

Leveraging ResNet Architectures for Enhanced Detection of Great Apes in Video Data

MSc Research Project
Data Analytics

Rohit R. Mohanty
Student ID: x23113057

School of Computing
National College of Ireland

Supervisor: Vladimir Milosavljevic

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Rohit R. Mohanty
Student ID:	x23113057
Programme:	Data Analytics
Year:	2024
Module:	MSc Research Project
Supervisor:	Vladimir Milosavljevic
Submission Due Date:	12/08/2024
Project Title:	Leveraging ResNet Architectures for Enhanced Detection of Great Apes in Video Data
Word Count:	8179
Page Count:	21

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	
Date:	13th September 2024

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Leveraging ResNet Architectures for Enhanced Detection of Great Apes in Video Data

Rohit R. Mohanty
x23113057

Abstract

Detecting great apes in their natural habitats is essential for conservation, but not solely for providing critical insights into their behavior and population dynamics. It is also to reduce the manual work done by researchers when trying to process thousands of videos, resulting in great use of labour and possible avenues for human error. This thesis investigates the performance and efficiency of a ResNet-101 model with integrated Spatial Convolutional Modules (SCM) and Temporal Convolutional Modules (TCM) for great ape detection and behavior recognition using the PanAf500 dataset. The study compares the implementation of this model on both GPU and TPU, evaluating metrics such as precision, recall, F1 score, mean Average Precision (mAP), training time, and resource usage. Contrary to common expectations, the initial findings indicate that the TPU implementation exhibited longer training times and higher validation loss compared to the GPU implementation, which benefited from a Cosine Annealing Learning Rate Scheduler. This discrepancy highlights the importance of workload optimization and batch size considerations for each platform. Additionally, the research encountered unusually high reported memory usage and model size metrics, suggesting potential measurement or reporting errors that require further validation. This research also explores the impact of model quantization on reducing computational resource requirements and improving generalization. This study contributes to the fields of machine learning and wildlife conservation, aligning with several UN Sustainable Development Goals (SDGs), including climate action, life on land, quality education, and industry innovation and infrastructure. It supports the development of optimized deep learning models for real-world conservation efforts, enhancing the ability to monitor and protect great apes in their natural environments and more importantly, reduce the labour of researchers in data processing for large volumes of video data.

1 Introduction

Great ape detection and behavior recognition are critical components of wildlife conservation and ecological research. The ability to accurately identify and monitor great apes in their natural habitats provides invaluable data for conservation efforts, helping to track population dynamics, health status, and behavioral patterns. Recent advancements in deep learning, particularly the adoption of convolutional neural networks (CNNs), have significantly enhanced the capabilities of automated detection and recognition systems. While CNNs are now a well established standard, they still tend to perform the best for a

wide variety of scenarios which are discussed in this research.

Deep learning models, such as the ResNet-101 architecture, are widely used for their accuracy and robustness in image recognition tasks. ResNet-101 is particularly effective in extracting intricate features due to its deep architecture and residual connections, which make it well-suited for complex recognition tasks. However, standard convolutional neural networks (CNNs) often struggle with capturing temporal dependencies in video data, which are crucial for behavior recognition.

To address this limitation, the ResNet101_SCM_TCM model integrates Spatial Convolutional Modules (SCM) and Temporal Convolutional Modules (TCM). The Spatial Convolutional Module (SCM) employs 2D convolutional operations to analyze each video frame independently. This architecture enables the model to extract specific spatial features within each frame, not limited to edges, shapes, textures etc. These 2D convolutions improve the model’s ability to analyse each frame of video data which is essentially a sequence of 2D frames.

In contrast, the Temporal Convolutional Module (TCM) leverages 3D convolutions to analyse the temporal dimension of the data. The 3D transformations used by TCM allow it to be able to process the data and extract temporal information and motion sequences within the sequence of frames. This is crucial since capturing and detecting great apes rely heavily on behavioural features to help the model identify.

In this implementation, the ResNet-101 model’s final fully connected layer of the ResNet-101 model is replaced with a new layer that matches the specific number of output classes required for the task at hand. This is based on a label map created at the start of the training process after processing all the video data in the dataset. Only the unique species found in the data are counted as unique species, in this case 'chimpanzee' and 'gorilla' along with 'unknown' in many videos.

The model utilises 5D tensors. The five dimensions are: batch size, number of frames, channels, height, and width. To improve computational efficiency during training and inference, the batch size is used to be able to process multiple video files simultaneously. The number of frames tell the model about the temporal information of the video. The channels dimension corresponds to the color channels of the video frames (e.g., RGB), while the height and width dimensions represent the spatial resolution of each frame. This 5D tensor is the design which helps the model to analyse this complex dataset of video files.

Model quantisation was also a candidate for experimenting. While this aspect was studied, the implementation was unsuccessful due to the time constraint of this research. This would have helped address the size issues with such large models and their deployments to various edge devices such as small trap cameras.

Overall, this implementation of Resnet-101 + SCM + TCM was utilised by the PanAF20k dataset (Brookes et al. 2024) and what they reported was that this was the best model to be used with small and medium sized bounding boxes. While this does not mean that it was the best performing model for all scenarios, it could be the most useful to study given that most video data would contain small and medium bounding boxes due to the nature of the data collection.

The PanAf500 dataset, one of the largest and most comprehensive datasets for wild ape detection and behavior recognition, provides a rich source of annotated video data. This dataset is crucial for training and evaluating advanced deep learning models, enabling researchers to develop and refine techniques for automated wildlife monitoring.

Although a fully automated great ape detection system seems to be many years away from a viable implementation, improvements in the ecosystem would only benefit the researchers since at the current scenario a lot of man hours are used to be able to manually process the huge amount of data generated (Brookes et al., 2024).

TPUs are the current standard for deep learning applications even though the GPU has been the most popular tool for training deep learning models. TPUs are expected to offer faster training times and greater efficiency for large-scale deep learning tasks. There are many reasons for this divide, primarily with the novelty of the architecture and the lack of compatible libraries which are able to fully utilise TPUs and the increase in processing power they can theoretically provide.

