# Building Footprint Segmentation from Satellite Imagery Using Convolutional and Transformer-Based Models

## Abstract

Accurate building footprint extraction from satellite imagery is essential for applications such as disaster response and urban planning, yet manual annotation remains costly and time-consuming. In this work, we evaluate automated building segmentation on the SpaceNet-2 Las Vegas dataset using two deep learning approaches: a convolutional baseline based on a ResNet-50 encoder and a transformer-based Seg-Former model. The ResNet-50 model leverages residual connections and a fully convolutional segmentation head to produce dense pixel-level predictions, achieving high recall but exhibiting over-segmentation in dense urban regions. SegFormer, which replaces convolutional feature extraction with transformer encoders, captures global spatial context across the image and produces cleaner boundaries with improved precision. Experimental results show that while ResNet-50 provides strong baseline performance, SegFormer achieves higher overall segmentation quality, highlighting the benefits of global context modeling for high-resolution urban imagery.

## 1. Introduction

Infrastructure data plays a critical role in urban planning, disaster response, and economic development (Taubenböck et al., 2012; Kuffer et al., 2016). However, many regions lack accurate and up-to-date information on core infrastructure, such as buildings, limiting effective planning, resource allocation, and recovery efforts (Kuffer et al., 2016). Traditional mapping methods rely heavily on manual surveys and annotation, which are costly, slow, and difficult to scale, particularly in rapidly changing or data-scarce regions (Taubenböck et al., 2012). At the same time, large volumes of high-resolution satellite imagery are continu-

. Correspondence to: Anonymous Author <anon.email@domain.com>.

ously collected but remain underutilized due to the difficulty of extracting structured information from them at scale (Kuffer et al., 2016).

Following the February 6, 2023, earthquakes in Turkey and Syria, emergency response teams relied heavily on satellite imagery to assess building damage; however, the lack of reliable pre-event building footprint maps significantly slowed the identification of affected structures and the estimation of losses (UNOSAT, 2023; World Bank, 2023). This gap highlights the need for automated, scalable approaches that can reliably extract infrastructure information directly from satellite imagery.

In this project, we focus on building footprint detection from high-resolution satellite images using supervised deep learning methods. Specifically, we evaluate convolutional and transformer-based semantic segmentation models on the SpaceNet2 dataset, which provides large-scale, well-labeled satellite imagery for building detection tasks. By framing building footprint extraction as a pixel-level segmentation problem, these models can learn both fine-grained local features, such as rooftop boundaries, and broader spatial patterns, such as neighborhood structure.

Convolutional Neural Networks (CNNs) have been widely used for building footprint segmentation due to their strong performance in visual feature extraction and spatial generalization from high-resolution remote sensing imagery (Zhu et al., 2017a; Audebert et al., 2017). By leveraging local receptive fields and hierarchical feature representations, CNN-based architectures are effective at capturing fine-grained urban structures such as building edges and rooftops across diverse environments.

To address this, we conduct a comparative evaluation of a ResNet-based fully convolutional network (ResNet-FCN) and a SegFormer model on SpaceNet2 building footprint data. We assess performance using standard segmentation metrics, including Precision, Recall, F1 score, and Intersection-over-Union (IOU), and analyze qualitative differences in model predictions. Through this comparison, our goal is to better understand the trade-offs between convolutional and transformer-based approaches for automated infrastructure mapping from satellite imagery.

## 2. Related Work

Building footprint extraction from high-resolution satellite imagery is a well-studied problem in remote sensing and computer vision, with applications in urban planning, disaster response, and infrastructure monitoring (UNOSAT, 2023; Ye et al., 2025). Early approaches relied on hand-crafted features and traditional image processing techniques, including rule-based segmentation and morphological operations (Pesaresi & Benediktsson, 2001; Mayer, 1999). More recent advances have demonstrated that deep learning methods significantly outperform these classical approaches on large-scale remote sensing benchmarks by learning hierarchical and spatially robust feature representations directly from data (Tsagkatakis et al., 2019; Ma et al., 2019).

Convolutional neural networks (CNNs) have become the dominant approach for building footprint segmentation due to their ability to learn hierarchical spatial features directly from image pixels (Zhu et al., 2017a; Ma et al., 2019). Encoder–decoder architectures and fully convolutional networks have demonstrated strong performance in capturing both fine-grained boundary details and broader spatial context in high-resolution remote sensing imagery (Badrinarayanan et al., 2017; Long et al., 2015). Prior studies report that CNN-based models trained on high-resolution satellite imagery can generalize effectively across diverse urban environments when sufficient labeled data is available (Audebert et al., 2017; Ye et al., 2025).

