



Application of Convolutional Neural Network in the Segmentation and Classification of High-Resolution Remote Sensing Images

By Dr. E. Kesavulu Reddy

S. V. University

Abstract- Numerous convolution neural networks increase accuracy of classification for remote sensing scene images at the expense of the models' space and time sophistication. This causes the model to run slowly and prevents the realization of a trade-off among model accuracy and running time. The loss of deep characteristics as the network gets deeper makes it impossible to retrieve the key aspects with a sample double branching structure, which is bad for classifying remote sensing scene photos. We suggest a dual branch inter feature dense fusion-based lightweight convolutional neural network to address this issue (BMDF-LCNN). In order to prevent the loss of shallow data due to network development, the network model can fully extricate the data from the current layer through 3 x 3 depthwise separable method is structured and 1 x 1 standard pooling layers, identity sections, and fusion with the extracted features out from preceding stage through 1 x 1 standard pooling layer.

Keywords: remote sensing, convolutional neural network, standard convolution, feature extraction.

GJCST-D Classification: DDC Code: 621.3678 LCC Code: G70.4



Strictly as per the compliance and regulations of:



Application of Convolutional Neural Network in the Segmentation and Classification of High-Resolution Remote Sensing Images

Dr. E. Kesavulu Rreddy

Abstract- Numerous convolution neural networks increase accuracy of classification for remote sensing scene images at the expense of the models' space and time sophistication. This causes the model to run slowly and prevents the realization of a trade-off among model accuracy and running time. The loss of deep characteristics as the network gets deeper makes it impossible to retrieve the key aspects with a sample double branching structure, which is bad for classifying remote sensing scene photos. We suggest a dual branch inter feature dense fusion-based lightweight convolutional neural network to address this issue (BMDF-LCNN). In order to prevent the loss of shallow data due to network development, the network model can fully extricate the data from the current layer through 3 x 3 depthwise separable method is structured and 1 x 1 standard pooling layers, identity sections, and fusion with the extracted features out from preceding stage through 1 x 1 standard pooling layer. Additionally, we suggest a down sampling structure that makes use of the pooled branching to down sample and the convolution branching to make up for the pooled features in order to extract the shallow characteristics of the network more effectively. Four open and difficult remote sensing images scene data sets were used for the experiments. The experimental findings show that the suggested method realises the trade-off among model accuracy and system running speed and has better classification accuracy and reduced model complexity than various state-of-the-art classification techniques.

Keywords: remote sensing, convolutional neural network, standard convolution, feature extraction.

1. INTRODUCTION

In applications including urban development, land-use planning, infrastructure construction management, natural disasters, and crisis management, urban land-use classification is crucial [1]. The rate of change in land usage increases with the nation's rate of growth. Costly, labor-intensive, and time-consuming are land-use surveys [2]. China conducts a national land-use survey every ten years. High-resolution remote sensing processing technologies are being developed, which may assist planners in quickly and affordably gathering comprehensive land-cover data [3]. Deep convolutional neural networks (DCNNs), for instance, might fully achieve the classification of urban land-use by

automatically extracting species-specific information from remote sensing photos. According to current criteria for land-use classification, one typical class may include more than one type of item. Each might also contain various objects that adhere to various standards. The Land-Use Standard of the 2nd and 3rd National Land-Use Resource Surveys, for instance, has various contents. Convolutional neural networks (CNN) trying to identify high-resolution remote sensing images are faced with significant difficulties by the complex spatial and textural patterns in one class [4]. Early FCN-based models had a limited ability to reconstruct spatial information despite acquiring rich contextual data and suffered from loss of high-frequency details, blurring boundaries, and difficulty identifying features. A skip connection was introduced to the networks to address this issue. By combining the multi-layer feature maps from the encoder with the decoder structure for incremental up sampling, Ronneberger et al U-Net.'s.

