

Spatial Spectral Transformer with Conditional Position Encoding for Hyperspectral Image Classification

Muhammad Ahmad, Muhammad Usama, Adil Mehmood Khan, Salvatore Distefano, Hamad Ahmed Altuwaijri, and Manuel Mazzara

Abstract—In Transformer-based Hyperspectral Image Classification (HSIC), predefined positional encodings (PEs) are crucial for capturing the order of each input token. However, their typical representation as fixed-dimension learnable vectors makes it challenging to adapt to variable-length input sequences, thereby limiting the broader application of Transformers for HSIC. To address this issue, this study introduces an implicit conditional PEs (CPEs) scheme in a Transformer for HSIC, conditioned on the input token's local neighborhood. The proposed SSFormer integrates spatial-spectral information and enhances classification performance by incorporating a CPE mechanism, thereby increasing the Transformer layers' capacity to preserve contextual relationships within the HSI data. Moreover, SSFormer ensembles the cross-attention between patches and proposed learnable embeddings. This enables the model to capture global and local features simultaneously while addressing the constraint of limited training samples in a computationally efficient manner. Extensive experiments on publicly available HSI benchmarking datasets were conducted to validate the effectiveness of the proposed SSFormer model. The results demonstrated remarkable performance, achieving classification accuracies of 97.7% on the Indian Pines dataset and 96.08% on the University of Houston dataset.

Index Terms—Spatial Spectral Transformer (SSFormer); Hyperspectral Image Classification (HSIC).

I. INTRODUCTION

HYPERSPECTRAL IMAGING (HSI) has gained significant attention due to its ability to capture fine-grained spectral information, providing a wealth of data across multiple domains [1]–[7]. However, accurately classifying HSI data poses challenges due to its high dimensionality, spectral variability, complex spatial patterns, and Hughes phenomenon [8]–[10]. Several efforts have focused on developing accurate and computationally efficient algorithms to extract rich spectral information [11]. Recent works have developed effective classification frameworks based on jointly exploiting spectral-spatial features [12]–[14].

M. Ahmad and M. Usama are with the Department of Computer Science, National University of Computer and Emerging Sciences, Chiniot 35400, Pakistan. e-mail: mahmad00@gmail.com

A. M. Khan is with the School of Computer Science, University of Hull, Hull HU6 7RX, UK. e-mail: a.m.khan@hull.ac.uk

S. Distefano and M. Ahmad are with Dipartimento di Matematica e Informatica—MIFT, University of Messina, Messina 98121, Italy. e-mail: sdistefano@unime.it

H.A. Altuwaijri is with the Department of Geography, College of Humanities and Social Sciences, King Saud University, Riyadh, 11451 Saudi Arabia. e-mail: Haaltuwaijri@ksu.edu.sa

M. Mazzara is with the Institute of Software Development and Engineering, Innopolis University, 420500 Innopolis, Russia. e-mail: m.mazzara@innopolis.ru

Recent advances in Transformer-based models have demonstrated remarkable success in HSIC by employing self-attention mechanisms to process image patches directly [15]–[22]. For instance, combining spectral-spatial kernels with improved SST enables the joint extraction of spectral-spatial features [23]. Lifan et al. [24] introduced a spatial-spectral hierarchical Vision Transformer (ViT), while X. He et al. [25] proposed an SST that jointly exploits spectral-spatial information through a dedicated module. Hyper-ES2T and CSiT achieved state-of-the-art performance on several datasets using Transformer-based models [26], [27], and a spatial-spectral Transformer has been utilized for HSI denoising [28]. Despite their effectiveness, these models often face computational constraints and challenges in handling variable-length input sequences. Transformers rely on predetermined positional encodings (PEs) to capture the sequential order of input tokens, typically represented as learnable fixed-dimensional vectors. However, they struggle to adapt to variable-length input sequences, which is crucial for generating meaningful outputs from HSI data. This limitation restricts the broader applicability of Transformers.

