
Deep Learning Techniques for Long-Tailed Image Classification: A Comparative Study of Loss Designs and Model Architectures

MASTER'S THESIS IN
ELECTRICAL ENGINEERING

By

Christine Annelise Midtgaard

AU521655

STUDY PROGRAM: ELECTRICAL ENGINEERING

SUPERVISOR

Kim Bjerger

ASSOCIATE PROFESSOR

kbe@ece.au.dk



M.Sc. IN ELECTRICAL ENGINEERING
DEPARTMENT OF ENGINEERING SCIENCE & TECHNOLOGY
AARHUS UNIVERSITY

JANUARY 3, 2025

Abstract

Long-tailed datasets, where a few classes dominate with abundant samples while many classes have sparse representation, pose significant challenges for traditional training methods. These imbalances often lead to models that perform well on majority classes but struggle to recognize or generalize to minority (tail) classes.

Acknowledgements

Contents

Abstract	iv
Acknowledgements	v
1 Introduction	1
1.1 Problem Definition	1
1.1.1 Goals of this thesis	2
1.1.2 Hypothesis	2
1.1.3 Approach	2
1.1.4 Scope of this thesis	3
1.2 Related Work	3
1.3 Reading Guide	4
2 Background	5
2.1 Long-Tailed Datasets	5
2.2 Model Architectures	7
2.2.1 Introduction to Deep Neural Networks	7
2.2.2 Convolutional Neural Networks	8
2.2.3 Vision Transformers	11
2.2.4 Transfer Learning	14
2.3 Class Re-balancing Methods for Long-Tailed Learning	14
2.3.1 Data Modification	15
2.3.2 Loss Modification	15
3 Methodology	25
3.1 Overview of Approach	25
3.2 Dataset Preparation and Specifications	26
3.2.1 Benchmark Dataset Selection	26
3.2.2 Data Characteristics: Class Distribution	26
3.2.3 CIFAR-100-LT	28
3.2.4 Data Augmentation	30
3.3 Long-tailed Learning Techniques	31
3.3.1 Model Selection	31
3.3.2 Selection of Loss Function	33
3.3.3 Optimizer	34
3.4 Evaluation Metrics	34

3.4.1	Baselines	34
3.5	Model Training and Evaluation	35
4	Experimental Setup	36
5	Results and Discussion	38
5.1	Main Findings	38
5.2	Overall Results	39
5.2.1	MobileNetV2	39
5.2.2	ResNet50	40
5.2.3	ViT-B/16	41
5.2.4	ConvNeXt Base	42
5.3	Benchmarks	43
5.4	Performance Comparisons	44
5.5	Summary and Discussion	47
6	Conclusion and Future Work	50
6.1	Revisiting the Goals of the Thesis	50
6.2	Future Work	51
	Bibliography	52
A	Results	57
A.1	MobileNetV2	57
A.2	ResNet50V2	58
A.3	ViT-B/16	60
A.4	ConvNeXt Base	61
B	Experimental Setup Details	63
B.1	Dataset Specifications	63
B.2	Data Preprocessing	64
B.3	Model Architecture Settings	64
B.4	Training Configurations	65
B.5	Evaluation Metrics	66
B.6	Hardware and Software Configurations	66
B.6.1	Hardware	66
B.6.2	Software	66
B.7	Reproducibility Considerations	66
B.8	Implementation Faults	67

Chapter 1

Introduction

TODO: Emphasize head and tail classes.

Image classification is one of the main challenges computer vision.

Deep learning has become a prominent solution in recent years for tackling image recognition tasks. With the availability of large datasets, i.e. ImageNet, along with GPUs, training of deep learning models have become easier, and have led to remarkable results **TODO: reference here**. The trained models have shown image classification with accuracies of over 80 % **TODO: reference here**, and hence the interest in deep learning for image recognition is high. However, the high accuracies are for on models trained on balanced datasets with thousands of samples per class **TODO: reference here**, while most real-world datasets are skewed in samples per class **TODO: reference here**.

