



Performance Improvements Using Deep Learning Based Object-Identification

[Shouket Abdulrahman Ahmed](#)^{*1}, [Hazry Desa](#)¹, [Abadal-Salam T. Hussain](#)¹²

¹ Centre of Excellence for Unmanned Aerial Systems (COEUAS), Universiti Malaysia Perlis, Malaysia.

² Department of Medical Instrumentation Engineering Techniques, Al-Kitab University, Iraq.

*Corresponding Author: shouketunimap@gmail.com

Citation: Ahmed S., Desa H., Hussain A. Performance Improvements Using Deep Learning Based Object-Identification. Al-Kitab Journal for Pure Sciences (2021); 6(1): 1-13. DOI: <https://doi.org/10.32441/kjps.06.01.p1>

Keyword

Deep Learning, Artificial Neural Networks, Object-Identification.

Article History

Received	15 Jan. 2022
Accepted	12 Feb. 2022
Available online	10 April 2022

©2021. Al-Kitab University. THIS IS AN OPEN ACCESS ARTICLE UNDER THE CC BY LICENSE
<http://creativecommons.org/licenses/by/4.0/>



Abstract:

Deep Learning incorporates numerous hidden layers and more rooted combinations that average Artificial Neural Networks (ANN), to produce more refined and better performing autonomy in learning algorithms. An incredible volume of literature details and improves upon Deep Learning related methods and their improvement through the years as well as their suitability in uses. Nevertheless, the primary focus of the literature review is not the enlistment of these techniques hence a concise overview will be provided over the mechanisms before delving into the intended applications.

Keywords: Deep Learning, Artificial Neural Networks, Object-Identification.

تحسينات الأداء باستخدام التعرف على المكونات القائمة على التعلم العميق

*، [جزري ديسي¹](#)، [عبد السلام طه حسين¹²](#)، [شوكت عبد الرحمن احمد¹](#)

المركز الاستشاري للطائرات المسيرة، جامعة بيرليس الماليزية، بيرليس، ماليزيا¹
قسم هندسة تقنيات الأجهزة الطبية، الكلية التقنية الهندسية، جامعة الكتاب، كركوك، العراق²

[*shouketunimap@gmail.com](mailto:shouketunimap@gmail.com)

الخلاصة:

يشتمل التعلم العميق على العديد من الطبقات المخفية ومجموعات أكثر تجذرًا من الشبكات العصبية الاصطناعية لإنتاج استقلالية أكثر دقة وأفضل أداءً في خوارزميات التعلم. حجم لا يصدق من التوصيلات والتفاصيل وتحسن أساليب التعلم العميق ذات الصلة وتحسينها خلال عملية التعلم وكذلك مدى ملاءمتها في الاستخدامات المختلفة. ومع ذلك، فإن التركيز الأساسي لمراجعة التوصيلات ليس هو تجنيد هذه التقنيات ومن ثم سيتم تقديم نظرة عامة موجزة على الآليات قبل الخوض في التطبيقات المقصودة.

الكلمات المفتاحية: التعلم العميق، الشبكات العصبية الاصطناعية، التعرف على الأشياء.

1. INTRODUCTION:

The performance of DNNs (Deep Neural Networks) in operations concerning the processing of images improved considerably in the past decade because of the abundance of samples (labeled examples) and increased computer functionality. Deep Neural Networks have shown promising results when utilized for information handling operations, but there is yet much unexplored regarding the capabilities and constraints that this approach may have. For this reason, many studies were devised to more accurately detail the perspective contributions that deep learning may provide in the field of image segmentation.