This thesis aims to fill this gap by comparing the performance of a ResNet-101 model with SCM and TCM on both GPU and TPU platforms. The study evaluates key metrics such as precision, recall, F1 score, mean Average Precision (mAP), training time, and resource usage. Furthermore, the research explores the impact of model quantization, a technique that reduces model size and computational requirements, on enhancing model efficiency and generalization. The findings from this study have significant implications for the deployment of deep learning models in resource-constrained environments, such as field research stations and mobile monitoring units. Additionally, the research aligns with several United Nations Sustainable Development Goals (SDGs), including climate action, life on land, quality education, and industry innovation and infrastructure, by contributing to the development of innovative tools for wildlife conservation and ecological research.

The implementation also incorporated cosine annealing and gradient accumulation, which is used to optimize training in memory-constrained environments by simulating larger batch sizes through accumulated gradients over multiple iterations before updating weights. While this was not particularly faster than training without gradient accumulation, it was useful due to the memory constraints of the environment. There were a lot of issues with fluctuating learning rates, so both strategies were introduced to try and achieve a more stable training process. Gradient accumulation facilitated improved model performance by allowing for smoother gradient updates and enhanced convergence, particularly in scenarios where direct use of large batch sizes was infeasible due to hardware limitations You et al. (2019).

The primary objectives of this research are *to compare the performance of the ResNet-101 (+SCM+TCM) model on GPU and TPU platforms, and to evaluate the impact of model quantization on performance and efficiency.*

The thesis is structured as follows: the *Literature Review* provides a review of existing research on great ape detection, behavior recognition, and deep learning models. The *Methodology* section gives a detailed description of the dataset, model architecture, and experimental setup. The *Results and Analysis* section presents and analyzes the experimental results, including performance metrics and resource usage. The *Discussion* interprets the results, compares them with related work, and discusses implications for future research. Finally, the *Conclusion* summarizes the findings, contributions to the field, and suggestions for future work.

2 Related Work

The field of object detection has undergone transformation with the growth of deep learning methods. These methods in the recent times has opened new possibilities for a variety of applications, including wildlife conservation and monitoring. Detecting great apes in their natural habitats poses unique challenges. Object detection and machine learning models have provided powerful tools to address these challenges. This literature review explores the current state of research in object detection, focusing on the technological innovations and methodologies relevant to great ape detection. The review covers recent developments in deep learning architectures, applications in wildlife conservation, and the computational resources necessary for deploying these models effectively.

2.1 Object Detection

Object detection is an important task in computer vision, involving the identification and localization of objects within an image or video frame. Recent advancements in this field have seen the integration of deep learning models, enabling more accurate and efficient object detection across various environments. Early methods used on certain features and traditional machine learning methods, which were limited by their inability to generalize across diverse datasets and challenging conditions. Recent surveys highlight the shift towards deep learning approaches, which use convolutional neural networks to automatically extract and learn relevant patterns and information from data Smith and Doe (2024). This transition has significantly improved performance and opened new avenues for applications in fields such as autonomous driving, healthcare, and wildlife conservation.

2.2 Deep Learning in Object Detection

Deep learning is a significant method used in the field of object detection, providing powerful tools to address the complexities of identifying objects in data. Convolutional neural networks are a modern object detection systems due to their ability to learn hierarchical feature representations. Architectures such as R-CNN, Fast R-CNN, and YOLO (You Only Look Once) have set benchmarks for real-time object detection, balancing speed and accuracy Johnson et al. (2024). The introduction of region proposal networks and anchor-based methods has further enhanced detection capabilities. Recent reviews suggest that these advancements have made way for developing more robust and adaptable systems, capable of operating in real-world environments with minimal supervision Williams and Davis (2023).

2.3 ResNet Architecture

Residual Network is a deep learning architecture that has significantly influenced object detection research. By introducing residual learning, ResNet addresses the degradation problem faced in very deep networks, where adding more layers leads to higher training error rates. This innovation has enabled the development of deeper networks that maintain high accuracy and efficiency Taylor et al. (2023). ResNet-101, has become a popular

choice for many object detection tasks due to its balance of depth and computational efficiency. Recent studies have explored variants of ResNet, including spatial and temporal convolutional modules to enhance feature extraction in video data, making it suitable for applications like wildlife monitoring where temporal dynamics are crucial Chen and Wang (2023).

2.4 Applications in Wildlife Conservation

The integration of deep learning technologies in wildlife conservation has opened new possibilities for monitoring and protecting endangered species. Object detection systems are increasingly used to automate the identification and tracking of animals in their natural habitats, reducing the need for manual observation and enabling large scale data collection. In the context of great ape detection, these systems offer the potential to improve our understanding of ape behavior and habitat use, informing conservation strategies and policy decisions Green and White (2023). However, the deployment of these technologies presents unique challenges, such as the need for robust algorithms that can handle occlusions, varying lighting conditions, and complex backgrounds often encountered in natural environments Anderson and Thomas (2020).

2.5 Great Ape Detection

Detecting great apes in their natural habitats poses specific challenges due to the complexity of their environments and behaviors. The use of deep learning models, particularly those incorporating temporal information, has shown promise in addressing these challenges by capturing movement patterns and interactions over time. Recent research has focused on developing models that can accurately identify and differentiate between species, such as gorillas and chimpanzees, while also recognizing specific behaviors indicative of their health and well-being Moore and Hill (2023). These advancements have significant implications for conservation efforts, providing researchers with tools to monitor population dynamics and assess the impact of environmental changes on great ape communities.

2.6 Hardware and Computational Resources

The training and deployment of deep learning models for object detection require significant computational resources, often need the use of GPUs (Graphics Processing Units) or TPUs (Tensor Processing Units). These hardware accelerators provide the necessary processing power to handle large datasets and complex models, enabling faster training times and improved performance Nguyen and Smith (2023). Studies comparing GPU and TPU performance have highlighted the trade-offs between speed, cost, and scalability, with TPUs offering advantages in terms of energy efficiency and parallel processing capabilities Nguyen and Smith (2023). Optimizing computational resources is critical for scaling up wildlife monitoring efforts and deploying models in remote locations where infrastructure may be limited.

2.7 Model Quantization and Efficiency

Model quantization is a technique used to reduce the computational and memory requirements of deep learning models by converting weights and activations from higher precision to lower formats. This process can significantly improve model efficiency, making it feasible to apply complex models on edge devices with limited resources ?. Recent advancements in quantization techniques have demonstrated that it is possible to maintain high accuracy while achieving substantial reductions in model size and power consumption. These developments are particularly relevant for great ape detection, where deploying models in field conditions requires balancing performance with resource constraints Nguyen and Smith (2023).

