# Object Detection Model, Image Data and Results from the "When Computers Dream of Charcoal: Using Deep Learning, Open Tools and Open Data to Identify Relict Charcoal Hearths in and Around State Game Lands in Pennsylvania" Paper

**JEFF BLACKADAR** (iD)
**BENJAMIN CARTER** (iD)
**WESTON CONNER** (iD)

*Author affiliations can be found in the back matter of this article*

]u[ ubiquity press

## ABSTRACT

These data were used to build an object detection model to locate Relict Charcoal Hearths (RCH) as described in the paper "When Computers Dream of Charcoal: Using Deep Learning, Open Tools and Open Data to Identify Relict Charcoal Hearths in and around State Game Lands in Pennsylvania" [1]. This is the second grouping of data for the paper above. The first grouping is also available in this journal, see "Geospatial and image data from the "When Computers Dream of Charcoal: Using Deep Learning, Open Tools and Open Data to Identify Relict Charcoal Hearths in and around State Game Lands in Pennsylvania" paper" [2].

**THESE FILES CONSIST OF:**

- JPEGs representing tiles of larger Slope TIFF files derived from LiDAR for the State Game Lands (SGL) of Pennsylvania, United States [3, 4, 5, 6]. A subset of these tiles was used to train the model.
- A Shapefile of points of known relict charcoal hearths (RCH).
- XML files representing the pixel points of known RCHs on JPEG files used for training.
- Jupyter notebooks of programs used to prepare data and train a Mask R-CNN model.
- The Mask R-CNN model H5 file.
- Shapefile and GeoJSON of object detection results from the model showing locations of possible RCH in all SGLs.
- XML files representing the pixel points of predicted RCH on JPEG files used for predictions.
- GeoJSON of results using cluster analysis.

These data are stored on Zenodo.org. The programs are stored on *Github.com*.

# (1) OVERVIEW

## CONTEXT

### Spatial coverage

Description: State Game Lands within the State of Pennsylvania, United States.

(WGS84):

North_Bounding_Coordinate: 42.269479

South_Bounding_Coordinate: 39.719860

East_Bounding_Coordinate: –74.689583

West_Bounding_Coordinate: –80.519349

### Temporal coverage

AD1700-AD1945

# (2) METHODS

## INTRODUCTION

These data were used to build an object detection model to locate Relict Charcoal Hearths (RCH) as described in the paper "When Computers Dream of Charcoal: Using Deep Learning, Open Tools and Open Data to Identify Relict Charcoal Hearths in and around State Game Lands in Pennsylvania" [1].

## STEPS

A Python program was used to split the Slope TIFF files into smaller 1024x768 pixel JPEG tiles. If a tile contained at least one location point of a known RCH it was also retained to train the model. For each training image the program generated a corresponding XML file with the pixel coordinates of the boundaries of known RCHs in the image. Mask R-CNN is used to train the model using the images and XML files [7].

The trained model H5 file was used to detect RCHs in all of the JPEG tiles. This produces a Shapefile of predicted RCHs. The initial predictions from the model were processed using cluster analysis to produce a final list of detected RCHs, stored in GeoJSON. (*Figure 1*) lists the detailed steps.

## SAMPLING STRATEGY

For training the model, 20% of the images across all SGLs were set aside for automated testing.

## QUALITY CONTROL

The quality control for model training and selection had these steps:

- Each model was assessed for its Average Precision. A score above 50% meant the model was a possible candidate for object detection.
- Candidate models were used to detect RCHs in 20 images. The results were manually inspected for accuracy.
- Models that passed visual inspection were formally scored using a set of 100 images selected at random.
- Models passing formal scoring were reviewed in depth using large sample size, spanning multiple SGLs.

## CONSTRAINTS

Predictions for State Game Lands 264 and 258 could not be processed due to problems in the source files.

# (3) DATASET DESCRIPTION

## OBJECT NAME

Since the dataset contains multiple objects, filenames are displayed in Zenodo.

