
Exploring the MLP-Mixer Architecture for Image Classification: Performance Evaluation, Tuning, and Ablation Study on MNIST Dataset

Sijan Poudel

The University of British Columbia
Vancouver, Canada
sijuichi@gmail.com

Abstract

The research paper "MLP-Mixer: An all-MLP Architecture for Vision" introduces a novel approach to computer vision tasks based entirely on multi-layer perceptrons (MLPs). This report reproduces and analyzes the MLP-Mixer architecture proposed in the paper, specifically training and testing the model on the MNIST dataset. In addition, model tuning and baseline comparisons were conducted with CNN and simple MLP, and an ablation study was performed to determine the impact of different components of the MLP-Mixer model. The results demonstrate that MLP-Mixer outperforms other baseline models and achieves state-of-the-art performance on the MNIST dataset. Furthermore, the study revealed that certain components, such as the channel mixing layer, have a greater impact on the model's performance.

1 Introduction

Convolutional Neural Networks (CNNs) have long been the favored architecture for image classification tasks. However, recent advances in all Multi-Layer Perceptron (MLP) based architectures have demonstrated remarkable performance across various vision tasks. One such architecture, MLP-Mixer, exclusively relies on MLPs for image processing, entirely foregoing convolution operations.

The MLP-Mixer architecture was proposed by Tolstikhin et al. in their paper "MLP-Mixer: An all-MLP Architecture for Vision" [1]. The authors argued that the existing CNN-based architectures are computationally inefficient and do not generalize well to new domains. In contrast, MLP-Mixer is a simple and efficient architecture that can capture spatial and channel-wise relationships in the image. MLP-Mixer has achieved state-of-the-art performance on various image classification benchmarks, including ImageNet, CIFAR-100, and STL-10 [1]. In addition to classification, MLP-Mixer has also shown promising results on other vision tasks such as object detection and segmentation.

This report presents a reproduction of the MLP-Mixer architecture, its performance evaluation on the MNIST dataset, and an exploration of model tuning, baseline model comparisons, and ablation studies to understand the impact of various design choices on the model's performance.

2 Related work

CNNs have been the preferred architecture for numerous computer vision tasks, including image recognition, object detection, and segmentation, due to their ability to extract significant features from input images. However, CNNs face limitations, including a large number of parameters, slow training and inference times, and overfitting susceptibility.

Alternative architectures have been proposed to address these limitations. One such architecture is the Spatial Transformer Network (STN) [2], introduced by Jaderberg et al. in 2015. STN is a differentiable module that can be incorporated into a CNN to learn spatial transformations of the input, making the network more robust to input variations. Nevertheless, STN still relies on convolutional layers and shares the same limitations as CNNs.

Recently, researchers have investigated the use of fully connected layers or Multi-Layer Perceptrons (MLPs) for image recognition tasks. The Vision Transformer (ViT) [3], proposed by Dosovitskiy et al. in 2020, is one such architecture. ViT replaces convolutional layers with self-attention layers, allowing the network to learn long-range dependencies between different parts of the image. ViT has achieved state-of-the-art performance on several benchmark datasets and is a promising alternative to traditional CNNs.

MLP-Mixer is another architecture that substitutes convolutional layers with MLPs. While ViT employs self-attention layers to capture global dependencies, MLP-Mixer utilizes a different approach called token mixing. The input image is first divided into tokens, which are then processed by MLPs. Another MLP mixes the outputs of these MLPs, resulting in a set of transformed tokens. This process is repeated multiple times to learn different abstraction levels. Finally, the transformed tokens are concatenated and passed through another MLP to generate the final output.

Compared to traditional CNNs and ViT, MLP-Mixer offers several advantages, including a simpler and more interpretable architecture, fewer parameters, and state-of-the-art performance on various benchmark datasets like ImageNet, CIFAR-100, and STL-10. MLP-Mixer’s unique architecture efficiently and interpretable captures both local and global dependencies, making it a promising alternative to traditional CNNs and ViT for image recognition tasks. The success of MLP-Mixer has spurred further research into MLPs for computer vision tasks [4], and it will be intriguing to observe how this architecture evolves in the future.

