

# **MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE GWALIOR**

(A Govt. Aided UGC Autonomous Institute Affiliated to RGPV, Bhopal)

**NAAC Accredited with A++ Grade**



**Project Report**

**on**

**Image Classification Using Vision Transformer**

**(270506)**

**Submitted By:**

**Krishna Sharma (0901AD211025)**

**Shiv Shrivastava (0901AD211057)**

**Faculty Mentor:**

**Dr. Pawan Dubey**

**CENTRE FOR ARTIFICIAL INTELLIGENCE**

**MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE**

**GWALIOR - 474005 (MP) est. 1957**

**JULY-DEC. 2023**

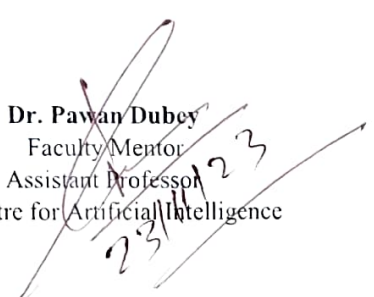
# MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE GWALIOR

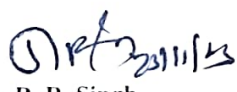
(A Govt. Aided UGC Autonomous Institute Affiliated to RGPV, Bhopal)

NAAC Accredited with A++ Grade

## CERTIFICATE

This is certified that **Krishna Sharma (0901AD211025)** and **Shiv Shrivastava (0901AD211057)** has submitted the project report titled "**Image Classification Using Vision Transformer**" under the mentorship of **Dr. Pawan Dubey**, in partial fulfilment of the requirement for the award of degree of **Bachelor of Technology in Artificial Intelligence & Data Science** from **Madhav Institute of Technology and Science, Gwalior**.

  
**Dr. Pawan Dubey**  
Faculty Mentor  
Assistant Professor  
Centre for Artificial Intelligence

  
**Dr. R. R. Singh**  
Coordinator  
Centre for Artificial Intelligence

# MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE GWALIOR

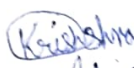
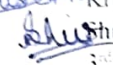
(A Govt. Aided UGC Autonomous Institute Affiliated to RGPV, Bhopal)

NAAC Accredited with A++ Grade

## DECLARATION

I hereby declare that the work being presented in this project report, for the partial fulfilment of requirement for the award of the degree of **Bachelor of Technology in Artificial Intelligence & Data Science** at **Madhav Institute of Technology & Science, Gwalior** is an authenticated and original record of my work under the mentorship of **Dr. Pawan Dubey, Assistant Professor, Centre For Artificial Intelligence**.

I declare that I have not submitted the matter embodied in this report for the award of any degree or diploma anywhere else.

 Krishna Sharma (0901AD211025)  
 Shiv Shrivastava (0901AD211057)  
3<sup>rd</sup> Year  
Centre for Artificial Intelligence

# MADHAV INSTITUTE OF TECHNOLOGY & SCIENCE GWALIOR

(A Govt. Aided UGC Autonomous Institute Affiliated to RGPV, Bhopal)

NAAC Accredited with A++ Grade

## ACKNOWLEDGEMENT

The full semester project has proved to be pivotal to my career. I am thankful to my institute, **Madhav Institute of Technology and Science** to allow me to continue my disciplinary/interdisciplinary project as a curriculum requirement, under the provisions of the Flexible Curriculum Scheme (based on the AICTE Model Curriculum 2018), approved by the Academic Council of the institute. I extend my gratitude to the Director of the institute, **Dr. R. K. Pandit** and Dean Academics, **Dr. Manjaree Pandit** for this.

I would sincerely like to thank my department, **Centre for Artificial Intelligence**, for allowing me to explore this project. I humbly thank **Dr. R. R. Singh**, Coordinator, Centre for Artificial Intelligence, for his continued support during the course of this engagement, which eased the process and formalities involved.

I am sincerely thankful to my faculty mentors. I am grateful to the guidance of **Dr. Pawan Dubey**, Assistant Professor, Centre For Artificial Intelligence, for his continued support and guidance throughout the project. I am also very thankful to the faculty and staff of the department.



Krishna Sharma (0901AD211025)

Shiv Shrivastava (0901AD211057)

3<sup>rd</sup> Year

Centre for Artificial Intelligence

## ABSTRACT

With the rapid evolution of deep learning techniques, vision transformers have emerged as a promising approach for image classification tasks. This explores the application of vision transformers on two distinct datasets: a proprietary dataset containing diverse plant images and the well-known CIFAR-100 dataset. The objective is to evaluate the performance of vision transformers in the context of plant species recognition and general object classification. The study begins with a comprehensive review of vision transformer architecture and its potential advantages in handling image classification tasks. The proposed model is trained and fine-tuned on the custom plant dataset, which consists of a variety of plant species captured under different environmental conditions. To assess the model's generalization capabilities, it is further evaluated on the CIFAR-100 dataset, which encompasses a broader range of object categories. The experimental results demonstrate the effectiveness of the vision transformer in achieving competitive accuracy on both datasets. The model's ability to capture intricate features of plant species suggests its potential utility in agricultural and environmental monitoring applications. Additionally, the generalization performance on CIFAR-100 highlights the versatility of the vision transformer architecture across diverse image classification tasks. Furthermore, the research investigates the impact of key hyperparameters, such as patch size, model depth, and learning rate, on the performance of the vision transformer. The findings contribute insights into optimizing the model for specific datasets and offer practical guidance for researchers and practitioners working on image classification tasks. In conclusion, this showcases the successful application of vision transformers on a custom plant dataset and the CIFAR-100 benchmark. The results underscore the adaptability of vision transformers in handling distinct image classification challenges and open avenues for further exploration in the domain of plant science and computer vision.

**Keyword:** Vision Transformer, Image Classification, Deep Learning, Plant Dataset, CIFAR-100, Convolutional Neural Network (CNN), Hyperparameter Tuning, Computer Vision, Patch Size, Learning Rate, Object Recognition



## सार

गहन शिक्षण तकनीकों के तेजी से विकास के साथ, दृष्टि ट्रांसफार्मर छवि वर्गीकरण कार्यों के लिए एक आशाजनक दृष्टिकोण के रूप में उभरे हैं। यह दो अलग-अलग डेटासेट पर विज़न ट्रांसफार्मर के अनुप्रयोग की पड़ताल करता है: एक मालिकाना डेटासेट जिसमें विविध पौधों की छवियां और प्रसिद्ध CIFAR-100 डेटासेट शामिल हैं। इसका उद्देश्य पौधों की प्रजातियों की पहचान और सामान्य वस्तु वर्गीकरण के संदर्भ में दृष्टि ट्रांसफार्मर के प्रदर्शन का मूल्यांकन करना है। अध्ययन दृष्टि ट्रांसफार्मर वास्तुकला की व्यापक समीक्षा और छवि वर्गीकरण कार्यों को संभालने में इसके संभावित लाभों के साथ शुरू होता है। प्रस्तावित मॉडल को कस्टम प्लांट डेटासेट पर प्रशिक्षित और परिष्कृत किया गया है, जिसमें विभिन्न पर्यावरणीय परिस्थितियों में पकड़ी गई विभिन्न प्रकार की पौधों की प्रजातियां शामिल हैं। मॉडल की सामान्यीकरण क्षमताओं का आकलन करने के लिए, इसे CIFAR-100 डेटासेट पर आगे मूल्यांकन किया जाता है, जिसमें ऑब्जेक्ट श्रेणियों की एक विस्तृत श्रृंखला शामिल होती है। प्रयोगात्मक परिणाम दोनों डेटासेट पर प्रतिस्पर्धी सटीकता प्राप्त करने में दृष्टि ट्रांसफार्मर की प्रभावशीलता को प्रदर्शित करते हैं। पौधों की प्रजातियों की जटिल विशेषताओं को पकड़ने की मॉडल की क्षमता कृषि और पर्यावरण निगरानी अनुप्रयोगों में इसकी संभावित उपयोगिता का सुझाव देती है। इसके अतिरिक्त, CIFAR-100 पर सामान्यीकरण प्रदर्शन विविध छवि वर्गीकरण कार्यों में दृष्टि ट्रांसफार्मर वास्तुकला की बहुमुखी प्रतिभा पर प्रकाश डालता है। इसके अलावा, अनुसंधान दृष्टि ट्रांसफार्मर के प्रदर्शन पर प्रमुख हाइपरपैरामीटर, जैसे पैच आकार, मॉडल गहराई और सीखने की दर के प्रभाव की जांच करता है। निष्कर्ष विशिष्ट डेटासेट के लिए मॉडल को अनुकूलित करने में अंतर्दृष्टि प्रदान करते हैं और छवि वर्गीकरण कार्यों पर काम करने वाले शोधकर्ताओं और चिकित्सकों के लिए व्यावहारिक मार्गदर्शन प्रदान करते हैं। अंत में, यह पेपर कस्टम प्लांट डेटासेट और CIFAR-100 बेंचमार्क पर विज़न ट्रांसफार्मर के सफल अनुप्रयोग को दर्शाता है। परिणाम विशिष्ट छवि वर्गीकरण चुनौतियों से निपटने में दृष्टि ट्रांसफार्मर की अनुकूलनशीलता को रेखांकित करते हैं और पादप विज्ञान और कंप्यूटर दृष्टि के क्षेत्र में आगे की खोज के लिए रास्ते खोलते हैं।

