

Release Statement

Gridded maps of residential/non-residential building classifications and associated building patterns for Nigeria (NGA), version 1.0

21 December 2022

This work is part of the GRID3 project with funding from the Bill and Melinda Gates Foundation (BMGF) and the United Kingdom's Foreign, Commonwealth & Development Office (INV-009579, formerly OPP1182425). Method, coding, modelling, and data production have been carried out by Christopher Lloyd. Oversight was provided by Andrew J. Tatem. These data may be distributed using a [Creative Commons Attribution NonCommercial ShareAlike 4.0 License](#). Contact release@worldpop.org for more information.

The whole WorldPop group are acknowledged for overall project support. The author acknowledges the use of the IRIDIS High Performance Computing Facility, and associated support services at the University of Southampton, in the completion of this work. Thanks go to Maksym Bondarenko (WorldPop) for support with online publication, and to Attila Lazar and Edith Darin for internal review of data and the release statement. The author acknowledges the efforts of WorldPop's partners (United Nations Population Fund (UNFPA), Center for International Earth Science Information Network (CIESIN) in the Earth Institute at Columbia University, and the Flowminder Foundation) in supporting access to the building footprints. The author thanks to Blair-Freese (formerly of BMGF) and Heather Chamberlain (WorldPop) for providing coordination between WorldPop and BMGF and Ecopia.

CITATION

Lloyd, C.T and Tatem, A.J. 2022. Gridded maps of residential/non-residential building classifications and associated building patterns for Nigeria (NGA), version 1.0. University of Southampton: Southampton, UK. doi:10.5258/SOTON/WP00753

RELEASE CONTENT

1. NGA_building_class_metrics_v1_0_classification_binary.tif
2. NGA_building_class_metrics_v1_0_residential_mean_score.tif
3. NGA_building_class_metrics_v1_0_residential_count.tif
4. NGA_building_class_metrics_v1_0_nonresidential_count.tif
5. NGA_building_class_metrics_v1_0_residential_density.tif
6. NGA_building_class_metrics_v1_0_nonresidential_density.tif

Rasters are provided for countries identified using [ISO3 country codes](#).

FILES DESCRIPTION

The geotiff rasters have a spatial resolution of approximately 100m (3 arc seconds). Their coordinate reference system is WGS84. Each building has been considered in the grid cell that contained the centroid* of its building footprint. NAs represent grid cells that contain no building footprint centroid*.

NGA_building_class_metrics_v1_0_classification_binary.tif

This raster contains the binary classification for all buildings per each grid cell (as a binary integer. 1=residential, 0=non-residential) derived from modelling.

NGA_building_class_metrics_v1_0_residential_mean_score.tif

This raster contains the mean percentage score (as a decimal number) for all buildings per each grid cell. A value of 1 indicates highest likelihood that all buildings within the grid cell are residential. A value of 0

indicates highest likelihood that all buildings within the grid cell are non-residential. This layer enables the user to pick a threshold for the residential/non-residential cut-off threshold other than that identified as optimal during modelling.

NGA_building_class_metrics_v1_0_residential_count.tif

This raster contains a count of residential buildings per each grid cell.

NGA_building_class_metrics_v1_0_nonresidential_count.tif

This raster contains a count of non-residential buildings per each grid cell.

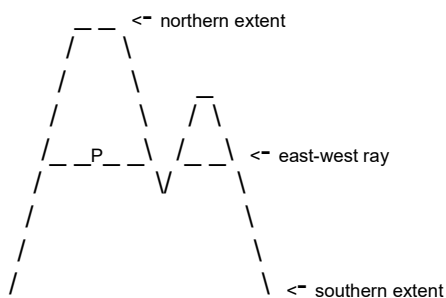
NGA_building_class_metrics_v1_0_residential_density.tif

This raster contains a measure of the number of residential buildings per grid cell area in square kilometres (km²), i.e. building count divided by the area in square kilometres of the grid cell. If needed, the grid cell area can be retrieved by dividing the count raster by the density raster.

NGA_building_class_metrics_v1_0_nonresidential_density.tif

This raster contains a measure of the number of non-residential buildings per grid cell area in square kilometres (km²), i.e. building count divided by the total number of square kilometres in the grid cell. If needed, the grid cell area can be retrieved by dividing the count raster by the density raster.

+ For simplicity, the term centroid is used throughout this document to indicate the centre point within each building. However, true building centroids are often located outside of the building polygon, perhaps because the building is L-shaped or because several polygons represent the building and so these are grouped into one geometry (i.e. multi-polygon geometry). It is undesirable for building centroids to be located outside of building polygons because allocation of building to a raster grid cell is less accurate. Hence, rather than use the `st_centroid` function to define the centre of a building we instead use the `st_point_on_surface` function (Pebesma, 2018). Building centre is accordingly defined as the point half-way along the longest segment of the east-west ray that intersects the building and lies half-way between the northern and southern extents of the building geometry (see illustration below).



