

Aerial2Poly: Aerial Imagery to Polygon Vectorization

Anonymous CVPR submission

Paper ID 4

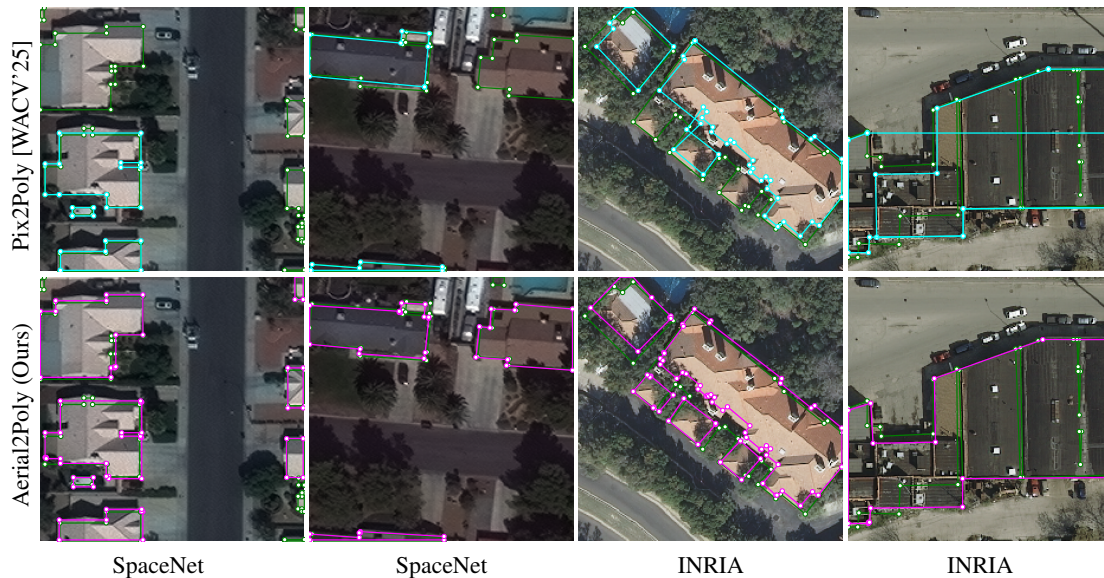


Figure 1. Aerial2Poly extends autoregressive polygon sequencing to aerial building footprint extraction and introduces Gaussian-biased multi-trial decoding to recover complete, geometrically faithful boundaries. Without task-specific modules or post-hoc vectorization, Aerial2Poly sets new state-of-the-art results on the challenging SpaceNet [6], INRIA [17], and WHU Buildings [11] benchmarks. We visualize our predicted polygons (magenta) against those from the strong baseline Pix2Poly [1] (cyan) on exemplar SpaceNet and INRIA validation images and show that our model’s predictions align more closely to the ground truth (green).

Abstract

001 Building footprint extraction from aerial imagery is often
 002 treated as a mask-first problem, followed by contour extrac-
 003 tion, simplification, or geometric post-processing to recover
 004 polygons. We instead formulate it as scene-level autore-
 005 regressive polygon generation. We present Aerial2Poly, an
 006 encoder-decoder model that predicts the polygon set of an
 007 image directly as a sequence of coordinate, separator, and
 008 class tokens. Aerial2Poly encodes connectivity in traversal
 009 order without relying on object prompts or auxiliary con-
 010 nectivity modules. To make this formulation effective and
 011 scalable in practice, we introduce Gaussian-biased decoding,
 012 a recall-oriented anti-EOS objective with noise polygons, and
 013 a vector-space trial consolidation procedure. Aerial2Poly

sets new state-of-the-art results on the challenging SpaceNet,
 INRIA, and WHU Buildings benchmarks, including a 10.36%
 relative C-IoU gain on WHU Buildings. These results show
 that a simple scene-level sequence interface is a strong foun-
 dation for direct polygon extraction from aerial imagery.