2. Deep Neural Networks Overview:

Deep Neural Networks (DNNs) is a concept of neural-networks that are made up of “neurons”, referred to as units, and contain a specific incitation (or activation) and specifications (or parameters) functioning as processes that convert input data such as UAV-imagery to the desired output, which in this case is scenario-based maps, simultaneously gaining a greater level of knowledge [1, 2]. This dynamic knowledge enhancement is carried in layers called the “hidden layers” that lie between the input layer and output layer [3]. Deep

Neural Networks exist as the most simplified Deep Learning approach, as such containing at the very least two hidden layers. They were introduced in the mid-20th century and the ideology of this concept was derived from AI (Artificial Intelligence) produced under the natural neural networks present in human biology. Deep Neural Networks have resurfaced in many important domains of science only after the improvements in computation technology and the rising abundance of sample image datasets. Due to the benefits gained in performance with the implementation of Deep Learning networks in image-processing operations, it has caught the eye of many researchers and scientists in the image segmentation field in the previous decade [4, 5].

The operation of a Deep Neural Network is identical to that of an Artificial Neural Network such that it can be instructed to evaluate specific input characteristics that are examined through various layers, and the required data is provided by a final output layer. However, there are many differences between the more conventional Artificial Neural Networks and Deep Neural Networks that need to be examined. One of the most prevalent pieces of literature in the scope of Deep Learning, details Deep Neural Networks as “Representation-Learning” algorithms with several bands of representation. The concept of Representation Learning is a prominent subject in Deep Learning as it provides the capability for the model to autonomously generate the required representations from raw data, commonly in text, image, and video formats.

Usually, Deep Neural Networks consist of activation instructions that are applied to the dense layers in the network as illustrated in **Figure 1** below. The purpose of these activation instructions is to evaluate the importance score of the input data and provided biases and accordingly activate (or not activate) a unit [6]. These instructions contain decision-making parameters that assist in studying inherent patterns [7] hence, this process is essential in allowing different units to interact and gain knowledge. Some Activation algorithms that are broadly implemented are Rectified Linear Unit (ReLU) and others based on its structure e.g., Leaky ReLU, Parametric Rectified Linear unit (PReLU), Exponential Linear Unit (ELU), max-out, linear, tanh, and sigmoid [8]. Out of these ReLU, a piecewise linearly structured algorithm is the most prevalent of the modern Deep Neural Network algorithms. Its popularity can be attributed to its lower processing-power requirements than other functions, effectively handles the vanishing-gradient issue, allows broader representation of information, and provides the capability to alter data structure. On the other hand, a novel approach for implementing activation algorithms has been introduced as Mish, an autonomously organized

and multi-parametric activation algorithm, which is providing encouraging results, as further studies are being implemented.

Besides the Activation instructions, other constituent operations of a layer present in a Deep Neural Networks model can be batch normalization, convolution and deconvolution, max pooling, dropout, encoding and decoding, and memory holding units among many others. Out of these, dropout and batch normalization will be discussed currently while others will be detailed in further sections of this literature. Due to the nature of Deep Learning models to inconsistently discard neurons and connections using a configured probability, dropout layers are introduced to provide normalization in this process. They assist the network in reducing overfitting through the removal of connections showing correlation, along with enhancing rationalization allowing for more refined and quicker learning processes [9, 10]. Similarly, batch normalization also assists in enhancing rationalization by operating as a catalyst and stabilizing the path of loss-gradient. It is often implemented on problems pertaining to covariance shift in land use maps [11]. The systemization of these various types of layers and their construction is a defining feature of this network.