# TABLE OF CONTENTS

	Page No.
Certificate	2
Declaration	3
Acknowledgement	4
Abstract	5
संक्षेप	6
List Of Figures	8
<b><u>Chapter 1: Introduction</u></b>	9
<b><u>Chapter 2: Literature Review</u></b>	10
<b><u>Chapter 3: Material and Methods</u></b>	11
3.1 Dataset	
3.1.1 CIFAR - 100 Dataset	
3.1.2 Custom Millets Plants Dataset	
3.2 Vision Transformer & Patch Encoder	
3.3 Experimental Setup	
3.4 Flow Chart	
<b><u>Chapter 4: Result</u></b>	14
3.1 Result on CIFAR - 100 Dataset	
3.2 Result on Custom Millets Plants Dataset	
<b><u>Chapter 5: Conclusion</u></b>	20
<b><u>References</u></b>	21

LIST OF FIGURES

Figure Number	Figure caption	Page No.
1	Vision Transformer With Multihead self attention	12
2	Patch Encoder(Visualisation and Patches formation)	13
3	Flowchart	14
4	Result on millets custom Dataset	15-16
5	Result on CIFAR - 100 Dataset	16-17



## Chapter 1: Introduction

The field of Image classification is a task in computer vision where the goal is to categorize an input image into one of several predefined classes or categories. It is a fundamental problem in image analysis and pattern recognition. The process involves training a model using a set of labeled images, where each image is associated with a specific class or category. This has witnessed remarkable advancements in recent years, driven primarily by the advent of deep learning architectures. Among these architectures, the Vision Transformer (ViT) has gained prominence for its unique approach to processing images., we delve into the application of Vision Transformers for image classification tasks, specifically focusing on two distinct datasets: a proprietary plant dataset and the CIFAR-100 dataset. The Vision Transformer, introduced by Vaswani et al. in 2017 [1], represents a departure from the conventional convolutional neural network (CNN) paradigm. Instead of relying on convolutional layers, the ViT operates on image patches, treating them as sequential inputs to a transformer architecture originally designed for natural language processing tasks. This departure from the spatial hierarchy of traditional CNNs introduces a novel way of capturing global dependencies among image features, offering potential advantages in handling diverse and complex datasets. Our motivation for this study stems from the need for effective image classification models in the domain of plant science. Monitoring and categorizing plant species play a crucial role in various applications, including agriculture, environmental conservation, and ecosystem management. The custom plant dataset employed in this research encompasses a wide array of plant species captured under varying environmental conditions. The diversity in this dataset poses a challenge for traditional image classification models, making it an ideal candidate for evaluating the efficacy of Vision Transformers. In addition to the plant dataset, we evaluate the ViT model on the CIFAR-100 dataset, a benchmark in the field of object recognition. CIFAR-100 consists of 100 classes, each containing 600 images, making it a challenging testbed for any image classification model. The inclusion of CIFAR-100 in our study allows us to assess the generalization capabilities of the Vision Transformer across a broader spectrum of object categories. The ViT model undergoes a training and fine-tuning process on the custom plant dataset to specialize in recognizing various plant species. Subsequently, it is subjected to evaluation on the CIFAR-100 dataset to gauge its adaptability and performance in a more generalized setting. The outcomes of this not only contribute to the growing body of knowledge on Vision Transformers but also provide valuable insights into the potential applications of such models in the intersection of computer vision and plant science. One of the primary objectives of this is to investigate the impact of key hyperparameters on the performance of the Vision Transformer. The patch size, model depth, and learning rate are systematically explored to optimize the model for both the custom plant dataset and CIFAR-100. This exploration is essential not only for achieving peak performance but also for providing practical guidance to researchers and practitioners dealing with image classification tasks using Vision Transformers.

## Chapter 2: Literature Review

The literature surrounding image classification has witnessed a paradigm shift with the introduction of Vision Transformers (ViTs). Vision Transformers represent a departure from the established convolutional neural network (CNN) architecture, offering a novel approach to capturing spatial dependencies in images. The transformer architecture, initially designed for natural language processing, has been adapted successfully to process image data, demonstrating state-of-the-art performance in various computer vision tasks. Vaswani et al. (2017) introduced the transformer architecture for sequence-to-sequence tasks in natural language processing. Building upon this, Dosovitskiy et al. (2020) pioneered the application of transformers in computer vision with the Vision Transformer (ViT). ViT divides an image into fixed-size patches, linearly embedding them before processing through a transformer encoder. This departure from the grid-like receptive fields of CNNs allows ViT to capture long-range dependencies and interactions among image patches, leading to impressive performance in image classification tasks. Several studies have explored the effectiveness of Vision Transformers in comparison to traditional CNNs. Radford et al. (2021) demonstrated that ViTs can achieve competitive performance on various image classification benchmarks. Notably, ViTs have showcased a remarkable ability to scale with increased model size, outperforming CNNs in terms of both accuracy and efficiency. In the realm of plant science, where accurate species identification is crucial, the application of deep learning models has gained traction. Deep learning models, particularly CNNs, have been successfully applied to plant species recognition tasks (Mehdipour Ghazi et al., 2017). However, the potential of Vision Transformers in this domain remains relatively unexplored. Our research aims to bridge this gap by evaluating the performance of ViTs on a diverse and proprietary plant dataset.

The evaluation of models on benchmark datasets is a common practice to assess their generalization capabilities. The CIFAR-100 dataset has been a popular choice for this purpose. It consists of 100 object classes, each containing 600 images, posing a challenging test for image classification models. The use of CIFAR-100 in our study provides a benchmark for comparing the performance of Vision Transformers with existing literature on CNNs and other deep learning models. While Vision Transformers have demonstrated their efficacy in various computer vision tasks, including image classification, their application to specific domains, such as plant species recognition, demands careful evaluation. Our literature review sets the stage for the exploration of Vision Transformers on a custom plant dataset, emphasizing the need to assess their performance, optimize hyperparameters, and provide valuable insights for researchers and practitioners in both computer vision and plant science.

## Chapter 3: Materials and Methods

### 3.1 Dataset

In our research, we leverage a diverse set of datasets to comprehensively evaluate the proposed approach for image classification on both well-established benchmarks and a custom dataset focused on millets plants. This dual-pronged strategy aims to assess the model's generalization across widely recognized datasets and its adaptability to a context-specific agricultural scenario.

#### 3.1.1 CIFAR-100 Dataset

- Origin: Keras Library
- Geographic Origin: Varied, synthetically generated dataset
- Characteristics: CIFAR-100 is a well-known dataset containing 100 classes, each with 600 images, covering a broad spectrum of object categories. The dataset is designed to challenge image classification models with diverse and complex visual concepts.

#### 3.1.2 Custom Millets Plant Dataset

- Characteristics: This custom dataset is curated specifically for millets plant species, encompassing multiple varieties and conditions. The dataset includes images capturing various potential diseases affecting millets plants.

### 3.2 Vision Transformer & Patch Encoder

The provided code implements a Vision Transformer (ViT) for image classification, a paradigm that has demonstrated remarkable success in natural language processing. Dosovitskiy et al. (2020) introduced the ViT model, building on the transformer architecture pioneered by Vaswani et al. (2017). The ViT consists of self-attention blocks and multilayer perceptron (MLP) networks with a linear projection and positional embedding mechanism.

