

# ENHANCING REMOTE SENSING IMAGE FUSION AND CLASSIFICATION ACCURACY USING DEEP LEARNING MODELS

G. Brindha

Department of Artificial Intelligence and Data Science, Dr. N.G.P. Institute of Technology, India

## Abstract

*Remote sensing imagery has become a pivotal source for land-use information at broad spatial scales due to advancements in satellite technology. However, challenges persist in accurately segmenting and classifying remote sensing data, particularly with high-resolution imagery. This paper proposes a novel hybrid deep learning model for spatiotemporal fusion to address these challenges, integrating SRCNN and LSTM models. The SRCNN enhances spatial details using MODIS-Landsat image pairs, while the LSTM learns phenological patterns in the enhanced images, facilitating dynamic agricultural system predictions. Evaluation comparing against benchmark fusion models. Implementation details are provided, including the use of loss functions for image segmentation and training specifics. Results demonstrate superior performance in land cover extraction accuracy compared to existing models, with an overall accuracy of 95.77% and a mean Intersection over Union (MIoU) of 82.23%. This study highlights the effectiveness of the proposed hybrid model in capturing both spatial and temporal dynamics, essential for applications ranging from land cover mapping to disaster assessment.*

## Keywords:

*Deep Learning, SRCNN, LSTM, MioU Scores*

## 1. INTRODUCTION

Currently, remote sensing pictures are considered the gold standard for doing land-use analysis on a wide scale [1]. There has been a substantial improvement in the spatial-temporal resolution of remote sensing imagery as a result of the continued development of satellite remote sensing technology and remote sensing platforms. When it comes to remote sensing imagegraphs, the data that is obtained by high-resolution satellite sensors for remote sensing delivers more clear ground geometry, spatial information, and rich texture information. The spatial resolution of hyperspectral remote sensing satellite data is frequently insufficient for use in remote sensing applications, despite the fact that the data contains a wealth of spectral information. Surface feature extraction, categorization, and application are all simplified as a result of the wealth of surface feature information that is present in high-resolution satellite remote sensing images. Due to the extremely robust feature extraction capabilities that image segmentation algorithms possess, an increasing number of academics are adding them into the process of remote sensing picture categorization [2].

GEographic Object-Based Picture Analysis (OBIA) is said to rely extensively on picture segmentation, as evidenced by [3]. Accurate remote sensing data segmentation could be of tremendous use for a wide variety of applications, including but not limited to mapping land cover and agricultural monitoring, urban development surveys, and disaster damage assessment [4]. Picture segmentation, which involves classifying images down to the pixel level, is an essential area of study for remote sensing image classification [5].

Image segmentation is a subfield of image classification. Pictures are often segmented using methods that entail the extraction of features and the classification of those features [6,7]. The tasks involved in remote sensing image categorization are as diverse as the fields that make use of and perform operations on these images. Throughout the course of history, the major basis for the categorization of remote sensing pictures has been the spectral disparities of ground features, with prior knowledge playing a secondary role. The differentiation of different ground features can be accomplished by the utilization of a wide range of spectral properties [8,9].

The classification of common ground characteristics such as water, agricultural land, and vegetation is accomplished by the utilization of normalized difference indices (NDWI, NDVI, NDBI, and so on) [10,11,12]. On the other hand, when we are presented with the phenomenon of “same spectral from different materials” and “same material with different spectral,” we may experience misclassification and accuracy concerns as a result of depending entirely on spectral information for the categorization of ground features. Within the realm of high-resolution remote sensing pictures, the spectral resolution is quite modest, while the spatial resolution is quite great.

## 2. DATASET

The Gaofen Image Dataset (GID) serves as a benchmark dataset that we employ in order to validate the methodologies that we have developed. For the purpose of training CNN models, the research investigation utilized partially processed GF-2 satellite images from the GID dataset, which is accessible to the general public. These images only comprise green, blue, and red bands. The fine land-cover classification set and the large-scale classification are the two fundamental components that make up the Geographic Information System (GID).

Annotations are made for five exemplary land-use categories in the case of the former. These categories include urban, agricultural, woodland, meadow, and water. These different land-use categories were represented by the colors red, green, blue, cyan, and yellow. These colors were used to identify the categories. Additionally, regions that do not fall into any of the five categories, as well as any areas that are considered to be clutter zones, are designated as backdrop and colored black. There are 15 subcategories that are included in the fine land-cover categorization set.

These subcategories include paddy field, irrigated farmland, dry cropland, garden, arbor forest, shrub, natural meadow, artificial meadow, industrial, urban, rural, residential, traffic, river, lake, and pond. The annotation of unknown regions is also performed in the event that these locations do not fall into any of the categories or if they were not manually identified.

