# Convolutional Neural Networks and Long Short-Term Memory for Hyperspectral Classification and Time Series Prediction

**Juan F. Ramirez Rochac[1] and Nian Zhang[2]**
[1] Department of Computer Science and Information Technology and [2] Department of Electrical and Computer Engineering
University of the District of Columbia, Washington, DC, 20008
[1] jrochac@udc.edu and [2] nzhang@udc.edu

## Hypothesis

Machine Learning and Deep Learning (ML&DL) are areas of active research with multiple expert-level performing algorithms, yet additional research is needed to reach this level of accuracy on highly imbalanced, non-linear datasets. This is evident in the presence of environment data, such as remotely-sensed hyperspectral images and time series data that describe streamflow and surface-water quantity. In this study, we propose different approaches, as follows: (1) The implementation of context-based feature augmentation (CFA) to tackle highly imbalanced data in hyperspectral classification using deep convolutional neuronets; (2) The implementation of Scharr-based adaptive filtering (SAF) to deal with non-linearities in time series prediction using deep recurrent neuronets; and finally, (3) The application of CFA and SAF to improve performance accuracy in hyperspectral classification and time series prediction using real-world environment datasets.

## Introduction

With Conv-based and LSTM-based approaches at expert-level performance, researchers have found great success across a range of applications, such as autonomous driving/piloting, healthcare, cybersecurity, speech and image recognition, natural language processing and financial markets. However, these DL approaches are data-depended. Since they use labeled data to learn patterns during the fitting phase and later apply this knowledge during on unlabeled data the predict phase, the level of quantity and quality of the input data greatly impacts the learning ability and performance. *Data Quantity* refers to the number of labeled data samples. Recently, there has been an exponential growth in data capturing in every possible area of life. All sectors are gathering, collecting, warehousing massive datasets from research institutions to government agencies to private corporations. Academia, industry and governments are creating new sources (Internet of Things) with greater detail (Browsing History) and finer granularity. *Data Quality* refers to how well the data samples model the entire dataset. In the presence of environmental data, that is data collected with sensors that measure any environmental variable, deep learning algorithms leverage from these oceans of data and hybrid models to address emerging challenges in highly imbalanced data samples and non-linearity of the systems. In hyperspectral datasets, classification algorithms need to also address the curse of dimensionality, high levels of imbalanced data distribution, and the presence of noise while in time series datasets, such as streamflow, prediction algorithms need to deal with the non-linearity of the systems.

## Evaluation Methods

In all the experiments, the proposed CNN-variants and LSTM-variants use the Adam optimizer and a learning rate of 0.00001 for 200 epochs. The test bench for all performance experiments was implemented using Python 3.7. The runtime environment consists of a 2-core CPU, Intel(R) Xeon(R) @ 2.20GHz, a single GPU, Tesla K80, 4992 CUDA, 12GB GDDR5, 13GB of RAM and 80GB of HDD on a Linux-based virtual machine.

**Fig 3. Confusion Matrix**

**Fig 4. Classification Report**

$$RMSE = \sqrt{\frac{\sum_{i=1}^{N}(y_i - \hat{y}_i)^2}{N}}$$

**Eq. 4: Root Mean Square Error,** where $N$ represents the total number of $i$-th testing samples, $y_i$ corresponds to the actual observations and $\hat{y}_i$ the predicted values.
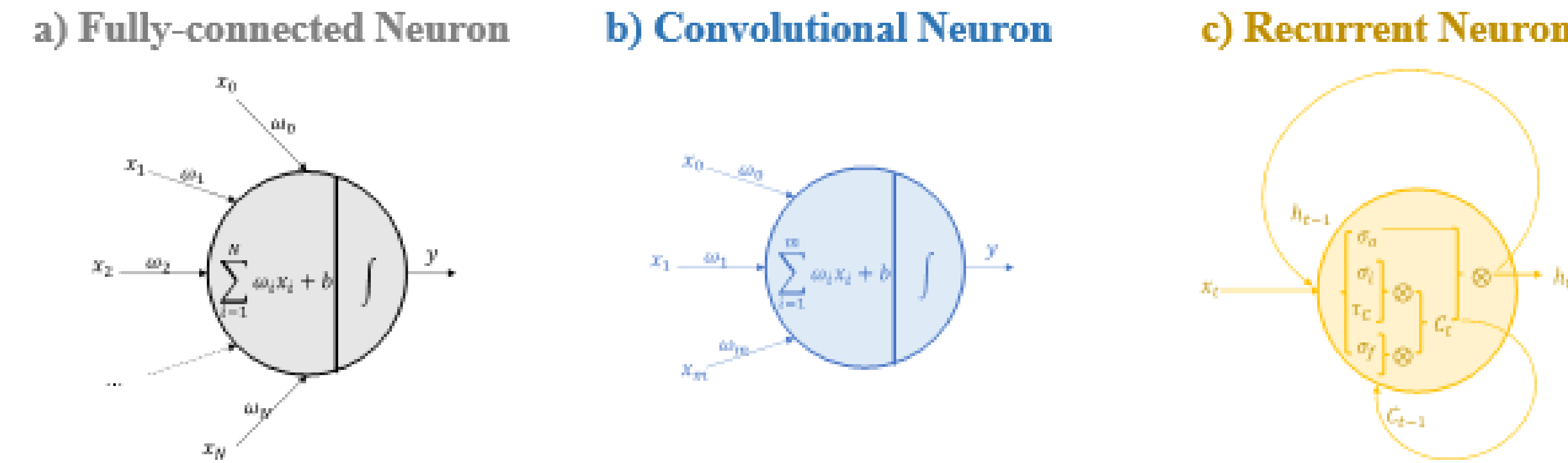
## Methodology