3 MLP-Mixer model

The MLP-Mixer architecture proposed by Tolstikhin et al. (2021) is an all-MLP (multi-layer perceptron) approach for vision tasks. It is designed to be more efficient and scalable than traditional CNN (convolutional neural network) architectures while achieving state-of-the-art performance on several image classification benchmarks.

The MLP-Mixer architecture consists of two main components: the mixer layers and the per-patch fully connected layers. The mixer layers are responsible for mixing the information across different spatial locations, while the per-patch fully connected layers allow for modeling interactions between different feature channels within each patch.

As shown in Figure 1, the input image is first divided into non-overlapping patches, each of which is then linearly projected into a higher dimensional space. These per-patch embeddings are then processed by a series of mixer layers, each of which consists of two sub-layers: the token mixing layer and the channel mixing layer.

The token mixing layer processes the spatial information by computing pairwise interactions between different patches along the spatial dimensions. Specifically, it applies a multi-layer perceptron to each patch independently to compute a set of key, value, and query vectors. These vectors are then used to compute a weighted sum of the values across all patches, where the weights are given by the softmax of the dot product between each patch’s query vector and the key vectors of all other patches.

The channel mixing layer, on the other hand, models the interactions between different feature channels within each patch. It applies a fully connected layer to each channel independently, followed by a channel shuffle operation that randomly permutes the feature channels across different patches. This allows for the model to capture a wider range of interactions between different channels, rather than being limited to modeling interactions within each patch independently.

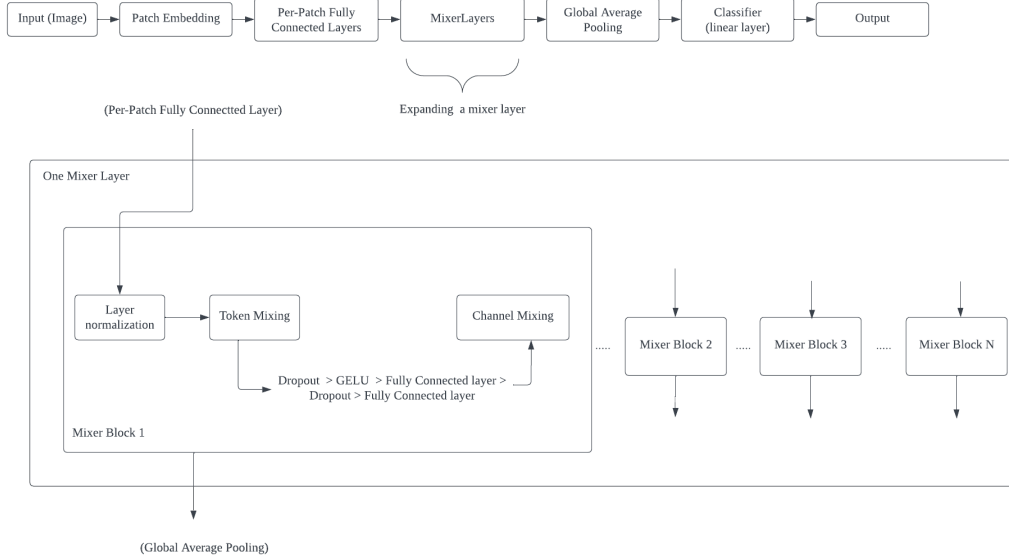


Figure 1: Graph of the MLP-Mixer architecture.

85 After passing through several mixer layers, the per-patch embeddings are aggregated by averaging
 86 them across all spatial locations, and then fed into a final linear layer to produce the output logits for
 87 classification.

88 Below are the equations to represent MLP-Mixer Architecture in detail:

89 The MLP-Mixer architecture is composed of a series of MixerBlocks, each of which consists of a
 90 Token mixing layer followed by a Channel mixing layer. The MixerBlock is defined as follows:

$$MixerBlock(x) = Channel_Mixer(Token_Mixer(LayerNorm(x))) \quad (1)$$

91 where x is the input tensor and $LayerNorm$ is a layer normalization layer.

