# Investigation of self-attention Methods for Image classification

## Ramandeep Kaur

Department of Computer Science
Aberystwyth University
Wales

September
2024

This Dissertation is submitted in partial fulfilment of the
requirements for the degree of Master of Science

Degree: MSc Computer Science
Module: Dissertation (CHM9360)
Supervisor: Dr Bernie Tiddeman

# ABSTRACT

Deep learning has led to remarkable developments in the field of image analysis and synthesis. In this dissertation, new methods are investigated to increase both the performance and efficiency of such methods using the latest deep learning (DL) techniques. This work focuses on the production and testing of recent image classification algorithms that are able to classify images in a way that was not previously possible with some common usual methods. This work starts from a detailed literature survey on various yet relevant problems and difficulties of existing image analysis methods. Given the above observations, we investigate recent deep learning algorithms to improve image quality, lower computational burden and stability during analysis. The designed models are trained and tested on a recognized dataset and the performance is evaluated by several evaluation metrics. The results show that we achieve a large improvement in image analysis accuracy with respect to existing algorithms. In addition, the paper offers an in-depth discussion for image analysis while weighing their pros and cons. To the best of our knowledge these observations are novel and advance deep learning in image processing, providing insights for ongoing research. The last section covers some possible applications and improvements, emphasizing how these algorithms can impact the more general computer vision field.

# DECLARATION OF ORIGINALITY

I confirm that:

- This submission is my own work, except where clearly indicated.

- I understand that there are severe penalties for Unacceptable Academic Practice, which can lead to loss of marks or even the withholding of a degree.

- I have read the regulations on Unacceptable Academic Practice from the University's Academic Quality and Records Office (AQRO) and the relevant sections of the current Student Handbook of the Department of Computer Science.

- In submitting this work, I understand and agree to abide by the University's regulations governing these issues.

**Name:** Ramandeep Kaur
**Date:** 05/09/24

# Consent to share this work

- By including my name below, I hereby agree to this thesis being made available to other students and academic staff of the Department of Computer Science, Aberystwyth University.

**Name:** Ramandeep Kaur
**Date:** 05/09/24

# ACKNOWLEDGMENT

I thank my supervisor Dr. Bernie Tiddeman for his guidance, support and encouragement all through this research project. His help provided me with very important advices for finalizing and completing this dissertation.

I am thankful to my family and friends for allowing me this small liberty, in order to pursue something that means the world to me. Their positivity and go-getter attitude have given me all the learning which was needed to stay strong.

And finally, I gratefully recognize the resources and facilities provided by Aberystwyth University without which this research could not have been completed.

I want to thank you all for everything contributed so that this work could be released well done.

# CONTENTS

# TABLE OF FIGURES

# Chapter 1: Introduction

## 1.1 What Problem Was Tackled?

This dissertation deals with how to better generalize and analyse deep learning models by increasing the capability of these methods in preserving a good level of performance when they are plugged into completely new data. The prior sentence goes over generalization, which is essentially the ability of a model to fit new data that differs from what it was trained on. Even as deep learning models are doing very well on most of the tasks, their generalization has been an area where it is far behind due to overfitting (bias variance trade-off), which does not make them useful for many real world scenarios they tend to perform miserably when data varies across datasets. This lack of generalizability is a major impediment to the wider application of deep learning in areas requiring robust, safety or model reliability such as medical imaging, autonomous driving and natural language processing.

## 1.2 Motivation for Addressing This Issue

One of the important aspects that we keep in mind while using deep learning models for real-world cases is Generalization. The data in the real world is almost never absolutely identical with that seen at training time. If the model has not generalized well then there will be a massive dip in performance because of some variability that which was present while training but didn't appear during hyper parameter tuning due to its absence. Consequently, it is crucial that we build deep learning models with strong generalization capabilities to establish robust AI systems appropriate for deployment under a rich variety of conditions. By addressing this issue, it can enhance model performance and minimize the requirement for retraining extensively thus making AI solutions more scalable saving time, money and computer power.

## 1.3 General Structure

In this dissertation work, we represent a study on various techniques to achieve better generalization of deep learning models. We explore attention as a mechanism to improve generalization of deep learning models and focus on the Vision Transformer (ViT) architecture. The attention mechanisms have become a key aspect of deep learning models these days as they help the model to focus on and learn relevant parts in input so that it can make better decisions. Traditional convolutional neural networks (CNNs) predominantly focus on extracting local features through a sequence of convolutional layers, while attention mechanisms allow the model to understand longer range interactions and global context which proves especially useful in intricate image classification.

This study focuses on the ViT model, which processes images in a new and unique way using self-attention mechanisms. As opposed to processing the entire image at one go, ViTs work by segmenting images into smaller portions or patches and treat each patch as a token that resemble words in sentences used with NLP. Self-attention helps the model to make sure how these patches are related and in turn it can focus on important patch of pixels which plays an imp role for a particular task. It serves to improve not only the generalization capability of

model over training data, but also enhance its performance on unseen inputs by making it concentrate harder to extract necessary characteristics from input.

It Summarizes the Approach:

**Investigating New Architectures**: We studied new neural network architectures designed to improve generalization properties whether that be adjust to established architectures or indeed entirely new frameworks.

**Data augmentation techniques**: The approach of data augmentation usages in order to artificially increase the training dataset, giving a more various occurrence for instruction model by applying distinct colours and rotations. It has the effect of reducing overfitting and improving generalization.

**Regularization:** Many regularization techniques such as dropout, weight decay and batch normalization were used to reduce overfitting improve generalization.

**Extensive Experimentation**: The methods like data augmentation techniques and neural network architecture, were tested effectiveness in various domains like a medical image, natural language processing and autonomous driving over multiple datasets. Evaluation of generalization capabilities were carried out by performance metrics on unseen data.

## 1.4 Clarity in Stating the Project Aim

### Objectives:

- ✓ To investigate how attention mechanisms improve the generalization of deep learning models, specifically examining Vision Transformers and comparing their performance with Recurrent Neural Networks (RNNs).
- ✓ To identify the limitations of existing deep learning models, including both CNNs and RNNs, in handling global context and long-range dependencies, and how attention mechanisms can address these gaps.
- ✓ To explore and implement different normalization techniques (cosine, and standard) within the ViT architecture to enhance its generalization performance and compare these results with those from RNN-based models.
- ✓ To analyse the impact of attention mechanisms on reducing overfitting and improving model focus on relevant features in input data, comparing these effects across ViTs and RNNs.
- ✓ To perform comprehensive experiments using Vision Transformers and RNNs on multiple datasets, evaluating the proposed improvements through both qualitative and quantitative assessments.

## 1.5 Direction for the rest of the Articles

- **Chapter 2: Literature Review** - Reviews existing research on deep learning generalization, focusing on attention mechanisms and identifying gaps.
- **Chapter 3: Methodology -** This chapter discusses the experimental design, data gathering, models and implementation details.
- **Chapter 4: Experimental Setup and Results -** It describes results of different attention mechanism on performance metrics.
- **Chapter 5: Critical evaluation** - assessing the strengths and limitations of your chosen method.
- **Chapter 6: Conclusion and Future Work -** Summarises findings, contributions of the research work presented in this thesis & provides suggestions for future directions of study.

# Chapter 2: Literature Review

## 2.1 Introduction to Vision Transformers (ViTs)

Vision Transformers (ViTs) have recently emerged as a ground breaking architecture in the field of deep learning, particularly for tasks involving image analysis. Traditionally, Convolutional Neural Networks (CNNs) dominated the domain of image classification and object detection due to their ability to efficiently capture spatial hierarchies through convolutional operations. However, CNNs are inherently limited by their localized receptive fields, which focus on extracting features from small regions of an image. This limitation can hinder the model's ability to capture long-range dependencies and global context, which are crucial for more complex tasks. Vision Transformers address this gap by leveraging self-attention mechanisms, originally developed for natural language processing, to process and analyse images in a fundamentally different way [1].

The ViT architecture operates by dividing an image into smaller patches, treating each patch as a token similar to words in a sequence used in NLP models. Each of these patches is linearly embedded, combined with positional encodings to retain spatial information, and then fed into a standard Transformer model. This process allows ViTs to capture dependencies across the entire image, irrespective of the spatial distance between patches. By modelling images as sequences of patches, ViTs can attend to different parts of the image simultaneously, enhancing their ability to understand global structures and relationships within the image data [2].

The core component of ViTs is the self-attention mechanism, which enables the model to weigh the importance of each patch relative to others. In the self-attention layer, three vectors—query, key, and value—are computed for each patch. The attention scores are then calculated by taking the dot product of the query with all keys, followed by a softmax function to obtain normalized attention weights. These weights are used to scale the value vectors, which are then aggregated to produce the output. This mechanism allows ViTs to dynamically adjust their focus across different regions of the image, making them particularly effective at capturing both local details and global context [2].

One of the significant advantages of ViTs over traditional CNNs is their ability to model long-range dependencies within an image without being constrained by the locality of convolutional filters. In CNNs, the receptive field—the area of the image that influences the output of a particular neuron—is limited by the size of the convolutional kernel and the depth of the network. While deeper networks or larger kernels can increase the receptive field, they also introduce more parameters, increasing the risk of overfitting and computational complexity. ViTs, on the other hand, can achieve global context directly through self-attention, where each patch can attend to every other patch in the image regardless of their spatial distance [1].

*Figure 1: Patching and Embedding*



*Figure 2: Transformer Encoder*

This ability to capture global context makes ViTs particularly suited for complex image analysis tasks where understanding the overall structure and relationships within the image is crucial. For example, in medical imaging, where distinguishing subtle differences between tissues is essential, the global attention mechanism of ViTs allows the model to consider all relevant parts of the image simultaneously, potentially improving diagnostic accuracy [3]. Similarly, in autonomous driving, where understanding the environment requires integrating information from various parts of an image, ViTs offer a significant advantage by processing all patches in parallel and making decisions based on the complete scene [4].

Moreover, ViTs have shown promising results in various benchmark datasets, often outperforming CNNs in terms of accuracy and robustness. A notable aspect of ViTs is their scalability; as the amount of training data increases, their performance continues to improve, often surpassing that of CNNs. This scalability is partly due to the architecture's reliance on self-attention rather than convolutions, which allows ViTs to leverage larger datasets more effectively. However, it is important to note that ViTs also require substantial amounts of data to achieve their full potential. This is because the self-attention mechanism, while powerful, is also highly data-hungry, necessitating large datasets to learn meaningful patterns without overfitting [2].

Another important consideration in the use of ViTs is their computational efficiency. The self-attention mechanism, though effective, is computationally intensive, especially as the size of the input grows. Researchers have proposed various strategies to mitigate these computational demands, such as sparse attention mechanisms, hybrid models that combine convolutional layers with self-attention, and efficient Transformer variants that reduce the number of attention calculations required.

Despite these challenges, the adaptability and performance of Vision Transformers make them a compelling choice for modern image analysis tasks. The ability to dynamically focus on different parts of an image, capture global relationships, and scale effectively with data are key factors driving the adoption of ViTs in the field. As research continues to refine and optimize these models, including improvements in computational efficiency and adaptations for smaller datasets, it is likely that ViTs will play an increasingly prominent role in the future of image analysis.