The Spatial-Spectral Transformer (SSFormer) proposed in this study addresses the aforementioned challenges by effectively integrating spatial-spectral patterns through Conditional Positional Encoding (CPE) and cross-attention mechanisms. The key contributions can be summarized as:

- 1) Unlike traditional fixed-dimension encodings, the CPE is conditioned on the input token's local neighborhood, enabling the model to process varying input sequence lengths. This provides the Transformer with crucial pixel arrangement and sequence knowledge, allowing for better differentiation of pixels based on their relative positions and preserving the underlying spatial associations.
- 2) Cross-attention mechanisms enhance the ability to capture spatial context and relationships between pixels by combining positional information with spectral data. This integration improves the overall accuracy and efficiency of the model.
- 3) SSFormer leverages a backbone CNN to extract joint spatial-spectral features from the HSI cube, which are then encoded in the Transformer alongside CPEs. This combination enhances the model's ability to exploit spatial-spectral patterns, leading to improved generalization performance.

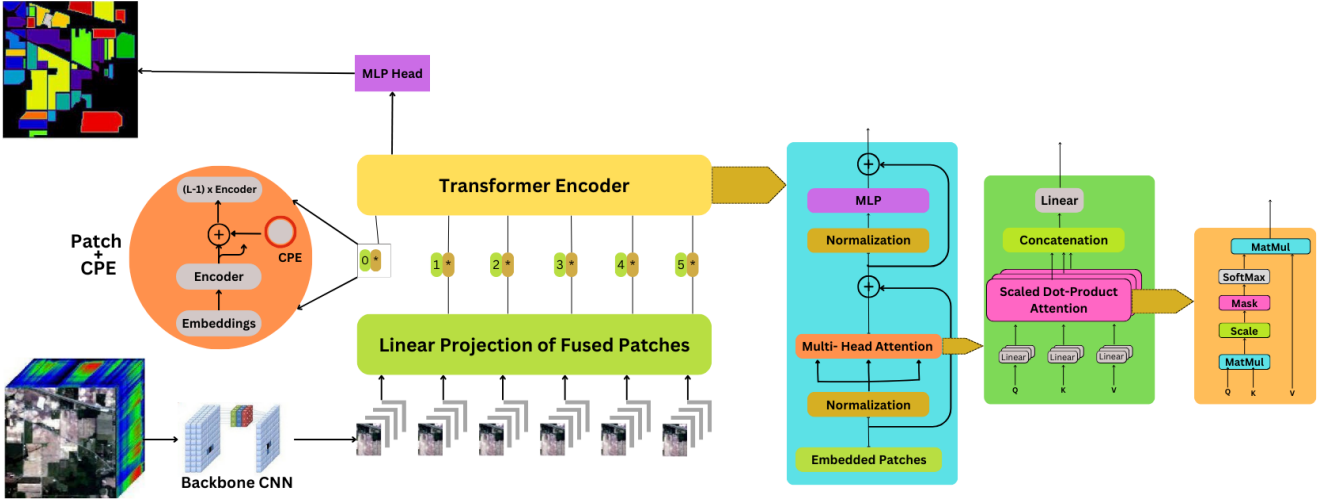


Fig. 1: SSFormer: A backbone CNN extracts joint spatial-spectral features from the HSI cube. These are encoded in the Transformer alongside CPEs, which selectively represent image information better than standard fixed-dimension encodings. The Transformer output then feeds into an MLP head for ground truth generation.

