CLOUD CLASSIFICATION WITH UNSUPERVISED DEEP LEARNING

Takuya Kurihana¹, Ian Foster^{1,2}, Rebecca Willett^{1,4}, Sydney Jenkins^{1,5}, Kathryn Koenig^{1,6}, Ruby Werman⁷, Ricardo Barros Lourenco¹, Casper Neo¹, Elisabeth Moyer³

Abstract—We present a framework for cloud characterization that leverages modern unsupervised deep learning technologies. While previous neural networkbased cloud classification models have used supervised learning methods, unsupervised learning allows us to avoid restricting the model to artificial categories based on historical cloud classification schemes and enables the discovery of novel, more detailed classifications. Our framework learns cloud features directly from radiance data produced by NASA's Moderate Resolution Imaging Spectroradiometer (MODIS) satellite instrument, deriving cloud characteristics from millions of images without relying on pre-defined cloud types during the training process. We present preliminary results showing that our method extracts physically relevant information from radiance data and produces meaningful cloud classes.

I. MOTIVATION

Clouds play a dominant role in the Earth's radiation budget, both reflecting sunlight and trapping infrared radiation. Their responses are the principal source of uncertainty in numerical simulations of future climate, because even state-of-the-art climate models cannot accurately resolve cloud formation and evolution on scales from sub-kilometers to thousands of kilometers [1]. NASA satellite instruments have observed cloud behavior for several decades, providing us with a rich dataset that can potentially inform understanding of cloud dynamics and feedbacks, but these large datasets have not yet been fully employed, in part because computing power has only recently approached the necessary scale. Clouds are therefore a timely target for large scale computational analyses that can automate

Corresponding author: I Foster, foster@uchicago.edu

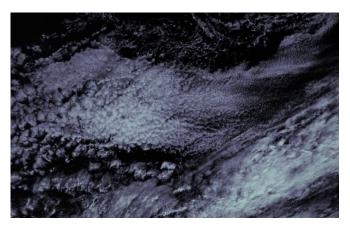


Fig. 1. MODIS Terra satellite visible imagery on December 1, 2015 over the East Pacific ocean, off the coast of California, capturing stratocumulus clouds at a variety of scales. Clouds at lower right are higher in altitude. We use this scene in all examples that follow.

detection of cloud attributes and identify scientifically relevant cloud classes.

Cloud classification effectively reduces the dimensionality of information in satellite images, rendering them tractable to analysis. Attempts to use neural network methods for this purpose date to the 1990s, when Lee et al. [2] used human-labeled images to train shallow, fully connected networks to recognize stratocumulus, cumulus, and cirrus clouds. Similar efforts have continued sporadically up to the present (e.g., [3-5]), all involving supervised learning based on images labeled with historically defined classes. However, none have produced an operational tool for automated analysis of cloud images. Supervised learning generally falls short because historical classes are artificial and are welldefined only for "classic" examples that make up a small fraction of total data. Furthermore, those classes do not distinguish scales of features that in nature may vary by an order of magnitude or more. For example, in the MODIS example image of Figure 1, stratocumulus clouds show a wide range of cell sizes, and cloud textures and patterns vary in complex ways.

Unsupervised learning methods may be a more ap-

¹Department of Computer Science, University of Chicago

²Data Science and Learning Division, Argonne National Lab

⁴Department of Statistics, University of Chicago

⁵Department of Physics, University of Chicago

⁶Harris School of Public Policy, University of Chicago

⁷College of Letters & Science, University of California, Berkeley

³Department of the Geophysical Sciences, University of Chicago

propriate means of making use of the complex information in large multi-spectral satellite datasets. Such methods allow novel data-driven cloud types to emerge from imagery data, and in principle can track changes in cloud textures and patterns over time, by identifying both changing frequencies of individual cloud classes and evolving characteristics within a class.

A scientifically useful operational classification system would:

- produce physically reasonable classes with scientifically relevant distinctions
- 2) capture information on cloud *spatial distributions*, i.e., be not reproducible using only mean properties over the target area
- produce classes that in high-dimensional space are cohesive within each class but separated between classes
- 4) be *rotationally invariant*, i.e., insensitive to the orientation of an image
- 5) be *stable*, i.e., produce similar or identical classes when different subsets of the data are used.