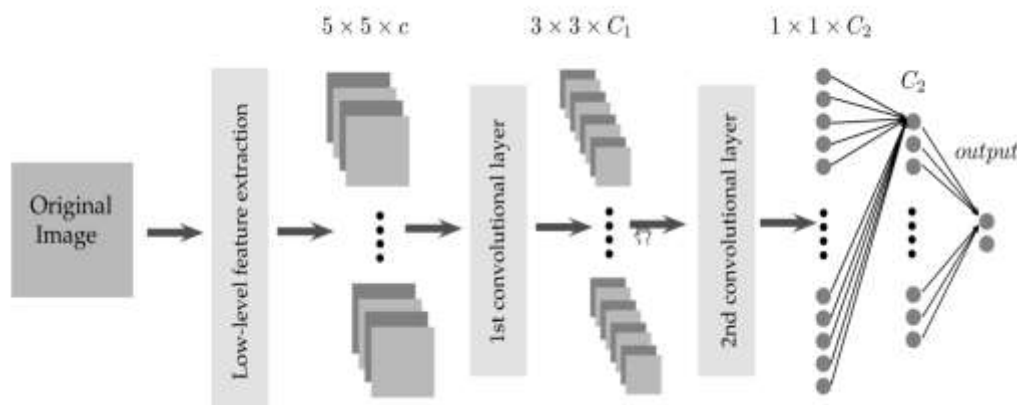


Figure 1: Deep learning model architecture

When assembling a network, its configuration requires some elementary pieces of information. The “optimizer” is one such piece that is applied to evaluate the rate of learning. A few reliable optimizer algorithms that are commonly used are momentum algorithm, Adam, Root Mean Squared Propagation (RMSprop), and Stochastic Gradient Descent (SGD). Choosing the most suitable optimizer for every application or architecture can assist greatly in refining the precision. The simplest of the above algorithms is SGD, in which the units are

combined and processed with the optimum-cost operation at a rate of a single sample-per-step. Another method, the Momentum algorithm aims to overcome the limitation caused by being bound to the local bottleneck by employing a transitory model. On the basis of gradient-based refinement strategies, RMSprop incorporates the momentum and AdaGrad (Adaptive Gradient) algorithms to apply an increasingly diminishing mean value of the gradients. The most prominent optimizer model, Adam, provides the capability to conjoin the adaptive-learning flow with the momentum algorithm, further examination of which can be found in [12]. These functions are essential in Deep Learning models as when used with a suitable loss algorithm, they can greatly enhance the precision of the network.

In the scope of optimization, the algorithm responsible for examining the network is the cost function, also commonly referred to as an objective function or loss function. The purpose of the loss function is to evaluate the capability of a network to provide a single-scalar score from the test information. Due to this often the focus of enhancing the learning ability of a network is directed towards varying the networks specifications to reduce the cost function. This assists in novel approaches to be scored and then benchmarked against different unit operations [13]. Cost functions utilize numerical statistics to perform their computation. This representation is analogous to the context of the issue, for instance, if a model is responsible for performing object classification or regression-based operations. In the case of object classification, the statistical loss can be evaluated using various approaches such as Poisson, cross-entropy, KL divergence (Kullback- Leibler), etc. In the case of regression-based operations, approaches such as Mean Squared Logarithmic Error (MSLE), Mean Absolute Error (MAE), Mean-Squared Error (MSE), Mean Absolute Percentage Error (MAPE), etc., are often applied. A more thorough approach to objective functions can be found in the literature by [4].

To assess the capabilities of a Deep Neural Network several criteria have been applied [5], as scientists often consult these evaluations. In the context of object classification, the general parameter used is accuracy, or by extension sensitivity or recall, other metrics such as F-Score, an area under the ROC curve path, precision, Intersection over Union can also be utilized to evaluate the efficiency of a model. An additional parameter is the “Kappa Co-efficient” which has been discredited as an accurate metric in the publication by [6]. Furthermore, in the context of regression-based operations, the preferred parameters are MAE, MSE, Mean Relative Error (MRE), Correlation Co-efficient (r), Root-Mean-Squared (RMS) Error, and others. Evaluation of such qualities is necessary to benchmark the

performance of sample and acquired data and compare novel models with existing ones [7]. Despite regression-based evaluation not being as prevalently applied in the examination of image-processing, it is still essential as UAV-related operations are often dependent on both processes.