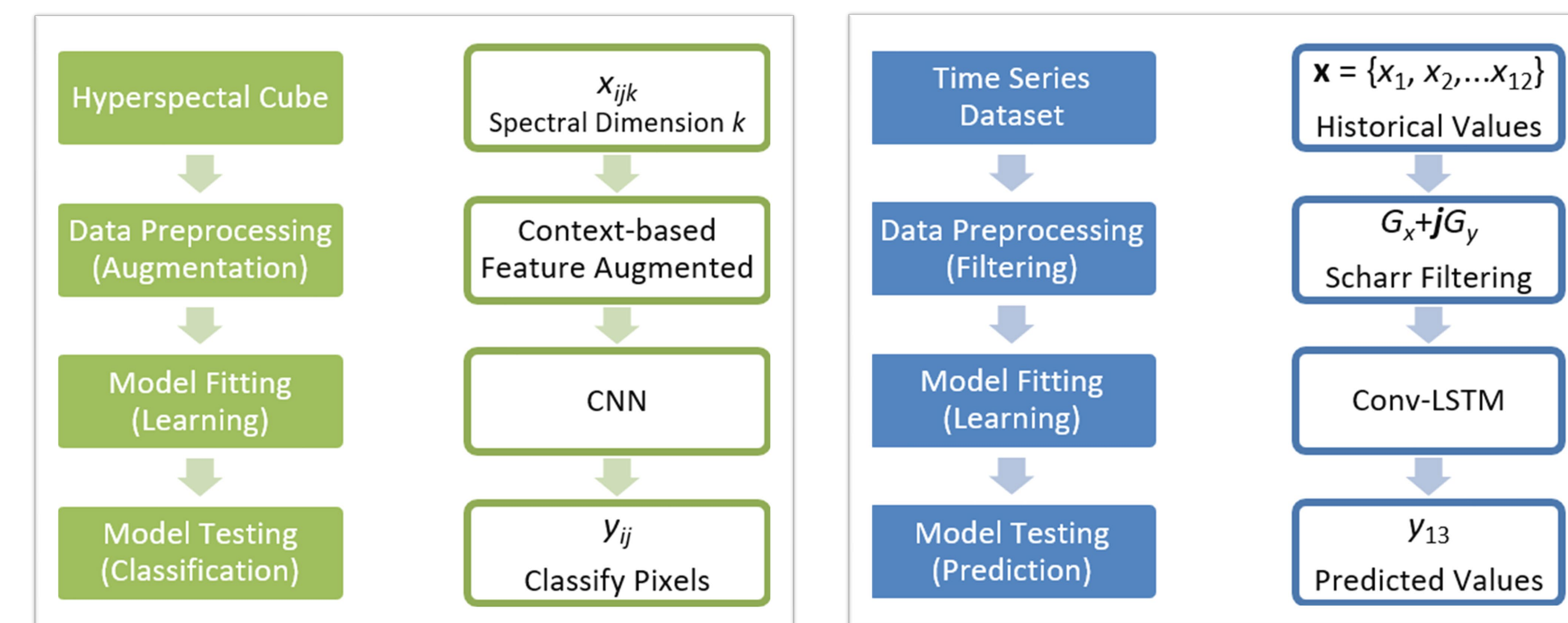**Fig 1:  Deep Learning Models: a) MLP, b) CNN and c) LSTM neurons**

**Fig 2:  Overview of the proposed hybrid DL-based frameworks for:**
**Hyperspectral Classification on the left and Streamflow Prediction on the right**

$$AGWN(r|\mu,\sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} exp\left[-\frac{(r-\mu)^2}{2\sigma^2}\right]$$