Architecture created high-resolution feature maps [5]. The classification effects of object boundaries are improved by the merging of high- and low-level semantic information. Later, Yu and Koltun added atrous convolution to fully convolutional networks (FCN), which could maintain the resolution of a featured image, expand the receptive field to capture multi-scale context information, and boost the precision of semantic segmentation using spatial information in the images [6]. Spatial Pyramid Pooling (SPP) [7] has been widely used to better capture information about the global context.

To take advantage of the potential of global context information, Zhao et al. used a pyramid pooling module to aggregate the context of several regions [8]. In order to gather multi-scale information, Chen et al. implemented pyramid-shaped atrous pooling in spatial dimensions [9] and piled up atrous convolution [10] with various atrous in cascade or in parallel. The resolution in the scale axis dimension was insufficient for Atrous Spatial Pyramid Pooling (ASPP) [9] to precisely extract target features from remote sensing images, therefore it still had certain drawbacks (RSIs). In order to effectively identify complicated situations while maintaining the model's size, Yang et al. introduced densely-connected Atrous Spatial Pyramid Pooling (DenseASPP) [12], which was able to cover a broader scale of the feature

Author: Assistant Professor Department of Computer Science S. V. University, Tirupati, Andhra Pradesh India.
e-mail: ekreddy@svuniversity.edu.in

map and acquire more intensive receptive field information. A labor-intensive foundation, including well-tagged remote sensing image labels for the most recent urban land cover types under distinct categorization standards, must be built in order to address the inherent difficulties in the present classification methods used to classify urban land use. Combining algorithms to produce higher-level semantic class images is another effective way to replace the original photos in labor-intensive tasks. We proposed a double-layer deep convolutional neural network called DUA-Net that mainly combined two networks with different advantages, U-Net and DenseASPP, into a parallel structure in order to take into account the characteristics of urban land-use types, which contain multiple elements in one type. The methodology utilised in this study can produce a larger, continuous block of land use classification for urban areas. It can considerably cut down on operation durations and manual interactions when using the image of this classification result as the input for artificial fine classification, which can increase efficiency. Additionally, with the aid of vector data, we can fully utilise the same standard to categorise the photographs taken at various points in time in order to study the changes in land types at various times.

High resolution remote sensing images are being used in a variety of applications, including the classification of remote sensing scenes [1], hyperspectral image classification [2], change detection [3, 4], geographic image classification, and the classification of land use [5, etc. However, image categorization presents significant challenges due to the intricate spatial patterns and geographical structure of remote sensing images. Therefore, it is crucial to effectively comprehend the semantic information of remote sensing photographs [6]. The goal of this study is to identify a straightforward and effective lightweight network model that is capable of quickly and reliably classifying remote sensing scene photos. Researchers have suggested a variety of techniques for efficiently extracting visual information. To begin with, manually created feature descriptors including colour histograms, texture descriptors, local binary mode, GIST, directional gradient histograms, bag-of-visual words (BOW), etc. were used to extract picture features. Researchers then proposed some unsupervised feature learning techniques that can automatically extract shallow detail features from images in order to address the drawbacks of the method of manually extracting features. These techniques include principal component analysis (PCA), sparse coding, autoencoder, Latent Dirichlet allocation, and probabilistic latent semantic analysis. For the extraction of shallow picture information, the two feature extraction techniques mentioned above work quite well. However, the extraction of high-level features from images using these techniques is challenging, which restricts the development of classification accuracy.