We describe here the construction of a prototype data-driven workflow for cloud classification based on unsupervised deep learning. In the following, we introduce our model architecture and clustering procedure, apply it to images from the MODIS satellite instrument, and evaluate its ability to meet some of the key criteria listed above.

II. METHOD

A. Model Architecture

We leverage recent work in self-supervised learning, in which an encoder-decoder network is trained to recover an input image. We use a deep convolutional autoencoder [6] to obtain dimensionally reduced information from input data. Autoencoders have been widely used to retrieve dimensionally reduced information from high-dimensional input data. The resulting lower-dimensional latent representations incorporate important input features, simplifying the classification task. In the general framework of an autoencoder, the learning process minimizes the loss function L:

$$\min L(x) = \min ||x - F(x)||_p$$
, (1)

where x is the input image; $F(\cdot)$ is a function which maps the input image on the dimension-reduced representation and then reconstructs the image from the intermediate information, meaning F(x) is the reconstructed input image; and $||\cdot||_p$ denotes the p-norm of the two images.

Our loss function combines four metrics: L1 and L2 loss, corresponding to p=1 and p=2 in Equation 1; the high frequency error norm after passing through the Sobel filter to detect edges of input clouds; and the multi-scale structure similarity index (MSSIM) [7], a multi-band version of SSIM [8], an index often used in computer vision to assess image similarity. We use the Adam optimizer [9], a combination of RMSprop and stochastic gradient descent with momentum, to optimize our loss function, with a learning rate of 10^{-4} .

We also include a convolution layer in our model in order to preserve spatial structure of the input image. The convolution operation implements a small-size filter to subset the entire image iteratively, with specified stride and width. The filter kernel operation extracts local features and parses the activation layer. The convolutional layer with activation function is described as

$$h^{l} = f\left(\sum_{i}\sum_{j}x_{(i+w-1)(j+s-1)}^{l}\otimes W_{ij}^{l} + b^{l}\right)$$
 (2)

where h^l is the lth layer's latent representation; f is a nonlinear activation function; $x_{(i+w)(j+s)}$ is a $w \times s$ domain for the convolutional filter; W_{ij} is the weight at the ith column and jth row; \otimes is the convolutional operation; and b denotes the bias. We set the filter size to 3×3 and use Leaky Rectified Linear Unit (Leaky ReLU) as the activation function $f(x) = \max(0.3x, x)$, as that performs better than common ReLU. Additionally, we build a residual connection every two convolutional layers to improve network performance, and add batch normalization after each residual connection. Between residual blocks, the size of an input image is scaled by a factor of two. In the encoder, the width and the stride are halved at each block, while the depth is doubled. In the decoder, these transformations are reversed, with the minor modification that we apply a transposed convolution kernel to map each input pixel to 3×3 pixels for up-sampling. The overall model architecture is illustrated in Figure 2.

We implement the convolutional autoencoder in the TensorFlow deep learning library [10] and use the Horovod framework [11] for data parallelization. Our encoder-decoder architecture stacks 20 convolutional layers and has 297 232 trainable parameters, and our latent representation has size $8\times8\times128$. Training took 100 000 steps and 17 hours to converge the loss function on four NVIDIA K80 GPUs on the University of Chicago's Midway compute cluster. We chose a batch size of 32 in accordance with common deep

CLOUD CLASSIFICATION WITH UNSUPERVISED DEEP LEARNING

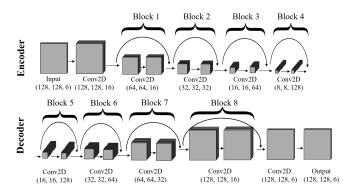


Fig. 2. Our autoencoder architecture. The encoder consists of four blocks, each with two convolutional layers activated by Leaky ReLU with residual connections; the decoder has the mirror structure of the encoder. The arrows represent the flow of input images, and each bracketed triple gives the height, width, and channel dimensions of the layer(s) above.

learning parameter settings and processed 789 GB of training images. We describe the input data and its preprocessing below.

