

---

# Convolutional neural networks and satellite imagery: How deep is necessary?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Applying off-the-shelf models (e.g., ResNet) to satellite imagery has become  
2 standard practice. While convolutional neural networks (CNNs) have been shown  
3 to outperform baseline methods in remote sensing prediction tasks, differences  
4 in satellite and natural images (i.e., images that comprise common datasets like  
5 ImageNet and CIFAR-10) may make ResNet-type models overkill for many satellite  
6 imagery tasks. In this paper, we present a comparison of off-the-shelf CNNs to  
7 a much smaller CNN over a range of satellite imagery tasks and show that a  
8 CNN with significantly fewer parameters performs on par with standard CNN  
9 architectures for five out of six tasks. Our findings are especially pertinent to those  
10 working with satellite imagery who face computational constraints.

## 11 1 Introduction

12 Machine learning has continually proven successful in informing sustainability-related tasks from  
13 satellite imagery. Example tasks include crop type mapping (1; 2), poverty prediction (3), and water  
14 quality monitoring (4; 5). It has become standard practice to apply off-the-shelf models (e.g., ResNet)  
15 to satellite imagery. While CNNs have been shown to outperform baseline methods in remote sensing  
16 tasks, there are substantial enough differences between natural images (i.e., images that comprise  
17 common datasets like ImageNet and CIFAR-10) and satellite images, that off-the-shelf CNNs may be  
18 unnecessarily large for many satellite imagery analyses.

19 Prior to deep learning, handcrafted features were common in satellite imagery analysis. Features  
20 generally consisted of low-level color and texture descriptors, such as color histograms (6). Recently,  
21 basic color descriptors have been shown to be highly effective in discriminating coffee/non-coffee  
22 scenes (6) and simple statistics over images have been shown to be informative in predicting bird  
23 distributions from satellite imagery (7). As filters in earlier layers of CNNs generally pick up on color  
24 and texture and the later layers are more representative of concepts (8), many satellite imagery tasks  
25 may not require significantly deep networks.

26 Parameter reduction in neural networks is a popular area of research (9; 10; 11; 12; 13). Not only are  
27 smaller models beneficial from a training perspective (e.g., time and carbon footprint), but in many  
28 cases smaller models are necessary when deployed on devices with limited computational resources.  
29 Methods for designing smaller neural networks generally involve compressing pretrained networks  
30 or designing smaller networks (10). MobileNet (10; 14) and SqueezeNet (9) are two common  
31 architectures optimized for fewer parameters but capable of achieving ResNet-level accuracy. While  
32 both of these networks are indeed smaller than ResNet18, they still contain roughly a million  
33 parameters. We hypothesize that common CNN architectures, even those optimized to have fewer  
34 parameters than ResNet, are overkill for many satellite imagery tasks.