Researchers have proposed convolutional neural networks, which have the ability to automatically extract significantly discriminative features from images, as a way to get around the limitations of existing methods. Since then, the model based on convolution neural networks has replaced other techniques as the industry standard for classifying remote sensing scene images. A lightweight convolution neural network may now achieve a balance between the speed of model operation and the precision of model classification thanks to advancements in convolution neural networks. Lightweight networks have currently been used for a variety of applications, including target recognition, image segmentation, and classification. The fire module, which separates the initial basic convolution layer into an extrusion layer and an expansion layer, was proposed by SqueezeNet. The extension layer is made up of a set of continuous 1×1 convolution and 3×3 convolution channels, whereas the extruded layer is made up of a continuous set of 1×1 convolution channels. The Google team's MobileNet has three iterations: V1, V2, and V3. In order to divide the regular convolution into depthwise convolution and 1×1 convolution, MobileNetV1 employs depthwise separable convolution. This significantly decreases the number of network parameters and, to some extent, increases accuracy. An inverse residual module and a linear bottleneck structure were presented by MobileNetV2. The convolution of 1×1 for ascending dimension, 3×3 depthwise separable convolution for feature extraction, and 1×1 convolution for dimension reduction were all applied to this bottleneck structure in that order. With the addition of the SE module and the use of neural structure search, MobileNet V3 examines the network's setup and parameters [10]. An extremely effective convolution neural network architecture called ShuffleNet was created for mobile devices with constrained processing resources. Compared to some sophisticated ones with comparable accuracy, the design only requires two operations—group convolution and channel mixing—which significantly lowers the computation time.

II. RELATED WORKS

Remote sensing picture databases are being produced in greater numbers. These datasets use a variety of land cover and land use categories, and Castillo-Navarro et al. [14] have developed datasets that cover various scenes to increase surface coverage. Additionally, the labels that have been applied to the datasets vary [15]. For instance, BigEarthNet [17] and SEN12MS [16] both give image-level labels and pixel-level labels, respectively, and both datasets with varying scene categories can only be used for particular semantic segmentation applications. For instance, LULC has hundreds of fine-grained semantic classifications that may be further broken down into

categories like highways, buildings, cars, the countryside, cities, etc.

The circumstances that can show the relationship between the content of interest and its surroundings are rarely taken into account, and many datasets simply ignore the relationships within and across semantic classes [18]. Rich and detailed geometric features, texture information, and geospatial data are all present in high-resolution RSIs [19]. For land-use classification, the features extracted from these images can be interpreted with high accuracy. Pixel-based image analysis, object-based image analysis, and pixel-level semantic segmentation have all been used to classify the land use of RSIs [20]. Low-resolution remote sensing photos have historically been classified primarily using spectral data from remote sensing photographs.

Because the spectral features of pixels, which lack textural features and structural data, are unable to fully capture the characteristics of land-use kinds, the classification results for complicated land-use types, such as residential land and wasteland, are frequently less than optimal [21]. Similar pixels in different land-use types on residential and industrial land may exist. Some strategies, like Transfer Learning [22], Active Learning [24], and others, have been developed with the goal of increasing the size and enhancing the effectiveness of training datasets. Ammour et al. [25] merged two asymmetric networks for data domain adaption and classification, projected the two networks to the same feature space, and performed post-training for the weight coefficient adjustment method of the two networks. They employed a pretraining network for feature extraction.

Migration tests were conducted by Zhou et al. [26] using data from the same sensor at various dates. Additionally, they created a very difficult migration experiment that tested the efficacy of feature extraction and migration structure and was performed on hyperspectral remote sensing data from various viewpoints. The object-oriented classification approach [27] takes into account the correlation information between pixels and the internal texture features of ground objects while leveraging the spectrum information of RSIs [28] to make up for the inadequacies of conventional pixel-based classification methods. However, feature descriptions are generally incomplete, and the data collected is frequently insufficient to assist the characterization and identification of ground objects.