1. Introduction

Building footprint extraction from aerial imagery is a core
 problem in remote sensing, with applications in mapping,
 urban analysis, disaster response, and GIS pipelines. The
 desired output is vector geometry rather than raster masks:
 polygons preserve boundaries, vertices, and topology, and
 they integrate directly with editors, spatial databases, and

027 cadastral records. Yet many current systems remain mask-
028 first, predicting dense segmentations and only later convert-
029 ing them into polygons through contour extraction, simpli-
030 fication, or geometric post-processing [7, 13, 18, 22]. This
031 conversion is lossy: small mask errors surface as shifted
032 vertices, irregular jagged outlines, and missing corners.
033 The final polygon quality is bounded as much by the post-
034 processing chain as by the recognition model itself.

035 Recent work moves toward direct polygon prediction
036 via connectivity graphs of building corners [1, 19, 23], but
037 important gaps remain. Decoupling recognition from con-
038 nectivity introduces a failure mode in which every corner
039 is well localized yet the polygon is topologically broken.
040 The result is the crossed edges, missing rings, and merged
041 neighbors visible in Figure 1. Even modern autoregressive
042 approaches [4] decode one instance at a time conditioned
043 on a previously detected box, rather than describing the full
044 scene in one pass. Aerial tiles vary widely in object count
045 and polygon complexity, so an ideal model should produce
046 the full polygon output of a tile in a single pass.

047 We address these gaps with Aerial2Poly, a scene-
048 level autoregressive model for aerial polygon extraction.
049 Aerial2Poly generates the polygon set of the entire scene
050 as a single sequence of coordinate, separator, and class to-
051 kens. Connectivity is encoded by the traversal order itself, so
052 there is no need for an auxiliary matching module, contour
053 extractor, or per-instance prompt. Moreover, our approach
054 efficiently stops the polygon decoding procedure once the
055 scene is described rather than running to a fixed maximum
056 length [3]. The design of Aerial2Poly comprises three sim-
057 ple yet effective components: a Gaussian coordinate bias for
058 locally coherent decoding, a recall-oriented anti-EOS objec-
059 tive trained with synthetic noise polygons, and a multi-trial
060 consolidation step that operates entirely in polygon space.

061 **Main Contributions.** (1) We present Aerial2Poly, a scal-
062 able, end-to-end, scene-level autoregressive model for aerial
063 polygon extraction. (2) Our framework adapts the autore-
064 gressive sequence formulation to efficiently work on dense
065 aerial scenes with the introduction of Gaussian-biased coor-
066 dinate decoding, an anti-EOS recall objective paired with a
067 noisy-polygon training procedure, and a vector-space multi-
068 trial consolidation scheme. (3) Aerial2Poly demonstrates
069 excellent performance and scalability by establishing state-
070 of-the-art results, with compelling margins, on the challeng-
071 ing INRIA, SpaceNet, and WHU Buildings benchmarks.

072 2. Related Work

073 **Polygons from Segmentation Masks.** Early work on
074 polygonal building extraction relied on pixel-perfect seg-
075 mentation maps learned by a convolutional neural network
076 (CNN) to produce high-quality vector shapes via a separate
077 post-processing step. Zhao et al. [21] proposed a bound-

ary regularization method to refine the building instances
078 predicted from Mask-RCNN [9]. Mask2Poly [22] trained
079 two separate CNNs to refine raster segmentation masks into
080 vector polygons by using a combination of adversarial, re-
081 construction, and regularized losses. FFL [7] imposed frame
082 fields on top of raster segmentation maps via a multi-task
083 objective, which were then used to generate building poly-
084 gons. HiSup [18] strengthened the multi-stage polygoniza-
085 tion pipeline with line-segment attraction fields and hierar-
086 chical supervision over masks, but polygon recovery still
087 occurs through sub-optimal geometric post-processing. 088

Polygons as Connectivity Graphs. A promising class of
089 end-to-end approaches formulates the polygonization task as
090 predicting vertices in a connectivity graph of edges [14–16].
091 PolyWorld [23] and TopDig [19] proposed to first detect
092 building corner candidates from encoded image features fol-
093 lowed by a matching network to connect polygon vertices.
094 More recently, Pix2Poly [1] extends this graph-based frame-
095 work by adopting an image-to-sequence transformer network
096 to detect building corners as a discrete sequence prediction
097 task [3]. Although Pix2Poly achieves strong polygonal re-
098 sults, topological errors and misalignment can still arise from
099 the auxiliary graph connectivity module. 100