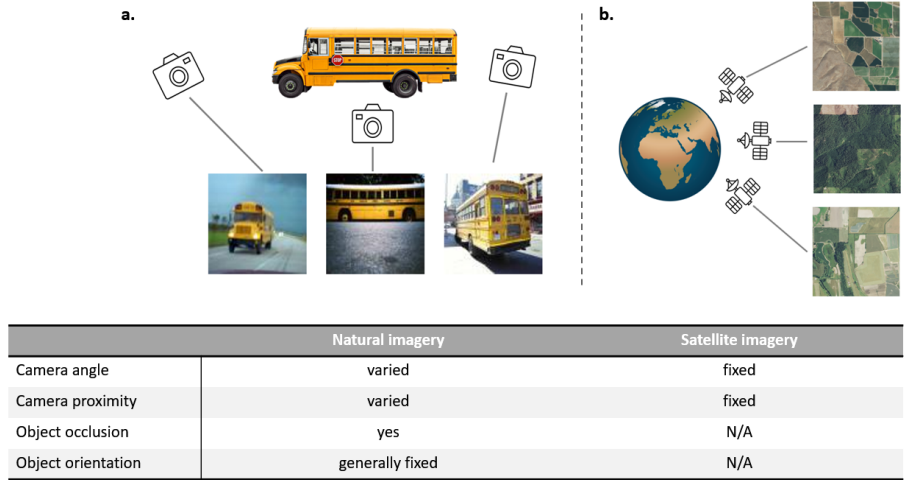


Figure 1: Natural versus satellite images. Left: three school buses captured at different viewing angles and proximities (15). There is no occlusion in the school bus images and the orientation of all school buses is the same (wheels on the ground, below the body of the buses). Right: three satellite images with the same camera angle, proximity, and orientation as images from the same satellite are captured at a fixed, bird’s-eye view.

35 We performed an analysis on CNN architectures to determine if larger models, such as ResNet and  
 36 ResNet alternatives, are necessary for satellite imagery tasks. Specifically, we compared several  
 37 off-the-shelf models to a simple, shallow CNN on multiple regression and classification tasks. Our  
 38 main contributions are: 1) an analysis of modern CNN architectures across several satellite imagery  
 39 tasks, and 2) results showing that a shallow CNN, with millions of fewer parameters than ResNet18,  
 40 is comparable to standard ResNet-type models for five out of six satellite imagery tasks.

## 41 2 Natural versus satellite images

42 Satellite images differ from natural images in several key ways. Most notably, satellite images are  
 43 captured in a significantly more structured manner compared to natural images. There are several  
 44 degrees of freedom when capturing natural images: 1) camera proximity from the subject, 2) camera  
 45 angle and 3) camera orientation relative to the subject, and 4) items occluding the subject (Figure 1a).  
 46 Satellite images within publicly available datasets, such as Landsat and Sentinel, are captured at a  
 47 relatively fixed distance and orientation to Earth and rarely suffer from occlusion (Figure 1b). This  
 48 fixed nature in which satellite images are captured enforces that images taken from the same satellite  
 49 have essentially the same scale. It is possible for atmospheric conditions (e.g., clouds) to obscure  
 50 satellite images; however, a common preprocessing step is to filter out days/images that contain heavy  
 51 atmospheric conditions. The fewer degrees of freedom in which satellite images are captured likely  
 52 simplifies the complexity of many prediction tasks.

53 CNNs for natural images must have incredibly large capacities (i.e., millions of parameters) to  
 54 represent classes well. Accounting for different camera viewing angles and distances, and the  
 55 potential of object occlusion makes object classification a difficult task. CNNs must learn many  
 56 representations of the same class (e.g., a school bus is a school bus no matter what angle or position  
 57 the image is captured from and whether or not the bus is partially occluded - Figure 1a). Apart from  
 58 the differences in how natural and satellite images are captured, there is evidence that color and  
 59 texture descriptors are effective features in satellite imagery analysis (6). Such features are detected  
 60 in early CNN layers, meaning relatively deep CNNs are potentially unnecessary. Below, we test if  
 61 these differences between domains does indeed simplify satellite imagery analysis and determine if  
 62 shallower CNN architectures are capable of performing as well as larger, off-the-shelf models.

Table 1: **Summary of architectures.** The number of parameters and parameter reductions compared to ResNet18 and memory required to train each architecture.

Architecture	# parameters	Reduction from ResNet18	Memory
ResNet18	11.2M	-	42.80MB
MobileNetV3	1.7M	7x	6.34MB
SqueezeNet	0.7M	16x	2.76 MB
ShallowCNN	0.1M	87x	0.49MB

