Kar, and Sanja Fidler. Efficient interactive annotation of segmentation datasets with polygon-rnn++. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [2] Anil Batra, Suriya Singh, Guan Pang, Saikat Basu, C.V. Jawahar, and Manohar Paluri. Improved road connectivity by joint learning of orientation and segmentation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 3
- [3] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. 2013. arXiv:1308.3432. 2
- [4] Mikhail Bessmeltsev and Justin Solomon. Vectorization of line drawings via polyvector fields. *ACM Trans. Graph.*, 2019. 3
- [5] Lluís Castrejón, Kaustav Kundu, Raquel Urtasun, and Sanja Fidler. Annotating object instances with a polygon-rnn. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [6] Bodhiswatta Chatterjee and Charalambos Poullis. On building classification from remote sensor imagery using deep neural networks and the relation between classification and reconstruction accuracy using border localization as proxy. In *2019 16th Conference on Computer and Robot Vision (CRV)*, pages 41–48, 2019. 8, 26
- [7] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. 2017. arXiv:1706.05587. 3, 5
- [8] Zhiqin Chen, Andrea Tagliasacchi, and Hao Zhang. Bspnet: Generating compact meshes via binary space partitioning. *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [9] Jakub Czakon, Kamil A. Kaczmarek, Andrzej Pyskir, and Piotr Tarasiewicz. Best practices for elegant experimentation in data science projects. *EuroPython*, 2018. 6, 7, 18, 20
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. 5
- [11] Olga Diamanti, Amir Vaxman, Daniele Panozzo, and Olga Sorkine-Hornung. Designing N -polyvector fields with complex polynomials. *Eurographics SGP*, 2014. 3
- [12] Lei Ding and Lorenzo Bruzzone. Diresnet: Direction-aware residual network for road extraction in vhr remote sensing images. 2020. arXiv:2005.07232. 3
- [13] David H. Douglas and Thomas K. Peucker. Algorithms for the reduction of the number of points required to represent a digitized line or its caricature. *The Canadian Cartographer*, 10(2):112–122, 1973. 1, 16
- [14] Jun Gao, Zian Wang, Jinchun Xuan, and Sanja Fidler. Beyond fixed grid: Learning geometric image representation with a deformable grid. In *ECCV*, 2020. 2
- [15] Nicolas Girard, Guillaume Charpiat, and Yuliya Tarabalka. Noisy Supervision for Correcting Misaligned Cadaster Maps Without Perfect Ground Truth Data. In *IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, 2019. 5, 17
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 7, 20
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. 2015. arXiv:1512.03385. 5
- [18] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9:1735–80, 12 1997. 2
- [19] Michael Kass, Andrew Witkin, and Demetri Terzopoulos. Snakes: Active contour models. *International Journal of Computer Vision (IJCV)*, 1(4):321–331, 1988. 2, 4
- [20] Muxingzi Li, Florent Lafarge, and Renaud Marlet. Approximating shapes in images with low-complexity polygons. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1, 2, 6, 7, 17, 18, 20
- [21] Zuoyue Li, Jan Dirk Wegner, and Aurélien Lucchi. Topological map extraction from overhead images. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 1, 2, 6, 7, 17, 18, 20
- [22] 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 David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *The European Conference on Computer Vision (ECCV)*, pages 740–755, Cham, 2014. Springer International Publishing. 6, 17
- [23] Huan Ling, Jun Gao, Amlan Kar, Wenzheng Chen, and Sanja Fidler. Fast interactive object annotation with curve-gcn. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2
- [24] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 7, 20
- [25] William Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. 1987. 1, 12
- [26] Emmanuel Maggiori, Yuliya Tarabalka, Guillaume Charpiat, and Pierre Alliez. Can semantic labeling methods generalize to any city? the inria aerial image labeling benchmark. In *IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, 2017. 5, 17
- [27] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. 2013. arXiv:1312.5602. 2
- [28] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, Feb. 2015. 2
- [29] Juan Nunez-Iglesias. skan: skeleton analysis in python. <https://github.com/jni/skan>, 2017. 13

- [30] OpenStreetMap contributors. Planet dump retrieved from <https://planet.osm.org>, 2017. **1, 17**
- [31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. **12**
- [32] Urs Ramer. An iterative procedure for the polygonal approximation of plane curves. *Computer graphics and image processing*, 1(3):244–256, 1972. **1, 16**
- [33] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. 2015. arXiv:1505.04597. **3, 5**
- [34] Sharada Prasanna Mohanty. Crowdai dataset. https://www.crowdai.org/challenges/mapping-challenge/dataset_files, 2018. **5, 6, 7, 16, 18, 20**
- [35] Sharada Prasanna Mohanty. Crowdai mapping challenge 2018: Baseline with mask rcnn. <https://github.com/crowdai/crowdai-mapping-challenge-mask-rcnn>, 2018. **7, 20**
- [36] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In *Deep learning in medical image analysis and multimodal learning for clinical decision support*, pages 240–248. Springer, 2017. **4**
- [37] Richard S. Sutton and Andrew G. Barto. Reinforcement learning: An introduction, 2020. **2**
- [38] Maria Taktasheva, Albert Matveev, Alexey Artemov, and Evgeny Burnaev. Learning to approximate directional fields defined over 2d planes. In Wil M. P. van der Aalst, Vladimir Batagelj, Dmitry I. Ignatov, Michael Khachay, Valentina Kuskova, Andrey Kutuzov, Sergei O. Kuznetsov, Irina A. Lomazova, Natalia Loukachevitch, Amedeo Napoli, Panos M. Pardalos, Marcello Pelillo, Andrey V. Savchenko, and Elena Tutubalina, editors, *Analysis of Images, Social Networks and Texts*, pages 367–374, Cham, 2019. Springer International Publishing. **3**
- [39] Meng Tang, Federico Perazzi, Abdelaziz Djelouah, Ismail Ben Ayed, Christopher Schroers, and Yuri Boykov. On regularized losses for weakly-supervised cnn segmentation. In *The European Conference on Computer Vision (ECCV)*, 2018. **1**
- [40] T. Tieleman and G. Hinton. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012. **15**
- [41] S. Tripodi, L. Duan, F. Trastour, V. Poux, Lionel Laurore, and Yuliya Tarabalka. Automated chain for large-scale 3d reconstruction of urban scenes from satellite images. *ISPRS - International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, XLII-2/W16:243–250, 09 2019. **2**
- [42] Amir Vaxman, Marcel Campen, Olga Diamanti, Daniele Panozzo, David Bommes, Klaus Hildebrandt, and Mirela Ben-Chen. Directional field synthesis, design and processing. *Computer Graphics Forum*, 35:545–572, 05 2016. **3**
- [43] T. Y. Zhang and C. Y. Suen. A fast parallel algorithm for thinning digital patterns. *Commun. ACM*, 27(3):236–239, Mar. 1984. **5, 13**
- [44] Stefano Zorzi, Ksenia Bittner, and Friedrich Fraundorfer. Machine-learned regularization and polygonization of building segmentation masks, 2020. arXiv:2007.12587. **2, 8, 18, 20, 26**
- [45] Stefano Zorzi and Friedrich Fraundorfer. Regularization of building boundaries in satellite images using adversarial and regularized losses. In *IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, 2019. **1**

Supplementary Materials

Contents

1. Frame field learning details	11
1.1. Model architecture	11
1.2. Losses	11
1.3. Handling numerous heterogeneous losses	11
1.4. Training details	11
2. Frame field polygonization details	12
2.1. Data structure	12
2.2. Active Skeleton Model	14
2.3. Corner-aware polygon simplification	16
2.4. Detecting building polygons in planar graph	16
3. Experimental setup details	16
3.1. Datasets	16
3.2. Metrics	17
4. Additional results	18
4.1. CrowdAI dataset	18
4.2. Inria OSM dataset	20
4.3. Inria polygonized dataset	20
4.4. Private dataset	21

1. Frame field learning details

1.1. Model architecture

We show in Fig. S1 how we add a frame field output to an image segmentation backbone. The backbone can be any (possibly pretrained) network as long as it outputs an F -dimensional feature map $\hat{y}_{backbone} \in \mathbb{R}^{F \times H \times W}$.

1.2. Losses

We define image segmentation loss functions below:

$$L_{BCE}(y, \hat{y}) = \frac{1}{HW} \sum_{x \in I} y(x) \cdot \log(\hat{y}(x)) + (1 - y(x)) \cdot \log(1 - \hat{y}(x)), \quad (9)$$

$$L_{Dice}(y, \hat{y}) = 1 - 2 \cdot \frac{|y \cdot \hat{y}| + 1}{|y + \hat{y}| + 1}, \quad (10)$$

$$L_{int} = \alpha \cdot L_{BCE}(y_{int}, \hat{y}_{int}) + (1 - \alpha) \cdot L_{Dice}(y_{int}, \hat{y}_{int}), \quad (11)$$

$$L_{edge} = \alpha \cdot L_{BCE}(y_{edge}, \hat{y}_{edge}) + (1 - \alpha) \cdot L_{Dice}(y_{edge}, \hat{y}_{edge}), \quad (12)$$

where $0 < \alpha < 1$ is a hyperparameter. In practice, $\alpha = 0.25$ gives good results.

For the frame field, we show a visualization of the L_{align} loss in Fig. S2.

