

SEMI-SUPERVISED DETECTION OF EXTREME WEATHER EVENTS IN LARGE CLIMATE DATASETS

Evan Racah¹, Christopher Beckham², Tegan Maharaj²

Prabhat¹, Christopher Pal²

¹ Lawrence Berkeley National Lab, Berkeley, CA, eracah, prabhat@lbl.gov.

² École Polytechnique de Montréal, firstname.lastname@polymtl.ca.

ABSTRACT

The detection and identification of extreme weather events in large scale climate simulations is an important problem for risk management, informing governmental policy decisions and advancing our basic understanding of the climate system. Recent work has shown that fully supervised convolutional neural networks (CNNs) can yield acceptable accuracy for classifying well-known types of extreme weather events when large amounts of labeled data are available. However, there are many different types of spatially localized climate patterns of interest (including hurricanes, extra-tropical cyclones, weather fronts, blocking events, etc.) found in simulation data for which labeled data is not available at large scale for all simulations of interest. We present a multichannel spatiotemporal encoder-decoder CNN architecture for semi-supervised bounding box prediction and exploratory data analysis. This architecture is designed to fully model multichannel simulation data, temporal dynamics and unlabelled data within a reconstruction and prediction framework so as to improve the detection of a wide range of extreme weather events. Our architecture can be viewed as a 3D convolutional autoencoder with an additional modified one-pass bounding box regression loss. We demonstrate that our approach is able to leverage temporal information and unlabelled data to improve localization of extreme weather events. Further, we explore the representations learned by our model in order to better understand this important data, and facilitate further work in understanding and mitigating the effects of climate change.

1 INTRODUCTION

Climate change is one of the most important challenges facing humanity in the 21st century and climate simulations are one of the only viable mechanisms for understanding the future impact of various carbon emission scenarios and intervention strategies. Large climate simulations produce massive datasets: a single 30-year run from a 25-km resolution model produces on the order of 10TB of multi-variate data. This scale of data makes post-processing such datasets to make quantitative assessments challenging and as a result, climate analysts and policy makers typically take global and annual averages of temperature or sea-level rise. These quantities are very coarse measurements and while they are perfect for public and media consumption, they ignore spatially (and temporally) resolved extreme weather events such as extra-tropical cyclones and tropical cyclones (hurricanes). Because the general public and policy makers are concerned about the local impacts of climate change, it is critical that we be able to analyze trends in extreme weather events which can have dramatic and tragic impacts on local and national populations and economies.

The task of finding extreme weather events in climate data has some similarities to the task of detecting objects and activities in video streams - a popular application area for deep learning techniques. An important difference is that in the case of climate data the 'video' has 16 or more 'channels' of information (such as water vapour, pressure and temperature), while conventional video only has 3 (RGB). In addition, these climate simulations do not share the same statistics as natural images. As a result, we cannot build off of some of the computer vision communities' efforts, such as pretrained weights from VGG or AlexNet with ImageNet.

Deep neural networks, especially deep convolutional neural networks, have enjoyed breakthrough success in recent recent years, achieving state-of-the-art results on many benchmark datasets (Krizhevsky et al. (2012); He et al. (2015); Szegedy et al. (2015)) and also compelling results on many practical tasks such as disease diagnosis (Hosseini-Asl et al. (2016)), facial recognition (Parkhi et al. (2015)), autonomous driving (Chen et al. (2015)), and many others. Furthermore, deep neural networks have also been very effective in the context of unsupervised and semi-supervised learning; some recent examples include variational autoencoders (Kingma & Welling (2013)), adversarial networks (Goodfellow et al. (2014); Makhzani et al. (2015)), and what-where autoencoders (Zhao et al. (2015)).

In this paper we present a 3D convolutional architecture for the task of semi-supervised climate event detection. Our contribution is the creation of a 3D spatiotemporal autoencoding CNN architecture for bounding box prediction and the exploration of its use for this important problem. Our architecture is capable of semi-supervised event detection, as well as exploratory data analysis through clustering representations learned using this approach. Our framework allows climate researchers to identify important climate events which have not yet been precisely defined, and could provide a basis for mining large climate simulation data archives and potentially near-real-time prediction of extreme weather events.

2 RELATED WORK

2.1 DEEP LEARNING FOR CLIMATE AND WEATHER DATA

The climate science community primarily relies on expert engineered systems and ad-hoc rules for characterizing climate and weather patterns. Prabhat et al. (2012; 2015) have implemented an assortment of heuristics in the TECA MapReduce framework to process large scale climate datasets on high-performance computing (HPC) platforms. Using the output of TECA analysis as ground truth, Liu et al. (2016) demonstrated for the first time that convolutional architectures could be successfully applied to predict the class label for two extreme weather event types. Their work considered the binary classification task on centered, cropped patches from 2D (single-timestep) multi-channel images. Like Liu et al. (2016) we use TECA's planetary-scale output as ground truth, but we build on this work by: 1) using uncropped images, 2) considering the temporal axis of the data 3) doing multi-class bounding box detection and 4) exploring a semi-supervised approach based on a hybrid predictive and reconstructive model.

Some recent work has applied deep learning methods to weather forecasting. Xingjian et al. (2015) have explored a convolutional LSTM architecture (described in 2.2 for predicting future precipitation on a local scale (i.e. the size of a city) using radar echo data. In contrast, we focus on extreme event detection on planetary-scale data. Our aim is to capture patterns which are very local in time (e.g. a hurricane may be present in half a dozen sequential frames) compared to the scale of our underlying climate data, consisting of global simulations over many years. As such we use full blown 3D CNNs as they make more sense for our detection application compared to LSTMs, whose strength is in capturing long-term dependencies.