Various variants of network models have been introduced in the past decade to advance and refine Deep Neural Networks by applying numerous types of layers, activation and optimizer algorithms, cost functions, level of depth, etc. Albeit the rise in Deep Neural Network's prevalence in the computer vision community can most accurately be attributed to the abundance of openly accessible data available to train these models. The prominent criterion among data-scientists suggested that at the very least 5000 samples should be available for every data class [8]. However, the direction taken by modern research in Deep Neural Networks focuses on achieving similar levels of feature detection using fewer samples than those suggested. Due to fewer expenditures and manpower needed to acquire samples, many applications with special requirements can profit from this. As such, despite efforts being made at this front, modern research in the computer-vision field is also advancing towards accommodating approaches for data augmentation, self-sufficiency, and unsupervised data-learning techniques, and many others. More details on these approaches are provided at the end of this literature, as such the study by [9] can also be consulted.

3. Image Segmentation Algorithms

Unsupervised image segmentation algorithms have progressed to the stage they can produce segmentations that match human intuition to a great extent. It is about time for segmentations to be implemented in the recognition of objects. Unsupervised segmentation can clearly be used to aid in the cueing and refinement of many recognition algorithms; however, a major stumbling block is the fact that it is not yet clear how efficiently these segmentation schemes can perform objectively. Most segmentation algorithms include cursory evaluations that only show visual clues of the segmentation process and rely on the instincts of the reader for decisions. Considering the persistent absence of numerical findings, it is tedious to determine the best segmentation scheme that will produce beneficial performances and in what scenarios. While appealing to human intuition is convenient, objective findings on huge datasets are necessary if the method is to be employed in an automated system.

4. Convolutional and Recurrent Neural Networks

Deep Neural Networks can be constructed using various models, where the sophistication of the architecture can be dependent on the number and constitution of layers in the networks and the mathematical processing that is applied. The various types of Deep Learning models that are used commonly are Recurrent Neural Networks (RNN) as show in **Figure 2**, Convolutional Neural Networks (CNN), Deep Belief Networks (DBN) [9], and Generative Adversarial Networks, that has been newly introduced. These Convolutional Neural Networks and Recurrent Neural Networks have been the most prominently used in the context of supervised networks [1].

Modern research in the context of image-processing and object-identification operations is primarily implemented with Convolutional Neural Network models. They are prominently recognized in the field of computer vision despite having recently come into the spotlight. Despite having been predicted to provide great enhancements in image-classification operations, its promise only became recognized in 2012, when [2] developed an approach incorporating CNN that severely outperformed all others in an image-classification contest. This model is referred to as AlexNet consisting of 8 total layers. The starting 5 were entirely convolutional based, of which a few were accompanied with max-pooling layers, and the final 3 were fully connected layers, all of which used the ReLu activation algorithm [5]. This architecture is currently thought of as a basic Deep Learning model and its efficiency is attributed to the comprehensiveness of its layers.

Convolutional Neural Networks (comprehensively illustrated in **Figure 2**) primarily consists of three discrete systemized webs of layers, i.e., convolutional, fully connected, and pooling layers [2], and have several assigned metrics such as biases, weights, the number of units and layers, stride, activation instructions, learning-rate, size of filters, etc. [11]. As the data passes through every layer several filters and incorporated biases are applied to allow the network to develop a land-use map [13]. Convolution processes utilize the size of filters to explore the correlation between pixels from the input data [11]. Due to the regular arrangement of pixels in multi-band image segmentation image-data Convolutional Neural Networks, that were initially introduced to compute data in the format of multiple arrays. As such this model has garnered a reputation as one of the most prevalent Deep Neural Network architectures currently [1] and has contributed greatly to the enhancement of numerous UAV-related image-processing operations.