Deep learning overcomes the limitations of artificial features, directs object categorization, and achieves pixel-level land-use classification of RSIs by mastering the shape and texture aspects of various objects. Deep learning has been used extensively in RSIs for land-use classification. To automatically train

the representative and discriminative features in a hierarchical way for land-use scene classification, deep filter banks were proposed to integrate multicolumn stacked denoising sparse autoencoders (SDSAE) with Fisher vectors (FV) [29]. A land-use classification framework for photographs (LUCFP) was presented by Xu et al., and it was effectively used to automate the verification of land surveys in China [30]. Adaptive hierarchical image segmentation improvement, multilevel extraction of features, and multiscale supervised deep - learning models were integrated to accurately produce detailed maps for disparate urban areas from the fusion of the UHSR ortho mosaic and digital elevation model, taking into account the high-level details in an ultrahigh-spatial-resolution (UHSR) unmanned aerial vehicle (UAV) dataset (DSM). Excellent potential was shown by this framework for the thorough mapping of varied urban areas [31]. Another cutting-edge hybrid approach is multi-temporal relearning using convolutional long short-term memory (LSTM) models. It integrates post-classification relearning with locational semantic segmentation and is effective at categorising complex LULC maps with multitemporal VHR pictures [32].

III. METHODOLOGY

a) *Proposed Architecture*

Figure 1 depicts the model's overall structure, which is broken down into nine sections. We suggest a feature extraction structure for the network's shallow layers in the first and second groups. The maximum pool layer is used for down sampling in the third group, where standard convolution and depth-wise separable convolution are combined. This reduces the spatial dimensions of the input images and lowers the danger of overfitting from irrelevant features. The majority of representative features from remote sensing image are extracted by groups 4 through 8. For the purpose of extracting richer feature information, Groups 4 through 7 use the proposed dual branch multi-level feature intense fusion method.

To extract deep-level features from Group 8, we sequentially applied 1 x 1 standard convolution, 3 x 3 standard convolution, and 3 x 3 depth wise separable convolution. The multilevel characteristics are fully exchanged and fused on the basis of double branch fusion, which not only increases classification accuracy but also significantly speeds up the network and achieves a balance between accuracy and speed. The number of convolution channels in Groups 5 and 8 is also increased to 256 and 512, respectively, in order to extract more features. The feature information generated by the final fusion is used to calculate the likelihood of each scene category, and Group 9 is used for classification. Each layer in deep feature extraction structures from Group 4 to Group 7 may fully extract the

data of the current layer through three branches, including Identity, 1×1 standard convolution, and 3×3 depth wise separable convolution. Additionally, the shallow information loss caused by network deepening can be successfully avoided by merging the features retrieved by 1×1 standard convolution with each prior layer. Batch normalization (BN) can speed up training and use a greater learning rate while reducing the network's reliance on parameter initialization. Additionally, there are far less remote sensing photos available for training compared to the natural image data collection.

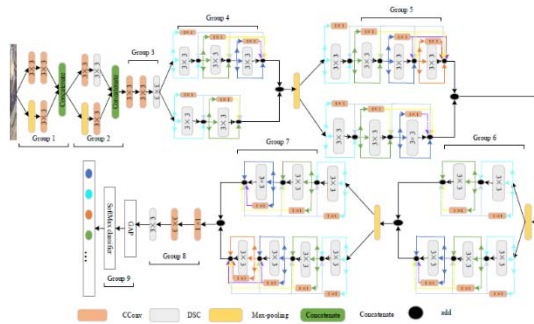


Figure 1: Proposed Architecture

In order to gather spatial information, boundary information, multi-scale contextual information, and global contextual information, our proposed model utilised parallel modules. As a result, it was able to reduce border ambiguity and class imbalance, address the inaccurate, fragmented single element classification in urban land-use semantic segmentation, and increase urban land-use classification accuracy. This section showed the DUA-suggested Net's architecture for classifying urban land use. The essential components of the suggested design, including the U-Net module, DenseASPP module, and Channel Attention Fusion module, were then thoroughly explained. In this study, U-Net and DenseASPP, two different DCNNs, were used to build the distributed system of DUA-Net, which fully utilised the various benefits of these two types of networks in the semantic segmentation of RSIs. The suggested framework has three components, as seen in Figure 1: a backbone network, a parallel extracting features module, and a feature fusion module. First, the VGG16 network is introduced as the foundation for extracting the features in U-Net and DenseASPP. Second, we use the U-Net module and DenseASPP module to simultaneously collect various semantic pieces of information due to the complexity of land-use type, structure, and geographic distribution of abnormality.