2.2 RELATED METHODS AND MODELS

Following the dramatic success of CNNs in static 2D images, a wide variety of CNN architectures have been explored for video, ex. (Karpathy et al., 2014; Yao et al., 2015; Tran et al., 2014). The details of how CNNs are extended to capture the temporal dimension are important and Karpathy et al. (2014) explored different strategies for fusing information from 2D CNN subcomponents. In contrast, Yao et al. (2015) created 3D volumes of statistics from low level image features.

CNNs have been combined with RNNs for modeling video and other sequence data and we briefly review some relevant video models here. The most common approach to modeling sequential images is to feed single-frame representations from a CNN at each timestep to an RNN. This approach has been examined for a number of different types of video (Donahue et al., 2015). Srivastava et al. (2015) have explored an LSTM architecture for the unsupervised learning of video representations using a pretrained CNN representation as input. Another popular model, also used on 1D data, is a convolutional RNN, wherein the hidden-to-hidden transition layer is 1D convolutional (i.e. the state

is convolved over time). Ballas et al. (2016) combine these ideas, applying a convolutional LSTM to frames processed by a (2D) CNN.

The 3D CNNs we use here are based on 3-dimensional convolutional filters, taking the height, width, and time axes into account for each feature map, as opposed to aggregated 2D CNNs. This approach was studied in detail in Tran et al. (2014). 3D convolutional neural networks have been used for various tasks ranging from human activity recognition (Ji et al., 2013), to large-scale YouTube video classification (Karpathy et al., 2014), and video description (Yao et al., 2015). Hosseini-Asl et al. (2016) use a 3D convolutional autoencoder for diagnosing Alzheimer’s disease through MRI - in this case, the 3 dimensions are height, width, and depth. Whitney et al. (2016) use 3D (height, width, depth) filters to predict consecutive frames of a video game for continuation learning. Some recent work has also examined ways to use CNNs to generate animated textures and sounds Xie et al. (2016). This work is similar to our approach in that it uses a 3D convolutional encoder, but they use a stochastic approach whereas we use a deterministic autoencoder with a 3D convolutional decoder for the unsupervised component of our overall objective function.

Stepping back, our approach and model is closely related to Tran et al. (2014) in that we use a 3D deconvolutional CNN approach to construct the autoencoding component of our model. Our approach is related to recent work from Zhang et al. (2016) (and others) in that we also use a hybrid prediction and autoencoding architecture, but our model is for multidimensional 3D data and yields multiple bounding box predictions. Our formulation for the bounding box prediction loss is inspired by the approach in Redmon et al. (2015) and extended in Ren et al. (2015), the single shot multiBox detector formulation used in Liu et al. (2015) and the seminal bounding box work in OverFeat Sermanet et al. (2013).

3 THE MODEL

We use a 3D convolutional encoder-decoder architecture (the code can be found here: <https://github.com/eracah/hur-detect/>); feature maps are 3D: time, height, and width. The encoder part of this architecture is described in Table 2. The decoder is the equivalent structure in reverse, using tied weights and deconvolutional layers. As we take a semi-supervised approach, the code (bottleneck) layer of the autoencoder is used as the input to the loss layers, which make predictions for bounding box location and size, the class associated with the bounding box, and the confidence (sometimes called ‘objectness’) of the bounding box. In training, we input one day’s simulation at a time (8 time steps, 16 variables) and reconstruct all 8 time steps, while only predicting the bounding box for 4 time steps, since only 4 out of 8 time steps are labelled, so a 1:1 ratio of unlabelled data to labelled data. The entire architecture is depicted in Figure 1.

For comparison, to evaluate how useful the time axis is to recognizing extreme weather events, we also run experiments with a 2D (width, height) version of this architecture.

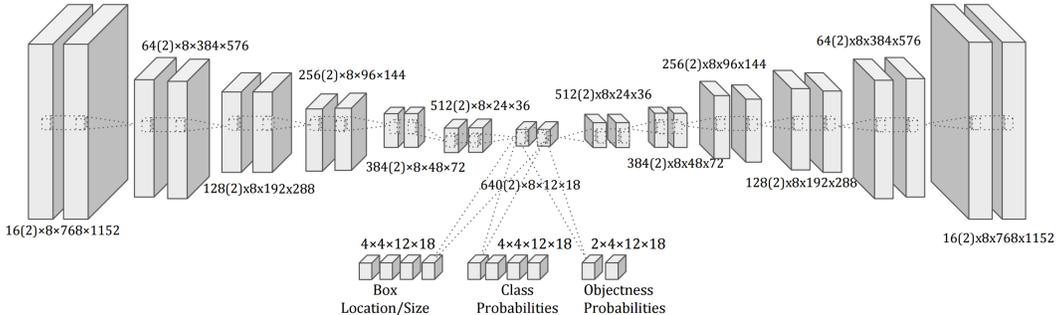


Figure 1: Diagram of the 3D semi-supervised architecture. Parentheses denote subset of total dimension shown. Only two (2) feature maps per layer shown for ease of visualization in encoder-decoder. All feature maps are shown in bounding box scoring portion of the network

Table 1: Details of the encoder part of the architecture used in our experiments. Input (F×T×H×W) = F feature maps, T time axes, H height, and W width. For convolutions, F(T,H,W) = the filter size for the time, height, and width axis respectively, S(T,H,W) = the stride amount for the time, height and width axis respectively, P(T,H,W) = the zero padding and N = the number of output feature maps. Leaky ReLU with 0.1 nonlinearities are used for all layers.