92 The Token mixing layer processes the tokens or patches individually while the Channel mixing layer
 93 processes the channels of each token or patch. The equations for Token mixing and Channel mixing
 94 layers are as follows:

95 The Token mixing layer is defined as follows:

$$Token_Mixer(x) = x + MLP_2(Dropout(GELU(MLP_1(x)))) \quad (2)$$

96 where x is the input tensor, MLP_1 and MLP_2 are multi-layer perceptrons with two linear layers
 97 each, $GELU$ is the GELU activation function, and $Dropout$ is a dropout layer with a dropout
 98 probability of 0.1.

99 The Channel mixing layer is defined as follows:

$$Channel_Mixer(x) = x + MLP_4(Dropout(GELU(MLP_3(LayerNorm(x)))))) \quad (3)$$

100 where MLP_3 and MLP_4 are multi-layer perceptrons with two linear layers each, $LayerNorm$
 101 is a layer normalization layer, $GELU$ is the GELU activation function, and $Dropout$ is a dropout
 102 layer with a dropout probability of 0.1.

103 Hence, the MLP-Mixer model is defined as follows:

$$MLPMixer(x) = Classifier(Pool(Network(ToPatchEmbedding(x)))) \quad (4)$$

where x is the input image, *ToPatchEmbedding* is a module that converts the image into a sequence of patches and applies a linear projection to obtain a patch embedding, *Network* is a series of MixerBlocks from (1), *Pool* is an adaptive average pooling layer that reduces the spatial dimensions of the output to 1, and *Classifier* is a linear layer that maps the output to the number of classes.

Having defined the MLP-Mixer architecture, in order to fully illustrate the mathematical computations of the MLP-Mixer architecture, let us define some dimensions and formulas used across the layers.

1. **Input image dimensions:** (H, W, C_{in}) where H is the height, W is the width, and C_{in} is the number of input channels.
2. **Patch dimensions:** (P, P) where P is the patch size (assuming square patches).
3. **Number of patches:** $S = \frac{H \times W}{P^2}$.
4. **Linear projection dimension:** C .
5. **Token mixing hidden dimension:** D_S .
6. **Channel mixing hidden dimension:** D_C .
7. **Number of mixer layers:** L .

Given these dimensions, we can now represent the dimensions of the layers and the formulas used across them:

1. The input image is divided into non-overlapping patches of size (P, P) , resulting in S patches.
2. Each patch of size (P, P, C_{in}) is linearly projected into a higher-dimensional space of size C . The resulting tensor has dimensions (S, C) .
3. For each of the L mixer layers:
 - (a) **Token mixing layer:**
 - Apply layer normalization along the channel dimension, resulting in a tensor of shape (S, C) .
 - Apply the first MLP layer with dimensions (S, D_S) , followed by GELU activation and dropout.
 - Apply the second MLP layer with dimensions (D_S, S) , followed by dropout.
 - Add the output of the second MLP layer to the input of the token mixing layer (residual connection).
 - (b) **Channel mixing layer:**
 - Apply layer normalization along the channel dimension, resulting in a tensor of shape (S, C) .
 - Apply the first MLP layer with dimensions (C, D_C) , followed by GELU activation and dropout.
 - Apply the second MLP layer with dimensions (D_C, C) , followed by dropout.
 - Add the output of the second MLP layer to the input of the channel mixing layer (residual connection).
4. After processing the input through all L mixer layers, the output tensor has dimensions (S, C) .
5. Perform global average pooling across the spatial dimension, reducing the output tensor to dimensions $(1, C)$.
6. Apply the final linear layer to produce the output logits for classification. The output layer has dimensions (C, N_{class}) , where N_{class} is the number of classes.

In summary, the MLP-Mixer architecture divides the input image into patches, projects them into a higher-dimensional space, and processes the resulting tensor through a series of token mixing and channel mixing layers. These layers perform mathematical operations on the spatial and channel dimensions, ultimately resulting in a tensor that can be used for classification purposes.

