
STCONVS2S: SPATIOTEMPORAL CONVOLUTIONAL SEQUENCE TO SEQUENCE NETWORK FOR WEATHER FORECASTING

ACCEPTED MANUSCRIPT

Rafaela Castro
CEFET/RJ
Rio de Janeiro, RJ, Brazil
rafaela.nascimento@eic.cefet-rj.br

Yania M. Souto
LNCC
Petrópolis, RJ, Brazil
yaniams@lncc.br

Eduardo Ogasawara
CEFET/RJ
Rio de Janeiro, RJ, Brazil
eogasawara@ieee.org

Fabio Porto
LNCC
Petrópolis, RJ, Brazil
fporto@lncc.br

Eduardo Bezerra
CEFET/RJ
Rio de Janeiro, RJ, Brazil
ebezerra@cefet-rj.br

ABSTRACT

Applying machine learning models to meteorological data brings many opportunities to the Geosciences field, such as predicting future weather conditions more accurately. In recent years, modeling meteorological data with deep neural networks has become a relevant area of investigation. These works apply either recurrent neural networks (RNN) or some hybrid approach mixing RNN and convolutional neural networks (CNN). In this work, we propose STConvS2S (Spatiotemporal Convolutional Sequence to Sequence Network), a deep learning architecture built for learning both spatial and temporal data dependencies using only convolutional layers. Our proposed architecture resolves two limitations of convolutional networks to predict sequences using historical data: (1) they violate the temporal order during the learning process and (2) they require the lengths of the input and output sequences to be equal. Computational experiments using air temperature and rainfall data from South America show that our architecture captures spatiotemporal context and that it outperforms or matches the results of state-of-the-art architectures for forecasting tasks. In particular, one of the variants of our proposed architecture is 23% better at predicting future sequences and five times faster at training than the RNN-based model used as a baseline.

Keywords Spatiotemporal data analysis · Sequence-to-Sequence models · Convolutional Neural Networks · Weather Forecasting

1 Introduction

Weather forecasting plays an essential role in resource planning in cases of severe natural phenomena such as heat waves (extreme temperatures), droughts, and hurricanes. It also influences decision-making in agriculture, aviation, retail markets, and other sectors, since unfavorable weather negatively impacts corporate revenues (Štulec et al., 2019). Over the years, with technological developments, predictions of meteorological variables are becoming more accurate. However, due to the stochastic behavior of the Earth systems, which is governed by physical laws, traditional forecasting requires complex, physics-based models to predict the weather (Karpatne et al., 2018). In recent years, an extensive volume of data about the Earth systems has become available. The remote sensing data collected by satellites provide meteorological data about the entire globe at specific time intervals (e.g., 6h or daily) and with a regular spatial resolution (e.g., 1km or 5km). The availability of historical data allows researchers to design deep learning models that can make more accurate predictions about the weather (Reichstein et al., 2019).

Even though meteorological data exhibits both spatial and temporal structures, weather forecasting can be modeled as a sequence problem. In sequence modeling tasks, an input sequence is encoded to map the representation of the sequence output, which may have a different length than the input. In Shi et al. (2015), the authors proposed the ConvLSTM architecture to solve the sequence prediction problem using a radar echo dataset for precipitation forecasting. They integrated the convolution operator, adopted by the convolutional neural network (CNN), into a recurrent neural network (RNN) to simultaneously learn the spatial and temporal context of input data to predict the future sequence. Although ConvLSTM architecture has been considered the potential approach to build prediction models for geoscience data (Reichstein et al., 2019), new opportunities have emerged from recent advances in deep learning. In Wang et al. (2017, 2019), authors proposed improved versions of the long short-term memory (LSTM) unit for memorizing spatiotemporal information.

RNN-based architectures may be ideal for multi-step forecasting tasks using spatiotemporal data (Shi et al., 2015; Wang et al., 2017, 2019), due to the ability to respect the temporal order (causal constraint) and predict long sequences. However, these architectures maintain the information from previous time steps to generate the output, which consequently leads to a high training time. Taking this as motivation, we address the spatiotemporal forecasting problem by proposing a new architecture using entirely 3D CNN. CNN are an efficient method for capturing spatial context and have attained state-of-the-art results for image classification using a 2D kernel (Krizhevsky et al., 2012). In recent years, researchers expanded CNN actuation field, such as machine translation (Gehring et al., 2017) using a 1D kernel, which is useful to capture temporal patterns in a sequence. 3D CNN-based models are commonly used for video analysis and action recognition (Yuan et al., 2018; Tran et al., 2018) or climate event detection (Racah et al., 2017). However, CNN-based models are generally not considered for multi-step forecasting tasks, because of two intrinsic limitations. They violate the temporal order, allowing future information during temporal reasoning (Singh and Cuzzolin, 2019), and they cannot generate a predictive output sequence longer than the input sequence (Bai et al., 2018). To tackle these limitations, we introduce STConvS2S (*Spatiotemporal Convolutional Sequence to Sequence Network*), a spatiotemporal predictive model for multi-step forecasting task. To our knowledge, STConvS2S is the first 3D CNN-based architecture built as an end-to-end trainable model, suitable to satisfy the causal constraint and predict flexible length output sequences (i.e., not limited to be equal to the input sequence length).

We compared STConvS2S to RNN-based architectures through experimental studies in terms of both predictive performance and time efficiency. The proposed architecture matches or outperforms state-of-the-art methods on meteorological datasets obtained from satellites and in-situ stations - CHIRPS (Funk et al., 2015), and climate model - CFSR (Saha et al., 2014).

The contributions of this paper are twofold. Firstly, we provide two variants of the STConvS2S architecture that satisfy the causal constraint. One adapts the causal convolution in 3D convolutional layers, and the other introduces a new approach that strategically applies a reverse function in the sequence. Secondly, we devise a temporal generator block designed to extend the length of the output sequence, which encompasses a new application of the transposed convolutional layers.

The rest of this paper is organized as follows. Section 2 discusses works related both to weather forecasting and spatiotemporal architectures. Section 3 presents the formulation of the spatiotemporal data forecasting problem. Section 4 describes our proposed deep learning architecture. Section 5 presents our experiments and results. Section 6 provides the conclusions of the paper.

2 Related work

Several statistical methods and machine learning techniques have been applied to historical data about temperature, precipitation, and other meteorological variables to predict the weather conditions. Auto-regressive integrated moving average (ARIMA) are traditional statistical methods for times series analysis (Babu and Reddy, 2012). Other studies have also applied artificial neural networks (ANN) to time series prediction in weather data, such as temperature measurements (Corchado and Fyfe, 1999; Baboo and Shereef, 2010; Mehdizadeh, 2018). Recently, some authors have been developing new approaches based on deep learning to improve time series forecasting results, in particular, using LSTM networks. Traffic flow analysis (Yang et al., 2019), displacement prediction of landslide (Xu and Niu, 2018), petroleum production (Sagheer and Kotb, 2019) and sea surface temperature forecasting (Zhang et al., 2017) are some applications that successfully use LSTM architectures. However, these approaches (addressed to time series) are unable to capture spatial dependencies in the observations.

Spatiotemporal deep learning models deal with spatial and temporal contexts simultaneously. In Shi et al. (2015), the authors formulate weather forecasting as a sequence-to-sequence problem, where the input and output are 2D radar map sequences. In addition, they introduce the convolutional LSTM (ConvLSTM) architecture to build an end-to-end model for precipitation nowcasting. The proposed model includes the convolution operation into LSTM network to

capture spatial patterns. Kim et al. (2019) also define their problem as a sequence task and adopt ConvLSTM for extreme climate event forecasting. Their model uses hurricane density map sequences as spatiotemporal data. The work proposed in Souto et al. (2018) implements a spatiotemporal aware ensemble approach adopting ConvLSTM architecture. Based on Shi et al. (2015), Wang et al. (2017) present a new LSTM unit that memorizes spatial and temporal variations in a unified memory pool. In Wang et al. (2019), they present an improved memory function within LSTM unit adding non-stationarity modeling. Although related to the use of deep learning for climate/weather data, our model adopts only CNN rather than a hybrid approach that combines CNN and LSTM.

Some studies have applied spatiotemporal convolutions (Yuan et al., 2018; Tran et al., 2018) for video analysis and action recognition. In Tran et al. (2018), the authors compare several spatiotemporal architectures using only 3D CNN and show that factorizing the 3D convolutional kernel into separate and successive spatial and temporal convolutions produces accuracy gains. A limitation of both 3D CNN or factorized 3D CNN (Tran et al., 2018) is the lack of causal constraint, violating the temporal order. Singh and Cuzzolin (2019) and Cheng et al. (2019) factorize the 3D convolution as Tran et al. (2018). Singh and Cuzzolin (2019) propose a recurrent convolution unit approach to address causal constraint in temporal learning for action recognition tasks, and Cheng et al. (2019) satisfy the causal constraint by adopting causal convolution in separate and parallel spatial and temporal convolutions. We also adopt a factorized 3D CNN, but with a different implementation, where Figure 2 highlights our approach. In contrast to Singh and Cuzzolin (2019), we use an entirely CNN approach, and to Cheng et al. (2019), besides not using parallel convolutions when adopting a causal convolution, we introduce a new method to not violate the temporal order (details in Section 4.2).

Following the success of 2D CNN in capturing spatial correlation in images, Xu et al. (2019) propose a model to predict vehicle pollution emissions using 2D CNN to capture temporal and spatial correlation separately. Racah et al. (2017) use a 3D CNN in an encoder-decoder architecture for extreme climate event detection. Their architecture consists of a downsampling path in the encoder using a stack of convolutional layers, and an upsampling path in the decoder using a stack of transposed convolutional layers. Their model adopts the typical use of transposed convolutional layers to reconstruct the output to match the entire input dimension. Instead, we use these layers to generate an output with a larger dimension, different from the dimensions of the input. Furthermore, unlike our work, they do not satisfy the causal constraint in their models.

