

ESTIMATING DISPLACED POPULATIONS FROM OVERHEAD

Armin Hadzic^{1,2}, Gordon Christie¹, Jeffrey Freeman¹, Amber Dismer³, Stevan Bullard⁴
Ashley Greiner³, Nathan Jacobs², Ryan Mukherjee¹

¹Johns Hopkins University Applied Physics Laboratory ²University of Kentucky
³Centers for Disease Control and Prevention ⁴Agency for Toxic Substances and Disease Registry

ABSTRACT

We introduce a deep learning approach to perform fine-grained population estimation for displacement camps using high-resolution overhead imagery. We train and evaluate our approach on drone imagery cross-referenced with population data for refugee camps in Cox’s Bazar, Bangladesh in 2018 and 2019. Our proposed approach achieves 7.02% mean absolute percent error on sequestered camp imagery. We believe our experiments with real-world displacement camp data constitute an important step towards the development of tools that enable the humanitarian community to effectively and rapidly respond to the global displacement crisis.

Index Terms— Deep Learning, Machine Learning, Regression, CNN, Population Estimation, Remote Sensing

1. INTRODUCTION

According to the United Nations High Commissioner for Refugees (UNHCR), there are currently 70.8 million people forcibly displaced worldwide [1], which is the largest human displacement crisis in history. Many displaced persons find themselves living in camps where health and other basic human services are constrained and the threat of violence a persistent concern. The humanitarian community has mobilized significant resources to address the displacement crisis. However, responses are inherently reactionary and there is a strong desire for tools that enable a rapid and accurate assessment of populations in need during times of crisis. Deep learning with overhead imagery offers a practical solution to improve response. In this paper, we detail an approach for estimating camp population using overhead imagery, which can often be obtained rapidly and throughout an event. Population estimates provide a good starting point because they can feed into many different types of analysis, such as determining whether sufficient water, sanitation, and hygiene (WASH) facilities are present in a camp.

While population estimation using overhead imagery is not novel, high-resolution (i.e., 30cm ground sample distance and better) estimates have only recently become possible due to improved sensor capabilities. Given the increasing availability of such imagery, along with the recent advancements

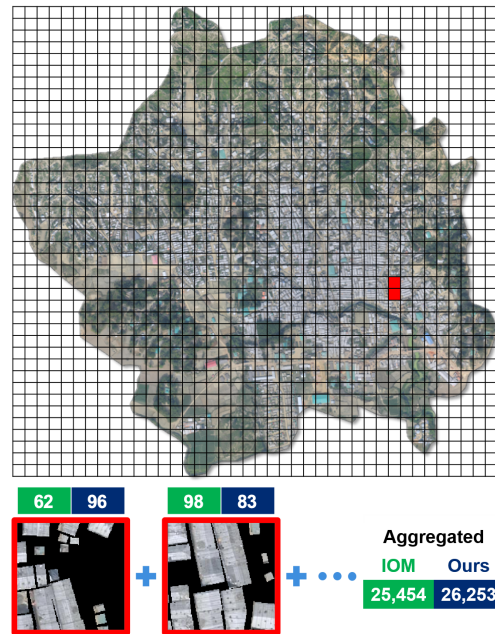


Fig. 1: Generating test images (top) using equally spaced and sized image chips. Population predictions (blue) and IOM-reported population labels (green) are shown per chip. The predicted population of the entire camp (blue, bottom right) is the sum of all of the image chip predictions.

in computer vision and machine learning, we believe the time is ripe to start training and deploying high-resolution models for the humanitarian community. Most previous work estimate population at a substantially coarser scale, leveraging approaches unlikely to transfer well to high-resolution camp imagery. [2] disaggregate US census tract data into fine-grain cells using weighted combinations of intersecting census blocks, similar to our disaggregation approach. However, [2] also discretize population counts into 17 bins to perform classification and operate on coarse 15m Landsat imagery. [3] perform population estimation using imagery of India, including rural areas that may be more visually similar to humanitarian camps, however they also use coarse resolution Landsat and Sentinel data and their approach struggles with fine-grain village-level ($\leq 20.25\text{km}^2$ area) population estimation. [4] introduce techniques for performing

pixel-level population estimation with 3m resolution Planet imagery. However, these techniques are only applied to well-organized US cities and it is unclear how well they might handle less-organized camp settings where the density and appearance of structures are drastically different.

Related approaches using sub-meter resolution imagery include [5] detecting impoverished settlements and [6] counting dwellings in Darfur camps, both using 30-50cm DigitalGlobe imagery but neither estimating population. On the other hand, [7] use structure areas extracted from Quickbird imagery to predict population, but their method is mostly manual and lacks large-scale evaluation.