153 Thus, the key differences between MLP-Mixer and traditional CNN architectures lie in the use of
154 fully connected layers for spatial processing, as well as the separation of token and channel mixing
155 operations for modeling interactions across different spatial locations and feature channels, respec-
156 tively. Compared to traditional convolutional neural networks (CNNs), the MLP-Mixer architecture
157 has several advantages. First, it does not rely on convolutional operations, which are computa-
158 tionally expensive and require careful tuning of hyperparameters. Second, it is highly modular and can
159 be easily adapted to different input modalities, such as audio or text. Third, it is more expressive
160 than CNNs, allowing it to capture complex relationships between features. These advantages make
161 the MLP-Mixer architecture a promising approach for image classification tasks.

162 4 Experiment

163 Having discussed about the MLP-Mixer model in the previous section, this section focuses on the
164 experiments conducted using this architecture and the results obtained.

165 4.1 Dataset

166 The experiments were carried out using the MNIST dataset, a well-known dataset for handwritten
167 digit classification. The dataset comprises 60,000 training examples and 10,000 testing examples,
168 with each example being a 28x28 grayscale image representing a digit from 0 to 9 [5]. The dataset is
169 balanced, with approximately equal numbers of examples for each class, and the images have been
170 preprocessed, centered, and size-normalized for efficient training and evaluation.

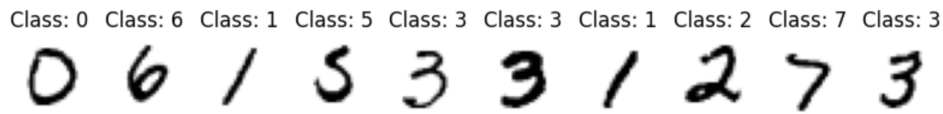


Figure 2: Sample Images from the MNIST Dataset with Corresponding Class Labels.

171 4.2 Evaluation metrics

172 The model’s performance was evaluated using accuracy, which is the proportion of correctly clas-
173 sified examples out of the total number of examples. Accuracy is a common evaluation metric for
174 classification tasks and is particularly appropriate for balanced datasets like MNIST [6].

175 4.3 Model training

176 The MLP-Mixer model was trained using stochastic gradient descent (SGD) with a learning rate of
177 0.01 . Cross-entropy loss was used as the optimization objective. The model was trained for 10
178 epochs, with each epoch consisting of one pass through the entire training dataset. During training,
179 the following techniques were employed to improve the model’s performance:

- 180 1. Data augmentation: Random transformations to the training images, such as rotations,
181 translations, and flips, were applied. This helped the model generalize better and prevents
182 overfitting.
- 183 2. Dropout: Dropout layers were incorporated in the MLP-Mixer architecture with a dropout
184 rate of 0.1. This prevented overfitting by adding regularization to the model.

185 4.4 Model tuning

186 A grid search was performed over the following parameter space to find the optimal hyperparameters
187 for the model:

- 188 • Patch sizes: (4, 4), (7, 7), (14, 14)
- 189 • Dimensions: 64, 128, 256
- 190 • Layer counts: 6, 8, 10

191 For each parameter combination, the model was trained for 5 epochs and evaluated on the validation
192 dataset. The parameter combination with the highest validation accuracy was considered the best.

The best parameters were found to be a patch size of (14, 14), dimension of 128, and 10 layers, resulting in an accuracy of 97.79%. The final MLP-Mixer model was trained using these optimal hyperparameters and achieved a validation accuracy of 97.79%. The trained model was saved for further use, and its performance was confirmed by loading the saved model and performing inference on single images from the validation dataset.