Layer	Output size
Input (16x8x768x1152)	16, 8, 768, 1152
Conv F(3,5,5),S(1,2,2),P(1,2,2),64	64, 8, 384, 576
Conv F(3,5,5),S(1,2,2),P(1,2,2),128	128, 8, 192, 288
Conv F(3,5,5),S(1,2,2),P(1,2,2),256	256, 8, 96, 144
Conv F(3,5,5),S(1,2,2),P(1,2,2),384	384, 8, 48, 72
Conv F(3,5,5),S(1,2,2),P(1,2,2),512	512, 8, 24, 36
Conv F(3,5,5),S(1,2,2),P(1,2,2),640	640, 8, 12, 18

Table 2: Details of the loss layers of the architecture. Note that the encoder output is given as input to each of the 3 ‘Scorer’ layers.

Name	Layer	Output size
Encoder Output	Input (640x8x12x18)	640, 8, 12, 18
Class Scorer	Conv + Softmax F(3,3,3),S(2,1,1),P(1,1,1),4	4, 4, 12, 18
Objectness Scorer	Conv + Softmax F(3,3,3),S(2,1,1),P(1,1,1),2	2, 4, 12, 18
B-box Location/Size Scorer	Conv + ReLU F(3,3,3),S(2,1,1), P(1,1,1),4	4, 4, 12, 18

In order to regress bounding boxes, we use some techniques from Redmon et al. (2015), Ren et al. (2015), and Liu et al. (2015). We split the original 768 by 1152 image into a 12x18 grid of 64x64 ‘‘anchor’’ boxes. We then make a guess for a box at each grid point by transforming the representation to 216 ‘‘scores’’. Each score encodes three pieces of information: how much the predicted box differs in size and location from the anchor box, how confident we are that an object of interest is in the predicted box, and the class probability distribution for that object. We compute each component of the score by several 3x3 convolutions applied to the 640 12x18 feature maps from output of the last layer of the encoder (‘‘code layer’’). Because each set of pixels in each feature map at a given x, y coordinate can be thought of as a learned representation of the climate data in a 64x64 patch of the input image, we can think of the 3x3 convolutions as using the representation of a 192 x 192 neighborhood from the input image as context to determine the box and object centered in the given 64x64 patch. Our approach is similar to Liu et al. (2015) and Sermanet et al. (2013), which use convolutions from small local receptive field filters to regress boxes. This choice is motivated by the fact that extreme weather events occur in relatively small spatiotemporal volumes, with the ‘background’ context being highly consistent across event types and between events and non-events. This is in contrast to that of Redmon et al. (2015), which uses a fully connected layer to consider the whole image as context. Their approach is appropriate to the objective of object identification in natural images, where there is often a strong relationship between background and object.

The bounding box regression loss is determined as follows:

$$L_{sup} = \frac{1}{N}(L_{box} + L_{conf} + L_{cls}), \quad (1)$$

where N is the number of time steps in the minibatch, and L_{box} is defined as:

$$L_{box} = \alpha \sum_i \mathbb{1}_i^{obj} R(u_i - u_i^*) + \beta \sum_i \mathbb{1}_i^{obj} R(v_i - v_i^*), \quad (2)$$

where $i \in [0, 216)$ is the index of the anchor box for the i th grid point (12x18 grid, so 216 grid points). where $\mathbb{1}_i^{obj}$ is 1 if an object is present at the i th grid point and 0 if not, $R(z)$ is the smooth L1 loss (Ren et al. (2015)), $u_i = (t_x, t_y)_i$ and $u_i^* = (t_x^*, t_y^*)_i$, $v_i = (t_w, t_h)_i$ and $v_i^* = (t_w^*, t_h^*)_i$

where t is the parametrization defined in (Ren et al. (2015)) such that:

$$t_x = (x - x_a)/w_a, t_y = (y - y_a)/h_a, t_w = \log(w/w_a), t_h = \log(h/h_a)$$

$$t_x^* = (x^* - x_a)/w_a, t_y^* = (y^* - y_a)/h_a, t_w^* = \log(w^*/w_a), t_h^* = \log(h^*/h_a),$$

where (x_a, y_a, w_a, h_a) is the center coordinates and height and width of the closest anchor box, (x, y, w, h) are the predicted coordinates and (x^*, y^*, w^*, h^*) are the ground truth coordinates.

L_{conf} is the weighted cross entropy of the object confidence (probability distribution for whether object is present in a grid cell or not):

$$L_{conf} = \sum_i \mathbb{1}_i^{obj} [(-\log(p(obj)_i))] + \gamma * \sum_i \mathbb{1}_i^{noobj} [-\log(p(\overline{obj})_i)] \quad (3)$$

L_{cls} is just the cross-entropy between the one-hot encoded class distribution and the softmax predicted class distribution, evaluated only for predicted boxes at the grid points containing a ground truth box:

$$L_{cls} = \sum_i \mathbb{1}_i^{obj} \sum_{c \in classes} -p^*(c) \log(p(c)) \quad (4)$$

This loss is similar in spirit to YOLO (Redmon et al. (2015)), but has a few differences. One, the object confidence and class probability terms in YOLO are squared-differences between ground truth and prediction, so we make the logical change to cross-entropy, which is more appropriate for probabilities and is used in the region proposal network from Faster R-CNN Ren et al. (2015) and the network from Liu et al. (2015) for the object probability term and the class probability term respectively.

In addition, there is a difference in the parametrization for the coordinates and the size of the bounding box.