For more specifics, the DenseASPP module aggregates semantic information at various scales to capture multi-scale contextual information and global contextual information, and the U-Net module fuses

high-level and low-level semantic information to improve the extraction of spatial and boundary information. Then, to address the issue of improper segmentation caused by comparable characteristics of similar categories, the feature maps output by the U-Net module and DenseASPP module were fused in the channel dimension through the attention mechanism in the Channel Attention Fusion module. The segmentation results were then produced by convolution with a convolution kernel size of 1×1 after the feature vectors had been mapped to the necessary number of classes. The DenseASPP module was introduced as the feature extractor in order to gather multi-scale contextual information as well as global contextual information in RSIs. In order to achieve integration at various levels with various dilation rates, DenseASPP employs the concept of dense connection and arranges numerous convolution layers in a cascade manner. Without significantly growing the model size, this organising approach not only covers a broader scale, but also covers it intensively. To gather semantic information from various scales, this work specifically exploited dense connections to send the output of each atrous convolution layer to all previously unvisited atrous convolution layers.

Additionally, the atrous convolution dilation rate at each layer increased layer by layer, enlarging the receptive field while maintaining the same level of feature map resolution. The layer with the lowest dilation rate among them was positioned in the lower layer, and the layer with the highest dilation rate was positioned in the upper layer. The outputted feature map from the multi-scale convolution process was the last step. The connection between feature channels is typically ignored by traditional techniques, and thus exhibit low sensitivity to critical information characteristics throughout the fusion process.

We used the channel attention method to successfully combine the feature maps from the U-Net module and the DenseASPP module. This fusion module achieved the automatic selection and weight assignment of attention regions, then increased output feature quality by using SENet to learn the correlation between various feature channels (and to boost the extraction of significant features). In particular, its primary operations were concatenation, squeeze, and excitation. The shallow features of the network are intended to be extracted by the first and second sets of down sampling structures. The impact of down sampling on network performance during the shallow feature extraction phase is significant. Down sampling is the process of scaling down the complex feature map to maintain the image's primary features while reducing the spatial size of the image. Having a greater number of pooling layers is one of the primary techniques for down sampling in deep convolution neural networks.

IV. RESULTS AND DISCUSSION

Three of the down sampling techniques indicated in Section IIB are utilised in the first and second levels of the network to verify the effectiveness of our suggested down sampling techniques. Two datasets, UC and RSSCN, were employed for the experiments, and the OA and Kappa were used as evaluation metrics. According to Figure 2, the first and third convolution steps for the Conv-Downsampling (CD) are 1, while the second and fourth convolution steps are 2. The convolution kernels for the pooling down sampling (Max pooling-Downsampling, MD) are all 3 x 3, with convolution steps of 1 x 1. The pooling step size and maximum pooling size are both 2. On the two datasets, pooling down sampling had poorer classification accuracy and Kappa values than convolution down sampling. Convolution down sampling in deep networks produces superior non-linear performance than pooled down sampling, which is the reason. On the 80/20UC and 50/50RSSCN datasets, the suggested down sampling methods have classification accuracy scores of 99.53% and 97.86%, respectively, and Kappa values of 99.50% and 97.50%, which are greater than those of the other two down sampling methods.

This demonstrates once again how much more accurately the multi-level features dense fusion technique can identify remote sensing scene photos. In this section, three types of visualisation, including grad cam, t-distribution random neighbour embedding (T-SNE), and randomly picked and tested are explained and examined in order to more clearly demonstrate the effectiveness of the suggested method. Through a visual thermal map, the grad cam presents the retrieved features in order of significance. The most comprehensive spatial and semantic information is found in the final layer of a convolution neural network.