3 Problem Statement

Spatiotemporal data forecasting can be modeled as a sequence-to-sequence problem. Thus, the observations of spatiotemporal data (e.g. meteorological variables) measured in a specific geographic region over a period of time serve as the input sequence to the forecasting task. More formally, we define a spatiotemporal dataset as $[\tilde{X}^{(1)}, \tilde{X}^{(2)}, \dots, \tilde{X}^{(m)}]$ with m samples of $\tilde{X}^{(i)} \in \mathbb{R}^{T \times H \times W \times C}$, where $1 \leq i \leq m$. Each training example is a tensor $\tilde{X}^{(i)} = [X_1^{(i)}, X_2^{(i)}, \dots, X_T^{(i)}]$, that is a sequence of T observations containing historical measurements. Each observation $X_j^{(i)} \in \mathbb{R}^{H \times W \times C}$, for $j = 1, 2, \dots, T$ (i.e. the length of input sequence), consists of a $H \times W$ grid map that determines the spatial location of the measurements, where H and W represent the latitude and longitude, respectively. In the observations, C represents how many meteorological variables (e.g. temperature, humidity) are used simultaneously in the model. This structure is analogous to 2D images, where C would indicate the amount of color components (RGB or grayscale).

Modeled as sequence-to-sequence problem in Equation 1, the goal of spatiotemporal data forecasting is to apply a function f that maps an input sequence of past observations, satisfying the causal constraint at each time step t , in order to predict a target sequence $\hat{X} \in \mathbb{R}^{H \times W \times C}$, where the length T'' of output sequence may differ from the length T of input sequence.

$$\hat{X}_{t+1}, \hat{X}_{t+2}, \dots, \hat{X}_{t+T''} = f(X_{t-T+1}, \dots, X_{t-1}, X_t) \quad (1)$$

4 STConvS2S architecture

STConvS2S is an end-to-end deep neural network suited for learning spatiotemporal predictive patterns, which are common in domains weather forecasting. Our approach makes multi-step (sequences) prediction without feeding the predicted output back into the input sequence. Figure 1 is an overview of our proposed deep learning architecture.

Although some methods for weather forecasting using a radar echo dataset apply a hybrid approach, combining 2D CNN (to learn spatial representations) and LSTM (to learn temporal representations) (Shi et al., 2015; Wang et al., 2017, 2019), our method uses only 3D convolutional layers to learn spatial and temporal contexts. Distinct from conventional convolution applied in some 3D CNN architectures (Tran et al., 2015; Tran et al., 2018; Racah et al., 2017), during

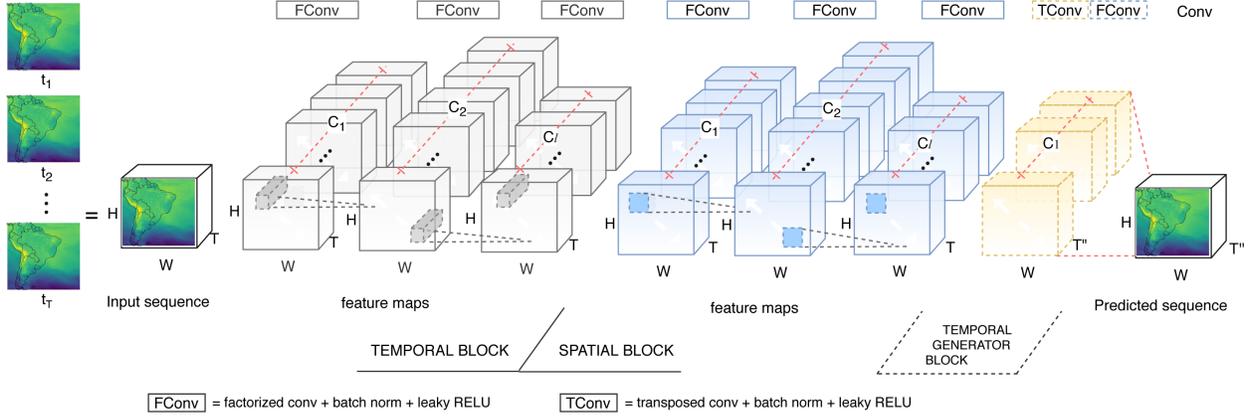


Figure 1: An illustration of STConvS2S architecture, which comprises three components: temporal block, spatial block, and temporal generator block. Each block is a set of layers. The temporal block learns a temporal representation of the input sequence, the spatial block extracts spatial features from the output of the previous block. On top of the spatial block, there is the temporal generator block designed to increase the sequence length T if the task requires a longer predictive horizon, where $T' \geq T$. Finally, the output of this block is further fed into a final convolutional layer to complete the prediction

temporal learning, STConvS2S takes care not to depend on future information, a crucial constraint on forecasting tasks. Another core feature of our designed network is the capability to allow flexible output sequence length, which means the possibility to predict many time-steps ahead, regardless of the fixed-length of the input sequence. In the following, we provide more details about the components which comprise our architecture.

4.1 Factorized 3D convolutions

Instead of adopting a conventional $t \times d \times d$ kernel for 3D convolutional layers, where d and t are the kernel size in space ($H \times W$) and time (T) dimensions, respectively, we use a factorized 3D kernel adapted from R(2+1)D network, proposed in Tran et al. (2018). The factorized kernel $1 \times d \times d$ and $t \times 1 \times 1$ split the convolution operation of one layer into two successive operations, named as a spatial convolution and a temporal convolution in their work. In our new architecture, we take a different approach: operations are not successive inside each convolutional layer. Instead, the factorized kernels are separated into two blocks, giving them specific learning skills. The temporal block applies the $t \times 1 \times 1$ kernel in its layers to learn only temporal dependencies, while the next component, the spatial block, encapsulates spatial dependencies using $1 \times d \times d$ kernel. Figure 2 schematically illustrates the difference between these three approaches.

Compared to the full 3D kernel applied in standard convolutions, the kernel decomposition used in STConvS2S offers the advantage of increasing the number of nonlinearities in the network (additional activation functions between factorized convolutions), which leads to an increase in the complexity of representable patterns (Tran et al., 2018). An advantage of our proposed approach over the (2+1)D block is flexibility since temporal and spatial blocks can have a distinct number of layers, facilitating their optimization.

4.2 Temporal Block

In STConvS2S, the temporal block is a stack of 3D convolutional layers which adopt $t \times 1 \times 1$ kernel during convolutions. Each layer receives a 4D tensor with dimensions $T \times H \times W \times C_{l-1}$ as input, where C_{l-1} is the number of filters used in the previous layer ($l-1$), T is the sequence length (time dimension), H and W represent the size of the spatial coverage for latitude and longitude, respectively. Within the block, the filters C are increased twice in the feature maps as the number of layers increases, but the final layer reduces them again to the number of filters initially defined.

In detail, this block uses batch normalization and leaky rectified linear unit (LeakyReLU) with a negative slope set as 0.01 after each convolutional layer. This block discovers patterns over the time dimension T exclusively. Besides, since we are using 3D convolutional layers to analyze historical series of events, we must prevent data leakage from happening. That is, the model should not violate the temporal order and should ensure that, at step t , the learning process uses no future information from step $t+1$ onward. To satisfy this constraint, we propose two variants of the temporal block.

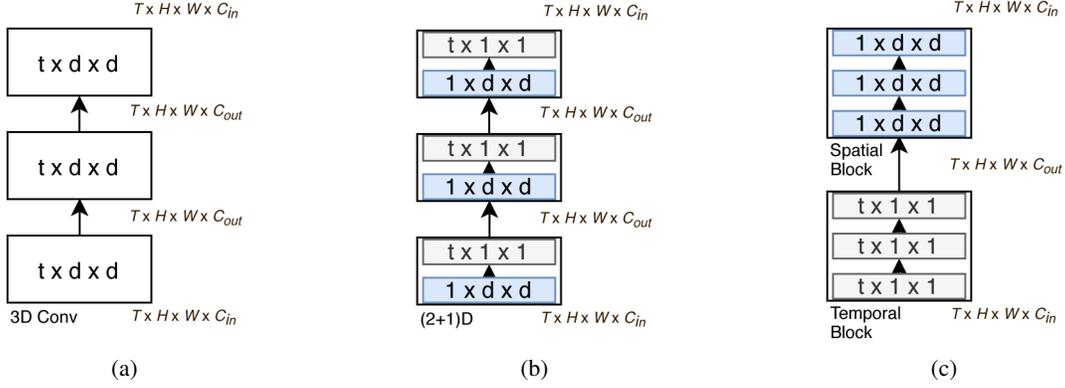


Figure 2: Comparison of convolution operations applied in three convolutional layers. The spatial kernel is defined as $1 \times d \times d$ and the temporal kernel as $t \times 1 \times 1$, where d and t are the kernel size in spatial ($H \times W$) and time (T) dimensions, respectively. (a) Representation of the standard 3D convolution operation using the $t \times d \times d$ kernel. (b) Factorized 3D kernels proposed in Tran et al. (2018) as successive spatial and temporal convolution operations in a unique block called (2+1)D. (c) Our proposal for the factorized 3D kernels usage is in separate blocks. First, the temporal block stacks three convolutional layers, each performing convolutions using only the temporal kernel. Likewise, the spatial block applies the spatial kernel to its layers.

Temporal Causal Block. We name our architecture as *STConvS2S-C* when it adopts this block to learn the temporal patterns. We apply causal convolutions within the block to incorporate the ability to respect the temporal order during learning in convolutional layers. Causal convolution was originally presented in WaveNet (van den Oord et al., 2016) for 1D CNN and applied with factorized 3D convolutions in Cheng et al. (2019). This technique can be implemented by padding the input by $k - 1$ elements, where k is the kernel size. Figure 3 shows the operation in details.

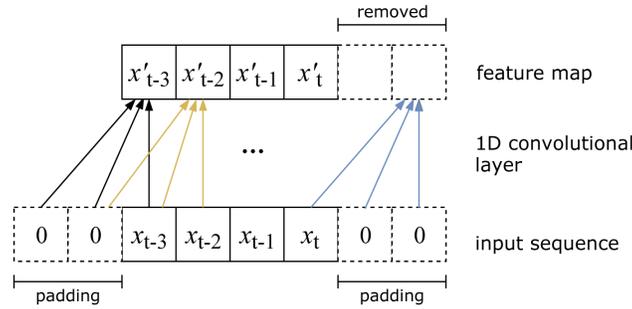


Figure 3: Causal convolution operation in a 1D convolutional layer (used to simplify the illustration) with $k = 3$ (kernel size). Input is padded by $k - 1$ elements to avoid learning future information. To ensure that the output feature map has the same length as the input, the last $k - 1$ elements are removed since they are related to the zeros added to the right of the input.