More recently, transformer-based architectures have been introduced for semantic segmentation, largely because self-attention allows models to capture long-range spatial relationships that are difficult for purely convolutional networks to model (Vaswani et al., 2017; Dosovitskiy et al., 2021). SegFormer adapts vision transformers for segmentation by using a hierarchical transformer encoder together with a lightweight decoder, allowing information from multiple spatial scales to be combined efficiently without relying on positional encodings (Xie et al., 2021).

Transformer-based segmentation models have shown strong results on standard natural image benchmarks, leading to increased interest in alternatives to traditional per-pixel classification approaches (Cheng et al., 2021). These architectures have also begun to appear in remote sensing applications, where access to global image context can be useful for interpreting large and complex urban scenes (Aleissaee et al., 2023; Noman et al., 2024).

However, comparisons between transformer-based models and established CNN baselines for building footprint extraction remain limited, and reported performance often depends strongly on dataset properties, image resolution, and training setup, with CNNs continuing to perform compet-

itively in many cases (Zhu et al., 2017b; Ma et al., 2019). In this work, we address this gap by comparing a ResNet-based fully convolutional network and a SegFormer model on the SpaceNet2 building footprint dataset using consistent preprocessing and evaluation metrics.

## 3. Dataset and Problem Setup

### 3.1. Dataset

We use the SpaceNet2 building footprint dataset, focusing on Area of Interest 2 (AOI), Las Vegas, which provides high-resolution satellite imagery with corresponding building footprint annotations. The dataset consists of pan-sharpened RGB satellite image tiles paired with polygon-based building labels provided as GeoJSON files. These define building footprints at pixel-level resolution, enabling supervised semantic segmentation.

Each satellite tile is associated with a corresponding set of building polygons, which we rasterize into binary segmentation masks aligned with the original imagery. Pixels belonging to building footprints are labeled as 1, while background pixels are labeled as 0. This formulation defines a binary building footprint segmentation task.

### 3.2. Data Preprocessing

The raw SpaceNet dataset is not organized in a format suitable for model training, as images and annotations are stored separately, and building labels are provided as vector polygons. To address this, we convert all building footprint polygons into pixel-aligned raster masks using the same spatial transform and coordinate reference system as the input images to ensure precise alignment between the images and labels.

All images are resized to a resolution of 256x256 pixels to enable sufficient batching during training. Input images are normalized as well to an intensity range to stabilize optimization. The segmentation masks are resized using nearest-neighbor interpolation to ensure that the binary label is preserved.

### 3.3. Train-Validation-Test Split

The dataset is split into training, validation, and test sets using an 80/10/10 ratio. The split is performed by randomly shuffling image-mask pairs with a fixed random seed. The training set is used for model optimization, the validation set for hyperparameter tuning and early stopping, and the test set for final evaluation.

### 3.4. Data Loading

To effectively feed data into models during training and evaluation, we implement a custom PyTorch dataset that loads paired-image mask samples and applies preprocessing transforms. We use PyTorch DataLoaders with mini-batch training, shuffling the training set each epoch while keeping validation and test sets fixed.
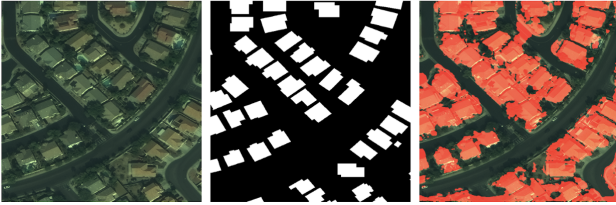


*Figure 1.* Overview of the SpaceNet2 dataset, showing example satellite image tiles and corresponding rasterized building footprint masks.

## 4. Methods (ResNet + SegFormer)

### 4.1. ResNet-FCN Baseline

As a convolutional baseline, we implement a deep convolutional network with a ResNet-50 backbone for binary building segmentation. ResNet-50 is composed of stacked residual blocks with identity skip connections, allowing the network to learn residual functions of the form F(x)+x. These residual connections improve gradient flow during backpropagation and enable stable optimization of deeper convolutional architectures.

In our implementation, ResNet-50 serves as the encoder, extracting hierarchical spatial features from the input RGB satellite tiles. The FCN segmentation head upsamples the encoder's final feature maps to produce dense, per-pixel predictions at the original image resolution. Since building footprint extraction is a binary task, the model outputs a single logit channel representing building versus background for each pixel.

Training is performed for five epochs using mini-batch gradient descent. During each iteration, the model processes batches of image, mask pairs, computes pixel-wise loss, and updates weights via backpropagation. Validation loss is computed after training to assess generalization performance and detect overfitting.