1.3. Handling numerous heterogeneous losses

We linearly combine our eight losses using eight coefficients, which can be challenging to balance. Because the losses have different units, we first compute a normalization coefficient $N_{\langle loss\ name \rangle}$ by computing the average of each loss on a random subset of the training dataset using a randomly-initialized network. Then each loss can be normalized by this coefficient. The total loss is a linear combination of all normalized losses:

$$\lambda_{int} \frac{L_{int}}{N_{int}} + \lambda_{edge} \frac{L_{edge}}{N_{edge}} + \lambda_{align} \frac{L_{align}}{N_{align}} + \lambda_{align90} \frac{L_{align90}}{N_{align90}} + \lambda_{smooth} \frac{L_{smooth}}{N_{smooth}} + \lambda_{int\ align} \frac{L_{int\ align}}{N_{int\ align}} + \lambda_{edge\ align} \frac{L_{edge\ align}}{N_{edge\ align}} + \lambda_{int\ edge} \frac{L_{int\ edge}}{N_{int\ edge}}, \quad (13)$$

where the $\lambda_{\langle loss\ name \rangle}$ coefficients are to be tuned. It is also possible to separately group the main losses and the regularization losses and have a single λ coefficient balancing the two loss groups:

$$\lambda L_{main} + (1 - \lambda) L_{regularization}, \quad (14)$$

with

$$L_{main} = \frac{L_{int}}{N_{int}} + \frac{L_{edge}}{N_{edge}} + \frac{L_{align}}{N_{align}}, \quad (15)$$

$$L_{regularization} = \frac{L_{align90}}{N_{align90}} + \frac{L_{smooth}}{N_{smooth}} + \frac{L_{int\ align}}{N_{int\ align}} + \frac{L_{edge\ align}}{N_{edge\ align}} + \frac{L_{int\ edge}}{N_{int\ edge}}. \quad (16)$$

In practice we started experiments with the single-coefficient version with $\lambda = 0.75$ and then used the multi-coefficient version to have more control by setting $\lambda_{int} = \lambda_{edge} = 10$, $\lambda_{align} = 1$, $\lambda_{align90} = 0.2$, $\lambda_{smooth} = 0.005$, $\lambda_{int\ align} = \lambda_{edge\ align} = \lambda_{int\ edge} = 0.2$.

1.4. Training details

We do not heavily tune our hyperparameters: once we find a value that works based on validation performance we keep it across ablation experiments. We employ early stopping for the U-Net16 and DeepLabV3 models (25 and 15 epochs, respectively) chosen by first training the full method on the training set of the *CrowdAI dataset*, choosing the epoch number of the lowest validation loss, and finally re-training the model on the train and validation sets for that number of total epochs.

Segmentation losses L_{int} and L_{edge} are both a combination of 25% cross-entropy loss and 75% Dice loss. To balance the losses in ablation experiments, we used the single-coefficient version with $\lambda = 0.75$. For our best performing model UResNet101 we used the multi-coefficients version to have more control by setting $\lambda_{int} = \lambda_{edge} = 10$,

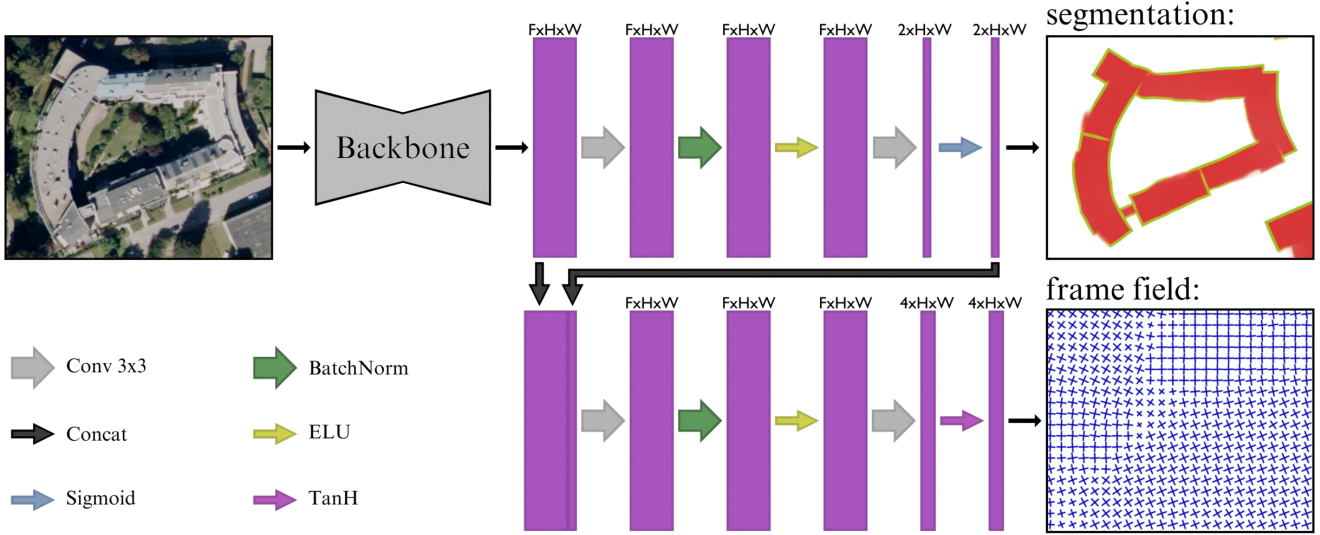


Figure S1: Details of our network’s architecture with the addition of the frame field output.

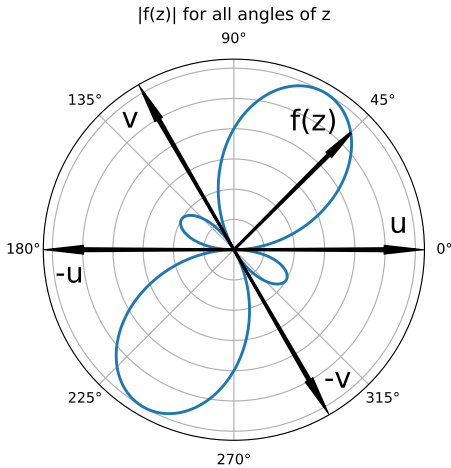


Figure S2: Visualization of the frame field align loss L_{align} (in blue) for a certain configuration of $\{-u, u, -v, v\}$ and all possible ground truth $z = e^{i\theta r}$ directions.

$\lambda_{align} = 1$, $\lambda_{align90} = 0.2$, $\lambda_{smooth} = 0.005$, $\lambda_{int align} = \lambda_{edge edge} = \lambda_{int edge} = 0.2$. The U-Net16 was trained on 4 GTX 1080Ti GPUs in parallel on 512×512 patches and a batch size of 16 per GPU (effective batch size 64). For all training runs, we compute for each loss its normalization coefficient $N_{(loss_name)}$ on 1000 batches before optimizing the network.

Our method is implemented in PyTorch [31]. On the *CrowdAI dataset*, training takes 2 hours per epoch on 4 1080Ti GPUs for the U-Net16 model and 3.5 hours per epoch for the DeepLabV3 backbone on 4 2080Ti GPUs.

Inference with the U-Net16 on a 5000×5000 image (requires splitting into 1024×1024 patches) takes 7 seconds on a Quadro M2200 (laptop GPU).

2. Frame field polygonization details

We expand here on the algorithm and implementation details of our frame field polygonization method.

2.1. Data structure

Our polygonization method needs to be initialized with geometry, which is then optimized to align to the frame field (among other objectives we will present later).

In the case of extracting individual buildings, we use the marching squares [25] contour finding algorithm on the predicted interior probability map y_{int} with an isovalue ℓ (set to 0.5 in practice). The result is a collection of contours $\{\mathcal{C}_i\}$ where each contour is a sequence of 2D points:

$$\mathcal{C}_i = ((r_0, c_0), (r_1, c_1), \dots, (r_{n_i-1}, c_{n_i-1})).$$

where $r_i, c_i \in \mathbb{R}$ correspond to vertex i ’s position along the row axis and the column axis respectively (they are not restricted to being integers). A contour is generally closed with $(r_0, c_0) = (r_{n_i-1}, c_{n_i-1})$, but it can be open if the corresponding object touches the border of the image (therefore start and end vertices are not the same).

In the case of extracting buildings with potential adjoining buildings sharing a common wall, we extract the skeleton graph of the predicted edge probability map y_{edge} . This skeleton graph is a hyper-graph made of nodes connected together by chains of vertices (i.e., polylines) called paths (see Fig. S4 for examples). To obtain this skeleton

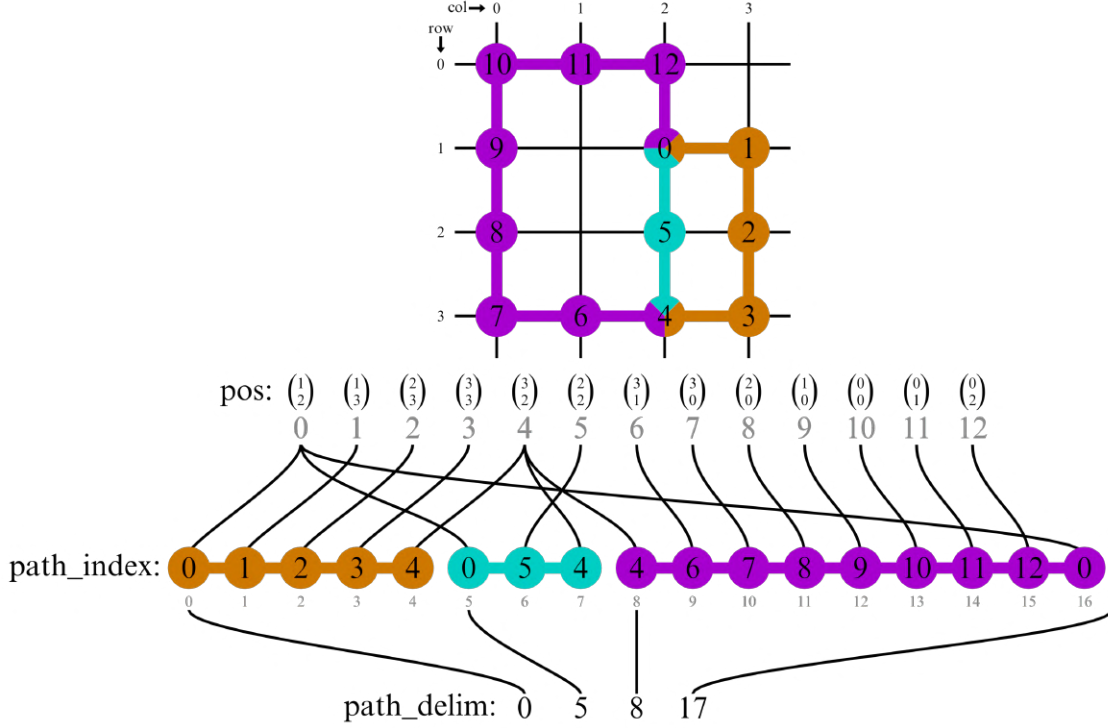


Figure S3: Our data structure of an example skeleton graph. It represents two buildings with a shared wall, necessitating 3 polyline paths. Here nodes 0 and 4 are shared among paths and are thus repeated in *path_index*. We can see *path_index* is a concatenation of the node indices in *pos* of the paths. Finally, *path_delim* is used to store the separation indices in *path_index* of those concatenated paths. Indices of arrays are in gray.