### 63 3 Experiments

64 We compared four CNN architectures on three regression and three classification tasks. We compared  
 65 a basic, shallow CNN, hereafter *ShallowCNN*, to three pretrained off-the-shelf CNNs: ResNet18,  
 66 MobileNetV3-small, and SqueezeNet. Our ShallowCNN has a straightforward (conv2d-batch norm-  
 67 relu-pooling) architecture (Table A1). For each of the tasks and models, with the exception of  
 68 Brazilian Coffee Scenes, we used 10-fold cross-validation to train 10 models in order to quantify  
 69 variation in model performance. We used 5-fold cross-validation on Brazilian Coffee Scenes as the  
 70 dataset came with 5 predefined folds. We fine-tuned all layers of the pretrained models and trained  
 71 ShallowCNN from scratch, for each task respectively. We evaluated the regression tasks with  $R^2$ ,  
 72 mean squared error (MSE), and mean absolute error (MAE) and evaluated the classification tasks  
 73 with overall accuracy, precision, and recall. All three classification tasks have roughly equal class  
 74 balance, therefore, we report the average overall accuracy, precision, and recall across classes. In  
 75 addition to comparing the CNN architectures across tasks, we evaluated the impact of dataset size on  
 76 model performance. Details on data preprocessing and model training are in the appendix.

77 We compared the four CNN architectures by modeling six satellite imagery tasks. In order to  
 78 generalize across tasks, we selected a set of problems that range in difficulty and dataset size. We  
 79 modeled three regression tasks from Rolf et al. (2021): percent forest cover, nighttime light intensity,  
 80 and elevation (Table A2, Figure A1) Additionally, we modeled three classification tasks: crop type,  
 81 Brazilian Coffee Scenes (6), and UCMerced Land-use (17) (Table A2, Figure A1). We selected  
 82 the Brazilian Coffee Scenes and UCMerced Land-use datasets as they have been commonly used  
 83 in previous remote sensing studies (18; 19; 20) and created our own crop type dataset. A detailed  
 84 description of the tasks can be found in the appendix.

### 85 4 Results and Discussion

86 ShallowCNN, with 87 times fewer parameters than ResNet18 (Table 1), was within two standard  
 87 deviations of ResNet18 in three of the six tasks and exceeded ResNet18 by over four standard  
 88 deviations in crop type mapping (Table 2). In addition to crop type mapping, ShallowCNN also  
 89 achieved the highest accuracy for coffee scene identification (Table 2b). The only case in which  
 90 ShallowCNN had significantly degraded performance is in the land use problem (Table 2b). In  
 91 comparing the images from all tasks (Figure A1), the land use images are visually the most similar  
 92 to ImageNet, while the other tasks are likely more reliant on color and texture descriptors. The  
 93 land use task may require higher-level features to distinguish classes. As the deeper layers of  
 94 CNNs are generally more representative of concepts (8), satellite imagery tasks which are visually  
 95 similar to ImageNet may benefit from larger CNNs. Further work is needed to understand which  
 96 satelliteimagery tasks are better suited for ShallowCNN versus ResNet18-type models.

97 There are computational benefits to ShallowCNN. While ShallowCNN takes almost as long to train  
 98 as it does to fine-tune the larger pretrained models, ShallowCNN requires significantly less memory  
 99 to train (Table 1). If computational resources are limited, training multiple larger models at the same  
 100 time may not be feasible, however, it may be possible to train several ShallowCNNs concurrently. In  
 101 addition to benefits in training, ShallowCNN’s smaller model size is beneficial if models are being  
 102 deployed externally and space is limited by hardware (e.g., analyzing images real time on a drone  
 103 or satellite). Further work should investigate the potential of a pretrained ShallowCNN for satellite  
 104 imagery.

Table 2: **Summary of model performances across tasks.** Regression tasks are evaluated on MSE and classification tasks on accuracy (other metrics showed similar trends). The reported accuracy is averaged across all classes since all classification tasks have good class balance. Reported results are averaged across cross-validation folds (10 folds for all but Coffee, which had 5), plus or minus their standard deviation. Models are ordered from greatest number of parameters (top) to fewest (bottom).

