

# Lunar Crater Identification via Deep Learning

Ari Silburt<sup>a,b,c,f</sup>, Mohamad Ali-Dib<sup>a,d,f</sup>, Chenchong Zhu<sup>b,d</sup>, Alan Jackson<sup>a,b,e</sup>,  
Diana Valencia<sup>a,b</sup>, Yevgeni Kissin<sup>b</sup>, Daniel Tamayo<sup>a,d</sup>, Kristen Menou<sup>a,b</sup>

<sup>a</sup>*Centre for Planetary Sciences, Department of Physical & Environmental Sciences, University of Toronto Scarborough, Toronto, Ontario M1C 1A4, Canada*

<sup>b</sup>*Department of Astronomy & Astrophysics, University of Toronto, Toronto, Ontario M5S 3H4, Canada*

<sup>c</sup>*Department of Astronomy & Astrophysics, Penn State University, Eberly College of Science, State College, PA 16801, USA*

<sup>d</sup>*Canadian Institute for Theoretical Astrophysics, 60 St. George St, University of Toronto, Toronto, ON M5S 3H8, Canada*

<sup>e</sup>*School of Earth and Space Exploration, Arizona State University, 781 E Terrace Mall, Tempe, AZ 85287-6004, USA*

<sup>f</sup>*These authors contributed equally to this work.*

---

Keywords: Moon; Crater Detection; Automation; Deep Learning

Corresponding authors:

- Ari Silburt, Department of Astronomy & Astrophysics, Penn State University, Eberly College of Science, State College, PA 16801, USA. e-mail: arisilburt@gmail.com .
- Mohamad Ali-Dib, Centre for Planetary Sciences, Department of Physical & Environmental Sciences, University of Toronto Scarborough, Toronto, Ontario M1C 1A4, Canada. e-mail: m.alidib@utoronto.ca .

## **Abstract**

Crater counting on the Moon and other bodies is crucial to constrain the dynamical history of the Solar System. This has traditionally been done by visual inspection of images, thus limiting the scope, efficiency, and/or accuracy of retrieval. In this paper we demonstrate the viability of using convolutional neural networks (CNNs) to determine the positions and sizes of craters from Lunar digital elevation maps (DEMs). We recover 92% of craters from the human-generated test set and almost double the total number of crater detections. Of these new craters, 15% are smaller in diameter than the minimum crater size in the ground-truth dataset. Our median fractional longitude, latitude and radius errors are 11% or less, representing good agreement with the human-generated datasets. From a manual inspection of 361 new craters we estimate the false positive rate of new craters to be 11%. Moreover, our Moon-trained CNN performs well when tested on DEM images of Mercury, detecting a large fraction of craters in each map. Our results suggest that deep learning will be a useful tool for rapidly and automatically extracting craters on various Solar System bodies. We make our code and data publicly available at <https://github.com/silburt/DeepMoon.git> and <https://doi.org/10.5281/zenodo.1133969>.

## **1. Introduction**

Craters formed by small impactors constitute an important surface property for many bodies in the Solar System. On airless bodies like the Moon,

Mercury, Ceres, and Vesta, weather based erosion, tectonics and volcanic activity have been largely non-existent resulting in the accumulation of impact craters over time. However, other eroding factors such as micrometeorite bombardment can affect smaller craters.

Crater densities permit the geological history of a body to be examined, and the relative chronology of a region to be assessed remotely. In addition, when in-situ samples are recovered from a body, absolute chronologies can be determined too. Inferred temporal variation in cratering rates have been used to make inferences about the dynamical history of the Solar System, including the (debated) possibility of a Late Heavy Bombardment, (e.g., Hartmann 1970; Ryder 2002; Gomes et al. 2005; Chapman et al. 2007; Bottke and Norman 2017). Crater records and chronology are thus central to any formation theory about the Solar System. In addition, the size distribution of craters directly probes the dynamics and size distribution of the impactor population (Strom et al., 2005). For example from the size distribution of craters on the Lunar highlands, Minton et al. (2015) argued that the impactor population contained comparatively fewer large bodies than the asteroid belt does today.

Traditionally, crater detection has been done manually via visual inspection of images. However this approach is not practical for the vast numbers of kilometre and sub-kilometre sized craters on the Moon (and other Solar System bodies), resulting in human-generated databases that are either spatially comprehensive but restricted to the largest craters, or size

comprehensive but limited to a very specific geographic region (Stepinski et al., 2012; Bandeira et al., 2012). In addition, manual crater counting by experts can yield disagreements as high as 40% (Greeley and Gault, 1970; Kirchoff et al., 2011; Robbins et al., 2014).

As a result, scientists have developed crater detection algorithms (CDAs) to automate the process of classifying craters. Such CDAs include edge detection (Emami et al., 2015), Hough transforms (Salamunićcar and Lončarić, 2010), support-vector machines (Wetzler et al., 2005), decision trees (Stepinski et al., 2012) and neural networks (Wetzler et al., 2005; Cohen et al., 2016; Palafox et al., 2017). Multi-step approaches have also been tried. For example, Di et al. (2014) used a boosting algorithm to box the crater-containing region, then a Hough transform to delimit the edges of the crater. Boukercha et al. (2014) used a similar approach where an initial detection algorithm provided crater candidates which were subsequently classified as true or false positives by a support-vector machine or polynomial classifier.

These CDAs tend to perform well on the datasets upon which they were trained, but not to generalize well on unseen patches or other bodies (see Stepinski et al. (2012) for a review and Chung et al. (2014) for a comparison between classical and machine learning based techniques.). The difficulty in designing robust CDAs stems from the complex nature of craters, having large variations in shape and illumination, orders of magnitude size differences, overlap and degradation. An algorithm capable of universally

identifying craters on Solar System bodies would be invaluable to the community.

In this work we train a deep learning architecture known as a convolutional neural network (CNN) to perform crater identification<sup>1</sup> on Lunar digital elevation map (DEM) images, and transfer-learn our Moon-trained CNN to identify craters on Mercury. There are numerous reasons for using CNNs to detect craters. First, CNNs have demonstrated impressive performance on a variety of computer vision problems and other datasets where features are correlated (e.g. Long et al., 2015), demonstrating their versatility. This includes, in addition to images, sounds and signals. Second, CNNs engineer their own representation features, alleviating the need for a human to develop sophisticated pre-processing algorithms and custom input features. Finally, CNNs have been able to successfully classify objects that appear at multiple scales in a single image (Zhang et al., 2016; Zeng et al., 2017), a property very relevant to crater counting.

## 2. Methods

The code to generate the data set (Section 2.1), train our model (Section 2.7), and extract the resulting crater distribution (Section 2.4 and Section 2.5) is available at <https://github.com/silburt/DeepMoon.git>. The

---

<sup>1</sup>By crater identification, we mean a) output the pixel locations of crater rims from DEM images via a CNN segmentation task, and b) extract the crater coordinates from these CNN outputs using a custom pipeline, explained in the Methods section.

data used to train, validate and test our model, the global DEM used to make our input images, our best model and final test set crater distribution can be found at <https://doi.org/10.5281/zenodo.1133969>. We use Keras (Chollet, 2015) version 1.2.2 with Tensorflow (Abadi et al., 2016) version 0.10 to build and train our model, but our code is also compatible with Keras 2.0.2 and Tensorflow 1.0.

### *2.1. Data Preparation*

Our input data was generated by randomly cropping digital elevation map (DEM) images from the Lunar Reconnaissance Orbiter (LRO) and Kaguya merged digital elevation model, which spans  $\pm 60$  degrees in latitude (and the full range in longitude) and has a resolution of 512 pixels/degree, or 59 m/pixel (Barker et al. 2016; available at LOLA Team and Kaguya Team 2015). This global grayscale map is a Plate Carree projection with a resolution of  $184320 \times 61440$  pixels and a bit depth of 16 bits/pixel; we downsampled it to  $92160 \times 30720$  pixels and 8 bits/pixel. We use an elevation map, rather than an optical one, because a crater's appearance in an elevation map is not affected by the direction of incident sunlight. This reduces variation in appearance between craters, making it easier to train a CNN to identify them.

Each input DEM image is generated by a) randomly cropping a square area of the global map, b) downsampling the cropped image to  $256 \times 256$  pixels, c) transforming the image to an orthographic projection using

the Cartopy Python package (UK Met. Office, 2015) in order to minimize image distortion, and d) linearly rescaling image intensity to boost contrast. The position of the cropped region in a) is randomly selected with a uniform distribution, and its length is randomly selected from a log-uniform distribution with minimum and maximum bounds of 500 and 6500 pixels (59 km and 770 km), respectively. The transformation in step c) for our input data often produces non-square images that are padded with zeros; these are the black bars on the sides of the Moon DEM in Figure 1.<sup>2</sup>

For each input image, we generate a corresponding ground-truth “output target” that is also  $256 \times 256$  pixels. Craters are encoded in the targets as rings with thicknesses of 1 pixel, and with radii and centers derived from craters’ physical locations and diameters in the catalog, described below. All target pixel values are binary, including at ring intersections. Any craters with diameter  $D_{\text{pix}} < 1$  pix are excluded. We experimented with other target formats, including density maps, and binary and non-binary filled circles. We found, however, that out of all of these formats binary ring targets were best reproduced by the CNN, particularly for situations with many overlapping craters. While the rings highlight crater edges, the CNN still uses information about the interior of craters to generate the ring masks. We tested this on a limited sample following the method of

---

<sup>2</sup>To check that this padding does not affect our results, we experimented with training our CNN on DEM images with all padding cropped out, and found our performance metrics (Table 3.1) differed by only a few percent. We thus conclude that it has a negligible effect.