B. Dataset

We test our workflow to perform cloud feature extraction on multi-spectral data from the MODIS instrument. Our input data is MODIS level 1B calibrated radiance imagery at 1 km resolution (MOD021KM; hereafter MOD02). This product has 36 spectral bands, from the visible to the thermal infrared; we work with bands 6, 7, 20, 28, 29, and 31, as these are the most important for the MODIS06 Level 2 algorithms that characterize cloud properties [12, 13]. Bands 6, 7, and 20 (1.6, 2.1, and 3.7 µm) are encoded in the algorithm to estimate cloud optical properties (e.g., optical thickness and effective radius), and the brightness temperatures at bands 28, 29, and 31 (7.3, 8.5, and 11 µm) are used in the separation of high and low clouds and the detection of cloud phase. Note that because we seek to discover aspects of these physics variables in our classification, we do not use derived properties such as brightness temperature, but instead input radiance data directly.

To enable efficient learning of cloud features, we define the unit of our input data, a "patch," as a small subset of a typical MODIS image: $128 \text{ km} \times 128 \text{ km} \times 6$ selected bands, out of an image of $2030 \text{ km} \times 1354 \text{ km} \times 36$ bands. To select input patches that contain clouds, we align the MOD02 data to its corresponding MODIS35 Level 2 cloud flag product ("MOD35"), and define a patch to be valid if more than 30% of the patch is comprised of cloud pixels as detected by MOD35. We then train the network using $\sim 1.01 \text{ million patches}$:

about 1% of the full 19-year dataset from a single MODIS satellite instrument.

C. Clustering

We use hierarchical agglomerative clustering (HAC) to merge data points by minimizing cluster variance, thus building a tree structure during the merging process. We choose HAC because it exhibits greater stability with respect to initialization than does k-means clustering. Our linkage metric is Ward's method, formulated as following

$$\delta \text{dist}(X_A, X_B) = \frac{n_A n_B}{n_A + n_B} ||C_A - C_B||^2 , \qquad (3)$$

where the distance between two clusters X_A and X_B is evaluated as the squared distance between the centroids of merged clusters C_A and C_B weighted by the number of patches in these clusters n_A and n_B .

To choose the number of clusters for an analysis, we would ideally determine the number for which clustering results are stable (allowing the permutation of clustering categories). As an approximate measure of stability, we measure the similarity of clusters using the Adjusted Mutual Information score (AMI). We first obtain pseudo ground-truth labels by applying clustering to ~320,000 patches, and then conduct tests with varying subsets of patches (chosen at random from the full dataset) and varying numbers of clusters, and compute the AMI score between the ground truth labels and subsets. The AMI score typically stabilizes at cluster numbers of about 10 and higher. Most demonstration analyses shown below use 12 clusters; a larger number is likely desirable for eventual science use.

Fig. 3 shows the complete pipeline from input data to resulting clusters.

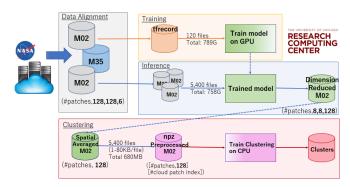


Fig. 3. Cartoon of the framework used in this work. M02 and M35 are the MOD02 and MOD35 satellite data products used as inputs. Orange and blue arrows show the paths taken by the training and test data, respectively; the red arrow depicts the clustering process.

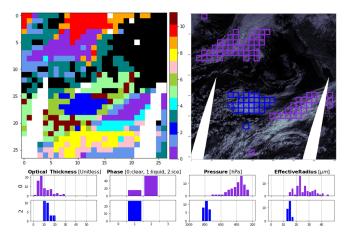


Fig. 4. Results of clustering results via autoencoder of MOD02 level 1B radiances in a representative MODIS swath image (December 1st, 2015, 11-44 N, 144-112 E; center of scene is shown in Fig. 1). Orbital coverage gaps leave missing data on sides of swath. Top left: labeled patches, classified into 12 clusters. Color bar shows cluster number in the 491 patches used; white indicates no data or invalid data; black indicates patches with <30% cloud pixels. Top right: raw visible image (band 1, which is not used as input to the autoencoder), with clusters #0 (violet) and #2 (blue) highlighted. Bottom: histograms of path-mean values of four derived cloud physics parameters in clusters #0 and #2: optical thickness, phase, cloud top pressure, and effective radius. Cluster #2 captures stratocumulus and Cluster #0 two instances of high-altitude cirrus.