Temporal Reversed Block. When dealing with historical data, respecting the temporal order (causal constraint) is an essential behavior of deep learning models. This because, in real applications, future information is not available in forecasting. The common approach in the literature for adapting convolutional layers to satisfy this constraint is through causal convolutions. Here, we introduce a better alternative to avoid violating the temporal order, applying a function ψ in the time dimension to reverse the sequence order. This function is a linear transformation $\psi : \mathbb{R}^{T \times H \times W \times C} \rightarrow \mathbb{R}^{T \times H \times W \times C}$. The architecture is named as *STConvS2S-R* when composed with this block. Formally, *STConvS2S-R* computes the output feature map R of a temporal reversed block using

$$R'_{1:l_r} = \begin{cases} g(W_u * \psi(I_u) + b_u), & \text{if } u = 1 \\ g(W_u * I_u + b_u), & \text{if } 2 \leq u \leq l_r \end{cases} \quad (2)$$

$$R = \psi(R'_{l_r}) \quad (3)$$

where $W_{1:l_r}$ and $b_{1:l_r}$ is the learnable weight tensor and bias term in l_r layers of this block, $*$ denotes a convolution operator and $g(\cdot)$ is a non-linear activation function. For the first layer of the temporal reversed block, I_1 is the input sequence \tilde{X} previously defined in Section 3, and for the subsequent layers, $I_{2:l_r}$ is the feature map calculated in the previous layer R'_{l_r-1} .

4.3 Spatial Block

The spatial block is built on top of the temporal block and has a similar structure with batch normalization and LeakyReLU as non-linearity. In contrast, each 3D convolutional layer of this block extracts only spatial representations since kernel decomposition allows us to analyze the spatial and temporal contexts separately. In the STConvS2S, each feature map generated has a fixed-length in $H \times W$ dimensions and, to ensure this, the input in the spatial block is padded following $p = \frac{k_s-1}{2}$, where k_s is the size of spatial kernel. This design choice differentiates our model from 3D encoder-decoder architecture (Racah et al., 2017), which needs to stack upsample layers after all convolutional layers, due to the downsampling done in the latter.

4.4 Temporal Generator Block

In addition to ensuring that our model satisfies the causal constraint, another contribution of our work is generating output sequences with longer lengths than the length of the input sequence. When CNNs are used for sequence-to-sequence learning, such as multi-step forecasting, the length of the output sequence must be the same size or shorter than the input sequence (Gehring et al., 2017; Bai et al., 2018). To tackle this limitation, we designed a component placed on top of the spatial block used when the task requires a more extended sequence (e.g., from the previous 5 grids, predict the next 15 grids). First, we compute the intermediate feature map G' :

$$G'_{1:l_{g_t}} = tconv(I_{1:l_{g_t}}) \quad (4)$$

where $l_{g_t} = \left\lceil \frac{T''-T}{2T} \right\rceil$ is the number of transposed convolutional layers ($tconv$) necessary to guarantee that G' has the size of time dimension $T_{G'} \geq T'' - T$. The kernel size, stride and padding of $tconv$ are fixed and extend the feature map by a factor of 2 in time dimension only. For the first layer, I_1 is the output of the spatial block S and for the other layers, $I_{2:l_{g_t}}$ is the feature map calculated in the previous layer $G'_{l_{g_t}-1}$. Follow, given G' and S , we can compute G'' :

$$G'' = \rho(S \oplus G') \quad (5)$$

In the equation above, \oplus denotes a concatenation operator in the time dimension and $\rho(\cdot)$ is a function to ensure that the feature map G'' matches exactly the length T'' of the desired output sequence. Finally, the output feature map G of this block can be defined as

$$G_{1:l_{g_c}} = g(W_{1:l_{g_c}} * I_{1:l_{g_c}} + b_{1:l_{g_c}}) \quad (6)$$

where $l_{g_c} = \left\lceil \frac{T''}{T} \right\rceil$ is the number of convolutional layers that use factorized kernels as in the spatial block. $W_{1:l_{g_c}}$ and $b_{1:l_{g_c}}$ is the learnable weight tensor and bias term in l_{g_c} layers, $*$ denotes a convolution operator and $g(\cdot)$ is a non-linear activation function. For the first convolutional layer, I_1 is G'' . Unlike the temporal and spatial block, where the number of layers is a defined hyperparameter to execute the model, in the temporal generator block l_{g_t} and l_{g_c} are calculated based on the length T'' of the desired output sequence and the length T of the input sequence (the size of time dimension).

5 Experiments

We perform experiments on two publicly available meteorological datasets containing air temperature and precipitation values to validate our proposed architecture. The deep learning experiments were conducted on a server with a single Nvidia GeForce GTX1080Ti GPU with 11GB memory. We executed the ARIMA methods on 8 Intel i7 CPUs with 4 cores and 66GB RAM. We start by explaining the datasets (Section 5.1) and evaluation metrics (Section 5.2). Further, we describe the main results of the experiments for each dataset (Section 5.3 and 5.4) and summarize the results of ablation studies (Section 5.5).

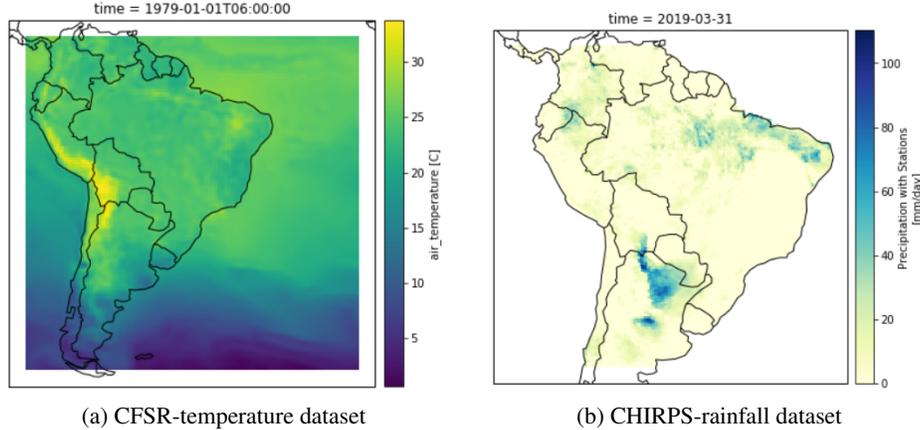


Figure 4: Spatial coverage of the datasets used in all experiments. (a) It shows the selected grid on January 1, 1979 with air temperature values. (b) It shows the selected grid of the sequence on March 31, 2019 with rainfall values.

5.1 Datasets

The CFSR¹ is a reanalysis² product that contains high-resolution global land and ocean data (Saha et al., 2014). The data contain a spatial coordinate (latitude and longitude), a spatial resolution of 0.5 degrees (i.e., $0.5^\circ \times 0.5^\circ$ area for each grid cell) and a frequency of 6 hours for some meteorological variables, such as air temperature and wind speed.

In the experiments, we use a subset of CFSR with the air temperature observations from January 1979 to December 2015, covering the space in 8°N - 54°S and 80°W - 25°W as shown in Figure 4 (a). As data preprocessing, we scale down the grid to 32×32 in the H and W dimensions to fit the data in GPU memory. The other dataset, CHIRPS³, incorporates satellite imagery and in-situ station data to create gridded rainfall times series with daily frequency and spatial resolution of 0.05 degrees (Funk et al., 2015). We use a subset with observations from January 1981 to March 2019 and apply interpolation to reduce the grid size to 50×50 . Figure 4 (b) illustrates the coverage space 10°N - 39°S and 84°W - 35°W adopted in our experiments.

Similar to Shi et al. (2015), we define the input sequence length as 5, which indicates the use of the previous five grids to predict the next T'' grids. Thus, the input data shapes to the deep learning architectures are $5 \times 32 \times 32 \times 1$ for CFSR dataset and $5 \times 50 \times 50 \times 1$ for CHIRPS dataset. The value 1 in both dataset shapes indicates the one-channel (in this aspect similar to a grayscale image), 5 is the size of the sequence considered in the forecasting task, and 32 and 50 represent the numbers of latitudes and longitudes used to build the spatial grid in each dataset.

We create 54,041 and 13,960 grid sequences from the temperature dataset and rainfall datasets, respectively. Finally, we divide both datasets into non-overlapping training, validation, and test set following 60%, 20%, and 20% ratio, in this order. The adoption of temperature and rainfall datasets in our experimental evaluation relies on the fact that they are the two main meteorological variables. Research about their spatiotemporal representation is relevant to short-term forecasting and improves the understanding of long-term climate variability (Rahman and Lateh, 2017). However, the proposed architecture is suitable for other meteorological variables or other domains, as long as the training data can be structured as defined in Section 3.

5.2 Evaluation metrics

To evaluate the proposed architecture, we compare our results against ARIMA models, traditional statistical approaches for time series forecasting, and state-of-the-art models for spatiotemporal forecasting. To accomplish this, we use the two evaluation metrics presented in Equation 7 and 8.

¹<https://climatedataguide.ucar.edu/climate-data/climate-forecast-system-reanalysis-cfsr>

²Scientific method used to produce best estimates (analyses) of how the weather is changing over time (Fujiwara et al., 2017).

³<https://chc.ucsb.edu/data/chirps>

RMSE, denoted as E_r , is based on MSE metric, which is the average of squared differences between real observation and prediction. The MSE square root gives the results in the original unit of the output, and is expressed at a specific spatiotemporal volume as:

$$E_r(T, H, W) = \sqrt{\frac{1}{N} \sum_{n=1}^N \sum_{t \in T} \sum_{h \in H} \sum_{w \in W} [x(t, h, w) - \hat{x}(t, h, w)]^2} \quad (7)$$

where N is the number of test samples, $x(t, h, w)$ and $\hat{x}(t, h, w)$ are the real and predicted values at the location h and w at time t , respectively.