## DATA TYPE

Primary data: Locations of known RCHs. Secondary data: JPEG tiles, XML files, Object Recognition File, ShapeFile and GeoJSON of results.

## FORMAT NAMES AND VERSIONS

TIFF, Shapefile, GeoJSON, JPEG, XML, H5, Jupyter Notebook

## CREATION DATES

01/01/2020 – 31/08/2020

## DATASET CREATORS

Jeff Blackadar, object recognition programmer, author, Carleton University. *https://orcid.org/0000-0002-8160-0942*

Ben Carter, archaeologist, author, data collection and verification, GIS expert, Muhlenberg College. *https://orcid.org/0000-0002-7464-0989*

Weston Conner, archaeologist, author, data collection and verification, GIS expert, Lehigh University. *https://orcid.org/0000-0001-9906-3762*

## LANGUAGE

English

## LICENSE

CC-BY.

## REPOSITORY LOCATION

Documentation: *https://zenodo.org/record/4766351*

Data for

1. Vector files resulting from manual identification of relict charcoal hearths (RCHs). These were used to train the Mask R-CNN model.
    a. Shapefile - *https://zenodo.org/record/4593605*
    b. GeoJSON - *https://zenodo.org/record/4593622*
2. Program for splitting slope analysis (in TIFF format) into smaller tiles (in JPEG format).
    a. *https://github.com/jeffblackadar/charcoalhearths/blob/master/0_split_tifs_refactored.ipynb*
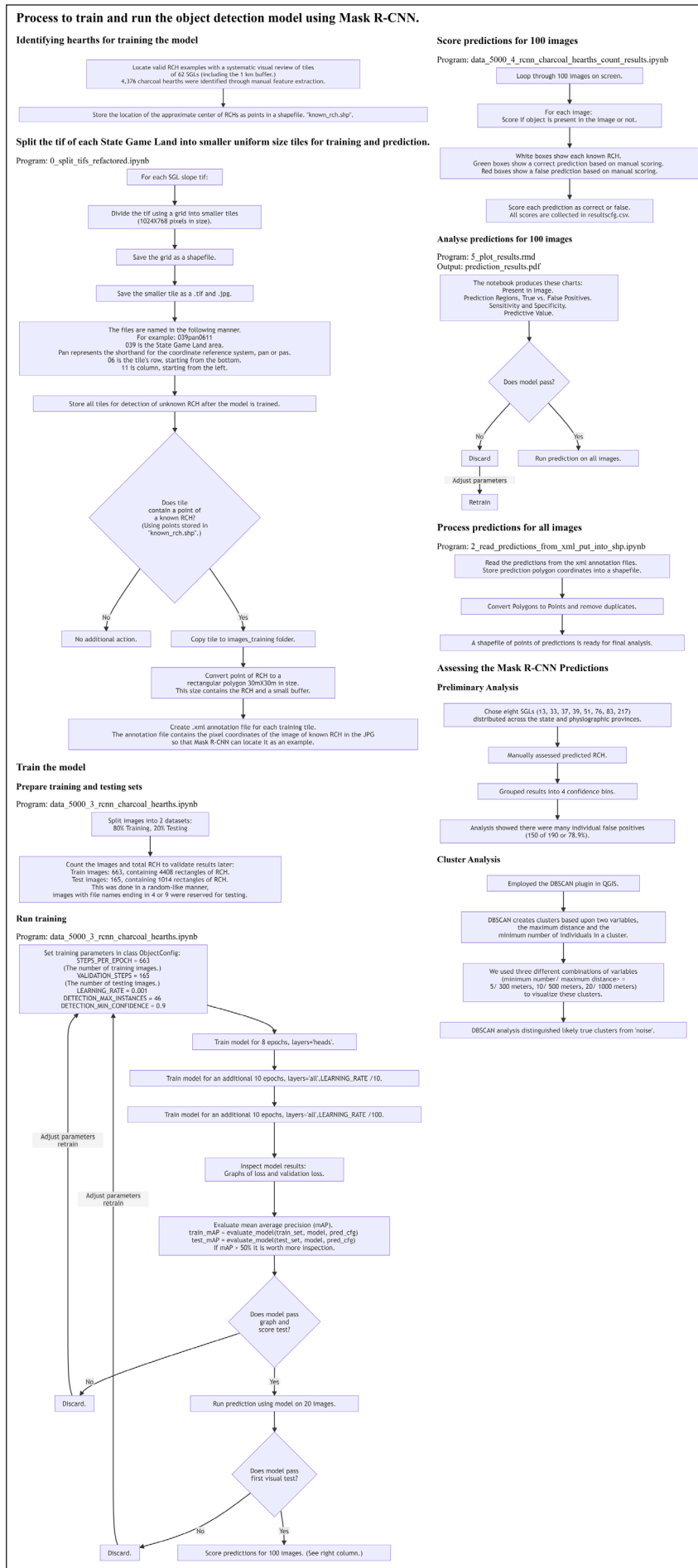3. Data files for Training:

