



UWS Academic Portal

Deformable patch-based-multi-layer perceptron mixer model for forest fire aerial image classification

Mittal, Payal; Sharma, Akashdeep; Singh, Raman

Published in:
Journal of Applied Remote Sensing

DOI:
[10.1117/1.JRS.17.022203](https://doi.org/10.1117/1.JRS.17.022203)

Published: 07/11/2022

Document Version
Peer reviewed version

[Link to publication on the UWS Academic Portal](#)

Citation for published version (APA):
Mittal, P., Sharma, A., & Singh, R. (2022). Deformable patch-based-multi-layer perceptron mixer model for forest fire aerial image classification. *Journal of Applied Remote Sensing*, 17(2), [022203].
<https://doi.org/10.1117/1.JRS.17.022203>

General rights

Copyright and moral rights for the publications made accessible in the UWS Academic Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

If you believe that this document breaches copyright please contact pure@uws.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.

Mittal, P., Sharma, A., & Singh, R. (2022). Deformable patch-based-multi-layer perceptron mixer model for forest fire aerial image classification. *Journal of Applied Remote Sensing*, 17(2), [022203]. <https://doi.org/10.1117/1.JRS.17.022203>

Deformable Patch based MLP Mixer Model for Forest Fire Aerial Image Classification

Payal Mittal,^{a*} Akashdeep Sharma,^a Raman Singh,^b

^aUniversity Institute of Engineering and Technology, CSE Department, Sector 25, Chandigarh, India, 160014

^bUniversity of the West of Scotland, Cyber Security Department, Lanarkshire, United Kingdom, ML11 7AA

Abstract. Unmanned Aerial Vehicles (UAVs), equipped with mounting camera sensors, facilitate a wide domain of applications deployed in the real-time world. The situational awareness for applications such as search and rescue in case of wildfires, estimation of endangered flora and fauna and emergency responses have seen paradigm shift due to UAVs capability of accessing in remote and challenging areas such as forests. The last decade has seen tremendous growth in CNN based methods for object classification, detection and segmentation tasks. Recently emerged Attention-based Transformer models have been trying to achieve state-of-the art results in predicting images. This paper proposed a novel MLP-Mixer architecture for classification of burned piles in dense forests. MLP mixer architecture tries to eliminate the shortcomings of convolutions and attention by merging them to obtain good performance. The shallow learning of CNN layers and fixed-size patch embedding in transformers have been eliminated by introducing a new module of DePatch in proposed MLP mixer model which divides the input images in a deformable pattern to detect forest fires at an early stage. On the pile images dataset obtained by drones during a burning pile of debris in an Arizona pine forest, our suggested classification algorithm has been tested. The performance of the proposed model has been compared with transformer models.

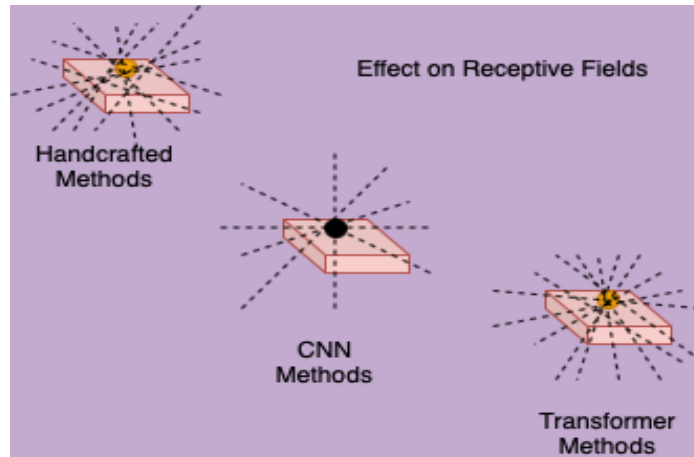
Keywords: Aerial Image Classification; Convolutional Neural Networks; Transformers; MLP Mixer; Computer Vision; Deep Learning

*First Author, E-mail: payalmittal6792@gmail.com

1. Introduction

The use of unmanned aerial vehicles (UAVs) as a remote sensing platform for a range of practical applications, such as traffic monitoring [1], wildfire detections [2], precision agriculture [3], and the processing of satellite data [4], has received a lot of attention recently. Due to its use in numerous civil and commercial applications, UAV growth has accelerated recently. The UAVs have characteristics of cost-effectiveness, high performance and low power consumption which make it possible to incorporate vision-based automatic applications [5]. Height of capturing, tilt camera angle, and light settings play a significant role in the information of aerial images. The quality of generated images is highly dependent on the altitude of the flying vehicle and the characteristics of the sensors used. Recent years have seen the emergence of high-performance deep learning-based categorization architectures. Applications for emergency response and catastrophe management can benefit from deep learning approaches to quickly retrieve vital information, improve response times in time-sensitive circumstances, and help in-the-loop decision-making processes. The challenges of densely located, occluded objects, noise and background clutter exist in aerial images [6]. There is a need for powerful classification algorithms for overcoming these challenges. CNNs became the mainstream for performing standard classification in computer vision. The effective deep learning classification architectures of VGG [7], GoogleNet [8], Inception [9], ResNet [10], DenseNet [11], Lite-HRNet [12], and EfficientNet [13] are based on architectural advances such as depth-wise and deformable convolutions. The CNNs continue to be the well-known foundational designs for computer vision applications, but

49 Transformer-like architectures have also shown great promise for unified modelling of vision and
50 language. In addition to efficiently capturing long-range dependency within the sequence and
51 extracting additional semantic data, freshly developed transformers outperform CNNs in this
52 regard. In terms of research discoveries and logical knowledge, transformers in vision assist in
53 bridging the gap between the NLP and computer vision communities. The transformers help to
54 increase the receptive field deprived from the resolutions. The effect of receptive fields on



55

56 **Fig. 1** Difference of receptive areas in handcrafted, CNN and Transformer based methods