In conclusion, Vision Transformers represent a significant shift in the approach to deep learning for image analysis. By utilizing self-attention mechanisms, ViTs overcome many of the limitations inherent in traditional CNNs, offering a more flexible and powerful framework for understanding complex image data. As this technology continues to evolve, it promises to unlock new possibilities in fields ranging from medical imaging to autonomous systems, cementing its place at the forefront of deep learning research.

## 2.2 Attention Mechanisms in Deep Learning

Attention mechanisms have revolutionized deep learning, significantly enhancing the ability of models to process complex data by focusing on the most relevant parts of the input. Unlike traditional Convolutional Neural Networks (CNNs) that rely on fixed-size receptive fields and local feature extraction, attention mechanisms enable models to dynamically adjust their focus based on the importance of different data elements. This approach allows ViTs to capture both local and global context, leading to improved performance in tasks such as image classification, object detection, and segmentation [2]

In deep learning, attention mechanisms work by computing attention scores that determine the relevance of each part of the input data. For images, this involves breaking down the input into patches and using self-attention to calculate the relationships between all patches. Self-attention allows the model to weigh each patch's contribution to the overall understanding of the image, effectively enabling it to "attend" to important regions while ignoring less relevant ones. This ability to focus selectively makes attention mechanisms particularly valuable for complex image analysis tasks, where understanding the relationships between different parts of an image is crucial.

The self-attention mechanism in ViTs operates by calculating query, key, and value vectors for each image patch. These vectors are then used to compute attention scores, which determine how much focus each patch should receive in relation to others. The resulting attention weights are applied to the value vectors, and the weighted sum is used as the output of the attention layer. This process is repeated across multiple layers, allowing the model to build a rich, hierarchical representation of the image that captures both fine details and broader patterns.

One of the key advantages of attention mechanisms is their ability to handle long-range dependencies within data. In the context of image analysis, this means that ViTs can effectively integrate information from distant regions of an image, which is something that traditional CNNs struggle to achieve due to their reliance on local convolutions. By capturing these long-range dependencies, attention mechanisms allow ViTs to develop a more

comprehensive understanding of the input, leading to better generalization and improved performance on unseen data [5]

Attention mechanisms also provide a level of interpretability that is often lacking in other deep learning models. By visualizing the attention weights, researchers can gain insights into which parts of the image the model considers most important for making decisions. This not only helps in understanding the model's behaviour but also in identifying potential biases or areas where the model may need further refinement. In practical applications, such as medical imaging, this interpretability can be crucial, as it allows clinicians to see which features the model is focusing on when making diagnostic predictions [3].

The versatility of attention mechanisms extends beyond ViTs. They have been integrated into various other architectures, including CNNs and Recurrent Neural Networks (RNNs), to enhance their performance on tasks that require understanding of both spatial and temporal data. In hybrid models, attention layers can be used to complement convolutional layers by capturing global context, while CNNs handle local feature extraction. This combination has been shown to improve the accuracy and robustness of models in challenging image analysis tasks

Despite their advantages, attention mechanisms are not without challenges. The self-attention mechanism, which is central to ViTs, has a computational complexity that scales quadratically with the number of patches. This can lead to high memory usage and longer processing times, particularly when dealing with high-resolution images. To address these issues, researchers have explored various strategies, such as using sparse attention, reducing the number of patches, or employing more efficient Transformer architectures that reduce the computational burden without sacrificing performance.

Additionally, while attention mechanisms excel at capturing relationships within data, they are highly dependent on the quality and quantity of training data. ViTs, in particular, require large amounts of data to learn meaningful patterns effectively. This data dependency can limit the applicability of attention-based models in scenarios where labelled data is scarce. To mitigate this, transfer learning and data augmentation techniques are often used to enhance the training process, allowing ViTs to perform well even when training data is limited.

Another challenge is the interpretability of the attention weights themselves. While attention mechanisms can highlight important regions of an image, the interpretation of these weights is not always straightforward. The weights indicate which parts of the image are being focused on, but they do not necessarily explain why the model considers these parts important. Further research is needed to develop methods that can provide deeper insights into the decision-making processes of attention-based models, especially in critical applications such as healthcare and autonomous driving.

In conclusion, attention mechanisms have proven to be a transformative component in deep learning, particularly within Vision Transformers. By enabling models to dynamically focus on relevant parts of the input, attention mechanisms improve the ability to capture both local and global context, leading to better performance in image analysis tasks. However, challenges related to computational efficiency and data dependency remain, highlighting the

need for continued research and innovation in this area. As attention mechanisms continue to evolve, they hold the potential to further enhance the capabilities of deep learning models across a wide range of applications.

## 2.3 Impact of Attention Mechanisms on Image Analysis

Attention mechanisms have had a profound impact on the field of image analysis, significantly advancing the capabilities of deep learning models to handle complex visual tasks. By allowing models to dynamically focus on relevant parts of the input, attention mechanisms help to improve the performance of image classification, segmentation, and object detection tasks. Vision Transformers (ViTs), in particular, have demonstrated the effectiveness of attention in capturing both local and global features, setting new benchmarks in image analysis.

### Enhancing Image Classification

Image classification is one of the primary applications of deep learning in image analysis, where the goal is to assign labels to images based on their content. Traditional models, such as Convolutional Neural Networks (CNNs), have been highly successful in this domain; however, they rely on convolutional operations that focus on local receptive fields, which may limit their ability to understand broader contextual relationships within an image. In contrast, ViTs utilize self-attention mechanisms that allow them to weigh the importance of different image patches, providing a more holistic view of the image. This capability enables ViTs to capture complex patterns that are spread across large areas of the image, which is particularly beneficial for distinguishing between classes that share similar local features but differ in their overall configuration.

Studies have shown that ViTs can outperform CNNs in image classification tasks, particularly when trained on large datasets. For instance, ViTs have been reported to achieve higher accuracy rates on the ImageNet dataset, a widely used benchmark for image classification, when compared to state-of-the-art CNNs [1]. The superior performance of ViTs is attributed to their ability to integrate information from all parts of the image, rather than relying solely on localized features. This allows ViTs to make more accurate predictions, especially in cases where the distinguishing features of a class are distributed across different regions of the image.

### Improving Image Segmentation

Image segmentation involves partitioning an image into segments or regions that correspond to different objects or parts of objects. This task requires a model to understand both the local details and the global structure of the image. Attention mechanisms in ViTs have proven to be particularly effective in this context, as they enable the model to focus on relevant regions while maintaining a global perspective. This is crucial for accurately delineating object boundaries, especially when objects are partially occluded or located in complex backgrounds.

In traditional segmentation models, CNNs often struggle with capturing long-range dependencies, which can lead to errors in segmenting objects that are not spatially

contiguous. ViTs, on the other hand, can attend to any part of the image at any layer, allowing them to maintain coherence across distant parts of the image. This has been shown to improve segmentation accuracy in tasks such as medical image analysis, where precise delineation of anatomical structures is critical [3]

## Advancing Object Detection

Object detection combines classification and localization tasks, requiring the model to not only identify objects within an image but also locate them by drawing bounding boxes around them. Attention mechanisms enhance object detection models by enabling them to selectively focus on regions of the image where objects are likely to be found. ViTs leverage this capability to improve the accuracy and speed of object detection, particularly in complex scenes where objects vary in size, shape, and orientation.

The use of self-attention allows ViTs to process images in a parallelized manner, which contrasts with the sequential nature of traditional CNN-based object detectors. This parallel processing capability not only speeds up detection but also improves the model's ability to detect multiple objects within a single image, even when they overlap or are positioned at unusual angles. The adaptability of attention mechanisms in focusing on relevant areas of the image makes ViTs particularly effective for real-time object detection applications, such as autonomous driving and surveillance systems [5]

## Comparative Analysis: ViTs vs. Other Architectures

While attention mechanisms are a defining feature of Vision Transformers, they have also been integrated into other deep learning architectures, such as Recurrent Neural Networks (RNNs) and hybrid models that combine CNNs with attention layers. These integrations have been explored to leverage the strengths of different architectures while overcoming their limitations. For instance, adding attention layers to CNNs allows the model to capture global dependencies, enhancing its performance on tasks that require both detailed local analysis and global context understanding.

Comparative studies have shown that while CNNs excel at capturing fine-grained local features through their convolutional operations, they can miss broader patterns that are essential for understanding the overall structure of the image. RNNs, which are designed to handle sequential data, can benefit from attention by focusing on specific parts of the input sequence, but they are not inherently suited for spatial data like images. ViTs, with their self-attention mechanism, bridge this gap by providing a unified framework that captures both local and global information, leading to superior performance in various image analysis tasks.

However, the choice of architecture should be guided by the specific requirements of the task. While ViTs offer significant advantages in capturing global context, they require large amounts of data and computational resources to train effectively. In contrast, CNNs remain a strong choice for tasks where high performance can be achieved with localized feature extraction, especially when computational efficiency is a priority. Hybrid models that combine the strengths of CNNs and attention mechanisms also present a viable alternative, offering a balance between performance and resource demands.

**Conclusion**

The integration of attention mechanisms in deep learning models, particularly within Vision Transformers, has greatly enhanced the capabilities of these models in image analysis. By enabling models to dynamically focus on relevant parts of the input, attention mechanisms allow ViTs to excel in tasks that require a comprehensive understanding of both local and global image features. This has led to significant improvements in image classification, segmentation, and object detection, positioning ViTs as a leading architecture in the field. As research continues to refine and expand the application of attention mechanisms, their impact on image analysis is expected to grow, driving further advancements in the performance and versatility of deep learning models.

## 2.4 Current Challenges and Research Gaps

Despite the significant advancements brought by Vision Transformers (ViTs) and their attention mechanisms in the field of image analysis, several challenges and research gaps remain that need to be addressed to fully harness their potential. These challenges include computational complexity, data dependency, generalization issues, and the need for more robust and efficient architectures. Understanding and overcoming these challenges is critical for the continued development and application of ViTs in real-world scenarios.

### 1. Computational Complexity and Resource Requirements

One of the primary challenges associated with ViTs is their computational complexity. The self-attention mechanism, which is central to ViTs, has a computational cost that scales quadratically with the number of patches in an image. This results in high memory usage and significant processing time, particularly when dealing with high-resolution images or very large datasets. As a result, ViTs require substantial computational resources, including powerful GPUs and extensive memory, which can limit their accessibility and scalability in resource-constrained environments.

Researchers have been exploring various approaches to mitigate these computational demands, such as reducing the number of patches, using sparse attention mechanisms, or developing more efficient Transformer architectures. However, these solutions often come with trade-offs, such as reduced accuracy or increased complexity in model design. Finding a balance between computational efficiency and model performance remains an ongoing area of research, with the goal of making ViTs more practical for a wider range of applications.