Zeiler and Fergus (2013); Shallue and Vanderburg (2018) by perturbing the pixels of input DEMs, and found that we could prevent a crater detection by only increasing the pixels within its interior, without modifying its rim. This demonstrates that the interiors of craters are being used by our network.

The data used to construct the targets was obtained by merging two human-generated crater catalogs. For 5 – 20 km craters we used the global crater dataset assembled by Povilaitis et al. (2017) using the LRO Wide Angle Camera (WAC) Global Lunar DEM at 100 m/pixel (303 pixels/degree) resolution (GLD100; Scholten et al. 2012), and for > 20 km craters we used the global crater dataset assembled by Head et al. (2010) using the LOLA DEM with a resolution of 64 pixels/degree (472 m/pixel). Merging these two datasets in this way was intended by Povilaitis et al. (2017), who explicitly designed their dataset as a continuation of that of Head et al. (2010) to smaller sizes. The methods used to assemble these datasets are described in Head et al. (2010) and Povilaitis et al. (2017) and are typical for human generated datasets.<sup>3</sup>

We split our input image/output target pairs into three separate datasets to be used for training, validating and testing our CNN (see Section 2.7 for their uses). The three datasets are sampled from equal sized and mutu-

---

<sup>3</sup>During initial testing we also used the LU78287GT Lunar Crater Catalog (Salamunićcar et al., 2014), which was generated by a Hough Transform-based CDA. We transitioned to solely using human-generated catalogs to prevent the CNN from inadvertently learning the biases of another CDA.



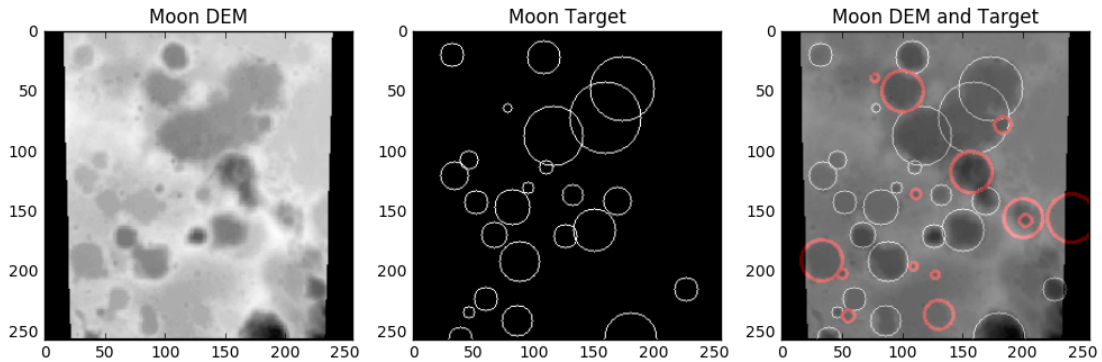


Figure 1: Sample Moon DEM image (left) and target (middle) from our dataset, with the two overlaid (right). Red circles (right panel) show features that appear to be craters but are absent from the Head et al. (2010) and Povilaitis et al. (2017) datasets, representing apparently missed classifications.

ally exclusive portions of the Moon, spanning the full range in latitude and  $-180^\circ$  to  $-60^\circ$ ,  $-60^\circ$  to  $60^\circ$  and  $60^\circ$  to  $180^\circ$  in longitude for the training, validation and test sets, respectively. Each dataset contains 30000 DEM images, and the median number of craters per DEM image is 21. Because we use a log-uniform sampling of crop lengths, all scales are equally represented in area, but the datasets contain far more DEM images with small crop lengths. An example input DEM image and target pair is shown in the left and middle panels of Figure 1.

Our human-generated crater dataset is incomplete and contains many apparent missed identifications. Fassett et al. (2012) estimated an incompleteness of 12% for the Head et al. (2010) dataset. The incompleteness

of the Povilaitis et al. (2017) dataset is unknown at this time, but appears to be comparable or higher. For the sample DEM image/target pair shown in Figure 1, we highlight features that appear to be craters but are missing from the Head et al. (2010) and Povilaitis et al. (2017) datasets as red circles in the right panel. It is unclear at this time why these craters were missed. In addition, the right panel of Figure 1 shows how our binary ring targets do not match the rims of non-circular craters. Together, these hinder the training of our CNN since genuine crater rims will be present in our dataset with no corresponding target rings, potentially confusing the CNN (see Section 4 for a discussion).

## 2.2. Convolutional Neural Networks (CNNs)

In this section we provide a brief, heuristic background to convolutional neural networks (CNNs) before introducing our network architecture. More in depth descriptions of the theory and mathematics of CNNs can be found in references such as Goodfellow et al. (2016) Chapter 9.

Machine learning algorithms in general can be thought of as universal function approximators, and neural networks (NNs) are a type of machine learning algorithm that uses “neurons” (loosely modelled after the human brain) to make such approximations. A neuron can be represented mathematically by:

$$y = f \left( \sum_j w_j x_j + b \right) \quad (1)$$

where  $x_j$  are a set of input values that are linearly combined with a set of weights  $w_j$ , and added to a bias offset  $b$ . This linear combination is then fed through a (typically non-linear) “activation function”  $f$ , which returns output  $y$ . Depending on the choice of  $f(z)$ , where  $z = \sum_j w_j x_j + b$ , a neuron is able to represent a number of simple functions, eg.  $f(z) = z$  for a line, or  $f(z) = 1 / (1 + \exp(-z))$  for a sigmoid function. Varying the weights  $w_j$  and bias  $b$  to best approximate a set of known  $y$  values given a corresponding set of  $x_j$  values is thus equivalent to linear regression when  $f(z)$  is linear, and logistic regression when  $f(z)$  is sigmoidal.

A neural network contains sets, or layers, of neurons  $\{y_i\}$ , where  $y_i$  represents the output of the  $i$ -th neuron in the layer. Neurons within each layer are independent of one another, but different layers are stacked on top of one another, with the output of one layer serving as the input to the next.<sup>4</sup> Input data is fed into the first layer, while the network’s “prediction”, or predicted output target, comes from the last. A network is “trained” by tuning the weights of all neurons in all layers so that the network best approximates a set of known, or ground-truth, targets given corresponding input. This tuning is typically done via backpropagation (Rumelhart et al., 1986), and goodness of fit is defined through a loss function, e.g. mean squared error between the known targets and network predictions. Much like adding terms to a Taylor series, increasing the number of layers and/or

---

<sup>4</sup>This is true of standard “feed-forward” neural networks. Other types of networks exist, e.g. recurrent neural networks, but are beyond the scope of this paper.

the number of neurons per layer allows the network to approximate functions of increasing complexity. However, this added complexity comes at the cost of more tunable parameters and the potential for errors to be amplified from one layer to the next; both make network optimization more difficult.

In computer vision problems, the input is typically pixel intensities from images, while the target is typically either a single class label for the image or another array of pixel intensities. In the latter case (which used in this work), the approximated function is a mapping from one type of image to another. Since the input and output are two-dimensional, we may represent the neuron in the  $i$ -th row and  $j$ -th column of one layer as

$$y_{ij} = f \left( \sum_k \sum_l w_{ij,kl} x_{kl} + b_{ij} \right) \quad (2)$$

where, for the NN's first layer,  $x_{kl}$  represents input image pixel intensities. Classical NNs, however, place no restrictions on  $w_{ij,kl}$ , and allow weights of different neurons in a layer to vary independently of one another. For images, this means any number of pixels in any orientation can be connected, and there is no guarantee the spatial information of the image is preserved from one layer to the next.

CNNs primarily differ from traditional NNs in that their neurons are only locally connected to the previous input volume (LeCun et al., 1989), and layers are structured to perform a discrete convolution between their

inputs and a kernel, or “filter”, which is represented by the weights. In the context of Equation 2, this means  $w_{ij,kl}$  is zero other than a few adjacent values of  $k, l$ , and the weights for one neuron are just index-shifted versions of the weights for another. Convolutional layers hence embed the spatial continuity of images directly into the architecture of the network. Also, since there are only a few non-zero weights that are shared between all the neurons, only a small number of weights need to be stored in memory and adjusted during training, simplifying network optimization. The weight-sharing also exploits the property that weights useful to one section of an image might also be useful to another.

The main components of CNNs relevant to our work are convolutional layers, pooling layers, and merge layers. Convolutional layers, the primary component of a CNN, contain neurons that convolve an input with a filter in the manner described above. The output of this convolution (i.e.  $y_{ij}$  in Equation 2) is called a “feature map”. In practice, a single convolutional layer can contain a large number of filters, each acting on the layer’s input to produce its own feature map. We make use of these “filter banks” in our CNN architecture. Pooling layers perform downsampling operations along the spatial dimensions, reducing the dimensionality of the image and thus the number of learnable weights needed for subsequent layers. Merge layers combine feature maps from convolutional layers of compatible spatial dimensions, facilitating complex connections within the CNN. Neither pooling nor merge layers have trainable parameters, and are not represented

by Equation 2.

In recent years CNNs have demonstrated impressive performance on image-related tasks, including classifying individual pixels (aka “segmentation” in computer vision terminology) (Long et al., 2015). A particularly successful network for pixel-wise classification (image to image mapping) is the UNET architecture (Ronneberger et al., 2015), which was originally designed for biomedical segmentation problems. A novel aspect of the UNET architecture is the use of numerous “skip connections” which merge deep and shallow layers together, providing both spatial and semantic classification information for future convolutions.