During inference, the ResNet-FCN model outputs raw per-pixel logits, which are passed through a sigmoid activation function to obtain probability maps. These probabilities are thresholded at 0.5 to generate binary segmentation masks aligned with the ground-truth building labels. The resulting predictions are used to compute pixel-level precision, recall, and F1 score on the held-out test set.

### 4.2. SegFormer

We use a pre-trained SegFormer encoder (MiT backbone) initialized with ImageNet weights. The encoder processes the input image at multiple resolutions, producing hierarchical feature maps that capture both local texture and global context. These features are fused by the SegFormer decoder, which outputs a dense per-pixel prediction map at one quarter of the input spatial resolution.

Due to SegFormer producing lower-resolution logits, we upsampled the output using bilinear interpolation to match the original image resolution before computing loss and evaluation metrics.

SegFormer was trained for five epochs using the Adam optimizer with a fixed learning rate. The model was optimized using binary cross-entropy loss with logits. Training was performed in mini-batches, and gradients were computed via backpropagation. Validation loss was evaluated after each epoch to determine convergence and detect overfitting.

During interference, the model outputs raw logits representing per pixel building likelihoods. A sigmoid activation function was applied to convert logits into probabilities. These probability maps were then thresholded at 0.5 to produce binary segmentation masks. The predicted masks were aligned with the ground truth labels by resizing outputs to the original tile resolution.

### 4.3. Evaluation Metrics

Model performance for both Res-Net-FCN and SegFormer was evaluated using pixel-level precision, recall, and F1 score. True positives, false positives and false negatives were computed by comparing predicted binary masks against ground-truth building masks on a per-pixel basis. These metrics quantify the difference between correctly identifying building pixels and avoiding false detections, thus providing a comprehensive assessment of segmentation quality.

### 4.4. Hypothesis

Motivated by prior work on transformer-based semantic segmentation Xie et al., 2021), we hypothesize that SegFormer will outperform the ResNet-FCN baseline in Intersection-over-Union (IOU) and F1 score due to its ability to capture long-range spatial dependencies through self-attention.

# 5. Experiment and Results

## 5.1. ResNet

Over five training epochs, the ResNet-50 FCN baseline showed clear convergence, with training loss decreasing from 0.2480 → 0.1142. As shown in Figure 2 (a), validation loss generally decreased early and reached its best value at Epoch 3 (Val Loss = 0.1608), then slightly increased afterward (Epochs 4–5), suggesting the start of overfitting or diminishing returns beyond the best checkpoint.

On the test set, ResNet-50 achieved Precision = 0.7864, Recall = 0.8667, F1 = 0.8246, and IoU = 0.7016, with pixel counts TP = 4,122,573, FP = 1,119,528, and FN = 633,864. The relatively high recall indicates the model successfully captures most building pixels (few missed building regions), while the lower precision reflects a noticeable number of false positives, consistent with a tendency to slightly over-segment or "spill" into background regions.
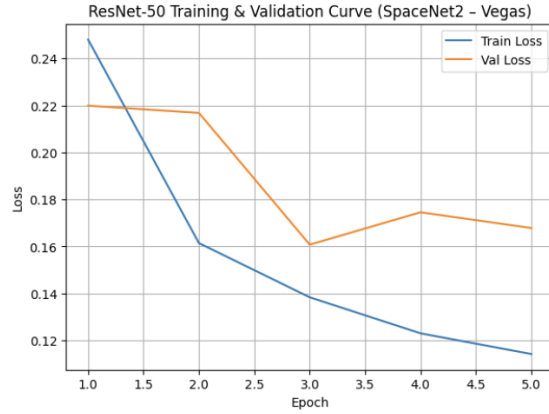
Qualitatively, Figure 2 (b) shows that ResNet predictions recover the overall building layout well, but building boundaries appear less crisp and can merge nearby structures or thicken edges, which aligns with the observed false-positive behavior and the lower IoU compared to an ideal boundary-aligned mask.
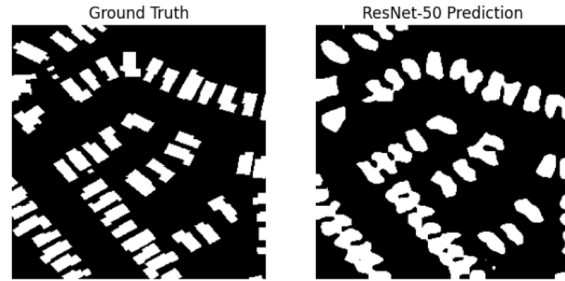
## 5.2. Segformer

Over 5 training epochs, Segformer displayed stable convergence, with training loss decreasing from 0.2305 to 0.1512. As shown in Figure 3 (a), training and validation loss followed similar downward trends, indicating effective learning without overfitting.