In YOLO, the parametrizations for x and y are equivalent to Faster R-CNN's t_x and t_y for an anchor box the same size as the patch it represents (64x64). For the w and h parametrizations in YOLO, they would only be equivalent to Faster-RCNN's t_h and t_w for a 64x64 anchor box if the anchor box had a height and width equal to the size of the whole image and if there were no log transform in the faster-RCNN's parametrization. We found in practice that both these differences in YOLO did not work well in practice. For one, without the log term and using a RELU transform as the nonlinearity with a standard weight initialization scheme sampled from a distribution centered around 0, most outputs (more than half) for the parametrization were zero or close to zero as expected (at least at the beginning of training); when the parametrization was converted to raw width and height, this resulted in 0 height and width boxes, so technically the default boxes were of 0 width and height. In practice, it took many epochs for the network to learn to resize 0 area boxes to the correct size due to all the other terms of the loss that were being optimized. Adding the log term, would in effect make the "default" box (an output of 0) equal to a box of height and width h_a and w_a , which in YOLO's case was the height and width of the entire image. This once again proved to be a problem as it took a long time to resize the boxes to the size of the much smaller ground truth boxes. Making h_a and w_a equal to 64x64 made training a lot more efficient, as the optimization was focused more on picking *which* box contained an object and not as much on what size the box should be. The last difference was YOLO used squared difference between predicted and ground truth for the coordinate parametrizations, as opposed to smooth L1. We used smooth L1 due its lower sensitivity to outlier predictions Ren et al. (2015).

The loss is a weighted combination of reconstruction error and bounding box regression loss:

$$L = L_{sup} + \lambda L_{rec}, \quad (5)$$

where L_{unsup} is the classic squared difference between input and reconstruction:

$$L_{rec} = \|X - X^*\|_2^2, \quad (6)$$

4 EXPERIMENTS AND DISCUSSION

4.1 CLIMATE DATA

The climate science community utilizes three flavors of global datasets: observational products (satellite, gridded weather station); reanalysis products (obtained by assimilating disparate observational products into a climate model) and simulation products. In this study, we analyze output from the third category: free-running climate simulations with prescribed initial and boundary conditions. In particular, we consider a 27-year run of the CAM5 (Community Atmospheric Model v5); configured at 25-km spatial resolution (Wehner et al. (2015)), with ground-truth labelling for four extreme weather events: Tropical Depressions (TD) Tropical Cyclones (TC), Extra-Tropical Cyclones (ETC) and Atmospheric Rivers (AR).

At 25-km resolution, each snapshot of the global atmospheric state in the CAM5 model output corresponds to a 768x1152 image. Each image is comprised of 16 variables (including surface temperature, surface pressure, precipitation, zonal wind, meridional wind, humidity, cloud fraction, water vapor, etc.) stored in double-precision floating point representation. While climate models are run on a 3D grid, with the vertical dimension corresponding to 30 levels; we only consider surface quantities (i.e. 2D data) in this study. Model output is stored every 3 hours.

The training data (which can be downloaded here: <http://portal.nersc.gov/project/dasrepo/climate/>) consists of 365 days of simulations from 1979 with 8 frames a day (every three hours), but labels are only every 6 hours, so one day for supervised is 4 frames and 8 for semi-supervised. The test set consisted of 365 days from 1984. Table 3 shows the breakdown of the dataset splits for each class.

TECA, the ground truth labelling framework, tries to implement the state of the art heuristics for determining 'ground truth' data for the four types extreme weather events. However, it is entirely possible there are errors in the labeling: for instance, there is little agreement in the climate community on a standard heuristic for capturing Extra-Tropical Cyclones (Neu et al. (2013)); Atmospheric Rivers have been extensively studied in the northern hemisphere (Lavers et al. (2012); Dettinger et al. (2011)), but not in the southern hemisphere; and spatial extents of such events not universally agreed upon. As such, there is potential for many false negatives, resulting in partially annotated images. This, in addition to lower representation for classes, like AR's and Tropical Depressions, is part of our motivation in exploring semi-supervised methods to better understand the features underlying extreme weather events. Moreover, although there is more than enough labelled data for our training set sizes, not all flavors of climate datasets can be labelled by TECA due to its "hard-coded" set up for a specific resolution simulation with specific parameters, so demonstrating the success of semi-supervised learning for climate pattern detection is important. Lastly, TECA only labels four types of events, so detection of other climate events would require hand-labelling, which would likely result in only 10's of labels, necessitating semi-supervised approaches.

Table 3: Class frequency breakdown for Tropical Cyclones (**TC**), Extra-Tropical Cyclones (**ETC**), Tropical Depressions (**TD**), and United States Atmospheric Rivers (**US-AR**). Raw counts in parentheses.

Data	TC %	ETC (%)	TD (%)	US-AR (%)
Train	42.32 (3190)	46.57 (3510)	5.74 (433)	5.36 (404)
Test	39.04 (2882)	46.47 (3430)	9.44 (697)	5.04 (372)

4.2 FRAMEWISE RECONSTRUCTION

As a simple experiment, we first train a 2D convolutional autoencoder on the data, treating each timestep as an individual training example (everything else about the model is as described in Section 3), in order to visually assess reconstructions and ensure reasonable accuracy of detection. Figure 2 shows the original and reconstructed feature maps for the 16 climate variables of one image in the training set. Though we are showing reconstructions on the training set, the reconstruction loss on the validation set was also similar.