This thesis focuses on the problem with long-tailed datasets. The problem with training a deep learning model on long-tailed datasets is that the model will effectively the data from the classes with most samples, and not the classes with few samples. The finished model will then not recognize an input from the tail classes. Most real-world datasets follows a long-tailed structure, hence the need for a reliable method to detect examples of tail-class data. The aim of this thesis is to try out some of the methods tackling the long-tailed problem for deep learning described in the paper *Deep Long-Tailed Learning: A Survey* by Zhang et al.[1] to find a method for long-tailed learning that works on a specific long-tailed dataset of images of moths taken around equator. The goal of the moth dataset is to identify species.

1.1 Problem Definition

TODO: Something about head and tail classes.

The goal of this project is to investigate deep learning models, like Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs), and methods, mainly class re-balancing through cost-sensitive learning, and analyze their performance on datasets representaing a long-tailed structure.

Four models are investigated, where three of them are CNNs, and one is a ViT. The CNNs investigated are the MobileNetV2, ResNet50v2, and ConvNeXt Base,

while the ViT investigated is the ViT-B/16. All four models are pretrained on ImageNet, and further trained on both a balanced version of CIFAR100 and an long-tailed version of CIFAR100, called CIFAR100-LT, and introducing the long-tailed technique as class re-balancing in the loss functions by using re-weighting. All models are trained with Softmax Cross-Entropy Loss, Weighted Softmax Cross-Entropy Loss, Focal Loss, Class-Balanced Loss, Balanced Softmax Loss, LDAM Loss, and Equalization Loss. The main purpose of this master's thesis is to compare class re-balancing methods on a well representing type of models used in deep learning image recognition tasks. Further comparison with state-of-the-art methods for deep long-tailed learning will be made. The motivation for this comparison is to find a way to train on a real-world long-tailed dataset of moths to correctly identify species.

1.1.1 Goals of this thesis

The goals of this thesis are to:

1. Investigate the efficacy of long-tailed learning methods by assessing their performance on tail classes without sacrificing accuracy on head classes.
2. Understand how model design affects the performance of deep long-tailed learning methods.
3. Provide a comprehensive insight in comparisons that inform the choice of methods for long-tailed distributions.

1.1.2 Hypothesis

Deep long-tailed learning methods, such as the carefully designed loss functions tailored for long-tailed distributed datasets, can improve the performance of underrepresented (tail) classes while maintaining the overall accuracy across diverse model architectures. The effectiveness of these loss functions is influenced by the choice of model architecture and the degree of imbalance of the dataset.

1.1.3 Approach

TODO: finish this section. This master's thesis consists of six steps, described below:

First step is to investigate the dataset used for training, testing and validation in *Deep Long-Tailed Learning: A Survey*, as the class re-balancing in this paper is used as inspiration for this thesis.

Second step is generating a long-tailed version of CIFAR100 that can be used for training and comparisons of methods.

Third step is the implementation of models and methods, combining them.

Fourth step is hyperparameter settings.

Fifth step is training and evaluation of the models with different loss functions. The training is both on the balanced and long-tailed version of CIFAR100.

Sixth step is a comparison of the methods.

Other steps could involve comparisons to related work and other long-tailed learning methods.

1.1.4 Scope of this thesis

This thesis focuses on applying deep long-tailed learning methods to image classification tasks with an emphasis on loss re-weighting as a solution to handle imbalanced datasets. The loss re-weighting methods include cross-entropy loss, focal loss, weighted cross-entropy loss, class-balanced loss, balanced softmax loss, equalization loss, and LDAM loss. The experiments are conducted on the CIFAR100 dataset, with both a balanced and synthetically generated long-tailed version. The evaluation is done with particular focus on top-1 accuracy, with attention paid to performance on head, middle, and tail classes. This thesis does not explore other long-tailed learning approaches such as re-sampling, information augmentation, or model architecture modifications. **TODO: not yet module improvement.**

1.2 Related Work

The challenge of long-tailed datasets has been extensively studied in the literature [1, 2].

TODO: Add references to re-weighting papers.

Data Re-sampling Data re-sampling techniques aim to modify the training dataset to mitigate class imbalance by simulating a balanced dataset. The most popular methods for re-sampling are random over-sampling (ROS) and random under-sampling (RUS) [3, 4]. ROS involves duplicating random samples from the minority classes, whereas RUS involves removing samples from the majority classes [1, 4]. However, these techniques have their weaknesses: over-sampling the minority classes will eventually lead to overfitting tail classes, and under-sampling can lead to underperformance on head classes [1]. Another re-sampling technique, called Synthetic Minority Over-Sampling Technique (SMOTE) [3], involves synthetic generation of instances of the minority classes based on the existing data.