### 2.3. CNN Architecture

In this work we implement a custom version of the UNET architecture (Ronneberger et al., 2015), shown in Figure 2.<sup>5</sup> This architecture consists of a contracting path (left side) and expansive path (right side), joined through multi-level skip connections (middle). Lunar DEM images are input to the contracting path and predictions are output from a final layer following the expansive path. Unless otherwise stated, all convolutional layers have banks of filters that each apply 3x3 padded convolutions followed by a rectified linear activation unit (ReLU; e.g. Goodfellow et al. 2016 Chapter 6.1), whose functional form is  $f(z) = \max(0, z)$ .

---

<sup>5</sup>In the early stages of this work, we attempted to count the craters of an image, rather than localize them, using a traditional CNN regressor. However, that model’s resulting accuracy was low, motivating our shift to the UNET-based one.

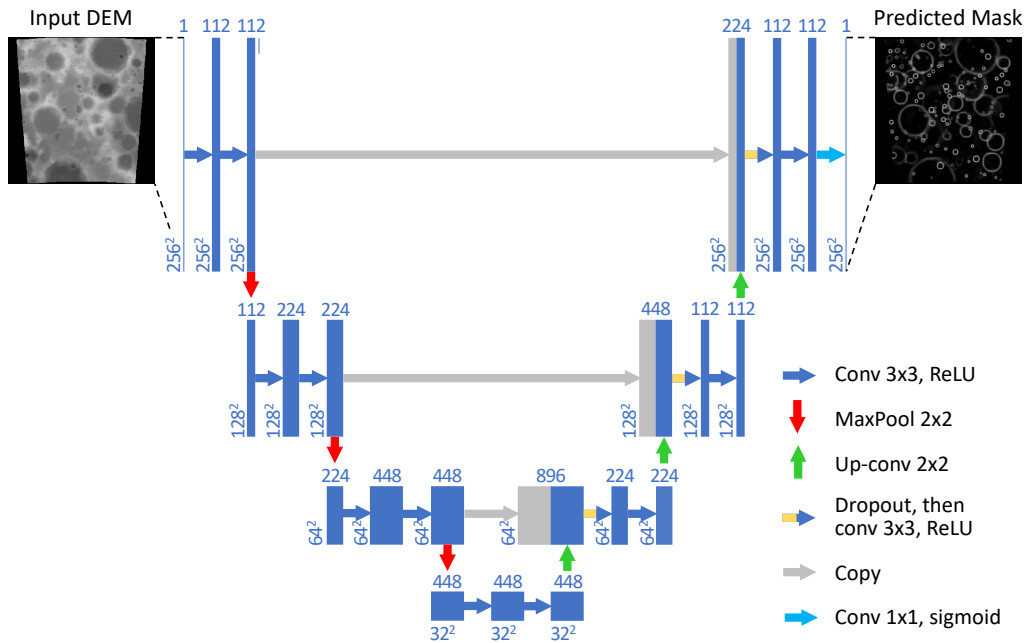


Figure 2: Convolutional neural network (CNN) architecture, based on UNET (Ronneberger et al., 2015). Boxes represent cross-sections of sets of square feature maps. For each set, its maps' dimensions are indicated on its lower left, and its number of maps are indicated above it. Half-grey boxes represent sets for which half of their maps are copied. The leftmost map is a  $256 \times 256$  grayscale image sampled from the digital elevation map, and the rightmost the CNN's binary ring mask prediction. Arrows represent operations, specified by the legend - notably, blue arrows represent convolutions, while gray ones represent copying (skip connections).

The contracting and expansive paths each contain 3 convolutional blocks. A block in the contracting path consists of two convolutional layers followed by a max-pooling layer with a 2x2 pool size. A block in the expansive path consists of a 2x2 upsampling layer, a concatenation with the corresponding block from the contracting path (i.e. a merge layer), a dropout layer (Srivastava et al., 2014), and two convolutional layers. The connecting path consists of two convolutional layers. Lastly, the final output layer is a 1x1 convolutional layer with a sigmoid activation and a single filter to output pixel-wise class scores. In the contracting path, each convolutional layer in blocks 1, 2 and 3 contain 112, 224 and 448 filters, respectively, while in the expansive path blocks 5, 6 and 7 contain 224, 122 and 122, respectively. Each convolutional layer in the connecting path contains 448 filters. Our CNN differs from the original UNET by the number of filters in each convolution layer (which we selected for the model to fit into GPU memory) and the use of dropout in the expansive path.

#### 2.4. Crater Extraction

A  $256 \times 256$  DEM image passed through the CNN will output a  $256 \times 256$  target with activated pixels corresponding to the locations of the crater rims. However, the CNN does not explicitly extract crater position and size from these rims. Instead, this must be done separately using a custom pipeline that relies heavily on the `match_template` algorithm from



`scikit-image` (Van der Walt et al., 2014)<sup>6</sup>. This algorithm iteratively slides generated rings through the targets and calculates a match probability at each  $(x, y, r)$  coordinate where  $(x, y)$  is the centroid of the generated ring and  $r$  is the radius.

Our custom crater extraction pipeline is as follows. For each CNN-predicted target we apply a binary threshold  $B$  such that pixel intensities greater than  $B$  are set to 1 and those otherwise are set to 0. We then apply Scikit’s `match_template` algorithm over a radius range  $r_{\min}$  to  $r_{\max}$  and classify any  $(x, y, r)$  ring with a match probability greater than  $P_m$  as a crater. Two craters  $i$  and  $j$  that fulfill the following criteria are flagged as duplicates if they satisfy both of the following conditions:

$$\begin{aligned} ((x_i - x_j)^2 + (y_i - y_j)^2) / \min(r_i, r_j)^2 &< D_{x,y} \\ \text{abs}(r_i - r_j) / \min(r_i, r_j) &< D_r \end{aligned} \tag{3}$$

where  $D_{x,y}$  and  $D_r$  are tunable hyperparameters. For duplicates we keep only the crater with the highest match probability  $P_m$ . This process is repeated for all CNN-predicted targets.

As is standard practice in machine learning, our hyperparameters are tuned by measuring the performance of various combinations on the vali-

---

<sup>6</sup>Although template matching is an expensive technique done in a brute force manner we found it far more accurate than others including the Hough transform (Duda and Hart, 1972) and Canny edge detection (Canny, 1986) with enclosed circle fitting.

validation data and picking the optimal set. After training our CNN (see Section 2.7), we perform a randomly sampled grid search of size 210 on the validation data over the following hyperparameter ranges:

$$B = [0.05, 0.1, 0.15]$$

$$P_m = [0.3, 0.4, 0.5, 0.6, 0.7]$$

$$D_{x,y} = [0.6, 0.8, 1.0, 1.2, 1.4, 1.6, 1.8, 2.0, 2.2]$$

$$D_r = [0.2, 0.4, 0.6, 0.8, 1.0, 1.2]$$

We find  $B = 0.1$ ,  $P_m = 0.5$ ,  $D_{x,y} = 1.8$  and  $D_r = 1.0$  yields the optimal  $F_1$  score (see Equation 9) of 0.74. We set  $r_{\min} = 5$  to minimize errors (see Section 2.8) and set  $r_{\max} = 40$ .

### 2.5. Post-Processing

Since our dataset contains DEM images with a log-uniform distribution of magnifications and uniform distribution of locations, a single crater will appear on average in  $120 \pm 30$  different DEM images. This increases the likelihood of detection but also yields many duplicates across targets which must be filtered. Therefore, in a final post-processing step we aggregate crater detections across all targets, convert from pixel coordinates to degrees and kilometers, and filter out duplicates.

Using the known properties of each DEM image, craters are converted from pixel coordinates  $(x, y, r)$  to degrees and kilometer coordinates  $(\mathcal{L},$

$L, R$ ):

$$\begin{aligned}
L - L_0 &= \frac{\Delta L}{\Delta H} (y - y_0) \\
\mathcal{L} - \mathcal{L}_0 &= \frac{\Delta L}{\cos\left(\frac{\pi L}{180^\circ}\right) \Delta H} (x - x_0) \\
R &= r \frac{C_{KD} \Delta L}{\Delta H},
\end{aligned} \tag{4}$$

where  $\mathcal{L}$  and  $L$  are the crater's longitude and latitude centroid, subscript 0 values are those for the center of the DEM image,  $\Delta L$  and  $\Delta H$  are the latitude and pixel extents of the DEM image along its central vertical axis (where  $\mathcal{L} = \mathcal{L}_0$ ), excluding any padding, and

$$C_{KD} = \frac{180^\circ}{\pi R_{\text{Moon}}} \tag{5}$$

is a kilometer-to-degree conversion factor, where  $R_{\text{Moon}}$  is the radius of the Moon in km.

We then employ a similar filtering strategy as Section 2.4, classifying craters  $i$  and  $j$  as duplicates if they satisfy both of the following conditions:

$$\begin{aligned}
\frac{((\mathcal{L}_i - \mathcal{L}_j)^2 \cos^2\left(\frac{\pi}{180^\circ} \langle L \rangle\right) + (L_i - L_j)^2)}{C_{KD}^2 \min(R_i, R_j)^2} &< D_{\mathcal{L},L} \\
\frac{\text{abs}(R_i - R_j)}{\min(R_i, R_j)} &< D_R.
\end{aligned} \tag{6}$$

where  $\langle L \rangle = \frac{1}{2}(L_i + L_j)$ .  $D_R$  and  $D_{\mathcal{L},L}$  are hyperparameters which, like in

Section 2.4, we tune by performing a grid search on the validation dataset after training our CNN and tuning its hyperparameters, this time sampling every combination from:

$$D_{\mathcal{L},L} = [0.6, 1.0, 1.4, 1.8, 2.2, 2.6, 3.0, 3.4, 3.8]$$

$$D_R = [0.2, 0.6, 1.0, 1.4, 1.8, 2.2, 2.6, 3.0, 3.4, 3.8]$$

We find  $D_{\mathcal{L},L} = 2.6$  and  $D_R = 1.8$  yields the optimal  $F_1$  score (see Equation 9) of 0.67.