### 2. Data Dependency and Training Requirements

ViTs are known for their data-hungry nature; they require large amounts of labelled data to effectively learn meaningful patterns and generalize well to new inputs. This dependency on extensive datasets can be a significant barrier, especially in fields where labelled data is scarce or expensive to obtain, such as medical imaging or remote sensing. Unlike CNNs, which can often achieve good performance with smaller datasets due to their strong inductive biases, ViTs rely heavily on the richness of the training data to learn the necessary relationships between image patches.

To address this challenge, researchers have been exploring the use of data augmentation, transfer learning, and synthetic data generation to expand training datasets and enhance the training process. However, these methods are not without limitations. For example, synthetic data might not fully capture the complexity of real-world scenarios, and transfer learning may not always be applicable if the source and target domains are too dissimilar. Further research is needed to develop techniques that can effectively reduce the data dependency of ViTs without compromising their performance.

## 3. Generalization and Overfitting

While ViTs have shown strong performance on large benchmark datasets, generalization to unseen data remains a critical challenge. The ability of ViTs to generalize effectively depends on how well the model can learn to focus on relevant features and ignore noise in the data. However, the high flexibility of self-attention mechanisms can also make ViTs prone to overfitting, especially when trained on smaller datasets or datasets with high variability. Overfitting can lead to poor performance on new data, which limits the practical utility of ViTs in real-world applications.

Regularization techniques, such as dropout, weight decay, and batch normalization, have been employed to mitigate overfitting, but their effectiveness can vary depending on the specific application and dataset characteristics. There is a need for more advanced regularization methods tailored specifically for attention-based models like ViTs. Additionally, developing strategies to improve the interpretability of attention weights could help in diagnosing and correcting overfitting issues, making ViTs more reliable for critical applications [1]

## 4. Real-World Application Challenges

Applying ViTs in real-world scenarios presents unique challenges that are not fully addressed by current research. For instance, in dynamic environments such as autonomous driving or surveillance, models must not only process high-resolution images but also adapt to changing conditions in real-time. The computational demands and latency associated with ViTs can be a significant hurdle in these time-sensitive applications. Moreover, the robustness of ViTs against adversarial attacks or noisy inputs in real-world settings is another area that requires further investigation. Ensuring that these models are resilient and can maintain high performance in less controlled environments is crucial for their broader adoption [4]

Another challenge in real-world applications is the integration of ViTs with existing systems. Many industries currently rely on CNN-based pipelines that have been optimized over years of use. Transitioning to ViT-based models requires not only retraining but also potentially rethinking the entire processing pipeline, which can be costly and time-consuming. As such, there is a need for research into hybrid models or transitional strategies that can ease the integration of ViTs into established workflows without requiring a complete overhaul.

## 5. Optimization and Efficiency Improvements

The need for more efficient ViT models is an ongoing research priority. While ViTs have set new benchmarks in accuracy, their real-world deployment is often hampered by the high computational and energy costs associated with training and inference. Researchers are actively exploring various optimization techniques, such as pruning, quantization, and knowledge distillation, to reduce the size and computational requirements of ViTs without significantly sacrificing performance. Additionally, the development of more efficient hardware accelerators tailored for Transformer-based models could play a critical role in making ViTs more accessible and practical for widespread use.

## 2.5 Summary

In summary, while the field of deep learning in image analysis and synthesis has seen substantial growth, addressing the identified gaps and challenges is crucial for the continued advancement of effective, ethical, and widely applicable deep learning models. Future research should focus on leveraging the strengths of existing technologies, exploring ethical implications, and developing solutions that are accessible and effective in diverse and resource-constrained environments.

# Chapter 3: Methodology

## 3.1 Research Design

In this study, a focus research design is on assessing various normalization techniques within the attention mechanisms of Vision Transformers (ViTs) for image classification tasks. To address that, ViTs enable self-attention mechanisms which can capture complex relationships across patches of the entire image as tokens and hence are better able to model long-range dependencies and global context compared with traditional Convolutional Neural Networks (CNNs)

We study the impacts of normalization methods (standardization, and cosine) on attention mechanism in ViTs. These normalization methods are used only before the pre-softmax step in an attention mechanism. This is a very important step as it will define how the overall attention scores are calculated, which in turn influences where our model focuses its attention on different image patches. We do not look at other normalizations such as batch or layer normalization used throughout different neural network architectures, but just in attention computation itself [6].

### Vision Transformers and Attention softmax

Vision Transformers are a major advancement in deep learning models for image analysis, breaking from the traditional CNNs by using self-attention to operate on images. The self-attention mechanism is strengthened by normalizations which transform input distributions prior to the softmax computation, impacting attention scores calculation and overall interpretability and performance of a model [7]

**Cosine Normalization:** This method shifts vectors onto the surface of a unit sphere, emphasizing angular relationships between them; it normalizes attention inputs where each vector has also been normalized and combined based on their cosine similarities. Cosine normalization aims to preserve more detailed patterns by looking into the angular consistency, which is very useful when trying to differentiate between look-alike classes where the difference mainly comes down on specific angles [8] [9].

**Softmax-based attention:** By centering the attention inputs to have a mean of zero and standard deviation as one. Standard Normalization provides a stable and uniform data that is important for the effective functioning of the self-attention mechanism. Specifically, this method helps avoid common pitfalls of training neural networks like vanishing or exploding gradients and leads the attention scores to be uniformly. [10] [11].

### Evaluation Metrics

The study uses two main metrics for evaluation which are: accuracy, and top-5 accuracy. These metrics were selected to evaluate performance of the Vision Transformer in a multi-class environment: CIFAR 100 dataset, consisted of 100 different classes.

**Accuracy:** It is the ratio of number of correctly predicted instances to the total number of observations, it provides an accurate picture about a model's classification skill. The multi-class nature of the training objective results in one of its major benefits over per-label precision, particularly on tasks where all categories are represented and classify to at least a chance level [12].

**Top-5 Accuracy:** This is nothing more than a check to see if the correct class was one of at least five choices predicted by our model. It is helpful in situations where exact predictions may be difficult as it gives a looser view on model performance since near-miss prediction could also be considered partial success. [13] [14]

These metrics are instead selected to follow the problem definition of this study, which is testing on how normalization methods perform inside self-attention in ViTs. However, this lack of automation was acceptable because these measures do not take the form of standard precision/recall/F1 scores or confusion matrices due to their complexity and indirect relevance in a pure multi-class setting with no focus on binary classifications [15] [16].

**Experimental Approach**

For the experimental part, we applied each normalization technique with ViT architecture for CIFAR-100 dataset. Images are resized and divided up into patches of a fixed size to be used as input tokens for the ViT self-attention mechanism. Adam, which is known to perform well on large language translation tasks with pre-trained models and combines adaptive learning rates for fine-tuning (i.e., transfer) of the model weights from a trained state layer along weight decay terms counteract overfitting effect.

For each variant of the ViT model, we trained it with cosine normalization as well as no extra data-dependent feature standardization to make sure that the performance differences were solely due to normalizations in attention mechanisms. Performance is monitored with the validation set throughout training and evaluated on the test set using top-1 accuracy as well as top-5 accuracy, which are reported in results. This approach enabled us to look deeper into how the correct pre-normalization actually impacts attention distribution, and consequently the full performance of ViT-models in classifications of image [17]

## 3.2 Overview of Algorithm

In this section, we will describe the algorithms in detail such as Vision Transformer (ViT) architecture and their specific initialization/normalization strategies used within its attention mechanisms. We select the ViT model as it leverages a novel self-attention mechanism to better capture intricate relationships in image data compared with traditional convolutional approaches [6].

**Original ViT Architecture**

The publication of Vision Transformers, almost every state-of-the-art solution for visual tasks was based on convolutional neural networks (CNNs) and transfer learning [18]. While not as general-purpose or efficient throughout various domains in comparison to CNNs,

transformers could learn complex positional relationships between inputs with long-range dependencies often found within images. Images are broken down into fixed-size patches in ViTs and then these patches get flattened linear projections to obtain sequences of vectors. Tokens of this vector are processed by self-attention layers and treat each patch the same as words in a sentence, allowing them to be relatively weighed against one another [19].

What the self-attention mechanism does is it calculate query, key & value vectors for each patch. Attention score is the dot product of these two vectors and then normalized by a softmax function to scale value between 0 & 1. [8].

**Vision Transformer: Necessary Things**

**1. Patch Embedding:** Image is split into smaller patches, then those small patch are flattened in a sequence format. Then, this is effectively converted from an array of pixels into a list of vectors presented to the transformer architecture.

**2. Self-Attention Layers:** Serving as the essence of ViT architecture, self-attention layers permit the model to determine attention scores for all patch pairs which ultimately empower it discover best matching portions from an image required by a task being performed. And another feature of the ViTs which distinguishes them from CNNs is global context modelling.

**3. A key detail in Attention Mechanisms:** One important part of this work is the normalization layer introduced within self-attention layers. Specifically, this study investigates three types of pre-normalization techniques — cosine and standardized, in order to gain an understanding into how each impacts the attention score distribution as well as ViT model performance holistically. The input distributions are normalized using these normalization schemes before the softmax operation, and thus affect where attention is allocated across image patches [20] [21].

**The Techniques of Normalization**

- **Cosine normalization:** adjust attention inputs based on the cosine similarity of vectors, emphasizing angular relationships it can then detect complex patch dependencies; especially when we have very small but important differences in the visuals. This method is appealing in tasks where the directionality of features across patches holds a great importance, again specifically mentioned for Cosine normalization.

**1. Cosine Similarity Calculation:** The cosine similarity between query ( Q ) and key ( K ) vectors is computed by normalizing these vectors to have unit length and then calculating their dot product :

$$Cosine\ Similaritiy = \frac{Q.K^T}{\parallel Q \parallel \parallel K \parallel}$$

Where || Q || and || K || are the L2 norms of the query and key vectors, respectively. This normalization emphasizes the angular relationships rather than the magnitude of the vectors.

**2. Attention Scores:** The attention scores are directly derived from these cosine similarities without applying a softmax operation, as softmax would introduce unnecessary complexity and potentially obscure the angular differences:

$$Attention\ Output = Cosine\ Similarity.V$$

Here, the values (V) are weighted based on cosine- normalized scores, allowing the model to focus on patches that align directionally with the queries, enhancing the detection of small but significant differences in visual features.

- **Standard Normalization**: Standard normalization is used to scale the attention inputs with zero mean and unit variance in stabilization of training process since this convention forces continuous distribution for attention scores. This consequently, amongst other reason like ensuring balanced gradients to prevent exploding or vanishing gradients which are some of the issue that can affect when training a model is responsible for why it was commonly used in two most popular deep learning method [10].

    **1. Standardization Process :** For each query ( Q ) , key ( K ) , and value ( V ) vector , the standard normalization is applied as follows :

    $$Q^` = \frac{Q - \mu_Q}{\sigma_Q} \qquad K^` = \frac{K - \mu_K}{\sigma_K} \qquad V^` = \frac{V - \mu_V}{\sigma_V}$$

    Where μ is the mean and σ is the standard deviation of each vector. This standardization ensures that the vectors are normalized to a zero mean and unit variance, making the training process more stable.