Quantitatively, SegFormer achieved a precision of 0.8571, a recall of 0.8137, and an F1 score of 0.8348, with an IOU of 0.7164. The relatively high precision indicates that the model produced few false positive building predictions, while the slightly lower recall suggests that smaller or irregularly shaped buildings were more likely to be missed. Overall, the F1 score reflects a strong balance between detection accuracy and coverage.

Figure 3 (b) displays a qualitative comparison between ground truth building masks and SegFormer predictions. The model produces clean and coherent building boundaries, particularly in dense urban regions where structures are closely packed. While most large and well-defined buildings are accurately segmented, some fine-grained details and small structures are missed, consistent with the observed recall behavior.
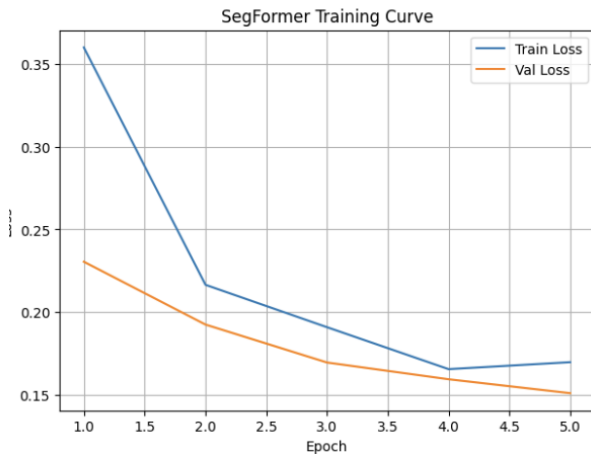


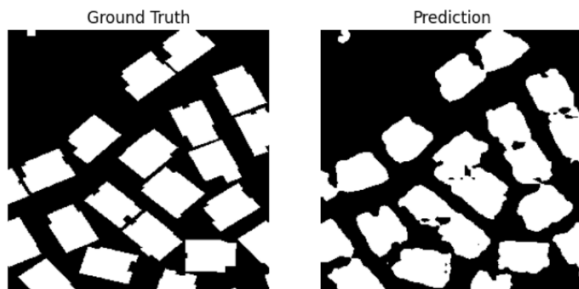*(a)* Training and validation loss curves for ResNet-50 FCN.



*(b)* Qualitative ResNet-FCN predictions vs. ground truth.

*Figure 2.* ResNet-50 FCN quantitative and qualitative results on SpaceNet2.

4

*(a)* Training and validation loss curves for SegFormer over five epochs.



*(b)* Qualitative SegFormer predictions compared to ground truth building footprint masks.

*Figure 3.* SegFormer quantitative and qualitative segmentation results on the SpaceNet2 dataset.

## 6. Conclusion

This project evaluated two deep learning approaches for building footprint segmentation on the SpaceNet-2 Las Vegas dataset: a ResNet-50 FCN convolutional baseline and a transformer-based SegFormer model. Both models successfully learned to segment buildings from high-resolution RGB satellite tiles, demonstrating that supervised semantic segmentation is a viable and scalable alternative to manual building map creation.

Overall, the ResNet-50 FCN provided a strong baseline with high recall (0.8667), indicating it detected most building regions, but its lower precision (0.7864) reflected a tendency to over-segment and blur boundaries in dense areas. In contrast, SegFormer achieved higher precision (0.8571) and a slightly higher F1 score (0.8348 vs. 0.8246) and IoU (0.7164 vs. 0.7016), producing cleaner and more coherent building boundaries by leveraging global context through self-attention. A Wilcoxon signed-rank test comparing per-tile performance confirmed that this difference isstatisti-

cally significant ($p \approx 3 \times 10^{-33}$), indicating that Seg-Former consistently outperformed the ResNet-FCN baseline across the test set. These results support our hypothesis that transformer-based segmentation can outperform CNN baselines on complex urban imagery, particularly when boundary quality and false positive control matter.

Future work could improve both approaches by incorporating stronger loss functions for class imbalance, applying data augmentation, training for additional epochs with early stopping, and evaluating generalization across multiple SpaceNet AOIs to test robustness beyond Las Vegas.

## References

Aleissaee, A. A., Kumar, A., Anwer, R. M., Khan, S., Cholakkal, H., Xia, G.-S., and Khan, F. S. Transformers in remote sensing: A survey. *Remote Sensing*, 15(7): 1860, 2023. doi: 10.3390/rs15071860. URL `https://www.mdpi.com/2072-4292/15/7/1860`.