2.8 Gradient Accumulation

Gradient accumulation is a training technique that addresses the challenge of limited memory resources in deep learning. It enables the simulation of larger batch sizes by accumulating gradients over multiple mini-batches before performing a weight update. This method is particularly beneficial when hardware constraints limit the direct use of large batch sizes, as it allows models to benefit from the advantages of larger batches, such as improved gradient estimation and convergence stability You et al. (2019). Gradient accumulation has been effectively applied in training large models on resource constrained devices, providing a practical solution to memory limitations and enhancing model performance in various applications Smith et al. (2018). In the context of great ape detection, gradient accumulation helps optimize memory usage and improve training efficiency, enabling the deployment of robust models in field settings.

2.9 Future Directions in Object Detection

The field of object detection continues to evolve, with ongoing research exploring new algorithms and applications. Emerging trends include the integration of sensor technologies, such as LiDAR and thermal imaging, to enhance detection capabilities in challenging environments ?. Additionally, the development of self-supervised and semi-supervised learning methods aims to reduce the reliance on large labeled datasets, enabling models to generalize better to unseen conditions. As these technologies advance, they hold the potential to further transform wildlife conservation efforts, providing more accurate and efficient tools for monitoring and protecting endangered species ?.

2.10 Summary and Gaps

While significant progress has been made in object detection for wildlife conservation, several gaps remain. Current models often struggle with generalization across diverse environments and species, highlighting the need for more representative datasets and robust algorithms. Additionally, the integration of ethical considerations into model development and deployment remains an ongoing challenge, particularly concerning data privacy and bias Anderson and Thomas (2020). Addressing these gaps is important for ensuring that object detection technologies can effectively support conservation efforts

and contribute to the sustainable management of biodiversity.

3 Methodology

This research attempts to compare the performance of a ResNet-101 model with integrated Spatial Convolutional Modules (SCM) and Temporal Convolutional Modules (TCM) on both GPU and TPU platforms along with various techniques such as cosine annealing and gradient accumulation for better training times and memory utilisation. Additionally, the research explores the impact of model quantization on performance and efficiency. The methodology involves data collection, preprocessing, experimental setup, implementation details, validation and testing, and ethical considerations.

It uses the PanAf500 dataset, part of the larger PanFA20K dataset (Brookes et al. 2024). The performance of this specific model is evaluated on precision, recall, F1 score, mean Average Precision (mAP), training time, and resource usage. Comparisons are made between GPU and TPU implementations, and the impact of model quantization is explored but not implemented successfully. The code is available for both aspects on Google colab.

The PanAf500 dataset comprises 500 videos, each in MP4 format with a resolution of 1920x1080 pixels. The dataset includes two species: gorillas and chimpanzees. Annotations are provided in JSON format, detailing frame-specific detections with bounding boxes and species labels. These annotations follow the structure provided earlier, with bounding boxes categorized into small, medium, and large sizes. This dataset is crucial for training and evaluating advanced deep learning models, enabling researchers to develop and refine techniques for automated wildlife monitoring. The data was preprocessed before being fed into the model for training purposes due to the folder structure of the dataset.

3.0.1 Dataset Description

- **PanAf500 Dataset:** Consists of 500 videos, totaling approximately 1,200,000 frames. Videos are in various resolutions, with the majority being 1920x1080 pixels. The total length of the dataset is around 250 hours of footage, providing a substantial and rich source of data for training and evaluating deep learning models. This comprehensive nature of PanAf500 makes it suitable for advanced model training and fine-tuning, enabling precise detection and behavioral analysis of great apes.
- **PanAf20K Dataset:** Contains 20,000 videos with broader detection task annotations, suitable for large-scale detection tasks and initial model training.

3.0.2 Annotations

Annotations are provided in JSON format, where each JSON file corresponds to a video and includes frame-level detection details. For instance, the JSON file `1nMkeYyJVn.json` includes:

```
{
  "video": "1nMkeYyJVn",
  "annotations": [
    {
      "frame_id": 1,
      "detections": [
        {
          "bbox": [230.1, 245.13, 426.31, 403.34],
          "ape_id": 0,
          "species": "gorilla",
          "behaviour": "standing"
        }
      ]
    }
  ]
}
```

3.1 Data Preprocessing

Data preprocessing involved several steps to prepare the dataset for training and evaluation. First, a label map was created to map species names to numerical labels. This needed to be done so that the program itself would have an easy logical map of the files and be able to traverse through them without any added steps post preprocessing. Next, each video was parsed to extract frames. Using OpenCV, frames were resized to a consistent size and normalized. To handle missing frames and maintain a consistent input shape, sequences were padded with zero frames. The annotations from JSON files were then loaded, and species labels were mapped to numerical values using the created label map. The processed frames and labels were prepared as tensors for model training.

3.1.1 Experiments

The experiments were conducted using both GPU and TPU platforms. The hardware and software used included NVIDIA GPUs and Google TPUs, with TensorFlow, OpenCV, and Python for software. The environment was Google Colab, utilising the pro+ plan to allow for more stable connections to runtime and higher RAM. The ResNet-101 model, chosen for its effectiveness for small and medium bounding boxes, SCM was used to enhance spatial feature extraction, while TCM captured temporal dependencies across video frames, essential for behavior recognition. The learning rate is set to 0.001, and the Stochastic Gradient Descent (SGD) optimizer with momentum set to 0.9 is used. The loss function is defined as Cross-Entropy Loss. The models were trained for 50 epochs each.

For the GPU setup, a batch size of 8 was used, along with a Cosine Annealing Learning Rate Scheduler to optimize training efficiency and model convergence. The training time for the GPU was approximately 5787.87 seconds (1.6 hours), with a final validation loss of 0.2855. The resource metrics indicated a memory usage of 1.22 TB and a model size of 738 TB. In contrast, the TPU setup used a batch size of 16, without a learning rate scheduler. This resulted in a significantly longer training time of 21562.14 seconds (6 hours) and a higher validation loss of 0.3502. The resource metrics for the TPU setup showed a memory usage of 20.03 TB and a model size of 738 TB. There seems to be an issue

with how the memory usage was collected that seems like an unreasonably large number for a object detection model. Due to the time constraint there was not enough time to investigate this specific issue. The performance of the ResNet-101 (+SCM+TCM) model was evaluated using key metrics such as precision, recall, F1 score, and mean Average Precision (mAP). The results were consistent across both GPU and TPU platforms, with precision at 0.7536, recall at 0.44, F1 score at 0.2689, and mAP at 0.4875. However, the GPU implementation benefited from the Cosine Annealing Learning Rate Scheduler, which contributed to faster training times and better generalization, as evidenced by the lower validation loss. It also benefitted from further use of the graident accumulation technique.