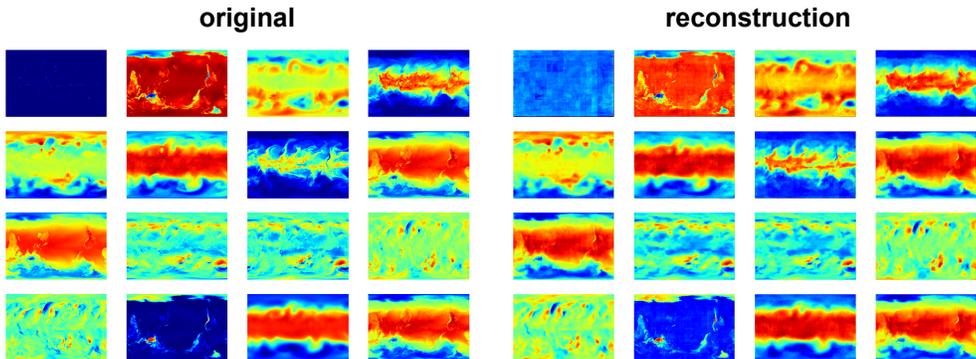


Figure 2: Feature maps for the 16 channels in a select image in the training set (left) versus their reconstructions from the 2D convolutional autoencoder (right).

As the reconstruction visualizations suggest, the 2D convolutional autoencoder architecture does a good job of encoding spatial information from climate images. In the upcoming subsections, we utilize the full 3D architecture for the detection of climate events.

4.3 COMPARING FULLY SUPERVISED AND SEMI-SUPERVISED DETECTION

We run two experiments of 3D supervised and semi-supervised bounding box prediction, where values for α, β, γ were selected with experimentation and some inspiration from (Redmon et al. (2015)) to be 5, 7 and 0.5 respectively.

A lower value for γ also can be interpreted as way to deal with ground truth false negatives, as we want more emphasis on pushing up the confidence of true positive examples in order to learn the features of the four events, as opposed to pushing down the confidence of false negative examples, which may interfere with generalization. We used the Adam optimizer (Kingma & Ba (2014)) with a base learning rate of 0.0001. We used a weight decay coefficient of 0.0005.

In Table 4 and Table 5, we present 2D and 3D supervised and semi-supervised results for various settings of λ . Because the 3D model has inherently higher capacity in terms of number of parameters than the 2D model, we also show some experiments with higher capacity 2D models by doubling the number of filters that are in each layer in Table 2. Average Precision (AP) was calculated in the manner of ImageNet (Russakovsky et al. (2015)), where we integrated the precision-recall curve for each class. Results are shown for two modes of evaluation: denoting a true positive as a bounding box prediction with IOU with the ground truth box of 0.1 and 0.5. Mean average precision was just the average of the AP over the four classes. Furthermore, in Figure 3 we provide bounding box predictions shown on 2 consecutive (6 hours in between) simulation frames comparing the 3D supervised vs 3D semi-supervised model predictions.

As we can see from Table 4, Table 5 and Figure 3, 3D and semi-supervised approaches help to some degree for rough localization of the weather events (IOU=0.1). Namely, semi-supervised approaches help the performance of the 3D model for IOU=0.1. In contrast, semi-supervised approach 2D does not have the same effect for IOU=0.1. It is important to note that a more thorough hyperparameter search for λ and other parameters may yield better results for semi-supervised. Also of note is that the 3D models perform significantly better than their 2D counterparts for ETC and TC (hurricane) classes. This potentially means that the time evolution of these weather events is an important criteria for discriminating them. In addition, the semi-supervised model significantly improves the ETC and TC performance, which suggests unsupervised shaping of the spatio-temporal representation is important for these events. It would be interesting to further explore and confirm the spatio-temporal representation learned for these weather events in the 3D models.

In contrast, Table 4 shows that for a more stringent, but standard detection metric (IOU=0.5), the semi-supervised models do much worse than their supervised counterparts. This could potentially be improved by using a dataset from a different year for obtaining unlabelled data, so as not to

have as much correlation with the labelled data. Moreover, the 3D supervised models perform worse than the 2D supervised models and 3D semi-supervised perform comparably to the corresponding 2D semi-supervised models. The potential reason for this can be explained to a certain degree with Figure 3. As shown in the figure, the 3D models do a good job roughly localizing the various events. However, as mentioned in Section 3, the network has a hard time adjusting the size of the boxes. As such, in this figure we see mostly boxes of size 64x64. For example, for TD's, which are usually much smaller than 64x64 and for AR's, which are always much bigger than 64x64, a 64x64 box roughly centered on the event is sufficient to count as a true positive with $\text{IOU}=0.1$, but not for the more stringent $\text{IOU}=0.5$. This lead to a large dropoff in performance for AR's and TD's and a sizable dropoff in the variable sized TC's as well as shown in Table 5. While the model has the capacity to adjust the sizes, it can take a while to converge on more different sizes; with longer training time we expect it (and have seen it for 2D) to be able to detect other size boxes, but due to the size of the model we did not have time for longer training. Hence, we believe the 2D model was able to do more size and location adjustments in its training and as a result, it's AP dropoff for the 4 classes was not as large. However, the 3D semi-supervised model's ability to better roughly localize events ($\text{IOU}=0.1$) is still useful for the application and shows promise toward future use in an improved system. Future work would involve harnessing anchor boxes of many shapes and sizes, like (Ren et al. (2015), Liu et al. (2015), which we believe would greatly increase the detection accuracy for $\text{IOU}=0.5$.

In Figure 4 we provide t-SNE visualisation of the first 7 days in the training sets for both 3D supervised (left) and semi-supervised (right) experiments. It is possible to make out some grouping of various spiral patterns in the latent space for the TC's and ETC's.