Audebert, N., Le Saux, B., and Lefèvre, S. Beyond rgb: Very high resolution urban remote sensing with multimodal deep networks. *arXiv preprint arXiv:1711.08681*, 2017. URL `https://arxiv.org/abs/1711.08681`.

Badrinarayanan, V., Kendall, A., and Cipolla, R. Segnet: A deep convolutional encoder–decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(12):2481–2495, 2017. doi: 10.1109/TPAMI.2016.2644615. URL `https://ieeexplore.ieee.org/document/7803544`.

Cheng, B., Schwing, A. G., and Kirillov, A. Per-pixel classification is not all you need for semantic segmentation. *arXiv preprint arXiv:2107.06278*, 2021. doi: 10.48550/arXiv.2107.06278. URL `https://arxiv.org/abs/2107.06278`.

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2021. doi: 10.48550/arXiv.2010.11929. URL `https://arxiv.org/abs/2010.11929`.

Kuffer, M., Pfeffer, K., and Sliuzas, R. Slums from space—15 years of slum mapping using remote sensing. *Remote Sensing*, 8(6):455, 2016. doi: 10.3390/rs8060455.

Long, J., Shelhamer, E., and Darrell, T. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3431–3440, 2015.

doi: 10.1109/CVPR.2015.7298965. URL https://ieeexplore.ieee.org/document/7298965.

Ma, L., Liu, Y., Zhang, X., Ye, Y., Yin, G., and Johnson, B. Deep learning in remote sensing applications: A meta-analysis and review. *ISPRS Journal of Photogrammetry and Remote Sensing*, 152:166–177, 2019. doi: 10.1016/j.isprsjprs.2019.04.015.

Mayer, H. Automatic object extraction from aerial imagery—a survey. *ISPRS Journal of Photogrammetry and Remote Sensing*, 54(2–3):95–106, 1999. doi: 10.1016/S0924-2716(99)00006-0.

Noman, M., Fiaz, M., Cholakkal, H., Narayan, S., et al. Remote sensing change detection with transformers trained from scratch. *IEEE Transactions on Geoscience and Remote Sensing*, 2024. doi: 10.1109/TGRS.2024.3383800.

Pesaresi, M. and Benediktsson, J. A. A new approach for the morphological segmentation of high-resolution satellite imagery. *IEEE Transactions on Geoscience and Remote Sensing*, 39(2):309–320, 2001. doi: 10.1109/36.905239.

Taubenböck, H., Esch, T., Felbier, A., Wiesner, M., Roth, A., and Dech, S. Monitoring urbanization in mega cities from space. *Remote Sensing of Environment*, 117:162–176, 2012. doi: 10.1016/j.rse.2011.09.015.

Tsagkatakis, G. et al. Survey of deep-learning approaches for remote sensing observation enhancement. *Sensors*, September 2019. URL https://pmc.ncbi.nlm.nih.gov/articles/PMC6767260/.

UNOSAT. Earthquake damage assessment: Türkiye and syria, 2023. URL https://unosat.org. Satellite-based damage assessment following the February 6, 2023 earthquakes.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017. doi: 10.48550/arXiv.1706.03762. URL https://arxiv.org/abs/1706.03762.

World Bank. Türkiye earthquakes 2023: Rapid damage and needs assessment, 2023. URL https://www.worldbank.org. Infrastructure damage and economic loss assessment.

Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., and Luo, P. Segformer: Simple and efficient design for semantic segmentation with transformers. *arXiv preprint arXiv:2105.15203*, 2021. doi: 10.48550/arXiv.2105.15203. URL https://arxiv.org/abs/2105.15203.

Ye, M. et al. A deep learning pipeline to power infrastructure detection in high-resolution satellite images. *International Journal of Digital Earth*, April 2025. URL https://www.tandfonline.com/doi/full/10.1080/20964471.2025.2490408.

Zhu, X. X., Tuia, D., Mou, L., Xia, G.-S., Zhang, L., Xu, F., and Fraundorfer, F. Deep learning in remote sensing: A comprehensive review. *IEEE Geoscience and Remote Sensing Magazine*, 5(4):8–36, 2017a. doi: 10.1109/MGRS.2017.2762307.

Zhu, X. X., Tuia, D., Mou, L., Xia, G.-S., Zhang, L., Xu, F., and Fraundorfer, F. Deep learning in remote sensing: A comprehensive review and list of resources. *IEEE Geoscience and Remote Sensing Magazine*, 5(4):8–36, 2017b. doi: 10.1109/MGRS.2017.2762307.