## 2.6. Accuracy Metrics

To train our network we use the pixel-wise binary cross-entropy,  $\ell$  (Abadi et al., 2016; Chollet, 2015), a standard loss function used for segmentation problems:

$$\ell_i = x_i - x_i z_i + \log(1 + \exp(-x_i)) \quad (7)$$

where  $z_i$  is the ground-truth output target value of pixel  $i$  and  $x_i$  is the CNN-predicted one.

To optimize the hyperparameters in our crater extraction and post-processing routines (Sections 2.4 and 2.5) we use the precision,  $P$ , and recall,  $R$ , to measure accuracy, calculated according to:

$$\begin{aligned}
P &= \frac{T_p}{T_p + F_p} \\
R &= \frac{T_p}{T_p + F_n}
\end{aligned}
\tag{8}$$

where  $T_p$  are true positives,  $F_p$  are false positives and  $F_n$  are false negatives. There is always a trade-off between precision and recall. For example, a machine that only classifies craters when extremely certain will have a high precision but low recall, while a machine that classifies craters when only moderately certain will have a higher recall but lower precision. A common single-parameter metric that balances precision and recall is the  $F_1$  score:

$$F_1 = 2 \frac{PR}{P + R}
\tag{9}$$

Implicitly encoded into our accuracy metrics is the assumption that our ground-truth datasets of Head et al. (2010) and Povilaitis et al. (2017) are complete. However, as mentioned in Section 2.1 this is incorrect. As a result, genuine new craters identified by our CNN will be interpreted as false positives, penalizing these loss functions vs. improving them. This is an unavoidable consequence of using an incomplete ground truth (see Section 4 for further discussion).

To measure the accuracy of identified craters (Sections 2.4 and 2.5) we calculate the fractional errors in longitude,  $\mathcal{L}$ , latitude,  $L$ , and radius,  $R$ , according to:

$$\begin{aligned}
d\mathcal{L}/R &= \text{abs}(\mathcal{L}_P - \mathcal{L}_G) \cos(\pi \langle L \rangle / 180^\circ) / (R_G C_{KD}) \\
dL/R &= \text{abs}(L_P - L_G) / (R_G C_{KD}) \\
dR/R &= \text{abs}(R_P - R_G) / R_G
\end{aligned}
\tag{10}$$

where subscript  $P$  corresponds to our CNN-predicted craters, subscript  $G$  corresponds to our ground-truth craters and  $\langle L \rangle = \frac{1}{2}(L_P + L_G)$ .

Finally, our pipeline discovers thousands of new craters (as will be shown in Section 3). We measure the new crater percentage,  $P$ , according to:

$$P = N / (N + G) \tag{11}$$

where  $N$  is the number of CNN-predicted craters without a corresponding match from the ground truth (i.e. they are either genuine new craters or false positives), and  $G$  is the number of ground-truth craters. A “match” between a CNN-predicted and ground-truth crater is determined via Eq. 3 and Eq. 6 for post-CNN (Section 2.4) and post-processed (Section 2.5) craters, respectively.

## 2.7. Training

A recurring theme in machine learning is “overfitting”, which occurs when a model latches onto overly-complex and/or irrelevant features during training. Overfit models typically achieve high accuracy on the training set but low accuracy (poor generalization) on new data. Many algo-

rithms control for overfitting by penalizing overly-complex models, retaining only the essential characteristics from the training data that will generalize to new examples. This penalization is generally mediated through the model’s hyperparameters, which control the model’s complexity. For our CNN, such hyperparameters include weight regularizations for each convolutional layer, the learning rate, dropout layers after each merge layer, filter size, and depth of the network. These hyperparameters are tuned on a separate validation dataset, forcing the model to achieve high accuracy on two different datasets.

After training the model and tuning the hyperparameters a final evaluation is conducted on the test set, another dataset distinct from both the training and validation datasets. If the model achieves comparable accuracy on the test set as the training and validation sets, it is likely that minimal overfitting has occurred and the model should generalize well to new examples. We also address overfitting through a custom image augmentation scheme that randomly flips, rotates and shifts DEM images (and their corresponding targets) before they are used to train the CNN. This augments the effective dataset size and minimizes the chance of the CNN generating features related to image orientation.

We tune the hyperparameters of our model by training 60 models with randomly chosen hyperparameters over standard ranges on the training set and selecting the model with the best binary cross-entropy score (Equation 7) on the validation set. Defining an “epoch” as a single pass through

the entire training set and “batch size” as the number of examples seen per backpropagation gradient update, each model is trained for 4 epochs with a batch size of 8 using the ADAM optimizer (Kingma and Ba, 2014). The hyperparameters of our best model are weight regularization =  $10^{-5}$ , learning rate =  $10^{-4}$ , dropout = 15%,  $3 \times 3$  filter sizes, and a depth of 3.

### 2.8. Errors

A few sources of error affect the final extracted  $(\mathcal{L}, L, R)$  coordinates of each detected crater. First, craters can only be detected in pixel increments, and converting from pixels to degrees yields a quantization error,  $E_q$ , of:

$$E_q = C_{\text{offset}} \frac{\Delta L}{\Delta H}, \quad (12)$$

where  $C_{\text{offset}} \leq 1$  is a constant of order unity representing typical sub-pixel offsets. Setting  $C_{\text{offset}} = 1$ , and considering our largest DEM images (6500 pixels) where  $\Delta L \approx 25^\circ$ , we find a maximum quantization error of  $E_q \approx 0.1^\circ$ , or  $\sim 3$  km. In principle this error could be reduced by increasing the pixel resolution of each DEM image, though doing so would be memory-intensive.

Second, objects within an orthographic projection become more distorted further from the central longitude and latitude  $(\mathcal{L}_0, L_0)$ , which changes the size of smaller craters, and introduces non-circular deformations in larger ones. Along the central vertical axis, the deviation of the distorted



radius from our estimated value using Equation 4 is

$$F_{rd} = \frac{R_{\text{distorted}}}{R} = \frac{1}{\cos\left(\frac{\pi}{180^\circ}(L - L_0)\right)}. \quad (13)$$

For our largest DEM images,  $L - L_0 \approx 12^\circ$ ,  $F_{rd} \approx 1.02$ , so deviations are at most 2% of a crater’s radius.

Third, the longitude and latitude estimates in Equation 4 neglect higher order terms (including the distortion described above) and cross-terms containing both  $\mathcal{L}$  and  $L$ . To quantify this effect we passed the crater pixel positions from the ground-truth test dataset through Equation 4 to obtain  $\mathcal{L}_C$  and  $L_C$ . We then subtracted the ground truth’s longitude and latitude values from  $\mathcal{L}_C$  and  $L_C$ , respectively, and normalized by the longitude/latitude extent of our DEM images. We found median relative offsets of 0.13% and 0.28% in longitude and latitude, but, for the largest DEM images, maximum relative offsets can reach 1.0% in longitude and 1.9% in latitude. For a 6500 pixel DEM image, this translates to  $0.25^\circ$  in longitude and  $0.5^\circ$  in latitude. We also calculated the fractional error using Equation 10, replacing  $\mathcal{L}_P$  and  $L_P$  with  $\mathcal{L}_C$  and  $L_C$ , and find median fractional errors of 3% in longitude and 5% in latitude.

To help offset these errors we impose a minimum search radius,  $r_{\min} = 5$ , for our crater extraction pipeline (Section 2.4). This prevents quantization and projection errors from ever being a significant fraction of the crater’s radius. This comes at a cost of not being able to probe the smallest

craters in each DEM image, yielding fewer new crater detections than we otherwise would obtain.

### 3. Results

#### 3.1. Crater Identification on the Moon

We apply our trained CNN and optimized crater identification pipeline on the test set and list our various accuracy metrics in Table 3.1 for the validation and test sets. “Post-CNN” statistics were generated on an image-by-image basis after Section 2.4 of the pipeline with an averaged mean and standard deviation taken across all predicted targets. “Post-processed” statistics were generated after Section 2.5 of the pipeline, and hence represent our final crater distribution after combining extracted craters from each target into one distribution and removing duplicates. Together, these statistics convey how our pipeline is performing at various stages.