MAE, denoted as E_m , is the average of differences between real observation and prediction, which measures the magnitude of the errors in prediction. MAE also provides the result in the original unit of the output, and is expressed at a specific spatiotemporal volume as:

$$E_m(T, H, W) = \frac{1}{N} \sum_{n=1}^N \sum_{t \in T} \sum_{h \in H} \sum_{w \in W} |x(t, h, w) - \hat{x}(t, h, w)| \quad (8)$$

where N, t, h, w are defined as shown in Equation 7.

5.3 CFSR Dataset: results and analysis

We first conduct experiments with distinct numbers of layers, filters, and kernel sizes to investigate the best hyperparameters to fit the deep learning models. As a starting point, we set the version 1 based on the settings described in Shi et al. (2015) with two layers, each containing 64 filters and a kernel size of 3^4 . To make fair comparisons using the chosen datasets, we explored variations of the hyperparameters for our architecture (STConvS2S-C and STConvS2S-R) and the following state-of-the-art methods: ConvLSTM (Shi et al., 2015), PredRNN (Wang et al., 2017), and MIM (Wang et al., 2019). Thus, for versions 2-4, we defined the number of layers (L), kernel size (K) and the number of filters (F) in a way that would help us understand the behavior of the models during the learning process by increasing L (versions 1 and 3), K (versions 2 and 4) or F (versions 2 and 3). In the training phase, we perform for all models mini-batch learning with 50 epochs, and RMSprop optimizer with a learning rate of 10^{-3} .

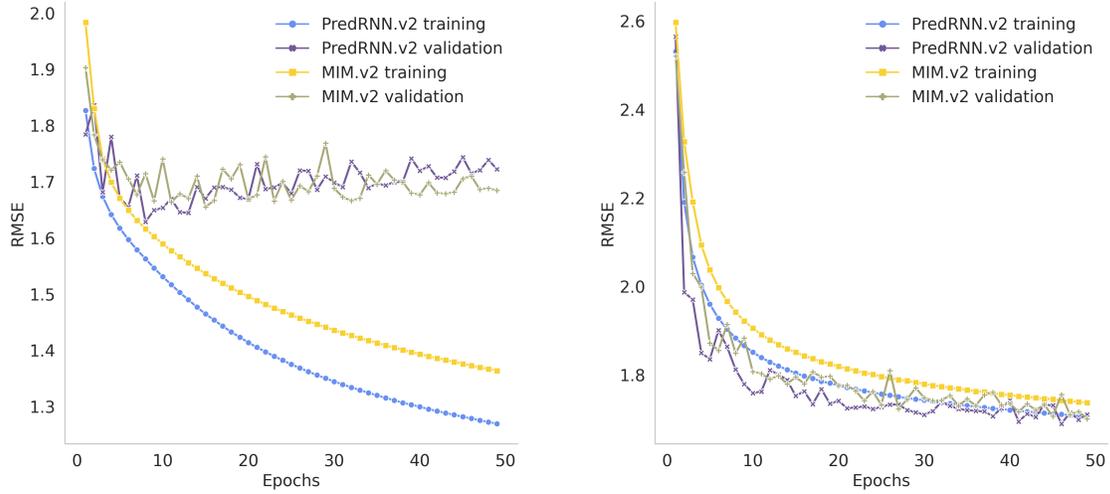
We applied dropout after convolutional layers during the training of PredRNN and MIM models⁵ to reduce the model complexity and avoid overfitting. Without dropout, these models do not generalize well for this dataset and make less accurate predictions in the validation set. We adopt 0.5 as the dropout rate for both models after evaluating the best rate employing the grid search technique, which performed several experiments changing the rate by $\{0.3, 0.5, 0.8\}$. Figure 5 (a) and (b) illustrate the differences in the learning curve, where the former shows a high error on the validation set early in the training stage for both models, and the latter indicates the learning curve with dropout applied.

As a sequence-to-sequence task, we use the previous five grids, as we established before in Section 5.1, to predict the next five grids (denoted as $5 \rightarrow 5$). Table 1 provides the models considered in our investigation with four different settings, the values of the RMSE metric on the test set, the training time, and the memory usage by GPU. The results present the superiority of version 4, which has the highest values of L and K , reaching the lowest RMSE for all models, except PredRNN, where version 2 is superior. Another aspect to note is that when increasing the number of filters (versions 2 and 3) the impact is more significant in the training time than when increasing the number of layers (versions 1 and 3) for state-of-the-art models, indicating that version 2 is faster for these RNN-based architectures. This impact is not seen in our models (CNN-based) compared to version 3, as in versions 1 and 2 they spend almost the same time during training, showing that STConvS2S models have more stability when increasing the hyperparameters.

To improve the comprehension of the analysis, Figure 6 highlights the differences between the performances of RMSE metric and training time for the models. As shown, STConvS2S-R and STConvS2S-C models perform favorably against the state-of-the-art models for the CFSR dataset in all versions, demonstrating that our architectures can simultaneously capture spatial and temporal correlations. Comparing the best version of each model, our models significantly outperform the state-of-the-art architectures for spatiotemporal forecasting. In detail, STConvS2S-R (version 4) takes only 1/4 memory space, is 5x faster in training, and has achieved a 23% improvement in RMSE over MIM (version 4), the RNN-based model with better performance. These results reinforce that our models have fewer parameters to optimize than MIM and PredRNN. Furthermore, our model can be completely parallelized, speeding up the learning process, as the output of the convolutional layers does not depend on the calculations of the previous step,

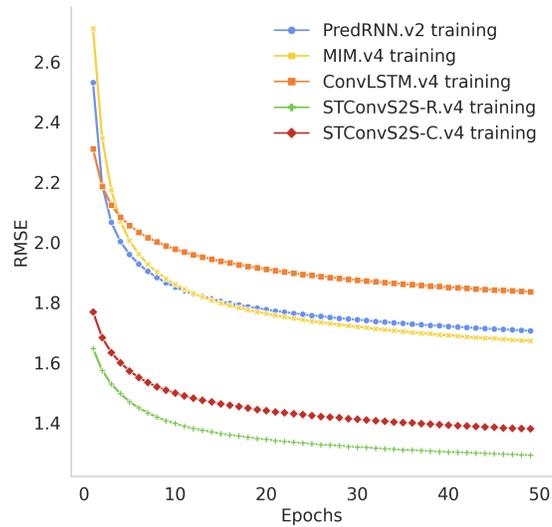
⁴ 3×3 kernel for ConvLSTM, PredRNN, and MIM. $3 \times 1 \times 1$ temporal kernel and $1 \times 3 \times 3$ spatial kernel for STConvS2S.

⁵STConvS2S and ConvLSTM models do not overfit during training on any version.



(a) Learning curve - overfitting

(b) Learning curve - dropout 0.5



(c) Training curve for all models

Figure 5: Learning curves after running 50 epochs on temperature dataset (CFSR). (a) To exemplify, we select version 2 to illustrate the overfitting observed when analyzing the training and validation curve of PredRNN and MIM models. (b) The same version and models using dropout to improve its generalization. (c) Comparison of training curve for the best version for each model. Our models (STConvS2S-R and STConvS2S-C) achieved lower RMSE and thus, better ability to learn spatiotemporal representations.

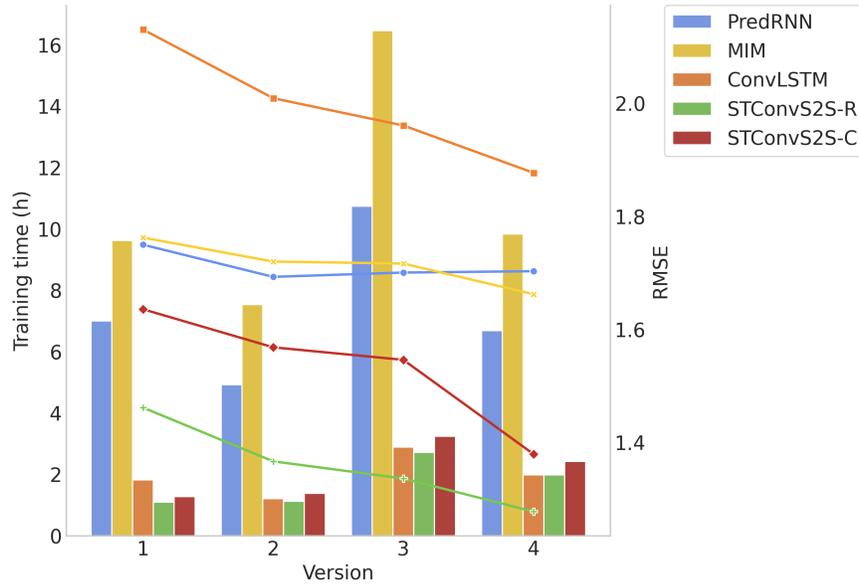


Figure 6: Comparison between training time in hours (bar plot) and RMSE (line plot) for each model version on temperature dataset (CFSR).

Table 1: Evaluation of different settings on the CFSR dataset for STConvS2S and state-of-the-art methods, where the best version has the lowest RMSE value.