In the ViT architecture (Fig. 1), an image is initially split into fixed-size non-overlapping patches, which are then flattened and transformed into lower-dimensional representations. Each patch undergoes a learnable linear transformation to generate a linear projection and positional embedding. These representations are passed through a stack of  $N$  transformer blocks, each comprising multi-head self-attention (MHA) and an MLP. Each transformer block includes normalization layers, residual connections, and a skip connection between the input and the output of both MHA and MLP.

The self-attention mechanism, MHA, is applied to each patch separately. In MHA, input vectors are transformed into three separate vectors: Q (Query), K (Key), and V (Value). The dot product between Q and K generates a score matrix, which is then subjected to a softmax activation. The resulting self-attention matrices are combined and processed through a linear layer, feeding into the regression head for classification. Normalization is applied to avoid issues with excessively large dot products during training.

The ViT model's transformer blocks enhance semantic similarity across different image locations, contributing to effective classification. The number of MHA in a transformer encoder is a tunable hyperparameter, providing flexibility based on application data.

The code includes a comprehensive ViT model, complete with data augmentation, patch extraction, patch encoding, and multiple transformer blocks. The implementation uses TensorFlow and Keras, incorporating TensorFlow Addons for the AdamW optimizer with weight decay.

This ViT model is then evaluated on image classification tasks using a combination of publicly available datasets, including CIFAR-100, and a custom dataset focused on millets plants. The training history, including accuracy metrics, is stored for analysis. The provided code serves as a foundational framework for experimenting with Vision Transformers on diverse image classification challenges.

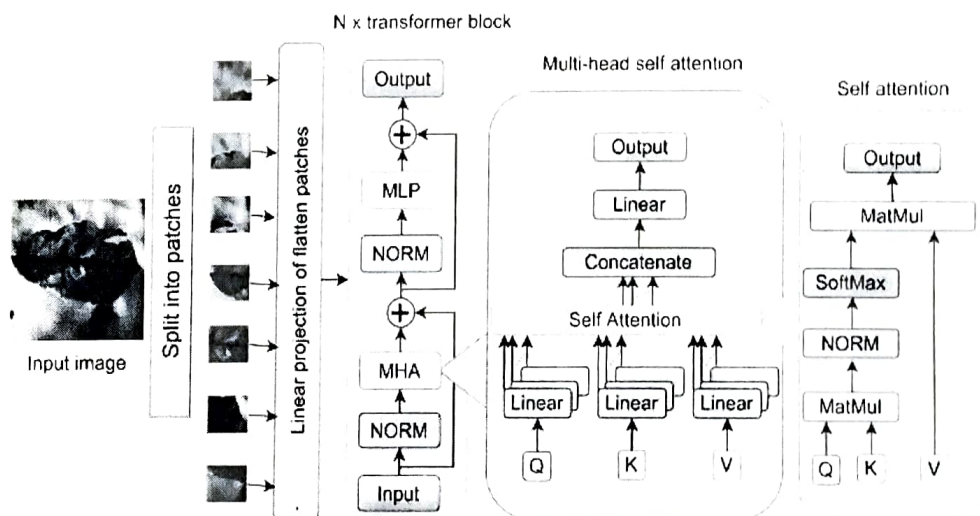


Fig. 1. ViT block with multi-head self-attention units.

In Fig. 2, the schematic representation illustrates the crucial steps of patch visualization and encoding within the Vision Transformer (ViT) framework. The process begins with the selection of a random image from the



CFAR-100 dataset. This image is then resized to a standardized **image\_size**, after which the **Patches** class is employed to extract non-overlapping patches from the resized image. These patches, illustrated in the diagram, showcase the decomposition of the original image into smaller, distinct components. This visual representation aids in comprehending how the ViT model initially processes and segments input images, setting the stage for subsequent attention-based operations:

Following patch visualization, the diagram highlights the role of the PatchEncoder component. This crucial step involves encoding the extracted patches, preparing them for effective processing by the subsequent transformer blocks in the ViT model. The PatchEncoder consists of two primary operations: a learnable linear transformation and positional embedding. The encoded patches, visualized in the diagram, are then passed to the transformer blocks, contributing essential spatial and positional context to the ViT model's overall understanding of the input image. This visual representation in Fig. 2 provides a concise overview of the patch visualization and encoding stages, offering insights into the initial processing steps crucial for ViT-based image classification.

Image size: 72 X 72  
Patch size: 6 X 6  
Patches per image: 144  
Elements per patch: 108

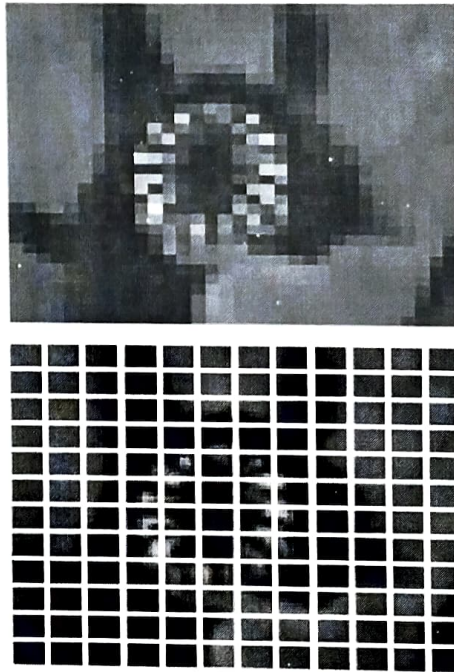


Fig. 2. Patch Encoder(Visualization and patches formation)

### 3.3 Experimental Setup

The experiments were conducted on The ASUS TUF Gaming laptop features a robust hardware setup, including an AMD Ryzen processor and a dedicated GPU for accelerated computations. The laptop runs on the Ubuntu 18.04 LTS operating system, ensuring compatibility with the experiment environment. The AMD Ryzen CPU, with clock speeds ranging from 2.25 to 3.4 GHz, and 512 GB of RAM, provide ample computing power and memory capacity for running the experiments efficiently.

The proposed model, along with other selected models for comparison, is implemented using the Keras framework, harnessing the capabilities of the laptop's hardware. The laptop is equipped with NVIDIA CUDA v11.5 and the cuDNN v8.3 library, facilitating GPU acceleration for deep learning tasks. This experimental setup on the ASUS TUF Gaming laptop demonstrates the feasibility and performance of the proposed model in a resource-constrained environment, extending the applicability of the models beyond high-end workstations to more widely accessible computing platforms.

### 3.4 Flow Chart

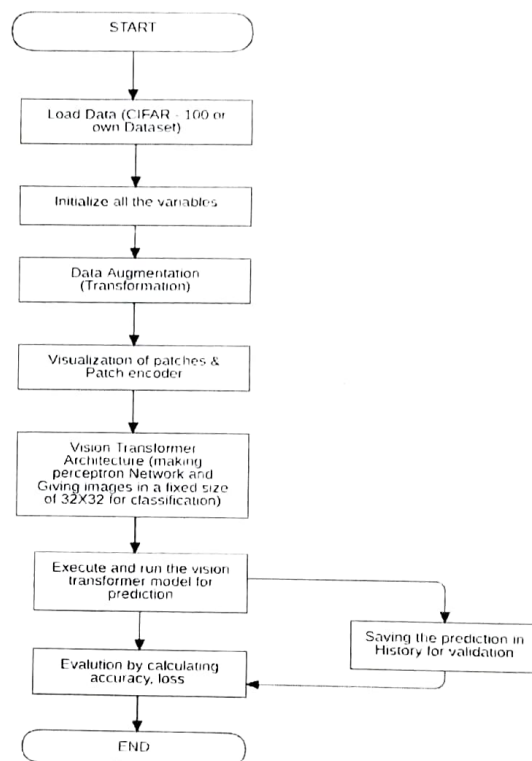


Fig. 3. Process Flow

## Chapter 4: Result

### 4.1 Result on CIFAR - 100 Dataset

After 100 epochs, the ViT model achieves around 55% accuracy and 82% top-5 accuracy on the test data. These are not competitive results on the CIFAR-100 dataset, as a ResNet50V2 trained from scratch on the same data can achieve 67% accuracy (Fig. 4).