**Eq. 1: Additive Gaussian White Noise,** where $r$ is a uniformly random-generated number between [0,1], $\mu$ corresponds to expected value zero and $\sigma$ corresponds to the noise variance.

$$SNR_{dB} = 20\log_{10}\frac{\mu_{precip}}{\sigma_{noise}}$$

**Eq. 2: Signal-to-Noise Ratio,** where $\mu_{precip}$ corresponds to mean value of all streamflow data points and $\sigma_{noise}$ corresponds to the noise variance.

$$(G_x + jG_y) = Amp * \begin{bmatrix} -3-3j & 0-10j & +3-3j \\ -10+0j & 0+0j & +10+0j \\ -3+3j & 0+10j & +3+3j \end{bmatrix}$$

**Eq. 3: Complex Scharr operator,** where $G_x$ and $G_y$ represent the horizontal and vertical gradients, respectively, and $Amp$ is a scalar factor.
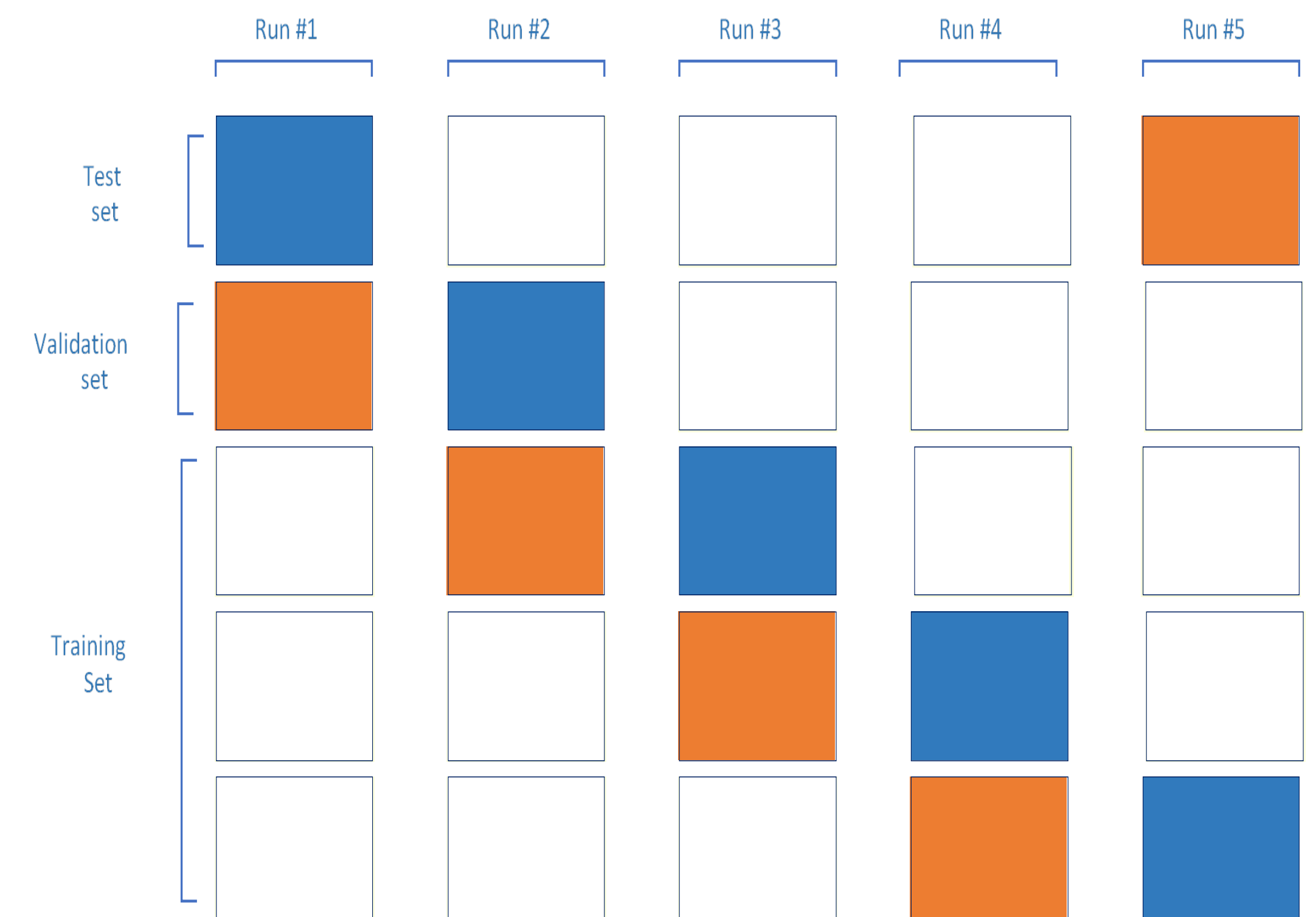
**Fig 5. Data Partitions for K-fold Cross Validation,** where the *Test set* consists of ⅕ of the data (blue squares), the *Validation set* also consists of ⅕ of the data (orange squares) and the *Training set* consists of the remaining ⅗ of the data (white squares). And, in every run, different partitions are used for training, validation and testing.