The similarity between our validation and test set statistics in Table 3.1 implies that little to no overfitting has occurred. Our post-processed test recall is 92%, recovering almost all craters from the test set. By comparison, our post-CNN test recall is lower at  $57\% \pm 20\%$ , meaning that (on average) our CNN detects only half of the craters per target. The drastic difference between post-processed and post-CNN recalls demonstrates the effectiveness of aggregating crater detections across multiple images and scales. A major reason for our low post-CNN recall is our CNN does not reliably detect craters with radii greater than  $\sim 15$  pixels (see Section 4 for

Accuracy Metric	Post-CNN (Validation)	Post- Processed (Validation)	Post-CNN (Test)	Post- Processed (Test)
Recall	56% $\pm$ 20%	92%	57% $\pm$ 20%	92%
Recall ( $r < 15$ pixels)	83% $\pm$ 16%	–	83% $\pm$ 13%	–
Precision	81% $\pm$ 16%	53%	80% $\pm$ 15%	56%
New Crater Percentage	12% $\pm$ 11%	45%	14% $\pm$ 13%	42%
False Positive Rate	–	–	–	11% $\pm$ 7%
Frac. longitude error	10% <sup>+2%</sup> <sub>–2%</sub>	13% <sup>+10%</sup> <sub>–7%</sub>	10% <sup>+2%</sup> <sub>–2%</sub>	11% <sup>+9%</sup> <sub>–6%</sub>
Frac. latitude error	10% <sup>+3%</sup> <sub>–2%</sub>	10% <sup>+8%</sup> <sub>–5%</sub>	10% <sup>+2%</sup> <sub>–2%</sub>	9% <sup>+7%</sup> <sub>–5%</sub>
Frac. radius error	8% <sup>+2%</sup> <sub>–2%</sub>	6% <sup>+5%</sup> <sub>–3%</sub>	8% <sup>+1%</sup> <sub>–1%</sub>	7% <sup>+5%</sup> <sub>–4%</sub>

Table 1: Accuracy metrics on the validation and test sets. “Post-CNN” statistics were generated after Section 2.4 of the pipeline with a mean and standard deviation taken across targets, while “post-processed” statistics were generated after Section 2.5 of the pipeline, after combining extracted craters into a single global distribution. Precision and recall are calculated according to Eq. 8, new crater percentage according to Eq. 11, fractional longitude, latitude and radius errors according to Eq. 10 with a median and interquartile range (IQR) taken across all detections. The false positive rate of new craters is estimated by four different scientists classifying 361 new craters and averaging the results. We note that precision drops at the post-processed stage because many new craters (absent in the ground truth) are identified.

a discussion). Restricting to craters with a pixel radius  $r$  less than 15 pixels, our post-CNN test recall improves to  $83\% \pm 16\%$ . Wetzler et al. (2005) estimated a human recall of 75% when re-classifying crater images, making our post-CNN recall consistent with human performance for  $r < 15$  pixels.

42% of post-processed test craters are new, almost doubling our catalog, with 15% of them having diameters under 5 km (i.e. below the limits of our ground-truth catalogs). Our estimated false positive rate of these new post-processed craters is  $11\% \pm 7\%$ , which was estimated by four scientists from our research group each classifying the same 361 new craters re-projected onto their original DEM images and averaging the results. This procedure allowed the scientists to classify the new craters under the same conditions as our identification pipeline. These manually classified craters along with their corresponding Moon DEMs, ground-truth targets and CNN predictions are publicly available at <https://doi.org/10.5281/zenodo.1133969>. Although individual false positive estimates differed between scientists, this is in line with previous research (e.g. Robbins et al., 2014) that large disagreements in human crater classification is common. For post-CNN,  $14 \pm 13\%$  of test craters per DEM image are new. As a result of these new crater detections, our post-CNN and post-processed precisions are low since new craters are interpreted by the precision metric exclusively as false positives (see Section 4 for a discussion).

Figure 3 compares our post-processed craters (top left) to the ground truth (top right) for a large swath of the Moon (bottom left) from the test

set. Blue circles represent post-processed craters that were successfully matched to the ground truth (and vice versa), red circles represent new crater detections from our pipeline (without a corresponding ground-truth match), and purple circles represent ground-truth craters missed by our pipeline. As can be visually seen, our pipeline recovers many more craters than the ground truth, with overall few false positives and duplicates. Our median post-processed and post-CNN fractional errors in longitude, latitude and radius are 11% or less, representing overall good agreement with the ground truth despite the sources of error mentioned in Section 2.8.

### 3.2. Lunar Crater Size Distribution

In Figure 4 we show the cumulative size-frequency distribution (CSFD) of the craters recovered by our CNN as compared with the CSFD of craters in the human counted set. This is constructed following the recommendations in NASA Technical Memorandum (79730, 1978). In this case the surface area,  $A$ , used to normalise the number of craters is the area of the Moon between latitudes of  $60^\circ$  North and South. We note that another widely used method to visualize crater distributions is R plots (as for example used in Head et al. (2010)), however we refrain from using these due to their significant dependency on data binning that can be completely avoided using CSFDs, where each crater is its own bin (Weaver et al., 2015).

As we can see, between roughly 5 and 20 km diameters the CNN de-

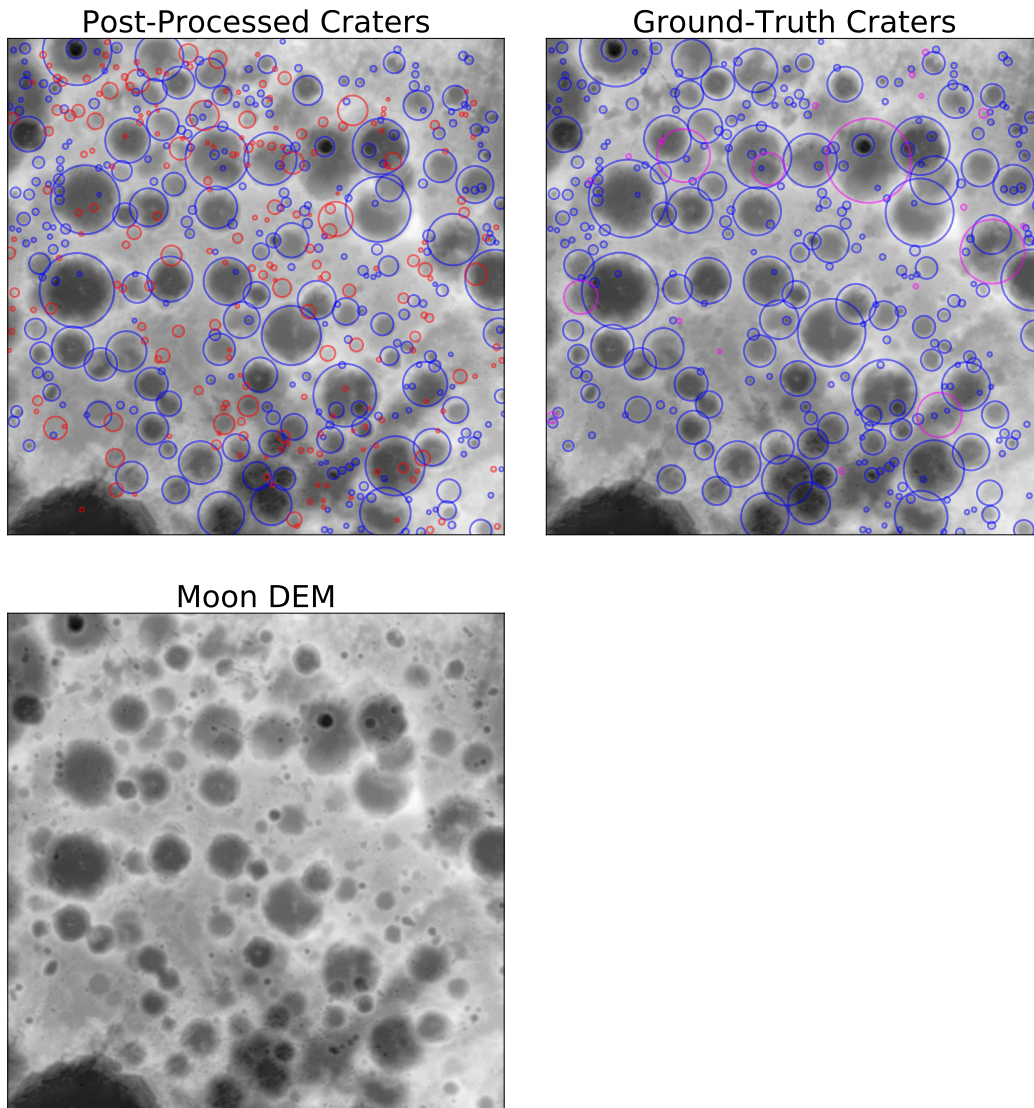


Figure 3: Sample patch of the Moon from the test set (lower left), with post-processed (top left) and ground-truth (top right) craters overlaid. Blue circles represent post-processed craters that were successfully matched to the ground truth (and vice versa), red circles represent new crater detections from our pipeline (without a corresponding ground-truth match), and purple circles represent ground-truth craters missed by our pipeline.