57 handcrafted, CNN and transformers based methods is observed in Fig. 1. The concentration of the
58 receptive field in case of handcrafted methods is more when compared with CNN methods but
59 prior methods need manual engineering. But, the transformers based methods contribute
60 significantly to detect and classify objects in an effective manner. Transformer-based classification
61 models employ a fixed-size patch embedding with the implicit presumption that images are
62 appropriate for the fixed image split design. However, such strong patch splits may cause semantic
63 disagreement between images and issues with the collapsing of local structures in an image. In
64 order to address the discussed shortcomings of CNNs and transformers, MLP Mixer architecture
65 is utilized for the categorization of fire photographs taken by drones during a burning pile of trash
66 in an Arizona pine forest [14]. Only straightforward matrix multiplication operations, changes to
67 the data layout, and scalar non-linearities are used in Mixer's architecture. The use of multi-layer
68 perceptrons is repeated across either feature channels or spatial locations. To the best of our
69 knowledge, the classification model we propose is the first MLP mixer that executes patch splitting
70 in a data-specific manner. In this study, we create the DePatch deformable patch, which adaptably
71 changes the scale and position of each patch. The DePatch module is easy to create and can be
72 used as a plug-and-play module. Our study is focused by proposing a novel deep learning based
73 classification model for aerial images. The proposed deep model contains a deformable patch for
74 MLP based mixer model for object classification. The proposed study aims to identify forest
75 wildfires by deep feature representation approach. Early wildfire detection is crucial because it has
76 the potential to cause serious harm to ecosystems, residential areas, forests, and wildlife habitats
77 in the past. Recent developments in aerial surveillance systems in particular can give operational
78 troops and first responders more precise information on the behaviour of fires for improved fire
79 management. The proposed research work aims to contribute to a modeling shift in achieving
80 strong performance on visual recognition tasks. The major offerings of this paper include:

81 a) Proposing a novel deformable patch based MLP mixer model for aerial image classification

82 b) Proposed deep learning-based classification model is an amalgamation of CNN and transformer
83 based methods

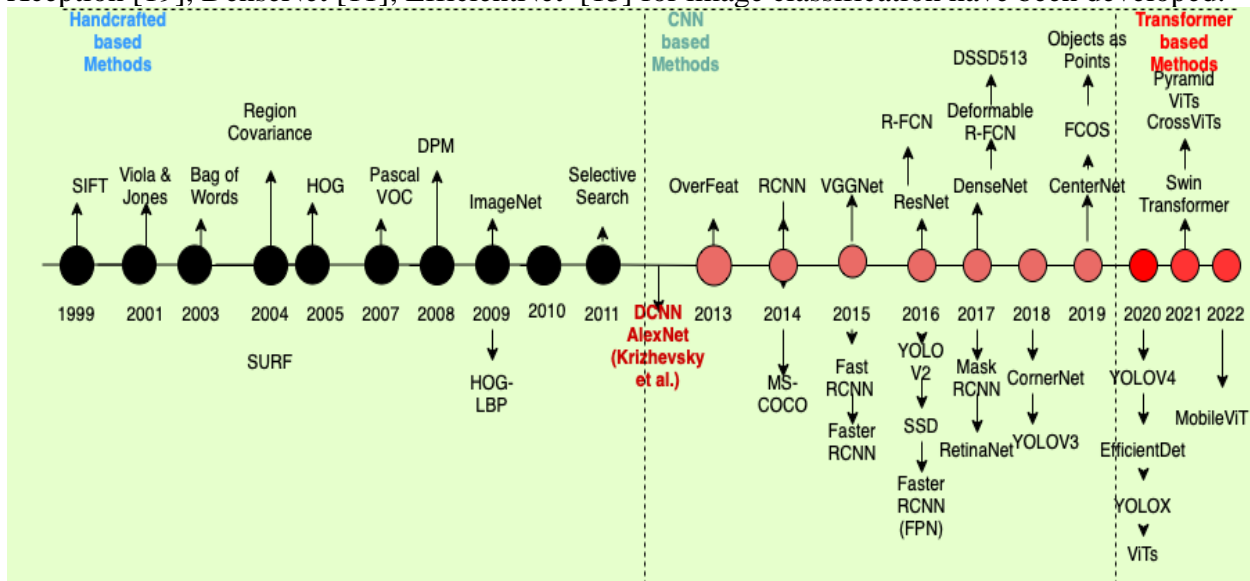
84 c) Proposed classification model has been evaluated on the pile imagery dataset collected by
85 drones.

86 The organization of the paper is as follows: Section 2 describes the related work of deep learning
87 object classification models based on CNNs and transformers. Section 3 presents the details of the
88 proposed deep learning-based object classification model in an elaborate manner. Section 4
89 discusses training experimental setup in which the proposed model requirements and evaluation
90 parameters are being analyzed. Section 4 also represents analysis of results obtained from the
91 proposed classification model. The last section concludes the results achieved with focus on future
92 directions.

93

94 2 Related Work

95 The growth of object detection algorithms and classification architectures is shown as a timeline
96 diagram in Fig. 2. On a year-by-year time period, the figure depicts the evolution of various hand-
97 crafted, deep learning, and transformer-based algorithms. Scale Invariant Feature Transform
98 (SIFT) [15], Histogram of Gradients (HOG) [16], and SURF [17] were among the manual feature
99 descriptors that dominated prior to 2012, but since then, CNN-based object detection methods have
100 emerged. A breakthrough happened in handcrafted-based object detection and classification
101 methods when deep learning-based convolutional architecture AlexNet [18] performed
102 significantly when compared with former approaches. In 2014, a powerful detector region-based
103 CNN method was developed through the combination of region proposals with CNNs. By
104 classifying object proposals using a deep convolution network, it obtains outstanding object
105 identification accuracy. After the proposal of the RCNN method, advancement in deep learning-
106 based object detectors began. After 2015, a number of powerful architectures such as VGG [7],
107 Xception [19], DenseNet [11], EfficientNet [13] for image classification have been developed.