**a. Regression tasks - MSE ( $\times 10^{-3}$ )**

Model	Percent forest cover	Nighttime light intensity	Elevation
ResNet18	<b><math>6.06 \pm 0.31</math></b>	$3.40 \pm 0.14$	$7.55 \pm 2.52$
MobileNetV3	$7.20 \pm 1.12$	<b><math>3.22 \pm 0.26</math></b>	<b><math>7.00 \pm 2.32</math></b>
SqueezeNet	$7.92 \pm 1.36$	$3.54 \pm 0.23$	$8.40 \pm 2.63$
ShallowCNN	$6.85 \pm 2.33$	$3.60 \pm 0.26$	$10.80 \pm 2.65$

**b. Classification tasks - overall accuracy (%)**

Model	Crop type	Land use	Coffee
ResNet18	$93.35 \pm 0.38$	<b><math>98.89 \pm 0.47</math></b>	$91.28 \pm 2.19$
MobileNetV3	$93.75 \pm 0.45$	$96.53 \pm 2.86$	$88.33 \pm 1.48$
SqueezeNet	$92.80 \pm 0.34$	$91.00 \pm 1.90$	$90.23 \pm 0.93$
ShallowCNN	<b><math>95.04 \pm 0.33</math></b>	$89.00 \pm 1.87$	<b><math>92.13 \pm 1.02</math></b>

105 A significant benefit to ShallowCNN is its performance  
 106 on small datasets, an issue common among  
 107 sustainability-related tasks. Although a common  
 108 principle in machine learning says that smaller models  
 109 should be favored for simpler prediction problems,  
 110 it is still common practice to apply ResNet-type  
 111 models to satellite imagery tasks. Until the training  
 112 size is reduced to 8k images, there is a small,  
 113 roughly linear increase in MSE across all models  
 114 (Figure 2). Once the training set size falls below  
 115 8k, there is an exponential increase in MSE for  
 116 ResNet and SqueezeNet and a much smaller increase  
 117 in MSE for ShallowCNN and MobileNet (Figure 2).  
 118 While MobileNet outperforms ShallowCNN on smaller  
 119 datasets, ShallowCNN’s smaller memory requirements  
 120 may outweigh the difference in performance for some  
 121 applications. Further work should investigate why  
 122 MobileNetV3, which has significantly more parameters  
 123 than ShallowCNN and SqueezeNet, outperforms both  
 124 methods in the small data regime.



Figure 2: MSE for nighttime light intensity as dataset size is reduced. Note that the x-axis is decreasing.

## 125 5 Conclusion

126 Evaluating the performance of smaller, non-standard deep architectures is generally underexplored  
 127 (21). This is especially of interest when applying CNNs to new domains where there are fundamental  
 128 differences between the domain images and images that comprise common computer vision datasets  
 129 (i.e., natural images). We compared a small three layer CNN, ShallowCNN, to much larger off-  
 130 the-shelf architectures and showed that in most cases ShallowCNN has comparable performance  
 131 to ResNet-type models. Our results align with other studies which have shown that simple CNN  
 132 architectures perform well on satellite imagery tasks (22; 23). Our study differs from the previous  
 133 studies in that it is the first to baseline the performance of modern pretrained architectures with  
 134 a smaller CNN across several remote sensing tasks, whereas the previous studies investigated the  
 135 viability of CNNs compared to other deep learning methods. Our results are of special concern to  
 136 those who have computational constraints, whether in model training or in model deployment.