To our knowledge, our approach is the first to perform learning-based population estimation using sub-meter overhead imagery (10cm GSD). We have publicly released our code to train and test population estimation models, as well as generate a dataset based on open source imagery and population data¹. Together, we believe these contributions constitute an important step towards the development of tools that enable the humanitarian community to effectively and rapidly respond to the global displacement crisis.

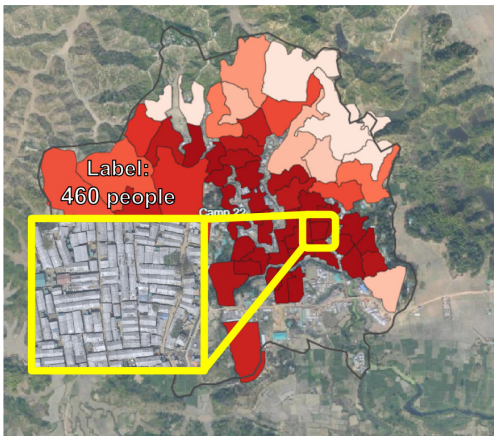


Fig. 2: Majhee block (population polygon) labels from Camp 22 in Cox’s Bazar, Bangladesh. Combining these labels produces an IOM-reported population label for an entire camp.

2. PROBLEM STATEMENT

Population data is traditionally provided in the form of census tracts, which vary in size and cover large spatial areas (Figure 2). However, fine-grain population mapping is desirable and beneficial for humanitarian efforts such as camp management. Pixel-level population labels are generally unavailable due to security concerns and annotation costs. We address this problem by using a simple area-based tract disaggregation method. Then to reduce the cost of conducting camp censuses altogether, we train a model using aerial imagery to directly predict camp population at high spatial resolution.

¹<https://github.com/JHUAPL/EstimatingDisplacedPopulations>

3. DATASET

Our dataset is comprised of (1) overhead drone imagery, (2) population polygons (majhee blocks), and (3) OpenStreetMap (OSM) structure segmentation masks, all from refugee camps in Cox’s Bazar, Bangladesh. Overhead images were tiled into square chips and paired with labels extracted from the population polygons. Population polygons correspond to majhee block data taken from routine International Organization for Migration (IOM) Bangladesh: Needs and Population Monitoring (NPM) site assessments. The data corresponding to 10% of the 34 camps were sequestered to the test split.

3.1. Overhead Imagery

Overhead imagery for our dataset was also sourced from NPM site assessments [8]. The imagery is comprised of georeferenced 10cm overhead drone images. Each of the 34 camps in the Cox’s Bazar region have up to nine overhead images, each averaging a 2.7km² area, and totaling 294 images. The nine possible images per camp correspond to different site assessments performed during different months and seasons across two years.

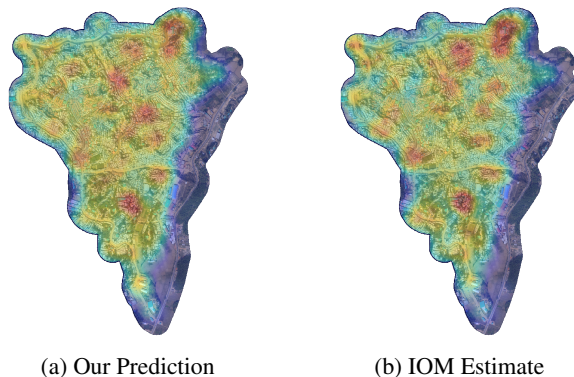


Fig. 3: Qualitative comparison of model predictions and IOM estimates from March 2019 for Camp 11 in Cox’s Bazar, Bangladesh. Red indicates relatively high population density, while blue indicates relatively low population density.

For each image in the training split, 200 square bounding boxes were extracted, stored as image chips, and used to train a convolutional neural network (CNN). Bounding boxes were both randomly positioned and sized, where a box’s edge length l is randomly chosen from the set $l \in \{224, 320, 480, 640, 720, 1024\}$ without exceeding the source image size. Randomly selecting bounding box locations and sizes provides several benefits, such as increased quantity of training data and improved generalizability by training on imagery with slightly varying effective spatial resolutions and population densities. Input image chips were then resized, not cropped, from their original bounding box size down to 224x224 before being passed through the model.

Conversely, test split image chips are all of size 224x224 and were sampled in a sliding-window fashion to ensure overhead test images were fully covered with no overlap. As Figure 1 illustrates, the sliding-window chip generation approach allows a model to make a prediction on each image chip such that the sum of all predictions reflects the total displacement camp population estimate.

OSM’s human-annotated structure masks were used to focus the model’s feature extraction by removing non-relevant image pixels and help reduce scene appearance bias for when the model is applied to new scenes. We accomplished this by using OSMNX [9] to retrieve and register structure annotations for each image chip (Figure 1). Note that this registration process does result in some misalignment. Test chips containing structures with 0 corresponding population were omitted during test time as the buildings were outside Majhee block boundaries.