II. PROPOSED METHODOLOGY

To transform a 3D HSI cube $X \in \mathcal{R}^{M \times N \times B}$ into small 3D cubes suitable for the Transformer, the model first divides HSI into overlapping small patches. So, the model reproduces X into a sequence of patches which are then linearly projected into lower-dimensional embedding using a learnable weight matrix $W \in \mathcal{R}^{(D \times N \times S \times S)}$, where D , H , and S are embedding dimension, liner projection, and patch size, respectively:

$$H(n, d, i, j) = \sum_m \sum_k \sum_l X_{patches} * W(d, n, k, l) \quad (1)$$

$$X_{patches} = \phi \left(b_{i,j} + \sum_{\tau=1}^{d_{l-1}} \sum_{\sigma=-\delta}^{\delta} \sum_{\lambda=-v}^v \sum_{\rho=-\gamma}^{\gamma} w_{i,j,\tau}^{\sigma,\rho,\lambda} \times v_{i-1,\tau}^{x+\sigma,y+\rho,z+\lambda} \right) \quad (2)$$

where d_{l-1} , $b_{i,j}$, and $w_{i,j}$ correspond to the number of feature maps, bias parameter, and depth of the kernel for the j^{th} feature map at the $(l-1)^{th}$ layer, respectively. Additionally, $2v+1$, $2\gamma+1$, and $2\sigma+1$ represent the depth, width, and height of the kernel of 3D CNN as the backbone model [29]. To better incorporate spatial information, CPEs are added, i.e., for each position (i, j) in the patch grid, we generate a CPE $\mathcal{F} \in \mathcal{R}^D$. The CPE is added element-wise (where deemed necessary) to the embeddings, resulting in $H_{pos} \in \mathcal{R}^{(M \times D \times (P//S) \times (Q//S))}$:

$$H_{pos}(m, d, i, j) = H(m, d, i, j) + \mathcal{F}(pos, i) \quad (3)$$

where pos is the position of the element in the sequence and i is the index of the dimension in the embedding vector. The \mathcal{F} provides a unique CPE value for each position and dimension within the input sequence. For instance, to create the reference points, features undergo boundary padding, as an illustrative example, when conducting convolutions on feature maps, the use of zero padding serves to denote

the position of the boundary point. The combination of the above CPEs with the patch data allows the SSFormer to learn and capture positional information, ensuring that it can effectively process sequences with different orders or positions. \mathcal{F} can be viewed as a local regularization applied to $X_{patches}$. As the class token does not incorporate position information, it remains unaltered. Consequently, the ultimate output is created by concatenating the unmodified class token and the regularized $X_{patches}$ along the last dimension. The SSFormer employs multi-head attention to capture contextual relationships within the patches. To capture both spatial-spectral information, cross-attention is performed between the output (O_{att}) and the original embeddings (H_{pos}).

$$Attention^c(Q_c, K_c, V_c) = \left(\frac{Q_c \times K_c^T}{\sqrt{D}} \right) \times V_c \quad (4)$$

where K_c , V_c , and Q_c are the query, key, and value matrices for cross-attention, respectively. The output of the cross-attention is represented as:

$$O_{cross-att} = Attention^c(O_{att}, H_{pos}, H_{pos}) \quad (5)$$

The cross-attention output is fed into a feed-forward neural network with two fully connected layers to further process the features. Global average pooling is applied along the spatial dimensions to obtain a fixed-size representation for each band. The output is reshaped back into patches and concatenated to form the final output Y . The softmax function generates class probabilities to produce ground truth maps. The SSFormer effectively handles the HSI cube with reduced complexity by integrating spatial-spectral information through attention mechanisms and projections, making it suitable for resource-constrained environments. The complete model is demonstrated in Figure 1.

$$Y(m, d, i, j) = FFN(O_{cross-att}(m, d, i, j)) \quad (6)$$

III. EXPERIMENTAL RESULTS AND DISCUSSION

SSFormer’s weights were initially randomized and then optimized over 50 epochs using the Adam optimizer (learning rate 0.0001) and softmax loss. Training occurred in mini-batches of 256 samples per epoch, enabling the SSFormer to learn patterns from repeated exposure to training data and adjust its parameters to reduce loss over multiple iterations.