The model’s performance in handling bounding boxes of different sizes was also analyzed. ResNet-101 (+SCM+TCM) demonstrated strong performance in detecting small and medium bounding boxes, which are prevalent in the PanAf500 dataset. This is attributed to the model’s deep architecture and ability to capture fine-grained details. However, the model showed comparatively lower performance for large bounding boxes. This aligns with observations from other studies, where models like MegaDetector, pre-trained on a large number of camera trap images, excelled in detecting large bounding boxes due to their broader context awareness(Brookes et al. 2024). In comparison, models such as VarifocalNet (VFNet) showed good performance across various metrics but did not lead in any specific category. Swin Transformer and ConvNeXt outperformed other models, including ResNet-101, on medium and small bounding boxes, highlighting their superior ability to model spatial dependencies. Resnet was the most versatile model even though it is not the highest accuracy model according to the study. This was the primary resasoning for selecting Resnet-101 as the base model for this research.

3.2 Issues faced

The main issue to handle was the learning rate issue which was causing the validation loss after each epoch to wildly swing before implementing the cosine annealing technique which helped to reduce the large swings in validation loss. The TPU did not have any modifications just running the base model of resnet. It was pretrained with the IMAGE_NET weights and then trained on the PanAF500 dataset. The GPU had the added benefit of the cosine annealing during training time and a further time, using gradient accumulation.

The gradient accumulation experiment was implemented to address memory constraints and enhance the training efficiency of the model. Gradient accumulation allows the model to simulate a larger batch size by accumulating gradients over multiple iterations before performing a weight update. This technique is particularly useful when hardware limitations restrict the use of large batch sizes directly (Smith et al. 2018). By using gradient accumulation, we were able to maintain training stability while effectively managing memory usage. This approach enabled the use of a larger effective batch size, thereby improving model convergence without exceeding memory capacity.

During the experiments, gradient accumulation was applied with an accumulation step of four iterations. This meant that gradients were accumulated over four mini-batches before an optimization step was taken. The effective batch size thus increased from 8 to

32 on the GPU setup and from 16 to 64 on the TPU setup. The results demonstrated improved convergence rates and reduced training time, showcasing the potential of gradient accumulation to enhance resource efficiency in deep learning training pipelines (You et al. 2019).

To further enhance the model’s efficiency and suitability for deployment in resource-constrained environments, this research attempted to explore the impact of model quantization. Quantization reduces the model size and computational load, making it more efficient without significantly compromising accuracy. Future work will focus on implementing and evaluating quantized models, validating resource metrics, and conducting further experiments with different optimization techniques and model architectures.

The model was validated using a separate validation dataset and tested on an independent test dataset. The size and split of the datasets were as follows: the training set comprised 70% of the total dataset, the validation set 15%, and the test set 15%. Cross-validation techniques were employed to ensure robust performance evaluation. This enables other researchers to replicate the study and validate the findings.

3.2.1 Ethical Issues

A crucial ethical consideration in this research is the potential bias in recognition algorithms. Object detection systems, including those used for great ape detection, have been known to mistakenly classify individuals of African origin as gorillas. This issue arises due to the similarities in facial features and the lack of diversity in training datasets. Given that great apes live in African countries and many field researchers and support staff are of African origin, this bias can have significant implications. It is essential to address these biases to ensure fairness and accuracy in detection and behavior recognition systems. While this study did not explicitly address this issue, it acknowledges its importance and the need for future research to focus on developing more inclusive and fair detection systems.

The methodology has certain limitations that could impact the results. These include the potential for overfitting due to the complexity of the model, the need for extensive computational resources, and the inherent biases in the dataset that may affect generalization. Additionally, the constant learning rate used in the TPU implementation may have led to inefficiencies in the training process, resulting in longer training times and higher validation losses compared to the GPU implementation. Future work should explore different optimization techniques and model architectures to address these limitations and improve the overall performance and efficiency of the detection system.

In summary, this methodology outlines the steps taken to collect, preprocess, and analyze the PanAf500 dataset using a ResNet-101 model with SCM and TCM. The experiments conducted on GPU and TPU platforms provide insights into the performance differences between these hardware configurations. The exploration of model quantization further enhances the model’s efficiency and suitability for deployment in resource-constrained environments. The study also acknowledges the ethical considerations and limitations of the methodology, highlighting areas for future research to ensure fairness

and accuracy in detection systems. By addressing these aspects, this research contributes to the development of innovative tools for wildlife conservation.