a. RCH Detection with Mask R-CNN Image Annotations. *https://zenodo.org/record/4575582* This file contains a collection of xml files that contain coordinates of the locations of known RCHs on images (from above). These files are known as annotations and are used by Mask R-CNN to identify objects to detect during training of a model.

b. RCH Detection with Mask R-CNN Training Images. *https://zenodo.org/record/4579935* This file contains all of the images used for training the Mask R-CNN model. Each image contains at least one known RCH.

C. Polygons for tiles of LiDAR data. *https://zenodo.org/record/4580726*

4. Program for Mask R-CNN training and prediction. Known as "data_5000_3_rcnn_charcoal_hearths. ipynb" [7]

*https://github.com/jeffblackadar/charcoalhearths/blob/ master/data_5000_3_rcnn_charcoal_hearths.ipynb*

5. Data files of the model and predictions

a. Resultant trained Mask R-CNN model. *https://zenodo.org/record/4579946*

b. Predictions from the model, in x, y coordinates (not geolocated) in XML format. *https://zenodo. org/record/4581281* The format of these files is similar to the training annotations.

c. RCH Detection with Mask R-CNN Images. *https://zenodo.org/record/4583945* This file contains all of the images representing tiles of LiDAR images of State Game Lands. These images are used for predictions to locate RCHs. (A subset was used for training. See 3b RCH Detection with Mask R-CNN Training Images above.)

6. Program that produces confidence scores for the predictions above. Known as "data_5000_4_rcnn_ charcoal_hearths_count_results.ipynb"

*https://github.com/jeffblackadar/charcoalhearths/blob/ master/data_5000_4_rcnn_charcoal_hearths_count_results. ipynb*

7. Program to remove duplicates (because some tiles included multiple SGLs and some tiles overlapped both the north and south LiDAR tile indices of the state), converts squares to their centroid for comparison and saves unique squares in geolocated vector files. Known as "2_read_predictions_from_ xml_put_into_shp.ipynb"

*https://github.com/jeffblackadar/charcoalhearths/blob/ master/2_read_predictions_from_xml_put_into_shp.ipynb*

8. Vector file of prediction results.

a. Shapefile- *https://zenodo.org/record/4593734*

b. GeoJSON- *https://zenodo.org/record/4593747*

9. Prediction results with additional variables (bins for assessment, ID of training data, cluster analysis and visual confirmation)

a. GeoJSON (no shapefile)- *https://zenodo.org/ record/4593767*

b. Variables:

i. **id** = unique identifier starting with 3-digit SGL number, PAN or PAS (projections) and, within those a unique four-digit identifier

ii. **score** = confidence score

iii. **SGL** = State Game Land number

iv. **SGLImage** = name of TIFF file of merged LiDAR tiles

V. **Confirm** = Whether the predicted hearth was determined, through visual inspection, to be a likely true positive (Y) or a false positive (N)

vi. **Bin#**- in assessing these predictions we "binned" the results based upon the confidence score.