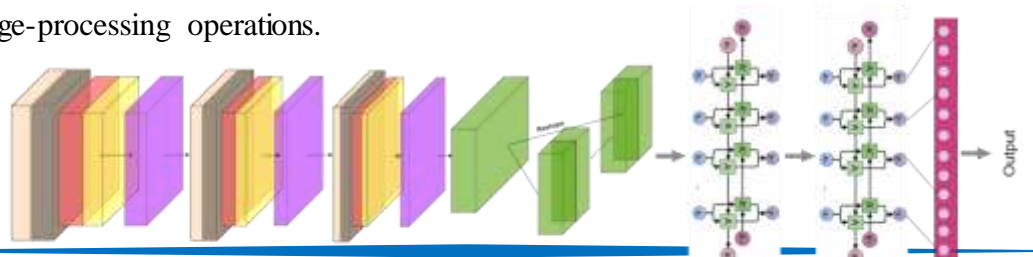


Figure 2: Convolutional Recurrent Neural Networks

Another architecture built upon the concept of Deep Learning is Recurrent Neural Networks, which implement a supervised-learning approach. Its advent in image segmentation applications is a recent development despite having been commonly utilized for numerous computer-vision-related operations. The Recurrent Neural Network architecture was introduced with the aim of handling discrete-sequence data examinations [1].

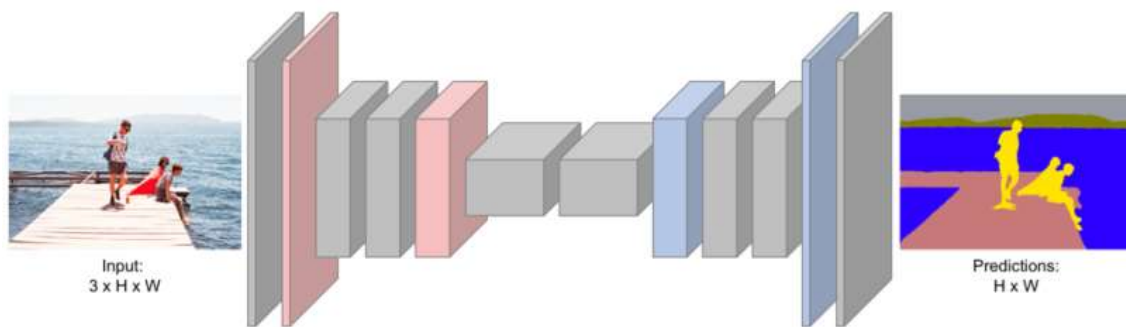


Figure 3: A Convolutional Neural Network-based model consisting of convolution layers and deconvolutional layers

The fundamental benefit received in using Recurrent Neural Networks is their ability to enhance their knowledge database through repeated examination of a specific setting or object, commonly pertaining to data of a time series-related format. A variant of Recurrent Neural Networks that is gaining importance and is being utilized with many applications is the Long Short-Term Memory model. They are useful in time-series-based operations as they overcome the diminishing gradient issue encountered with Recurrent Neural Networks. To accomplish these specific parameters and functions are included to smoothen the generated gradients [3]. Neurons found in an LSTM-based architecture consist of a cell along with an input and output gate and a forget gate. These gates exist to assist the unit in retaining and discarding information at discretely specified time-intervals.

Within the domain of image segmentation, Recurrent Neural Network architectures have been implemented in operations requiring handling and examination of time-series image-data intending to, for instance, generate land-use maps [4,5]. When working with pixel-related time-series image-data to classify categories of winter-time herbage density utilizing the SAR Sentinel-1 model, traditional Machine Learning networks paled in comparison to the performance demonstrated by the Recurrent Neural Network architecture. Furthermore, a modern model built for performing more precise herbage-density mapping been introduced,

incorporated a multi-domain Convolutional Neural Network to derive spatial-characteristics from UAV-acquired RGB-based image data and then processed through an attention focused Recurrent Neural Network to examine the flowing reliance withing multi-transient characteristics. The collective spatial temporal properties are utilized to deduce the class of vegetation. Examples such as these prove the promise of applying Recurrent Neural Networks to image segmentation information. Another promising model is the CNN-LSTM Architecture (illustrated in **Figure 4**). In this model, the convolution layers are applied to the image-data for feature-extraction which are then input into the LSTM-model. Very few implementations of this architecture can be found in literature as it caters to distinct applications such as, operations involving multi-temporal data.