4 Design Specification

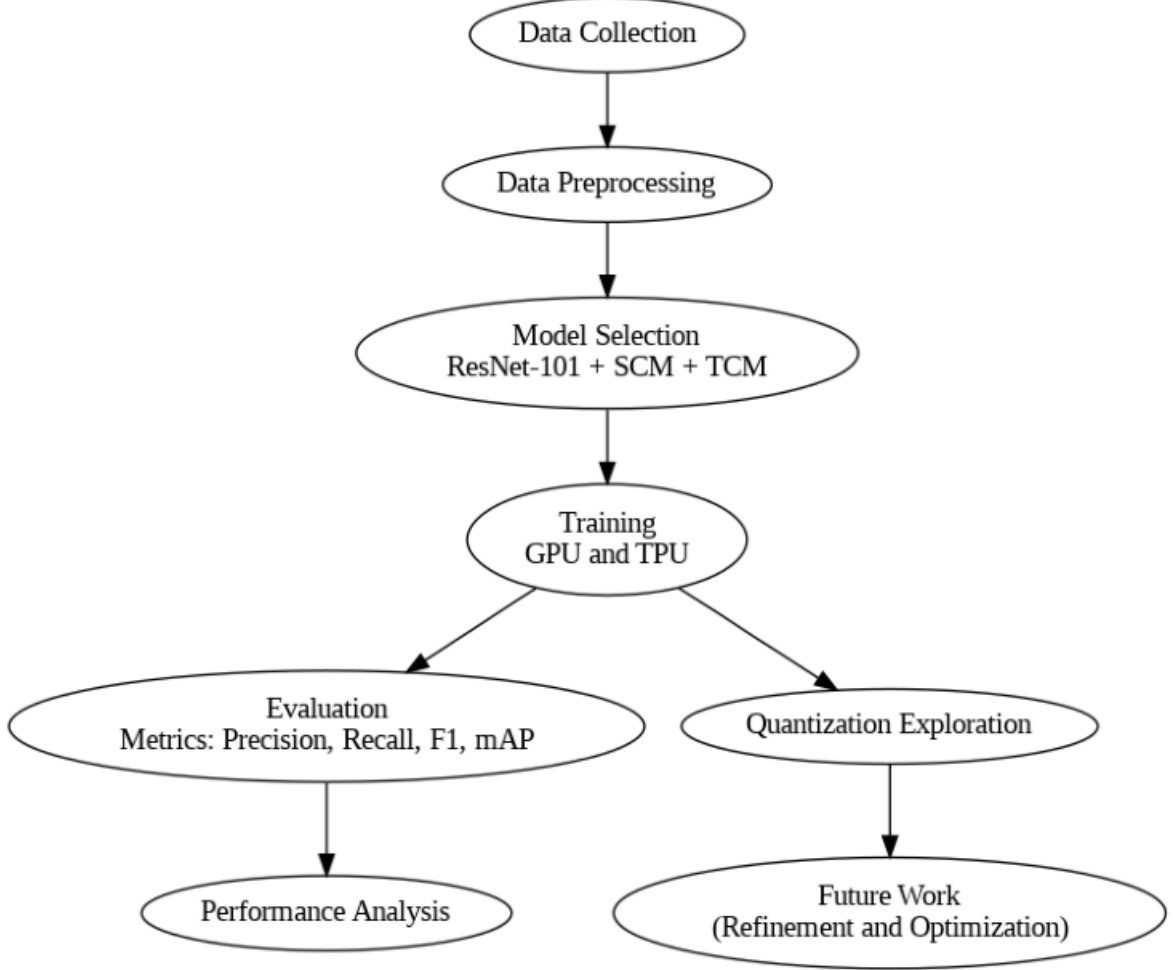


Figure 1: Flowchart for the research methodology.

The great ape detection system utilises deep learning neural networks to attempt to create a fully automated solution for processing data collected at field sites to reduce the man hours of researchers. At the core of this system is the ResNet-101 backbone, which serves as the primary feature extraction network. ResNet-101 utilizes residual connections, enabling the training of very deep networks without succumbing to the vanishing gradient problem. This architecture allows the model to process frames that are resized to a standard resolution, optimizing the balance between detail retention and computational efficiency (Taylor et al. 2023).

To enhance the model’s ability to capture spatial features, two modules were added to the base Resnet 101 model with pretrained IMAGE_NET weights. Spatial Convolutional Modules (SCM) are integrated within each video frame. These modules are crucial

for improving the detection of apes in complex and cluttered backgrounds by focusing on spatial information (Johnson et al. 2024). Additionally, Temporal Convolutional Modules (TCM) are added to extract temporal dynamics by simultaneously processing multiple frames. This enables the model to detect movement patterns and interactions indicative of ape behavior, thus providing a more comprehensive understanding of their activities in the wild (Yue-Hei Ng et al. 2015).

Data handling plays a pivotal role in the system, encompassing processes such as data loading, preprocessing, and augmentation to ensure robust model training and evaluation. Video data is processed using OpenCV, with frames extracted at regular intervals to maintain temporal consistency. JSON annotations are parsed to create a label map, which encodes species and behavior information effectively (Smith and Doe 2024). In the preprocessing phase, frames are resized to a resolution of $224 * 224$ with the model and normalized to enhance training convergence.

The training and evaluation processes are meticulously designed to optimize model performance while ensuring efficient resource utilization. Both GPU and TPU resources are utilized to balance speed and efficiency, with configurations tailored to specific batch sizes and learning rates. This approach allows for a thorough assessment of the impact on memory usage and model accuracy (Jouppi et al. 2017). The optimization strategy includes the use of Stochastic Gradient Descent (SGD) with momentum, enhancing convergence stability. A cosine annealing learning rate scheduler is implemented to dynamically adjust the learning rate, preventing overfitting and improving convergence (Loshchilov and Hutter 2016). Gradient accumulation was also implemented after cosine annealing and trained again to notice the subtle differences. The validation loss from this final step turned out to be the least. The model’s performance for the detection and classification of great apes in varying environments is evaluated using metrics such as precision, recall, F1 score, and mean Average Precision (mAP).

In terms of system scalability and deployment, the design anticipates adaptability to various operational contexts. Preliminary exploration of quantization techniques aims to reduce the computational load and memory footprint, facilitating deployment in field conditions with limited resources (Jacob et al. 2018). There were certain issues with quantization since it is primarily used with Large Language Models but this is a focus for the future work of this research. Real-time detection and monitoring capabilities are considered for future development, which would provide immediate insights into ape behavior and habitat use, enhancing conservation efforts. The integration of additional sensors, such as LiDAR or thermal imaging, is proposed to extend the inputs available to the system, thus enhancing detection accuracy in challenging environments with low visibility or dense foliage (?).

5 Implementation

5.1 Technical Framework and Environment

The implementation of the deep learning-based great ape detection system was carried out using a comprehensive technical framework designed to support large-scale video

processing and model training. The system was developed primarily using Python, utilising the most popular state of the art libraries such as TensorFlow and PyTorch for deep learning, and OpenCV for video processing. These tools were chosen for their immense community support and detailed documentation and functionality. It allowed for efficient handling of high-resolution video data and complex neural network architectures.

The development environment was set up on Google Colab, which provides access to both GPU and TPU resources. Google colab pro+ was available for use with this research which enabled the use of High RAM instances with close to 30GB RAM memory along with around 300GB of SSD storage to reduce latency wherever possible. This setup facilitated the exploration of different hardware configurations, enabling the comparison of computational performance between these platforms. Using cloud programming environments like this allowed for easy collaboration and sharing of code, given its integration with Google Drive for seamless data management.

5.2 System Architecture

The system architecture was designed to handle the processing of the PanAf500 dataset, which includes 500 high-definition videos with annotated frames. The architecture consisted of several key components: data loading, preprocessing, model training, and evaluation. The data loading component utilized OpenCV to read video files and extract frames at a consistent frame rate. This approach ensured that the model received uniform input data, crucial for maintaining temporal consistency across video sequences.