III. EVALUATION

We report on three initial evaluations of our framework's capabilities.

As a first, we evaluate the physical reasonableness of assigned cluster labels in our full workflow. That is, we ask whether clusters are associated with reasonable distributions of patch-mean values of physical variables. Fig. 4 shows results for the representative MODIS swath also shown in Figure. 1. Left panel shows the cluster labels assigned to each patch in the image; right panel shows the raw visible image (band 1) and highlights patches assigned to two selected clusters (#0 and #2); and bottom panels show the distributions for these patches of four derived cloud physics parameters. Clustering is clearly correlated with meaningful physical cloud attributes: cluster #2 (blue) is stratocumulus and cluster #0 is cirrus, likely convective outflow.

We then conduct a simple test of whether our clustering via autoencoder captures richer and more meaningful information on cloud distributions and properties than can be provided by the deterministic algorithms used to produce derived cloud parameters. To have scientific value, our framework must produce information beyond that encoded in MOD06 products. We apply agglomerative clustering directly to MOD06 physics parameters and evaluate how well patches are classified

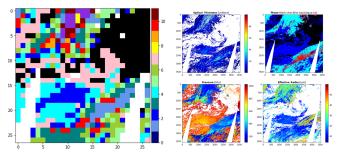


Fig. 5. Results of clustering based on patch-mean values of five MOD06 parameters (the 4 shown in Fig. 5 plus cloud water path) in the same swath as in Fig. 4. Left: labeled patches, classified into 12 clusters, with same figure conventions as in Fig. 5. Right: spatial distribution (heat maps) of values of four MOD06 cloud physics parameters. Clustering based on patch-mean parameter value produces less spatially coherent assigned classes than in Fig. 4 and does not capture important physical gradients, e.g. the sharp transition in effective radius at lower right.

without the guidance of dimension-reduced MOD02 radiance information. Results suggest that clustering via autoencoder produces classes that are spatially more cohesive and that better capture important physical transitions. (Compare left panels of Figures 4 and 5.)

Finally, we examine the spatial distribution of the latent representation itself using t-Distributed Stochastic Neighbor Embedding (t-SNE). This nonlinear dimensionality reduction technique maps each point in a high-dimensional space to a two-dimensional point such that similar objects are placed near to each other and dissimilar objects far apart, with high probability [14]. Resulting patch clusters are cohesive and distinct, suggesting that agglomerative clustering within our latent representation meaningfully separates different patch types (Figure 6).

IV. CONCLUSIONS

We describe here a prototype application of unsupervised learning to the problem of automated classification of clouds in multi-spectral satellite imagery. Our convolutional autoencoder generates a latent representation that, when clustered, yields physically meaningful cloud classes that pass a number of requisite tests for a scientifically useful tool. Assigned classes appropriately produce spatially coherent classifications, and capture meaningful aspects of cloud physics without being reproducible from mean values of physics parameters alone. This work supports the possibility of using unsupervised data-driven frameworks for automated cloud classification and pattern discovery without requiring the prior hypothesis of ground-truth labeled data. While results here are preliminary, they suggest that similar

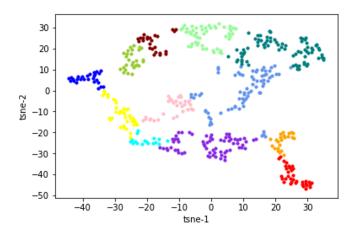


Fig. 6. A t-SNE visualization of the latent representations of the MOD02 patches of Figure 4, with cluster assignments represented by the same colors as in that figure. Patches in each cluster are projected near to each other in the t-SNE map.

frameworks can be used to analyze multi-year, global data to track short-term evolution during cloud lifecycles and long-term trends in the distribution of cloud features and characteristics.