3.2. Population Polygons

Population estimates for refugee camps in Cox’s Bazar, Bangladesh are currently performed using majhee blocks, named after the majhee community leaders responsible for portions of each camp. Majhee blocks are represented as irregular polygons in the NPM shapefile dataset available on the Humanitarian Data Exchange [8], as shown in Figure 2. It is worth noting that manual population estimation is a challenging and time-consuming task. The majhee block system is not perfect due to the existence of some gaps in camp coverage and unquantified errors in population count.

Nonetheless, image chip population labels (y_{GT}) are calculated using the weighted sum of majhee block polygons that overlap with a chip’s bounding box, as described by Equation 1 and similar to the methodology used by [2]. The square bounding box ($b_{i,j,l}$) is defined by its bottom-left corner pixel coordinate (i,j), and by length l . The spatial area of a bounding box or polygon is represented by A_p and $c(z_k)$ is the population of a given polygon. Additionally, z_k is the k th polygon from the set of polygons n in the shapefile that most closely matches the date and time of when the overhead image was captured, as follows:

$$y_{GT}(b_{i,j,l}) = \sum_{k=0}^n c(z_k) \frac{A_p(z_k \cap b_{i,j,l})}{A_p(b_{i,j,l})}. \quad (1)$$

4. METHOD

We perform displacement camp population estimation using a convolutional neural network (CNN). Our model approximates the population density \hat{f} , which can be represented as $\hat{f}(m(I_{i,j,l}) \wedge I_{i,j,l}, GSD; \Theta) \approx f(I_{i,j,l})$, where $I_{i,j,l}$ is the input image chip, $m(I_{i,j,l}) \wedge I_{i,j,l}$ is the image chip overlaid by the mask (m) chip, GSD is the ground sample distance of the image chip, A_c is the spatial area depicted in the chip,

Θ is the set of model parameters, and f is the known chip population density. The population density of the camp depicted in the chip is then translated into a chip scale population prediction, $\hat{P}(I_{i,j,l})$, using the spatial area of the chip, $\hat{P}(I_{i,j,l}) = \hat{f}A_c \approx P(I_{i,j,l})$. During inference all chips are uniformly sized ($l = 224$ for $I_{i,j,l}$) and do not overlap, so the chips corresponding to a single full overhead image can be aggregated to obtain the total predicted camp population.

5. EXPERIMENTAL RESULTS

Quantitative performance is evaluated on entire camps since chip-scale population labels are unavailable. However, we perform a chip-level qualitative evaluation to assess our model’s fine-grain prediction performance.

As a baseline, which we call OSM-ONLY, we used a linear Huber regressor that predicts the population of an image chip using the structure area (A_s) of the corresponding chip as an input. Similar to the CNN, the baseline performs fine-grained predictions on chips, sums the predictions, and yields the full camp population estimation. Let n be the number of structures (s_r) in a given structure mask chip. The total structure area (A_t) is calculated using the structure segmentation mask chip ($m(I_{i,j,l})$) of the corresponding image chip ($I_{i,j,l}$), as follows:

$$A_t(m(I_{i,j,l})) = \sum_{r=0}^n A_s(s_r). \quad (2)$$

Metrics: Similar to [2], we use Mean Absolute Error (MAE) to reflect the total number of individuals unaccounted for, and Mean Absolute Percent Error (MAPE) to generally reflect model performance.

Implementation Details: The proposed CNN model was built from a ResNet50 architecture pre-trained on ImageNet and modified for regression. The image chip’s GSD was concatenated with the bottleneck features of the network attached to the regression head. Ground truth was represented as density (i.e., population normalized by area) instead of raw population to enable more robust training and testing of models with imagery of varying GSD and chip sizes. Training was performed with the Adam optimizer and a default initialized triangular cyclic learning Rate [10]. Huber loss was used and found to work better than mean squared error or log-cosh loss. Additionally, the following image augmentations [11] were applied to improve model generalization: vertical flip, random 90 deg rotation, CLAHE, random brightness, random gamma, hue saturation, and random contrast.

Quantitative Evaluation: The quantitative results are 3,704 MAE (10.00% MAPE) and 3,341 MAE (7.02% MAPE) for OSM-ONLY and our approach, respectively. These results demonstrate that the CNN’s image-based features are better-suited for predictive displacement camp populations over structure area alone.