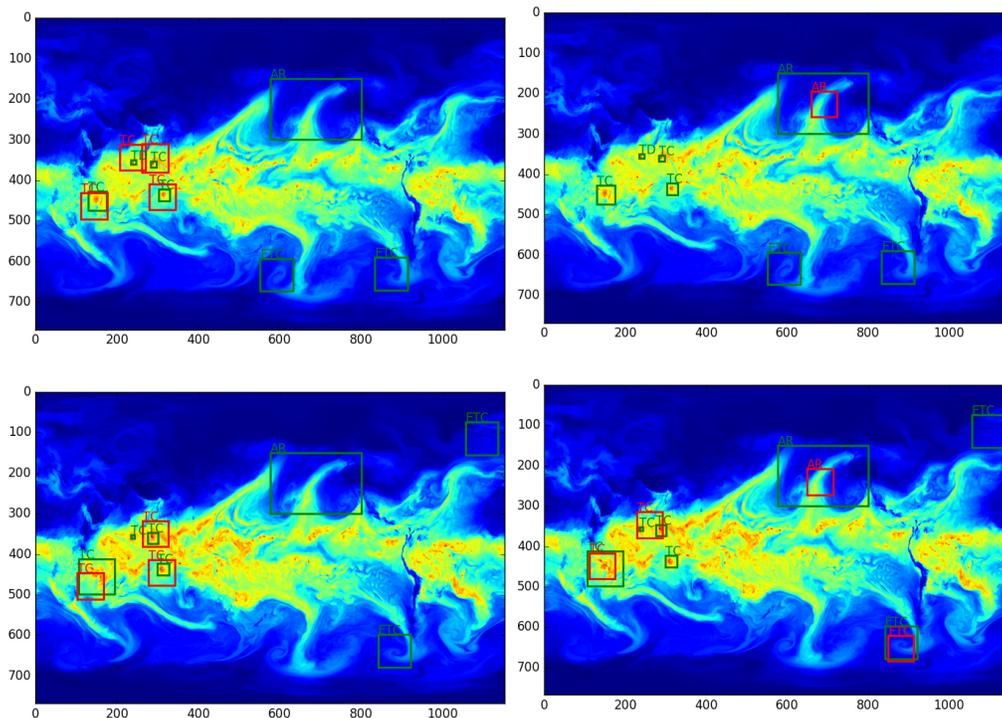


Figure 3: Bounding box predictions shown on 2 consecutive (6 hours in between) simulation frames (integrated water vapor column). Green = ground truth, red = high confidence predictions (confidence above 0.8, (Left) 3D Fully supervised model, (Right) 3D Semi-supervised model.