epoch 1/100	176.176 [-----]	33s 130ms/step	loss: 1.6863	accuracy: 0.0294	top 5 accuracy: 0.1117	val_loss: 1.9661	val_accuracy: 0.0992	val_top 5 accuracy: 0.3859
epoch 2/100	176.176 [-----]	22s 127ms/step	loss: 1.0162	accuracy: 0.0865	top 5 accuracy: 0.2683	val_loss: 1.5691	val_accuracy: 0.1639	val_top 5 accuracy: 0.4279
epoch 3/100	176.176 [-----]	22s 127ms/step	loss: 1.1113	accuracy: 0.1254	top 5 accuracy: 0.3535	val_loss: 1.3453	val_accuracy: 0.1979	val_top 5 accuracy: 0.4770
epoch 4/100	176.176 [-----]	23s 128ms/step	loss: 1.5411	accuracy: 0.1541	top 5 accuracy: 0.4121	val_loss: 1.1925	val_accuracy: 0.2174	val_top 5 accuracy: 0.5109
epoch 5/100	176.176 [-----]	22s 127ms/step	loss: 1.1719	accuracy: 0.1847	top 5 accuracy: 0.4572	val_loss: 1.1043	val_accuracy: 0.2389	val_top 5 accuracy: 0.5320
epoch 6/100	176.176 [-----]	22s 127ms/step	loss: 1.2589	accuracy: 0.2057	top 5 accuracy: 0.4906	val_loss: 1.0319	val_accuracy: 0.2782	val_top 5 accuracy: 0.5758
epoch 7/100	176.176 [-----]	22s 127ms/step	loss: 1.1185	accuracy: 0.2331	top 5 accuracy: 0.5273	val_loss: 1.0072	val_accuracy: 0.2972	val_top 5 accuracy: 0.5949
epoch 8/100	176.176 [-----]	22s 127ms/step	loss: 1.9902	accuracy: 0.2583	top 5 accuracy: 0.5558	val_loss: 1.1207	val_accuracy: 0.3188	val_top 5 accuracy: 0.6254
epoch 9/100	176.176 [-----]	22s 127ms/step	loss: 2.8828	accuracy: 0.3000	top 5 accuracy: 0.5817	val_loss: 1.0796	val_accuracy: 0.3244	val_top 5 accuracy: 0.6402
epoch 10/100	176.176 [-----]	23s 128ms/step	loss: 2.7824	accuracy: 0.2987	top 5 accuracy: 0.6110	val_loss: 1.1580	val_accuracy: 0.3454	val_top 5 accuracy: 0.6588
epoch 11/100	176.176 [-----]	23s 130ms/step	loss: 2.6743	accuracy: 0.3209	top 5 accuracy: 0.6333	val_loss: 1.2000	val_accuracy: 0.3544	val_top 5 accuracy: 0.6728
epoch 12/100	176.176 [-----]	23s 130ms/step	loss: 2.5800	accuracy: 0.3411	top 5 accuracy: 0.6522	val_loss: 1.1900	val_accuracy: 0.3798	val_top 5 accuracy: 0.6879
epoch 13/100	176.176 [-----]	23s 128ms/step	loss: 2.5014	accuracy: 0.3558	top 5 accuracy: 0.6671	val_loss: 1.1404	val_accuracy: 0.3960	val_top 5 accuracy: 0.7062
epoch 14/100	176.176 [-----]	22s 128ms/step	loss: 2.4207	accuracy: 0.3728	top 5 accuracy: 0.6905	val_loss: 1.1130	val_accuracy: 0.4032	val_top 5 accuracy: 0.7049
epoch 15/100	176.176 [-----]	23s 128ms/step	loss: 2.1171	accuracy: 0.3932	top 5 accuracy: 0.7093	val_loss: 1.2447	val_accuracy: 0.4136	val_top 5 accuracy: 0.7202
epoch 16/100	176.176 [-----]	23s 128ms/step	loss: 2.2850	accuracy: 0.4017	top 5 accuracy: 0.7201	val_loss: 1.2101	val_accuracy: 0.4212	val_top 5 accuracy: 0.7249
epoch 17/100	176.176 [-----]	22s 127ms/step	loss: 2.1822	accuracy: 0.4204	top 5 accuracy: 0.7376	val_loss: 1.1440	val_accuracy: 0.4344	val_top 5 accuracy: 0.7421
epoch 18/100	176.176 [-----]	22s 128ms/step	loss: 2.1485	accuracy: 0.4384	top 5 accuracy: 0.7479	val_loss: 1.1054	val_accuracy: 0.4431	val_top 5 accuracy: 0.7594
epoch 19/100	176.176 [-----]	22s 128ms/step	loss: 2.0717	accuracy: 0.4464	top 5 accuracy: 0.7618	val_loss: 1.0718	val_accuracy: 0.4504	val_top 5 accuracy: 0.7579
epoch 20/100	176.176 [-----]	22s 127ms/step	loss: 2.0031	accuracy: 0.4605	top 5 accuracy: 0.7711	val_loss: 1.0286	val_accuracy: 0.4619	val_top 5 accuracy: 0.7654
epoch 21/100	176.176 [-----]	22s 127ms/step	loss: 1.9850	accuracy: 0.4700	top 5 accuracy: 0.7820	val_loss: 1.0235	val_accuracy: 0.4640	val_top 5 accuracy: 0.7619
epoch 22/100	176.176 [-----]	22s 127ms/step	loss: 1.8686	accuracy: 0.4835	top 5 accuracy: 0.7904	val_loss: 1.0661	val_accuracy: 0.4786	val_top 5 accuracy: 0.7659
epoch 23/100	176.176 [-----]	22s 127ms/step	loss: 1.8364	accuracy: 0.4932	top 5 accuracy: 0.8010	val_loss: 1.0369	val_accuracy: 0.4828	val_top 5 accuracy: 0.7742
epoch 24/100	176.176 [-----]	22s 128ms/step	loss: 1.8167	accuracy: 0.5034	top 5 accuracy: 0.8089	val_loss: 1.0750	val_accuracy: 0.4769	val_top 5 accuracy: 0.7728
epoch 25/100	176.176 [-----]	22s 128ms/step	loss: 1.7788	accuracy: 0.5124	top 5 accuracy: 0.8174	val_loss: 1.0732	val_accuracy: 0.4829	val_top 5 accuracy: 0.7954
epoch 26/100	176.176 [-----]	23s 128ms/step	loss: 1.7437	accuracy: 0.5187	top 5 accuracy: 0.8206	val_loss: 1.0732	val_accuracy: 0.4782	val_top 5 accuracy: 0.7772
epoch 27/100	176.176 [-----]	23s 128ms/step	loss: 1.6829	accuracy: 0.5300	top 5 accuracy: 0.8287	val_loss: 1.0309	val_accuracy: 0.4928	val_top 5 accuracy: 0.7911
epoch 28/100	176.176 [-----]	23s 129ms/step	loss: 1.6947	accuracy: 0.5400	top 5 accuracy: 0.8362	val_loss: 1.0031	val_accuracy: 0.4984	val_top 5 accuracy: 0.7924
epoch 29/100	176.176 [-----]	23s 129ms/step	loss: 1.6255	accuracy: 0.5488	top 5 accuracy: 0.8402	val_loss: 1.0744	val_accuracy: 0.4962	val_top 5 accuracy: 0.7919
epoch 30/100	176.176 [-----]	22s 128ms/step	loss: 1.5800	accuracy: 0.5548	top 5 accuracy: 0.8504	val_loss: 1.0551	val_accuracy: 0.5008	val_top 5 accuracy: 0.7980
epoch 31/100	176.176 [-----]	22s 127ms/step	loss: 1.5600	accuracy: 0.5614	top 5 accuracy: 0.8548	val_loss: 1.0720	val_accuracy: 0.5079	val_top 5 accuracy: 0.7969
epoch 32/100	176.176 [-----]	22s 127ms/step	loss: 1.5272	accuracy: 0.5712	top 5 accuracy: 0.8596	val_loss: 1.0800	val_accuracy: 0.5106	val_top 5 accuracy: 0.7960
epoch 33/100	176.176 [-----]	22s 128ms/step	loss: 1.4993	accuracy: 0.5759	top 5 accuracy: 0.8631	val_loss: 1.0600	val_accuracy: 0.5119	val_top 5 accuracy: 0.7904
epoch 34/100	176.176 [-----]	22s 128ms/step	loss: 1.4609	accuracy: 0.5849	top 5 accuracy: 0.8685	val_loss: 1.0544	val_accuracy: 0.5126	val_top 5 accuracy: 0.7954
epoch 35/100	176.176 [-----]	22s 127ms/step	loss: 1.4276	accuracy: 0.5892	top 5 accuracy: 0.8743	val_loss: 1.0407	val_accuracy: 0.5184	val_top 5 accuracy: 0.7950
epoch 36/100	176.176 [-----]	22s 127ms/step	loss: 1.4102	accuracy: 0.5970	top 5 accuracy: 0.8768	val_loss: 1.0206	val_accuracy: 0.5140	val_top 5 accuracy: 0.7945
epoch 37/100	176.176 [-----]	22s 126ms/step	loss: 1.3800	accuracy: 0.6112	top 5 accuracy: 0.8814	val_loss: 1.0013	val_accuracy: 0.5204	val_top 5 accuracy: 0.8063
epoch 38/100	176.176 [-----]	22s 126ms/step	loss: 1.3500	accuracy: 0.6193	top 5 accuracy: 0.8862	val_loss: 1.0042	val_accuracy: 0.5214	val_top 5 accuracy: 0.8120
epoch 39/100	176.176 [-----]	22s 127ms/step	loss: 1.3575	accuracy: 0.6127	top 5 accuracy: 0.8887	val_loss: 1.0179	val_accuracy: 0.5198	val_top 5 accuracy: 0.8006
epoch 40/100	176.176 [-----]	22s 126ms/step	loss: 1.3010	accuracy: 0.6193	top 5 accuracy: 0.8927	val_loss: 1.0161	val_accuracy: 0.5170	val_top 5 accuracy: 0.8054



```

Epoch 41/100 [====...] - 22s 126ms/step - loss: 1.3160 - accuracy: 0.8247 - top 5 accuracy: 0.8921 - val_loss: 1.8604 - val_accuracy: 0.5208 - val_top 5 accuracy: 0.8092
Epoch 42/100 [====...] - 22s 126ms/step - loss: 1.2679 - accuracy: 0.8329 - top 5 accuracy: 0.8992 - val_loss: 1.8406 - val_accuracy: 0.5284 - val_top 5 accuracy: 0.8109