graph, we first compute the skeleton image using the thinning method [43] on the binary edge mask (computed by thresholding y_{edge} with $\ell = 0.5$). It reduces binary objects to a one-pixel-wide representation. We then use the Skan [29] Python library to convert this representation to a graph representation connecting those pixels. The resulting graph is a collection of paths that are polylines connecting junction nodes together. We use an appropriate data structure only involving arrays (named tensors in deep learning frameworks) so that it can be manipulated by the GPU. We show in Fig. S3 an infographic of the data structure. A sequence of node coordinates “*pos*” holds the location of all nodes $i \in [0 \dots n - 1]$ belonging to the skeleton:

$$pos = ((r_0, c_0), (r_1, c_1), \dots, (r_{n-1}, c_{n-1}))$$

where n is the total number of skeleton pixels and $(r_i, c_i) \in [0 \dots H - 1] \times [0 \dots W - 1]$ correspond to the row number and column number, respectively (of skeleton pixel i). The skeleton graph connects junction nodes through paths, which are polylines made up of connected vertices. These paths are represented by the “paths” binary matrix $P_{p,n}$ where element (i, j) is one if node j is in path i . This $P_{p,n}$ is sparse, and, thus, it is more efficient to use the CSR

(compressed sparse row) format, which represents a matrix by three (one-dimensional) arrays respectively containing nonzero values, the column indices and the extents of rows. As $P_{p,n}$ is binary we do not need the array containing non-zeros values. The column indices array, which we name “*path_index*” holds the column indices of all “on” elements:

$$path_index = (j_0, j_1, \dots, j_{n - n_{junctions} + n_{degrees\ sum} - 1}),$$

where $n_{junctions}$ is the total number of junction nodes, $n_{degrees\ sum}$ is the sum of the degrees of all junction nodes and $\forall k \in [0 \dots n - n_{junctions} + n_{degrees\ sum} - 1], j_k \in [0 \dots n - 1]$. The extents of rows array which we name “*path_delim*” holds the starting index of each row (it also contains an extra end element which is the number of non-zeros elements n for easier computation):

$$path_delim = (s_0, s_1, \dots, s_p).$$

Thus, in order to get row i of $P_{p,n}$ we need to look up the slice (s_i, s_{i+1}) of *path_index*. In the skeleton graph case, this representation is also easily interpretable. Indices of path nodes are all concatenated in *path_index* and *path_delim* is used to separate those concatenated paths.

And finally a sequence of integers “degrees” stores for each node the number of nodes connected to it:

$$degrees = (d_0, d_1, \dots, d_{n-1}).$$

As a collection of contours is a type of graph, in order to use a common data structure in our algorithm, we also use the skeleton graph representation for the contours $\{C_i\}$ given by the marching squares algorithm (note we could use other contour detection algorithms for initialization). Each contour is thus an isolated path in the skeleton graph.

In order to fully leverage the parallelization capabilities of GPUs, the largest amount of data should be processed concurrently to increase throughput, i.e., we should aim to use the GPU memory at its maximum capacity. When processing a small image (such as 300×300 pixels from the *CrowdAI dataset*), only a small fraction of memory is used. We thus build a batch of such small images to process them at the same time. As an example, on a GTX 1080Ti, we use a polygonization batch size $B = 1024$ for processing the *CrowdAI dataset*, which induces a significant speedup. Building a batch of images is very simple: they can be concatenated together along an additional batch dimensions, i.e., B images $I_i \in \mathbb{R}^{3 \times H \times W}$ are grouped in a tensor $\mathbf{I} \in \mathbb{R}^{B \times 3 \times H \times W}$. This is the case for the output segmentation probability maps as well as the frame field. However, it is slightly more complex to build a batch of skeleton graphs because of their varying sizes. Given a collection of skeleton graphs $\{(pos_i, degrees_i, path_index_i, path_delim_i)\}_{i \in [0..B-1]}$, all pos_i and $degrees_i$ are concatenated in their first dimension to give batch arrays:

$$pos_{batch} = [pos_0, pos_1, \dots, pos_{B-1}],$$

and:

$$degrees_{batch} = [degrees_0, degrees_1, \dots, degrees_{B-1}].$$

All $path_index_i$ need their indices to be shifted by a certain offset:

$$offset_i = \sum_{k=0}^{i-1} |pos_k|,$$

with $|pos_k|$ the number of points in pos_k , so that they point to the new locations in pos_{batch} and $degree_{batch}$. They are then concatenated in their first dimension: $path_index_{batch} = [path_index_0 + offset_0, \dots, path_index_{B-1} + offset_{B-1}]$. In a similar manner, we concatenate all $path_delim_i$ into $path_delim_{batch}$ while taking care of adding the appropriate offset. We then obtain a big batch skeleton graph which is represented in the same way as a single skeleton graph. In order to later recover individual skeleton graphs in the batch, similar to $path_delim$, we need a $batch_delim$ array that stores the starting index of each individual skeleton

graph in the $path_delim$ array (it also contains an extra end element which is the total number of paths in the batch for easier computation). While we apply the optimization on the batched arrays pos_{batch} , $path_index_{batch}$, and so on, for readability we will now refer to them as pos , $path_index$ and so on. Note that in the case of big images (such as 5000×5000 pixels from the *Inria dataset*), we set the batch size to 1, as the probability maps, the frame field, and the skeleton graph data structure fills the GPU’s memory well.

At this point the data structure is fixed, i.e., it will not change during optimization. Only the values in pos will be modified. This data structure is efficiently manipulated in parallel on the GPU. All the operations needed for the various computations performed in the next sections are run in parallel on the GPU.

We compute other tensors from this minimal data structure which will be useful for computations:

- $path_pos = pos[path_index]$ which expands the positions tensor for each path (junction nodes are thus repeated in $path_pos$).
- A $batch$ tensor which for each node in pos_{batch} stores the index $i \in [0..B-1]$ of the individual skeleton this node belongs to. This is used to easily sample the batched segmentation maps and the batched frame fields at the position of a node.

2.2. Active Skeleton Model

We adapt the formulation of the Active Contours Model (ACM) to an Active Skeleton Model (ASM) in order to optimize our batch skeleton graph. The advantage of using the energy minimization formulation of ACM is to be able to add extra terms if needed (we can imagine adding regularization terms to, e.g., reward 90° corners, uniform curvature, and straight walls).

Energy terms will be parameterized by the node positions $\mathbf{p} \in pos$, which are the variables being optimized. The first important energy term is $E_{probability}$ which aims to fit the skeleton paths to the contour of the building interior probability map $y_{int}(v)$ at a certain probability level ℓ (which we set to 0.5 in practice, just like the isovalue used to initialize the contours by marching squares):

$$E_{probability} = \sum_{\mathbf{p} \in pos} (y_{int}(\mathbf{p}) - \ell)^2.$$

The value $y_{int}(\mathbf{p})$ is computed by bilinear interpolation so that gradients can be back-propagated to \mathbf{p} . Additionally, $y_{int}(\mathbf{p})$ implicitly entails using the $batch$ array to know which slice in the batch dimension of $y_{int} \in \mathbb{R}^{B \times 1 \times H \times W}$ to sample \mathbf{p} from. This will be the case anytime batched image-like tensors are sampled at a point \mathbf{p} . In the case of the marching squares initialization, this $E_{probability}$ energy is