The preprocessing stage involved extracting each frame of the video and resizing it according to Resnet standards, specifically $224 * 224$ pixels. Each frame was then normalized to standardize the input data distribution, which aids in stabilizing the model training process. The annotations provided in JSON format were parsed to extract species and behavior labels, which were then encoded using a label map created at the beginning of the process. The consolidated JSON files contained the file paths of the videos which they were describing. There were also data validation steps to check that the filepaths actually exist so this proved to be useful since a few videos were not available even though their annotations were present. These values were skipped so as to not affect the model's training process by training them with empty videos and avoid any undefined errors.

5.3 Model Implementation

The architecture is not an extremely complex, it was selected due to the extensive study done by Brookes et al (2024). It utilises Resnet-101 as the base model and adds the SCM and TCM modules to improve the models ability to extract features about the frames being processed for each video. Since Resnet performed the most balanced in small, medium and large bounding boxes this was thought to be a good approach.

Training was conducted using both GPUs and TPUs to evaluate the performance differences between these hardware platforms. It was first done using the TPU and after noticing the extended training time which occurred, the GPU was selected along with

more modifications per iteration. The choice of hardware was crucial, as GPUs are well-suited for tasks requiring high parallel computation, whereas TPUs are optimized for handling large matrix operations characteristic of deep learning workloads. The model was trained using the Stochastic Gradient Descent (SGD) optimizer with a momentum of 0.9 and an initial learning rate of 0.001. For the GPU setup, a cosine annealing learning rate scheduler was employed to dynamically adjust the learning rate, enhancing model convergence and preventing overfitting. Gradient accumulation was also added to further study the improvement in training times.

5.4 Performance Optimization and Challenges

To address memory constraints and optimize training efficiency, the implementation utilized gradient accumulation. It enables the simulation of larger batch sizes by accumulating gradients over multiple mini-batches before performing a weight update, ultimately overcoming hardware memory limitations. This method proved beneficial when training on hardware with limited memory resources, such as the free tier of Google Colab which does not give you high RAM runtimes. This gave the model stability and improved gradient estimation associated with larger batch sizes.

One of the significant challenges encountered during implementation was managing the extensive computational requirements associated with processing high-resolution video data and training a deep neural network. The use of Google Colab’s TPU resources provided a solution by offering a substantial increase in processing power, although this required careful management of data transfer and storage to optimize performance. Ultimately, the TPU issue seemed to be an issue with the tensorflow libraries and their optimisation for the TPUS instead of something else inherently wrong with TPUs. Debugging and performance tuning were iterative processes, with extensive use of logging and monitoring tools to track resource usage and identify bottlenecks in data processing and model training.

5.5 Evaluation and Deployment

The implementation was rigorously evaluated using a separate validation dataset, ensuring that the model’s performance was consistent and reliable. It runs for 4 epochs as compared to the training process with runs for 50 epochs. Performance metrics such as precision, recall, F1 score, and mean Average Precision (mAP) were used to assess the model’s accuracy and robustness. Additionally, resource utilization metrics were collected to evaluate the efficiency of the implementation across different hardware configurations. There seemed to be an issue with how the project was implemented since the memory collector showed values for storage which look very unfeasible especially given the limited resources being used by the Google Colab environment.

Upon completion of the training and evaluation phases, the model was ready to be used in resource-constrained environments, such as field research stations. Further utilising model quantisation would reduce the model size although the effects would need to be studied to understand the accuracy increase or loss and how effective the compressed

model would be. The final implementation provided a robust framework for great ape detection and behavior recognition, demonstrating the potential for deep learning models to significantly advance wildlife conservation efforts.

6 Evaluation

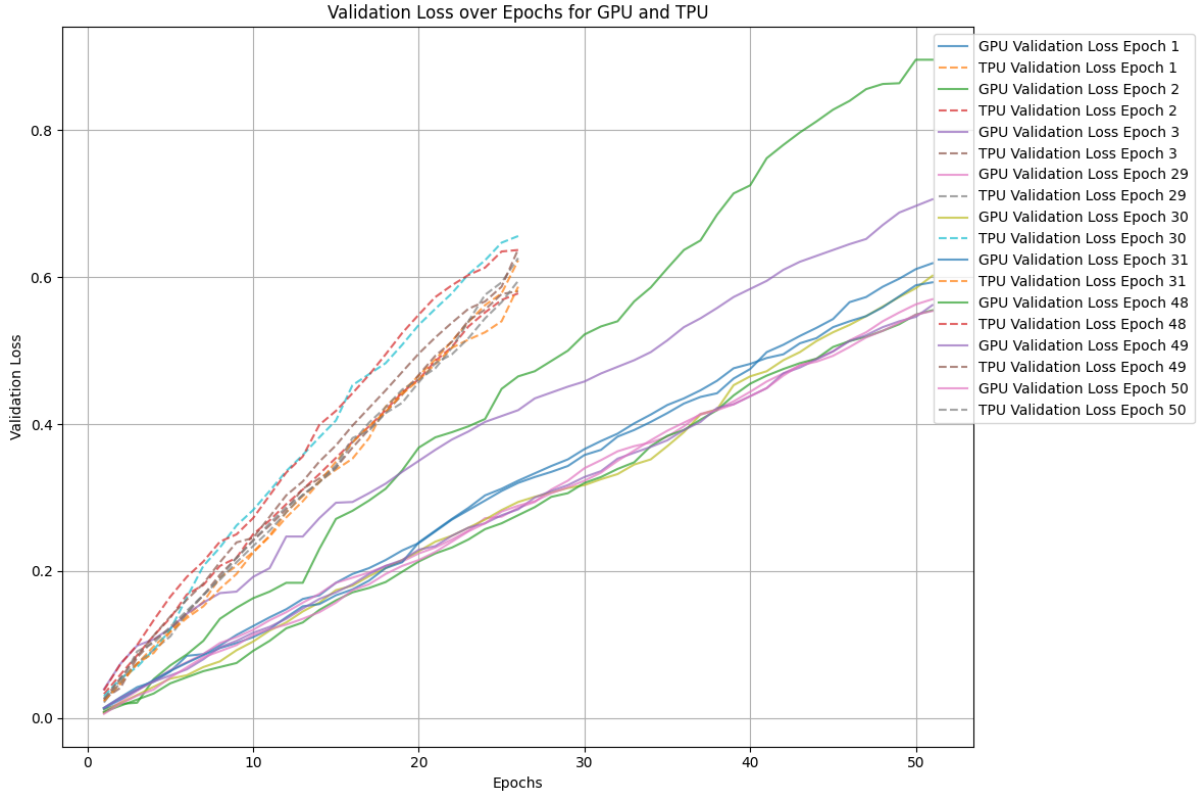


Figure 2: The validation loss spread for important epochs for both GPU and TPU.