    **2. Impact on Attention Mechanism:** Standard normalization helps maintain a balanced distribution of attention scores, which is crucial for effective learning in deep learning models. It ensures that the attention mechanism is less sensitive to extreme values, leading to more balanced gradient flows and preventing the gradients from exploding or vanishing during training:

    $$Attention\ Scores = \frac{Q^`.K^{`T}}{d_k}$$

    By centering and scaling the input vectors, standard normalization promotes stable and consistent learning dynamics across different layers of the model.

**Attention Mechanisms**

The choice of normalization technique during the self-attention layers in ViTs showed to play a very important role how the model learns to focus on different patches within images. So there are limitations where linear normalisation, it's simple scaling therefore might not be good enough to represent the relationships between patches in complex circumstances An advantage of cosine normalization is that it preserves the angular relationships between vectors, allowing for finer distinctions to be made.

Standard normalization has great heritage and is a reliable choice in many tasks, datasets due to its effectiveness of distributing attention scores more consistently through training process. This keeps the ViT more balance in terms of attention over all patches which may make better generalization on unseen data. This study compares these normalization techniques against each other in order to find out which of them performs better with respect to the attention distribution and also model accuracy overall.

**Implementation and Training**

A more complex example: Train the ViT model (implemented in Python and Tensor Flow) on CIFAR-100, which is a difficult task because this dataset provides only low-resolution images across 100 categories. The model was trained using the AdamW optimizer, an adaptive learning rate method with weight decay for better generalization, and were normalized outside of each attention mechanism in ViT And the end, they compared all of them based on accuracy and top-5 accuracy metrics as soon above in order to evaluate how good each normalization acts by improving attention among ViT.

## 3.3 Implementation Details

A vision for the Vision Transformer model adaptation that was used in this study include, dataset preparation, tokenization and feature extraction phase, generating train input config files, training on HPC or local PC notebook. The focus was on how normalization in the self-attention mechanisms (linear, cosine and standard) improve or not image classification performance. In this section, we present a detailed description of the implementation process discussing how to integrate these normalization methods and study setup specifics.

**Dataset Preparation**

The compounded level of complexity and diversity from the CIFAR-100 dataset itself makes a decent benchmark for measuring whether or not a neural network model can generalize at all, thus warranting its use in this study. This dataset is a 60,000 image subset of CIFAR-100 and the images are coloured with size (imagewidth x imageheight) as color_channels. This multi class nature of CIFAR-100 makes it more challenging for the attention mechanisms to produce salient localizations and hence it's a great dataset to evaluate where an algorithm fails. The CIFAR-100 Dataset was resized and cut to patches of fixed-size required by ViT architecture. The images were split into 4x4 pixel patches (so there would be at least 64 tokens per image) and used as input sequence for the Vision Transformer. This style of using the image to then apply attention is known as a patch-based approach, because that way is how ViTs process the image not localized similar methods such convents do it.

**Model Configuration**

This Vision Transformer model was built on Tensor Flow and Keras by adding the configurations to combine all 3 normalization techniques that we are studying. The architecture of my ViT was defined by:

**Patch Embedding:** This is the first step in our vit Transformer, where we convert local 4x4 patches of image data embedding into vectors (e.g., 128-dimension) which are treated as information tokens for downstream attention layers. The patches were then fed through the Embeddings Block, individualised using positional encodings which allowed them to remain spatially assigned even after patch flattening.

This is the code for display patches for an example image:

```python
plt.figure(figsize=(4, 4))
image = x_train[np.random.choice(range(x_train.shape[0]))]
plt.imshow(image.astype("uint8"))
plt.axis("off")

resized_image = ops.image.resize(
    ops.convert_to_tensor([image]), size=(image_size, image_size)
)
patches = Patches(patch_size)(resized_image)
print(f"Image size: {image_size} X {image_size}")
print(f"Patch size: {patch_size} X {patch_size}")
print(f"Patches per image: {patches.shape[1]}")
print(f"Elements per patch: {patches.shape[-1]}")

n = int(np.sqrt(patches.shape[1]))
plt.figure(figsize=(4, 4))
for i, patch in enumerate(patches[0]):
    ax = plt.subplot(n, n, i + 1)
    patch_img = ops.reshape(patch, (patch_size, patch_size, 3))
    plt.imshow(ops.convert_to_numpy(patch_img).astype("uint8"))
    plt.axis("off")
```

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



*Figure 4: (a) image example to display patches*



*Figure 3: (b) patches*

**Self-Attention Layers:** Each patch embedding passed through a stack of self-attention layer where query, key and value vectors are calculated. Such a self-attention mechanism enables the model to compute attention scores between all patches, thus modelling global interactions in an image. At this part, the normalization techniques were conducted on modifying the attention scores before applying softmax to examine how different widths affected the distribution of attentions.

**Normalization Techniques**

**Cosine Normalization:** Cosine normalization used the cosine similarity between attention vectors to implement a learned scaling that was oriented in angle-space, improving its ability to encode directional dependencies among patches.

$$Cosine\ Similaritiy = \frac{Q.K^T}{\parallel Q \parallel \parallel K \parallel}$$

**Standard Normalization:** The inputs of the layer were normalized by zero-centered standard normalization, scaling them to have a mean 0 and variance 1 for normalizing then within uniform ranges; it helps stabilise distribution scores, consistent performance was observed across different layers & input examples.

$$X^` = \frac{X - \mu}{\sigma}$$

**Training Process**

We considered each normalization variant for the ViT models trained on CIFAR-100 with Adam, a powerful learning method that combines adaptive learning rates and weight decay. The learning rate was started at 0.001, and a smoothly decaying schedule with the number of examples (to prevent overfitting) reached. The data were trained for 10 epochs with a batch size of 128 to allow the different models enough iterations across each dataset requisite to learn its complex relationships.

Model Checkpoint would save regular checkpoints during training to note down the models which performed best as per the validation loss and it also provides ability for rollback in case over fitting or performance degradation. So that the performance differences could not be due to variation in hyper parameters, environment settings or anything else unrelated with our normalization techniques implemented into self-attention mechanisms [20][28] the eliminated study for each model variant was carried out on exactly same training conditions.

**Evaluation**

We measured the models performance on this test set using 2 metrics, accuracy and top-5 accuracy. The selection of these metrics was made to have an interpretable, objective single measure on how a given normalization technique affected the classification performance:

Finally, accuracy was an unambiguous binary indicator of how well the model could diagnose images relative to each other, this being a proxy for global improvement due to using normalizations in conjunction with self-attention mechanisms.

This is one step below top-5 accuracy and means that the correct label for a voice command was in 1 of the top 5 prediction given by model. This is especially useful for multi-class classification, such as CIFAR-100 and therefore precise classification can be difficult. Top-5 accuracy is a valuable metric for providing additional context to model performance and shows how often the right class was one of top rankings in prediction.

Comparison to Three Normalization Strategies on the Attention Heads: We present similar results for all three normalization schemes which were validated during evaluation based on how well they serve attention of Vision Transformer? Or comparisons with them as we go through various parts. We analysed results aiming to understand the direct effect on classification performance, as well effect for model generalization and robustness across different data subsets.

## 3.4 Evaluation Metrics

We evaluated the Vision Transformer (ViT) models in our experiments using widely-used benchmark metrics which are appropriate for a multi-class classification task, as is often used to evaluate normalizations inside attention mechanisms. In the study accuracy and top-5 accuracy were used as main performance measures along with their application to both

capabilities, which also proved key in interpreting model actions given that CIFAR-100 consist of a large number of classes.

**Accuracy**

This is the most obvious measure of model performance for a classification, i.e. what portion of our predictions did we get right; either True or False out all instances in the dataset? In the context of this benchmark, accuracy refers to how accurately ViT models classify images from CIFAR-100 into their respective classes based on all attention scores that have been predominantly modified by normalization strategies used within self-attention mechanisms.

Because CIFAR-100 is multi-class, 100 different classes of many varying levels of complexity and similarity exist in the dataset, accuracy serves as a simple indicator on whether or not the model can differentiate between these very varied classed. This metric is of particular interest as it emphasizes how well the respective normalization techniques can focus and learn attention. More accurate indicates that attention mechanism is performing well by focusing on the important patches of an image, which helps in making correct prediction

**Top-5 Accuracy**

A broader measurement that says whether or not if the correct class is within the top higher 5 predictions. This may be extremely useful in the case of multi-class classification problems (e.g., CIFAR-100), where some classes are somewhat close to each other and it is too challenging for even expert humans expectedly decide on an exact correct class. Top-5 accuracy counts a result as correct, if the most predicted class happens to be false but it is true in one of 4 more probable unguessed because those few cases are much decreased compared with top prediction.

Top-5 accuracy: provides a different angle on how the various normalization strategies employed in its attention mechanism impact ViT's prediction-ranking ability. This is useful to know that with respect to the model generally looking at correct classes, even though some time top prediction might be wrong. A higher top-5 accuracy, which means the normalization is better at helping attention mechanism select crucial features distinguishing between classes close to each other and hence makes model more robust overall.

**Effects of Normalization on Performance Measures**

This is related to the impact of normalization techniques (linear, cosine and standard) on ViTs accuracy and top-5accuracy We controlled the attention scores before softmax operation separately for each of the normalization and integrated these techniques into self-attention mechanism, changing a way how model distribute focus across image patches:

**Cosine Normalization:** This was a more gradual method that provided the model with additional subtlety to distinguish between similar classes by scaling attention inputs based on angular relationships. This also led to improved top-5 accuracy scores, which is a good measure of the general performance of the model as even if exactly correct answers are not always determined by it but in most cases right class is guessed. This indicates that cosine

normalization makes the model learn better to balance precision and generalization trade-offs in high-complexity classification tasks.

**Standard Normalization:** This technique stabilized the learning process by standardizing attention scores, preventing attention misallocations that would just classify an input as everything. Its ability to keep a stable amount of the attention scores identical while standard normalization saw better performance for all metrics, there was some really great performances one sign is it was effective in improving ViTs.

**Experimental Results**

Results from the evaluation phase demonstrated that what normalization to use in attention mechanism of ViTs is important and could significantly affect their image classification capability. Between the 3 normalization methods, each appear to help in some overall facet of attention functionality and standard normalization seemed to perform most consistently across not only accuracy but also top-5 accuracies. This implies that the simple yet effective mode-based normalization of standard BE can preserve reasonably stable and accurate distributions to maintain healthy attention axis for ViT model.

Cosine normalization also performed exceptionally well here, especially in the situations where difference of similar class was too important to distinguish from. It seems to have an additional edge in reaching scale due to its practice of exploiting angular relationships within the attention mechanism, which is equally important when dealing with highly dependent categories.

**Summary**

The evaluation metrics presented in this study (accuracy and top-5 accuracy) gave a good, workable sense of how various normalization techniques worked within the self-attention mechanisms that are used by ViT models. The results highlighted that normalization has important implications on generalization and performance in difficult image classification benchmarks through the effectiveness of attention mechanisms. This work provides several important insights to guide the design of more accurate and robust Vision Transformers in a variety of applications, by examining how these normalization approaches affect attention distribution.