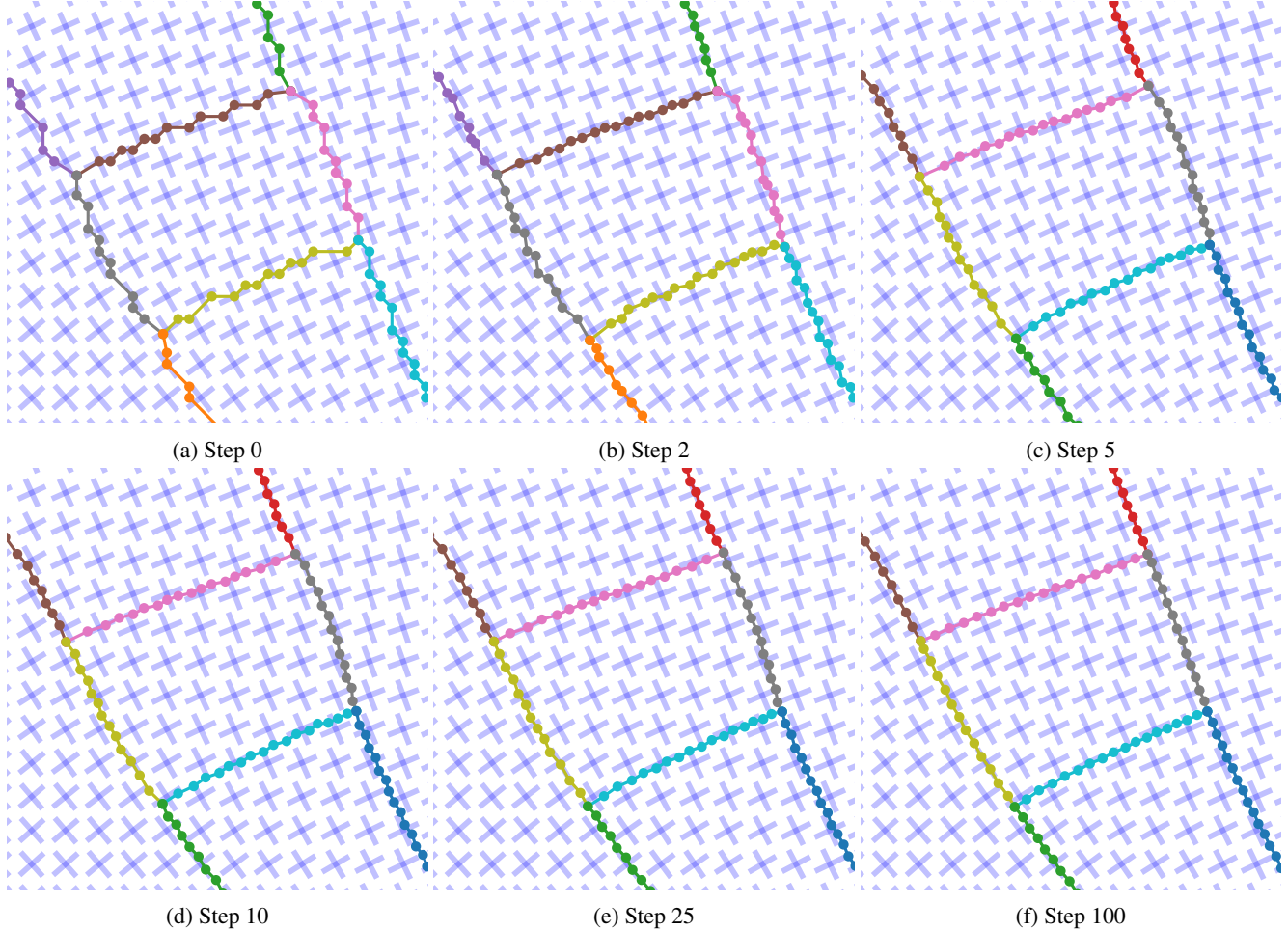


Figure S4: ASM optimization steps (zoomed example). Frame field in blue crosses.

actually zero at the start of optimization, since the initialized contour already is at isovalue ℓ . For the skeleton graph initialization, paths that trace inner walls between adjoining buildings will not be affected since the gradient is zero in a neighborhood of homogeneous values (i.e., $y_{int} = 1$ inside buildings).

The second important energy term is $E_{frame\ field\ align}$ which aligns each edge of the skeleton paths to the frame field. Edge vectors are computed in parallel as:

$$\mathbf{e} = \text{path_pos}[1:] - \text{path_pos}[:-1],$$

while taking care of removing from the energy computation “hallucinated” edges between paths (using the *path_delim* array). For readability we call E the set of valid edge vectors. For each edge vector $\mathbf{e} \in E$, we refer to its direction as $\mathbf{e}_{dir} = \frac{\mathbf{e}}{\|\mathbf{e}\|}$. We also refer to its center point as $\mathbf{e}_{center} = \frac{1}{2}(\text{path_pos}[1:] + \text{path_pos}[:-1])$. The frame

field align term is defined as:

$$E_{frame\ field\ align} = \sum_{\mathbf{e} \in E} |f(\mathbf{e}_{dir}; c_0(\mathbf{e}_{center}), c_2(\mathbf{e}_{center}))|^2.$$

This is the term that disambiguates between slanted walls and corners and results in regular-looking contours.

The last important term is the internal energy term E_{length} which ensures node distribution along paths remains homogeneous as well as tight:

$$E_{length} = \sum_{\mathbf{e} \in E} |\mathbf{e}|^2.$$

All energy terms are then linearly combined:

$$E_{total} = \lambda_{probability} E_{probability} + \lambda_{frame\ field\ align} E_{frame\ field\ align} + \lambda_{length} E_{length}.$$

In practice, the final result is robust to different values of coefficients for each of these three energy terms, and we determine them using a small cross-validation set. The total energy is minimized with the RMSprop [40] gradient descent method with a smoothing constant $\gamma = 0.9$ with an

initial learning rate of $\eta = 0.1$ which is exponentially decayed. The optimization is run for 300 iterations to ensure convergence. Indeed since the geometry is initialized to lie on building boundaries, it is not expected to move more than a few pixels and the optimization converges quickly. See Fig. S4 for a zoomed example of different stages of the ASM optimization.

2.3. Corner-aware polygon simplification

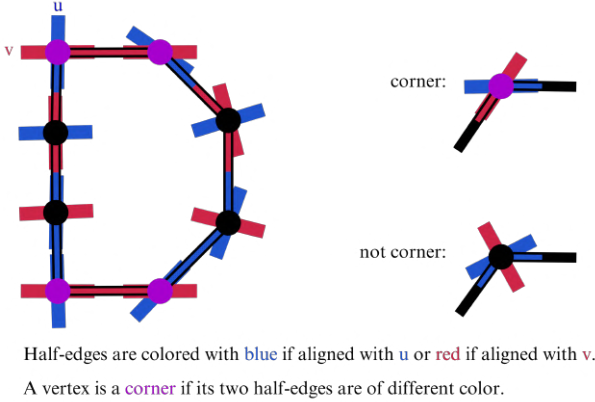


Figure S5: Corner detection using the frame field. For each vertex, the frame field is sampled at that location (with nearest neighbor) and represented by the $\{\pm u, \pm v\}$ vectors.

We now have a collection of connected polylines that forms a planar skeleton graph. As building corners should not be removed during simplification, only polylines between corners are simplified. For the moment our data structure encodes a collection of polyline paths connecting junction nodes in the skeleton graph. However, a single path can represent multiple walls. It is the case for example of an individual rectangular building: one path describes its contour while it has 4 walls. In order to split paths into sub-paths each representing a single wall we need to detect building corners along a path and add this information to our data structure. This is another reason to use a frame field input, as it implicitly models corners: at a given building corner, there are two tangents of the contour. The frame field learned to align one of u or $-u$ to the first tangent and one of v or $-v$ to the other tangent. Thus when walking along a contour path if the local direction of walking switches from $\pm u$ to $\pm v$ or vice versa, it means we have come across a corner, see Fig. S5 for an infographic for corner detection. Specifically for each node i with position $\mathbf{p} = \text{path_pos}[i]$ its preceding and following edge vectors are computed as: $\mathbf{e}_{prev} = \text{path_pos}[i] - \text{path_pos}[i-1]$ and $\mathbf{e}_{next} = \text{path_pos}[i+1] - \text{path_pos}[i]$. As the frame field is represented by the coefficients $\{c_0, c_2\}$ at each pixel, we first need to convert it to its $\{u, v\}$ representation with the simple formulas of eq. 17.

$$\begin{cases} c_0 = u^2 v^2 \\ c_2 = -(u^2 + v^2) \end{cases} \iff \begin{cases} u^2 = -\frac{1}{2} \left(c_2 + \sqrt{c_2^2 - 4c_0} \right) \\ v^2 = -\frac{1}{2} \left(c_2 - \sqrt{c_2^2 - 4c_0} \right) \end{cases} \quad (17)$$

	$ \langle \mathbf{e}_{prev}, \mathbf{u}_{\mathbf{p}} \rangle < \langle \mathbf{e}_{prev}, \mathbf{v}_{\mathbf{p}} \rangle $	$ \langle \mathbf{e}_{prev}, \mathbf{v}_{\mathbf{p}} \rangle < \langle \mathbf{e}_{prev}, \mathbf{u}_{\mathbf{p}} \rangle $
$ \langle \mathbf{e}_{next}, \mathbf{u}_{\mathbf{p}} \rangle < \langle \mathbf{e}_{next}, \mathbf{v}_{\mathbf{p}} \rangle $	False	True
$ \langle \mathbf{e}_{next}, \mathbf{v}_{\mathbf{p}} \rangle < \langle \mathbf{e}_{next}, \mathbf{u}_{\mathbf{p}} \rangle $	True	False

Table S1: Summary table for deciding if node i with position $\mathbf{p} = \text{path_pos}[i]$ is a corner (True) or not (False).

The frame field is sampled at that position \mathbf{p} : $\mathbf{u}_{\mathbf{p}} = u(\mathbf{p})$ and $\mathbf{v}_{\mathbf{p}} = v(\mathbf{p})$. Alignment between $\mathbf{e}_{prev}, \mathbf{e}_{next}$ and $\pm \mathbf{u}_{\mathbf{p}}, \pm \mathbf{v}_{\mathbf{p}}$ is measured with the absolute scalar product so that it is agnostic to the sign of u and v . For example alignment between \mathbf{e}_{prev} and $\pm \mathbf{u}_{\mathbf{p}}$ is measured by $|\langle \mathbf{e}_{prev}, \mathbf{u}_{\mathbf{p}} \rangle|$ and if $|\langle \mathbf{e}_{prev}, \mathbf{u}_{\mathbf{p}} \rangle| < |\langle \mathbf{e}_{prev}, \mathbf{v}_{\mathbf{p}} \rangle|$ then \mathbf{e}_{prev} is aligned to $\pm v$ and not $\pm u$. The same is done for \mathbf{e}_{next} . Finally if \mathbf{e}_{prev} and \mathbf{e}_{next} do not align to the same frame field direction, then node i is a corner. As a summary for corner cases we refer to Table. S1.

Because the path positions are concatenated together in the path_pos tensor, some care must be taken for nodes at the extremities of paths (i.e., junction nodes) as they do not have both preceding and following edges. The path_delim tensor is used to mark those nodes as not corners. Once corners are detected we obtain a tensor $\text{is_corner_index} = \{i \mid \text{node } i \text{ is a corner}\}$ which can be used to separate paths into sub-paths each representing a single wall by merging is_corner_index with the path_delim tensor through concatenation and sorting.