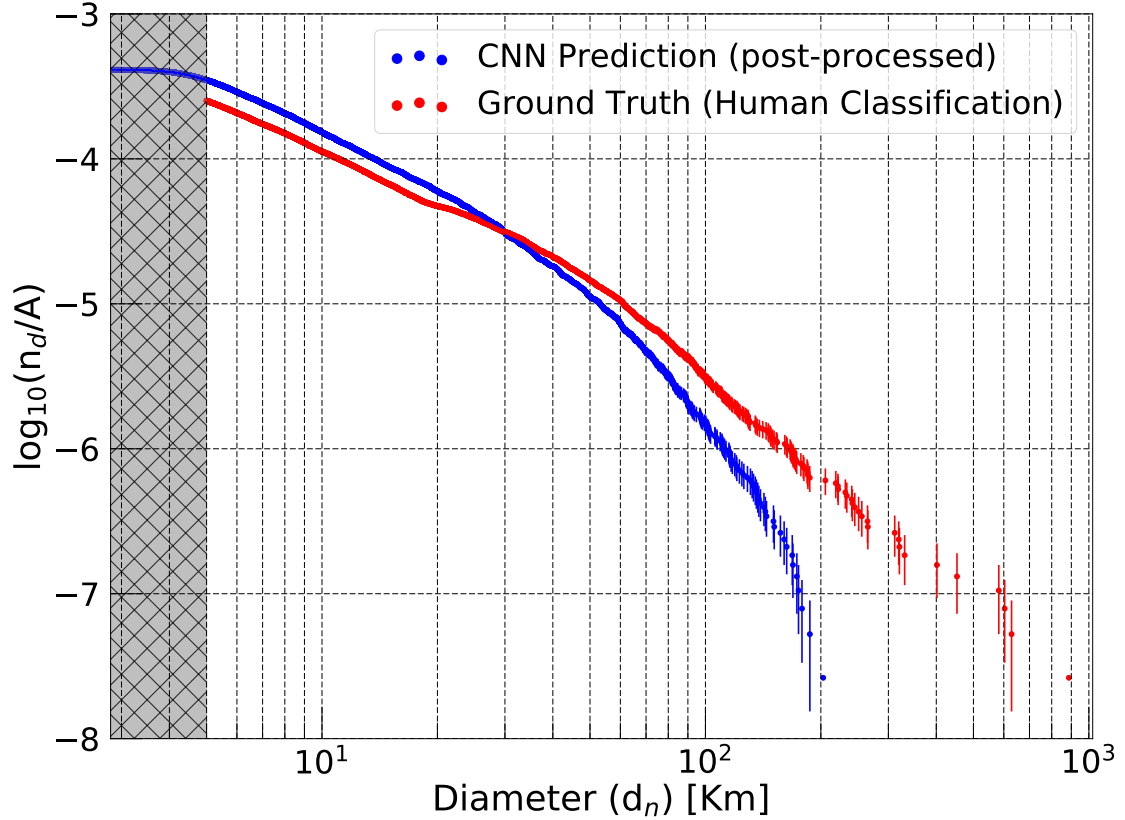


Figure 4: Lunar crater size-frequency distributions represented as CSFD plots. Red is the human-generated test set (the ground truth), compared to what our CNN predicts post-processing (blue). Our CNN recover the same slope as the GT for craters smaller than 20 km while detecting significantly more craters in this range. The shaded region inside 5 km however should not be physically interpreted due to data incompleteness.

rived CSFD is systematically higher than the human derived CSFD, but runs essentially parallel to it. This indicates that while the CNN finds a substantially larger number of craters than the human crater counters, the craters newly identified by the CNN follow the same size distribution as the human identified craters, which reassures us that they are likely to be real. At 20 km however there is a slight kink and upturn in the human dataset whereas the CNN prediction continues smoothly. The CNN prediction then curves downward gradually before cutting off more sharply at diameters of around 200 km, resulting in increasing divergence between the CNN prediction and the human classifications for these larger craters.

There are several factors that might contribute to the divergence between the CSFDs of the CNN prediction and the human classified dataset at large crater sizes. Firstly, as we describe in detail in Section 4, our CNN can struggle with identifying craters that have a radius of more than around 15 pixels in the input images. While this is somewhat ameliorated by sampling images at many different magnification levels (such that an  $r > 15$  pix crater in one image will have  $r < 15$  pix in another) it can still slightly hinder the detection of craters with diameters larger than around 50 km. This is likely responsible for some of the downward curvature of the CNN CSFD at larger sizes and is probably the cause of the cut-off at  $\sim 200$  km.

The location of the kink in the human classified CSFD at which the mismatch between the CNN prediction and the human classification begins is also notable. This kink appears at the boundary between the two compo-



ment datasets of our human classified ground truth, that of Povilaitis et al. (2017) and Head et al. (2010). It is possible that there is a systematic difference between the crater classifications of the two groups. Indeed some level of difference would be unsurprising given the findings of Greeley and Gault (1970), Kirchoff et al. (2011) and Robbins et al. (2014). The systematic offset between the CSFDs of our CNN predictions and the human classifications in the range where the two are parallel already indicates that the human datasets are not complete, as was also found by Fassett et al. (2012) for the Head et al. (2010) dataset. Since the Povilaitis et al. (2017) dataset contains around four times as many craters it is inevitable that if there are any systematic differences in the human classification that lead to size-dependent completeness effects, the CNN will tend to follow Povilaitis et al. (2017).

Another consideration is that 20 km is roughly the diameter at which craters on the Moon transition from simple to complex (e.g. Pike, 1980; Grieve, 1981; Stoffler, 2006). Above this size crater morphologies change, with crater floors becoming flatter and the appearance of features like central peaks. It is conceivable that this morphological transition could hinder the detection of larger craters by the CNN. In particular, since smaller, simple craters dominate the human dataset, if the morphological differences make it difficult to simultaneously optimise to detect both simple and complex craters it is inevitable that the CNN will favour simple craters, just as it would favour Povilaitis et al. (2017) over Head et al. (2010) in the case

of differences in their counting methods. In Section 4 we discuss how one could attempt to disentangle these possible effects.

Below around 5 km diameter the CNN prediction begins to roll over. This is due to incompleteness in the sampling of the lunar surface at the smallest scales/largest image magnifications. Despite incompleteness in the sampling and the lack of craters in the training set at these small sizes the CNN still finds many craters  $<5$  km diameter.

### *3.3. Transfer Learning on Mercury*

Domain shift is a common problem in machine learning that arises when a model is used to predict on data with different properties than its training set, and typically results in decreased performance. We briefly evaluate the sensitivity to domain shift for our network by taking our Moon-trained CNN and transfer learning to Mercury. Mercury has different properties than the Moon, including a different gravitational acceleration, surface composition, terrain, and impact history. In addition, we also use the Mercury MESSENGER Global DEM with a resolution of 64 pixel/degree, or 665 m/pixel (Becker et al. 2016; available at USGS Astrogeology Science Center 2016), which has different image properties than our Moon DEM. All these affect the distribution and appearance of impact craters.

To evaluate our CNN on Mercury, we prepared DEM images from the Mercury MESSENGER Global DEM in a similar manner as described in Section 2.1 (except that we do not use a corresponding human-generated

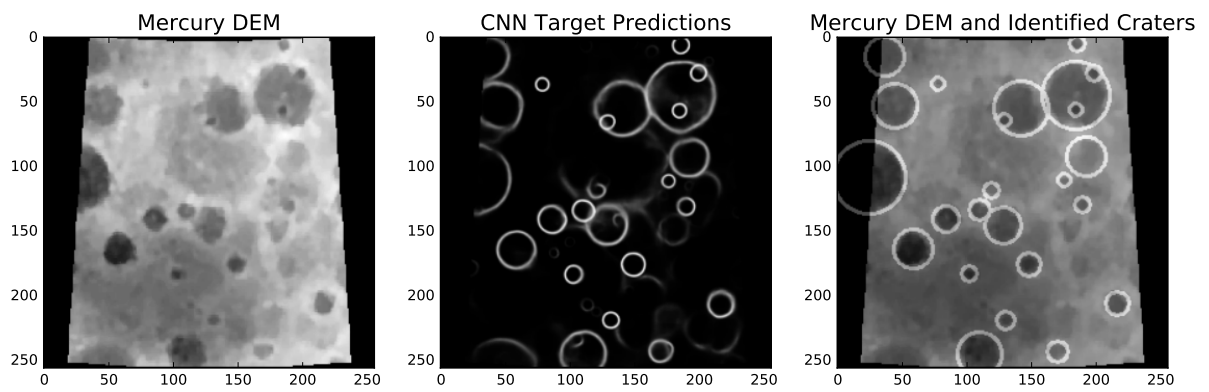


Figure 5: A sample Mercury DEM (left), CNN target predictions (middle), and post-processed identified craters overlaid on the original DEM image (right).

crater catalog). We then passed these DEM images through our Moon-trained CNN with no alterations to the architecture or weights. Figure 5 shows a sample input DEM image of Mercury (left), CNN target predictions (middle), and post-processed identified craters overlaid on top of the original DEM image (right). Comparing the left panels of Figure 1 and Figure 5, some differences are visible between Moon and Mercury craters, yet our CNN is able to detect both types. In addition, as shown in Figure 5, the CNN correctly identified a significant number of the crater-like features on Mercurys surface. Moreover, manual inspection on that patch of Mercury shows that almost all of the inferences made by the CNN do seem to be real craters. This demonstrates its efficiency in distinguishing craters from other terrain features. Simple edge detection techniques would not make such a distinction. While humans are often very good at transfer learning, and thus it may seem simple from a human perspective, transfer learning has commonly been challenging in machine learning. That our model seems to transfer well is thus greatly encouraging. While this demonstrates successful generalization, we leave a thorough analysis of transfer learning to future work.

#### **4. Discussion**

There are many reasons to believe our CNN has learned the complex features that define a crater. First, despite the Moon’s large structure variations across its surface our CNN was able to recover 92% of craters on a

face previously unseen by our CNN. Second, the similarity in accuracy metrics between the validation and test sets in Table 3.1 implies that minimal overfitting has occurred, and the CNN has indeed learned useful, generalizable features. Third, 42% of extracted Lunar craters are new, and from human validation of a subset most appear to be genuine matches. Fourth, our Moon-trained CNN successfully detected craters on Mercury, a surface completely distinct from any specific region on the Moon. Finally, while simple edge detection techniques would activate non-crater features like mountains, ridges, etc., our CNN almost exclusively activates crater rims (e.g. see middle panel of Figure 5).

As mentioned in Section 2.1 and shown in Figure 1, our training data is incomplete, containing many missed craters as well as target rings that differ from true crater rims. Despite these shortcomings our CNN was still able to understand our training objective and correlate the binary ring targets with the true rims of the craters. Proof of this can be seen in the middle panel of Figure 5, where some CNN-predictions are non-circular and better match the true crater rims than a circular ring could<sup>7</sup>. Together, these highlight the robustness and flexibility of deep learning solutions.