Autoregressive Polygon Sequencing. An emerging line
101 of work treats visual geometry predictions as a series of
102 sequence generations. Pix2Seq [3] casts object detection
103 as autoregressive prediction over quantized coordinates and
104 class tokens with a generic encoder-decoder transformer
105 network. Pix2Seq v2 [4] extends the same interface to mul-
106 tiple tasks, including polygon-based instance segmentation.
107 These approaches establish a shared token vocabulary for
108 geometric outputs, but they do not study scene-level polygon
109 decoding for remote sensing applications such as building
110 extraction, map generation, and aerial object detection. 111

Recent HiT [20] and GeoFormer [12] approaches have
112 adapted autoregressive polygon sequencing to aerial imagery,
113 but not without limitations. HiT requires an auxiliary re-
114 gion proposal network for building detection together with
115 polygon serialization. GeoFormer requires multiple forward
116 passes to extract all buildings and is not scalable in dense
117 urban scenes. By contrast, our method encodes connectivity
118 directly in autoregressive traversal order and efficiently de-
119 codes the polygon set of the full image in a single pass with-
120 out task-specific auxiliary modules or object-level prompts. 121

122 3. The Aerial2Poly Framework

We cast building footprint extraction as scene-level pixel-
123 to-sequence prediction. Given an aerial image, the decoder
124 emits the polygon set directly as coordinate, separator, and
125 class tokens. The remainder of this section describes se-
126 quence construction, the model architecture and training 127

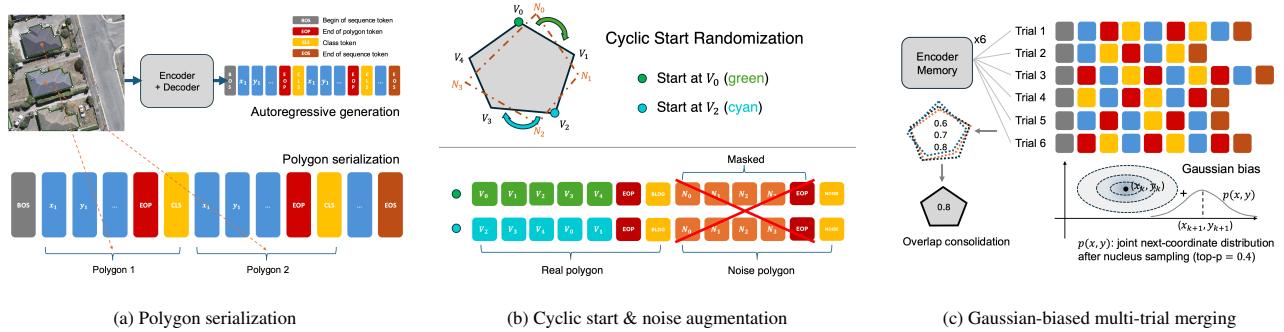


Figure 2. **Overview of Aerial2Poly.** (a) Building footprint polygons are serialized into a flat token sequence of quantized coordinates, polygon separators (EOP), class tokens, and a shared BOS/EOS framing. A ViT encoder and a Transformer decoder predict this sequence autoregressively. (b) During training, cyclic start randomization lets any vertex begin a polygon’s token subsequence (top), and synthetic noise polygons are appended with a dedicated `noise` class and masked coordinate losses to improve recall (bottom). (c) At inference, multiple stochastic trials are decoded with nucleus sampling; a Gaussian spatial bias on coordinate logits encourages locally coherent boundaries, and the resulting polygon sets are merged in vector space via confidence filtering and overlap consolidation.

128 objective, and inference, with the main serialization and
129 decoding steps illustrated in Figure 2.