108
109

Fig. 2 Timeline of various object detection and classification algorithms

110 The history of computer vision demonstrates that the availability of larger datasets along with
111 increased computational capacity frequently leads to a paradigm shift. Applications for computer
112 vision are increasingly using CNNs and other deep learning methods. In comparison to CNN-
113 based neural networks, the recently developed attention-based transformer models represent a

114 paradigm change for the middle of the 2020s [20]. When compared to CNNs, the recovered
115 features from transformers can more accurately reflect long-range dependency within the
116 sequence, and they also carry more semantic information.

117 In the next sections, we will be having a brief idea about self-attention based transformer based
118 models. ViTs and Swin transformer based architectures help in improving the model
119 interpretability as it relies on attention mechanisms which makes a prediction.

120 *2.1 ViT based Classification Model*

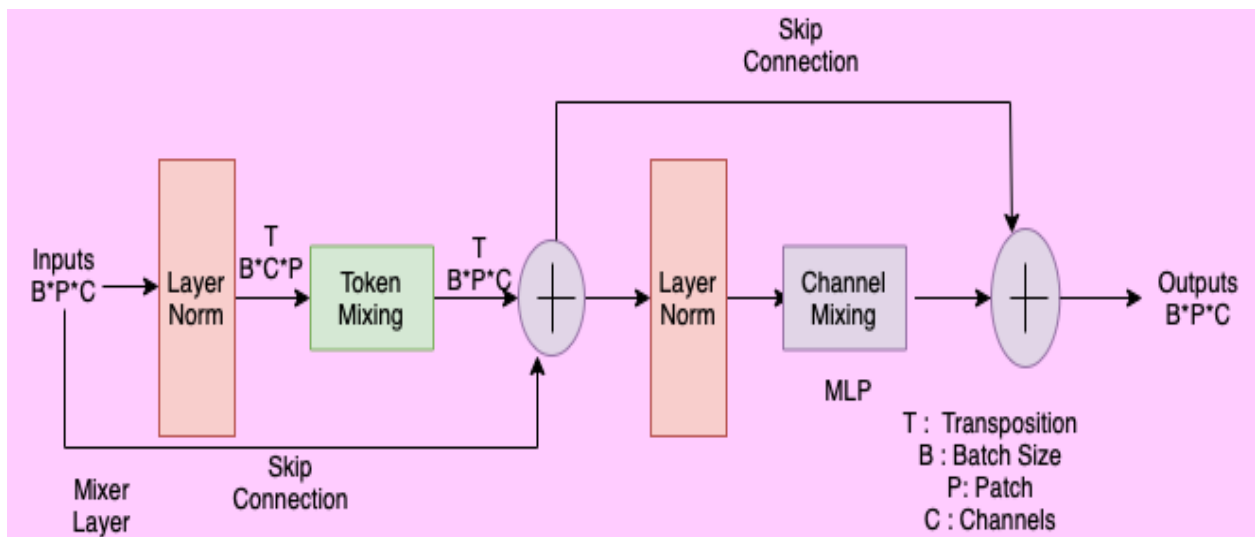
121 The ViT transformers for image classification were unveiled at the end of 2020. ViT carried on
122 the time-consuming process of learning from unprocessed data while eradicating arbitrary visual
123 features and inductive biases from models. The ViT transformer consists of three parts: a patch
124 embedding module, multi-head self-attention blocks, and feed-forward multi-layer perceptrons.
125 To develop features for image classification, the ViT divided an image into 16x16 patches and
126 sent the image patch sequence through the transformer architecture. After the patch embedding
127 module has turned the input image into a list of tokens, the network alternately stacks multi-head
128 self-attention blocks and MLPs to get the final representation. The patch embedding module
129 separates images into set sizes and positions before embedding a linear layer into each defined
130 patch. With the use of substantial training data, ViT obtained outcomes comparable to those of
131 conventional CNN designs. Despite great progress, most architectures still lost information since
132 they divided the input image according to a preset pattern without taking the input's content or
133 geometric variations into account [21].

134 *2.2 Swin Transformer*

135 Swin transformer functions as a general-purpose backbone for computer vision by converting
136 traditional multi-head attention to shifted window attention based models. Swin Transformer's
137 shift of the window partition between subsequent self-attention layers is a fundamental component
138 of its design, and it greatly outperformed ViT and ResNeXt models with comparable latency on
139 the tasks. Local multi-headed self-attention modules based on alternate shifting patch windows in
140 succeeding blocks make up the Swin Transformer block. The Swin transformer design consists of
141 a patch splitting module, similar to ViT, that divides an input RGB image into non-overlapping
142 patches. Each patch is handled as a token, and its feature is configured as a concatenation of the
143 RGB values of the individual pixels. The different transformer blocks with updated self-attention
144 computation are applied to these patch tokens, which keep the token count, along with the linear
145 embedding. The different Swin Transformer blocks are then applied to the patches in four stages,
146 progressively reducing the number of patches to maintain hierarchical representation. As the
147 network becomes deeper, patch merging layers reduces the number of tokens to provide a
148 hierarchical representation. After concatenating the features of each set of two adjacent patches,
149 the first patch combining layer applies a linear layer to the four-dimensional concatenated features.
150 It can model information at various scales and has a linear computational complexity with respect
151 to image size [22]. Swin Transformer obtained the best results on the MS-COCO dataset [23],
152 although it utilizes a lot more parameters than convolutional models. Transformers offers a
153 paradigm change away from CNN-based neural networks. Convolution may very likely be
154 replaced by it in these tasks, even if it is still in the early stages of use in vision. The fixed scale of
155 the tokens in the present transformer-based models precludes their usage in vision applications.
156 Another distinction is that text passages have a significantly lower word resolution than pixels in
157 graphics. Because the semantics of objects are destroyed by the Swin Transformers [22] and
158 MobileViTs [24] models, which are still in development, they are still having trouble.