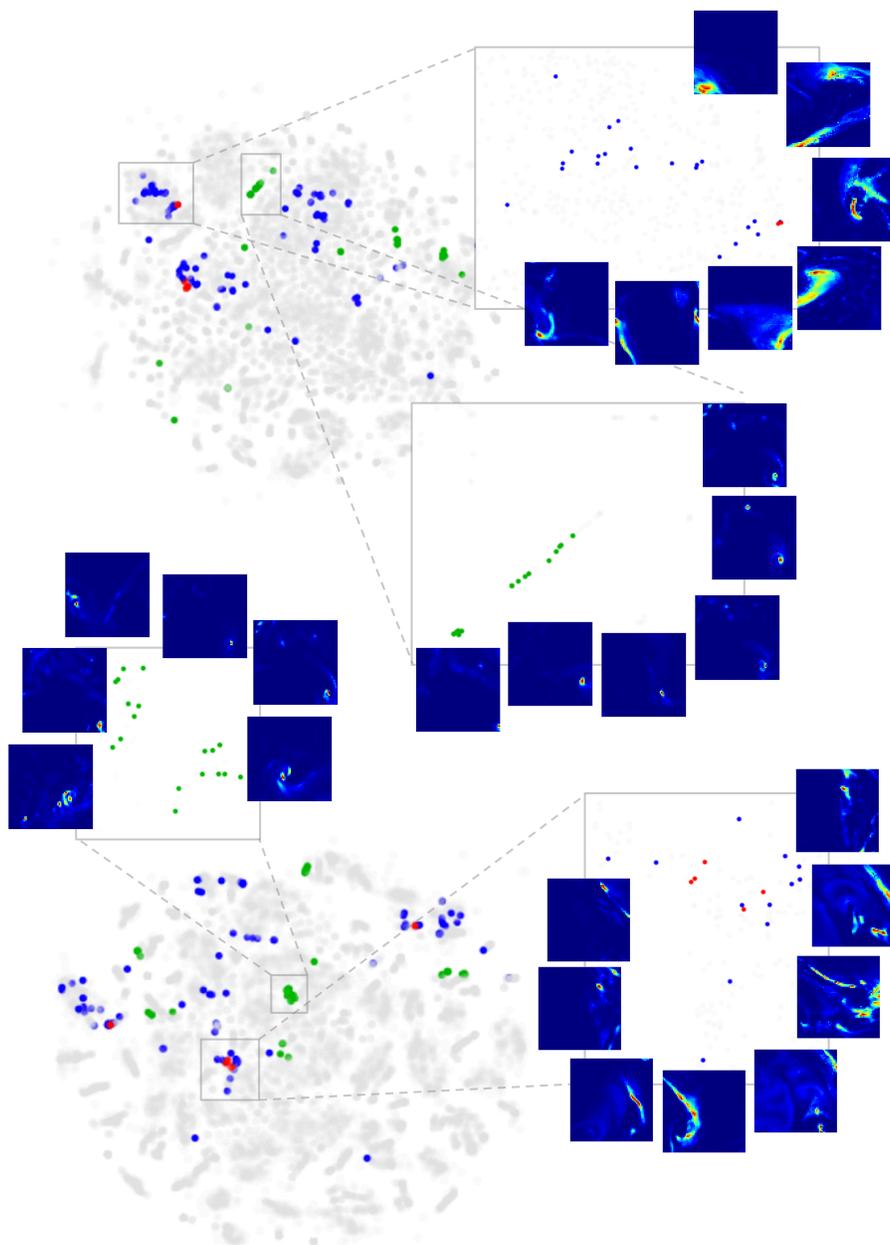


Figure 4: t-SNE visualisation of the first 7 days in the training set for both 3D supervised (top) and semi-supervised (bottom) experiments. Each frame (i.e. time step) in the 7 days has $12 \times 18 = 216$ vectors of length 640 (i.e. the number of feature maps in the final encode layer) in the coding layer (where each pixel in the 12×18 patch corresponds to a 64×64 patch in the original frame). These vectors are then fed to t-SNE and reduced to two dimensions for visualization. For both supervised and semi-supervised, we have zoomed into two clusters and sampled 64×64 patches from those clusters. Grey = unlabelled, yellow = tropical depression (not shown), green = hurricane, blue = ETC, red = AR. (The same t-SNE parameters were used for both experiments.)

5 CONCLUSIONS AND FUTURE WORK

We have explored semi-supervised methods for object detection and bounding box prediction using 3D CNNs. These architectures and approaches are motivated by finding extreme weather patterns;

Table 4: Accuracy Results: Mean Average Precision (mAP) for the models

Model	Mode	Parameters (millions)	λ	mAP (%)	
				(IOU=0.1)	(IOU=0.5)
2D	Supervised	66.53	0	51.42	16.98
2D	Semi-Supervised	66.53	10	48.85	9.24
2D	Semi-Supervised	66.53	1	51.11	6.21
2D	Supervised	16.68	0	49.21	15.49
2D	Semi-Supervised	16.68	1	44.01	7.71
3D	Supervised	50.02	0	51.00	11.60
3D	Semi-Supervised	50.02	1	52.92	7.31

Table 5: Accuracy Results By Class. Average Precision for each class. Frequency of each class in the test set shown in parentheses. First number is at IOU=0.1 and second number after the semicolon is at IOU=0.5 as criteria for true positive

Model	Mode	Parameters (millions)	λ	ETC (46.47%)		TC (39.04 %)		TD (9.44 %)		AR (5.04 %)	
				AP (%)	AP (%)	AP (%)	AP (%)	AP (%)	AP (%)		
2D	Sup	66.53	0	21.92; 14.42	52.26; 9.23	95.91; 10.76	35.61; 33.51				
2D	Semi	66.53	1	18.05; 5.00	52.37; 5.26	97.69; 14.60	36.33; 0.00				
2D	Semi	66.53	10	15.57; 5.87	44.22; 2.53	98.99; 28.56	36.61 ; 0.00				
2D	Sup	16.68	0	13.90; 5.25	49.74; 15.33	97.58; 7.56	35.63; 33.84				
2D	Semi	16.68	1	15.80; 9.62	39.49; 4.84	99.50 ; 3.26	21.26; 13.12				
3D	Sup	50.02	0	22.65; 15.53	50.01; 9.12	97.31; 3.81	34.05; 17.94				
3D	Semi	50.02	1	24.74 ; 14.46	56.40 ; 9.00	96.57; 5.80	33.95; 0.00				

a meaningful and important problem for society. Thus far, the climate science community has used hand-engineered criteria to characterize patterns. Our results indicate that there is much promise in considered deep learning based approaches. In addition, this work represents our first attempt at applying semi-supervised learning and exploratory data analysis (i.e. clustering in a learned representation) to climate data. In the future, we think it would be interesting to delve deeper into interpreting and visualizing what the network has learned. This would not only help open the "black box" of deep learning for skeptics in the climate community, but also potentially help resolve ambiguities in the definition of extreme weather events in addition to highlighting unknown underlying factors in extreme weather events. On a similar note, another goal would be to be able to segment every pixel/voxel in the multi-variate spatio-temporal volume as belonging to a particular weather pattern. Once we are able to decompose the climate system along various patterns, we will be able to conduct sophisticated conditional analysis: how much of global precipitation can be attributed to various patterns (cyclones vs. atmospheric rivers vs. fronts). Moreover, we think it would be interesting to apply the models learned from this CAM5 model output to satellite, reanalysis and other simulation model output. However, these datasets have different spatial resolution, variables, and missing values (in the case of satellite observations). Lastly, applying this extreme weather detection framework to simulations of years in the future would be a very informative exercise for the climate community and the world.

ACKNOWLEDGMENTS

This research used resources of the National Energy Research Scientific Computing Center (NERSC), a DOE Office of Science User Facility supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231. Code relies on open-source deep learning frameworks Theano (Bergstra et al.; Team et al., 2016) and Lasagne (Team, 2016), whose developers we gratefully acknowledge. We would also like to thank Yunjie Liu and Michael Wehner for providing access to the climate datasets; Alex Lamb and Thorsten Kurth for helpful discussions; and Samsung for their support of this research.