### 3. HYBRID DEEP LEARNING MODEL

In this study, we have created a new hybrid deep learning model for spatiotemporal fusion in order to better the prediction of spatially detailed information and variable phenological changes in dynamic agricultural systems. This was done in order to improve the accuracy of our predictions. The SRCNN model is trained to enhance the spatial details using MODIS-Landsat picture pairs as part of the hybrid deep learning model. Meanwhile, the LSTM model is trained to recognize patterns of phenological change in these improved images. Both models are trained to improve spatial details.

In addition, we develop three scenarios that depict different degrees of phenological changes that occur over the course of time in the fusion images. The phenological transition dates of the crop are used to determine if these scenarios involve rapid, moderate, or minimal phenological changes. Following that, we put the hybrid deep learning model through its paces by putting it through its paces under a variety of phenological change thresholds. The final decision is to select three benchmark image fusion models, namely STARFM, FSDAF, and STFDCNN, in order to conduct a more comprehensive evaluation of the performance of the model.

During the process of spatiotemporal image fusion, it is possible to retrieve both the spatial linkage between coarse MODIS and fine Landsat pictures as well as the temporal relationship between images that were taken on different dates. It is required to retrieve spatial and temporal linkages with a high level of accuracy to develop fusion models that are useful. During our investigation, we successfully developed a hybrid architecture for deep learning models that is capable of managing both kinds of interactions. The SRCNN and LSTM models are what make up the hybrid deep learning model. This model is a blend of the two. It is decided to employ the SRCNN model because of its convolutional operations' ability to restore the spatial information that has been degraded in coarse images and to register the reflectance of both coarse and fine images.

SRCNN is an effective model for mapping spatial properties in satellite data. This is mostly due to the fact that its construction is quite lightweight. After that, the LSTM model's one-of-a-kind recurrent network structure is utilized in order to carry out the process of learning the temporal phenological changes that occur among the SRCNN-derived pictures. This is accomplished with the assistance of the registered reflectance and the restored spatial information.

#### 3.1 HYBRID DEEP LEARNING - LSTM

With the help of the long short-term memory (LSTM) component of the hybrid deep learning model, the objective is to acquire knowledge of phenological patterns that have developed over time from a collection of imagegraphs. In addition to resolving the issue of disappearing gradients that is associated with RNN, LSTM improves the structure of the network cells by incorporating a gating mechanism. This mechanism enables the storage of information in memory for longer periods of time. It selectively keeps specific data in memory and adds new data to represent the patterns of change in the sequential data in order to regulate the flow of information that is changing over time.

This allows it to control the flow of communication that is occurring. By utilizing a mix of gates and memory cells, long-term short-term memory (LSTM) models are able to learn the characteristics that vary over time from a collection of images. This allows them to make predictions about future time series. Consequently, the peculiar design of LSTM offers significant promise for describing complex temporal phenological fluctuations across satellite data.

A cell state, an input gate, a forget gate, and an output gate are the components that make up a conventional LSTM cell unit. These components are used to regulate the transfer of data. The regulation of these gates makes it possible to bring about the gradual addition, deletion, or updating of information that is stored in the state of the cell.

### 4. IMPLEMENTATION

The Keras framework is utilized in conjunction with the Tensorflow backend in order to accomplish the creation of the hybrid deep learning model. Testing and training are carried out with the assistance of a GPU accelerator manufactured by NVIDIA and referred to as a GK110 "Kepler" K20X. The optimizer that the hybrid model uses is called Adam, and it is a system that makes use of adaptive learning rates. For determining the initial learning rate, an empirical value of 0.001 is utilized. The sub-images will be extracted from the original images on a regular basis by SRCNN.

10000 of these sub-images will be used for training purposes at random, and two thousand and five hundred will be reserved for testing. During the training phase of the Long Short-Term Memory (LSTM) approach, a random selection of 150,000 picture pixels is utilized, whereas the testing phase utilizes 30,000 pixels. To make use of the graphics processing unit's random-access memory (RAM), the mini-batch size has been set to 128.

Table.1. Effect of different loss functions on road segmentation results.

Loss Function	OA (%)	MioU (%)	Kappa
Softmax Loss	95.24	71.96	0.64
Dice Loss	95.20	74.27	0.68
BCE Loss	95.24	72.27	0.64
Dice Loss + BCE Loss	95.27	74.30	0.68

The DICE and BCE loss functions are combined in this study in order to increase the detect accuracy for surface characteristics such as highways. Other surface features, on the other hand, continue to employ the softmax loss function. An increase in the total classification accuracy is achieved by the complete application of these loss functions. We analyze the performance of the softmax and DICE-BCE loss functions in road picture segmentation on the basis of their respective capabilities. A breakdown of the results of the classification accuracy is presented in Table.1. The results for road ground characteristics are improved when the DICE-BCE loss function is applied, as opposed to simply applying the softmax loss function.