Grad Cam creates an attention map to highlight key portions of an image by fully utilizing the features of the last layer of convolution. In this experiment, some remote sensing scene photos from the RSSCN collection of "Industries," "Fields," "Residences," "Grass," and "Forests" are randomly chosen. Figure 2 compares the thermal diagram visualization outcomes of the enhanced BMDf-LCNN approach with the baseline LCNN-BFF method.

Figure 2 shows that, for "Industries" scenarios, the LCNN-BFF approach transfers the attention to the highway rather than accurately focusing on the factory region, whereas the proposed BMDf-LCNN method accurately focuses on the industrial area. While the BMDf-LCNN approach is well focused on the target region, the LCNN-BFF model's focused areas for the "Fields" and "Grass" scenarios both showed a partial deviation, ignoring the similar surrounding targets and searching with few objects. Additionally, the LCNN-BFF

method's restricted coverage and inability to fully extract the target for scenario regions like "Residence" and "Forests" has an impact on the classification accuracy. However, in these cases, the suggested BMDf-LCNN approach can get a comprehensive region of interest. Next, we use t-distribution random neighbor embedding to illustrate the classification results on the UC (8/2) and RSSCN (5/5) datasets (T-SNE). High-latitude characteristics are mapped by T-SNE to two- or three-dimensional space for visualization, which is a very effective way to assess the classification effect of the model.

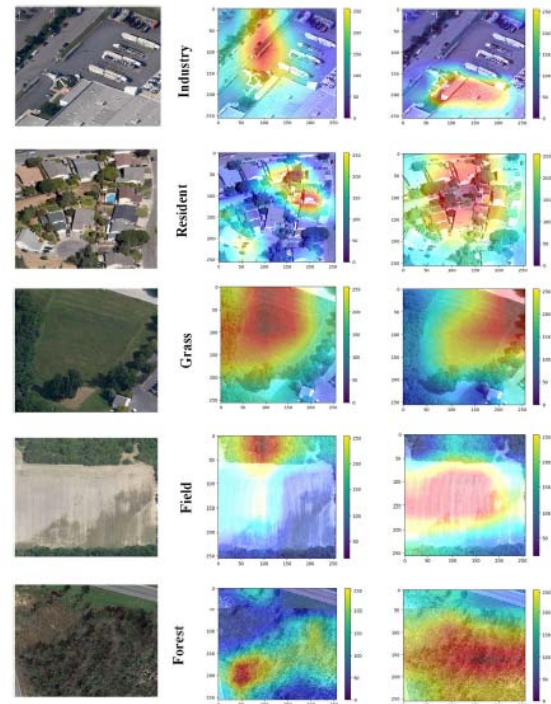


Figure 2: Results of Convolutional Neural Network Segmentation and Classification

V. CONCLUSION

A lightweight network based on the dense fusion of dual-branch, multi-level features is proposed for the categorization of remote sensing scene photos. A fresh down sampling technique was also developed to gather more accurate feature data. The information of the current layer can be fully extracted and fused with the features extracted by 1x1 standard convolution in the previous layer using the three branches of 3 3 depth wise separable convolution, 1 x 1 standard convolution, and identity in the network. This effectively realizes the information interaction between different levels of features and improves the classification performance and computational speed of the model. The suggested model still requires development. Due to the generation of certain redundant data during multi-level feature heavy fusion, the computational complexity rises. Future

research should discover a technique that can fuse data selectively, limit the production of redundant data, and further develop a lightweight model that combines speed and accuracy.