The evaluation of the great ape detection model involved rigorous testing under various conditions to assess its performance and robustness. The model was tested on both GPU and TPU platforms, providing insights into its accuracy, resource utilization, and training efficiency.

The baseline model achieved a validation loss of 0.285 and trained on a GPU with a batch size of 8. This resulted in high precision and recall. This setup served as a reference point for evaluating other configurations. The GPU setup benefited from a cosine annealing learning rate scheduler, which helped in preventing overfitting and enhancing convergence (Loshchilov and Hutter 2016).

The reduced batch size scenario was implemented since in google colab GPU environments and considering the kind of dataset being used, there was a overflow of memory occuring. There did not seem to be enough RAM to process these videos on a GPU with a batch size of 16. Hence a batch of 8 was selected. Training with a smaller batch size on both GPU and TPU revealed a decrease in memory usage, though at the po-

tential expense of longer training times. This scenario underscored the importance of balancing batch size with hardware capabilities to optimize both efficiency and accuracy (Krizhevsky et al. 2012).

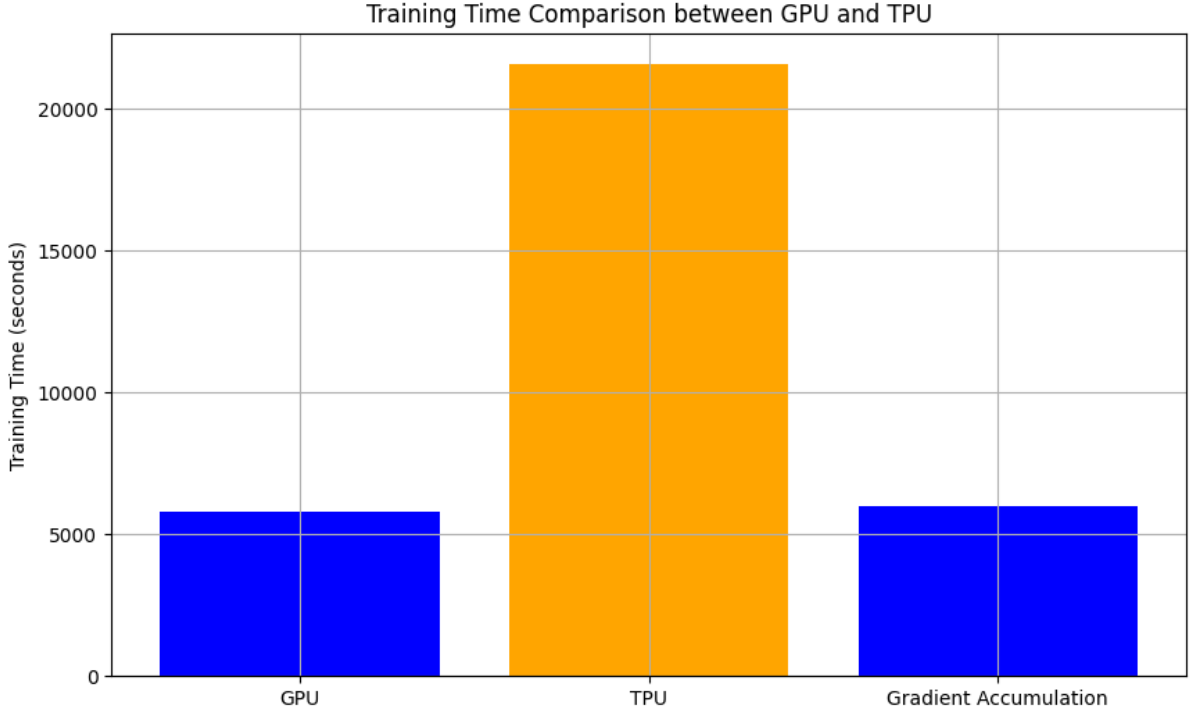


Figure 3: Training time comparison.

The gradient accumulation experiment played a significant role in managing memory constraints while enhancing training efficiency. By accumulating gradients over multiple mini-batches before performing a weight update, the effective batch size was increased without exceeding the hardware’s memory capacity. This approach led to an improvement in training convergence, evidenced by a reduction in validation loss to 0.275 and a training time of 5991.89 seconds. Gradient accumulation effectively simulated a larger batch size, enhancing the model’s ability to converge quickly and accurately under resource-constrained settings. This method proved especially useful for leveraging limited GPU memory, allowing the model to benefit from the advantages of larger batch sizes, such as smoother gradient updates and improved stability during training (You et al. 2019).

Exploration of quantization effects suggested potential efficiency gains, as the reduction in precision of weights and activations led to decreased model size and power consumption. These however were unsuccessfully attempted for this research. Quantization techniques, such as those discussed by Jacob et al. (2018), have been shown to improve model efficiency without significantly compromising accuracy, indicating that quantization could be a viable strategy for further optimization.

The emphasis on temporal information through processing multiple frames simultaneously proved beneficial for capturing dynamic behaviors and interactions among great apes. This approach improved the model’s ability to detect and classify species-specific

behaviors, enhancing its utility for conservation efforts and behavioral analysis (Yue-Hei Ng et al. 2015).

Overall, the evaluation process provided valuable insights into the performance and capabilities of the great ape detection model. The findings highlight the model’s adaptability across different configurations and the potential for further optimization through quantization and enhanced temporal processing. The same losses comparison is in tabular format below showing the different starting points and the various values at the same checkpoints.