Model	Version	Setting	5 → 5		
			RMSE	Training time	Memory usage (MB)
ConvLSTM (Shi et al., 2015)	1	L=2, K=3, F=64	2.1306	01:49:16	1119
	2	L=3, K=3, F=32	2.0090	01:12:16	920
	3	L=3, K=3, F=64	1.9607	02:53:00	1358
	4	L=3, K=5, F=32	1.8770	01:58:52	922
PredRNN (Wang et al., 2017)	1	L=2, K=3, F=64	1.7497	06:59:45	3696
	2	L=3, K=3, F=32	1.6928	04:55:19	2880
	3	L=3, K=3, F=64	1.7004	10:44:25	5242
	4	L=3, K=5, F=32	1.7028	06:41:14	2892
MIM (Wang et al., 2019)	1	L=2, K=3, F=64	1.7623	09:37:07	4826
	2	L=3, K=3, F=32	1.7199	07:31:59	4124
	3	L=3, K=3, F=64	1.7163	16:27:42	7789
	4	L=3, K=5, F=32	1.6621	09:49:40	4145
STConvS2S-C (ours)	1	L=2, K=3, F=64	1.6355	01:16:40	991
	2	L=3, K=3, F=32	1.5681	01:22:42	1021
	3	L=3, K=3, F=64	1.5459	03:14:25	1554
	4	L=3, K=5, F=32	1.3791	02:25:26	1040
STConvS2S-R (ours)	1	L=2, K=3, F=64	1.4614	01:05:48	880
	2	L=3, K=3, F=32	1.3663	01:07:33	891
	3	L=3, K=3, F=64	1.3359	02:42:53	1283
	4	L=3, K=5, F=32	1.2773	01:58:39	895

Table 2: Performance results for temperature forecasting using the previous five observations (grids) to predict the next five observations ($5 \rightarrow 5$), and the next 15 observations ($5 \rightarrow 15$).

Model	$5 \rightarrow 5$				
	RMSE	MAE	Memory usage (MB)	Mean training time	Training time/epoch
ARIMA	2.1880	1.9005	—	—	—
ConvLSTM (Shi et al., 2015)	1.8555 ± 0.0033	1.2843 ± 0.0028	922	02:38:27	00:02:21
PredRNN (Wang et al., 2017)	1.6962 ± 0.0038	1.1885 ± 0.0020	2880	06:59:34	00:05:52
MIM (Wang et al., 2019)	1.6731 ± 0.0099	1.1790 ± 0.0055	4145	11:05:37	00:10:43
STConvS2S-C (ours)	1.3699 ± 0.0024	0.9434 ± 0.0020	1040	03:34:52	00:02:48
STConvS2S-R (ours)	1.2692 ± 0.0031	0.8552 ± 0.0018	895	03:15:12	00:02:13
Model	$5 \rightarrow 15$				
	RMSE	MAE	Memory usage (MB)	Mean training time	Training time/epoch
ARIMA	2.2481	1.9077	—	—	—
ConvLSTM (Shi et al., 2015)	2.0728 ± 0.0069	1.4558 ± 0.0076	1810	5:29:30	00:07:32
PredRNN (Wang et al., 2017)	2.0237 ± 0.0067	1.4311 ± 0.0149	7415	11:45:48	00:17:03
MIM (Wang et al., 2019)	2.0287 ± 0.0361	1.4330 ± 0.0250	10673	19:19:00	00:31:19
STConvS2S-C (ours)	1.8739 ± 0.0107	1.2946 ± 0.0061	1457	03:12:24	00:05:17
STConvS2S-R (ours)	1.8051 ± 0.0040	1.2404 ± 0.0068	1312	03:15:42	00:05:03

as occurs in recurrent architectures. Figure 5 (c) illustrates that STConvS2S-R has a lower training error in 50 epochs compared to other models, including STConvS2S-C, proving to be a better alternative to make CNN-based models respect the temporal order.

To further evaluate our models, we chose the most efficient version for each model to perform new experiments. For STConvS2S-R, STConvS2S-C, ConvLSTM and MIM models, version 4 was chosen with 3 layers, 32 filters, and kernel size of 5, and for PredRNN, version 2 with the same number of layers and filters, but with kernel size of 3. We also included a comparison with ARIMA methods to serve as baseline for deep learning models, since they are a traditional approach to time series forecasting. The experiment for the baseline takes into account the same temporal pattern and spatial coverage. Thus, predictions were performed throughout all the 1,024 time series, considering in each analysis the previous 5 values in the sequence.

In this phase, we have not defined a specific number of epochs for each deep learning model’s execution. Therefore to avoid overfitting during the training of models, we apply the early stopping technique with patience hyperparameter set to 16 on the validation dataset. As the models run with different numbers of epochs, we include the training time/epoch to be able to compare the time efficiency of the models. We train and evaluate each deep learning model 3 times and compute the mean and the standard deviation of RMSE and MAE metrics on the test set. This time, we evaluate the models in two horizons: 5-steps ($5 \rightarrow 5$) and 15-steps ($5 \rightarrow 15$). These experiments are relevant to test the capability of our model to predict a long sequence.

As shown in Table 2, STConvS2S-R and STConvS2S-C perform much better than the baseline on RMSE and MAE metrics, indicating the importance of spatial dependence on geoscience data since ARIMA models only analyze temporal relationships. They also outperform the state-of-the-art models in these evaluation metrics in both horizons, demonstrating that our models can be efficiently adopted to predict future observations. However, beyond that, the designed temporal generator block in STConvS2S architecture can convincingly generate a more extended sequence regardless of the fixed-input sequence length. In a closer look at the best CNN-based architecture and RNN-based architecture in the task $5 \rightarrow 15$, STConvS2S-R takes less memory space and is faster than PredRNN. To provide an overview, Figure 7 illustrates the cumulative error based on both horizons.

5.4 CHIRPS Dataset: results and analysis

Similar to what we did with CFSR dataset, we divide the experiments into two phases. The first phase aims to investigate the best hyperparameter settings to adjust the models. In the second, we take into account the best version of each model and perform experiments with different initialization values for weight and bias to consolidate the analysis. In detail, we first set the hyperparameters as in the previous experiments on CFSR dataset. However, as the CHIRPS dataset is

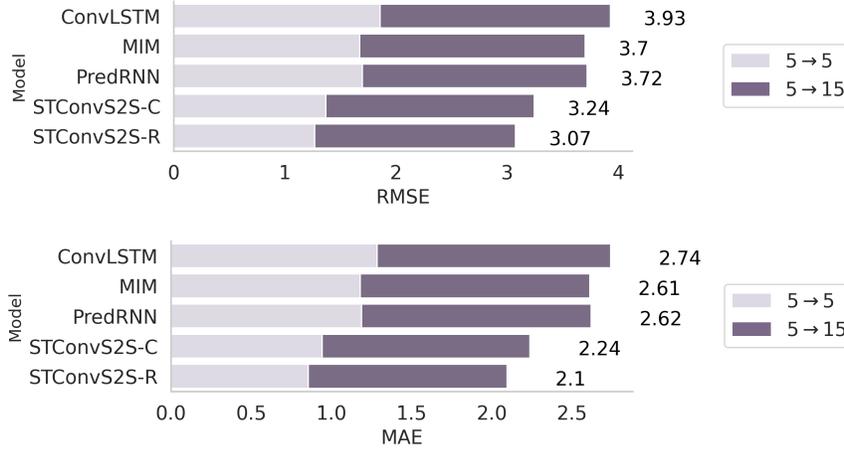


Figure 7: Cumulative error based on both horizons ($5 \rightarrow 5$ and $5 \rightarrow 15$) using temperature dataset (CFSR). Evaluations on RMSE and MAE metrics.

almost 4x smaller, all models overfit with those configurations. Thus, as an initial method to address this problem, we reduce the models' complexity by decreasing the number of layers (L) and the number of filters (F). However, we ensure fair comparability in the way we analyze the learning process when changing L (versions 1 and 3), K (versions 2 and 4) or F (versions 2 and 3). Again, all models were trained mini-batch learning with 50 epochs, and RMSprop optimizer with a learning rate of 10^{-3} .

Although we reduced the complexity, the overfitting problem remained with PredRNN and MIM models. We apply dropout to improve their performance in comparison with our proposed models. As before, we apply a search to find the best dropout rate among 0.2, 0.5, 0.8. Figure 8 (a) and (b) show the learning curves of these models with overfitting and with a dropout rate of 0.5 applied, respectively. Table 3 shows the experimental results of predicting five grids into the future by observing five grids ($5 \rightarrow 5$). For all models, version 1 with the fewest layers has the lowest memory usage and was faster in training than the other versions. Another notable analysis is that version 2 and 4 consume the GPU memory equally, except for the STConvS2S-C model, thus increasing the kernel size affects only computational time.

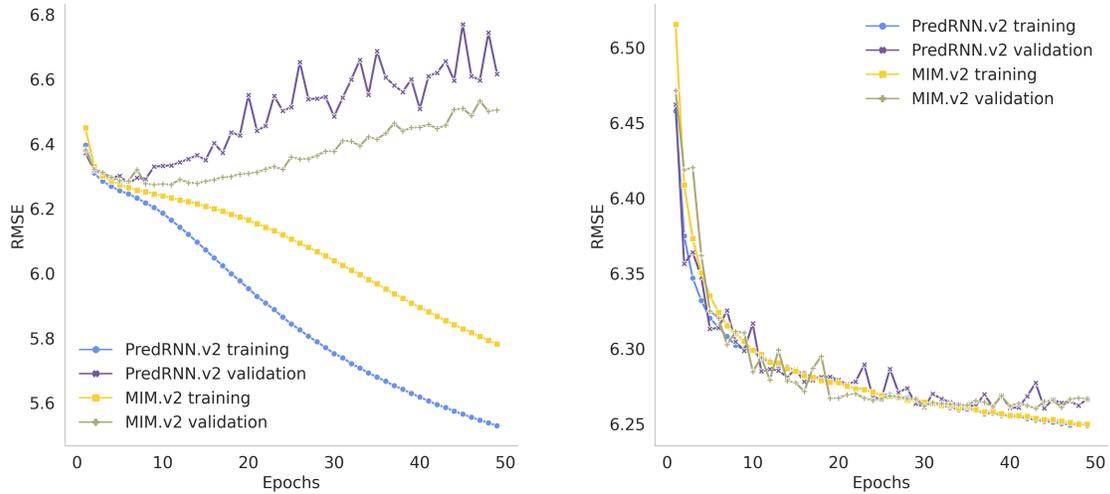
Figure 9 compares these results version by version. STConvS2S models outperform ConvLSTM with compatible training time on all versions. Comparing the best version of each model (version 4 for STConvS2S models and ConvLSTM, and version 2 for PredRNN and MIM), STConvS2S-R has the lowest prediction error and, compared to PredRNN, it is 3x faster and occupies only 1/3 of memory space. Besides, Figure 8 (c) illustrates training in 50 epochs and indicates that STConvS2S-R learns the spatiotemporal representation of rainfall better than RNN-based architectures.