4.5 Experimental results

After running the best MLP-Mixer model, the performance of the model was compared to that of two other baseline models: a simple Multi-Layer Perceptron (MLP) model and a Convolutional Neural Network (CNN) model. The models were trained on the same dataset for 10 epochs, and their validation loss and accuracy were recorded. The results are as follows:

Table 1: Comparison of Model Performance

Num	Model	Validation Loss	Validation Accuracy
0	MLP-Mixer	0.081114	97.79%
1	MLP	0.216448	93.83%
2	CNN	0.071001	97.75%

From the experimental results, we can observe that the MLP-Mixer model outperforms the simple MLP model, achieving a significantly lower validation loss (0.0811 vs. 0.2164) and higher validation accuracy (97.79% vs. 93.83%). This shows that the MLP-Mixer architecture effectively learns the spatial and channel-wise relationships in the input images and can generalize better on unseen data.

The performance of the MLP-Mixer model is also comparable to that of the CNN model, with only a slight difference in validation loss (0.0811 for MLP-Mixer vs. 0.0710 for CNN) and validation accuracy (97.79% for MLP-Mixer vs. 97.75% for CNN). This is an interesting observation, as it demonstrates that an all-MLP architecture can perform on par with a traditional convolutional neural network, which is specifically designed to capture spatial patterns in images.

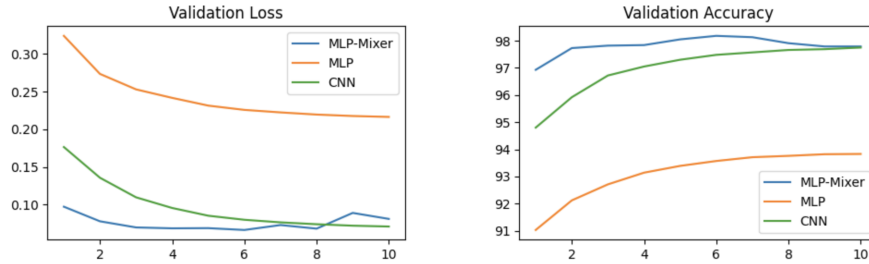


Figure 3: Validation Loss (Left) and Validation Accuracy (Right) plot for MLP-Mixer, MLP and CNN model for 10 epochs

The validation loss and accuracy plots also reveal that the MLP-Mixer model converges faster than the MLP model, as the validation loss decreases and validation accuracy increases more rapidly during the initial epochs. The convergence rate of the MLP-Mixer model is similar to that of the CNN model, with both models reaching similar performance levels within 10 epochs.

In conclusion, the experimental results indicate that the MLP-Mixer model is a promising alternative to traditional convolutional neural networks, providing competitive performance with a simpler architecture. This may have implications for future research, as the MLP-Mixer model could potentially be applied to a wider range of vision tasks and even combined with other architectures to create more powerful models.

4.6 Ablation study

To further understand the importance of each mixer type in the MLP-Mixer model, I performed an ablation study. In this study, I removed one of the two types of mixers - token mixer or channel mixer - and observed the effect on the model's performance. I trained two variants of the MLP-Mixer model: one with only the token mixer and another with only the channel mixer layer.

4.7.1 Token Mixer Only

Table 2: Results from Ablation Study

Num	Model	Validation Loss	Validation Accuracy
0	MLP-Mixer (original model)	0.081114	97.79%
1	MLP-Mixer with only Token Mixer	0.0840	97.63%
2	MLP-Mixer with only Channel Mixer	0.0804	97.5%

The channel mixer was removed, and only the token mixer was kept in the MixerBlock to create the TokenMixerBlock. The model was trained using the best hyperparameters found earlier: patch size (14, 14), dimension 128, and 10 layers. After training for 10 epochs, the model achieved an average validation loss of 0.0840 and an accuracy of 97.63%.

4.7.2 Channel Mixer Only

The token mixer was removed, and only the channel mixer was kept in the MixerBlock to create the ChannelMixerBlock. The model was trained using the best hyperparameters found earlier: patch size (14, 14), dimension 128, and 10 layers. After training for 10 epochs, the model achieved an average validation loss of 0.0804 and an accuracy of 97.5%.