More generally, unsupervised learning methods have broad potential applicability in the Earth sciences. In the satellite era, a primary challenge to environmental scientists is not gathering data but finding meaning in overwhelmingly large datasets. Unsupervised learning has the potential to reveal patterns directly learned from observation, which can then provide new insights and help diagnose drivers of system behavior.

ACKNOWLEDGMENTS

This work was supported by the Center for Robust Decision-making on Climate and Energy Policy (RD-CEP), NSF award SES-1463644, and used computers at the U.Chicago Research Computing Center and the Argonne Leadership Computing Facility, a DOE Office of Science User Facility, contract DE-AC02-06CH11357.

REFERENCES

- [1] M. Yoshiaki, K. Yoshiyuki, Y. Ryuj, Y. Tsuyoshi, Y. Hisashi, and T. Hirofumi, "Deep moist atmospheric convection in a subkilometer global simulation," Geophysical Research Letters, vol. 40, no. 18, pp. 4922-4926, 2013.
- [2] L. R. Jonathan, C. Weger, K. S. Sailes, and M. W. Ronaldo, "A neural network approach to cloud classification," IEEE Transactions on Geoscience and Remote Sensing, vol. 28, no. 5, pp. 846-855, 1990.
- [3] T. Bin, A. S. Mukhtiar, R. A. S. Mahmood, H. V. H. Thomas, and L. R. Donald, "A study of cloud classification with neural networks using spectral and textural features," IEEE Transactions on Neural Networks, vol. 10, no. 1, pp. 846-855, 1999.

- [4] W. Robert and L. H. Dennis, "Spatial variability of liquid water path in marine low cloud: The importance of mesoscale cellular convection," Journal of Climate, vol. 19, no. 9, pp. 1748-1764, 2005.
- [5] Z. Jinglin, L. Pu, Z. Feng, and S. Qianqian, "CloudNet: Ground-based cloud classification with deep convolutional neural network," Geophysical Research Letters, vol. 45, no. 16, pp. 8665-8672, 2018.
- [6] G. E. Hinton and S. Z. Richard, "Autoencoders, minimum description length and Helmholtz free energy," in Advances in Neural Information Processing Systems 6, pp. 3-10, Morgan-Kaufmann, 1994.
- [7] Z. Wang, E. P. Simoncelli, and A. C. Bovik, "Multiscale structural similarity for image quality assessment," in 37th Asilomar Conference on Signals, Systems & Computers, vol. 2, pp. 1398-1402, Ieee, 2003.
- [8] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli, et al., "Image quality assessment: from error visibility to structural similarity," IEEE Transactions on Image Processing, vol. 13, no. 4, pp. 600-612, 2004.
- [9] D. P. Kingma and J. L. Ba, "Adam: A method for stochastic optimization," International Conference on Learning Representations, 2015.
- [10] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al., "Tensorflow: A system for large-scale machine learning," in 12th USENIX Symposium on Operating Systems Design and Implementation, pp. 265-283, 2016.
- [11] A. Sergeev and M. Del Balso, "Horovod: Fast and easy distributed deep learning in Tensorflow," arXiv preprint arXiv:1802.05799, 2018.
- [12] S. Platnick, K. G. Meyer, M. D. King, B. Marchan, T. G. Arnold, Z. Zhang, P. A. Hubanks, R. E. Holz, P. Yang, W. L. Ridgway, and J. Riedi, "The MODIS cloud optical and microphysical products: Collection 6 updates and examples from Terra and Aqua," IEEE Transactions on Geoscience and Remote Sensing, vol. 55, no. 1, pp. 502-525, 2017.
- [13] B. A. Baum, P. W. Menzel, R. A.Frey, D. C. Tobin, R. E. Holz, S. A. Ackerman, A. K. Heidinger, and P. Yang, "MODIS cloud-top property refinements for collection 6," Journal of Applied Meteorology and Climatology, vol. 51, no. 6, pp. 1145-1163, 2012.
- [14] L. van der Maaten and G. Hinton, "Visualizing highdimensional data using t-SNE," Journal of Machine Learning Research, no. 9, pp. 2579-2605, 2008.