For the second phase, we train the models in the same set up as previously indicated for the CFSR dataset. We also include ARIMA methods as a baseline and evaluate the proposed architectures and state-of-the-art models in two tasks: feeding only five observations (grids) into the network and predicting the next 5 and 15 observations, denoted as $5 \rightarrow 5$ and $5 \rightarrow 15$, respectively. For ARIMA, predictions were performed throughout all the 2,500 time series, considering the previous five values in the sequence in each analysis. Results of Table 4 demonstrate that STConvS2S-R achieves a better trade-off between computational cost and prediction accuracy than state-of-the-art models in both tasks.

Figure 10 summarizes these results in an overview of the cumulative error based on the two forecast horizons. Trained on the rainfall dataset, our proposed architecture equipped with the temporal reversed block achieves performance comparable to RNN-based architecture. Besides, it can predict short and even long sequences in the spatiotemporal context. Such a statement is confirmed by Figure 11, which shows the observations at each time step for STConvS2S-R and PredRNN models. STConvS2S-R can predict in the long-term without many distortions since it presents a predictive result similar to the PredRNN. Given the high variability of rainfall, both models have difficulties in making an accurate forecast.

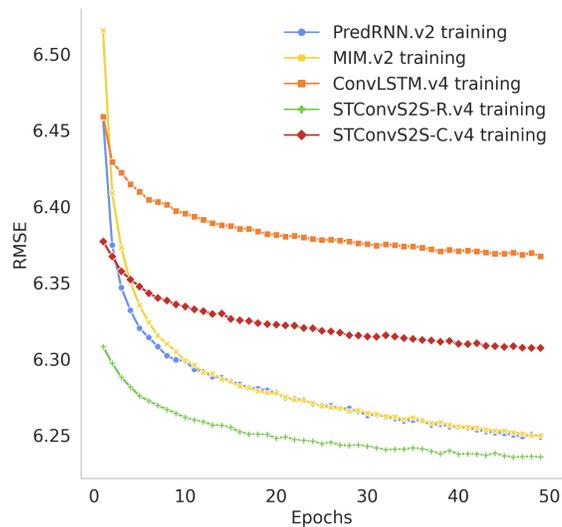
5.5 Ablation Study

We conduct a series of ablation studies, where the goal is to understand our architecture by removing/changing its main components and observing the impact on the evaluation metrics. For fair comparisons, we trained all models with the



(a) Learning curve - overfitting

(b) Learning curve - dropout 0.5



(c) Training curve for all models

Figure 8: Learning curves after running 50 epochs on the rainfall dataset (CHIRPS). (a) To exemplify, we select version 2 to illustrate the overfitting observed when analyzing the PredRNN and MIM models' training and validation curves. (b) The same version and models using dropout to improve its generalization. (c) Comparison of training curve for the best version for each model. STConvS2S-R achieved lower RMSE and, thus, better ability to learn spatiotemporal representations.

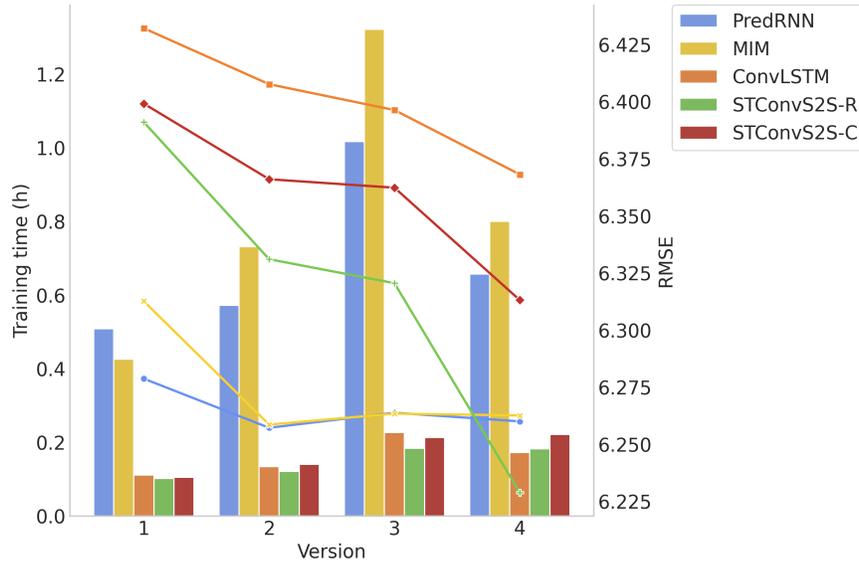


Figure 9: Comparison between training time in hours (bar plot) and RMSE (line plot) for each model version on the rainfall dataset (CHIRPS).

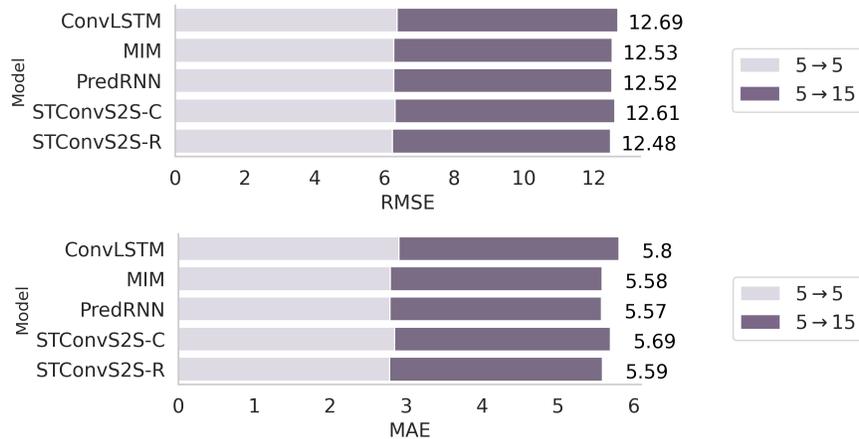
Table 3: Evaluation of different settings on the CHIRPS dataset for STConvS2S and state-of-the-art methods, where the best version has the lowest RMSE value.

Model	Version	Setting	5 → 5		
			RMSE	Training time	Memory usage (MB)
ConvLSTM (Shi et al., 2015)	1	L=1, K=3, F=16	6.4321	00:06:39	746
	2	L=2, K=3, F=8	6.4076	00:08:01	756
	3	L=2, K=3, F=16	6.3963	00:13:35	989
	4	L=2, K=5, F=8	6.3681	00:10:19	756
PredRNN (Wang et al., 2017)	1	L=1, K=3, F=16	6.2787	00:30:31	1673
	2	L=2, K=3, F=8	6.2572	00:34:19	1740
	3	L=2, K=3, F=16	6.2638	01:01:00	2775
	4	L=2, K=5, F=8	6.2600	00:39:24	1740
MIM (Wang et al., 2019)	1	L=1, K=3, F=16	6.3126	00:25:32	1447
	2	L=2, K=3, F=8	6.2586	00:43:52	2231
	3	L=2, K=3, F=16	6.2634	01:19:18	3521
	4	L=2, K=5, F=8	6.2626	00:48:00	2231
STConvS2S-C (ours)	1	L=1, K=3, F=16	6.3991	00:06:16	609
	2	L=2, K=3, F=8	6.3660	00:08:23	654
	3	L=2, K=3, F=16	6.3623	00:12:45	807
	4	L=2, K=5, F=8	6.3131	00:13:16	662
STConvS2S-R (ours)	1	L=1, K=3, F=16	6.3910	00:06:05	584
	2	L=2, K=3, F=8	6.3310	00:07:14	616
	3	L=2, K=3, F=16	6.3205	00:11:01	735
	4	L=2, K=5, F=8	6.2288	00:10:55	616

Table 4: Performance results for rainfall forecasting using the previous five observations (grids) to predict the next five observations ($5 \rightarrow 5$), and the next 15 observations ($5 \rightarrow 15$).

Model	$5 \rightarrow 5$				
	RMSE	MAE	Memory usage (MB)	Mean training time	Training time/epoch
ARIMA	7.4377	6.1694	—	—	—
ConvLSTM (Shi et al., 2015)	6.3666 ± 0.0019	2.9074 ± 0.0185	752	00:15:15	00:00:13
PredRNN (Wang et al., 2017)	6.2625 ± 0.0039	2.7880 ± 0.0110	1740	00:39:59	00:00:43
MIM (Wang et al., 2019)	6.2621 ± 0.0051	2.7900 ± 0.0178	2231	00:52:13	00:00:52
STConvS2S-C (ours)	6.3091 ± 0.0029	2.8487 ± 0.0280	662	00:15:54	00:00:15
STConvS2S-R (ours)	6.2248 ± 0.0006	2.7821 ± 0.0261	616	00:16:48	00:00:13

Model	$5 \rightarrow 15$				
	RMSE	MAE	Memory usage (MB)	Mean training time	Training time/epoch
ARIMA	7.9460	5.9379	—	—	—
ConvLSTM (Shi et al., 2015)	6.3244 ± 0.0025	2.8972 ± 0.0264	1308	00:44:30	00:00:33
PredRNN (Wang et al., 2017)	6.2600 ± 0.0013	2.7850 ± 0.0067	4115	01:53:30	00:01:58
MIM (Wang et al., 2019)	6.2722 ± 0.0020	2.7935 ± 0.0246	5276	02:48:38	00:02:34
STConvS2S-C (ours)	6.2962 ± 0.0039	2.8452 ± 0.0130	916	00:35:43	00:00:41
STConvS2S-R (ours)	6.2590 ± 0.0023	2.8054 ± 0.0175	912	00:39:35	00:00:39

Figure 10: Cumulative error based on both horizons ($5 \rightarrow 5$ and $5 \rightarrow 15$) using the rainfall dataset (CHIRPS). Evaluations on RMSE and MAE metrics.

same settings of version 4 for CFSR (see Table 1) and CHIRPS (see Table 3) datasets. Besides analyzing the structure of our model, we also compare it with three models: vanilla 3D CNN, 3D Encoder-Decoder architecture (Racah et al., 2017), and a CNN model using (2+1)D blocks (Tran et al., 2018). In Table 5 and 6, the results of these comparisons are shown on rows 1-3, rows 4-11 show the ablation experiments and on rows 12-13 the proposed models in this work. Following, we discuss the experimental results in detail.