A fundamental difficulty when using an incomplete dataset is tuning the hyperparameters. Under this regime genuine new crater detections are

---

<sup>7</sup>To be clear, the  $256 \times 256$  pixel CNN target predictions can produce non-circular ring boundaries, as shown in the middle panel of Figure 5. However, extracted post-processed craters (Section 2.4 and Section 2.5) do not retain non-circularity information, as shown in the right panel of Figure 5.

interpreted as false positives, penalizing the precision metric and artificially lowering the  $F_1$  score (which we are trying to maximize). Since thousands of new craters were detected, the  $F_1$  score, which favors hyperparameters that yield the fewest new crater detections whilst still maintaining a high recall, is reduced. As a result, our tuning procedure yields hyperparameters that are likely conservative compared to if we had focused on finding new craters or had a more complete ground truth. The same principle applies when using the binary cross-entropy loss to train our CNN, yielding a final model that likely generates more conservative crater predictions. Followup work that uses a more complete ground-truth crater distribution would presumably yield improved results.

Our CNN robustly detects craters from each DEM image with radii below 15 pixels, but tends to miss larger craters. An example of this can be seen in Figure 5, where a large crater with coordinates ( $x = 130, y = 78, r = 42$  pix) is only partially recognized by the CNN and missed by our crater extraction pipeline. We believe that this largely stems from the scale that is imposed when using small 3x3 filters in our convolutional chain, yielding a receptive field that is too small for large craters. Larger filter sizes were attempted, but this dramatically increases both the number of trainable parameters and network size, making model-optimization more difficult. Dilated Convolutions (Yu and Koltun, 2015), larger convolution strides, and/or deeper networks are possible avenues for improvement. However, increasing the receptive field likely accomplishes the same effect as reduc-

ing the magnification of an image, which we have already shown to be a successful remedy, achieving a post-processed recall of 92% (see Table 3.1).

In addition, as we noted in Section 3.2, there are potential external effects that may be feeding into the apparent size dependence of the ability of the CNN to recover larger craters. One possibility is that there is a systematic difference in the two human counted datasets that we have stitched together to form our ground truth. Another is that the physical transition between simple and complex craters makes it difficult for the CNN to adapt to detect both with equal efficiency. In both cases the CNN would tend to adapt itself to better identify craters in the larger population, which is the Povilaitis et al. (2017) dataset that consists of craters below the  $\sim 20$  km simple to complex transition. To disentangle these two possible effects requires a uniformly generated human dataset of both simple and complex craters. We would then be able to train the CNN on DEMs of solely simple craters, solely complex craters, or both and then test it on DEMs of all physical scales/crater morphologies. This requires a substantial new human counting effort and so we leave it for future work.

Our estimated post-processed false positive rate of new craters is  $11\% \pm 7\%$ , which, although generally low, is likely too high for our catalog to be used to produce high-precision crater catalogs. Our primary false positives are a) ambiguous circular looking depressions that may or may not be true craters (further analysis required), and b) overlapping craters that activate enough pixels in the region to breach the match probability threshold  $P_m$ ,

creating a “ghost” detection in our crater extraction pipeline (Section 2.4). In addition, Table 3.1 shows that roughly 25% of post-processed craters have coordinates that differ from the ground truth by 20% or more, and examples of this can be seen in Figure 3. This higher-error tail is not present in the post-CNN errors, so they arise from the post-processing methods, whose sources of error are detailed in Section 2.8. These issues indicate the need for further refinements to our overall crater identification pipeline in order to produce precision crater catalogs, which we save for future work.

## 5. Conclusions and Future Work

In this work we have demonstrated the successful performance of a convolutional neural network (CNN) in recognizing Lunar craters from digital elevation map images. In particular, we recovered 92% of craters from our test set, and almost doubled the number of total crater identifications. Furthermore, we have shown that our Moon-trained CNN can accurately detect craters on the substantially different DEM images of Mercury. This implies that the CNN has learned to robustly detect *craters*, and not features particular to the planetary surface on which it was trained.

Two primary advantages of a deep learning solution over human crater identification are consistency and speed. A CNN will classify an image identically each time, but the same is not true for humans (Wetzler et al., 2005). In addition, different humans will use slightly different criteria, which adds to the error budget (Robbins et al., 2014). Once trained, our



CNN greatly increases the speed of crater identification, taking minutes to generate predictions for tens of thousands of Lunar DEMs and a few hours to extract a post-processed crater distribution from those DEMs. This is of course all done passively, freeing the scientist to do other tasks. Our CNN could also be used to assist human experts, generating initial suggestions for the human expert to verify.

DEMs are available for many other Solar System bodies, including Mercury (Becker et al., 2016), Venus (Magellan Science Team, 2014), Mars (Fergason et al., 2017), Vesta (Preusker et al., 2014) and Ceres (Preusker et al., 2016). It will be interesting to study to what extent our CNN can transfer-learn to other Solar System bodies with a DEM, possibly facilitating a systematic, consistent, and reproducible crater comparison across Solar System bodies. While we have successfully shown transfer-learning from our Moon-trained CNN to Mercury, a detailed analysis for Mercury has been left to future work.

Our current work detected craters down to roughly 3 km diameter, but since our CNN accepts images of arbitrary magnification we can transfer-learn to kilometer and sub-kilometer scales on the Moon. We anticipate that the uncharted territory of systematic small-size crater identification will provide important new information about the size distribution of Lunar impactors and the formation history of the Moon.

## Acknowledgments

This work was supported by the Natural Sciences and Engineering Research Council of Canada. Computations were performed on the P8 supercomputer at the SciNet HPC Consortium. SciNet is funded by: the Canada Foundation for Innovation under the auspices of Compute Canada; the Government of Ontario; Ontario Research Fund - Research Excellence; and the University of Toronto. The authors are grateful to François Chollet for building Keras. The authors thank R. Povilaitis and the LROC team for providing the 5-20 km dataset and useful discussions regarding the input data, R. Malhotra and R. Strom for providing helpful comments on a draft of the manuscript, and N. Hammond for useful discussions on craters distributions. Finally, the authors thank David Minton and an anonymous referee for their constructive comments during the review process that significantly improved the quality of this manuscript.

## References

Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mane, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viegas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., Zheng, X., Mar. 2016. TensorFlow: Large-

Scale Machine Learning on Heterogeneous Distributed Systems. ArXiv e-prints.

Bandeira, L., Ding, W., Stepinski, T. F., Jan. 2012. Detection of sub-kilometer craters in high resolution planetary images using shape and texture features. *Advances in Space Research* 49, 64–74.

Barker, M. K., Mazarico, E., Neumann, G. A., Zuber, M. T., Haruyama, J., Smith, D. E., Jul. 2016. A new lunar digital elevation model from the Lunar Orbiter Laser Altimeter and SELENE Terrain Camera. *Icarus* 273, 346–355.

Becker, K. J., Robinson, M. S., Becker, T. L., Weller, L. A., Edmundson, K. L., Neumann, G. A., Perry, M. E., Solomon, S. C., Mar. 2016. First Global Digital Elevation Model of Mercury. In: *Lunar and Planetary Science Conference*. Vol. 47 of Lunar and Planetary Inst. Technical Report. p. 2959.

Bottke, W. F., Norman, M. D., Aug. 2017. The Late Heavy Bombardment. *Annual Review of Earth and Planetary Sciences* 45, 619–647.

Boukercha, A., Al-Tameemi, A., Grumpe, A., Wöhler, C., Mar. 2014. Automatic Crater Recognition Using Machine Learning with Different Features and Their Combination. In: *Lunar and Planetary Science Conference*. Vol. 45 of Lunar and Planetary Inst. Technical Report. p. 2842.

- Canny, J., 1986. A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence* (6), 679–698.
- Chapman, C. R., Cohen, B. A., Grinspoon, D. H., Jul. 2007. What are the real constraints on the existence and magnitude of the late heavy bombardment? *Icarus* 189, 233–245.
- Chollet, F., 2015. Keras. <https://github.com/fchollet/keras>.
- Chung, C. P., Chung, O. C., Yam, C. H., 2014. Lunar Crater Identification from Machine Learning Perspective. *Acta Futura* 9, 41–45.  
URL <http://www.esa.int/gsp/ACT/publications/ActaFutura/issues.html>
- Cohen, J. P., Lo, H. Z., Lu, T., Ding, W., Mar. 2016. Crater Detection via Convolutional Neural Networks. In: *Lunar and Planetary Science Conference*. Vol. 47 of *Lunar and Planetary Inst. Technical Report*. p. 1143.
- Di, K., Li, W., Yue, Z., Sun, Y., Liu, Y., 2014. A machine learning approach to crater detection from topographic data. *Advances in Space Research* 54 (11), 2419 – 2429.  
URL <http://www.sciencedirect.com/science/article/pii/S0273117714005304>
- Duda, R. O., Hart, P. E., Jan. 1972. Use of the hough transformation to detect lines and curves in pictures. *Commun. ACM* 15 (1), 11–15.  
URL <http://doi.acm.org/10.1145/361237.361242>

- Emami, E., Bebis, G., Nefian, A., Fong, T., 2015. Automatic Crater Detection Using Convex Grouping and Convolutional Neural Networks. Springer, Cham, pp. 213–224.  
URL [https://doi.org/10.1007/978-3-319-27863-6\\_20](https://doi.org/10.1007/978-3-319-27863-6_20)
- Fassett, C. I., Head, J. W., Kadish, S. J., Mazarico, E., Neumann, G. A., Smith, D. E., Zuber, M. T., Feb. 2012. Lunar impact basins: Stratigraphy, sequence and ages from superposed impact crater populations measured from Lunar Orbiter Laser Altimeter (LOLA) data. *Journal of Geophysical Research (Planets)* 117, E00H06.
- Ferguson, R., Hare, T., Laura, J., 2017. HRSC and MOLA blended digital elevation model at 200m. [http://bit.ly/HRSC\\_MOLA\\_Blend\\_v0](http://bit.ly/HRSC_MOLA_Blend_v0).
- Gomes, R., Levison, H. F., Tsiganis, K., Morbidelli, A., May 2005. Origin of the cataclysmic Late Heavy Bombardment period of the terrestrial planets. *Nat.*435, 466–469.
- Goodfellow, I., Bengio, Y., Courville, A., 2016. *Deep Learning*. MIT Press, <http://www.deeplearningbook.org>.
- Greeley, R., Gault, D. E., Sep. 1970. Precision Size-Frequency Distributions of Craters for 12 Selected Areas of the Lunar Surface. *Moon* 2, 10.
- Grieve, R. A. F., May 1981. Impact cratering. *Nat.*291, 16.
- Hartmann, W. K., Sep. 1970. Note: Lunar Cratering Chronology. *Icarus*13, 299–301.