A. Patch Size and Train-Validation-Test Samples

Two key parameters, train/test split, and patch size, were meticulously examined. Various combinations of these parameters were tested and evaluated to determine the optimal configuration for achieving the best results. The model’s performance with different training and test percentages is tested as is shown in Table I, with a fixed patch size of 9×9 . Furthermore, Table II details the results for using various patch sizes.

TABLE I: We assess the performance of the proposed model using a 9×9 patch size with various Train/Validation/Test (Tr/Val/Te) splits.

Tr/Val/Te	AA	OA	κ	Tr Time	Te Time
Indian Pines Dataset					
5/5/90	64.44	74.60	70.82	54.49	5.89
10/10/80	77.12	81.92	79.24	66.75	5.92
15/15/70	82.42	86.41	84.42	80.31	3.10
20/20/60	87.97	89.78	88.33	149.64	3.62
Pavia University Dataset					
5/5/90	93.71	96.05	94.74	153.73	21.89
10/10/80	96.46	97.75	97.02	210.88	21.53
15/15/70	97.28	98.26	97.70	204.77	11.22
20/20/60	97.89	98.64	98.20	270.83	8.50
University of Houston Dataset					
5/5/90	94.93	96.08	95.77	46.87	6.2470
10/10/80	96.68	97.38	97.17	65.45	11.2144
15/15/70	97.63	98.09	97.94	154.74	6.2593
20/20/60	98.09	98.68	98.57	109.03	3.7363

TABLE II: Proposed model’s performance using a fixed Training(40%)/Validation(40%)/Test(20%) split and different patch sizes.

Patch Size	AA	OA	κ	Tr Time	Te Time	Parameters
Indian Pines Dataset						
3×3	85.35	88.29	86.68	184.39	1.53	784,016
5×5	92.98	93.31	92.37	210.57	2.04	799,376
7×7	94.90	94.09	93.26	211.32	2.03	822,416
9×9	96.53	96.14	95.60	175.90	1.53	853,136
Pavia University Dataset						
3×3	95.52	96.76	95.71	213.72	21.21	782,217
5×5	96.04	97.33	96.46	154.42	11.82	797,577
7×7	96.41	97.51	96.71	149.17	11.69	820,617
9×9	98.39	99.06	98.76	574.58	3.26	851,337
University of Houston Dataset						
3×3	97.46	97.70	97.52	197.73	1.57	779,217
5×5	99.06	99.26	99.20	214.27	1.58	799,119
7×7	99.34	99.40	99.35	212.12	2.04	822,159
9×9	99.13	99.36	99.31	214.47	2.30	852,879

The results indicate that SSFormer excels at capturing minute details and local patterns with smaller patch sizes. However, sensitivity to noise and fluctuations may lead to overfitting, especially with smaller sample sizes. In contrast, larger patch sizes enable SSFormer to encapsulate global features and contextual information, exhibiting enhanced robustness against noise perturbations. As sample sizes increase, SSFormer’s robustness and generalization capabilities improve. The optimal patch size depends on the specific

sample size, with smaller sample sizes benefiting from larger patch sizes to capture essential features, and larger sample sizes leveraging smaller patch sizes for precise pattern extraction. Increasing the training set size improved performance, particularly with larger patch sizes, by enabling the model to learn complex patterns and reduce overfitting. A moderate validation sample size is crucial for proper hyperparameter tuning. Too small a validation set might lead to suboptimal parameter choices, while too large a set might increase computation time without significant performance gains. Additionally, extremely small test sets might not provide a reliable estimation of SSFormer’s generalization ability.

B. Experimental Results with CNN-based Models

A consistent experimental methodology is crucial for fair model comparison. In this paper, we utilized unique sample distributions, geographical locations, and a controlled 9×9 pixel patch size across models. This controlled experimental setup allows for a fair comparison of competing CNN architectures. Figure 2 shows loss and accuracy trends on the University of Houston dataset. The ground truth maps for comparative and SSFormer are presented in Figures 3, 4, and 5.