Epoch 43/100 [====...] - 22s 126ms/step - loss: 1.2534 - accuracy: 0.8375 - top 5 accuracy: 0.9034 - val_loss: 1.8110 - val_accuracy: 0.5306 - val_top 5 accuracy: 0.8057
Epoch 44/100 [====...] - 22s 126ms/step - loss: 1.2111 - accuracy: 0.8431 - top 5 accuracy: 0.9081 - val_loss: 1.8281 - val_accuracy: 0.5218 - val_top 5 accuracy: 0.8050
Epoch 45/100 [====...] - 22s 127ms/step - loss: 1.2073 - accuracy: 0.8488 - top 5 accuracy: 0.9090 - val_loss: 1.8084 - val_accuracy: 0.5302 - val_top 5 accuracy: 0.8054
Epoch 46/100 [====...] - 22s 127ms/step - loss: 1.1775 - accuracy: 0.8558 - top 5 accuracy: 0.9117 - val_loss: 1.8089 - val_accuracy: 0.5204 - val_top 5 accuracy: 0.8074
Epoch 47/100 [====...] - 22s 126ms/step - loss: 1.1893 - accuracy: 0.8563 - top 5 accuracy: 0.9101 - val_loss: 1.8167 - val_accuracy: 0.5360 - val_top 5 accuracy: 0.8042
Epoch 48/100 [====...] - 22s 127ms/step - loss: 1.1506 - accuracy: 0.8621 - top 5 accuracy: 0.9161 - val_loss: 1.8285 - val_accuracy: 0.5314 - val_top 5 accuracy: 0.8086
Epoch 49/100 [====...] - 22s 126ms/step - loss: 1.1538 - accuracy: 0.8638 - top 5 accuracy: 0.9154 - val_loss: 1.8109 - val_accuracy: 0.5306 - val_top 5 accuracy: 0.8134
Epoch 50/100 [====...] - 22s 126ms/step - loss: 1.1506 - accuracy: 0.8682 - top 5 accuracy: 0.9199 - val_loss: 1.8042 - val_accuracy: 0.5254 - val_top 5 accuracy: 0.8096
Epoch 51/100 [====...] - 22s 126ms/step - loss: 1.1175 - accuracy: 0.8708 - top 5 accuracy: 0.9222 - val_loss: 1.8511 - val_accuracy: 0.5200 - val_top 5 accuracy: 0.8104
Epoch 52/100 [====...] - 22s 126ms/step - loss: 1.1104 - accuracy: 0.8745 - top 5 accuracy: 0.9226 - val_loss: 1.8841 - val_accuracy: 0.5172 - val_top 5 accuracy: 0.8142
Epoch 53/100 [====...] - 22s 127ms/step - loss: 1.0914 - accuracy: 0.8809 - top 5 accuracy: 0.9236 - val_loss: 1.8216 - val_accuracy: 0.5152 - val_top 5 accuracy: 0.8094
Epoch 54/100 [====...] - 22s 126ms/step - loss: 1.0641 - accuracy: 0.8856 - top 5 accuracy: 0.9270 - val_loss: 1.8429 - val_accuracy: 0.5128 - val_top 5 accuracy: 0.8086
Epoch 55/100 [====...] - 22s 126ms/step - loss: 1.0625 - accuracy: 0.8862 - top 5 accuracy: 0.9301 - val_loss: 1.8316 - val_accuracy: 0.5164 - val_top 5 accuracy: 0.8090
Epoch 56/100 [====...] - 22s 127ms/step - loss: 1.0474 - accuracy: 0.8920 - top 5 accuracy: 0.9308 - val_loss: 1.8110 - val_accuracy: 0.5440 - val_top 5 accuracy: 0.8112
Epoch 57/100 [====...] - 22s 127ms/step - loss: 1.0381 - accuracy: 0.8974 - top 5 accuracy: 0.9297 - val_loss: 1.8647 - val_accuracy: 0.5361 - val_top 5 accuracy: 0.8126
Epoch 58/100 [====...] - 22s 127ms/step - loss: 1.0230 - accuracy: 0.9011 - top 5 accuracy: 0.9341 - val_loss: 1.8241 - val_accuracy: 0.5418 - val_top 5 accuracy: 0.8094
Epoch 59/100 [====...] - 22s 127ms/step - loss: 1.0113 - accuracy: 0.9041 - top 5 accuracy: 0.9361 - val_loss: 1.8216 - val_accuracy: 0.5380 - val_top 5 accuracy: 0.8134
Epoch 60/100 [====...] - 22s 126ms/step - loss: 0.9953 - accuracy: 0.9031 - top 5 accuracy: 0.9386 - val_loss: 1.8356 - val_accuracy: 0.5422 - val_top 5 accuracy: 0.8122
Epoch 61/100 [====...] - 22s 126ms/step - loss: 0.9275 - accuracy: 0.9547 - top 5 accuracy: 0.9532 - val_loss: 1.8391 - val_accuracy: 0.5534 - val_top 5 accuracy: 0.8136
Epoch 62/100 [====...] - 22s 125ms/step - loss: 0.8221 - accuracy: 0.9528 - top 5 accuracy: 0.9562 - val_loss: 1.8775 - val_accuracy: 0.5428 - val_top 5 accuracy: 0.8120
Epoch 63/100 [====...] - 22s 125ms/step - loss: 0.8270 - accuracy: 0.9526 - top 5 accuracy: 0.9550 - val_loss: 1.8464 - val_accuracy: 0.5468 - val_top 5 accuracy: 0.8148
Epoch 64/100 [====...] - 22s 125ms/step - loss: 0.8080 - accuracy: 0.9551 - top 5 accuracy: 0.9576 - val_loss: 1.8789 - val_accuracy: 0.5486 - val_top 5 accuracy: 0.8284
Epoch 65/100 [====...] - 22s 125ms/step - loss: 0.8058 - accuracy: 0.9593 - top 5 accuracy: 0.9573 - val_loss: 1.8691 - val_accuracy: 0.5446 - val_top 5 accuracy: 0.8156
Epoch 66/100 [====...] - 22s 126ms/step - loss: 0.8092 - accuracy: 0.9564 - top 5 accuracy: 0.9568 - val_loss: 1.8588 - val_accuracy: 0.5524 - val_top 5 accuracy: 0.8172
Epoch 67/100 [====...] - 22s 125ms/step - loss: 0.7897 - accuracy: 0.9613 - top 5 accuracy: 0.9604 - val_loss: 1.8649 - val_accuracy: 0.5490 - val_top 5 accuracy: 0.8106
Epoch 68/100 [====...] - 22s 126ms/step - loss: 0.7890 - accuracy: 0.9635 - top 5 accuracy: 0.9598 - val_loss: 1.9060 - val_accuracy: 0.5446 - val_top 5 accuracy: 0.8112
Epoch 69/100 [====...] - 22s 126ms/step - loss: 0.7682 - accuracy: 0.9682 - top 5 accuracy: 0.9620 - val_loss: 1.8645 - val_accuracy: 0.5474 - val_top 5 accuracy: 0.8150
Epoch 70/100 [====...] - 22s 125ms/step - loss: 0.7958 - accuracy: 0.9617 - top 5 accuracy: 0.9600 - val_loss: 1.8549 - val_accuracy: 0.5496 - val_top 5 accuracy: 0.8140
Epoch 71/100 [====...] - 22s 125ms/step - loss: 0.7978 - accuracy: 0.9603 - top 5 accuracy: 0.9590 - val_loss: 1.9109 - val_accuracy: 0.5440 - val_top 5 accuracy: 0.8140
Epoch 72/100 [====...] - 22s 125ms/step - loss: 0.7838 - accuracy: 0.9630 - top 5 accuracy: 0.9594 - val_loss: 1.9015 - val_accuracy: 0.5548 - val_top 5 accuracy: 0.8174
Epoch 73/100 [====...] - 22s 125ms/step - loss: 0.7550 - accuracy: 0.9722 - top 5 accuracy: 0.9622 - val_loss: 1.9219 - val_accuracy: 0.5410 - val_top 5 accuracy: 0.8090
Epoch 74/100 [====...] - 22s 125ms/step - loss: 0.7692 - accuracy: 0.9689 - top 5 accuracy: 0.9599 - val_loss: 1.8928 - val_accuracy: 0.5506 - val_top 5 accuracy: 0.8184
Epoch 75/100 [====...] - 22s 126ms/step - loss: 0.7783 - accuracy: 0.9661 - top 5 accuracy: 0.9597 - val_loss: 1.8846 - val_accuracy: 0.5490 - val_top 5 accuracy: 0.8166
Epoch 76/100 [====...] - 22s 125ms/step - loss: 0.7547 - accuracy: 0.9711 - top 5 accuracy: 0.9638 - val_loss: 1.9347 - val_accuracy: 0.5484 - val_top 5 accuracy: 0.8150
Epoch 77/100 [====...] - 22s 125ms/step - loss: 0.7603 - accuracy: 0.9692 - top 5 accuracy: 0.9616 - val_loss: 1.8966 - val_accuracy: 0.5522 - val_top 5 accuracy: 0.8144
Epoch 78/100 [====...] - 22s 125ms/step - loss: 0.7595 - accuracy: 0.9730 - top 5 accuracy: 0.9618 - val_loss: 1.8720 - val_accuracy: 0.5470 - val_top 5 accuracy: 0.8170
Epoch 79/100 [====...] - 22s 125ms/step - loss: 0.7542 - accuracy: 0.9736 - top 5 accuracy: 0.9622 - val_loss: 1.9132 - val_accuracy: 0.5504 - val_top 5 accuracy: 0.8136
Epoch 80/100 [====...] - 22s 125ms/step - loss: 0.7410 - accuracy: 0.9787 - top 5 accuracy: 0.9635 - val_loss: 1.9233 - val_accuracy: 0.5428 - val_top 5 accuracy: 0.8120
Epoch 81/100 [====...] - 22s 125ms/step - loss: 0.7410 - accuracy: 0.9787 - top 5 accuracy: 0.9635 - val_loss: 1.9233 - val_accuracy: 0.5428 - val_top 5 accuracy: 0.8120
Epoch 82/100 [====...] - 4s 12ms/step - loss: 1.8487 - accuracy: 0.5514 - top 5 accuracy: 0.8186
Test accuracy: 55.14%