130 **Sequence Construction.** Let I denote the input aerial
131 image. We serialize its annotation as a single token sequence.
132 The i th polygon ring of I is encoded as

$$133 P_i = [x_1^i, y_1^i, \dots, x_{n_i}^i, y_{n_i}^i, \text{EOP}, c_i], \quad (1)$$

134 where coordinates are uniformly quantized into B bins, c_i
135 is a class token, and in all experiments B equals the in-
136 put image resolution. Rings are concatenated between BOS
137 (begin-of-sequence) and EOS (end-of-sequence) in a shared
138 vocabulary over coordinate, class, and special tokens. Be-
139 cause the polygon set is unordered and each ring is cyclic,
140 we randomize ring order and circularly shift each ring’s start
141 vertex during training [4]. For annotations with holes, ex-
142 terior rings keep the object class and interior rings use a
143 dedicated interior token.

144 **Architecture and Objective.** The image encoder maps
145 the pixels of I to a sequence of visual features and can
146 be instantiated as a pre-trained CNN (e.g., ResNet [8]) or
147 Vision Transformer [5]. A Transformer decoder predicts
148 the polygon sequence $\mathbf{y} = (y_1, \dots, y_T)$ autoregressively.
149 Task structure is encoded by the sequence itself rather than
150 specialized prediction heads. Conditioned on the encoded
151 image features, the decoder models

$$152 P(\mathbf{y} | I) = \prod_{t=1}^T P(y_t | \mathbf{y}_{<t}, I), \quad (2)$$

153 and is trained with teacher forcing under the weighted cross-
154 entropy objective (optionally with label smoothing):

$$155 \mathcal{L} = \frac{\sum_{t=1}^T \mathbf{w}_t \text{CE}_t}{\sum_{t=1}^T \mathbf{w}_t} + \lambda_{\text{ae}} \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} P(\text{EOS} | \mathbf{y}_{<t}, I), \quad (3)$$

156 where $\text{CE}_t = -\log P(y_t | \mathbf{y}_{<t}, I)$, \mathcal{T} denotes all non-
157 padding, non-EOS positions, and \mathbf{w}_t is a token-specific
158 weight. We set the weights of coordinate, EOP, and class
159 tokens to 1.0, the weight of EOS to 0.5, and $\lambda_{\text{ae}} = 0.1$. The
160 lower EOS weight and anti-EOS term reduce premature ter-
161 mination by penalizing excessive EOS probability.

162 We also append synthetic noise polygons during training
163 and assign them a dedicated `noise` class token. Each noise
164 polygon is derived from a randomly selected ground-truth
165 ring by jittering and randomly dropping vertices. Their
166 coordinate and EOP labels are masked while the `noise` class
167 token remains supervised, as illustrated in Figure 2(b).

168 **Inference.** Prior autoregressive approaches either decode to
169 a fixed maximum length [3] or condition polygon generation
170 on per-instance prompts [4]. We instead decode at the scene
171 level: one trial produces a polygon set for the full image,
172 and the final prediction repeats this prompt-free decoding K
173 times with nucleus sampling [10], as shown in Figure 2(c).
174 Each trial is variable-length and stops at EOS, optionally with
175 a small EOS-logit penalty to reduce premature termination.
176 The decoding budget is therefore scene-adaptive, which suits
177 aerial tiles that range from a few sparse buildings to hundreds
178 in dense urban scenes, and the number of trials is indepen-
179 dent of the number of objects in the tile. All experiments
180 use $K = 6$ trials, top- p sampling with $p = 0.4$, and the
181 same decode-time EOS setting across datasets. Raw trial
182 sequences are converted to polygons without per-trial NMS
183 and then passed to a shared merge-and-filter pipeline.

184 We add a Gaussian spatial bias to coordinate logits during
185 decoding. For a coordinate bin k , the biased logit is

$$186 \tilde{\ell}_k = \ell_k + \alpha \exp\left(-\frac{(k - \hat{k})^2}{2\sigma^2}\right), \quad (4)$$