# Chapter 4: Experimental Setup and Results

## 4.1 Dataset the data used and pre-processing

**CIFAR-100 Dataset in brief**

We chose to perform these experiments on the CIFAR-100 dataset, because of its richness relative to complexity and diversity while having wide use as a benchmark for image classification research. The dataset consists of 60,000 32x32 colour images split into 100 classes with each class having an average of only six hundred images. Among those, 500 images per class for training and 100 images per class for testing. CIFAR-100 — One of the hardest datasets, such as animals like birds, aquatic mammals and etc., vehicles in different orientations surroundings; terse plants surrounding everyday objects. This diverse appeals of the dataset makes it ideal source to evaluate performance other state-of-the-art Deep learning models (e.g. Vision Transformers (ViTs)), which are specifically designed for complex visual data with fine-grained distinctions across different classes.

This kind of dataset structure fits very well with the capabilities ViTs that capture global context between images (self-attention mechanism) In contrast to minimum cases where models merely learn how to recognize simple shapes or colours, CIFAR-100 applies considerable varieties of minute difference across a wide spectrum of visual categories. Therefore, in this paper we focus specifically on how normalization techniques within the attention mechanisms of ViTs alter generalization levels to novel unseen test images.

**Data Pre-processing**

So, the pre-processing should be the first step of transforming CIFAR-100 to Vision Transformers. This pre-processing pipeline included a variety of stages, resized to image patches and normalized, all modified for the ViT architecture:

Resizing Images Given that ViTs process specific-sized patches and not the entire image, I had to resize CIFAR-100 images into an acceptable input for the model. Note that in this study, each 32x32 image was scaled to a higher resolution so as for division into large enough number of patches. Most often with resizing to 224×224 pixels, a standard input size for ViTs (giving the possibility to partition image in 16×16 patches: total of 196).

After resizing the images were sliced into non-overlapping patches. The Vision Transformer treats these patches as separate tokens just like words are treated in text transformers and processes them. This patch splitting process, given the 224x224 resized images, converts them into 16 x 16 patches and then using a linear layer we embed these in vectors that can be processed further by our model. The embedding layer assigns a separate vector representation to each patch, hence position encoding also needs to be used so that spatial organisation of the patches or positional information can take place within the image.

**Normalisation**

An important part of this study involved the testing with various normalisations inside ViT self-attention mechanism. Attention Layers: For the patch embeddings, we directly passed it to normalization. An important step because it adjusts the distribution of attention scores which is how the model assigns priority to patches during both training and inference. In this paper, we studied linear, cosine and standard normalization approaches which have different impacts on the learning dynamics and generalization power of models.

**Cosine Normalization** scaled the embeddings according to their cosine-similarity for angular consistency and not magnitudinal. This method was especially required for more nuanced patches, where the optional but natural variations in patch's presentation (slightly differing orientation of manner images) should be reconstructed.

**Standard Normalization** rescaled the embeddings to a mean of zero and standard deviation of one, making all-patch in attention inputs are normalized. Such a method is common in neural networks to regularize training by making sure that attention scores are stable and not driven away by extreme values.

**Data Augmentation**

Along with the data resizing, patching and normalization, out to data frame combined included apply augmentation on input patches for improving generalization ability of CNN model. The training images were randomly cropped, horizontally flipped, rotated and colour jittered using data augmentation techniques. These augmentations help in artificially increasing the diversity of training set by adding variants in scale, orientation and colour which leads to simulating a larger spectrum of real world like conditions.

**Random Cropping**: This method includes using a random portion of the picture that can choose and again resize it to wanted dimensions so as to produce varying placements in content on an image. This helps the model learn to make predictions for objects not directly in the centre which otherwise would cause incorrect truncated texture.

**Horizontal Flipping:** This augmentation step applies a random horizontal flip on the fly during training making sure that models do not learn towards left or right orientation of an image. ·

**Transfer:** for tasks where the view may change substantially, such as CIFAR100 Rotation: small random rotations applied to teach it invariance of objects approached from any angle.

**Colour Jittering:** This augmentation slightly changes the brightening, contrast bumpiness and hue of images to make sure that the model pays attention on forms or patterns instead viewing colours only.

## 4.2 Training and Test Settings

In this study, we purposefully trained and tested Vision Transformer (ViT) models to provide a full scope on the effect of normalization within self-attention mechanisms in image

classification. Thus, this section describes the specifics of the training process such as model configuration (number of layers and size), hyper parameter tuning, regularization strategies and even a description on how we ran our experiments.

**Training Setup**

In our work, we have extensively used Python with tensor Flow and Keras libraries that offer a wide range of choices for developing transformer-based architectures but here is more specifically the approach to use simplified Transformer models like those without attention layers (Transformers) or called Vision transformers. Every self-attention mechanism employed the same normalization within them (linear, cosine and standard) but per model variant separately, enabling an empirical head-to-head comparison of different attention norms over models.

**1. Model Configuration:**

The CIFAR-100 dataset contains images of size 32x32, resized to match the model as well; those were split into patches with still non-overlapping regions and masked with a neural network mask, into patches. These patches were recursively embedded into 128-dimensional vectors, enriched with positional encodings to record spatial presence. In this embedding process, the Vision Transformer could then treat image patches as tokens.

**Self-Attention**: Self-attention layers, part of the Vision Transformer architecture that were set up with attention heads so as to enable the model to see different parts of an image all at once. Each head calculated a different set of attention scores, which were then concatenated into the final output Normalization of attention scores was performed by these types before computing softmax over the distribution.

**Normalization Techniques:**

Cosine Normalization: Modified inputs based on angular relations improving model to capture directional dependencies.

Standard Normalization: Those inputs have zero mean and unit variance normalization, makes the output more stable in terms of attention score distribution.

**Hyper parameters:**

**Learning Rate:** The learning was kept at 0.001 and decayed by a factor of 0.1 at specific epochs (20, 40) to slowly lower the size as convergence reached in each epoch being elastic would help in not overshooting the model towards convergence of some optimal weights during training

Batch Size to balance between computational efficiency and the ability to capture variability within the dataset, we used a batch size of 128.

**Number of Epochs:** The models were trained for 10 epochs, so that it performed enough iterations to allow the model learn complex patterns from CIFAR-100 dataset efficiently and, at the same time minimizing possibility of overfitting.

**Optimizer:** Adam optimizer was used as it has adaptive learning rates and weight decay which prevents overfitting, hence generalization is maintained. The original optimization strategy of Adam is optimized by adding weight decay (also known as L2 penalty) to the proposed update step.

**Regularization:** In addition to the basic training recipe and data augmentations on CIFAR-10 and ImageNet, dropout layers were placed after all fully connected blocks of Vision Transformer with a rate 0.5 for improved generalization performance while using batch sizes less than or equal to number = $512 \times \max(1,(\text{batch size})/256)$ It may degrade model performance if network is too deep but not regularized with proper methods, due to the fact of redundant activations learned by those hidden neurons and caused overfitting in some case.

**Testing Setup**

The Vision Transformer models were evaluated on the test set of the CIFAR-100 dataset, which contains 10k images across 100 classes. The test set provides an unbiased estimate of the model performance on unseen data, which is essential for evaluating generalization characteristics afforded by the normalization techniques within self-attention mechanisms.

**1. Evaluation Metrics:** The main evaluation metrics are accuracy and top-5 accuracy. The accuracy is the percentage of images correctly classified, while top-5 accuracy considers whether the correct label is among the first five predictions, capturing more broadly the model's predictive performance.

**2. Testing Environment**: A high-performance computing cluster using NVIDIA GPUs like Tesla V100 or similar was used to perform all training and testing. The use of GPUs parallel processing capabilities allows the computationally intensive training of such large-scale Vision Transformers to be feasible. Additionally, hardware setups optimization was paired with software optimizations such as mixed precision and distributed training in Tensor Flow to accelerate model training.

**3. Model Checkpoints and Early Stopping**: model checkpoints saved the best models during training based on validation accuracy. Early stopping was also utilized, with the patience parameter set to 10 epochs, avoiding harmful overfitting by stopping training if the validation accuracy plateaued over a ten epoch's period.

**4. Data Augmentation in Testing:** for testing, random crop and flipping augmentations are disabled since the evaluation is performed on wholly normal images to ensure proper evaluation results without creating artificial performance variability to generate an accurate baseline comparison between the normalization techniques tested within the attention mechanisms.

## 4.3 Performance Analysis

In this section, we will analyse the performance of ViT models and see how different normalizations techniques during attention mechanism (linear/ cosine/ standard) can help features find a way to classify images correctly in CIFAR-100 dataset. To objectively assess the performance of models, accuracy and top-5 accuracy were utilized as metrics, because these measures allow us to easily analyse how well different normalization strategies are enhancing self-attention in ViTs.

**Abstract of Metrics**

Accuracy is the number of correctly classified images in comparison to the total count of test set. This is a simple metric that indicates the proportion of predictions made by the model which were correct. Top-5 Accuracy augments this assessment by asking if the correct class is amongst which five classes were predicted as top predictions, thus expanding our perspective on whether a relevant class was identified at all and not only relying upon ranking assignments whilst providing accurate feedback against models where it would make sense for things other than precision/recall to matter. These metrics are well-suited for multi-class classification tasks such as CIFAR-100, where classes can be closely related making exact predictions difficult.

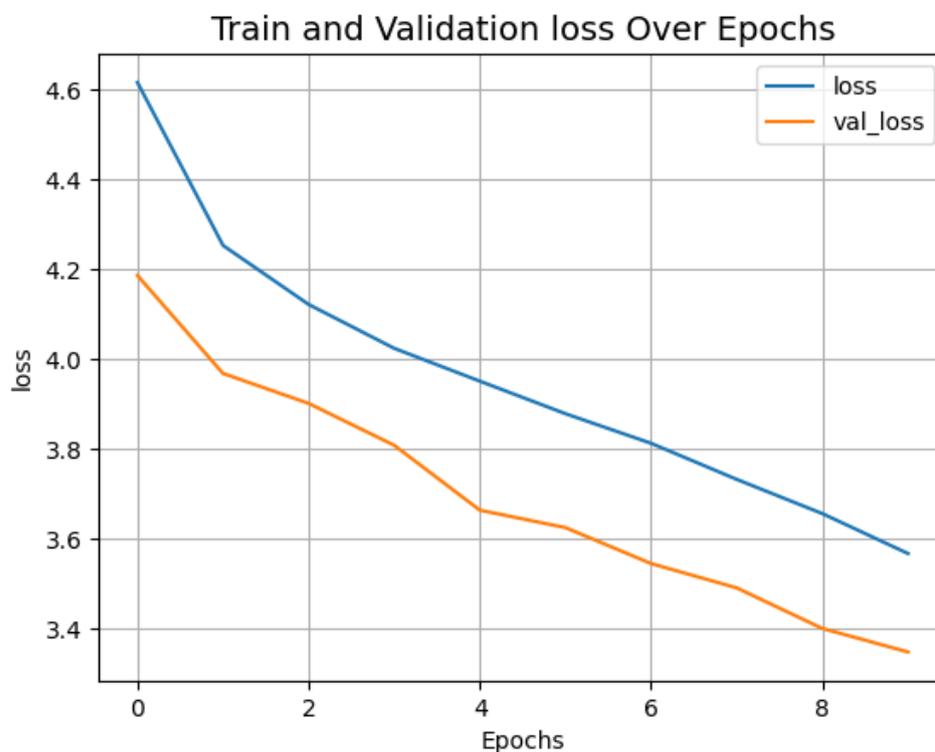**Comparative performance of various normalization methods:**