vii. **Bin_select** = 1 if this record (predicted RCH) was selected for assessment within that bin

viii. **TrainID** = Original ID of the training data (only training data that matched with a prediction are included).

ix. **Clusters5_300** = resultant clusters from DBSCAN where minimum cluster size = 5 and maximum distance = 300 meters

x. **Clusters10_500** = resultant clusters from DBSCAN where minimum cluster size = 10 and maximum distance = 500 meters

xi. **Clusters20_1000** = resultant clusters from DBSCAN where minimum cluster size = 20 and maximum distance = 1000 meters

xii. **CLUSTERCT** = How many of the above clusters included the predicted RCH (0–3). Derived from the previous three variables.

xiii. **3Cluster** = whether or not this predicted RCH was included in all three clusters.

10. False Negatives for the tiles around SGL 43 after close visual inspection (at 1:1000 scale).

a. GeoJSON = *https://zenodo.org/record/4758647*

## PUBLICATION DATE
11/03/2021

## (4) REUSE POTENTIAL

The model may be reused to detect similar looking objects in other landscapes. The shapefile and images may be used to train an improved prediction model. Also, the images and programs may be re-used to train a different object detection model to locate other types of objects in SGLs.

## COMPETING INTERESTS

The authors have no competing interests to declare.

## AUTHOR AFFILIATIONS

**Jeff Blackadar** orcid.org/0000-0002-8160-0942
Carleton University, CA

**Benjamin Carter** orcid.org/0000-0002-7464-0989
Data collection and verification, GIS expert, Muhlenberg College, US

**Weston Conner** orcid.org/0000-0001-9906-3762
Data collection and verification, GIS expert, Lehigh University, US

## REFERENCES

1. **Carter BP, Blackadar JH, Conner WLA.** "When Computers Dream of Charcoal: Using Deep Learning, Open Tools, and Open Data to Identify Relict Charcoal Hearths in and around State Game Lands in Pennsylvania." *Advances in Archaeological Practice*. 2021; 9(4): 257–71. Cambridge University Press. DOI: *https://doi.org/10.1017/aap.2021.17*

2. **Conner W, Carter B, Blackadar J.** "Geospatial and Image Data from the "When Computers Dream of Charcoal: Using Deep Learning, Open Tools and Open Data to Identify Relict Charcoal Hearths in and Around State Game Lands in Pennsylvania" Paper." *Journal of Open Archaeology Data*. Ubiquity Press. 2021; 7. DOI: *https://doi.org/10.5334/joad.80*

3. **PAMAP Program, PA Department of Conservation and Natural Resources, Bureau of Topographic and Geologic Survey.** PAMAP Program LAS Files (LiDAR Data of Pennsylvania). n.d. [Online]. Available: *https://www.pasda.psu.edu/uci/FullMetadataDisplay.aspx?file=pamap_lidar_LAS.xml* [Accessed 30 April 2020].

4. **PAMAP Program, Bureau of Topographic and Geologic Survey, PA Department of Conservation and Natural Resources.** PAMAP Tile Index - North 2006; 2006. [Online]. Available: *https://www.pasda.psu.edu/uci/DataSummary.aspx?dataset=266* [Accessed 30 April 2020].

5. **PAMAP Program, Bureau of Topographic and Geologic Survey, PA Department of Conservation and Natural Resources.** PAMAP Tile Index – South 2006. 2006. [Online]. Available: *https://www.pasda.psu.edu/uci/DataSummary.aspx?dataset=267* [Accessed 30 April 2020].

6. **Pennsylvania Game Commission., PGC State Game Lands.** 2018. [Online]. Available: *http://www.pasda.psu.edu/uci/DataSummary.aspx?dataset=86* [Accessed 30 April 2020].

7. **Abdulla W.** Mask R-CNN for Object Detection and Instance Segmentation on Keras and TensorFlow. n.d. [Online]. Available: *https://github.com/matterport/Mask_RCNN* [Accessed 30 April 2020].