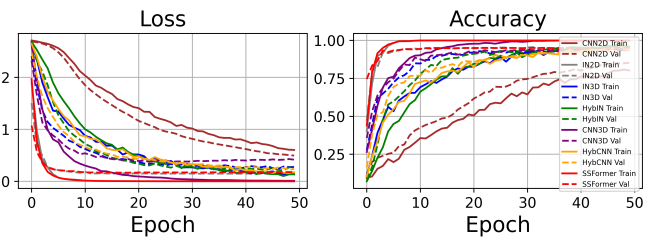
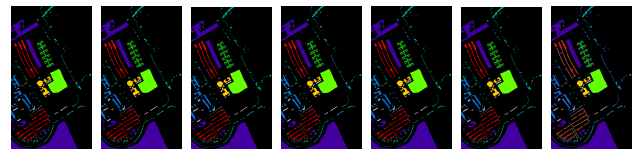


Fig. 2: The comparison above shows the accuracy and loss trends of CNN-based methods on the Houston dataset. The SSFormer (red line) demonstrates quicker convergence compared to traditional CNN models.



(a) [30] (b) [31] (c) [32] (d) [33] (e) [29] (f) [34] (g) SSF

Fig. 3: **Indian Pines Dataset:** The proposed SSFormer achieves OA=97.07% showing competitive performance.



(a) [30] (b) [31] (c) [32] (d) [33] (e) [29] (f) [34] (g) SSF

Fig. 4: **Pavia University Dataset:** The proposed SSFormer achieves OA=99.39% showing competitive performance.

Models evaluated include 3D CNN [29], Hybrid Inception Net [32], 3D Inception Net [31], 2D Inception Net [30], 2D CNN [33], and Hybrid CNN [34]. The findings highlight the advantage of decoupling spatial and spectral information over