159 *2.3 MLP Mixer*

160 Transformers have the aforementioned drawbacks; thus, this paper incorporated the recently
 161 created MLP Mixer architecture, which comes from non-local blocks. Because it doesn't require
 162 self-attention or convolutional layers, the MLP-Mixer is a novel design in computer vision that
 163 varies from earlier, successful models like CNNs and transformers. The MLP-blocks in Mixer are
 164 designed to immediately process embeddings of these patches after converting images into a series
 165 of patches. The design of Mixer is influenced by more contemporary transformer-based systems.
 166 It utilizes conventional regularization and optimization methods, relies on token and channel-
 167 mixing MLPs, and is scalable to big data sets successfully. The competitive approach behind the
 168 MLP-Mixer architecture is built exclusively on multi-layer perceptrons rather than convolutions
 169 or self-attention. These perceptrons depend on simple matrix multiplication operations,
 170 adjustments to the data layout, and scalar non-linearities. They are repeatedly used across either
 171 feature channels or spatial locations. The MLP-Mixer architecture is based on multi-layer
 172 perceptrons and contains two different layers as described in Fig. 3. These two layers were made
 173 up of two MLPs: one that is applied individually to picture patches in order to mix location-specific
 174 features, and the other that is applied across patches in order to blend spatial data. Convolutions
 175 with small kernels of 11 are used in Mixer to transform convolutions into typical dense matrix
 176 multiplications (channel-mixing MLPs), which apply independently to each spatial point [25].

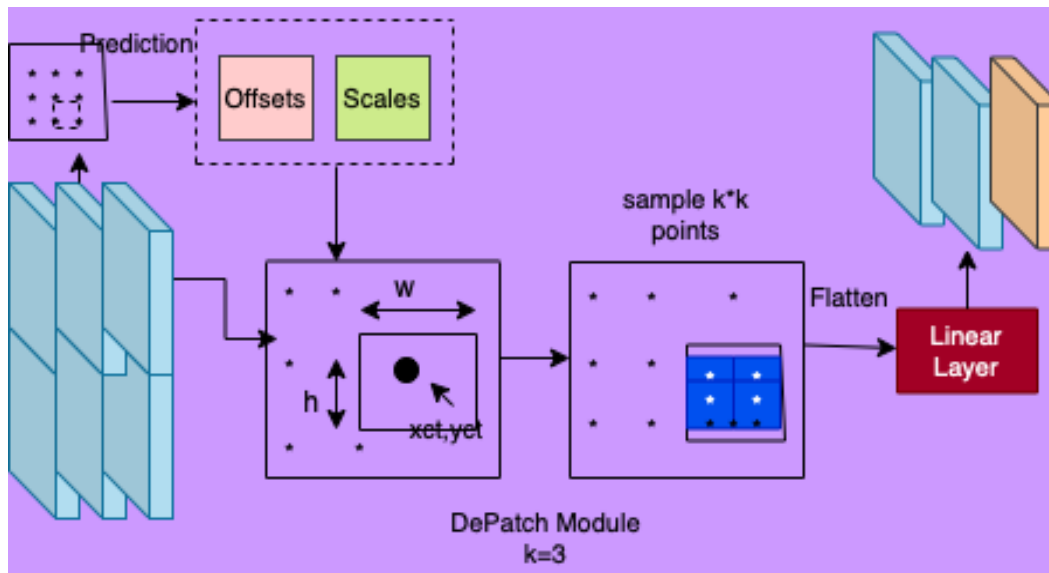


177
 178 **Fig.3** Mixer architecture containing token-mixing and channel-mixing MLP

179 As a result, geographic information cannot be aggregated; as a workaround, dense matrix
 180 multiplications, or token-mixing MLPs, are performed to each feature across all spatial locations.
 181 The competitive scores that the MLP mixer model achieved on benchmarks for image
 182 classification served as the basis for the proposed classification method. When trained on large
 183 datasets or using contemporary regularization techniques, these models have pre-training and
 184 inference costs that are comparable to those of cutting-edge models. The MLP-Mixer Architecture
 185 uses skip-connections and normalization layers, and each layer accepts an input of the same size.
 186 Matrix multiplications are then applied to the patches and features input table. As opposed to ViTs,
 187 Mixer uses a token-mixing layer to combine spatial information rather than position embeddings.
 188 MLPs are sensitive to the order of the input tokens and use a conventional classification head with
 189 a global average pooling layer, followed by a linear classifier.
 190

191 **3 De-Patch based MLP Mixer Architecture**

192 The existing CNN and transformer models for object classification pose serious challenges to
 193 aerial images. The transformer-based methods utilized a fixed-size patch embedding which might
 194 destroy the semantics of aerial objects. Further, the hard patch splits of CNNs brought two
 195 problems related to collapse of local structures and having semantic inconsistency across aerial
 196 images. Because scale-variance items can be seen in a variety of aerial photographs, it is difficult
 197 to capture the entire object-related local structure in a $16 * 16$ regular patch [26]. The same item
 198 may appear differently geometrically in many aerial photographs, depending on the scale, rotation,
 199 etc. The fixed method of picture splitting may capture contradictory information for the same
 200 object in various photographs. These updated patches run the risk of erasing semantic data, which
 201 would reduce classification accuracy. To address the aforementioned difficulties, we suggest
 202 DePatch, a novel module in the MLP mixer architecture that learns to adaptively partition the
 203 images into patches with varying positions and scales in a data-driven way as opposed to using
 204 predetermined fixed patches. The semantics in patches may be effectively preserved using our
 205 suggested classification architecture by integrating DePatch into MLP mixer design.