Accordingly, besides Convolutional and Recurrent Neural Networks, many other models of Deep Neural Networks are being introduced for image-data processing applications. Among these Generative Adversarial Networks are considered to be the most innovatory approach in the context of un-supervised Deep Learning architectures. Generative Adversarial Networks consist of two models, namely “generative” and “discriminative”, that are designed to compete with each other. The generative model focuses on deriving the required properties from a specific data-format, for example image-data, while on the other hand the discriminative model discriminates between the reference, also known as real or ground-truth, data and the data provided by the generative network i.e., the fake data) [6,8]. Currently, GAN-based Deep Learning models that assist the process of image-data computation such as object-classification of image segmentation imagery and image to image transcription operations are providing promising outcomes [9].

Originally, the majority of DNN-based architecture types, such as RNN, CNN, and CNN-LSTM were introduced to handle a distinct problem. These models can be conveniently categorized according to image-classification operations, for instance setting-specific classification, object sensing, semantic-segmentation, pixel-based instance-segmentation, and regression operations. These operations are followed through and detailed further in consequent sub-sections. As such the upcoming subsections will demonstrate the way in which these technologies are being used in image-processing operations and how they aid in overcoming problems encountered with previously applied models.

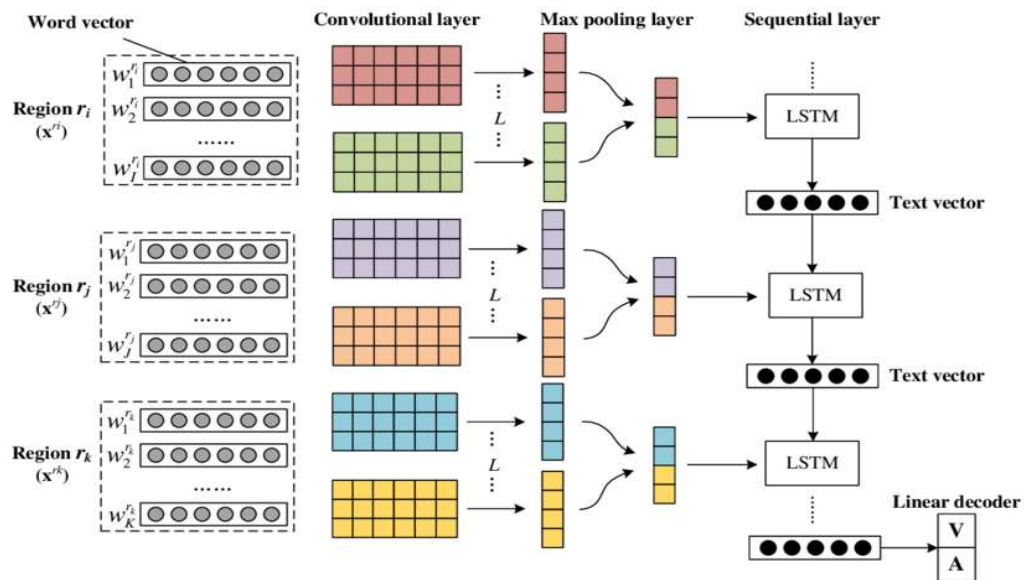


Figure 4: A sample of an architecture of type CNN-LSTM

5. Classification and Regression Approaches

In the context of implementing Deep Learning models to compute image segmentation image-data, the most common operations performed are, setting-specific classification, semantic-segmentation, instance-segmentation, and object-identification. The aim of setting-specific image-classification is to appoint a class-label to every input image, while the process of object detection is focused on building boundary-boxes that outline the detected image in an image and then assigning labels to these boundary-boxes. As such the process of identification is thought to be a more complex operation as it performs both detection and classification tasks. Additionally, identification can also be done with outlining sections or patches around an object rather than boundary-boxes, which specifies the category of an object at the pixel level. This operation is referred to as semantic segmentation. A drawback of this method is that it is unable to discriminate between objects of the same class owing to the capacity of a pixel to hold a single class-label [5]. As a means of overcoming this limitation a new approach called instance-segmentation was introduced that incorporated concepts of both semantic segmentation and object detection. It allows to perform identification of different objects in pixel level clouds, and these clouds are then labelled with corresponding class-labels [10].