REFERENCES

- [1] M. H. F. Butt, H. Ayaz, M. Ahmad, J. P. Li, and R. Kuleev, "A fast and compact hybrid cnn for hyperspectral imaging-based bloodstain classification," in *2022 IEEE Congress on Evolutionary Computation (CEC)*. IEEE, 2022, pp. 1–8.
- [2] M. H. Khan, Z. Saleem, M. Ahmad, A. Sohaib, H. Ayaz, M. Mazzara, and R. A. Raza, "Hyperspectral imaging-based unsupervised adulterated red chili content transformation for classification: Identification of red chili adulterants," *Neural Computing and Applications*, vol. 33, no. 21, pp. 14507–14521, 2021.
- [3] M. Zulfiqar, M. Ahmad, A. Sohaib, M. Mazzara, and S. Distefano, "Hyperspectral imaging for bloodstain identification," *Sensors*, vol. 21, no. 9, 2021.
- [4] M. H. Khan, Z. Saleem, M. Ahmad, A. Sohaib, H. Ayaz, and M. Mazzara, "Hyperspectral imaging for color adulteration detection in red chili," *Applied Sciences*, vol. 10, no. 17, 2020.
- [5] Z. Saleem, M. H. Khan, M. Ahmad, A. Sohaib, H. Ayaz, and M. Mazzara, "Prediction of microbial spoilage and shelf-life of bakery products through hyperspectral imaging," *IEEE Access*, vol. 8, pp. 176986–176996, 2020.
- [6] H. Ayaz, M. Ahmad, M. Mazzara, and A. Sohaib, "Hyperspectral imaging for minced meat classification using nonlinear deep features," *Applied Sciences*, vol. 10, no. 21, 2020.
- [7] H. Ayaz, M. Ahmad, A. Sohaib, M. N. Yasir, M. A. Zaidan, M. Ali, M. H. Khan, and Z. Saleem, "Myoglobin-based classification of minced meat using hyperspectral imaging," *Applied Sciences*, vol. 10, no. 19, 2020.
- [8] M. Ahmad, S. Shabbir, S. K. Roy, D. Hong, X. Wu, J. Yao, A. M. Khan, M. Mazzara, S. Distefano, and J. Chanussot, "Hyperspectral image classification—traditional to deep models: A survey for future prospects," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 968–999, 2021.
- [9] R. Hanachi, A. Sellami, I. Farah, and M. Dalla Mura, "Multi-view graph representation learning for hyperspectral image classification with spectral–spatial graph neural networks," *Neural Computing and Applications*, vol. 36, 12 2023.
- [10] X. Liao, B. Tu, J. Li, and A. Plaza, "Class-wise graph embedding-based active learning for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–13, 2023.
- [11] C. Pan, X. Jia, J. Li, and X. Gao, "Adaptive edge preserving maps in markov random fields for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 10, pp. 8568–8583, 2020.
- [12] N. Wambugu, Y. Chen, Z. Xiao, K. Tan, M. Wei, X. Liu, and J. Li, "Hyperspectral image classification on insufficient-sample and feature learning using deep neural networks: A review," *International Journal of Applied Earth Observation and Geoinformation*, vol. 105, p. 102603, 2021.
- [13] U. Ghous, M. S. Sarfraz, M. Ahmad, C. Li, and D. Hong, "Exnet: (2+1)d extreme xception net for hyperspectral image classification," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 17, pp. 5159–5172, 2024.
- [14] M. Ahmad and M. Mazzara, "Scsnet: Sharpened cosine similarity-based neural network for hyperspectral image classification," *IEEE Geoscience and Remote Sensing Letters*, vol. 21, pp. 1–4, 2024.
- [15] M. H. F. Butt, J. P. Li, M. Ahmad, and M. A. F. Butt, "Graph-infused hybrid vision transformer: Advancing geoi for enhanced land cover classification," *International Journal of Applied Earth Observation and Geoinformation*, vol. 129, p. 103773, 2024.
- [16] X. Chen, S.-I. Kamata, and W. Zhou, "Hyperspectral image classification based on multi-stage vision transformer with stacked samples," in *TENCON 2021-2021 IEEE Region 10 Conference (TENCON)*. IEEE, 2021, pp. 441–446.
- [17] J. Zou, W. He, and H. Zhang, "Lessformer: Local-enhanced spectral-spatial transformer for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–16, 2022.
- [18] M. Li, Y. Fu, and Y. Zhang, "Spatial-spectral transformer for hyperspectral image denoising," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 1, 2023, pp. 1368–1376.
- [19] P. Tang, M. Zhang, Z. Liu, and R. Song, "Double attention transformer for hyperspectral image classification," *IEEE Geoscience and Remote Sensing Letters*, vol. 20, pp. 1–5, 2023.
- [20] H. Yan, E. Zhang, J. Wang, C. Leng, A. Basu, and J. Peng, "Hybrid conv-vit network for hyperspectral image classification," *IEEE Geoscience and Remote Sensing Letters*, 2023.