REFERENCES RÉFÉRENCES REFERENCIAS

- Hu, F.; Xia, G.S.; Hu, J.; Zhang, L. Transferring deep convolutional neural networks for the scene classification of high-resolution remote sensing imagery. *Remote Sens.* 2015, 7, 14680–14707.
- Liu, Q.; Zhou, F.; Hang, R.; Yuan, X. Bidirectional-convolutional LSTM based spectral-spatial feature learning for hyperspectral image classification. *Remote Sens.* 2017, 9, 1330.
- Lu, X.; Yuan, Y.; Zheng, X. Joint dictionary learning for multispectral change detection. *IEEE Trans. Cybern.* 2017, 47, 884–897.
- Li, Y.; Peng, C.; Chen, Y.; Jiao, L.; Zhou, L.; Shang, R. A deep learning method for change detection in synthetic aperture radar images. *IEEE Trans. Geosci. Remote Sens.* 2019, 57, 5751–5763.
- Peng, C.; Li, Y.; Jiao, L.; Chen, Y.; Shang, R. Densely based multiscale and multi-modal fully convolutional networks for high-resolution remote-sensing image semantic segmentation. *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.* 2019, 12, 2612–2626.
- Ghamisi, P.; Maggiori, E.; Li, S.; Souza, R.; Tarabla, Y.; Moser, G.; Chen, Y. New frontiers in spectral-spatial hyperspectral image classification: The latest advances based on mathematical morphology, Markov random fields, segmentation, sparse representation, and deep learning. *IEEE Geosci. Remote Sens. Mag.* 2018, 6, 10–43.
- Swain, M.J.; Ballard, D.H. Color indexing. *Int. J. Comput. Vis.* 1991, 7, 11–32.
- Ojala, T.; Pietikainen, M.; Maenpää, T. Multiresolution gray-scale and rotation invariant texture classification with local binary patterns. *IEEE Trans. Pattern Anal. Mach. Intell.* 2002, 24, 971–987.
- Song, C.; Yang, F.; Li, P. Rotation invariant texture measured by local binary pattern for remote sensing image classification. In *Proceedings of the 2nd International Workshop on Education Technology and Computer Science, ETCS, Wuhan, China, 6–7 March 2010; Volume 3, pp. 3–6.*
- Oliva, A.; Antonio, T. Modeling the shape of the scene: A holistic representation of the spatial envelope. *Int. J. Comput. Vis.* 2001, 42, 145–175.
- Dalal, N.; Triggs, B. Histograms of oriented gradients for human detection. In *Proceedings of the CVPR, IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Boston, MA, USA, 7–12 June 2015; pp. 886–893.*
- Sivic, J.; Zisserman, A. Video Google: A text retrieval approach to object matching in videos. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV), Nice, France, 13–16 October 2003; p. 1470.*
- Zhou, Y.; Liu, X.; Zhao, J.; Ma, D.; Yao, R.; Liu, B.; Zheng, Y. Remote sensing scene classification based on rotation invariant feature learning and joint decision making. *EURASIP J. Image Video Process.* 2019, 2019, 1–11.
- Wang, C.; Lin, W.; Tang, P. Multiple resolution block feature for remote-sensing scene classification. *Int. J. Remote Sens.* 2019, 40, 6884–6904.
- Wilhelm, T.; Koßmann, D. Land Cover Classification from a Mapping Perspective: Pixelwise Supervision in the Deep Learning Era. In *Proceedings of the 2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS, Brussels, Belgium, 11–16 July 2021; pp. 2496–2499.*
- Schmitt, M.; Hughes, L.H.; Qiu, C.; Zhu, X.X. SEN12MS—A Curated Dataset of Georeferenced Multi-Spectral Sentinel-1/2 Imagery for Deep Learning and Data Fusion. *arXiv* 2019, arXiv:1906.07789.
- Sumbul, G.; Charfuelan, M.; Demir, B.; Markl, V. Bigearthnet: A Large-Scale Benchmark Archive for Remote Sensing Image Understanding. In *Proceedings of the IGARSS 2019—2019 IEEE International Geoscience and Remote Sensing Symposium, Yokohama, Japan, 28 July–2 August 2019; pp. 5901–5904.*
- Long, Y.; Xia, G.S.; Li, S.; Yang, W.; Yang, M.Y.; Zhu, X.X.; Zhang, L.; Li, D. On Creating Benchmark Dataset for Aerial Image Interpretation: Reviews, Guidances, and Million-AID. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 2021, 14, 4205–4230.
- Richards, J.A.; Jia, X. *Remote Sensing Digital Image Analysis: An Introduction*, 4th ed.; Springer: Berlin, Germany, 2005; ISBN 978-3-662-02464-5.
- Neupane, B.; Horanont, T.; Aryal, J. Deep Learning-Based Semantic Segmentation of Urban Features in Satellite Images: A Review and Meta-Analysis. *Remote Sens.* 2021, 13, 808.
- Pan, X.; Zhao, J. High-Resolution Remote Sensing Image Classification Method Based on Convolutional Neural Network and Restricted Conditional Random Field. *Remote Sens.* 2018, 10, 920.
- Ma, L.; Liu, Y.; Zhang, X.; Ye, Y.; Yin, G.; Johnson, B.A. Deep Learning in Remote Sensing Applications: A Meta-Analysis and Review. *ISPRS J. Photogramm. Remote Sens.* 2019, 152, 166–177.
- Lyu, H.; Lu, H.; Mou, L.; Li, W.; Wright, J.; Li, X.; Li, X.; Zhu, X.; Wang, J.; Yu, L.; et al. Long-Term Annual Mapping of Four Cities on Different Continents by Applying a Deep Information