	INRIA (155) val										
	Method	Encoder	Decoder	GPU Hrs	IoU (%) \uparrow	C-IoU (%) \uparrow	MTA ($^\circ$) \downarrow	PoLiS \downarrow	IoU ^{topo} (%) \uparrow	F1 ^{topo} (%) \uparrow	PA ^{topo} (%) \uparrow
	FFL [CVPR'21] [7]	-	-	-	68.30	49.80	35.62	2.865	43.38	58.78	89.67
	HiSup [ISPRS'23] [18]	-	-	-	74.90	66.10	43.86	2.438	53.51	67.94	93.20
	Pix2Poly [WACV'25] [1]	ViT-S/8	6L-256d	-	79.46	71.73	34.31	1.914	61.08	74.29	94.37
	Aerial2Poly (Ours)	ViT-S/8	6L-256d	12.93	85.02	78.55	34.74	1.944	72.75	81.22	96.32
	SpaceNet Vegas val										
	Method	Encoder	Decoder	GPU Hrs	IoU (%) \uparrow	C-IoU (%) \uparrow	MTA ($^\circ$) \downarrow	PoLiS \downarrow	IoU ^{topo} (%) \uparrow	F1 ^{topo} (%) \uparrow	PA ^{topo} (%) \uparrow
	FFL [CVPR'21] [7]	-	-	-	76.00	57.60	36.29	2.398	49.46	65.00	91.10
	HiSup [ISPRS'23] [18]	-	-	-	82.10	75.20	33.89	1.722	59.56	73.43	93.80
	Pix2Poly [WACV'25] [1]	ViT-S/8	6L-256d	-	81.81	75.05	33.40	1.717	60.31	74.20	93.80
	Aerial2Poly (Ours)	ViT-S/8	6L-256d	13.75	84.32	77.56	33.10	1.751	67.40	77.81	95.63
	WHU Buildings Test										
	Method	Encoder	Decoder	GPU Hrs	IoU (%) \uparrow	C-IoU (%) \uparrow	MTA ($^\circ$) \downarrow	PoLiS \downarrow	IoU ^{topo} (%) \uparrow	F1 ^{topo} (%) \uparrow	PA ^{topo} (%) \uparrow
	FFL [CVPR'21] [7]	-	-	-	77.64	54.52	35.79	1.747	56.56	70.01	94.02
	HiSup [ISPRS'23] [18]	-	-	-	87.12	79.62	34.75	1.158	72.11	82.47	96.80
	Pix2Poly [WACV'25] [1]	ViT-S/8	6L-256d	-	89.15	81.63	31.64	1.082	75.38	84.96	97.14
	Aerial2Poly (Ours)	ViT-S/8	6L-256d	12.85	93.79	90.09	31.58	0.993	87.17	91.68	98.65

Table 1. Main results on building footprint prediction. GPU hours report total training compute to 240k iterations on 8 NVIDIA L4 GPUs as a proxy for training cost. Results for methods not proposed in this work are quoted from Pix2Poly.

	INRIA (155) test				
	Method	Encoder	Decoder	IoU (%) \uparrow	Acc (%) \uparrow
	Mask2Poly [22]	-	-	74.40	96.10
	FFL [7]	-	-	74.80	95.96
	HiSup [18]	-	-	75.53	96.27
	Pix2Poly [1]	ViT-S/8	6L-256d	75.87	96.37
	Aerial2Poly (Ours)	ViT-S/8	6L-256d	76.88	96.53

Table 2. Official INRIA test-set results from the online evaluator. Decoder is reported as depth-hidden size. Results for methods not proposed in this work are quoted from Pix2Poly.

where \hat{k} is the most recently decoded coordinate bin on the same axis within the current polygon: the previous x -bin when the next token is x , or the previous y -bin when the next token is y . The bias is inactive for the first coordinate on each axis and reset after every EOP. In all experiments, we set the spatial bias strength to $\alpha = 1.5$ and $\sigma = 40$ coordinate bins, equivalent to pixels at our resolution. Applied only to coordinate tokens, the bias encourages local boundary coherence without altering separator or class prediction.

After discarding polygons predicted as noise along with very small or low-confidence ones, we consolidate the K trial outputs per class in vector space. The trials are split into two groups, an IoU-medoid is selected per overlap cluster within each group, and the survivors are merged across groups using a shared polygon-IoU threshold of 0.45. Interior polygons are consolidated as a separate class and subtracted from exterior polygons only when fully contained.