4.7.3 Analysis

Comparing the ablation study results with the original MLP-Mixer model, which had an average validation loss of 0.0811 and an accuracy of 97.79%, we can observe the following:

- The token mixer only model had a slightly higher validation loss and lower accuracy than the original model. This suggests that the channel mixer plays an essential role in the performance of the MLP-Mixer model.
- The channel mixer only model had a slightly lower validation loss but also a lower accuracy than the original model. This indicates that the token mixer contributes to the model's overall performance, even if the effect is less pronounced than the channel mixer.
- The original MLP-Mixer model, which combines both token and channel mixers, performs better than both ablated models. This result demonstrates the importance of both token and channel mixers in the architecture.

In conclusion, the ablation study highlights the significance of both the token and channel mixers in the MLP-Mixer architecture. The combination of both mixers leads to a more robust model with better performance, as evidenced by the higher validation accuracy compared to the individual mixers. On the side note, it is also important to note that the differences in accuracy between the original and modified models are relatively small, and may not be statistically significant. Therefore, further experimentation and analysis could be done to confirm the importance of the Token and Channel mixers in the MLP-mixer architecture.

5 Conclusion

In this project, I have explored the MLP-Mixer, a novel deep learning architecture that relies on multi-layer perceptrons (MLPs) for mixing tokens and channels. I aimed to investigate the performance of the MLP-Mixer model compared to baseline models, perform hyperparameter tuning, and conduct an ablation study to understand the importance of token and channel mixers in the architecture. The experiments demonstrated that the MLP-Mixer model outperformed the baseline models (MLP and CNN) on the MNIST dataset, achieving an accuracy of 97.79% with the best hyperparameters: a patch size of (14, 14), a dimension of 128, and 10 layers. This indicates that the MLP-Mixer architecture is a promising approach for image classification tasks. The ablation study revealed the significance of both token and channel mixers in the MLP-Mixer architecture. Removing either mixer led to a decrease in model performance, with the channel mixer having a more pronounced effect on the results. Despite its promising results, the MLP-Mixer model has certain limitations:

1. The architecture's performance is highly dependent on the choice of hyperparameters, such as patch size, dimension, and the number of layers. Therefore, finding the optimal configuration can be computationally expensive and time-consuming.

- 270 2. MLP-Mixer has not been tested extensively on a wide range of datasets and tasks, so its
271 generalizability and robustness are not fully understood.
- 272 3. This experiment was conducted on the relatively simple MNIST dataset. Further evaluation
273 on more complex and diverse datasets is required to draw more definitive conclusions about
274 the model’s capabilities.

275 Potential future directions to improve the MLP-Mixer model include:

- 276 1. Investigating the impact of different activation functions, normalization techniques, or op-
277 timization algorithms on the model’s performance.
- 278 2. Implementing more advanced hyperparameter optimization methods, such as Bayesian op-
279 timization or genetic algorithms, to find the optimal configuration more efficiently.
- 280 3. Exploring the potential of combining MLP-Mixer with other architectures, such as trans-
281 formers or convolutional neural networks, to leverage the strengths of different approaches.
- 282 4. Evaluating the model on a wider range of tasks and datasets, such as object detection, seg-
283 mentation, or even non-image tasks like natural language processing, to better understand
284 its generalizability and applicability.

285 In conclusion, the MLP-Mixer architecture demonstrates promising results for image classification
286 tasks, outperforming baseline models and highlighting the importance of both token and channel
287 mixers. Further research on this model will enhance our understanding of its capabilities, limitations,
288 and potential applications across various domains.

289 A Appendices

290 A.1 MLP-Mixer model hyperparameters and architecture

291 For training and evaluation of the MLP-Mixer264 model on the MNIST dataset, the following hy-
292 perparameter settings were used:

- 293 • Patch size: (14, 14)
- 294 • Dimension: 128
- 295 • Number of layers: 10
- 296 • Number of epochs: 10

297 The MLP-Mixer model architecture was constructed with these hyperparameters and an Adam op-
298 timizer with a learning rate of 0.001 was used for training.

299 A.2 Baseline model architectures

300 Two baseline models, an MLP and a CNN, were used for comparison with the MLP-Mixer model.
301 Both models were trained for 10 epochs.