## Experimental Results

**Dataset #1:** The first dataset was collected by the ROSIS optical sensor over an urban area centered at the University of Pavia, Italy. The flight was operated by the Deutschen Zentrum for Luftund Raumfahrt (DLR, the German Aerospace Agency) in the framework of the HySens project, managed and sponsored by the European Union. In Figure 6, the collected cube size in pixels is 610 × 340 by 115 channels.

**Dataset #2:** The second, real-world dataset consists of monthly adjusted river streamflow in cubic feet per second (cfs) collected at the Potomac Basin above Little Falls, which is near Washington DC, USA. Figure 7 shows the raw streamflow values, from March 1930 to December 2021. These adjusted values were collected by the National Water Information System (NWIS) at the U.S. Geological Survey, Station No. 01646502.

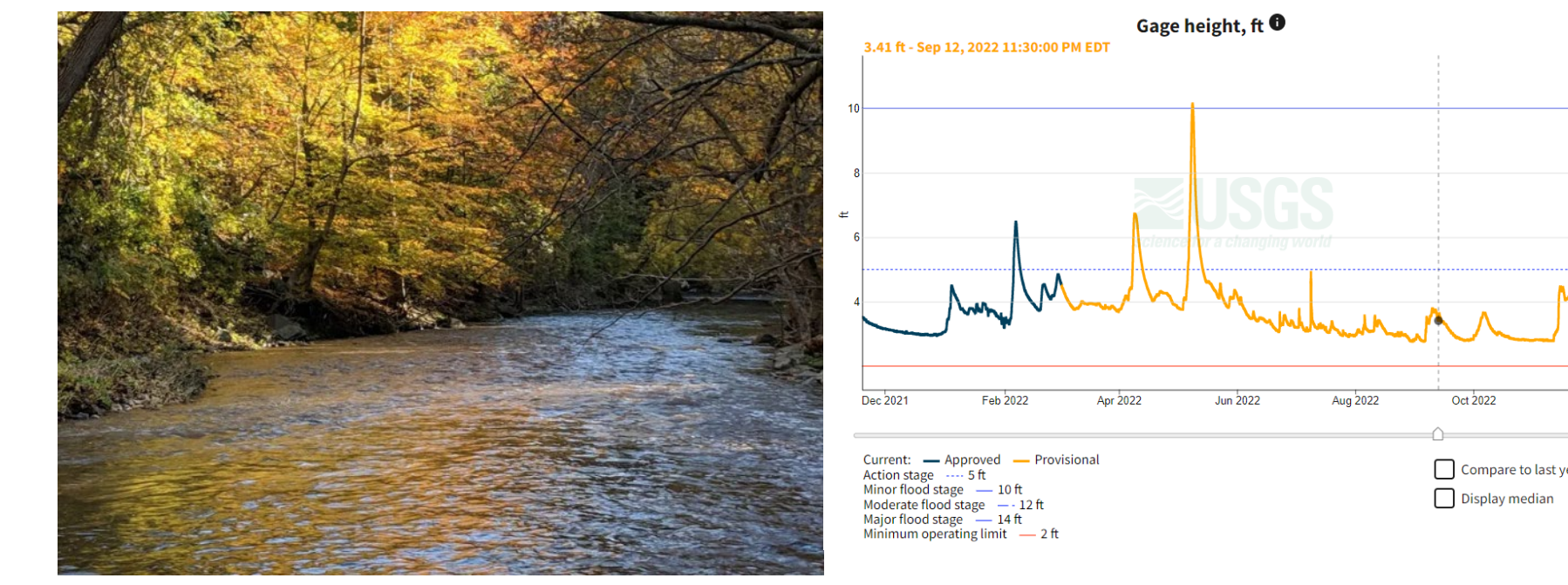**Fig. 6 Hyperspectral Cube Collection**

**Fig. 7 Streamflow Time Series Collection**

**Results**: Figure 8 and Table 1 summarize our remote sensing classification results published in [1], whereas Figure 9 and Table 2 summarize our streamflow quantity prediction results published in [2]. The ideal confusion matrix will only have nonzero values on the main diagonal. The ideal Predicted vs Actual plot will overlap perfectly generating an Absolute Error plot equal to the constant zero (a horizontal line). The perfect accuracy will have a value of 100.0 ± 0.0.