4. Experiments

Experimental Setup. We evaluate Aerial2Poly under the shared protocol for polygonal extraction following previous work [1, 2, 18, 23]. For a direct comparison with the Pix2Poly baseline, we use the same ViT-S/8 image encoder at 224×224 resolution and a Transformer decoder of match-

ing scale with 6 layers and hidden size 256. We train on the INRIA, SpaceNet Vegas, and WHU Buildings benchmarks following the same data splits and patchification procedure as HiSup and Pix2Poly. Unless noted otherwise, our model and training hyperparameters are fixed across datasets. We also fix the inference recipe across datasets: $K = 6$ trials, top- p sampling with $p = 0.4$, Gaussian bias with spatial bias strength $\alpha = 1.5$ and $\sigma = 40$ coordinate bins, raw-trial polygon extraction without per-trial NMS, and vector-space trial consolidation at polygon-IoU threshold 0.45, followed by the same final cross-class and area filters.

Metrics. Following the established evaluation protocol above, we report polygon quality using IoU, C-IoU (which penalizes vertex-count mismatch in addition to overlap), MTA, PoLiS, and the topology-aware IoU/F1/PA boundary metrics. For INRIA, we also report the official online test-set IoU and Accuracy metrics from the public evaluator.

4.1. Main Results

Table 1 shows that under the matched setting, Aerial2Poly outperforms Pix2Poly and competing methods on IoU, C-IoU, and the topology-aware IoU/F1/PA metrics across all three datasets. On WHU Buildings, C-IoU improves to 90.09%, yielding a +8.46 point absolute gain. This remarkable improvement suggests that our model recovers more geometrically aligned polygons rather than locally accurate but topologically broken outlines.

Official INRIA Test Set. Table 2 reports the pixel-level test scores returned by the INRIA online evaluator with undisclosed ground truth. Following the standard INRIA test-time protocol, we apply the model in a sliding-window fashion over each whole test image and combine the resulting pixel-level predictions into a test submission. Under this official evaluation, Aerial2Poly attains 76.88% IoU and 96.53% Accuracy, outperforming the best test results of Pix2Poly and

244 other alternative methods.

245 5. Next Steps

246 This work introduced Aerial2Poly, a simple and scalable
247 formulation for direct building footprint vectorization from
248 aerial imagery. By predicting the full-tile polygon set directly
249 in vector space, Aerial2Poly avoids mask-to-polygon recov-
250 ery, object prompts, and auxiliary connectivity heads. Cou-
251 pled with Gaussian-biased decoding, recall-oriented training,
252 and vector-space trial consolidation, our approach estab-
253 lishes new state-of-the-art results on the INRIA, SpaceNet
254 Vegas, and WHU Buildings benchmarks.

255 While this workshop paper focuses on building footprint
256 extraction, the broader takeaway is that full-tile polygon
257 generation is a viable interface for aerial perception. We
258 hope these results encourage future work on richer label
259 spaces, larger aerial benchmarks, and more general vector-
260 native prediction from imagery.

261 References

- 262 [1] Yeshwanth Kumar Adimoolam, Charalambos Poullis, and
263 Melinos Averkiou. Pix2Poly: A Sequence Prediction Method
264 for End-to-End Polygonal Building Footprint Extraction from
265 Remote Sensing Imagery. In *WACV*, 2025. 1, 2, 4
- 266 [2] Janja Avbelj, Rupert Müller, and Richard Bamler. A Metric
267 for Polygon Comparison and Building Extraction Evaluation.
268 *IEEE Geoscience and Remote Sensing Letters*, 2015. 4
- 269 [3] Ting Chen, Saurabh Saxena, Lala Li, David J. Fleet, and
270 Geoffrey Hinton. Pix2Seq: A Language Modeling Framework
271 for Object Detection. In *ICLR*, 2022. 2, 3
- 272 [4] Ting Chen, Saurabh Saxena, Lala Li, Tsung-Yi Lin, David J.
273 Fleet, and Geoffrey E. Hinton. A Unified Sequence Interface
274 for Vision Tasks. In *NeurIPS*, 2022. 2, 3
- 275 [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,
276 Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
277 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl-
278 vain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is
279 Worth 16x16 Words: Transformers for Image Recognition at
280 Scale. In *ICLR*, 2021. 3
- 281 [6] Adam Van Etten, Dave Lindenbaum, and Todd M. Bacastow.
282 SpaceNet: A Remote Sensing Dataset and Challenge Series.
283 <https://arxiv.org/abs/1807.01232>, 2018. 1
- 284 [7] Nicolas Girard, Dmitriy Smirnov, Justin Solomon, and Yuliya
285 Tarabalka. Polygonal Building Segmentation by Frame Field
286 Learning. In *CVPR*, 2021. 2, 4
- 287 [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
288 Deep Residual Learning for Image Recognition. In *CVPR*,
289 2016. 3
- 290 [9] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Gir-
291 shick. Mask R-CNN. In *ICCV*, 2017. 2
- 292 [10] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin
293 Choi. The Curious Case of Neural Text Degeneration. In
294 *ICLR*, 2020. 3
- 295 [11] Shunping Ji, Shiqing Wei, and Meng Lu. Fully Convolutional
296 Networks for Multisource Building Extraction From an Open