206

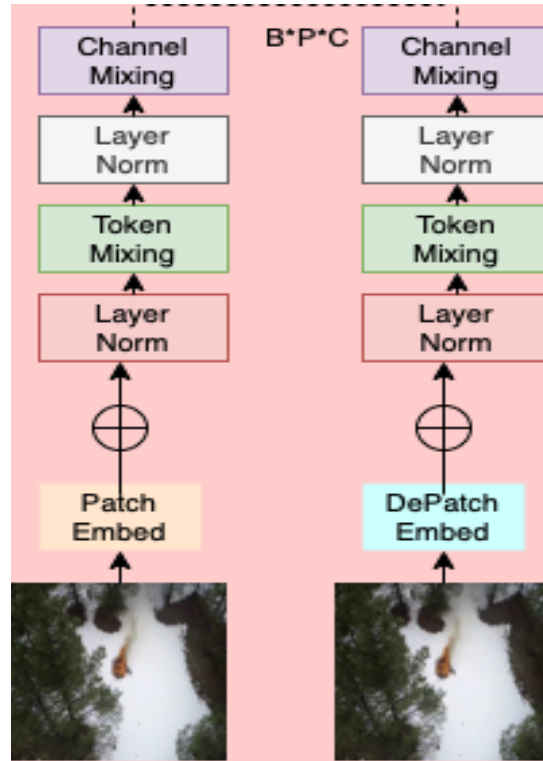
207

Fig. 4 DePatch module with offsets and scales within local features

208 The proposed DePatch based MLP Mixer architecture is shown in Fig. 4. The rectangle region is
 209 immutable for each patch, as illustrated in the figure, because the coordinates (x_{ct}, y_{ct}) , and size s
 210 of the patch are fixed. The interior pixels of the patch being used directly depict its feature. We
 211 relax these requirements to build our deformable patch embedding module, DePatch, which can
 212 better locate key structures and handle geometric deformation. Based on the contents of the input,
 213 projected parameters include the position and size of each patch. We estimate an offset (x, y) that
 214 will allow the location to move away from the original center. In terms of scale, all we do is swap
 215 out the fixed patch size s for the predictable s_h and s_w . In this way, we can determine a new
 216 rectangle region, and denote its left-top corner as (x_1, y_1) and right- bottom corner as (x_2, y_2) . We
 217 emphasize that $(\delta x, \delta y, s_w, s_h)$ can be different even for patches in a single image as shown in eq.
 218 1 and 2:

219
$$x_1 = x_{ct} + \delta x - \frac{S_w}{2}, y_1 = y_{ct} + \delta y - \frac{S_h}{2} \quad (1)$$

220
$$x_2 = x_{ct} + \delta x - \frac{S_w}{2}, y_2 = y_{ct} + \delta y - \frac{S_h}{2} \quad (2)$$



221
222 **Fig. 5** Left: Original patch based MLP architecture. Right: Modified architecture equipped with DePatch module

223 The architectural details of the proposed DePatch based MLP Mixer have been discussed. To retain
 224 semantics in the aerial photos, we integrated the DePatch module into the MLP Mixer framework.
 225 The MLP mixer architecture, as depicted in Fig. 5, receives the DePatch module as an input. The
 226 locations of the token and channel mixers are the same as in the original MLP mixer. When
 227 compared to CNNs and transformers, the performance of the suggested classification model shows
 228 that MLP Mixer is a superior option for aerial data. Our suggested classification approach is
 229 capable of handling tasks that demand pixel-level predictions without the use of transformers, as
 230 well as photos with significantly higher resolution. In this study, we use the suggested
 231 categorization model to locate slash piles that have been burned over the winter in high-elevation
 232 forests in the Southwest. The proposed model is able to combat challenges of aerial image
 233 classification such as lack of context information and imbalance of fore-ground and background
 234 training examples.

235
236 **4 Experimental Analysis**

237 Our proposed DePatch based MLP Mixer architecture performed pile aerial imagery classification
 238 in an effective manner. In this section, details about the aerial pile imagery dataset, training
 239 methodology and results for the proposed framework have been discussed.

241 *4.1 Aerial Pile Imagery Dataset*

242 The proposed classification technique has been evaluated on a fire image dataset collected through
243 drones during a burning piled detritus in Arizona pines forest. This dataset contains annotated
244 drone-based images and videos shot from infrared cameras for executing fire related detection,
245 classification and segmentation problems. The fire based classification and segmentation studies
246 can be evaluated on this dataset. For "Fire" vs. "Non-Fire," a total of 39,375 labelled frames were
247 used in the training phase, and 8,617 frames were used for the test data. Early wildfire detection is
248 crucial because it has the potential to cause serious harm to ecosystems, residential areas, forests,
249 and wildlife habitats in the past. These troubling statistics spur scientists to look for fresh
250 approaches to early fire detection and classification. A deep classification model can learn features
251 more effectively by being trained to execute fire picture classification tasks.

252 *4.2 Training Methodology*

253 The large-scale images dataset ImageNet has been deployed for initial training of feature extractor
254 for the proposed classification technique. The ImageNet dataset consists of 1.28M images
255 belonging to 1000 categories. We performed fine-tuning on a fire images dataset to detect wildfires
256 at an early stage. The dataset images are resized into 224×224 for training of proposed
257 classification technique. To broaden the variety of categorized items, advanced data augmentation
258 techniques like Mix up, CutMix, label smoothing, and Rand-Augment were used. The suggested
259 classification method has been trained using a 32-person batch size over 100 iterations. There is
260 no sign of convergence after 100 epochs. The training approach used the optimizer RMSProp with
261 a weight decay for non-bias parameters of 0.05 and an initial learning rate of 1 103. For fair
262 comparison, all of these settings are maintained with the MLP Mixer architecture.