- [21] M. Ahmad, U. Ghous, M. Usama, and M. Mazzara, "Waveformer: Spectral–spatial wavelet transformer for hyperspectral image classification," *IEEE Geoscience and Remote Sensing Letters*, vol. 21, pp. 1–5, 2024.
- [22] D. Hong, Z. Han, J. Yao, L. Gao, B. Zhang, A. Plaza, and J. Chanussot, "Spectralformer: Rethinking hyperspectral image classification with transformers," *IEEE Trans. Geosci. Remote Sens.*, vol. 60, pp. 1–15, 2022, doi: 10.1109/TGRS.2021.3130716.
- [23] A. Wang, S. Xing, Y. Zhao, H. Wu, and Y. Iwahori, "A hyperspectral image classification method based on adaptive spectral spatial kernel combined with improved vision transformer," *Remote Sensing*, vol. 14, no. 15, 2022.
- [24] L. Ji, Y. Shao, J. Liu, and L. Xiao, "Spatial-spectral hierarchical vision permutator for hyperspectral image classification," *European Journal of Remote Sensing*, vol. 56, no. 1, p. 2153747, 2023.
- [25] X. He, Y. Chen, and Z. Lin, "Spatial-spectral transformer for hyperspectral image classification," *Remote Sensing*, vol. 13, no. 3, 2021.
- [26] W. Wang, L. Liu, T. Zhang, J. Shen, J. Wang, and J. Li, "Hyperes2t: Efficient spatial–spectral transformer for the classification of hyperspectral remote sensing images," *International Journal of Applied Earth Observation and Geoinformation*, vol. 113, p. 103005, 2022.
- [27] W. He, W. Huang, S. Liao, Z. Xu, and J. Yan, "Csit: A multiscale vision transformer for hyperspectral image classification," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 9266–9277, 2022.
- [28] M. Li, Y. Fu, and Y. Zhang, "Spatial-spectral transformer for hyperspectral image denoising," 2022.
- [29] M. Ahmad, A. M. Khan, M. Mazzara, S. Distefano, M. Ali, and M. S. Sarfraz, "A fast and compact 3-d cnn for hyperspectral image classification," *IEEE Geoscience and Remote Sensing Letters*, 2020.
- [30] Z. Xiong, Y. Yuan, and Q. Wang, "Ai-net: Attention inception neural networks for hyperspectral image classification," in *IGARSS 2018-2018 IEEE International Geoscience and Remote Sensing Symposium*. IEEE, 2018, pp. 2647–2650.
- [31] X. Zhang, "Improved three-dimensional inception networks for hyperspectral remote sensing image classification," *IEEE Access*, vol. 11, pp. 32648–32658, 2023.
- [32] H. Firat, M. E. Asker, M. İ. Bayındır, and D. Hanbay, "Hybrid 3d/2d complete inception module and convolutional neural network for hyperspectral remote sensing image classification," *Neural Processing Letters*, vol. 55, no. 2, pp. 1087–1130, 2023.
- [33] X. Yang, Y. Ye, X. Li, R. Y. Lau, X. Zhang, and X. Huang, "Hyperspectral image classification with deep learning models," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 56, no. 9, pp. 5408–5423, 2018.
- [34] S. Ghaderizadeh, D. Abbasi-Moghadam, A. Sharifi, N. Zhao, and A. Tariq, "Hyperspectral image classification using a hybrid 3d-2d convolutional neural networks," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 14, pp. 7570–7588, 2021.
- [35] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," 2023.
- [36] X. Yang, W. Cao, Y. Lu, and Y. Zhou, "Hyperspectral image transformer classification networks," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2022.
- [37] L. Sun, G. Zhao, Y. Zheng, and Z. Wu, "Spectral–spatial feature tokenization transformer for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–14, 2022.
- [38] H. Guo and W. Liu, "S3l: Spectrum transformer for self-supervised learning in hyperspectral image classification," *Remote Sensing*, vol. 16, no. 6, 2024.
- [39] X. Zheng, H. Sun, X. Lu, and W. Xie, "Rotation-invariant attention network for hyperspectral image classification," *IEEE Transactions on Image Processing*, vol. 31, pp. 4251–4265, 2022.
- [40] J. Feng, Q. Wang, G. Zhang, X. Jia, and J. Yin, "Cat: Center attention transformer with stratified spatial–spectral token for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 62, pp. 1–15, 2024.
- [41] Z. Xue, Q. Xu, and M. Zhang, "Local transformer with spatial partition restore for hyperspectral image classification," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 15, pp. 4307–4325, 2022.
- [42] S. Jia, Y. Wang, S. Jiang, and R. He, "A center-masked transformer for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 62, pp. 1–16, 2024.