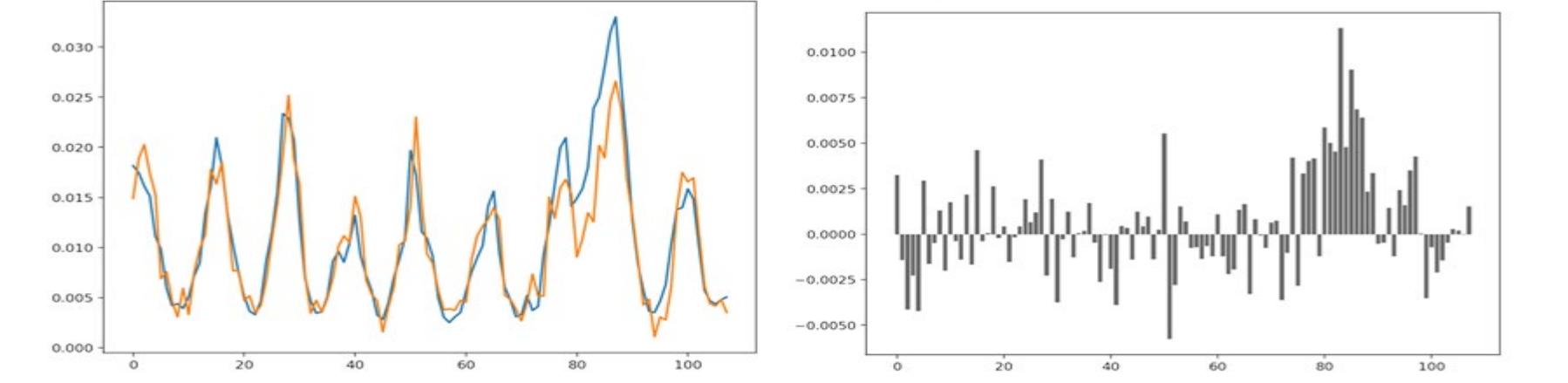**Fig. 8 Confusion Matrices for MLP, CNN and LSTM**

**Fig. 9 Predicted vs Actual and Absolute Error**

**Table 1: Hyperspectral Classification Accuracy**

| Technique | MLP | CNN | LSTM |
|---|---|---|---|
| No preprocess | 78.0 ± 9.0 | 80.7 ± 9.1 | 44.0 ± 9.2 |
| PCA | 77.6 ± 6.9 | 84.0 ± 6.5 | 60.0 ± 6.7 |
| CFA | 93.6 ± 4.7 | 99.0 ± 2.9 | 96.0 ± 3.8 |

**Table 2: Time Series Prediction Accuracy**

| Technique | LSTM | Bi-LSTM | ConvLSTM |
|---|---|---|---|
| No preprocess | 0.108 ± 0.005 | 0.118 ± 0.015 | 0.106 ± 0.003 |
| VMD | 0.057 ± 0.003 | 0.064 ± 0.002 | 0.067 ± 0.002 |
| SAF | 0.0043 ± 0.0003 | 0.0034 ± 0.0005 | 0.0028 ± 0.0003 |

## Conclusion

This research work proposes CFA-assisted hybrid approach for hyperspectral classification and a SAF-assisted hybrid approach for streamflow prediction. The experimental results demonstrate that the proposed CFA-assisted hybrid approach can effectively improve the overall classification accuracy for real-world data, and the proposed SAF-assisted hybrid approach also obtained competitive and even better performance compared with several state-of-the-art methods. In addition, our proposed design achieves comparable performance in terms of prediction time. The future work may include applying the proposed hybrid approaches into other environmental datasets, such as wind prediction.

## Acknowledgement

## References

[1] JF Ramirez Rochac, N Zhang, LA Thompson and T Deksissa, "A Robust Context-based Deep Learning Approach for Highly Imbalanced Hyperspectral Classification" *Computational Intelligence and Neuroscience*, ID 9923491, 2021.

[2] JF Ramirez Rochac, N Zhang, T Deksissa, J Xu and LA Thompson, "A Hybrid ConvLSTM Deep Neural Network for Noise Reduction and Data Augmentation for Prediction of Non-linear Dynamics of Streamflow," *10th Workshop on Data Mining in Earth System Science* (DMESS 2022) at the *IEEE International Conference on Data Mining* (ICDM 2022)