- Aerial and Satellite Imagery Data Set. *IEEE Transactions on
Geoscience and Remote Sensing*, 2019. 1 297
298
- [12] Maxim Khomiakov, Michael Riis Andersen, and Jes Frelsen.
GeoFormer: A Multi-Polygon Segmentation Transformer. In
BMVC, 2024. 2 299
300
301
- [13] Muxingzi Li, Florent Lafarge, and Renaud Marlet. Approxi-
mating Shapes in Images with Low-Complexity Polygons. In
CVPR, 2020. 2 302
303
304
- [14] Weijia Li, Wenqian Zhao, Huaping Zhong, Conghui He, and
Dahua Lin. Joint Semantic-Geometric Learning for Polygonal
Building Segmentation. In *AAAI*, 2021. 2 305
306
307
- [15] Zuoyue Li, Jan Dirk Wegner, and Aurélien Lucchi. Topo-
logical Map Extraction from Overhead Images. In *ICCV*,
2019. 308
309
310
- [16] Huan Ling, Jun Gao, Amlan Kar, Wenzheng Chen, and Sanja
Fidler. Fast Interactive Object Annotation with Curve-GCN.
In *CVPR*, 2019. 2 311
312
313
- [17] Emmanuel Maggiori, Yuliya Tarabalka, Guillaume Charpiat,
and Pierre Alliez. Can Semantic Labeling Methods General-
ize to Any City? The Inria Aerial Image Labeling Benchmark.
In *IEEE International Geoscience and Remote Sensing Sym-
posium (IGARSS)*, 2017. 1 314
315
316
317
318
- [18] Bowen Xu, Jiakun Xu, Nan Xue, and Gui-Song Xia. HiSup:
Accurate Polygonal Mapping of Buildings in Satellite Im-
agery with Hierarchical Supervision. *ISPRS Journal of Pho-
togrammetry and Remote Sensing*, 2023. 2, 4 319
320
321
322
- [19] Bingnan Yang, Mi Zhang, Zhan Zhang, Zhili Zhang, and Xi-
angyun Hu. TopDiG: Class-Agnostic Topological Directional
Graph Extraction from Remote Sensing Images. In *CVPR*,
2023. 2 323
324
325
326
- [20] Mingming Zhang, Qingjie Liu, and Yunhong Wang. HiT:
Building Mapping With Hierarchical Transformers. *IEEE
Transactions on Geoscience and Remote Sensing*, 2024. 2 327
328
329
- [21] Kang Zhao, Jungwon Kang, Jaewook Jung, and Gunho Sohn.
Building Extraction from Satellite Images Using Mask R-
CNN with Building Boundary Regularization. In *CVPR
Workshops*, 2018. 2 330
331
332
333
- [22] Stefano Zorzi, Ksenia Bittner, and Friedrich Fraundorfer.
Machine-Learned Regularization and Polygonization
of Building Segmentation Masks. In *ICPR*, 2020. 2, 4 334
335
336
- [23] Stefano Zorzi, Shabab Bazrafkan, Stefan Habenschuss, and
Friedrich Fraundorfer. PolyWorld: Polygonal Building Ex-
traction with Graph Neural Networks in Satellite Images. In
CVPR, 2022. 2, 4 337
338
339
340