Transfer Learning Transfer learning techniques address the issue of class imbalance by transferring learned features from head classes to tail classes. Examples

include transferring the intra-class variance [5] and transferring semantic deep features [6].

Data Augmentation Data augmentation is a deep learning technique used to expand the training dataset by applying transformations to samples or features, enhancing model accuracy while reducing overfitting [7, 8]. Common augmentation methods include Mixup [9], which combines pairs of samples to create interpolated examples; CutMix [10], which replaces regions of one image with patches from another; and UniMix [11], which focuses on calibrating the feature space for long-tailed distributions. Rare-class Sample Generator (RSG) [12] focuses on enhancing tail-class performance by transferring knowledge from head classes.

Decoupled Training Decoupled training [13] separates representation learning and classifier training. First, the model is trained to extract features from the entire dataset, afterwards, a re-balancing technique is applied during classification. Despite the simple framework, decoupling showed state-of-the-art performance on long-tailed benchmarks.

Ensemble Learning Ensemble learning is a technique that combines multiple network modules (i.e. experts) to solve long-tailed problems. By doing so, the predictions of the different models can be combined, ultimately outperforming a single model [14, 15].

TODO: Margin Loss

Margin loss. The hinge loss is often used to obtain a “max-margin” classifier, most notably in SVMs [Suykens and Vandewalle, 1999]. Recently, Large-Margin Softmax [Liu et al., 2016], Angular Softmax [Liu et al., 2017a], and Additive Margin Softmax [Wang et al., 2018a] have been proposed to minimize intra-class variation in predictions and enlarge the inter-class margin by incorporating the idea of angular margin. In contrast to the class-independent margins in these papers, our approach encourages bigger margins for minority classes. Uneven margins for imbalanced datasets are also proposed and studied in [Li et al., 2002] and the recent work [Khan et al., 2019, Li et al., 2019]. Our theory put this idea on a more theoretical footing by providing a concrete formula for the desired margins of the classes alongside good empirical progress. [16]

1.3 Reading Guide

TODO: Mention what each chapter will cover and how they relate to each other.

Chapter 2

Background

This chapter presents the different background topics of the thesis work, which are the long-tailed datasets, model architectures *Convolutional Neural Networks (CNN)* and *Vision Transformers (VT)*, the deep long-tailed learning methods *Class Re-balancing (CR)*, *Information Augmentation (IA)*, and *Module Improvement (MI)*. These topics will be explained for the reader.

TODO: Mention image classification, as it is the primary goal of this thesis.

2.1 Long-Tailed Datasets

Long-tailed datasets pose significant challenges in deep learning, as they represent an extreme form of class imbalance. Addressing these challenges is central to this thesis, which explores methods to improve model performance on underrepresented classes. This section outlines the structure of long-tailed distributions and their implications.

A balanced dataset is one where all classes are evenly represented, whereas imbalanced datasets feature varying sample sizes across classes. Long-tailed datasets are characterized by a significant class imbalance, where a few dominant classes account for most samples (head classes), while the majority of classes are underrepresented (tail classes) as depicted in Figure 2.1. This distribution is common for real-world datasets [17, 6]. For example, the iNaturalist, a popular benchmark for image classification, exhibits a long-tailed distribution of species [18]. Other benchmarks are constructed by sampling from datasets such as ImageNet [19] by using a Pareto distribution, which simulates long-tailed class distributions with a power-law decay [1, 20, 16].

CIFAR100-LT [16], derived from the CIFAR-100 dataset [21], serves as the primary dataset for the experiments conducted in this thesis. CIFAR-100 is a widely used benchmark in classification research due to its diverse class representation and manageable size. It consists of 60,000 32x32 color images divided into 100 classes, each with 600 samples. These are further split into 500 training images and 100 testing images per class. CIFAR100-LT is created by reducing the number of samples in certain classes of CIFAR-100 following an exponential decay in sample sizes, given by:

$$n_i = n_{max} \cdot \text{IR}^{\frac{i-1}{C-1}} \quad (2.1)$$

Where n_i is the number of samples in class i , n_{max} is the number of samples in the most frequent class, IR is the imbalance ratio, and C is the total number of classes [22].