- **Removal of the factorized convolutions.** In these experiments, there is no separation in temporal and spatial blocks in our architecture, since this split is only possible due to factorized convolutions. These models in relation to layers are similar to 3D CNN (see Figure 2 (a) and (c)). Removing factorized convolutions from the models underperform the proposed models in RMSE and MAE metrics (see rows 4 and 12; rows 5 and 13). To reinforce that the performance gain comes from design options rather than increased model parameters, we also performed experiments removing the progressive filters from our models (rows 10-11) for comparison. The models without factorized convolutions still performed worse in most cases.

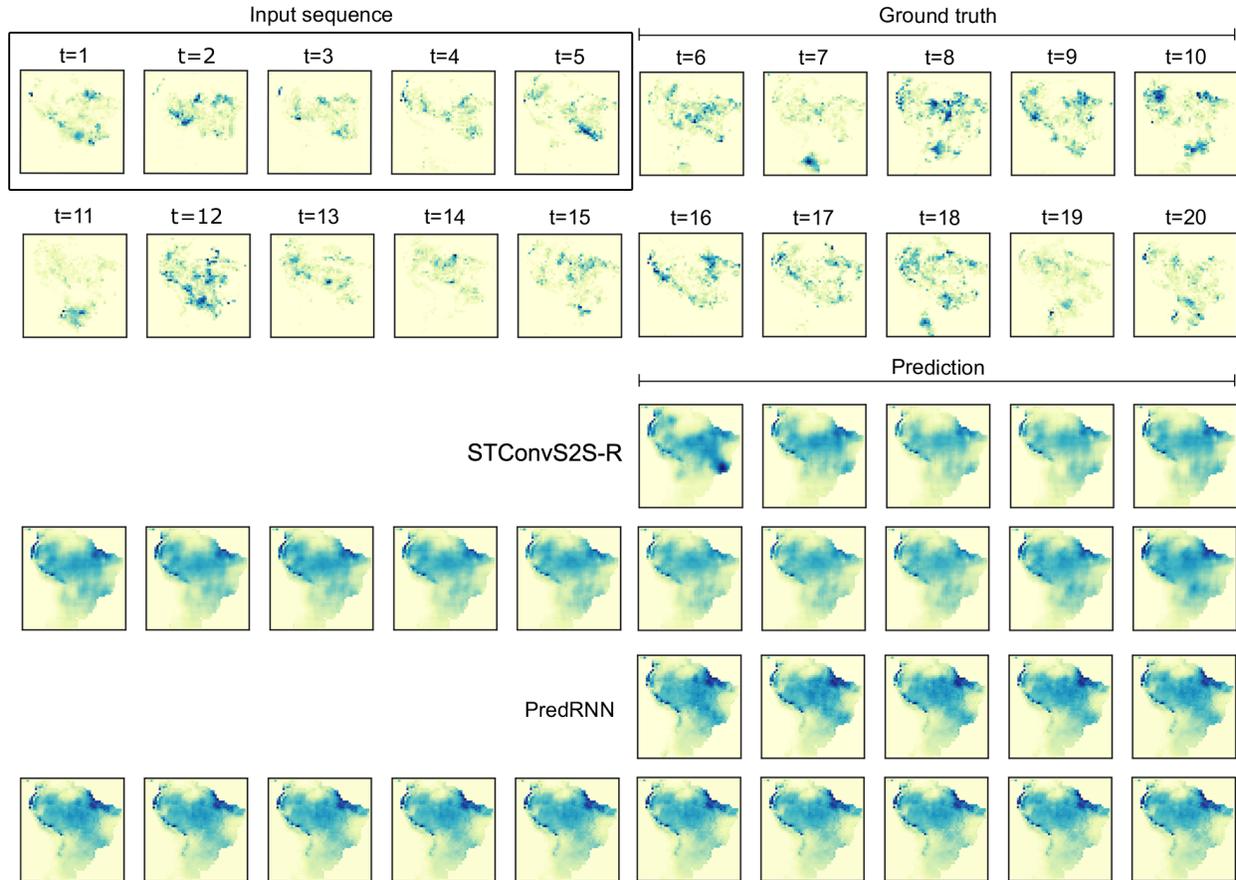


Figure 11: Prediction example on test set of rainfall dataset (CHIRPS). Comparison between the best CNN-based and RNN-based models: STConvS2S-R and PredRNN, respectively

- Removal of the causal constraint.** In general, respecting the temporal order is more a restriction of the problem domain than an additional feature to improve models' performance. However, the STConvS2S-R model slightly improves the results, at least in one of the evaluation metrics on both datasets (comparison between rows 6 and 13). This enhancement can also be observed in 3D CNN and "not factorized" STConvS2S-R (rows 1 and 5), as both use the same layers, but differ concerning the causal constraint. On the other hand, the STConvS2S-C model results do not show the same contribution to performance.
- Removal of the temporal block.** This experiment analyzes the importance of this component in our network since we propose two variations of STConvS2S based on the temporal block adopted. The results on row 7 indicate that although faster than the proposed models, this removal has a critical impact on RMSE and MAE, especially for the CFSR dataset.
- Inverted blocks.** Understanding the influence of the temporal and spatial blocks on each other is not straightforward. Thus, to analyze the model structure, we change the blocks from (temporal \Rightarrow spatial) to (spatial \Rightarrow temporal). Both GPU memory usage and training time are very similar, comparing each STConvS2S model with the respective inverted version. Regarding the evaluation metrics, the inverted versions have the worst performance on both datasets, except for MAE on CHIRPS using STConvS2S-R.
- Comparison with baselines methods.** There is no significant differentiation among STConvS2S models and the baselines on the CHIRPS dataset concerning memory usage and training time. These metrics almost increase twice on the CFSR dataset compared to 3D CNN and 3D Encoder-Decoder, but as a result, STConvS2S-R achieves a 4% improvement in RMSE over those same baselines. STConvS2S-R had a favorable or matching performance in the results of the experiments on both datasets, except for the MAE metric on the CHIRPS dataset comparing against (2+1)D Conv. STConvS2S-C does not perform as well as STConvS2S-R in these comparisons.

Table 5: Quantitative comparison of ablation experiments, baseline methods using 3D convolutional layers, and our proposed models on temperature dataset (CFSR) for 5 → 5 task.

Model	Factorized Conv.	Causal const.	RMSE	MAE	Memory usage (MB)	Training time
1. 3D CNN	—	—	1.3307	0.9015	578	01:13:54
2. 3D Encoder-Decoder	—	—	1.3327	0.9291	544	01:11:26
3. (2+1)D Conv	✓	—	1.2944	0.8763	847	01:51:24
4. STConvS2S-C	—	✓	1.4450	1.0068	605	01:39:47
5. STConvS2S-R	—	✓	1.3215	0.8958	580	01:16:53
6. STConvS2S	✓	—	1.2811	0.8645	884	02:00:52
7. STConvS2S*	✓	—	1.6780	1.1828	740	01:14:47
8. STConvS2S-C**	✓	✓	1.4152	0.9750	1000	02:15:46
9. STConvS2S-R**	✓	✓	1.3044	0.8796	895	01:56:05
10. STConvS2S-C***	✓	✓	1.4218	0.9821	698	01:04:58
11. STConvS2S-R***	✓	✓	1.3234	0.8966	649	00:57:11
12. STConvS2S-C	✓	✓	1.3791	0.9492	1040	02:25:26
13. STConvS2S-R	✓	✓	1.2773	0.8646	895	01:58:39

* No temporal block

** Inverted (spatial ⇒ temporal)

*** No filter increase

Table 6: Quantitative comparison of ablation experiments, baseline methods using 3D convolutional layers, and our proposed models on rainfall dataset (CHIRPS) for 5 → 5 task.

Model	Factorized Conv.	Causal const.	RMSE	MAE	Memory usage (MB)	Training time
1. 3D CNN	—	—	6.2519	2.8519	534	00:09:27
2. 3D Encoder-Decoder	—	—	6.2540	2.7977	513	00:09:26
3. (2+1)D Conv	✓	—	6.2323	2.7243	660	00:12:09
4. STConvS2S-C	—	✓	6.3310	2.9161	553	00:12:38
5. STConvS2S-R	—	✓	6.2510	2.8082	534	00:09:55
6. STConvS2S	✓	—	6.2281	2.8134	609	00:10:19
7. STConvS2S*	✓	—	6.3539	2.8980	572	00:06:57
8. STConvS2S-C**	✓	✓	6.3255	2.8594	656	00:12:10
9. STConvS2S-R**	✓	✓	6.2397	2.7971	616	00:10:56
10. STConvS2S-C***	✓	✓	6.3171	2.8418	591	00:10:23
11. STConvS2S-R***	✓	✓	6.2434	2.7829	565	00:09:25
12. STConvS2S-C	✓	✓	6.3131	2.8327	662	00:13:16
13. STConvS2S-R	✓	✓	6.2288	2.8060	616	00:10:55

* No temporal block

** Inverted (spatial ⇒ temporal)

*** No filter increase

- **Comparison of strategies to satisfy the causal constraint.** The results show the superiority of STConvS2S-R compared to STConvS2S-C in the evaluation metrics and time performance in all the experiments. The hypothesis of STConvS2S-R to better predict is that this model adds less zero-padding before performing the convolution operation on each layer inside the temporal block. Concerning time efficiency, in STConvS2S-R, the reverse function is performed only twice in the temporal reversed block, making it faster than STConvS2S-C. The latter performs its operations on each layer within the temporal causal block (Section 4.2). This study demonstrates the relevance of our original method to make convolutional layers satisfy the causal constraint during the learning process.