- Learning Method to Landsat Data. *Remote Sens.* 2018, 10, 471.
24. Liu, P.; Zhang, H.; Eom, K.B. Active Deep Learning for Classification of Hyperspectral Images. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 2017, 10, 712–724.
 25. Ammour, N.; Bashmal, L.; Bazi, Y.; Al Rahhal, M.M.; Zuair, M. Asymmetric Adaptation of Deep Features for Cross-Domain Classification in Remote Sensing Imagery. *IEEE Geosci. Remote Sens. Lett.* 2018, 15, 597–601.
 26. Zhou, X.; Prasad, S. Deep Feature Alignment Neural Networks for Domain Adaptation of Hyperspectral Data. *IEEE Trans. Geosci. Remote Sens.* 2018, 56, 5863–5872.
 27. Walter, V. Object-Based Classification of Remote Sensing Data for Change Detection. *ISPRS J. Photogramm. Remote Sens.* 2004, 58, 225–238.
 28. Chen, W.; Li, X.; He, H.; Wang, L. A Review of Fine-Scale Land Use and Land Cover Classification in Open-Pit Mining Areas by Remote Sensing Techniques. *Remote Sens.* 2018, 10, 15.
 29. Wu, H.; Liu, B.; Su, W.; Zhang, W.; Sun, J. Deep Filter Banks for Land-Use Scene Classification. *IEEE Geosci. Remote Sens. Lett.* 2016, 13, 1895–1899.
 30. Xu, S.; Zhang, S.; Zeng, J.; Li, T.; Guo, Q.; Jin, S. A Framework for Land Use Scenes Classification Based on Landscape Photos. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 2020, 13, 6124–6141.
 31. Gibril, M.B.A.; Kalantar, B.; Al-Ruzouq, R.; Ueda, N.; Saeidi, V.; Shanableh, A.; Mansor, S.; Shafri, H.Z.M. Mapping Heterogeneous Urban Landscapes from the Fusion of Digital Surface Model and Unmanned Aerial Vehicle-Based Images Using Adaptive Multiscale Image Segmentation and Classification. *Remote Sens.* 2020, 12, 1081.
 32. Zhu, Y.; Geis, C.; So, E.; Jin, Y. Multitemporal Relearning with Convolutional LSTM Models for Land Use Classification. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 2021, 14, 3251–3265.