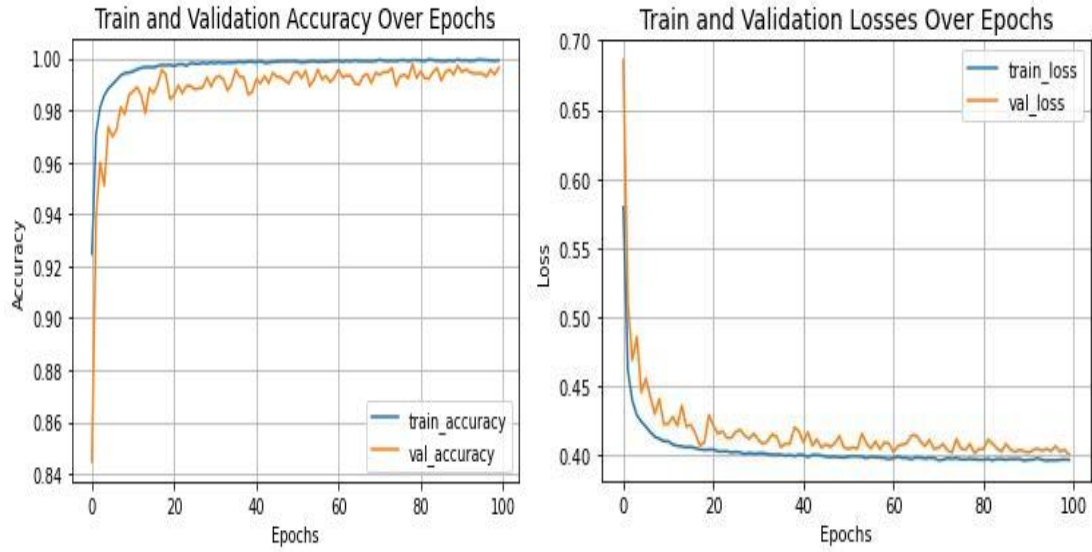
263 *4.3 Evaluation Parameters*

264 The popular classification-related evaluation metrics had been used to evaluate the effectiveness
265 of the suggested deep learning-based classification technique in order to effectively depict the
266 findings. While recall provides the percentage of True Positives (TP) from the whole quantity of
267 TP and False Negatives (FN), the metric precision refers to the TP fraction from the total sum of
268 TP and False Positives (FP). The true predictions from class one with the highest probability make
269 up the accuracy score.

270

271 **5 Results**

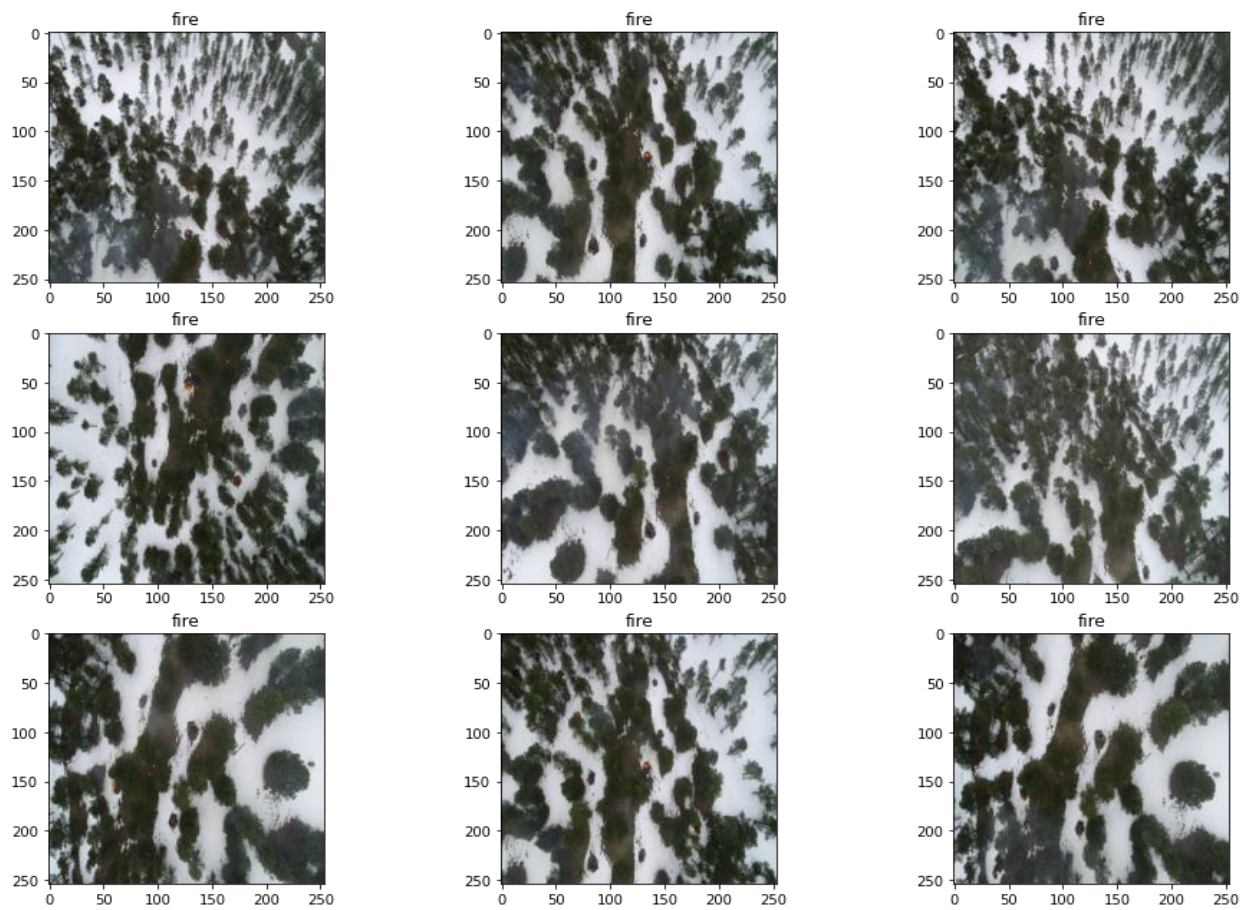
272 In this section, an analysis of the proposed optimized classification technique is provided by
273 considering the accuracy of chosen fire based aerial dataset. The plots related to accuracy and loss
274 for the proposed classification technique have been mentioned. The training accuracy value started
275 from 0.925 and goes up to 1.00 over 100 epochs. The validation accuracy values initially showed
276 bumpy behavior for starting 30 epochs, after which it started converging with the training values.
277 The starting value for the validation accuracy was 0.85 and goes up to 0.99. The training was
278 stopped as no further convergence took place. The training loss value started from 0.85 and goes
279 up to 0.40 over 100 epochs. The validation loss values initially showed a steep change in behavior
280 from .70 to 0.48 value in starting epochs, after which it started converging with the training loss
281 values. The settled value for the validation loss was 0.48 and goes up to 0.41 for 100 epochs. The
282 training was stopped as no further convergence took place. The predictions for the proposed
283 classification technique have been illustrated in Figure. The fire class was properly recognized by
284 DePatch-based MLP mixer technique.