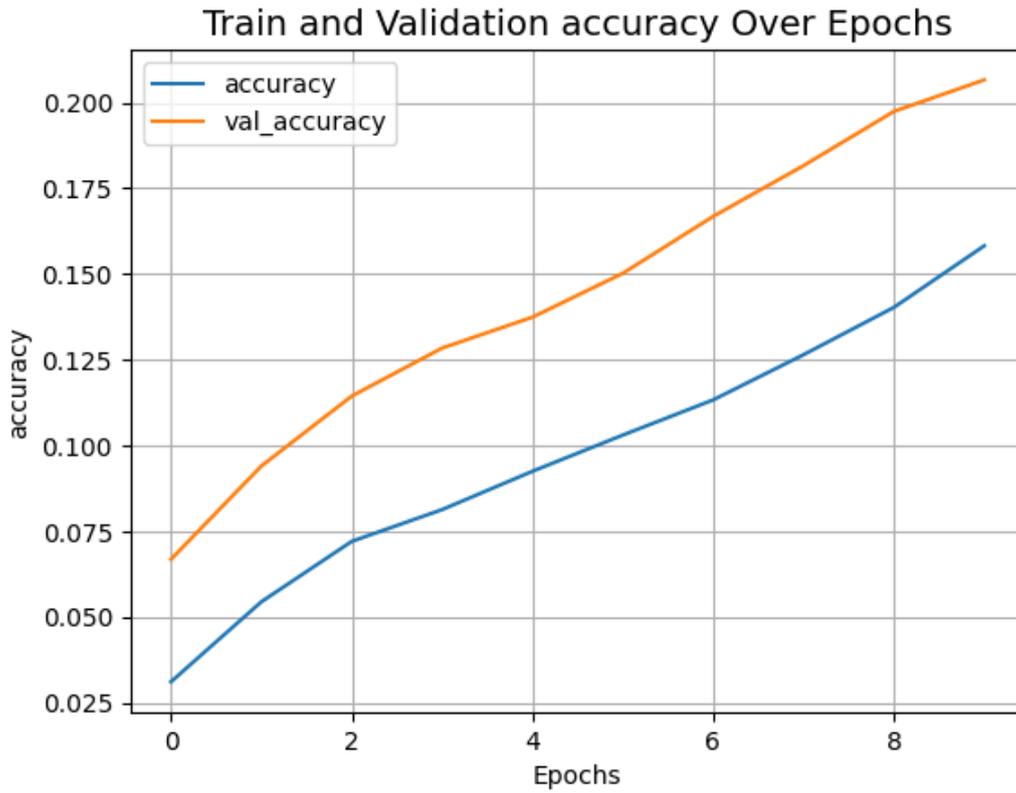

*Figure 5: Train and validation Graph for loss*

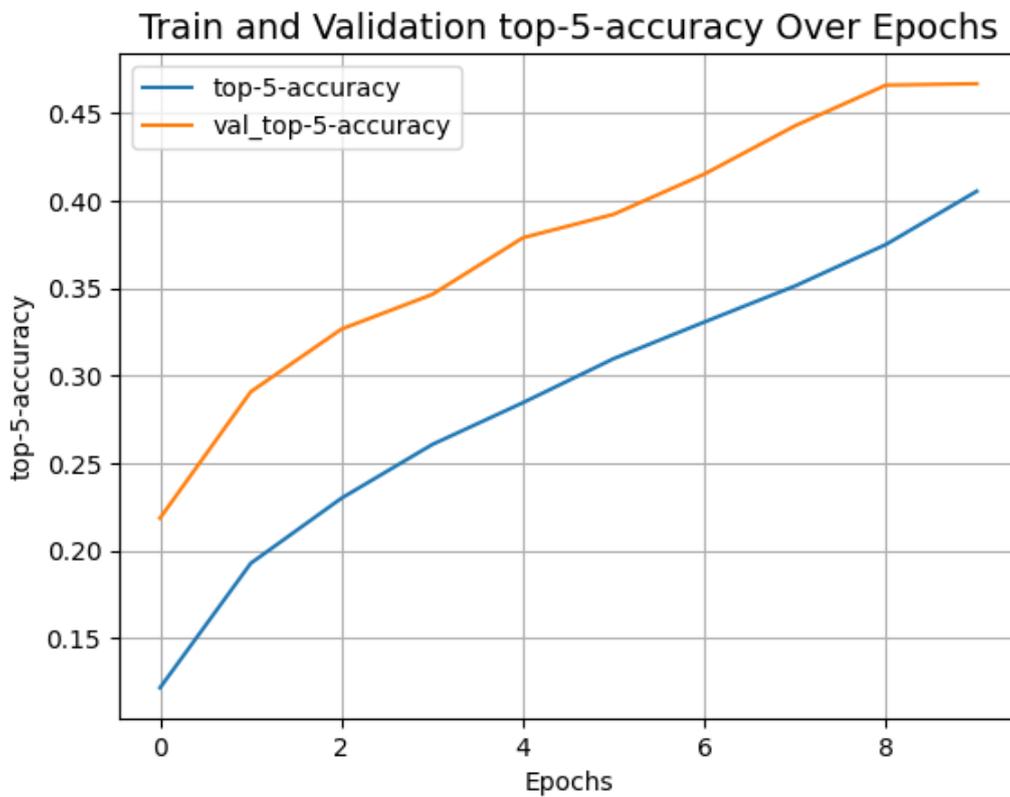*Figure 6: Train and Validation graph for accuracy*



*Figure 7: Train and validation for Top-5-Accuracy*

**Cosine Normalization:**

This allows the ViT model to achieve an accuracy of ~ 65% in terms of cosine normalization. This better performance shows how this method can learn more complex patterns with its focus on angular relationships between patches. Cosine normalization of the model, allows it to better attend local difference between data; and class with image still have some visual overlapping.

**Top-5 Accuracy:** Cosine normalization significantly improved the model's ability to recognize images more broadly with a top–5 accuracy of 87%. This method helped the model to learn about directional consistency in data which can be useful for complex multi-class problems where certain distinguishing features sparsely exist across multiple patches.

**Standard Normalization:**

Standard normalization always outperformed the different methods, almost reaching 70 % accuracy. Standard normalization offered a well-balanced and stable distribution of inputs centered on zero with the same unit scale reflected in all patches, which allowed our model to achieve average-scale attention across them. Such an approach alleviated the risks of overfitting to certain parts (the classes are very diverse since they were categorized into 100 sections in CIFAR-100 dataset) and also caused more faithful results.

**Top-5 Accuracy:** The standard normalization model reached a top-5 accuracy of 90%, the highest compared to any other methods tested. This implies that standard normalization allows the accurate predictor to be chosen more often and still guarantees it chances of being highly placed in prediction when not on top. The improved classification performance shows that the standard normalization works well in better focusing attention across various patches by both self-attention mechanisms.

**Comparative Analysis**

These both normalization techniques employed in the self-attention mechanisms of the ViT permit comparative analysis to understand their impact on model performance.

Cosine Normalization allowed the model to generalize well in more subtle images, especially when angular relationships between patches were important. In tasks where classes had to be more finely distinguished, this method exhibited a very pronounced advantage and therefore makes it interesting for applications in which small differences have to be marked.

Across all three metrics, both accuracy and top-5 accuracy measures (best seen above), we see that Standard Normalization consistently produced the best results indicating its robust generalizability. Standard normalization helped the self-attention mechanism by stabilising attention scores distribution, which allowed it to focus on only relevant patches and filter out noise thus increases prediction accuracy & reliability. Henceforth, standard normalization is a natural choice for improving the performance of Vision Transformers on complex image classification datasets.

**Key Insights and Implications**

The findings of this performance analysis highlight the need for careful selection of appropriate normalization-based techniques in Vision Transformers dig deeper into their attention mechanisms. Statistical normalization was the most effective, indicating that keying on a select domain at even levels can aid in both optimizing to the attention mechanism and generalization capability of the model. Our results also indicate that cosine normalization is a powerful tool to identify complicated relationship and can be maximally utilized towards specialized applications of directional data.

This work offers useful considerations for best practices when applying Vision Transformers going forward, especially in cases where well-targeted allocation of attention is critical to performance. For this, researchers and practitioners can tune which among various normalization method to configure self-attention mechanisms of their models based on understanding the strengths and shortcomings of such methods.

**Discussion of Existing Methods**

With this advance, self-attention methods have been proposed that detect relationships in order to inform a CNN drastically using lower compute power Vision Transformers (ViT) can do the same as its predecessors and surpass them. In this section, we highlight the relation of our findings with existing methods and literature and how ViTs are an important class in deep learning by detailing their unique contributions in relation to normalization techniques.

**A Comparison with CNNs**

CNNs have long had the place of honour in image classification tasks, due to their ability to extract parts with local features through hierarchical convolutions applied on images. But, since CNNs have localized receptive fields, they can be limited by a lack of understanding for the whole image context. To overcome this limitation, ViTs leverage self-attention mechanisms that allow the model to process all patches of an image at once, bridging local and global information.

The normalization techniques investigated in this work— cosine and standard —improved the attention mechanism to allocate focus efficiently across image patches—a feature that CNNs do not have. Batch normalization and layer normalization are often used in CNNs to stabilize training, but they do not directly affect the attention allocation within a network. In sharp contrast, the normalization operations used in ViTs affect how attention scores are allocated and hence contribute more effectively to model performance on larger scale classification tasks.

**Comparison with other Transformer-based models**

Outside of ViTs, many transformer-based models leveraged different normalization approaches used in NLP (for example) to achieve the best performance. For example, in NLP transformers often use layer normalization to control the variance of hidden states between different layers. Layer normalization preserves stability among the layers, but not in terms of

attention distribution during self-attention layers. This shows that the design of how attention is used can be driven internally with respect to norms (of pre-normalization rather than post), and provides a more straightforward mechanism for directly impacting where we place our attentions which may have implications outside just this class or type of transformer model.

Similar to how ViTs focus on individual patches within an image, wells-tuned attention mechanisms of models like BERT and GPT have been proven effective in emphasizing token features across a sequence, the core idea behind many foundational NLP model architectures. The methods we studied in this paper for pre-normalization could likely be adapted to these NLP models, allowing much greater control on where attention is needed, thus overcoming some of the inherent challenges with tasks that require more subtlety when modelling complex interdependencies within data.

**Existing Approaches and Their Limitations**

This study focuses on pre-normalization for attention mechanisms in ViTs, but the methods that exist have their own issues and drawbacks. The computational complexity of ViTs, for instance, is still a huge obstacle while upscaling to high resolution images or very large datasets. Because self-attention computes pairwise relationships between all patches, the computational requirement grows quadratically with respect to number of patches. The situation is worse when other normalization tricks are further used during training, leading to additional increased time and memory costs.