Now that each sub-path polyline represents a single wall between two corners, we apply the Ramer-Douglas-Peucker [32, 13] simplification algorithm separately on all sub-path polylines. As explained in the related works, the simplification tolerance ε represents the maximum Hausdorff distance between the original polyline and the simplified one.

2.4. Detecting building polygons in planar graph

To obtain our final output of building polygons, the collection of polylines is polygonized by detecting connecting regions separated by the polylines. A list of polygonal cells that partition the entire image is thus obtained. The last step computes a building probability value for each polygon using the predicted interior probability map and removes low-probability polygons (in practice those that have an average probability less than 50%).

3. Experimental setup details

3.1. Datasets

CrowdAI dataset. The *CrowdAI dataset* [34] originally has 280741 training images, 60317 validation images, and 60697 test images. All images are 300×300 pixels with unknown ground sampling distance, although they are aerial images. As the ground truth annotations of the test set are

unreleased because of the challenge, we use the original validation set as our test set and discard the original test images as is commonly done by other methods comparing themselves with that dataset [21, 20]. We then use 75% of the original training images as our initial training set and 25% for validation. Our final models are then trained on the entire original training set with hyperparameters selected using our validation test.

Inria dataset. The *Inria dataset* [26] has 360 aerial images of 5000×5000 pixels each with a Ground Sampling Distance of 30 cm. In total, 10 cities from Europe and the USA are represented, each city having 36 images. Each image is accompanied by its building ground truth mask with an average of a few thousand buildings per image. This dataset provides building ground truth in the form of binary mask images for each image. However, our method requires the ground truth annotations to be in vector format (polygons) so that the ground truth for the frame field can be computed: the tangent angle θ_τ used in L_{align} . We thus build two dataset variants with vector annotations.

The first variant is the *Inria OSM dataset* for which we discard completely the original ground truth masks and instead download annotations from Open Street Map (OSM) [30]. Because the OSM annotations are not always aligned, we align them using [15]. We randomly split the images into train (50%), validation (25%), and test (25%) sets. Because the OSM annotations have a lot of missing buildings in certain images, our test results on this dataset are somewhat skewed. Thus, for the test images, we manually select those with few missing buildings in the annotations, giving us 54 test images in total.

The second variant is the *Inria Polygonized dataset* for which we take the original ground truth masks and convert them to polygon format with our polygonization method. In this setting, the input to our network (we used the small U-Net16) is just the binary mask and the output a frame field. In order to train this model, we need a dataset of (binary masks, θ_τ) pairs. We used the OSM annotations of the *Inria OSM dataset*, which we rasterized to obtain the input binary masks and which we used to compute θ_τ . After our model finished training, we applied our frame field polygonization method on the original binary masks of the *Inria dataset* and their predicted frame fields. The new *Inria polygonized dataset* is thus made of (RGB image, polygonized annotations) pairs. We thus obtain the same ground truth as the original dataset but in vector format. This allows us to only use the same ground truth data as the other competitors of the Inria Aerial Image Labeling challenge and thus we can directly compare our method to them. Thus we keep the original train and test splits which do not have any cities overlap and tests cross-city generalization (the principal aim of the associated challenge). We then split the original train

split into our train (75%) and validation (25%) splits.

Private dataset. The *private dataset* is a large-scale dataset of satellite images built by a company we collaborate with. The images in this dataset were acquired using three types of satellites (Pleiades, WorldView, and GeoEye) over different types of cities (dense, industrial, residential areas, and city centers). We uniformized the image sampling at 50 cm/pixel spatial resolution, with 3-band RGB images. 57 images of 30 cities across 5 continents are present in the training dataset. The size of images varies from around 2000×2000 pixels to 20000×20000 pixels. The total dataset covers an area spanning around 700 km². The building outline polygons were manually labeled precisely by an expert. Satellite images are more challenging than aerial images (such as the CrowdAI and Inria images) because they are less clear due to atmospheric effects. This dataset also contains much more varied images compared to CrowdAI and Inria, making up for its smaller size. We preprocess the training images by splitting them into smaller 512×512 pixel patches. We then keep 90% of patches for training and 10% for validation.

3.2. Metrics

IoU, AP and AR. The usual metric for the image segmentation task is Intersection over Union (IoU) which computes the overlap between a predicted segmentation and the ground truth annotation. The IoU is then used to compute other metrics such as the MS COCO [22] Average Precision (AP and its variants AP_{50} , AP_{75} , AP_S , AP_M , AP_L) and Average Recall (AR and its variants AR_{50} , AR_{75} , AR_S , AR_M , AR_L) evaluation metrics. Precision and recall are computed for a certain IoU threshold: detections with an IoU above the threshold are counted as true positives while others are false positives and ground truth annotations with an IoU below the threshold are false negatives. Each object is also given a score value representing the model’s confidence in the detection. In our case, it is the mean value of the interior probability map inside the detection. The Precision-Recall curve can be obtained by varying the score threshold that determines what is counted as a model-predicted positive detection. Average Precision (AP) is the average value of the precision across all recall values and Average Recall (AR) is the maximum recall given a fixed number of detections per image (100 in our case). Finally, the mean Average Precision (mAP) is calculated by taking the mean AP over multiple IoU thresholds (from 0.50 to 0.95 with a step of 0.05). Likewise for the mean Average Recall (mAR). Following MS COCO’s convention, we make no distinction between AP and mAP (and likewise AR and mAR) and assume the difference is clear from context. The AP_{50} variant is AP computed with a single IoU threshold of 50% (similarly for AP_{75} , AR_{50} , and AR_{75}). The AP_S , AP_M and AP_L

variants are AP computed for small ($area < 32^2$), medium ($32^2 < area < 96^2$) and large ($area > 96^2$) objects respectively (like-wise for the AR equivalents).

Max tangent angle error. We introduce a max tangent angle error metric between predicted polygons and the ground truth to capture the regularity of the predicted contours. A max tangent angle scalar error is computed for each predicted contour. Only predicted contours with at least 50% overlap with the ground truth are selected, so that their measure makes sense. Each predicted contour is first sampled homogeneously with points $\{P_i\}_{i \in [1..n]}$ (specifically a point is sampled every 0.1 pixel). Then the P_i points are projected to the ground truth, meaning for each P_i we find the closest point Q_i belonging to the ground truth annotation. For both sequences of points P_i and Q_i , corresponding normed tangent directions are computed as:

$$T(P_i) = \frac{P_{i+1} - P_i}{\|P_{i+1} - P_i\|} \quad \text{and} \quad T(Q_i) = \frac{Q_{i+1} - Q_i}{\|Q_{i+1} - Q_i\|}.$$

The angle differences between the two are computed from the scalar product:

$$\Delta\theta_i = \cos^{-1}(\langle T(P_i), T(Q_i) \rangle).$$

Before computing the maximum angle error $\max_i \Delta\theta_i$ along the whole contour, some angle errors $\Delta\theta_i$ need to be filtered out as they are invalid. Angle error invalidity is due to the projection step. Indeed around ground truth corners, part of the predicted contour will be squashed to be zero-length for example. Another issue is when P_i and P_{i+1} are projected to two different ground truth polygon sides: the projected edge $P_{i+1} - P_i$ does not represent a ground truth tangent anymore. We thus filter out tangents whose projection is stretched more than a factor of 2, i.e., we keep all $\Delta\theta_j, \forall j \in V$ where $V = \{j \mid j \in [1..n], \frac{1}{2} < \frac{\|Q_{i+1} - Q_i\|}{\|P_{i+1} - P_i\|} < 2\}$. The final max tangent angle error for that contour is then:

$$E_{max \text{ tangent angle}} = \max_{j \in V} \Delta\theta_j.$$

As each contour gives a scalar error, we aggregate all the errors for a certain dataset by averaging this max tangent angle error metric.

4. Additional results

4.1. CrowdAI dataset

4.1.1 Complexity vs. fidelity

For the polygon complexity/fidelity trade-off ablation study we plot the AP and AR scores for difference simplification tolerance values ε on the *CrowdAI dataset*.

Method	Mean max angle error ↓
UResNet101 (no field), simple poly.	51.9°
UResNet101 (with field), simple poly.	45.1°
U-Net variant [9], ASIP poly. [20]	44.0°
U-Net variant [9], UResNet101 our poly.	36.6°
Zorzi et al. [44] poly.	36.8°
UResNet101 (no L_{smooth}), our poly.	33.6°
UResNet101 (no $L_{int \ align}$ and $L_{edge \ align}$), our poly.	33.5°
UResNet101 (no $L_{int \ edge}$), our poly.	33.4°
UResNet101 (no $L_{align90}$), our poly.	33.2°
PolyMapper [21]	33.1°
UResNet101 (with field), our poly.	31.9°

Table S2: Mean *max tangent angle errors* over all the original validation polygons of the *CrowdAI dataset* [34].

We perform an analysis of the polygonization complexity/fidelity trade-off by changing the tolerance value ε of the baseline simplification method and our corner-aware method. Fig. S6a shows that preventing the removal of building corners ensures key points of the contours and the global shape of the building remain intact even with extreme simplification tolerance values. We also plot the AP and AR values of both methods while increasing the tolerance value ε in Fig. S6. As expected the score of our method does not drop, unlike the simple polygonization method.

Our polygonization method allows the complexity-to-fidelity ratio to be tuned with the easy-to-interpret tolerance value ε of the Ramer-Douglas-Peucker algorithm, unlike ASIP [20], which uses a non-intuitive parameter λ to balance complexity and fidelity energies during polygonal partition optimization. Finally, PolyMapper [21] does not have the ability to tune the complexity-to-fidelity ratio.