285

286

Fig. 6 Accuracy and loss plots for the proposed classification model



287

288

Fig. 7 Correct Predictions for the proposed classification model

289

290 *5.1 Comparison with Other Approaches*

291 To compare the proposed classification technique, the promising transformer models of Vision
292 and Swin transformers have been chosen, shown in Table 1. ViTs used learning from unprocessed
293 data and inductive biases from models. By limiting self-attention computation to non-overlapping
294 local windows while simultaneously allowing for cross-window connections, Swim Hierarchical
295 Transformers were constructed with shifted windows and improved efficiency. These transformer
296 models lacked speed and accuracy while exhausting more hardware resources. The ViT-B/16
297 version obtained an accuracy score of 46.18 while exhausted 86.4 million parameters and 55.5 G
298 FLOPs. The better version of transformers i.e. Swin Transformers had also attained 61.35 and
299 82.48 as top-1 and top-5 accuracy scores. The classification results of Swin Transformers were
300 better than vision transformers but could not hold better position in front of our proposed
301 classification model. Transformers' subpar performance is caused by the absence of inductive bias
302 and fixed post-training weights. The drawbacks of transformers include the fact that they can only
303 compute global self-attention for non-overlapping local windows, the difficulty of pixel-level
304 predictions for transformers, and the fact that the computational complexity of their self-attention
305 is quadratic to image size. The input image size for vanilla MLP Mixer model is 224*224 and
306 utilized 19 million parameters through which it achieved top-1 accuracy 74.67 and 87.54 top-5
307 accuracy. The best performance was recorded by our proposed DePatch based MLP mixer model
308 which outperformed employed other classification models. The proposed classification technique
309 achieved top-1 accuracy of 77.23 and top-5 accuracy of 93.45 which outperformed employed
310 transformer models. The number of employed parameters were also minimum in case of DePatch-
311 based MLP mixer model.

312 **Table 1** Comparison of the proposed classification model with other models
313

Model	Resolution	Param	FLOPs	Top-1 (%)	Top-5 (%)
ViT-B/16	384*384	86.4M	55.5G	46.18	73.63
Swin	224*224	29M	4.5G	61.35	82.48
MLP Mixer	224*224	19M	2.2G	74.67	87.54
DePatch based MLP (ours)	224*224	18M	2.0G	77.23	93.45

314
315 A new module of DePatch had been proposed in MLP mixer model which divides the input images
316 in a deformable pattern to detect forest fires at an early stage. This DePatch module in the chosen
317 MLP mixer based classification model incorporated the awareness of input images and geometric
318 variations. The improved classification model for forest pile burn images of Arizona forest has
319 been proposed which can further be extended for performing multitude of applications. These
320 aerial applications provide the methods for performing smart agriculture, defense missions and
321 industry related activities.

322
323 **6 Conclusion**

324 This paper proposed an aerial scene classification of forest fire situations from drone images using
325 a novel multi-layer perceptron based network model. The hard patch split of CNNs brought two
326 problems related to collapse of local structures and having semantic inconsistency across images.
327 Transformers would be unable to provide pixel-level predictions on high-resolution aerial photos
328 because the computational difficulty of its self-attention scales quadratically with the size of the
329 image. In MLP Mixer model, multi-layer perceptron blocks turn pictures into a series of patches

330 and process embeddings of these patches directly. To efficiently analyze massive amounts of data,
331 it relied on token and channel-mixing MLPs as well as conventional regularization and
332 optimization approaches. In order to maintain the semantics of the aerial images, DePatch module
333 had been included into the MLP Mixer framework. Instead of utilising pre-determined fixed
334 patches, this DePatch module adaptively divides the photos into patches with varying positions
335 and scales. A drone-collected fire image dataset from a smouldering pile of debris in an Arizona
336 pine forest was used to evaluate the suggested deep learning-based categorization technique. The
337 suggested classification method outperforms transformer models and the standard MLP-Mixer
338 model with top-1 accuracy of 77.23 and top-5 accuracy of 93.45. The DePatch-based MLP Mixer
339 model likewise used the fewest possible parameters.

340

341 *Acknowledgements*

342 We acknowledge DIC, Panjab University Chandigarh for funding a workstation that enabled us to
343 perform experiments on NVIDIA TITAN XP GPUs. This research work has been done under UGC
344 NET JRF scholarship, New Delhi, India.

345

346 *Funding*

347 This research did not receive any specific grant from funding agencies in the public, commercial,
348 or not-for-profit sectors.