6 Conclusion

Predicting future information many steps ahead can be challenging, and a suitable sequence-to-sequence architecture to better represent spatiotemporal data for this purpose is still open for research. However, RNN-based models have been widely adopted in these cases (Shi et al., 2015; Wang et al., 2017, 2019). Previously to this work, CNN-based architectures were not considered for this task, due to two limitations. Firstly, they do not respect the temporal order in the learning process. They also cannot generate an output sequence of length that is higher than the input sequence. Considering these limitations, we proposed STConvS2S, an end-to-end trainable deep learning architecture. STConvS2S can do spatiotemporal data forecasting by only using 3D convolutional layers.

First, we address the problem of causal constraint, proposing two variations of our architecture, one using the temporal causal block and the other, the temporal reversed block. The former adopts causal convolution, which is commonly used in 1D CNN. We also introduced a new technique in the temporal reversed block that makes it not violate the temporal order by applying a reverse function in the sequence. These implementations are essential for a fair comparison with state-of-the-art methods, which are causal models due to the chain-like structure of LSTM layers. To overcome the sequence length limitation, we designed a temporal generator block at the end of the architecture to extend the spatiotemporal data only in the time dimension. We further compared our models with state-of-the-art models through experimental studies in terms of both performance and time efficiency on meteorological datasets. The results indicate that our model manages to analyze spatial and temporal data dependencies better since it has achieved superior performance in temperature forecasting, and comparable results in the rainfall forecasting, with the advantage of being up to 5x faster than RNN-based models. Thus, STConvS2S could be a natural choice for sequence-to-sequence tasks when using spatiotemporal data. We evaluate our architecture in the weather forecasting problem, but it is not limited to this domain. We expect that the results presented in this work will foment more research with comparisons between convolutional and recurrent architectures.

For future work, we will search for ways to decrease rainfall dataset error. Directions may include applying preprocessing techniques to sparse data and adding data from other geographic regions. Besides, we will investigate more architectures for spatiotemporal data forecasting.

Computer Code Availability

We implemented the deep learning models presented in this paper using PyTorch 1.0, an open-source framework. Our source code is publicly available at <https://github.com/MLRG-CEFET-RJ/stconvs2s>

Data Availability

In this paper, spatiotemporal datasets in NetCDF format were used and can be downloaded at <http://doi.org/10.5281/zenodo.3558773>, an open-source online data repository.

Acknowledgment

The authors thank CNPq, CAPES, FAPERJ, and CEFET/RJ for partially funding this research.

References

- Baboo, S., Shereef, K., 2010. An efficient weather forecasting system using artificial neural network. *International Journal of Environmental Science and Development* 1, 321–326. doi:10.7763/IJESD.2010.V1.63.
- Babu, C.N., Reddy, B.E., 2012. Predictive data mining on Average Global Temperature using variants of ARIMA models. *IEEE-International Conference on Advances in Engineering, Science and Management, ICAESM-2012*, 256–260.
- Bai, S., Zico Kolter, J., Koltun, V., 2018. An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling. arXiv:1803.01271v2.
- Cheng, C., Zhang, C., Wei, Y., Jiang, Y.G., 2019. Sparse Temporal Causal Convolution for Efficient Action Modeling, in: *Proceedings of the 27th ACM International Conference on Multimedia (MM'19)*, ACM. pp. 592–600. doi:10.1145/3343031.3351054.
- Corchado, J.M., Fyfe, C., 1999. Unsupervised neural method for temperature forecasting. *Artificial Intelligence in Engineering* 13, 351–357. doi:10.1016/S0954-1810(99)00007-2.

- Fujiwara, M., Wright, J.S., Manney, G.L., Gray, L.J., Anstey, J., Birner, T., Davis, S., Gerber, E.P., et al., 2017. Introduction to the sparc reanalysis intercomparison project (s-rip) and overview of the reanalysis systems. *Atmospheric Chemistry and Physics* 17, 1417–1452. doi:10.5194/acp-17-1417-2017.
- Funk, C., Peterson, P., Landsfeld, M., Pedreros, D., Verdin, J., Shukla, S., Husak, G., Rowland, J., Harrison, L., Hoell, A., Michaelsen, J., 2015. The climate hazards infrared precipitation with stations - A new environmental record for monitoring extremes. *Scientific Data* 2. doi:10.1038/sdata.2015.66.
- Gehring, J., Auli, M., Grangier, D., Yarats, D., Dauphin, Y.N., 2017. Convolutional sequence to sequence learning, in: *International Conference on Machine Learning*, pp. 2029–2042.
- Karpatne, A., Ebert-Uphoff, I., Ravela, S., Babaie, H.A., Kumar, V., 2018. Machine Learning for the Geosciences: Challenges and Opportunities. *IEEE Transactions on Knowledge and Data Engineering PP*, 1. doi:10.1109/TKDE.2018.2861006.
- Kim, S., Kim, H., Lee, J., Yoon, S., Kahou, S.E., Kashinath, K., Prabhat, 2019. Deep-hurricane-tracker: Tracking and forecasting extreme climate events, in: *Proceedings of the IEEE Winter Conference on Applications of Computer Vision, WACV 2019*, pp. 1761–1769. doi:10.1109/WACV.2019.00192.
- Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. ImageNet Classification with Deep Convolutional Neural Networks, in: *Proceedings of the 25th International Conference on Neural Information Processing Systems (NeurIPS)*, Curran Associates Inc.. pp. 1097–1105.
- Mehdizadeh, S., 2018. Assessing the potential of data-driven models for estimation of long-term monthly temperatures. *Computers and Electronics in Agriculture* 144, 114–125. doi:10.1016/j.compag.2017.11.038.
- van den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A.W., Kavukcuoglu, K., 2016. Wavenet: A generative model for raw audio. arXiv:1609.03499.
- Racah, E., Beckham, C., Maharaj, T., Ebrahimi Kahou, S., Prabhat, M., Pal, C., 2017. ExtremeWeather: A large-scale climate dataset for semi-supervised detection, localization, and understanding of extreme weather events, in: *Proceedings of the 30th International Conference on Neural Information Processing Systems (NeurIPS)*, Curran Associates, Inc.. pp. 3402–3413.
- Rahman, M.R., Lateh, H., 2017. Climate change in Bangladesh: a spatio-temporal analysis and simulation of recent temperature and rainfall data using GIS and time series analysis model. *Theoretical and Applied Climatology* 128, 27–41. doi:10.1007/s00704-015-1688-3.
- Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., Prabhat, 2019. Deep learning and process understanding for data-driven Earth system science. *Nature* 566, 195–204. doi:10.1038/s41586-019-0912-1.
- Sagheer, A., Kotb, M., 2019. Time series forecasting of petroleum production using deep lstm recurrent networks. *Neurocomputing* 323, 203 – 213. doi:10.1016/j.neucom.2018.09.082.
- Saha, S., Moorthi, S., Wu, X., Wang, et al., 2014. The NCEP Climate Forecast System Version 2. *Journal of Climate* 27, 2185–2208. doi:10.1175/JCLI-D-12-00823.1.
- Shi, X., Chen, Z., Wang, H., Yeung, D.Y., Wong, W.k., Woo, W.c., 2015. Convolutional LSTM network: A machine learning approach for precipitation nowcasting, in: *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 802–810.
- Singh, G., Cuzzolin, F., 2019. Recurrent Convolutions for Causal 3D CNNs, in: *Proceedings of the IEEE International Conference on Computer Vision (ICCV) Workshops*, pp. 1–10.
- Souto, Y.M., Porto, F., Moura, A.M., Bezerra, E., 2018. A Spatiotemporal Ensemble Approach to Rainfall Forecasting, in: *Proceedings of the International Joint Conference on Neural Networks*, pp. 574–581. doi:10.1109/IJCNN.2018.8489693.
- Štulec, I., Petljak, K., Naletina, D., 2019. Weather impact on retail sales: How can weather derivatives help with adverse weather deviations? *Journal of Retailing and Consumer Services* 49, 1–10. doi:10.1016/j.jretconser.2019.02.025.
- Tran, D., Bourdev, L., Fergus, R., Torresani, L., Paluri, M., 2015. Learning spatiotemporal features with 3d convolutional networks, in: *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 4489–4497. doi:10.1109/ICCV.2015.510.
- Tran, D., Wang, H., Torresani, L., Ray, J., LeCun, Y., Paluri, M., 2018. A closer look at spatiotemporal convolutions for action recognition, in: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, IEEE Computer Society. pp. 6450–6459. doi:10.1109/CVPR.2018.00675.

- Wang, Y., Long, M., Wang, J., Gao, Z., Yu, P.S., 2017. Predrnn: Recurrent neural networks for predictive learning using spatiotemporal lstms, in: Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS), Curran Associates, Inc.. pp. 879–888.
- Wang, Y., Zhang, J., Zhu, H., Long, M., Wang, J., Yu, P.S., 2019. Memory in memory: A predictive neural network for learning higher-order non-stationarity from spatiotemporal dynamics, in: Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition, IEEE Computer Society. pp. 9146–9154. doi:10.1109/CVPR.2019.00937.
- Xu, S., Niu, R., 2018. Displacement prediction of Baijiabao landslide based on empirical mode decomposition and long short-term memory neural network in Three Gorges area, China. Computers and Geosciences 111, 87–96. doi:10.1016/j.cageo.2017.10.013.
- Xu, Z., Cao, Y., Kang, Y., 2019. Deep spatiotemporal residual early-late fusion network for city region vehicle emission pollution prediction. Neurocomputing 355, 183–199. doi:10.1016/j.neucom.2019.04.040.
- Yang, B., Sun, S., Li, J., Lin, X., Tian, Y., 2019. Traffic flow prediction using LSTM with feature enhancement. Neurocomputing 332, 320–327. doi:10.1016/j.neucom.2018.12.016.
- Yuan, Y., Zhao, Y., Wang, Q., 2018. Action recognition using spatial-optical data organization and sequential learning framework. Neurocomputing 315, 221–233. doi:10.1016/j.neucom.2018.06.071.
- Zhang, Q., Wang, H., Dong, J., Zhong, G., Sun, X., 2017. Prediction of Sea Surface Temperature Using Long Short-Term Memory. IEEE Geoscience and Remote Sensing Letters 14, 1745–1749. doi:10.1109/LGRS.2017.2733548.