To develop a regression-based model, the network requires readjustment in that its end layer i.e., the fully connected layer is modified to handle regression operations in place of the more regularly implemented image-classification tasks. Due to this modification, sequential data is evaluated with different criteria from that of classification operations. The

implementation of regression-based operations utilizing Deep Learning models is not as common as that of classification-task, albeit recent research has demonstrated the promise shown by using this combination in processing of image segmentation related data. One such research is by [9] in which an in-depth investigation of deep-regression models was carried out and was shown that commonly used refined networks such as ResNet-50 [6] and VGG-16 [11] showed incredibly promising outcomes. One disadvantage of these models is that they have been tailored to perform distinct functions, meaning they are not that well suited for general purpose problems. Another drawback is that deep-regression methods are not successful every time. A useful technique is to categorize the output-space and to allow the model to handle the operation as an image-classification task. In the context of UAV-based image segmentation operations, it is preferred to utilize popular models. Models other than the “ResNet-50” and VGG-16” [11], include VGG-11 and AlexNet, used by [10]. A prospect for future research-considerations in retrospect of the required implementations is optimizer functions. Currently models that incorporate dynamic learning-rates, for instance RMSProp, AdaGrad, Adam, and AdaDelta are the most prominently recognized.

6. Scene-Wise Classification, Object Detection, and Segmentation

Setting-specific image-classification or scene-recognition are operations that assign a theme or label for an image, also referred to as a patch, according to other sample-images, for instance in agricultural, beach, and urban settings, among others [2,4]. Simple Deep Neural Network models were introduced for these operations of which there are many that are broadly used for conventional image-classification tasks. In the scope of image segmentation implementations, it is not common to employ a scene wise classification approach. Rather, applications in this field would find it more advantageous to adopt object-identification and instance-segmentation methods. In terms of scene wise categorization, the approach requires only a labeling of the class-label present in the image, whereas for object-identification every object present in the image is required to be outlined by a boundary-box, hence increasing the cost required to construct such sample data clusters. The abundance reduces even further for data clusters related to instance-segmentation as a “mask” needs to be constructed around every pixel the object is present in, meaning more accuracy is required for such an operation. **Figure 5** illustrates the operation of the annotation step for both object identification and pixel-wise semantic segmentation processes.

7. Conclusions

Object-identification models can be classified in two elementary approaches: regression-based mechanisms and region-proposal based mechanisms, also referred to as one stage and two stage detectors respectively [3-6]. Utilizing the common two stage object detector approach requires construction of probable regularly shaped boundary-boxes on the landscape map. Each object is then categorized with a class-label and the successful detections are reinforced through boundary-box regression. A broadly implemented technique utilized in many studies to compute region-proposals was through the Faster-RCNN model combined with a Region-Proposal-Network (RPN) [7]. Several other modern models are also available such as Cascade-RCNN [9-13]. Dynamic-RCNN [4], DetectoRS [6]. In the case of one stage detectors, they omit the region-proposal task and immediately outline the position of objects and label them accordingly. The omission of this task improves the speed of the identification process, but the accuracy is compromised as a result. This technique is commonly referred to as region free detection due to the nature of the model to utilize the image grid to segment the image and classify the objects with a class-label. Accordingly, there are a few detectors present that can perform both regression-based and region-proposal based approaches.