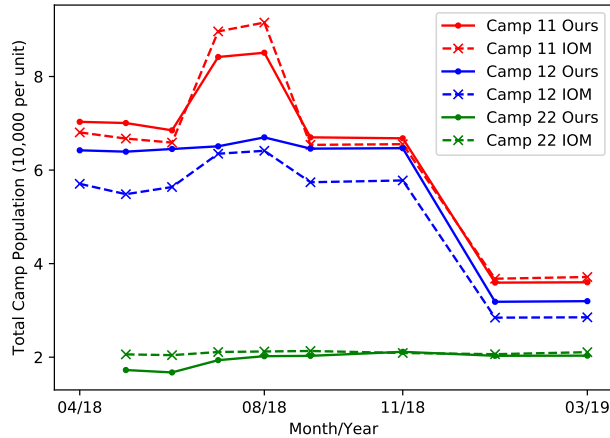


Fig. 4: CNN model (solid line) performance on the three held-out camps (11, 12, and 22) compared to IOM-reported population (dashed).

Qualitative Evaluation: Qualitative evaluation of our model is shown in Figure 3. Overall, our model’s predicted population distribution closely follows the IOM-reported counts. Our model tends to struggle with chips that are densely populated (dark red regions) as opposed to sparsely-populated regions (shades of blue). For fair comparison between the reported and predicted population, areas that do not overlap with majhee blocks are not shown. While our approach is capable of performing predictions in these areas, accuracy cannot be evaluated in said areas. Additionally, Figure 4 shows the CNN total camp population performance compared with IOM-reported population data over time². Note that imagery and population counts were recorded on a monthly basis for 8 months, and our method was able to predict an aggregated camp population that followed the contour of the reported population for all three camps. Additionally, our model did not strictly overpredict or underpredict for two of the three camps, suggesting our model is not unidirectionally biased.

6. CONCLUSIONS

We introduce a new dataset for developing fine-grained population estimation models using high-resolution drone imagery of refugee camps in Cox’s Bazar, Bangladesh. Using this dataset, we developed a novel approach capable of achieving 7.02% mean absolute population estimation error on a sequestered test set. The success of our approach demonstrates that structure masks can be used to encourage networks to learn correlations between building image features and population estimates. Future work could incorporate: methods for structure segmentation [12] and more rigorous disaggregation techniques [4]. Nevertheless, we believe this approach and dataset constitute an important step towards the development of tools that enable the humanitarian community to effectively and rapidly respond to the global displacement crisis.

²Camp 22 did not have recorded data for 04/18. Data for months 10/18, 12/18, and 02/19 was incorporated into the following month’s data.

7. REFERENCES

- [1] International Organization for Migration, UN High Commissioner for Refugees, UN Resident Coordinator for Bangladesh, Inter Sector Coordination Group, “2019 Joint Response Plan for Rohingya Humanitarian Crisis (January-December),” .
- [2] Caleb Robinson, Fred Hohman, and Bistra Dilkina, “A deep learning approach for population estimation from satellite imagery,” in *SIGSPATIAL GeoHumanities*, 2017.
- [3] Wenjie Hu, Jay Harshadbhai Patel, Zoe-Alanah Robert, Paul Novosad, Samuel Asher, Zhongyi Tang, Marshall Burke, David B. Lobell, and Stefano Ermon, “Mapping missing population in rural india: A deep learning approach with satellite imagery,” in *AIES*, 2019.
- [4] Nathan Jacobs, Adam Kraft, Muhammad Usman Rafique, and Ranti Dev Sharma, “A weakly supervised approach for estimating spatial density functions from high-resolution satellite imagery,” in *ACM SIGSPATIAL*, 2018.
- [5] Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee, “Functional map of the world,” in *CVPR*, 2018.
- [6] T. Kemper, M. Jenerowicz, M. Pesaresi, and P. Soille, “Enumeration of dwellings in darfur camps from geoeye-1 satellite images using mathematical morphology,” *JSTARS*, 2011.
- [7] F Galeon, “Estimation of population in informal settlement communities using high resolution satellite image,” in *XXI ISPRS Congress, Commission IV.*, 2008.
- [8] IOM Bangladesh, “IOM Bangladesh - Needs and Population Monitoring (NPM) UAV imagery and GIS package by camp (March 2019),” .
- [9] Geoff Boeing, “Osmnx: New methods for acquiring, constructing, analyzing, and visualizing complex street networks,” *Computers Environment and Urban Systems*, 2017.
- [10] Leslie N Smith, “Cyclical learning rates for training neural networks,” in *WACV*, 2017.
- [11] A. Buslaev, A. Parinov, E. Khvedchenya, V. I. Iglovikov, and A. A. Kalinin, “Albumentations: fast and flexible image augmentations,” *arXiv*, 2018.
- [12] Ziran Ye, Yongyong Fu, Muye Gan, Jinsong Deng, Alexis Comber, and Ke Wang, “Building extraction from very high resolution aerial imagery using joint attention deep neural network,” *Remote Sensing*, 2019.