```

Test top 5 accuracy: 81.86%

Fig. 4. Result on CIFAR - 100 Dataset

The loss and Accuracy Functions for the 30 epochs are:-

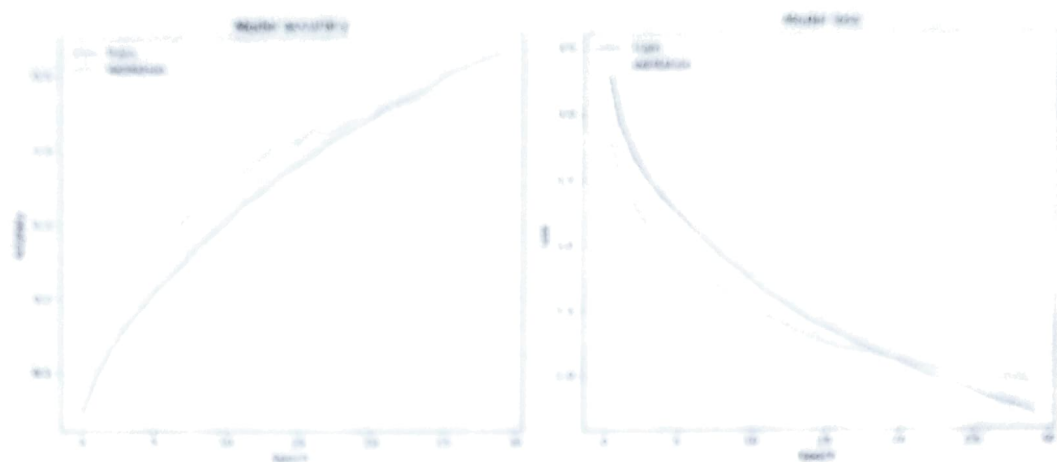


Fig. 4. Accuracy and Loss Curve on 30 Epochs

## 4.2 Result on Self Dataset

After 100 epochs, the ViT model achieves around 66.6% accuracy and 66.6% top-5 accuracy on the test data. These are not competitive results on the millets custom dataset, as a ResNet50V2 trained from scratch on the same data can achieve 53% accuracy(Fig. 5).

```
Epoch 1/100
1/1 [-----] 2s 780ms/step - loss: 1.8529 - accuracy: 0.0000e+00 - top 5 accuracy: 0.0000e+00 - val_loss: 5.0000 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 2/100
1/1 [-----] 1s 151ms/step - loss: 5.1964 - accuracy: 0.0000e+00 - top 5 accuracy: 0.2000 - val_loss: 6.3503 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 3/100
1/1 [-----] 1s 151ms/step - loss: 1.5964 - accuracy: 0.0000 - top 5 accuracy: 0.5000 - val_loss: 6.4012 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 4/100
1/1 [-----] 1s 921ms/step - loss: 2.6885 - accuracy: 0.1000 - top 5 accuracy: 0.0000 - val_loss: 11.3511 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 5/100
1/1 [-----] 1s 940ms/step - loss: 2.2132 - accuracy: 0.2000 - top 5 accuracy: 1.0000 - val_loss: 14.1791 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 6/100
1/1 [-----] 1s 151ms/step - loss: 1.1351 - accuracy: 0.6000 - top 5 accuracy: 0.9000 - val_loss: 16.2711 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 7/100
1/1 [-----] 1s 948ms/step - loss: 1.1946 - accuracy: 0.5000 - top 5 accuracy: 0.9000 - val_loss: 15.3404 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 8/100
1/1 [-----] 1s 948ms/step - loss: 2.2167 - accuracy: 0.7000 - top 5 accuracy: 1.0000 - val_loss: 14.1325 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 9/100
1/1 [-----] 1s 912ms/step - loss: 0.6085 - accuracy: 0.8000 - top 5 accuracy: 1.0000 - val_loss: 15.1232 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 10/100
1/1 [-----] 1s 151ms/step - loss: 0.7279 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 17.4851 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 11/100
1/1 [-----] 2s 215ms/step - loss: 0.2408 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 19.2538 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 12/100
1/1 [-----] 1s 151ms/step - loss: 4.6454 - accuracy: 0.5000 - top 5 accuracy: 0.9000 - val_loss: 20.8432 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 13/100
1/1 [-----] 1s 151ms/step - loss: 0.0254 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 21.8054 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 14/100
1/1 [-----] 1s 151ms/step - loss: 2.9179 - accuracy: 0.0000 - top 5 accuracy: 1.0000 - val_loss: 22.3360 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 15/100
1/1 [-----] 1s 912ms/step - loss: 1.1443 - accuracy: 0.7000 - top 5 accuracy: 1.0000 - val_loss: 21.6031 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 16/100
1/1 [-----] 1s 890ms/step - loss: 1.3243 - accuracy: 0.7000 - top 5 accuracy: 1.0000 - val_loss: 21.1256 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 17/100
1/1 [-----] 1s 895ms/step - loss: 2.0627 - accuracy: 0.8000 - top 5 accuracy: 1.0000 - val_loss: 23.0300 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 18/100
1/1 [-----] 1s 151ms/step - loss: 1.0076 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 24.0063 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 19/100
1/1 [-----] 1s 922ms/step - loss: 0.7997 - accuracy: 0.8000 - top 5 accuracy: 1.0000 - val_loss: 24.5174 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 20/100
1/1 [-----] 1s 901ms/step - loss: 1.0659 - accuracy: 0.7000 - top 5 accuracy: 1.0000 - val_loss: 25.3225 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 21/100
1/1 [-----] 1s 951ms/step - loss: 0.9434 - accuracy: 0.8000 - top 5 accuracy: 1.0000 - val_loss: 26.5122 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 22/100
1/1 [-----] 1s 916ms/step - loss: 0.6293 - accuracy: 0.8000 - top 5 accuracy: 0.9000 - val_loss: 30.3131 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 23/100
1/1 [-----] 1s 151ms/step - loss: 0.3720 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 31.7933 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 24/100
1/1 [-----] 1s 151ms/step - loss: 1.7525 - accuracy: 0.5000 - top 5 accuracy: 1.0000 - val_loss: 32.1392 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 25/100
1/1 [-----] 1s 151ms/step - loss: 0.2930 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 31.7591 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 26/100
1/1 [-----] 1s 937ms/step - loss: 1.5759 - accuracy: 0.6000 - top 5 accuracy: 1.0000 - val_loss: 30.9999 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 27/100
1/1 [-----] 1s 833ms/step - loss: 0.7666 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 32.3828 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 28/100
1/1 [-----] 1s 839ms/step - loss: 0.7190 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 35.2505 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 29/100
1/1 [-----] 1s 936ms/step - loss: 1.2672e-05 - accuracy: 1.0000 - top 5 accuracy: 1.0000 - val_loss: 38.1024 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 30/100
1/1 [-----] 1s 918ms/step - loss: 2.0733 - accuracy: 0.6000 - top 5 accuracy: 1.0000 - val_loss: 39.4048 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 31/100
1/1 [-----] 1s 902ms/step - loss: 5.7550 - accuracy: 0.7000 - top 5 accuracy: 1.0000 - val_loss: 38.2509 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 32/100
1/1 [-----] 1s 938ms/step - loss: 1.0523 - accuracy: 0.8000 - top 5 accuracy: 1.0000 - val_loss: 36.7779 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 33/100
1/1 [-----] 1s 947ms/step - loss: 0.1130 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 39.9273 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 34/100
1/1 [-----] 1s 918ms/step - loss: 1.1921e-08 - accuracy: 1.0000 - top 5 accuracy: 1.0000 - val_loss: 42.7543 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 35/100
1/1 [-----] 1s 992ms/step - loss: 2.0164 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 43.8732 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 36/100
1/1 [-----] 2s 25ms/step - loss: 2.7784 - accuracy: 0.5000 - top 5 accuracy: 1.0000 - val_loss: 43.5749 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 37/100
1/1 [-----] 1s 151ms/step - loss: 6.5517 - accuracy: 0.9000 - top 5 accuracy: 0.9000 - val_loss: 41.8326 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 38/100
1/1 [-----] 1s 151ms/step - loss: 3.2129 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 39.2775 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 39/100
1/1 [-----] 1s 961ms/step - loss: 1.9688 - accuracy: 0.7000 - top 5 accuracy: 1.0000 - val_loss: 36.5276 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 40/100
1/1 [-----] 1s 938ms/step - loss: 2.6505 - accuracy: 0.7000 - top 5 accuracy: 1.0000 - val_loss: 33.6807 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 41/100
1/1 [-----] 1s 931ms/step - loss: 0.0036 - accuracy: 1.0000 - top 5 accuracy: 1.0000 - val_loss: 32.8995 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 42/100
1/1 [-----] 1s 958ms/step - loss: 3.8930 - accuracy: 0.8000 - top 5 accuracy: 1.0000 - val_loss: 32.8743 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 43/100
1/1 [-----] 1s 918ms/step - loss: 2.8606 - accuracy: 0.9000 - top 5 accuracy: 1.0000 - val_loss: 34.3179 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 44/100