Object-identification related approaches can be thought to be built from three constituents namely:

- The backbone that is focused on deriving the inherent characteristics from the images.
- The neck, a mediatory element of the structure that lies in the middle of the head and backbone, its purpose being to refine the information generated by the backbone and lastly.
- The head, the structure responsible for constructing the boundary-boxes by performing the identification and classification operations on the image.

8. REFERENCES

- [1] Adão, T., Hruška, J., Pádua, L., Bessa, J., Peres, E., Morais, R., & Sousa, J. J. *Hyperspectral imaging: A review on UAV-based sensors, data processing and applications for agriculture and forestry. Remote Sensing*, 9(11), 1110. (2017).
- [2] Adayel, R., Bazi, Y., Alhichri, H., & Alajlan, N. *Deep open-set domain adaptation for cross-scene classification based on adversarial learning and pareto ranking. remote sensing*, 12(11), 1716. (2020).

- [3] Al-Najjar, H. A., Kalantar, B., Pradhan, B., Saeidi, V., Halin, A. A., Ueda, N., & Mansor, S. *Land cover classification from fused DSM and UAV images using convolutional neural networks. Remote Sensing*, 11(12), 1461. (2019).
- [4] Ammour, N., Alhichri, H., Bazi, Y., Benjdira, B., Alajlan, N., & Zuair, M. *Deep learning approach for car detection in UAV imagery. Remote Sensing*, 9(4), 312. (2017).
- [5] Ampatzidis, Y., & Partel, V. *UAV-based high throughput phenotyping in citrus utilizing multispectral imaging and artificial intelligence. Remote Sensing*, 11(4), 410. (2019).
- [6] Apolo-Apolo, O. E., Martínez-Guanter, J., Egea, G., Raja, P., & Pérez-Ruiz, M. *Deep learning techniques for estimation of the yield and size of citrus fruits using a UAV. European Journal of Agronomy*, 115, 126030. (2020).
- [7] Audebert, N., Le Saux, B., & Lefèvre, S. *Deep learning for classification of hyperspectral data: A comparative review. IEEE geoscience and remote sensing magazine*, 7(2), 159-173. (2019).
- [8] Bachman, P., Hjelm, R. D., & Buchwalter, W. *Learning representations by maximizing mutual information across views. arXiv preprint arXiv:1906.00910. (2019).*
- [9] Badrinarayanan, V., Kendall, A., & Cipolla, R. *Segnet: A deep convolutional encoder-decoder architecture for image segmentation. IEEE transactions on pattern analysis and machine intelligence*, 39(12), 2481-2495. (2017).
- [10] Bah, M. D., Hafiane, A., & Canals, R. *Deep learning with unsupervised data labeling for weed detection in line crops in UAV images. Remote sensing*, 10(11), 1690. (2018).
- [11] Ball, J. E., Anderson, D. T., & Chan Sr, C. S. *Comprehensive survey of deep learning in remote sensing: theories, tools, and challenges for the community. Journal of Applied Remote Sensing*, 11(4), 042609. (2017).
- [12] Banerjee, K., Gupta, R. R., Vyas, K., & Mishra, B. *Exploring Alternatives to Softmax Function. arXiv preprint arXiv:2011.11538. (2020).*
- [12] Barbedo, J. G. A., Koenigkan, L. V., Santos, P. M., & Ribeiro, A. R. B. *Counting cattle in UAV images dealing with clustered animals and animal/background contrast changes. Sensors*, 20(7), 2126. (2020).

Acknowledgement

The authors Mr. Shouket Abdulrahman Ahmed, Prof Dr. Hazry Desa, and Ass. Prof Dr. Abadal-Salam T. Hussain of University Malaysia Perlis (Malaysia) would like to acknowledge and thanks Al-Kitab university, Kirkuk, Iraq and its scientific journal “Al-Kitab Journal for Pure Sciences” for acceptance and support our research titled “Performance Improvements Using Deep Learning Based Object-Identification”