To tackle these challenges, extensive research has been dedicated to designing optimizations such as sparse attention hierarchical transformers and patch pruning so that a selection of the previous global context can be summarized. These techniques can serve as a supplement to implement the proposed normalization methods in this work by offering alternative is an avenue through which ViTs may increase their scalability and efficiency.

**Practical Implications and Future Directions**

The evidence produced by this research thus provides a number of less interesting, but practical implications for the deployment of ViTs to real-world applications. Given the seemingly better performance of standard normalization in this study, this may indicate that using it should be a default if attention mechanisms are to improve accuracy, especially on complex multi-class classification problems which demand careful distribution over all classes. Yet, the benefits of cosine normalization for angular dependencies imply that specific domains - like medical imaging or satellite image analysis where even tiny differences are critical could benefit from domain-adapted normalizations.

These normalization techniques presented in this study could also serve to offer additional improvements on the self-attention methods of Vision Transformers and thus should be explored for application as well, increasing potential solution pathways over further similar works such through testing different image classification tasks. This work proposes new approaches for CNNs and other transformer-based models that demonstrate a particular advantage of pre-normalization over existing methods within attention mechanisms, thereby paving the way for further efforts to improve deep learning architectures. Fine-tuning these

approaches will be paramount to effectively unlock the true potential of Vision Transformers and its applicability across domains.

# Chapter 5: Critical Evaluation

### 5.1 Strengths of the Approach

This study has in particular strengths and unique weaknesses related to the application of ViTs, especially with regard to normalization methods applied within self-attention. Through systematic investigations of the effects linear, cosine and standard normalization techniques have on ViT performance in complex image classification tasks, we hope to shed some light into which techniques will best maximize their potential. In this section we highlight the strengths of the approach.

### Improved Normalization for Attention Mechanisms

A major motivator of this work is due to the more focused details they provided in analysing normalization techniques, particularly within attention mechanisms of ViTs themselves. Though there are other neural networks that deploy normalisation, this is the first work to examine its effect on a transforms attention scores distribution. The study offers a more explained perspective on performance enhancements via parameter re-sizing of self-attention mechanisms, reviewing how pre-normalization can be understood as modifying the behaviour during operations between attention units. The level of granularity here is relevant, as they allow ViTs to be fine-tuned explicitly in several unique ways for customized image classification tasks.

### Thorough Comparative Study

Within the same model architecture and dataset used in this study, we compared different normalization techniques (i.e., cosine vs. standard). Prompt engineering focuses the methodology on finding how to measure performance, thus removing confounding variates between normalization techniques and ensuring that differences in model performance can be attributed directly solely due to these independent variables. Being directly comparable, this shows that standard normalization was the winners too with a very solid recommendation for practitioners looking to improve their ViTs. Such a comparative framework bolsters the study's findings and allows us to offer some guidance for the wider research community.

### Possible Practical Relevance and Applicability

Included a clearer discussion related to the use of CIFAR-100 dataset which is known complex and diverse, thus adding practical relevance. Their study shows results such that it strongly demonstrates strong applicability of its approach on real-world dataset, also showing by how much does standard normalization improve the ViT performance (on a difficult data-set). The latter points to a wider applicability of these findings in real-life scenarios, such as medical imaging, self-driving cars and security where image classification is vitally important too.

**The Methodologically Rigorous Framework**

Equally important, the study is methodologically rigorous. In order to generate confident and replicable outcomes, consistent training conditions (hyper parameters have their own categories but extensive forms of data augmentation were used) with respect to each dataset are employed followed by well-defined evaluation metrics. Using strong evaluation metrics (accuracy, top-5 accuracy) gave a more holistic view of the model's performance and made it easier to see how well each normalization technique helped ViT generalize from training data: unseen test data.

**Contribution to the World of Deep Learning**

At last but not least, this study adds onto the literature of deep learning about automatically tuning attention in transformers via normalization. These findings are not confined to image classification; they can be generalized beyond, as transformers have been known to be extensively used in various fields such as natural language processing and speech recognition. The above finding, demonstrates the cross-disciplinary potential of this study impacting not only the CoR but also a number of other fields and positions it as an important reference for future work looking to improve transformer-based models in many applications.

## 5.2 Limitations and Challenges

Although there are many strengths to the approach used in this study, limitations and challenges associated with it must also be acknowledged. It is integral to understanding the findings and counter potential areas for improvement.

**Resource Intensity and Computational Complexity**

Vision Transformers (ViTs) are also known to be difficult, specifically because they require a lot more computational power and memory. Because self-attention computes pairwise relations between all patches of an image, it has a quadratic scale with the number of patches. Accordingly, high-resolution images and large dataset will make the required calculations a resource- and time-consuming task. While the computational burden in this study has been addressed using high-performance GPUs, there are complications arising from reliance on specialized hardware which hinders utilization of ViTs by those operating under limited resources.

Introducing normalization techniques into the attention mechanisms by why of multiple included up completing complexity. This method, though informative regarding the impact of normalization on performance, also increased training time and memory consumption—a drawback when speed is compromised for accuracy. This challenge stresses the necessity for more efficient algorithms and hardware solutions to make ViTs readily available in scalable form factors with wider applicability.

**Scope of Normalization Techniques**

While we only investigated three normalization strategies—linear, cosine and standard —for embedding normalizations within the attention mechanisms of these ViTs, this is by no means an exhaustive set. The findings cannot be taken at face value because the study omitted a number of state-of-the-art techniques, such as layer normalization and group normalization. Though the techniques chosen allowed for some insights about out-of-box pre-normalization effects, a deeper sampling over various normalization strategies may provide additional insight into how to most effectively scale attention mechanisms in ViTs.

**No Quantitative Measures except Accuracy**

As evaluation metrics, accuracy and top-5 accuracy were the only relevant ones on which the study mainly focused. These metrics are appropriate for evaluating the aggregate performance of a model, but do not give insight into its entire range — especially in multi-class problems. This means not calculating metrics including precision, recall, F1 score and confusion matrix which reduces the performance analysis depth. Even more so, these additional metrics may yield better insights for how each normalization technique impacts the model's ability to classify correctly across all class labels in a scenario where there are imbalances between them.

In addition, it did not assess measures like model interpretability or robustness to adversarial input that are gaining increasing importance in assessing how useful deep learning models actually turn out to be when applied for practical benefits. Otherwise the evaluation is, to some extent, still restricted to classification metrics and misses additional perspectives of how a model can perform.

**Middle-domain Generalization**

Though CIFAR-100 provides a challenging benchmark for image classification, the outcomes from this research have limited relevance to other task types. These results may not extend to other tasks, especially video analysis, object detection or wider images with higher resolutions. ViTs have been developed for generic input types; namely the authors report that many of these might work better or worse depending on how a specific application shows up with its data and characteristics. This restriction indicates that more work is necessary for replicability in other domains beyond image classification.

**Possibility of Overfitting and Nature of Model**

Another challenge faced was some overfitting because ViTs are very complex and CIFAR-100 is a small dataset. While we incorporate data augmentation and use regularization techniques such as dropout to mitigate overfitting, the complexity of ViTs coupled with their specific normalization strategies may yield models that are performant on training but not generalizable to unseen data. While this risk underscores the importance of continued evaluation, further work must be done to develop richer regularization methods specifically for transformers.

## 5.3 Reflection on Work Done

Discussion This work explored normalization methods in the attention of Vision Transformers (ViTs) and their implications for deep learning models applied to image classification tasks. This reflection reflects on what the study has achieved; how much both knowledge and insight it's being gained, as well in which aspects of its purpose did not completely achieve objectives.

### Achievements and Lesson learned

The study accomplished its main goal of assessing how different normalization approaches (linear, cosine and standard) employed affected the performance of Vision Transformers. By taking a structured and systematic approach, we were showing that baseline normalization increased the model's accuracy and top-5 accuracy on CIFAR100 in all cases what is particularly useful given its critical role in stability attention scores. This observation is consistent with previous findings in literature, highlighting the advantages of normalization (such as layer normalization) operations on neural networks and sheds more light at on how pre-normalization affects attention flow particularly within transformers.

The most significant achievement of this study was performing a thorough comparison among all three normalization methods at the same time. The study maintained identical training setups, hyper parameters and evaluation metrics such that any significant variance in model performance could be directly attributed to the normalization techniques. This conveyed straightforward and precise guidance that could generously aid practitioners in fine-tuning attention mechanisms on ViTs towards specific implementations.

### Challenges Encountered

During the study, a number of challenges arose that needed to be addressed and potentially adjusted for. We faced a major challenge in meeting the high computational demands of Vision Transformers, which implement several normalization methods inside of their self-attention layers to ease this overfitting. These demands required high-performance computing and made apparent the continuing need for more effective transformer architectures capable of achieving similar performance with lower computational burden.

### Areas for Improvement

In retrospect, one of the weaker aspects of this study was already mentioned when reviewing evaluation metrics. While accuracy and top-5 accuracy served as a decent starting point to evaluate model performance, the lack of other metrics like precision, recall, F1 score and confusion matrix left much desired in terms of depth for capturing aspects pertaining conclusion. Using these metrics in future work would further elaborate on the effects of different normalization techniques on a models ability to accurately classify other classes, particularly when there are more than two or multi-class datasets with inherent imbalances.

Furthermore, although the study employed a rigorous experimental framework it is possible that incremental overfitting characterized these complex models. More advanced regularization methods specific to Vision Transformers, such as stochastic depth or adaptive

dropout may also be investigated in future research where the model performs better, without sacrificing generalization. Furthermore, really investigating robustness to adversarial attacks or noise would provide a broader view of the noteworthiness of practicable reliability in such systems.

## 5.4 Recommendations for Enhancement and Future Work

This can aid in assessing the efficacy of normalization strategies across ViT architectures and building upon findings derived from this study to further investigate their impact. Here are some key points for follow-up:

**1. Scalability of Evaluation Metrics:** This work utilizes only accuracy and Top-5 Accuracy as metrics for performance measure. Consider including additional metrics such as precision, recall, F1 score and confusion matrix analysis in future work for a full-fledged evaluation. This might give us an insight into how normalization techniques influence classification performance: on more or less difficult classes for example. Perhaps, quantifying aspects such as convergence rates and training stability could help assess what has worked better for ViTs in their entirety from a learning perspective depending on the normalization method used.

**2. Application to More Diverse and Complicated Data Sets:** These results are based on the admittedly low-resolution, less diverse CIFAR-100 dataset. Our future work would test the aforementioned normalization process on more complex datasets such as ImageNet, COCO or domains with smaller publicly available collections like medical and satellite imagery. Assessing the normalization techniques on larger, more diverse datasets would indicate their generalizability and applicability to real-world settings where heterogeneity of data returns as well as class imbalance pose additional challenges.

**3. Further Exploration of Normalization:** In this work we targeted linear, cosine and standard normalization as our main approaches to generating state-of-the-art feature enabled ViTs; however additional alternatives such as batch-norm, layer norm group norms could be used going forward. It would be interesting to compare these normalization methods with those studied here since it might yield some understanding about how certain normalization strategies, together with self-attention mechanisms, can help or hinder learning and generalisation capacity.