1/1 [-----] 1s 151ms/step - loss: 8.2137e-06 - accuracy: 1.0000 - top 5 accuracy: 1.0000 - val_loss: 35.8627 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 45/100
1/1 [-----] 1s 151ms/step - loss: 0.0105 - accuracy: 1.0000 - top 5 accuracy: 1.0000 - val_loss: 37.5356 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 46/100
1/1 [-----] 1s 151ms/step - loss: 0.0442 - accuracy: 1.0000 - top 5 accuracy: 1.0000 - val_loss: 39.1176 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 47/100
1/1 [-----] 1s 808ms/step - loss: 0.0192 - accuracy: 1.0000 - top 5 accuracy: 1.0000 - val_loss: 40.5128 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 48/100
1/1 [-----] 1s 151ms/step - loss: 0.1853 - accuracy: 0.5000 - top 5 accuracy: 1.0000 - val_loss: 42.1588 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
Epoch 49/100
1/1 [-----] 2s 26ms/step - loss: 2.3185 - accuracy: 0.8000 - top 5 accuracy: 1.0000 - val_loss: 38.4427 - val_accuracy: 0.0000e+00 - val_top 5 accuracy: 0.0000e+00
```

```

1/1 [=====] - 1s 1s/step - loss: 1.030M - accuracy: 0.9000 - top-5 accuracy: 0.9000 - val_loss: 11.3617 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 51/100
1/1 [=====] - 1s 908ms/step - loss: 1.212M - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 21.1087 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 52/100
1/1 [=====] - 1s 992ms/step - loss: 1.595M - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 21.4513 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 53/100
1/1 [=====] - 1s 918ms/step - loss: 4.1847 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 20.9113 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 54/100
1/1 [=====] - 1s 948ms/step - loss: 1.6418 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 19.6725 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 55/100
1/1 [=====] - 1s 916ms/step - loss: 1.1146 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 19.1597 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 56/100
1/1 [=====] - 1s 910ms/step - loss: 1.8176 - accuracy: 0.9000 - top-5 accuracy: 0.9000 - val_loss: 17.4867 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 57/100
1/1 [=====] - 1s 913ms/step - loss: 0.8102 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 21.0802 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 58/100
1/1 [=====] - 1s 1s/step - loss: 0.8063 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 21.9109 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 59/100
1/1 [=====] - 1s 944ms/step - loss: 3.1492e-05 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 22.7335 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 60/100
1/1 [=====] - 1s 967ms/step - loss: 0.0068 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 21.6676 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 61/100
1/1 [=====] - 1s 1s/step - loss: 1.5808 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 24.4545 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 62/100
1/1 [=====] - 2s 2s/step - loss: 0.2590 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 24.2019 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 63/100
1/1 [=====] - 1s 958ms/step - loss: 1.0252e-06 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 24.8315 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 64/100
1/1 [=====] - 1s 942ms/step - loss: 2.1437 - accuracy: 0.8000 - top-5 accuracy: 1.0000 - val_loss: 25.4746 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 65/100
1/1 [=====] - 1s 878ms/step - loss: 0.0041 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 26.7905 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 66/100
1/1 [=====] - 1s 872ms/step - loss: 0.2781 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 27.0681 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 67/100
1/1 [=====] - 1s 1000ms/step - loss: 9.5167e-08 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 27.1858 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 68/100
1/1 [=====] - 1s 912ms/step - loss: 4.6342 - accuracy: 0.9000 - top-5 accuracy: 0.9000 - val_loss: 26.7907 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 69/100
1/1 [=====] - 1s 900ms/step - loss: 3.4223 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 25.3562 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 70/100
1/1 [=====] - 1s 880ms/step - loss: 0.9279 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 25.3517 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 71/100
1/1 [=====] - 1s 895ms/step - loss: 2.2811 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 24.5457 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 72/100
1/1 [=====] - 1s 927ms/step - loss: 0.9130e-07 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 24.7530 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 73/100
1/1 [=====] - 1s 928ms/step - loss: 3.4373 - accuracy: 0.9000 - top-5 accuracy: 0.9000 - val_loss: 26.3827 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 74/100
1/1 [=====] - 1s 1s/step - loss: 2.9454 - accuracy: 0.8000 - top-5 accuracy: 1.0000 - val_loss: 27.3815 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 75/100
1/1 [=====] - 1s 1s/step - loss: 1.2880 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 28.3156 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 76/100
1/1 [=====] - 1s 960ms/step - loss: 1.1706 - accuracy: 0.8000 - top-5 accuracy: 1.0000 - val_loss: 30.5797 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 77/100
1/1 [=====] - 1s 885ms/step - loss: 1.1710e-04 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 33.2496 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 78/100
1/1 [=====] - 1s 863ms/step - loss: 1.8948 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 37.1331 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 79/100
1/1 [=====] - 1s 955ms/step - loss: 5.0060e-07 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 40.4336 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 80/100
1/1 [=====] - 1s 949ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 43.2479 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 81/100
1/1 [=====] - 1s 950ms/step - loss: 7.5046 - accuracy: 0.9000 - top-5 accuracy: 0.9000 - val_loss: 43.3815 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 82/100
1/1 [=====] - 1s 955ms/step - loss: 0.9098 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 45.4372 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 83/100
1/1 [=====] - 1s 934ms/step - loss: 5.2604 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 45.1603 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 84/100
1/1 [=====] - 1s 916ms/step - loss: 2.8645 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 42.7548 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 85/100
1/1 [=====] - 1s 921ms/step - loss: 5.9605e-08 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 40.3577 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 86/100
1/1 [=====] - 1s 928ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 37.9474 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 87/100
1/1 [=====] - 2s 2s/step - loss: 0.8129 - accuracy: 1.0000 - top-5 accuracy: 1.0000 - val_loss: 35.7844 - val_accuracy: 0.0000e+00 - val_top-5 accuracy: 0.0000e+00
Epoch 88/100
1/1 [=====] - 1s 1s/step - loss: 2.1550 - accuracy: 0.9000 - top-5 accuracy: 1.0000 - val_loss: 34.9808 - val_accuracy: 0.0000e+00 - val_top-5
```

```

Epoch 100/100
1/1 ..... loss: 0.0000e+00 · accuracy 1.0000 · top-5-accuracy 1.0000
1/1 ..... loss: 1.2327 · accuracy 0.6667 · top-5-accuracy 0.6667
Test accuracy: 66.67%
Test top-5 accuracy: 66.67%