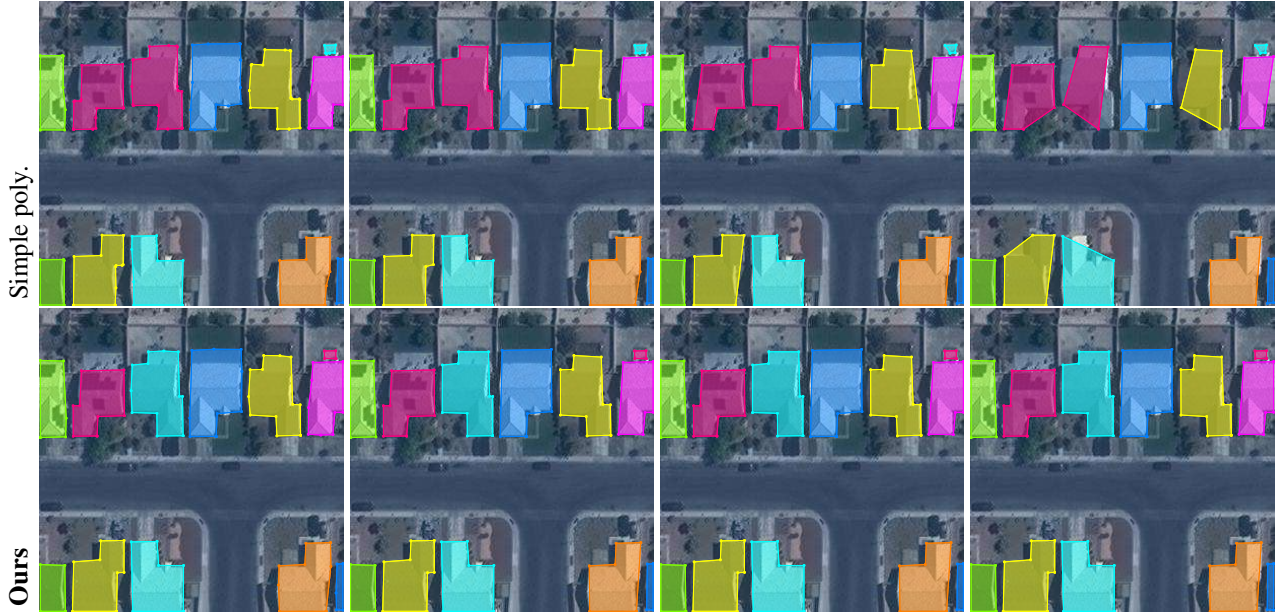
4.1.2 Ablation study

We visualize the predicted classification maps from each ablation study for an example test sample in Fig. S7. Both for the U-Net16 and DeepLab101 backbones, the (full) method yields more regular classification maps with sharper corners compared to (no field). Additionally, only learning the frame field with (no coupling losses) is insufficient, as can be seen in Fig. S7d.

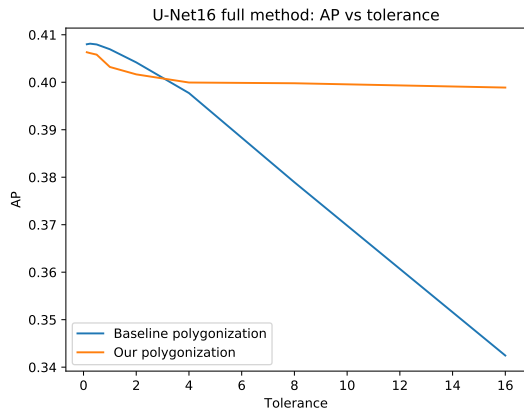
We observe the effect of only optimizing for IoU when removing coupling losses: we see that it does not impact AP and AR metrics in Table S3, while in Fig. S7 the (full) segmentations are clearly sharper compared to the (no coupling losses) ones.

In terms of AP and AR metrics, adding a frame field improves the final score (full) compared to (no field) for all backbones: U-Net16, DeepLab101 and UResNet101 (see Table S3).

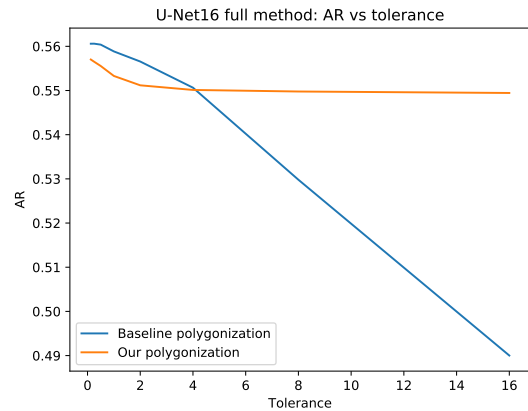
We also visually compare our frame field polygonization method with the simple baseline polygonization algorithm (both when the frame field is computed and when it is not) in Fig. S8. The UResNet101 without frame field learning



(a) Effect of increasing the simplification tolerance value ε from 0.5 px (left), then 2 px, then 8 px and 16 px (right).



(b) AP for both our corner-aware method and the simple polygonization for various tolerance value ε .



(c) AR for both our corner-aware method and the simple polygonization for various tolerance value ε .

Figure S6: Comparison between the baseline simplification algorithm with our corner-aware one. Both take the same classification map as input, but the baseline does not use the frame field. The corner-aware simplification guarantees that no corners will be simplified, regardless of the tolerance value ε .

and whose results are polygonized with the simple method performs the worst (see Fig. S8c), with the UResNet101 with frame field learning and whose results are polygonized with the simple method performs already much better (see Fig. S8b). Our UResNet101 with frame field learning and whose results are polygonized with our frame field polygonization method performs the best, with better corners using fewer vertices (see Fig. S8a). We can see our method provides the missing information needed to resolve ambiguous cases for polygonization and outputs more regular polygons.

Finally Table S2 and S3 also hold results for additional experiments of the ablation study which each remove a loss during training. We observe that removing one of those losses does not impact the AP or AR result of the final polygonization. However, if one of those loss is removed we observe a performance drop in terms of *max tangent angle errors*, with result polygons for all such experiments having a mean error slightly higher than PolyMapper (at 33.1°) while our full method achieves a mean error of 31.6°.

Method	AP \uparrow	AP $_{50}$ \uparrow	AP $_{75}$ \uparrow	AP $_S$ \uparrow	AP $_M$ \uparrow	AP $_L$ \uparrow	AR \uparrow	AR $_{50}$ \uparrow	AR $_{75}$ \uparrow	AR $_S$ \uparrow	AR $_M$ \uparrow	AR $_L$ \uparrow
U-Net16 (no field), mask	50.9	74.3	59.5	24.5	65.6	66.3	55.9	77.9	64.7	29.8	71.2	74.6
U-Net16 (no field), simple poly.	50.5	76.6	59.1	22.6	66.2	69.3	54.8	78.5	63.5	26.8	71.2	75.2
U-Net16 (no coupling losses), mask	53.7	77.7	62.8	25.7	69.0	68.9	57.7	79.2	66.4	31.0	73.4	74.4
U-Net16 (with field), mask	53.6	77.8	62.8	25.1	69.4	69.5	57.6	79.0	66.4	29.7	74.1	75.2
U-Net16 (with field), simple poly.	49.6	73.8	58.1	21.2	65.5	67.0	53.8	75.6	62.2	25.5	70.5	72.5
U-Net16 (with field), our poly.	50.5	76.6	59.3	20.4	67.4	69.0	55.3	78.1	64.0	25.7	72.8	75.0
DeepLab101 (with field)	54.9	78.1	64.9	25.6	71.2	76.8	58.7	79.8	68.1	29.5	75.8	81.6
DeepLab101 (no field)	52.8	75.2	61.8	26.1	67.7	75.0	57.8	78.4	66.7	30.3	73.7	81.8
UResNet101 (no field), mask	62.4	86.7	72.7	36.2	76.3	81.1	67.5	90.5	77.4	46.8	79.5	86.5
UResNet101 (no field), simple poly.	61.1	87.4	71.2	35.1	74.5	82.3	64.7	89.4	74.1	41.7	77.9	85.7
UResNet101 (with field), mask	64.5	89.3	74.6	40.3	76.6	84.0	68.1	91.0	77.7	47.5	80.0	86.7
UResNet101 (with field), simple poly.	61.7	87.7	71.5	35.8	74.9	83.0	65.4	89.9	74.6	42.6	78.6	85.9
UResNet101 (with field), our poly.	61.3	87.5	70.6	34.0	75.1	83.1	65.0	89.4	73.9	41.2	78.7	86.0
UResNet101 (no $L_{align90}$), mask	64.2	88.6	74.6	40.0	76.4	83.7	67.8	90.9	77.5	47.1	79.7	86.4
UResNet101 (no $L_{align90}$), simple poly.	61.4	87.7	71.4	35.4	74.5	82.7	65.0	89.7	74.4	42.1	78.2	85.6
UResNet101 (no $L_{align90}$), our poly.	61.1	87.5	70.6	34.1	74.9	82.8	64.7	89.3	73.8	41.2	78.4	85.6
UResNet101 (no $L_{int edge}$), mask	63.8	88.5	74.4	39.6	75.9	83.3	67.3	90.7	77.0	46.6	79.3	86.2
UResNet101 (no $L_{int edge}$), simple poly.	61.0	87.6	70.6	35.2	74.1	82.4	64.6	89.5	74.0	41.7	77.8	85.3
UResNet101 (no $L_{int edge}$), our poly.	60.9	87.4	70.5	33.7	74.4	82.5	64.4	89.1	73.4	40.7	78.1	85.4
UResNet101 (no $L_{int align}$ and $L_{edge align}$), mask	64.7	89.3	74.7	40.5	76.7	84.2	68.2	91.0	77.9	47.6	80.1	86.8
UResNet101 (no $L_{int align}$ and $L_{edge align}$), simple poly.	61.8	87.7	71.5	35.8	74.9	83.3	65.4	89.9	74.7	42.5	78.6	86.0
UResNet101 (no $L_{int align}$ and $L_{edge align}$), our poly.	61.5	87.5	71.3	34.2	75.2	83.4	65.0	89.5	74.0	41.3	78.8	86.1
UResNet101 (no L_{smooth}), mask	64.2	88.6	74.6	40.1	76.5	83.5	67.8	90.8	77.5	47.2	79.8	86.1
UResNet101 (no L_{smooth}), simple poly.	61.6	87.7	71.5	35.7	74.8	82.6	65.2	89.7	74.5	42.3	78.4	85.4
UResNet101 (no L_{smooth}), our poly.	61.3	87.5	70.7	34.1	75.0	82.7	64.8	89.3	73.9	41.1	78.6	85.5
U-Net variant [9], UResNet101 our poly.	67.0	92.1	75.6	42.1	84.2	92.7	73.2	93.5	81.1	48.8	87.3	95.4
Mask R-CNN [16] [35]	41.9	67.5	48.8	12.4	58.1	51.9	47.6	70.8	55.5	18.1	65.2	63.3
PANet [24]	50.7	73.9	62.6	19.8	68.5	65.8	54.4	74.5	65.2	21.8	73.5	75.0
PolyMapper [21]	55.7	86.0	65.1	30.7	68.5	58.4	62.1	88.6	71.4	39.4	75.6	75.4
U-Net variant [9], ASIP poly. [20]	65.8	87.6	73.4	39.3	87.0	91.9	78.7	94.3	86.1	57.2	91.2	97.6