## References

- [1] Sherrie Wang, Stefania Di Tommaso, Joey Faulkner, Thomas Friedel, Alexander Kennepohl, Rob Strey, and David B Lobell. Mapping crop types in southeast india with smartphone crowdsourcing and deep learning. *Remote Sensing*, 12(18):2957, 2020.
- [2] Gabriel Tseng, Hannah Kerner, Catherine Nakalembe, and Inbal Becker-Reshef. Learning to predict crop type from heterogeneous sparse labels using meta-learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 1111–1120, June 2021.
- [3] Neal Jean, Marshall Burke, Michael Xie, W. Matthew Davis, David B. Lobell, and Stefano Ermon. Combining satellite imagery and machine learning to predict poverty. *Science*, 353(6301):790–794, 2016. doi: 10.1126/science.aaf7894. URL <https://www.science.org/doi/abs/10.1126/science.aaf7894>.
- [4] Nima Pahlevan, Brandon Smith, John Schalles, Caren Binding, Zhigang Cao, Ronghua Ma, Krista Alikas, Kersti Kangro, Daniela Gurlin, Nguyn Hà, et al. Seamless retrievals of chlorophyll-a from sentinel-2 (msi) and sentinel-3 (olci) in inland and coastal waters: A machine-learning approach. *Remote Sensing of Environment*, 240:111604, 2020.
- [5] Jeremy Kravitz, Mark Matthews, Lisl Lain, Sarah Fawcett, and Stewart Bernard. Potential for high fidelity global mapping of common inland water quality products at high spatial and temporal resolutions based on a synthetic data and machine learning approach. *Frontiers in Environmental Science*, 9, 2021. ISSN 2296-665X. doi: 10.3389/fenvs.2021.587660. URL <https://www.frontiersin.org/articles/10.3389/fenvs.2021.587660>.
- [6] Otávio AB Penatti, Keiller Nogueira, and Jefersson A Dos Santos. Do deep features generalize from everyday objects to remote sensing and aerial scenes domains? In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 44–51, 2015.
- [7] Laurel M Hopkins, Tyler A Hallman, John Kilbride, W Douglas Robinson, and Rebecca A Hutchinson. A comparison of remotely sensed environmental predictors for avian distributions. *Landscape Ecology*, 37(4):997–1016, 2022.
- [8] Bolei Zhou, David Bau, Aude Oliva, and Antonio Torralba. Interpreting deep visual representations via network dissection. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2131–2145, 2018.
- [9] Forrest N. Iandola, Song Han, Matthew W. Moskewicz, Khalid Ashraf, William J. Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 1/10 model size, 2016. URL <https://arxiv.org/abs/1602.07360>.
- [10] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- [11] Gao Huang, Shichen Liu, Laurens Van der Maaten, and Kilian Q Weinberger. Condensenet: An efficient densenet using learned group convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2752–2761, 2018.
- [12] Bichen Wu, Alvin Wan, Xiangyu Yue, Peter Jin, Sicheng Zhao, Noah Golmant, Amir Ghohaminejad, Joseph Gonzalez, and Kurt Keutzer. Shift: A zero flop, zero parameter alternative to spatial convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9127–9135, 2018.
- [13] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6848–6856, 2018.
- [14] A. Howard, M. Sandler, B. Chen, W. Wang, L. Chen, M. Tan, G. Chu, V. Vasudevan, Y. Zhu, R. Pang, H. Adam, and Q. Le. Searching for mobilenetv3. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1314–1324, Los Alamitos, CA, USA, nov 2019. IEEE Computer Society. doi: 10.1109/ICCV.2019.00140. URL <https://doi.ieeecomputersociety.org/10.1109/ICCV.2019.00140>.