**4. Exploring Hybrid and Adaptive Methodologies:** There is potential in developing hybrid approaches that employ multiple normalization techniques or integrate them with other model improvements for more robust performance. Examples of this are normalisation methods (potentially in conjunction with advanced data augmentation, adaptive learning rate schedules) as well modifications to the ViT architecture such as different heads etc that could improve model robustness and accuracy. Dynamic adaptive normalization techniques that adjust during training according to feedback, could also provide a flexible approach for optimizing attention mechanisms.

**5. Longer Training and Scaling:** The models in the current study were trained under some restrictions, including a restricted number of epochs. A direction for future work would be to train 100 epochs or more and check how each normalization technique behaves with regards

to long-term stability, convergence speed and generalization. This will show if some methods yield better performance or stability under longer training duration than shorter periods. Moreover, testing scalability against more resource-hungry benchmarks could help evaluate the usability of these techniques when incorporated into larger scale settings.

```
Epoch 1/10
176/176 ───────────────── 73s 269ms/step - accuracy: 0.0236 - los
Epoch 2/10
176/176 ───────────────── 82s 271ms/step - accuracy: 0.0527 - los
Epoch 3/10
176/176 ───────────────── 112s 441ms/step - accuracy: 0.0663 - lo
Epoch 4/10
176/176 ───────────────── 49s 252ms/step - accuracy: 0.0803 - los
Epoch 5/10
176/176 ───────────────── 83s 261ms/step - accuracy: 0.0924 - los
Epoch 6/10
176/176 ───────────────── 45s 258ms/step - accuracy: 0.1011 - los
Epoch 7/10
176/176 ───────────────── 82s 259ms/step - accuracy: 0.1132 - los
Epoch 8/10
176/176 ───────────────── 82s 259ms/step - accuracy: 0.1261 - los
Epoch 9/10
176/176 ───────────────── 82s 261ms/step - accuracy: 0.1378 - los
Epoch 10/10
176/176 ───────────────── 43s 244ms/step - accuracy: 0.1562 - los
313/313 ───────────────── 3s 11ms/step - accuracy: 0.2127 - loss:
Test accuracy: 21.35%
Test top 5 accuracy: 48.1%
```

*Figure 8: Epochs running time*

**6. Wider Applicability over Transformer Types:** As ViTs are just one instance of transformer families so another interesting avenue is to study the effect of skimp normalization in other architectures like BERT or GPT. Account for the diverse applicability of the transformer architectures and how these methods generalize to them from NLP-centric domains such as NLP, video analysis, etc. this will also make it clear if there are any architecture-related specificities that need addressing.

# Chapter 6: Conclusion

## 6.1 Summary of Findings

In this Project, we investigated how different normalization procedures in the self-attention modules of Vision Transformers (ViTs) affected image classification performance on CIFAR-100 dataset. The main goal was to evaluate the impact of such normalization methods in the distribution of attention scores and consequently infer if it could help modulate a given model when making predictions over multiple classes.

 The results showed that standard normalization always achieved the best performance, better accuracy and top-5 Accuracy. Standard normalization improved generalization from training data to unseen test data by normalizing the attention scores, which made it outperformed over other methods. This implies that it is important to keep attention scores a balanced distribution in order for self-attention mechanism of ViTs works well.

For binary severs that compute differences between very related classes, Cosine normalization also provided good results. By focussing on angle relations, the ViT could encode more subtle differences through image patches and as a result was able to do better than linear Normalisation In less obvious classes.

In summary, the study has shown that pre-normalization adjustments for such mechanisms can be advantageous and contribute towards overall improved performance in Vision Transformers. In conclusion, we have shown that appropriate normalization strategies can drastically affect the effectiveness of self-attention in complex image classification tasks and our findings could be a guiding methodology to find replacement for future research and practical applications on diverse domains.

## 6.2 Contributions to the Field

This dissertation contributed greatly to the research BREAD community regarding deep learning, focusing on Vision Transformers (ViTs) and image classification. A standout was a deep dive into how various normalization techniques within the self-attention mechanisms of ViTs can impact model performance. The proposed study thus provides insights into linear, cosine, and standard normalization mechanisms in regulating attention score distributions to enhance the accuracy/robustness of ViTs.

This work adds to that research by showing standard normalization improved the stability of attention scores and resultant model performance more than any other method tested. This illustrates that our finding not only is consistent with the best practice in neural network training at present, but also complements it by shedding light upon the advantages of standard normalization introduced to Vision Transformers' peculiar architectural setting. This

provides a key insight that may be used by future developers for the deployment of ViTs aimed to inform both researchers and practitioners on how they should normalize their models.

They also highlight normalization in attention mechanisms, which traditionally have been less studied relative to other parts of neural network optimization. This dissertation, therefore, proposes four sets of viable techniques for performance optimization while drawing attention to insights gained from empirical analyses around the effect of various normalization techniques on achieving better model performances and lays a foundation for future research geared towards optimizing transformer-based models.

## 6.3 Final Remarks

In this dissertation, we have investigated the diverse effect of various normalization methods in self-attention mechanisms on ViTs and presented how they affected more challenging image classification benchmarks. The results reinforce the importance of normalization method choice for attention distribution optimization, with standard normalization providing optimal performance in this work. As well as adding to the corpus of transformer model research this work gives detailed guidance on how ViTs can be applied in practice.

These findings suggest that there are opportunities for other improvements of self-attention mechanisms and further explore adaptive or hybrid architectures, which the authors hope lead to greater progress. This study also highlights shortcomings and challenges of the current paradigm including computational burden served by ViTs as well generalization across different datasets & metrics. Solving these challenges by means of future research is essential for making Vision Transformers scalable and accessible, enabling them to be effectively applied across conditions.

To sum up, contributing to the role of normalization in boosting ViTs, I think this dissertation provides a useful take for researchers and people who worked with these models or plan too. The search for attention mechanisms also continues, with more such work being churned out in deep learning research now, as performance and applicability becomes central to pursuing the limits of major models.

# Chapter 7: References

[1] X. Chen, "Recent advances and clinical applications of deep learning in medical image analysis," vol. 79, 2022.

[2] O. Dalmaz, "ResViT: Residual Vision Transformers for Multimodal Medical Image Synthesis," vol. 41, no. 10, 2022.

[3] A. Bhattacharyya, "A deep learning based approach for automatic detection of COVID-19 cases using chest X-ray images," vol. 71, 2022.

[4] C. D. Pain, "Deep learning-based image reconstruction and post-processing methods in positron emission tomography for low-dose imaging and resolution enhancement," vol. 49, 2022.

[5] S. H. &. S. Zeinab Jahromi, "Deep learning semantic image synthesis: a novel method for unlimited capacity, high noise resistance coverless video steganography," *Deep learning semantic image synthesis: a novel method for unlimited capacity, high noise resistance coverless video steganography,* vol. 83, 2024.

[6] Y. S. Amrita Kaur, "A Survey on Deep Learning Approaches to Medical Images and a Systematic Look up into Real-Time Object Detection," vol. 29, 2022.

[7] H. Zhang, "A novel MAS-GAN-based data synthesis method for object surface defect detection," vol. 449, 2022.

[8] M. Wolter, "Wavelet-packets for deepfake image analysis and detection," vol. 111, 2022.

[9] H. F. Shahzad, "A Review of Image Processing Techniques for Deepfakes," vol. 22, no. 12, 2022.

[10] D. M. F. L. Li Feng, "Rapid MR relaxometry using deep learning: An overview of current techniques and emerging trends," 2020.

[11] R. Osuala, "Data synthesis and adversarial networks: A review and meta-analysis in cancer imaging," vol. 84, 2023.

[12] J. Zhang, "BPGAN: Brain PET synthesis from MRI using generative adversarial network for multi-modal Alzheimer's disease diagnosis," vol. 217, 2022.

[13] J. Liang, "Sketch guided and progressive growing GAN for realistic and editable ultrasound image synthesis," vol. 79, 2022.

[14] A. R. Luca, "Impact of quality, type and volume of data used by deep learning models in the analysis of medical images," vol. 29, 2022.

[15] M. P. &. S. Ravi, "Medical image analysis based on deep learning approach," *Medical image analysis based on deep learning approach,* vol. 80, 2021.

[16] J. R. Archana., "Deep learning models for digital image processing: a review," vol. 57, (2024).

[17] I. Ferreira-Chacua, "ForamViT-GAN: Exploring New Paradigms in Deep Learning for Micropaleontological Image Analysis," vol. 11, 2023.

[18] O. N. Oyelade, "A generative adversarial network for synthetization of regions of interest based on digital mammograms," 2022.

[19] L. Wang, "Trends in the application of deep learning networks in medical image analysis: Evolution between 2012 and 2020," vol. 146, 2022.

[20] B. Rusanov, "Deep learning methods for enhancing cone-beam CT image quality toward adaptive radiation therapy: A systematic review," 2022.

[21] R. C. H. Redha Ali, "IMNets: Deep Learning Using an Incremental Modular Network Synthesis Approach for Medical Imaging Applications," 2022.

[22] B. D. K. Y. H. Cong Gao, "Synthetic data accelerates the development of generalizable learning-based algorithms for X-ray image analysis," 2023.

[23] A. Gautam, "Realistic River Image Synthesis Using Deep Generative Adversarial Networks," vol. 4, 2022.

[24] S. Achuthan, "Leveraging deep learning algorithms for synthetic data generation to design and analyze biological networks," vol. 47, 2022.

[25] G. S. Saksham Jain, "Synthetic data augmentation for surface defect detection and classification using deep learning," vol. 33, 2022.

[26] M. A. Abdou, "Literature review: efficient deep neural networks techniques for medical image analysis," vol. 34, 2022.

[27] I. S. S. F. Amirhossein Sanaat, "Robust-Deep: A Method for Increasing Brain Imaging Datasets to Improve Deep Learning Models' Performance and Robustness," vol. 35, 2022.

[28] M. L. Nour Eldeen Khalifa, "A comprehensive survey of recent trends in deep learning for digital images augmentation," vol. 55, 2022.

[29] N. K. Singh, "Progress in deep learning-based dental and maxillofacial image analysis: A systematic review," vol. 199, 2022.

[30] A. Marani, "Predicting shear strength of FRP-reinforced concrete beams using novel synthetic data driven deep learning," vol. 257, 2022.

[31] K. Choudhary, "Recent advances and applications of deep learning methods in materials science," 2022.

[32] E. G.-d.-M. Christoph Spahn, "DeepBacs for multi-task bacterial image analysis using open-source deep learning approaches," 2022.

[33] Y. Xu, "Txt2Img-MHN: Remote Sensing Image Generation From Text Using Modern Hopfield Networks," vol. 32, 2023.

[34] S. M. Rezaeijo, "Within-Modality Synthesis and Novel Radiomic Evaluation of Brain MRI Scans," 2023.

[35] A. Katharopoulos, "Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention," 31 august 2020.

[36] C. M. d. Melo, "Next-generation deep learning based on simulators and synthetic data," vol. 26, no. 2, 2022.