349

350 *References*

351

- 352 [1] J.S. Zhang, J. Cao, and B. Mao, "Application of deep learning and unmanned aerial vehicle
353 technology in traffic flow monitoring," in *Machine Learning and Cybernetics (ICMLC),
354 2017 International Conference on*, 2017, vol. 1, pp. 189–194.
- 355 [2] S. Treneska and B. R. Stojkoska, "Wildfire detection from UAV collected images using
356 transfer learning."
- 357 [3] Y.-C. C. and P.Y. W. Pei-Chun Chen, "Imaging Using Unmanned Aerial Vehicles for
358 Agriculture Land Use Classification," *Agriculture*, vol. 10, no. 9, pp. 1–14, 2020.
- 359 [4] M. Kraft, M. Piechocki, B. Ptak, and K. Walas, "Autonomous, Onboard Vision-Based Trash
360 and Litter Detection in Low Altitude Aerial Images Collected by an Unmanned Aerial
361 Vehicle," *Remote Sens.*, vol. 13, no. 5, p. 965, 2021.
- 362 [5] P. Q. and T. P. E. Wang, W. Shu, J. Zhu, S. Xu, "Low-altitude UAV Recognition and
363 Classification Algorithm Based on Machine Learning," in *IEEE 16th Conference on
364 Industrial Electronics and Applications (ICIEA)*, 2021, pp. 1136–1141.
- 365 [6] X. Zhang, E. Izquierdo, and K. Chandramouli, "Dense and Small Object Detection in UAV
366 Vision Based on Cascade Network," in *Proceedings of the IEEE International Conference
367 on Computer Vision Workshops*, 2019, p. 0.
- 368 [7] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image
369 recognition," *arXiv Prepr. arXiv1409.1556*, 2014.
- 370 [8] M. Al-Qizwini, I. Barjasteh, H. Al-Qassab, and H. Radha, "Deep learning algorithm for
371 autonomous driving using googlenet," in *2017 IEEE Intelligent Vehicles Symposium (IV)*,
372 2017, pp. 89–96.
- 373 [9] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and
374 the impact of residual connections on learning," in *Thirty-First AAAI Conference on*

- 375 *Artificial Intelligence*, 2017.
- 376 [10] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in
377 *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp.
378 770–778.
- 379 [11] F. Iandola, M. Moskewicz, S. Karayev, R. Girshick, T. Darrell, and K. Keutzer, “Densenet:
380 Implementing efficient convnet descriptor pyramids,” *arXiv Prepr. arXiv1404.1869*, 2014.
- 381 [12] C. Yu *et al.*, “Lite-hrnet: A lightweight high-resolution network,” in *Proceedings of the*
382 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10440–
383 10450.
- 384 [13] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural
385 networks,” in *International conference on machine learning*, 2019, pp. 6105–6114.
- 386 [14] A. Shamsoshoara, F. Afghah, A. Razi, L. Zheng, P. Z. Fulé, and E. Blasch, “Aerial imagery
387 pile burn detection using deep learning: The FLAME dataset,” *Comput. Networks*, vol. 193,
388 p. 108001, 2021.
- 389 [15] H. Zhou, Y. Yuan, and C. Shi, “Object tracking using SIFT features and mean shift,”
390 *Comput. Vis. image Underst.*, vol. 113, no. 3, pp. 345–352, 2009.
- 391 [16] Y. Pang, Y. Yuan, X. Li, and J. Pan, “Efficient HOG human detection,” *Signal Processing*,
392 vol. 91, no. 4, pp. 773–781, 2011.
- 393 [17] H. Bay, T. Tuytelaars, and L. Van Gool, “Surf: Speeded up robust features,” in *European*
394 *conference on computer vision*, 2006, pp. 404–417.
- 395 [18] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep
396 convolutional neural networks,” in *Advances in neural information processing systems*,
397 2012, pp. 1097–1105.
- 398 [19] F. Chollet, “Xception: Deep Learning With Depthwise Separable Convolutions,” in *The*
399 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- 400 [20] H. Chefer, S. Gur, and L. Wolf, “Transformer interpretability beyond attention
401 visualization,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and*
402 *Pattern Recognition*, 2021, pp. 782–791.
- 403 [21] A. Dosovitskiy *et al.*, “An image is worth 16x16 words: Transformers for image recognition
404 at scale,” *arXiv Prepr. arXiv2010.11929*, 2020.
- 405 [22] Z. Liu *et al.*, “Swin transformer: Hierarchical vision transformer using shifted windows,”
406 in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp.
407 10012–10022.
- 408 [23] T.-Y. Lin *et al.*, “Microsoft coco: Common objects in context,” in *European conference on*
409 *computer vision*, 2014, pp. 740–755.
- 410 [24] S. Mehta and M. Rastegari, “Mobilevit: light-weight, general-purpose, and mobile-friendly
411 vision transformer,” *arXiv Prepr. arXiv2110.02178*, 2021.
- 412 [25] I. O. Tolstikhin *et al.*, “Mlp-mixer: An all-mlp architecture for vision,” *Adv. Neural Inf.*
413 *Process. Syst.*, vol. 34, pp. 24261–24272, 2021.
- 414 [26] Z. Chen *et al.*, “Dpt: Deformable patch-based transformer for visual recognition,” in
415 *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 2899–
416 2907.

417
418
419
420



Payal is awarded with Junior Research Fellow (JRF) award from University Grant Commission (UGC), New Delhi in 2016. She has published various papers in international journals and conferences. Presently, she is pursuing full time Ph.D. in Computer Science and Engineering at University Institute of Engineering. & Technology, Panjab University, Chandigarh, India. Her research interests include unmanned aerial vehicles, machine and deep learning based image processing.



Akashdeep is currently working as an assistant professor in Computer Science and Engineering at University Institute of Engineering. & Technology, Panjab University, Chandigarh, India. He has chaired 5+ sessions in national and international conferences. He has published 22 research papers in international journals and conferences. He has 12+ years of teaching and research experience. His research interests include wireless networks, soft computing and video analytics, moving object detection and tracking, traffic sensing and classification.



Raman Singh is working as Lecturer in Cyber Security, University of the West of Scotland. He has completed Ph.D. from Institute of Engineering and Technology, Panjab University Chandigarh. He has published 14 research papers in international journals and has 6 years of teaching and research experience. He is a Microsoft Certified Technology Specialist (MCTS) and TechNet Certified Technology Expert. His area of interest includes Intrusion Detection, Network Security, Cyber Security, Autonomous Driving and Machine Learning.