- 188 [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-  
189 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*  
190 *recognition*, pages 248–255. IEEE, 2009.
- 191 [16] E. Rolf, J. Proctor, T. Carleton, I. Bolliger, V. Shankar, M. Ishihara, B. Recht, and S. Hsiang. A  
192 generalizable and accessible approach to machine learning with global satellite imagery. *Nature*  
193 *Commun*, 12(4392), 2021. doi: 10.1038/s41467-021-24638-z.
- 194 [17] Yi Yang and Shawn Newsam. Bag-of-visual-words and spatial extensions for land-use clas-  
195 sification. In *Proceedings of the 18th SIGSPATIAL International Conference on Advances in*  
196 *Geographic Information Systems, GIS '10*, page 270–279, New York, NY, USA, 2010. Associa-  
197 tion for Computing Machinery. ISBN 9781450304283. doi: 10.1145/1869790.1869829. URL  
198 <https://doi.org/10.1145/1869790.1869829>.
- 199 [18] Marco Castelluccio, Giovanni Poggi, Carlo Sansone, and Luisa Verdoliva. Land use clas-  
200 sification in remote sensing images by convolutional neural networks, 2015. URL <https://arxiv.org/abs/1508.00092>.  
201
- 202 [19] Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification:  
203 Benchmark and state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017.
- 204 [20] Keiller Nogueira, Otávio AB Penatti, and Jefersson A Dos Santos. Towards better exploiting  
205 convolutional neural networks for remote sensing scene classification. *Pattern Recognition*, 61:  
206 539–556, 2017.
- 207 [21] Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding  
208 transfer learning for medical imaging. *Advances in neural information processing systems*, 32,  
209 2019.
- 210 [22] Saikat Basu, Sangram Ganguly, Supratik Mukhopadhyay, Robert DiBiano, Manohar Karki,  
211 and Ramakrishna R. Nemani. Deepsat - A learning framework for satellite imagery. *CoRR*,  
212 <abs/1509.03602>, 2015. URL <http://arxiv.org/abs/1509.03602>.
- 213 [23] Yanfei Zhong, Feng Fei, Yanfei Liu, Bei Zhao, Hongzan Jiao, and Liangpei Zhang. Satcnn:  
214 satellite image dataset classification using agile convolutional neural networks. *Remote Sensing*  
215 *Letters*, 8(2):136–145, 2017. doi: 10.1080/2150704X.2016.1235299. URL [https://doi.](https://doi.org/10.1080/2150704X.2016.1235299)  
216 [org/10.1080/2150704X.2016.1235299](https://doi.org/10.1080/2150704X.2016.1235299).
- 217 [24] National Geospatial Data Asset (NGDA) NAIP Imagery. [http://gis.apfo.usda.gov/](http://gis.apfo.usda.gov/arcgis/rest/services/NAIP)  
218 [arcgis/rest/services/NAIP](http://gis.apfo.usda.gov/arcgis/rest/services/NAIP), March 2015.
- 219 [25] Collin Homer, Jon Dewitz, Limin Yang, Suming Jin, Patrick Danielson, George Xian, John  
220 Coulston, Nathaniel Herold, James Wickham, and Kevin Megown. Completion of the 2011  
221 national land cover database for the conterminous united states-representing a decade of land  
222 cover change information. *Photogrammetric Engineering and Remote Sensing*, 81(5):345–354,  
223 May 2015.
- 224 [26] Roozbeh Valavi, Jane Elith, José J. Lahoz-Monfort, and Gurutzeta Guillera-Arroita. blockcv: An  
225 r package for generating spatially or environmentally separated folds for k-fold cross-validation  
226 of species distribution models. *Methods in Ecology and Evolution*, 10(2):225–232, 2019. doi:  
227 [10.1111/2041-210X.13107](https://doi.org/10.1111/2041-210X.13107).

228 **Appendix**

229 **A1 ShallowCNN**

230 Our shallowCNN is comprised of three convolutional layers. The weights in conv1 are taken from  
231 the first layer of ResNet18 pretrained on ImageNet as using pretrained weights in the initial layer can  
232 lead to faster convergence (21). To match ResNet18, the first convolutional layer has 7x7 filters. The  
remaining two layers have 3x3 filters. See Table A1 for the full architecture.

Table A1: **ShallowCNN architecture.** For the last layer we used a multiclass softmax activation for the classification tasks and a linear activation for the regression tasks.