`p` = Centre of building as defined by `st_point_on_surface` function

RELEASE HISTORY

Version 1.0 (21 December 2022)

- This is the original release of the data

SOURCE DATA

The classified building data have been modelled using building footprints (polygon features) provided by the Digitize Africa project (Ecopia.AI and Maxar Technologies, 2020/2021), and building footprint and highway data provided by OpenStreetMap (© 2020-2022 OpenStreetMap contributors; geofabrik.de). Digitize Africa is a two-year project funded by the Bill and Melinda Gates Foundation to map buildings and roads in 51 countries across sub-Saharan Africa using satellite imagery and artificial intelligence (AI) to support humanitarian assistance and sustainable development. Maxar provided their Vivid satellite imagery mosaics (50 cm resolution) and Ecopia.AI generated the building footprints using their artificial intelligence-assisted feature extraction techniques (Ecopia, 2021).

The building classification model utilizes impervious surface data provided by US NASA (SEDAC) (Global Man-made Impervious Surface data from Landsat; Brown de Colstoun et al., 2017).

The model has been trained using building label data provided by OpenStreetMap (OSM). In addition, building label data for Nigeria have been sourced from household surveys. The Oak Ridge National Laboratory (ORNL) contributed to the first household survey round. ORNL designed, and eHealth Africa implemented, data collection during 2016–17. WorldPop (University of Southampton) designed, and eHealth Africa implemented, household survey data collection during 2018–2019.

METHODS OVERVIEW

The classified building layers are generated in two steps: (1) The modelling of individual buildings into residential and non-residential classes, and then (2) the aggregation of these results at the grid cell-level.

1. Model input data are preprocessed using GIS, geospatial, database, and statistical software (OSGeo - GDAL, GRASS GIS, Spatialite, and R), discussed in detail in Lloyd et al. (2020). The modelling follows the methodology of Sturrock, et al. (2018), but with countrywide coverage. To handle potential building heterogeneity, the classification model is run separately for urban and rural areas, and outputs combined. US NASA (SEDAC) impervious surface data (Brown de Colstoun et al., 2017) represent urban areas for input to the model. Elsewhere, rural areas are demarcated by an absence of impervious surface data. Building classification predictions are combined with known building labels (used to train and test the model), with the latter taking precedence, to produce model outputs.
2. The raster datasets have been produced using a modified version of the code used to produce 'Gridded maps of building patterns throughout sub-Saharan Africa', version 2.0 (Dooley, Leasure, Boo, and Tatem. 2021). All data processing has been carried out using R (R Core Team, 2013). The classified building footprint polygons (model output) have been converted into centroid+ points in UTM projection using the *st_point_on_surface* function in the *sf* R package (Pebesma, 2018). The points have then been re-projected to WGS84 using *sf's st_transform* function so that their corresponding cell IDs could be identified in the WGS84 projected raster master grid. The WorldPop master grid national boundary (WorldPop et al. 2018) defines the extent, resolution, and coordinate reference system of the output building footprint metric layers. Building centroid+ cell IDs have been obtained using the *cellFromXY* function in the *raster* R package (Hijmans & van Etten, 2012). The 'settled' grid cells in the rasters have at least one building footprint centroid+ within the pixel area. No areas have been masked out.

The binary classification raster (**_classification_binary.tif*) represents the mode binary building classification per grid cell, derived from the model. Where, for a given grid cell, there are equal numbers of residential and non-residential buildings, the grid cell is classified as residential.

The mean score raster (**_residential_mean_score.tif*) represents the mean percentage likelihood of a grid cell containing residential buildings. The modelling technique invariably uses unequal amounts of residential and non-residential data in model training and testing. Hence, the model uses a country specific cut-off threshold at which equally good performance is achieved when classifying residential and non-residential

structures. This threshold is applied by the model to produce spatially comparable binary building classifications (Sturrock et al. 2018), which are subsequently averaged per grid cell to produce the binary classification raster (see above). The mean score raster allows the user to pick a different threshold for the residential/non-residential cut-off threshold than the one identified as optimal during modelling.

Building counts per grid cell have been calculated by simply summing the number of residential or non-residential building centroids* (i.e. binary building classification) for each cell ID. Building density per grid cell have been calculated by dividing the building count in a cell by the area of the cell. Grid cell area have been calculated in square kilometres using the *area* function in the *raster* R package.

NGA MODEL METRICS

57,887,322 buildings total
55,184,335 residential
2,702,987 non-residential
116,385 labelled buildings used in training and testing

Urban modelling

0.9278047 AUC
0.754 res/non-res cut-off threshold
84.6% residential correctly classified in the testing dataset
84.5% non-residential correctly classified in the testing dataset

Rural modelling

0.965716 AUC
0.967 res/non-res cut-off threshold
89.8% residential correctly classified in the testing dataset
89.6% non-residential correctly classified in the testing dataset

LIMITATIONS

In urban locations where building density is greatest, there will typically be greater numbers of residential than non-residential buildings per ~100m pixel in raster output. This is because buildings with mixed use are included in the residential classification for the purpose of modelling, and because residential buildings tend to be smaller than non-residential (e.g warehouses, universities, factories, malls) buildings. Hence, the binary classification raster will be more likely to reflect a residential classification in urban pixels, and in count rasters the count of residential buildings is likely to be greater than non-residential in urban pixels.