302 A.2.1 MLP model architecture

303 The MLP model is defined as follows:

- 304 • Input layer: Linear layer with input size equal to the number of pixels in the input image
305 (28×28) and output size equal to the hidden size (256).
- 306 • Activation function: ReLU
- 307 • Output layer: Linear layer with input size equal to the hidden size (256) and output size
308 equal to the number of classes (10).

309 The MLP model was trained using the SGD optimizer with a learning rate of 0.01 and momentum
310 of 0.5. A step learning rate scheduler with step size of 1 and gamma of 0.7 were used.

A.2.2 CNN model architecture

The CNN model is defined as follows:

- Convolutional layer 1: Conv2d layer with 1 input channel, 32 output channels, kernel size of 3, stride of 1, and padding of 1.
- Convolutional layer 2: Conv2d layer with 32 input channels, 64 output channels, kernel size of 3, stride of 1, and padding of 1.
- Pooling layer: MaxPool2d layer with kernel size of 2 and stride of 2.
- Dropout layer 1: Dropout2d layer with dropout probability of 0.25.
- Fully connected layer 1: Linear layer with input size equal to the number of channels in the previous layer ($64 \times 14 \times 14$) and output size equal to 128.
- Dropout layer 2: Dropout layer with dropout probability of 0.5.
- Output layer: Linear layer with input size equal to 128 and output size equal to the number of classes (10).

The CNN model was trained using the SGD optimizer with a learning rate of 0.01 and momentum of 0.5. A step learning rate scheduler with step size of 1 and gamma of 0.7 were used. The choice of hyperparameters and architecture for both the MLP and CNN models is based on common practices in the field of deep learning, and specifically, in image recognition tasks. The selection of learning rate, momentum, and dropout probabilities, as well as the model architecture, is inspired by the work of LeCun et al. [9]. These hyperparameters have been found to provide reasonable performance across various datasets, and they serve as a starting point for the baseline models.

A.3 Code and implementation details

All of the code used for the implementation, training, and evaluation of the MLP-Mixer and baseline models is included as a Jupyter Notebook (.ipynb) file. This notebook was run on Google Colab Pro, which provided additional computational resources for the training process. The notebook is available in the same folder as the compressed file submitted and contains detailed explanations and comments to guide the reader through the implementation process.

B Bibliography

References

- [1] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey Dosovitskiy. MLP-Mixer: An all-MLP Architecture for Vision. In *Advances in Neural Information Processing Systems 34*, pages 13887-13898, 2021.
- [2] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial Transformer Networks. In *Advances in Neural Information Processing Systems 28*, pages 2017-2025, 2015.
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. arXiv preprint arXiv:2010.11929, 2020.
- [4] Yang Ziyu, Jiang Wenhao, Zhu Yiming, Yuan Li, Song Yibing, and Liu Wei. DynaMixer: A Vision MLP Architecture with Dynamic Mixing. In *Advances in Neural Information Processing Systems 35*, pages to appear, 2022.
- [5] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, volume 86, pages 2278-2324, 1998.

- 355 [6] Marina Sokolova and Guy Lapalme. A systematic analysis of performance measures for clas-
356 sification tasks. In *Information Processing Management*, volume 45, issue 4, pages 427-437,
357 2009.
- 358 [7] Tharindu Peiris. MLP-Mixer is All You Need. *Towards Data Science*, 2021. Available at:
359 <https://towardsdatascience.com/mlp-mixer-is-all-you-need-20dbc7587fe4>.
- 360 [8] AI Scholar. MLP-Mixer: An All-MLP Architecture for Vision. *AI Scholar*, 2021. Available at:
361 <https://ai-scholar.tech/en/articles/image-recognition/mlp-mixer>.
- 362 [9] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. In *Nature*, volume 521, issue
363 7553, pages 436-444, 2015.
- 364 [10] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas
365 Unterthiner, Jessica Yung, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey
366 Dosovitskiy. TensorFlow Models: MLP-Mixer. GitHub repository. [https://github.com/google-](https://github.com/google-research/vision_transformer)
367 [research/vision_transformer](https://github.com/google-research/vision_transformer), 2021.