Layer name	
conv1	7x7, 64 conv; batch norm; ReLU; max pooling
conv2	3x3, 64 conv; batch norm; ReLU; max pooling
conv3	3x3, 128 conv; batch norm; ReLU; max pooling
avg_pool	average pooling
fc	128x64, ReLU, 64xnum_classes
activation	classification: softmax; regression: linear

233

234 **A2 Tasks**

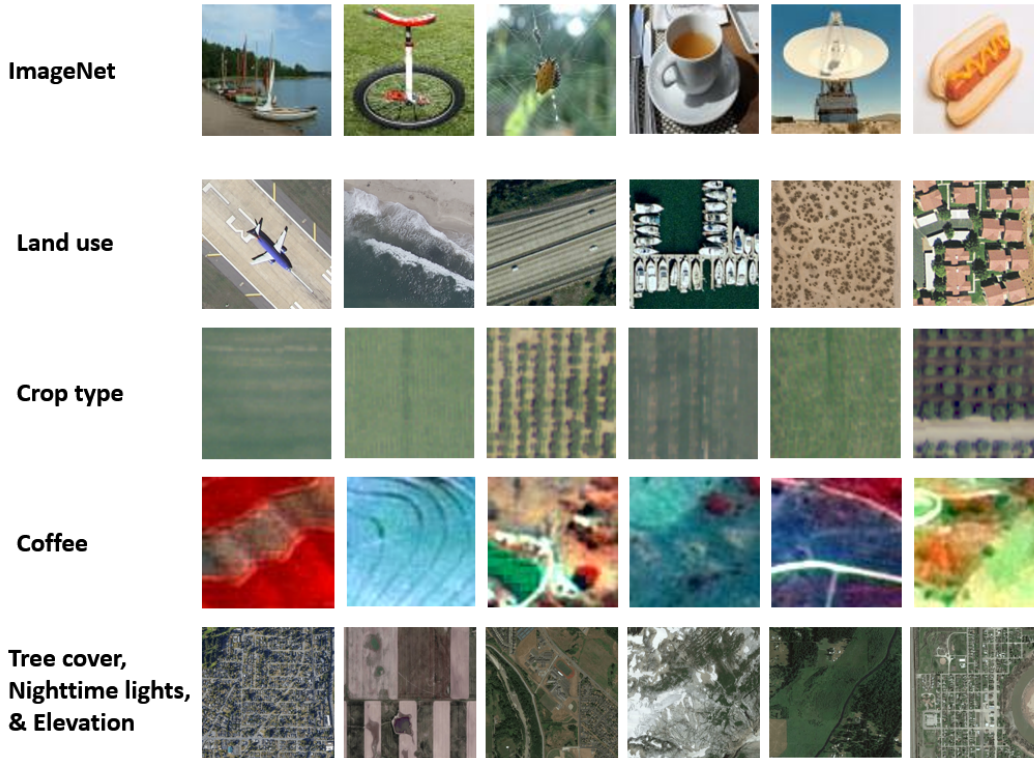


Figure A1: Sample images from ImageNet (15), UC Merced Land-use (17), crop type, Brazilian Coffee Scenes (6) and percent forest/nighttime light intensity/elevation tasks (16). Color and textural information appear more indicative of class in the satellite images as opposed to those of ImageNet.

### 235 **A2.1 Percent forest cover, nighttime light intensity, and elevation**

236 We purposefully selected a set of satellite imagery tasks that range in difficulty and dataset size.  
237 Table A2 outlines the tasks and dataset sizes. We selected three regression tasks from Rolf et al.  
238 (2021) with increasing complexity: percent forest cover, nighttime light intensity, and elevation. Forest  
239 cover is directly observable from satellite imagery and, therefore, should be the most straightforward  
240 to predict. Nighttime light intensity itself is not observable from daytime satellite images, however,  
241 proxies for nighttime light intensity (e.g., dense urban areas) are observable. Elevation on the other  
242 hand, is much more difficult to estimate solely from a satellite image. Many images may have similar  
243 appearances but dramatically different elevations. Images for all three regression tasks were collected  
244 from the contiguous United States based on the sampling schemes of Rolf et. al (2021). See Figure  
245 A1 for sample images.