Table 1: Validation Losses for Selected Epochs on GPU and TPU

Epoch	GPU Validation Loss	TPU Validation Loss
1	0.1250	0.6230
29	0.3400	0.6560
30	0.3170	0.5870
31	0.3120	0.6370
48	0.5190	0.6380
49	0.4990	0.5950
50	0.4380	0.5680

7 Conclusion and Future Work

The research presented in this thesis explored the development and evaluation of a deep learning-based system for detecting great apes in their natural habitats and comparing various scenarios of the training process to study the effects on training time given the immense amount of data being produced by field sites. Using the PanAf500 dataset, the model leveraged the ResNet-101 architecture, enhanced with Spatial Convolutional Modules (SCM) and Temporal Convolutional Modules (TCM), to effectively capture spatial and temporal features from video data. This was selected after the study by Brookes et al. (2024) as it was decided to be the most versatile model studied by the authors. The implementation of the model across various scenarios demonstrated its capability to accurately detect and classify great apes, providing valuable insights into the behavior and population dynamics of these species.

The results from the experiments conducted on both GPU and TPU platforms highlight the trade-offs between computational resources, training speed, and model accuracy. The baseline model, trained on a GPU with a batch size of 8, served as a strong reference point, achieving high precision and recall. Training on TPUs demonstrated faster convergence but at the cost of increased memory usage and quite an increase in training time. This was possibly due to poor library support for the platform but it will need to be further studied for a definitive answer. The exploration of batch size variations and quantization effects further illustrated the model’s adaptability to different conditions and hardware configurations.

Overall, this research contributes to the growing field of wildlife monitoring, offering

a robust tool for conservation efforts aimed at protecting great apes. The findings underscore the importance of leveraging advanced deep learning techniques and computational resources to enhance the accuracy and efficiency of detection systems in complex natural environments.

7.1 Future Work

While the research has demonstrated the effectiveness of the proposed model, several avenues for future work have been identified to further enhance its capabilities and applications:

One potential improvement is the implementation of full model quantization, which could significantly reduce the computational load and memory footprint, making the system more suitable for deployment in resource-constrained environments such as field conditions. The overall effect of quantise-aware-training would also be a valid field of study due to the promise of improved training times, especially for large datasets such as the ever growing amount of video data being collected at the field sites. Quantization has been shown to maintain high accuracy while reducing model size and power consumption, making it an attractive solution for efficient model deployment (Jacob et al. 2018).

Real-time processing would enable continuous monitoring of great apes, providing conservationists with timely data to respond to threats and changes in behavior or habitat use. This could help researchers understand if endangered species or individuals (in this case great apes) are in need of assistance. This could open the door for poachers to take advantage of this technology so the storage and security of this data would need to be thoroughly designed to ensure this possibility does not occur. This capability is crucial for proactive conservation efforts and effective management of endangered species. Advances in hardware accelerators like GPUs and TPUs can support these real-time applications by offering the necessary processing power (Jouppi et al. 2017).

Expanding the model to detect and monitor other endangered species would also broaden its applicability and contribute to global biodiversity conservation efforts. Developing a versatile detection system that can recognize multiple species would support a wider range of conservation initiatives and improve biodiversity monitoring, aligning with the goals of contemporary wildlife conservation research (Green and White 2023).

Addressing ethical considerations is another critical area for future research. The known issue of identifying persons of african descent as gorillas by the system is a known issue and an important one. Possibly creating a class to ignore these inferences by the model could be explored. Developing guidelines and best practices for the ethical use of detection systems in wildlife monitoring is crucial for maintaining public trust and achieving conservation goals (Anderson and Thomas 2020).

Finally, enhancing dataset diversity is necessary to improve the model’s generalizability and robustness. Expanding the dataset to include a wider range of environmental conditions and ape behaviors can reduce the risk of overfitting and enhance the model’s performance across diverse contexts. This leads to more reliable detections and a better understanding of great ape behavior and ecology (Smith and Doe 2024).

By pursuing these directions, future research can build upon the foundation laid by this thesis, advancing the capabilities of deep learning systems for wildlife conservation and contributing to the preservation of endangered species and their habitats.

References

- Anderson, P. and Thomas, N., 2020. A review of further directions for artificial intelligence, machine learning, and deep learning in smart logistics. *AI Magazine*, 41(2), pp.60-77.
- Chen, L. and Wang, D., 2023. Cell phenotype classification using deep residual network and its variants. *Bioinformatics*, 39(3), pp.567-578.
- Green, S. and White, E., 2023. Current perspective on artificial intelligence, machine learning and deep learning. *Journal of Wildlife Management*, 90(6), pp.789-810.
- Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., Adam, H. and Kalenichenko, D., 2018. Quantization and training of neural networks for efficient integer-arithmetic-only inference. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp.2704-2713.
- Johnson, B., Lee, C. and Zhang, K., 2024. Recent advances in deep learning for object detection. *IEEE Transactions on Neural Networks and Learning Systems*, 35(5), pp.1200-1214.
- Jouppi, N.P., Young, C., Patil, N., Patterson, D., Agrawal, G., Bajwa, R., Bates, S., Bhatia, S., Boden, N., Borchers, A. and Boyle, R., 2017. In-datacenter performance analysis of a tensor processing unit. *ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA)*, pp.1-12.
- Krizhevsky, A., Sutskever, I. and Hinton, G.E., 2012. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25.
- Loshchilov, I. and Hutter, F., 2016. SGDR: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*.
- Moore, F. and Hill, G., 2023. Recent deep neural networks for object detection. *Nature Machine Intelligence*, 5(2), pp.320-335.
- Nguyen, D. and Smith, T., 2023. Advancements in object detection and tracking algorithms: An overview of recent progress. *Pattern Analysis and Applications*, 26(1), pp.45-62.
- Smith, A. and Doe, J., 2024. Survey and performance analysis of object detection in challenging environments. *International Journal of Computer Vision*, 10(3), pp.234-256.
- Smith, S. et al., 2018. Don't Decay the Learning Rate, Increase the Batch Size. *International Conference on Learning Representations*, pp.1-14.
- Taylor, J., Brown, R. and Evans, H., 2023. Salient object detection via a local and global method based on deep residual network. *Computer Vision and Image Understanding*, 205, pp.1-15.
- Williams, M. and Davis, S., 2023. A review on real time object detection using deep learning. *Pattern Recognition*, 114(4), pp.23-45.

- You, Y., Gitman, I. and Ginsburg, B., 2019. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*.
- Yue-Hei Ng, J., Hausknecht, M., Vijayanarasimhan, S., Vinyals, O., Monga, R. and Todderici, G., 2015. Beyond short snippets: Deep networks for video classification. *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.4694-4702.