REFERENCES

- Nicolas Ballas, Li Yao, Chris Pal, and Aaron Courville. Delving deeper into convolutional networks for learning video representations. *In the Proceedings of ICLR. arXiv preprint arXiv:1511.06432*, 2016.
- James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio. Theano: A cpu and gpu math compiler in python.
- Chenyi Chen, Ari Seff, Alain Kornhauser, and Jianxiong Xiao. Deepdriving: Learning affordance for direct perception in autonomous driving. *In Proceedings of the IEEE International Conference on Computer Vision*, pp. 2722–2730, 2015.
- Michael D. Dettinger, Fred Martin Ralph, Tapash Das, Paul J. Neiman, and Daniel R. Cayan. Atmospheric rivers, floods and the water resources of california. *Water*, 3(2):445, 2011. ISSN 2073-4441. URL <http://www.mdpi.com/2073-4441/3/2/445>.
- Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. *In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *In Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *In Proceedings of the IEEE International Conference on Computer Vision*, pp. 1026–1034, 2015.
- Ehsan Hosseini-Asl, Georgy Gimel’farb, and Ayman El-Baz. Alzheimer’s disease diagnostics by a deeply supervised adaptable 3d convolutional network. 2016.
- Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu. 3d convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(1):221–231, 2013.
- Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. *In Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 1725–1732, 2014.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *In Advances in neural information processing systems*, pp. 1097–1105, 2012.
- David A. Lavers, Gabriele Villarini, Richard P. Allan, Eric F. Wood, and Andrew J. Wade. The detection of atmospheric rivers in atmospheric reanalyses and their links to british winter floods and the large-scale climatic circulation. *Journal of Geophysical Research: Atmospheres*, 117 (D20):n/a–n/a, 2012. ISSN 2156-2202. doi: 10.1029/2012JD018027. URL <http://dx.doi.org/10.1029/2012JD018027>. D20106.
- Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, and Scott Reed. Ssd: Single shot multibox detector. *arXiv preprint arXiv:1512.02325*, 2015.
- Yunjie Liu, Evan Racah, Prabhat, Joaquin Correa, Amir Khosrowshahi, David Lavers, Kenneth Kunkel, Michael Wehner, and William Collins. Application of deep convolutional neural networks for detecting extreme weather in climate datasets. 2016.

- Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, and Ian J. Goodfellow. Adversarial autoencoders. *CoRR*, abs/1511.05644, 2015. URL <http://arxiv.org/abs/1511.05644>.
- Urs Neu, Mirseid G. Akperov, Nina Bellenbaum, Rasmus Benestad, Richard Blender, Rodrigo Caballero, Angela Coccozza, Helen F. Dacre, Yang Feng, Klaus Fraedrich, Jens Grieger, Sergey Gulev, John Hanley, Tim Hewson, Masaru Inatsu, Kevin Keay, Sarah F. Kew, Ina Kindem, Gregor C. Leckebusch, Margarida L. R. Liberato, Piero Lionello, Igor I. Mokhov, Joaquim G. Pinto, Christoph C. Raible, Marco Reale, Irina Rudeva, Mareike Schuster, Ian Simmonds, Mark Sinclair, Michael Sprenger, Natalia D. Tilinina, Isabel F. Trigo, Sven Ulbrich, Uwe Ulbrich, Xiaolan L. Wang, and Heini Wernli. Imilast: A community effort to intercompare extratropical cyclone detection and tracking algorithms. *Bulletin of the American Meteorological Society*, 94(4):529–547, 2013. doi: 10.1175/BAMS-D-11-00154.1.
- Omkar M Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. In *British Machine Vision Conference*, volume 1, pp. 6, 2015.
- Prabhat, Oliver Rubel, Surendra Byna, Kesheng Wu, Fuyu Li, Michael Wehner, and Wes Bethel. Teca: A parallel toolkit for extreme climate analysis. *ICCS*, 2012.
- Prabhat, Surendra Byna, Venkatram Vishwanath, Eli Dart, Michael Wehner, and William D. Collins. Teca: Petascale pattern recognition for climate science. *CAIP*, 2015.
- Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. *CoRR*, abs/1506.02640, 2015. URL <http://arxiv.org/abs/1506.02640>.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. 2015.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *arXiv preprint arXiv:1312.6229*, 2013.
- Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. Unsupervised learning of video representations using lstms. *CoRR*, abs/1502.04681, 2, 2015.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–9, 2015.
- Lasagne Development Team. Lasagne, 2016. URL <https://lasagne.readthedocs.io/en/latest/>.
- The Theano Development Team, Rami Al-Rfou, Guillaume Alain, Amjad Almahairi, Christof Angermueller, Dzmitry Bahdanau, Nicolas Ballas, Frédéric Bastien, Justin Bayer, Anatoly Belikov, et al. Theano: A python framework for fast computation of mathematical expressions. *arXiv preprint arXiv:1605.02688*, 2016.
- Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. 2014.
- Michael Wehner, Prabhat, Kevin A. Reed, Dáithí Stone, William D. Collins, and Julio Bacmeister. Resolution dependence of future tropical cyclone projections of cam5.1 in the u.s. clivar hurricane working group idealized configurations. *Journal of Climate*, 28(10):3905–3925, 2015. doi: 10.1175/JCLI-D-14-00311.1.
- William F. Whitney, Michael Chang, Tejas Kulkarni, and Joshua B. Tenenbaum. Understanding visual concepts with continuation learning. 2016.

- Jianwen Xie, Song-Chun Zhu, and Ying Nian Wu. Synthesizing dynamic textures and sounds by spatial-temporal generative convnet. 2016.
- Shi Xingjian, Zhoung Chen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in Neural Information Processing Systems*, pp. 802–810, 2015.
- Li Yao, Atousa Torabi, Kyunghyun Cho, Nicolas Ballas, Christopher Pal, Hugo Larochelle, and Aaron Courville. Describing videos by exploiting temporal structure. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4507–4515, 2015.
- Yuting Zhang, Kibok Lee, and Honglak Lee. Augmenting supervised neural networks with unsupervised objectives for large-scale image classification. *arXiv preprint arXiv:1606.06582v1*, 2016.
- Junbo Zhao, Michael Mathieu, Ross Goroshin, and Yann Lecun. Stacked what-where auto-encoders. *arXiv preprint arXiv:1506.02351*, 2015.