### 246 **A2.2 Crop type**

247 We developed our own crop type dataset from images collected from three regions within the Central  
248 Valley of California. We collected National Agricultural Imagery Program (NAIP) imagery (24) from  
249 2012 and derived labels from the National Land Cover Database (NLCD) (25) for the same year.  
250 We subset the data to only include the three most commonly occurring crops: tomatoes, almonds,  
251 and alfalfa. When assigning labels, we only included images which contained more than 60% of the  
252 majority label in the image. In total the dataset consists of 36,000 images with a roughly even split  
253 across the three classes.

### 254 **A2.3 Brazilian Coffee Scenes**

255 SPOT satellite images were collected in 2005 over four counties in Brazil. Images were labeled by  
256 agricultural experts and labeled coffee if more than 85% of the pixels contained coffee and non-coffee  
257 if less than 10% of the pixels contained coffee. The dataset consists of 2876 images with an equal  
258 split of coffee and non-coffee (6).

### 259 **A2.4 UC Merced Land-use**

260 The UC Merced Land-use dataset consists of aerial images from 21 different land use classes. The  
261 classes span categories such as beach, parking lot, buildings, forest, and overpass. The images were  
262 collected from 20 cities across the United States. The images were manually annotated and each  
263 class contains 100 images (17).

Table A2: **Description of tasks.**

Prediction task	Type	Classes	Dataset size	Image size	Spatial res.
Percent forest cover (16)	regression	1	100k	256x256	~ 4 m
Nighttime light intensity (16)	regression	1	100k	256x256	~ 4 m
Elevation (16)	regression	1	100k	256x256	~ 4 m
Crop type	classification	3	36k	48x48	1 m
Brazilian Coffee Scenes (6)	classification	2	2876	64x64	-
UCMerced Land-use (17)	classification	21	2100	256x256	0.3 m

## 264 **A3 Data pre-processing**

265 Spatial autocorrelation is a common issue in spatial data and can be problematic as it can artificially  
266 overestimate the predictive power of models by having highly correlated datapoints (i.e., datapoints  
267 close in geographical space) in both the training and testing sets. To help address spatial autocorrela-  
268 tion, we used the blockCV R package (26) to split datasets with geographic location information into  
269 spatial blocks. We then used the spatial blocks to assign data points into 10 folds for cross validation.  
270 For the tasks without location information, we randomly split the data into 10 cross validation folds.  
271 All three regression tasks and the crop type task have location information (i.e., they were split



272 spatially). We randomly assigned splits for the land use task and used the predefined splits for the  
273 coffee dataset.

274 For image preprocessing we scaled pixel values to be in the range of  $[0, 1]$  and subtracted the channel  
275 means. During training, we augmented images by performing random horizontal and vertical flips  
276 and random rotations in increments of  $90^\circ$ .

## 277 **A4 Training**

278 We fine-tuned the three off-the-shelf models (ResNet18, MobileNetV3-small, and SqueezeNet) and  
279 trained ShallowCNN from scratch. For the three off-the-shelf architectures, we used pretrained  
280 weights and updated the fully-connected layer to match the number of outputs for the given task. For  
281 all models and all tasks, we performed hyperparameter tuning on the learning rate, weight decay, and  
282 batch size. Prior to comparing ShallowCNN to the other models, we experimented with different  
283 numbers of convolutional layers and different numbers of convolutions per layer. Across tasks, we  
284 found a three layer model with 64, 64, and 128 convolutions to perform the best. We transferred  
285 pretrained weights from the first convolutional layer of ResNet18 for the first layer of ShallowCNN as  
286 using pretrained weights in the initial layer can lead to faster convergence (21). To match ResNet18,  
287 the first convolutional layer has  $7 \times 7$  filters. The remaining two layers have  $3 \times 3$  filters and the weights  
288 were randomly initialized. See Table A1 for the full architecture.