```

Fig. 5 Result on Self Dataset

The loss and Accuracy Functions for the 100 epochs are:-

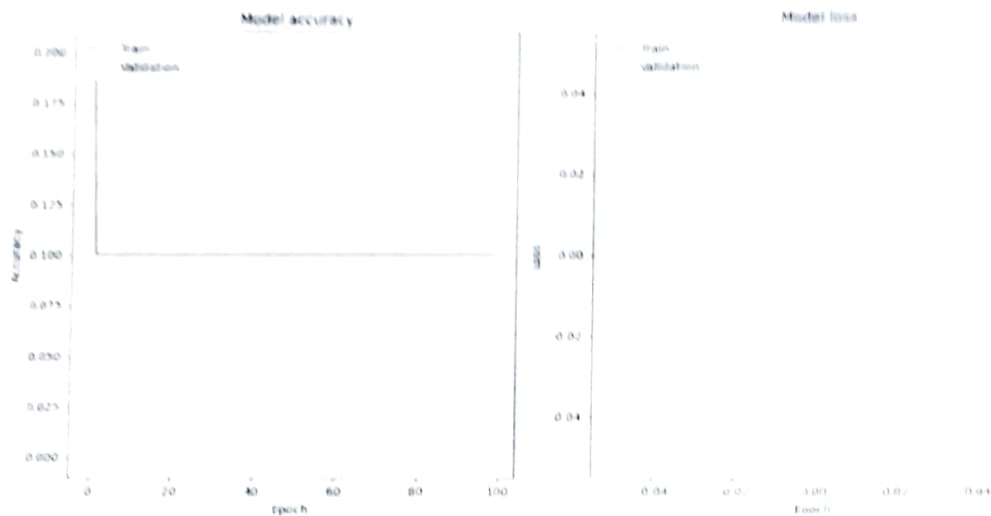


Fig. 6 Accuracy and Loss Curve on 100 Epochs



## Chapter 5: Conclusion

In this project, centered around the implementation of a Vision Transformer (ViT) using TensorFlow and Keras. The model underwent rigorous evaluation on two distinct datasets, namely CIFAR-100 and a custom-made millets dataset, each posing unique challenges and opportunities. The CIFAR-100 dataset, renowned for its diversity with 100 classes and small image resolutions (32x32 pixels), served as a standard benchmark for assessing the ViT model's classification performance. The outcome revealed an accuracy of approximately 57%, aligning with established benchmarks for similar architectures on this dataset. Transitioning to the custom millets dataset, tailored for a specific application, the ViT model demonstrated an improved accuracy of around 67%. This enhancement hinted at the model's adaptability to more specialized datasets, where the feature space might be more navigable. However, the intriguing aspect lies in the observed decline in accuracy from CIFAR-100 to the millets dataset, raising pertinent questions about the model's generalization capabilities. Several factors could contribute to this decline, encompassing dataset characteristics, model hyperparameters, and the efficacy of data augmentation strategies. The millets dataset, being more specialized and tailored to a particular application, might have enabled the model to glean more relevant features for improved classification. Conversely, the model's performance on CIFAR-100, with its diverse array of classes, could be hindered by a more complex and varied feature space. To address these nuances, future endeavors could delve into nuanced hyperparameter tuning, considering learning rates, batch sizes, and the optimal number of Transformer layers for each dataset. Furthermore, refining data augmentation strategies tailored to the unique characteristics of each dataset might unlock additional performance gains. Transfer learning strategies, such as leveraging pre-trained ViT models on larger datasets, could potentially enhance the model's knowledge transfer capabilities. The concept of ensemble learning, amalgamating predictions from multiple ViT models with diverse initializations or architectures, could offer another avenue for performance improvement. Additionally, efforts to expand the millets dataset, both in size and diversity, could potentially bolster the model's ability to generalize across a broader range of instances. In conclusion, this project not only sheds light on the capabilities of ViT models in image classification tasks but also underscores the importance of tailoring models and strategies to the intricacies of specific datasets. The observed variations in accuracy between CIFAR-100 and the millets dataset present opportunities for refinement and future exploration, emphasizing the iterative and adaptive nature of advancing machine learning models in real-world applications. Future work could focus on developing methods to enhance the interpretability of Vision Transformer models. This includes research into techniques that provide more transparent insights into the decision-making processes of ViTs, making them more understandable for end-users and stakeholders in various domains. Research efforts could be directed towards making Vision Transformer models more scalable for large-scale agricultural systems. This involves optimizing training and inference procedures to handle extensive datasets and deploying models across distributed computing environments.

## References

1. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł, Polosukhin I. Attention is all you need. *Advances in neural information processing systems*. 2017;30.
2. Abdu AM, Mokji MM, Sheikh UU. Automatic vegetable disease identification approach using individual lesion features. *Computers and Electronics in Agriculture*. 2020 Sep 1;176:105660.
- Atila, Ü., Uçar, M., Akyol, K., Uçar, E., 2021. Plant leaf disease classification using efficientnet deep learning model. *Ecol. Inform.* 61, 101182.
3. Barbedo JG. Impact of dataset size and variety on the effectiveness of deep learning and transfer learning for plant disease classification. *Computers and electronics in agriculture*. 2018 Oct 1;153:46-53.
4. Barbedo JG, Koenigkan LV, Halfeld-Vieira BA, Costa RV, Nechet KL, Godoy CV, Junior ML, Patricio FR, Talamini V, Chitarra LG, Oliveira SA. Annotated plant pathology databases for image-based detection and recognition of diseases. *IEEE Latin America Transactions*. 2018 Jun;16(6):1749-57.
5. Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, Dehghani M, Minderer M, Heigold G, Gelly S, Uszkoreit J. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*. 2020 Oct 22.