Modelling is most successful where data on structure type are available for a reasonable subset of buildings with a reasonable spatial distribution within a given country. In some modelling instances, this information is available for only a relatively small number of structures. In such instances, predictions are extrapolated in neighbouring regions. To handle potential building heterogeneity, the classification model is run separately for urban and rural areas per country, and outputs combined.

There may be some error in OSM or other label dataset attributes or bias could be introduced if building type information is only available for certain categories of structures. Ground truth data (i.e. household surveys) provide valuable additional building label data in some countries modelled.

This modelling is only really useful where every structure has been mapped or where the completeness can be quantified. If gaps exist of unknown size, as is often the case, estimates of numbers of residential structures will be low. While OSM does not currently have a mechanism to estimate the completeness of the building enumeration data, efforts are underway to provide this information (Sturrock et al. 2018).

Building footprint accuracy and completeness are discussed in Lloyd et al (2020). For ORNL and eHealth Africa building labels, the margin of error of building label positions, as sampled during original fieldwork surveys, is estimated to be no worse than approximately 450 m (WorldPop, personal communication).

The limitations of GMIS impervious data are discussed in Brown de Colstoun et al. (2017) and Gutman et al. (2013). As elucidated in Lloyd et al (2020), some limitations include those of the source GLS 2010 imagery, which has ~30m spatial resolution for the year 2010 and contains residual cloud covered areas and gaps caused by the Landsat 7 Scan Line Corrector (SLC) failure. Some of these gaps have not been filled. It is also possible that small areas with impervious cover have been removed, or small areas of bare soil within cities have a non-zero impervious cover. These errors are due to limitations of processing by the GMIS project (2017).

REFERENCES

- Brown de Colstoun E.C, Huang C, Wang P, Tilton J.C, Tan B, Phillips J, Niemczura S, Ling P.-Y, Wolfe R.E. 2017. Global Man-Made Impervious Surface (GMIS) Dataset From Landsat; NASA Socioeconomic Data and Applications Center (SEDAC): Palisades, NY, USA. <https://sedac.ciesin.columbia.edu/data/set/ulandsat-gmis-v1>
- Brown de Colstoun E.C, Huang, C, Wang P, Tilton J.C, Tan B, Phillips J, Niemczura S, Ling P.-Y, Wolfe R.E. 2017. Documentation for Global Man-made Impervious Surface (GMIS) Dataset From Landsat, v1 (2010); NASA Socioeconomic Data and Applications Center (SEDAC): Palisades, NY, USA.
- Dooley C.A, Leasure D.R, Boo G, and Tatem A.J. 2021. Gridded maps of building patterns throughout sub-Saharan Africa, version 2.0. University of Southampton: Southampton, UK. Source of building footprints “Ecopia Vector Maps Powered by Maxar Satellite Imagery”© 2020/2021. doi:10.5258/SOTON/WP00712.
- Ecopia.AI and Maxar Technologies. 2020/2021. Digitize Africa building footprint data. <https://www.ecopiatech.com/commerical-landing-africa>
- Ecopia.AI and Maxar Technologies. 2021. The Most Comprehensive Map of Buildings in the USA. <https://www.ecopiatech.com/post/the-most-comprehensive-map-of-buildings-in-the-usa>
- Gutman G, Huang C, Chander G, Noojipady P, Masek J.G. 2013. Assessment of the NASA–USGS Global Land Survey (GLS) datasets. *Remote Sens. Environ.* 134, 249–265.
- Hijmans R.J. & van Etten J. 2012. raster: Geographic analysis and modeling with raster data. R package version 2.0-12. <http://CRAN.R-project.org/package=raster>
- Lloyd C.T, Sturrock H.J.W, Leasure D.R, Jochem W.C, Lázár A.N, Tatem A.J. 2020. Using GIS and Machine Learning to Classify Residential Status of Urban Buildings in Low and Middle Income Settings. *Rem Sens* 12 (23). 3847. <https://doi.org/10.3390/rs12233847>
- Pebesma, E. 2018. Simple Features for R: Standardized Support for Spatial Vector Data. *The R Journal*, 10(1), 439–446. doi: 10.32614/RJ-2018-009, <https://doi.org/10.32614/RJ-2018-009> .
- R Core Team. 2013. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.
- Sturrock H.J.W, Woolheater K, Bennett A.F, Andrade-Pacheco R, Midekisa, A. Predicting residential structures from open source remotely enumerated data using machine learning. *PLoS ONE* 2018, 13, e0204399. <https://doi.org/10.1371/journal.pone.0204399>
- WorldPop (www.worldpop.org - School of Geography and Environmental Science, University of Southampton; Department of Geography and Geosciences, University of Louisville; Departement de Geographie, Universite de Namur) and Center for International Earth Science Information Network (CIESIN), Columbia University. 2018. Global High Resolution Population Denominators Project - Funded by The Bill and Melinda Gates Foundation (OPP1134076). <https://dx.doi.org/10.5258/SOTON/WP00651>