Table S3: AP and AR results on the *CrowdAI dataset* [34] for all polygonization experiments. (with field) refers to models trained with our full frame field learning method. (no field) refers to models trained without any frame field output. “mask” refers to the output raster segmentation mask of the network, “our poly.” refers to our frame field polygonization method, and “simple poly.” refers to the baseline polygonization of marching squares followed by Ramer-Douglas-Peucker simplification.

Method	Time (sec) \downarrow	Hardware
PolyMapper [21]	0.38	GTX 1080Ti
ASIP [20]	0.15	Laptop CPU
Zorzi et al. [44]	0.11	GTX 1080Ti
Ours	0.04	GTX 1080Ti

Table S4: Average time to extract buildings from a 300×300 pixel patch. **Ours** refers to UResNet101 (with field), our poly. ASIP’s time does not include model inference.

4.1.3 Additional runtimes

We report here the average runtimes for a 300×300 pixel patch of the different steps of the building polygonization pipeline of Zorzi et al. [44] along with corresponding GPU memory allocation (GTX 1080Ti):

1. segmentation: 0.152s with 20% GPU memory,
2. regularization: 0.269s with 12% GPU memory,
3. mask2poly: 0.257s with 19% GPU memory.

As we optimized our own method for maximum throughput, we want to compare to previous methods assuming perfect parallelization (as is done for the ASIP method in the main paper). Zorzi et al. would then get these runtimes:

1. segmentation: 0.0304s,
2. regularization: 0.03228s,
3. mask2poly: 0.04883s,

for a total of 0.11151s. For comparison we include the runtimes of all methods in Table S4, where we observe our method being competitive compared to previous works.

4.2. Inria OSM dataset

We show bigger crops of the result of our frame field polygonization in Fig. S9, S10, and S11. We observe the ability of our method to separate adjoining buildings, handle complex shapes with big buildings having non-rectangular shapes and possibly holes.

4.3. Inria polygonized dataset

We show a larger result comparison to other methods on the *Inria Polygonized dataset* in Fig S12, including the two best methods on the public leaderboard². While the result from “Eugene Khvedchenya” and ICTNet achieve an mIoU over 80%, they detect buildings with segmentation masks

²<https://project.inria.fr/aerialimagelabeling/leaderboard/>



Figure S7: Classification predictions on a test sample for all training ablation studies.

that need polygonization. We thus used the simple polygonization method which follows the boundaries in the segmentation raster image. Their results have blob-like features, with rounded corners and non-regular contours.

In order to compare to the ASIP polygonization method, we started to run the ASIP algorithm on the 180 output probability maps of our network, corresponding to the 180 test images. However the ASIP method is not well-suited for such big images (5000×5000 pixels) with thousands of buildings, requiring a very high number of iterations (that we set to 10000). The runtime of ASIP varies greatly depending on the building density of images. For the most dense ones, it did not finish within a day of computation, making it impractical to run on the whole test dataset. As such we compare to the ASIP method only on the *CrowdAI dataset*.

4.4. Private dataset

Because training on the *private dataset* must be done on a restricted computer with limited access, we only train two models: U-Net16 (full) and U-Net16 (no field) until validation loss converges (around 1500 epochs). First we display segmentation raster outputs in Fig S13, S14 and S15 and final polygonal buildings in Fig S16, S17 and S18. Satellite images being more challenging than aerial images, non-regularized segmentations (no field) appear even more rounded than usual. However, our frame field learning and polygonization (with field) still outputs clean, regular geometry, and separates adjoining buildings.



(a) **Ours**: UResNet101 with frame field learning (full) and frame field polygonization



(b) UResNet101 with frame field learning and simple polygonization



(c) UResNet101 (no field) learning and simple polygonization

Figure S8: Extracted polygons with: our (full) frame field learning and polygonization method; our frame field learning and simple polygonization method; the (no field) learning and simple polygonization baseline. A low tolerance of $\varepsilon = 0.125$ pixel was chosen to compare precise contours.

U-Net16 (no field), simple poly.



Ours: U-Net16 (with field), our poly.



Figure S9: Crop of results on the “sfo19” image from the *Inria OSM dataset*.

U-Net16 (no field), simple poly.



Ours: U-Net16 (with field), our poly.



Figure S10: Crop of results on the “innsbruck19” image from the *Inria OSM dataset*.

U-Net16 (no field), simple poly.



Ours: U-Net16 (with field), our poly.

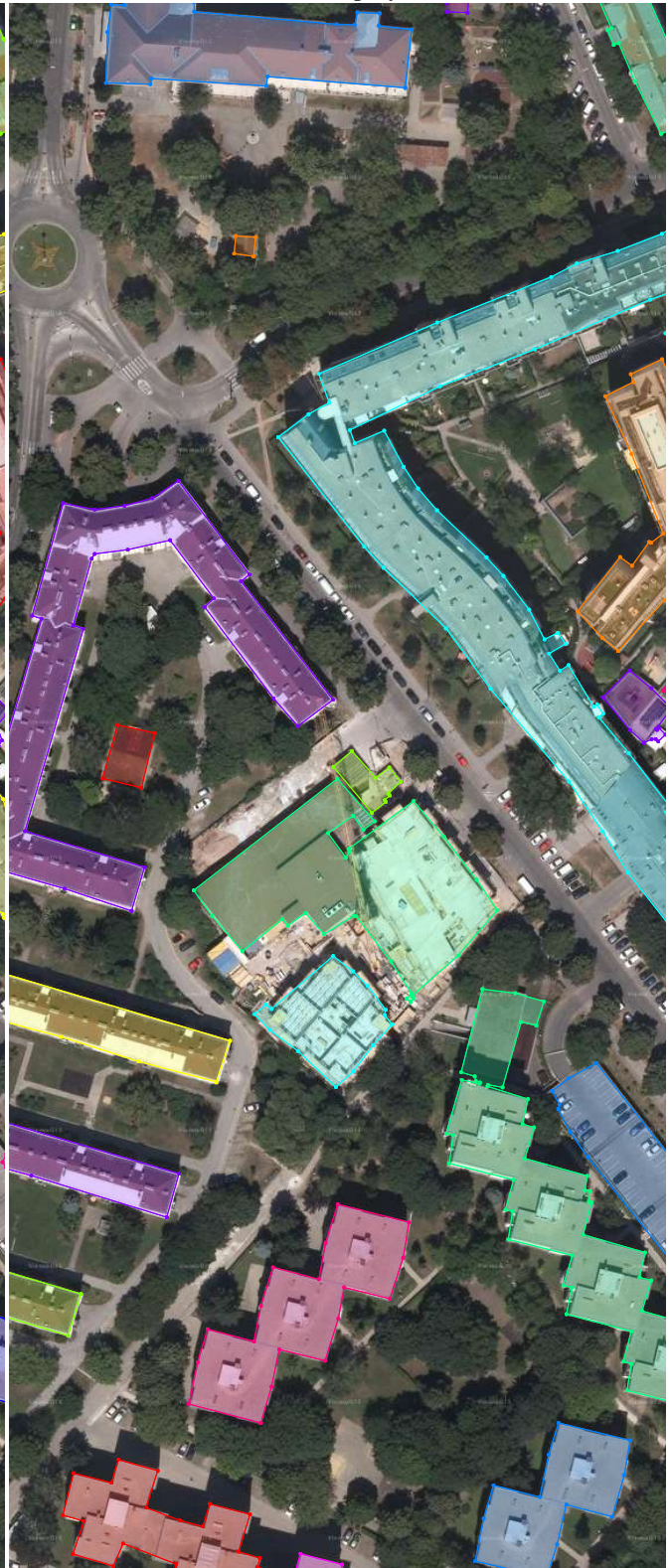


Figure S11: Crop of results on the “vienna36” image from the *Inria OSM dataset*.

Eugene Khvedchenya², simple poly.



ICTNet [6], simple poly.



Zorzi et al. [44] poly.