- Head, J. W., Fassett, C. I., Kadish, S. J., Smith, D. E., Zuber, M. T., Neumann, G. A., Mazarico, E., Sep. 2010. Global Distribution of Large Lunar Craters: Implications for Resurfacing and Impactor Populations. *Science* 329, 1504.
- Kingma, D. P., Ba, J., Dec. 2014. Adam: A Method for Stochastic Optimization. ArXiv e-prints.
- Kirchoff, M., Sherman, K., Chapman, C., Oct. 2011. Examining lunar impactor population evolution: Additional results from crater distributions on diverse terrains. In: EPSC-DPS Joint Meeting 2011. p. 1587.
- LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., Jackel, L. D., 1989. Backpropagation applied to handwritten zip code recognition. *Neural Computation* 1 (4), 541–551.  
URL <https://doi.org/10.1162/neco.1989.1.4.541>
- LOLA Team, Kaguya Team, 2015. LRO LOLA and Kaguya Terrain Camera DEM merge 60N60S 512ppd (59m). [https://astrogeology.usgs.gov/search/map/Moon/LRO/LOLA/Lunar\\_LRO\\_LrocKaguya\\_DEMmerge\\_60N60S\\_512ppd](https://astrogeology.usgs.gov/search/map/Moon/LRO/LOLA/Lunar_LRO_LrocKaguya_DEMmerge_60N60S_512ppd).
- Long, J., Shelhamer, E., Darrell, T., June 2015. Fully convolutional networks for semantic segmentation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

Magellan Science Team, 2014. Magellan global topography 4641m.  
[https://astrogeology.usgs.gov/search/map/Venus/Magellan/  
RadarProperties/Venus\\_Magellan\\_Topography\\_Global\\_4641m\\_v02.](https://astrogeology.usgs.gov/search/map/Venus/Magellan/RadarProperties/Venus_Magellan_Topography_Global_4641m_v02)

Minton, D. A., Richardson, J. E., Fassett, C. I., Feb. 2015. Re-examining the main asteroid belt as the primary source of ancient lunar craters. *Icarus*247, 172–190.

Palafox, L. F., Hamilton, C. W., Scheidt, S. P., Alvarez, A. M., 2017. Automated detection of geological landforms on mars using convolutional neural networks. *Computers & Geosciences* 101 (Supplement C), 48 – 56.

URL [http://www.sciencedirect.com/science/article/pii/  
S0098300416305532](http://www.sciencedirect.com/science/article/pii/S0098300416305532)

Pike, R. J., Jul. 1980. Formation of complex impact craters - Evidence from Mars and other planets. *Icarus*43, 1–19.

Povilaitis, R., Robinson, M., van der Bogert, C., Hiesinger, H., Meyer, H., Ostrach, L., 2017. Crater density differences: Exploring regional resurfacing, secondary crater populations, and crater saturation equilibrium on the moon. *Planetary and Space Science*.

URL [http://www.sciencedirect.com/science/article/pii/  
S0032063316303798](http://www.sciencedirect.com/science/article/pii/S0032063316303798)

Preusker, F., Scholten, F., Matz, K.-D., Elgner, S., Jaumann, R., Roatsch, T.,

- Joy, S. P., Polanskey, C. A., Raymond, C. A., Russell, C. T., Mar. 2016. Dawn at Ceres – Shape Model and Rotational State. In: Lunar and Planetary Science Conference. Vol. 47 of Lunar and Planetary Inst. Technical Report. p. 1954.
- Preusker, F., Scholten, F., Matz, K.-D., Roatsch, T., Jaumann, R., Raymond, C. A., Russell, C. T., Feb. 2014. Global Shape of (4) Vesta from Dawn FC Stereo Images. In: Vesta in the Light of Dawn: First Exploration of a Protoplanet in the Asteroid Belt. Vol. 1773 of LPI Contributions. p. 2027.
- Robbins, S. J., Antonenko, I., Kirchoff, M. R., Chapman, C. R., Fassett, C. I., Herrick, R. R., Singer, K., Zanetti, M., Lehan, C., Huang, D., Gay, P. L., May 2014. The variability of crater identification among expert and community crater analysts. *Icarus* 234, 109–131.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation. Springer International Publishing, Cham, pp. 234–241.  
URL [https://doi.org/10.1007/978-3-319-24574-4\\_28](https://doi.org/10.1007/978-3-319-24574-4_28)
- Rumelhart, D. E., Hinton, G. E., Williams, R. J., Oct. 1986. Learning representations by back-propagating errors. *Nat.* 323, 533–536.
- Ryder, G., Apr. 2002. Mass flux in the ancient Earth-Moon system and benign implications for the origin of life on Earth. *Journal of Geophysical Research (Planets)* 107, 6–1.



- Salamunićcar, G., Lončarić, S., 2010. Method for crater detection from digital topography data: interpolation based improvement and application to Lunar SELENE LALT data. In: 38th COSPAR Scientific Assembly. Vol. 38 of COSPAR Meeting. p. 3.
- Salamunićcar, G., Lončarić, S., Grumpe, A., Wöhler, C., Jun. 2014. Hybrid method for crater detection based on topography reconstruction from optical images and the new LU78287GT catalogue of Lunar impact craters. *Advances in Space Research* 53, 1783–1797.
- Scholten, F., Oberst, J., Matz, K.-D., Roatsch, T., Wählisch, M., Speyerer, E. J., Robinson, M. S., Mar. 2012. GLD100: The near-global lunar 100 m raster DTM from LROC WAC stereo image data. *Journal of Geophysical Research (Planets)* 117, E00H17.
- Shallue, C. J., Vanderburg, A., Feb. 2018. Identifying Exoplanets with Deep Learning: A Five-planet Resonant Chain around Kepler-80 and an Eighth Planet around Kepler-90. *Astron. J.*155, 94.
- Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., Salakhutdinov, R., 2014. Dropout: a simple way to prevent neural networks from overfitting. *Journal of machine learning research* 15 (1), 1929–1958.
- Stepinski, T., Ding, W., Vilalta, R., 2012. Detecting impact craters in planetary images using machine learning. *Information Science Reference, Hershey, PA*, pp. 146–159.

- Stoffler, D., Jan. 2006. Cratering History and Lunar Chronology. *Reviews in Mineralogy and Geochemistry* 60, 519–596.
- Strom, R. G., Malhotra, R., Ito, T., Yoshida, F., Kring, D. A., Sep. 2005. The Origin of Planetary Impactors in the Inner Solar System. *Science* 309, 1847–1850.
- UK Met. Office, 2015. Cartopy: a cartographic python library with a matplotlib interface. <http://scitools.org.uk/cartopy/index.html>.
- USGS Astrogeology Science Center, 2016. Mercury MESSENGER global DEM 665m (64ppd) v2 oct. 2016. [https://astrogeology.usgs.gov/search/map/Mercury/Topography/MESSENGER/Mercury\\_Messenger\\_USGS\\_DEM\\_Global\\_665m](https://astrogeology.usgs.gov/search/map/Mercury/Topography/MESSENGER/Mercury_Messenger_USGS_DEM_Global_665m).
- Van der Walt, S., Schönberger, J. L., Nunez-Iglesias, J., Boulogne, F., Warner, J. D., Yager, N., Gouillart, E., Yu, T., 2014. scikit-image: image processing in python. *PeerJ* 2, e453.
- Weaver, B. P., Robbins, S. J., Plesko, C. S., Riggs, J. D., May 2015. On the Binning and Associated Uncertainty of Crater Diameter Size-Frequency Distributions. In: *Issues in Crater Studies and the Dating of Planetary Surfaces*. Vol. 1841 of LPI Contributions. p. 9054.
- Wetzler, P., Honda, R., Enke, B., Merline, W., Chapman, C., Burl, M., 2005. Learning to detect small impact craters. In: *7th IEEE Workshop on Application of Computer Vision*. Vol. 1. pp. 178–184.

Yu, F., Koltun, V., Nov. 2015. Multi-Scale Context Aggregation by Dilated Convolutions. ArXiv e-prints.

Zeiler, M. D., Fergus, R., Nov. 2013. Visualizing and Understanding Convolutional Networks. ArXiv e-prints.

Zeng, L., Xu, X., Cai, B., Qiu, S., Zhang, T., 2017. Multi-scale convolutional neural networks for crowd counting. CoRR abs/1702.02359.

URL <http://arxiv.org/abs/1702.02359>

Zhang, Y., Zhou, D., Chen, S., Gao, S., Ma, Y., June 2016. Single-image crowd counting via multi-column convolutional neural network. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).