Table.2. Performance Analysis

Model Category	SRCNN		LSTM SRCNN	
	CA (%)	IoU (%)	CA (%)	IoU (%)
Desert	95.85	93.90	96.46	95.18
Cotton	90.78	86.28	92.88	88.89
Roads	72.68	56.15	95.27	74.30
Water	89.09	83.68	91.40	87.84
Wetland	90.38	86.17	92.98	89.83
Uncultivated land	83.50	78.38	89.13	84.67
Jujube trees	84.27	75.51	87.72	83.48
Populus euphratica	83.68	73.70	85.95	79.79
Buildings	84.19	80.89	89.95	86.82
Woodland	82.64	79.87	85.91	81.58
pear trees	89.39	79.31	90.83	85.46
backgrounds	94.97	90.61	95.73	92.44
OA (%)	93.62		95.77	
MIOU (%)	80.38		85.77	
Kappa coefficient	0.92		0.94	

By employing indicators such as OA, MIOU, Kappa coefficient, CA, and IoU, we are able to objectively and accurately evaluate the impact that each model has on the extraction accuracy of various land covers by making use of the test set. The results presented in Table.2 demonstrate that the model achieves an overall accuracy of 86.16% while it is being pre-trained on the COCO dataset.

## 5. CONCLUSION

The purpose of this work is to provide a comprehensive technique for evaluating remote sensing images, with a particular emphasis on the development of a hybrid deep learning model for spatiotemporal fusion. To appropriately segmenting and categorizing high-resolution remote sensing data, the suggested method combines SRCNN and LSTM models. This allows the method to overcome challenging situations. The intensive testing and evaluation that was performed on the dataset demonstrates that the model is capable of accurately representing both spatial and temporal dimensions of dynamics. When compared to benchmark fusion models, fusion models with high mean Intersection over Union (MIOU) scores and overall accuracy exhibit superior operational performance. There are many different applications for remote sensing that can benefit from the insights acquired from this study. Some of these applications include mapping land cover, monitoring urban development, and evaluating the damage caused by disasters. With additional applications and modifications of the hybrid deep learning model,

there is the potential for a significant improvement in the global knowledge and management of dynamic environmental systems.

## REFERENCES

- [1] D.A. Reynolds, "Experimental Evaluation of Features for Robust Speaker Identification," *IEEE Transactions on Speech and Audio Processing*, Vol. 2, pp. 639-643, 1994.
- [2] Mahesh P.K. and M.N. Shanmukhaswamy, "Comprehensive Framework to Human Recognition Using Palmprint and Speech Signal", *Communications in Computer and Information Science In Springer-Verlag Berlin Heideberg*, Vol. 131, pp. 368-377, 2011.
- [3] Farid Ahmed and Ira S. Moskowicz, "Correlation-Based Watermarking Method for Image Authentication Applications", *Optical Engineering*, Vol. 43, No. 08, pp. 1833-1838, 2004.
- [4] Bernard Achermann and Horst Bunke, "Combination of Face Classifiers for Person Identification", *Proceedings of International Conference on Pattern Recognition*, pp. 416-420, 1996.
- [5] S. Vasuki and V. Vaidehi, "Identification of Human Faces using Orthogonal Locality Preserving Projections", *Proceedings of International Conference on Computer Design and Applications*, pp. 718-722, 2009.
- [6] Rafael C. Gonzalez and Richard E. Woods, "*Digital Image Processing*", Prentice Hall, 2005.
- [7] P. Jagalingam and A.V. Hegde, "A Review of Quality Metrics for Fused Image", *Aquatic Procedia*, Vol. 4, No. 2, pp. 133-142, 2015.
- [8] Betsy Samuel and N. Vidya, "Full Reference Image Quality Assessment for Biometric Detection", *International Journal on Modern Trends in Engineering and Research*, Vol. 2, No. 6, pp. 1-13, 2015.
- [9] Mayuresh Gulame, K.R. Joshi and Kamthe, "A Full Reference Based Objective Image Quality Assessment", *International Journal on Advanced Electrical and Electronics Engineering*, Vol. 2, No. 6, pp. 123-129, 2013.
- [10] M. Amin-Naji and A. Aghagolzadeh, "Multi-Focus Image Fusion in DCT Domain using Variance and Energy of Laplacian and Correlation Coefficient for Visual Sensor Networks", *Journal of AI and Data Mining*, Vol. 6, No. 2, pp. 233-250, 2018.
- [11] Veerpal Kaur and Jaspreet Kaur, "Frequency Partitioning Based Image Fusion for CCTV", *International Journal on Computer Science and Information Technologies*, Vol. 6, No. 4, pp. 3968-3972, 2015.
- [12] C. Rama Mohan, S. Kiran and R. Pradeep Kumar Reddy, "Multi-Focus Image Synthesis based on DWT and Texture with Sharpening", *Pezzottaite Journals*, Vol. 4, No. 4, pp. 1662-1670, 2015.