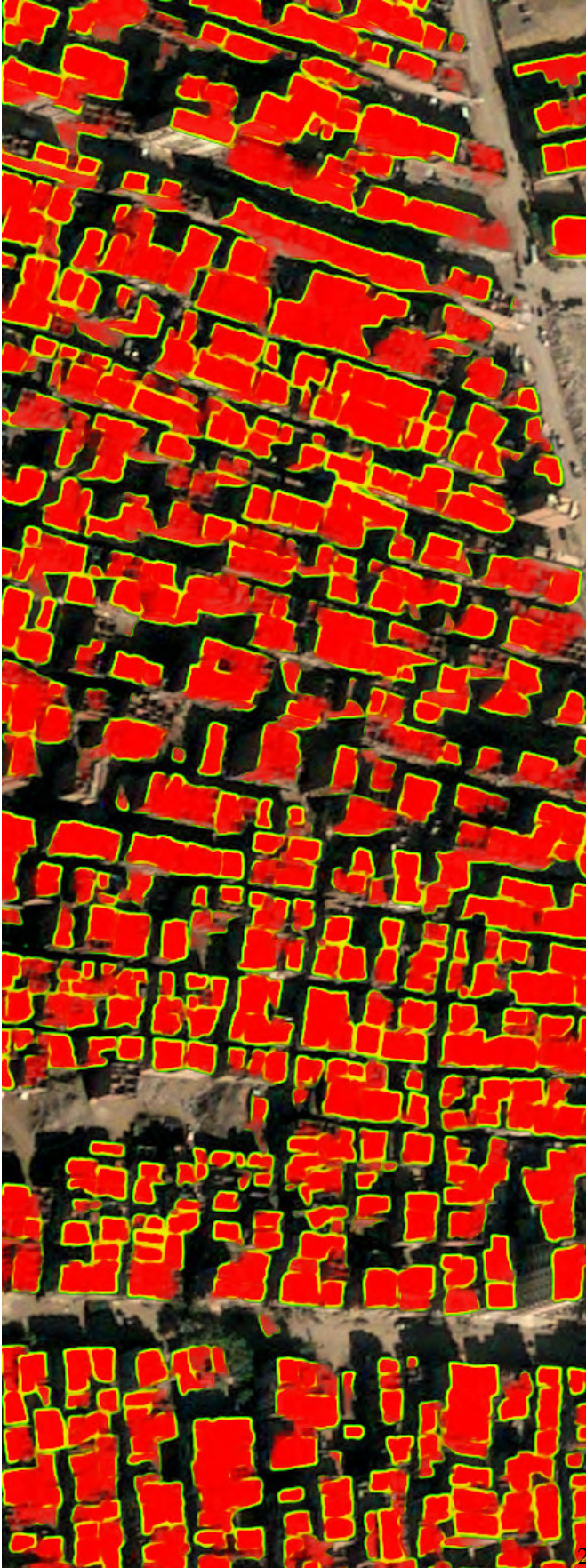


Ours: UResNet101 (with field), our poly.



Figure S12: Crop of results on an *Inria Polygonized dataset* test image.

U-Net16 (no field)



Ours: U-Net16 (with field)

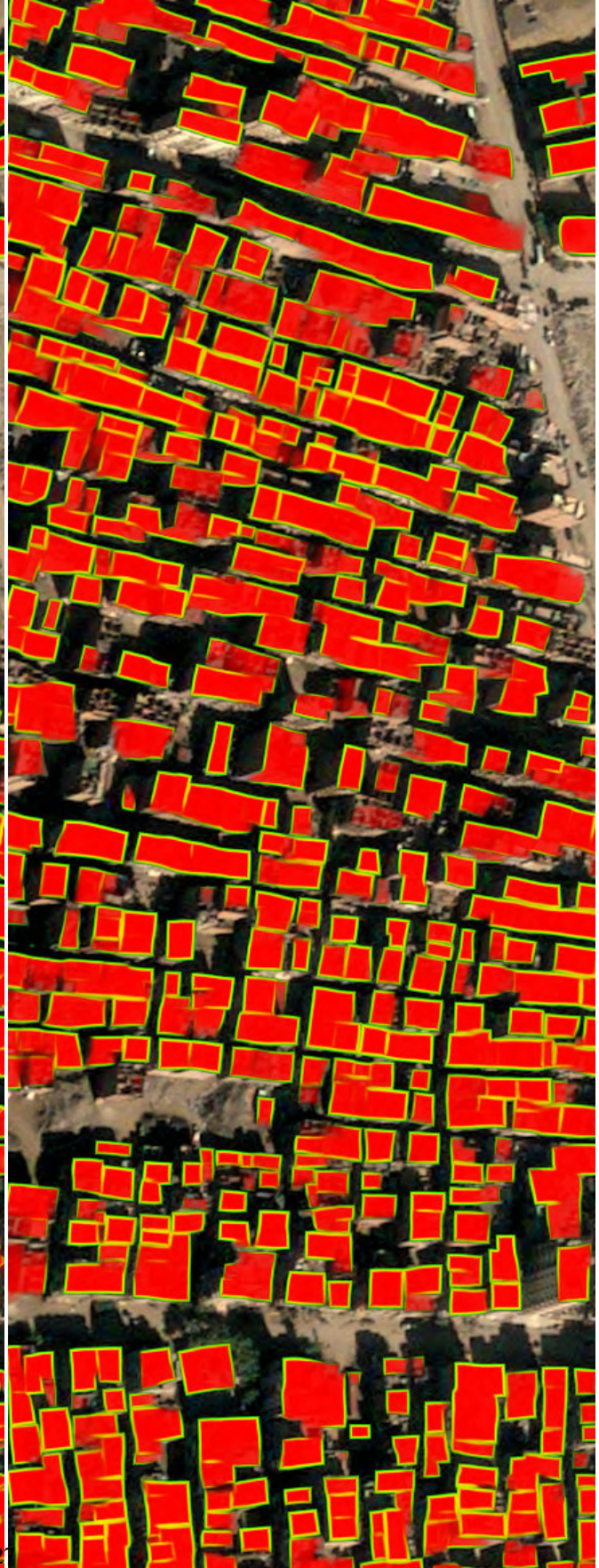


Figure S13: Crop results on the “Egypt” test image of the private dataset.

U-Net16 (no field)

Ours: U-Net16 (with field)



Figure S14: Crop results on the “Bangkok” test image of the private dataset.

U-Net16 (no field)

Ours: U-Net16 (with field)

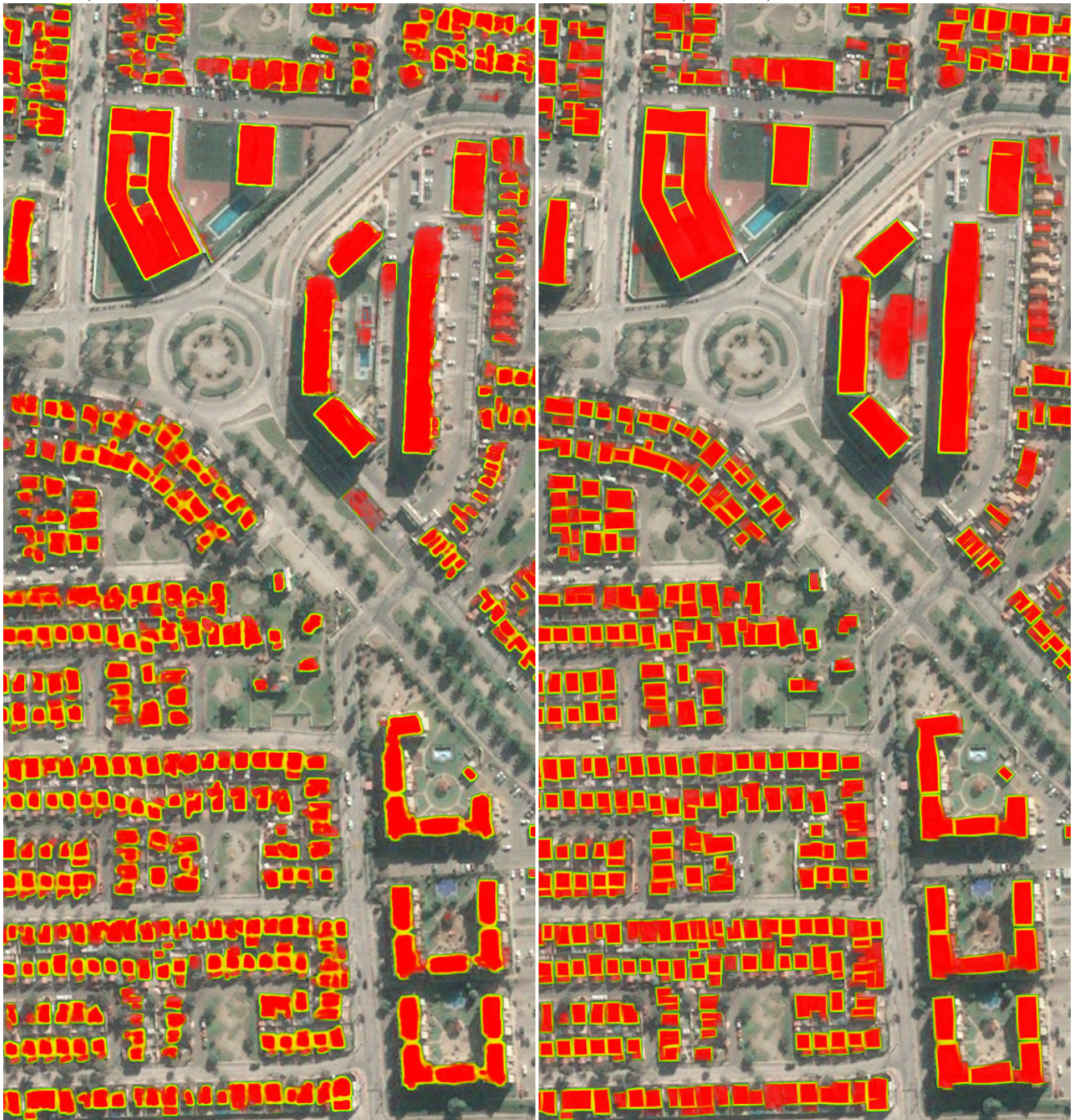


Figure S15: Crop results on the “Chile” test image of the private dataset.

U-Net16 (no field), simple poly.



Ours: U-Net16 (with field), our poly.



Figure S16: Crop results on the “Egypt” test image of the *private dataset*.

U-Net16 (no field), simple poly.

Ours: U-Net16 (with field), our poly.



Figure S17: Crop results on the “Bangkok” test image of the *private dataset*.

U-Net16 (no field), simple poly.

Ours: U-Net16 (with field), our poly.



Figure S18: Crop results on the “Chile” test image of the *private dataset*.