Other long-tailed datasets follow a Pareto distribution with number of samplers per class as followed [6]:

$$f(C) = \frac{\alpha C_m^\alpha}{C^{\alpha+1}}, \quad C \geq C_m, \quad \alpha > 0 \quad (2.2)$$

Here, α is the shape parameter, C_m is the scale parameter representing the minimum possible class index, and C is the class index. A larger α results in a more severe imbalance.

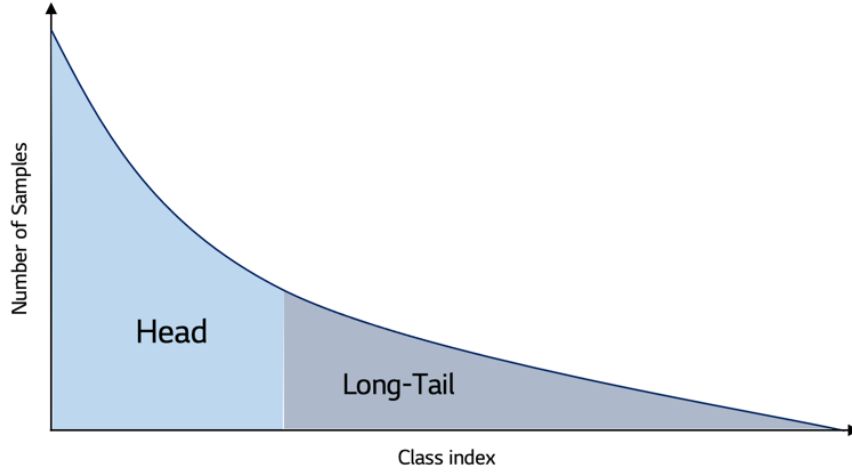


Figure 2.1: Illustration of a long-tailed distribution. Figure from [23].

Class imbalance has a profound impact on model performance compared to evenly distributed datasets [24, 25]. Deep networks trained on long-tailed datasets often exhibit biased performance, favoring head classes while performing poorly on tail classes [1]. Zhang et al. (2023) provide a comprehensive survey of methods addressing this challenge, categorizing current approaches into three main groups: class re-balancing, information augmentation, and module improvement. These methods will be further explored in section 2.3.

TODO: Write something of the following: "Let $\{x_i, y_i\}_{i=1}^n$ be the long-tailed training set, where each sample x_i has a corresponding class label y_i . The total number of training samples over K classes is $n = \sum_{k=1}^K n_k$, where n_k denotes the number of samples in class k . Let π denote the vector of label frequencies, where $\pi_k = \frac{n_k}{n}$ indicates the label frequency of class k . Without loss of generality, a common assumption in long-tailed learning [31, 32] is that the classes are sorted by cardinality in decreasing order (i.e., if $i_1 < i_2$, then $n_{i_1} \geq n_{i_2}$, and $n_1 \gg n_K$). The imbalance ratio is then defined as $\frac{n_1}{n_K}$." [1]

2.2 Model Architectures

Deep learning has revolutionized image classification by introducing models capable of learning complex patterns and representations from data. Among these, Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs) are chosen as the primary architectures used in this thesis due to their performance on image classification tasks. This section provides a theoretical foundation for these models, focusing on the specific architectures utilized: MobileNetV2 [26], ResNet50V2 [he2016], and ConvNeXt Base [27] as the CNN architectures, and ViT-B/16 [28] as the ViT architecture.

TODO: A brief summary of the relevance of the chosen architectures so that the reader knows why they are included.

2.2.1 Introduction to Deep Neural Networks

Before the introduction of Convolutional Neural Networks (CNNs) and, more recently, Vision Transformers (ViTs), the standard approach for image classification involved flattening a two-dimensional image matrix into a one-dimensional array and passing it through a Multilayer Perceptron (MLP), also known as a feed-forward neural network. MLPs are fully connected neural networks composed of an input layer, output layer, and one or more hidden layers. Being fully connected means that each neuron in a given layer is connected to all neurons in the next layer, forming a dense network. These connections are associated with weights and biases, which the network learns during training. Input features are fed into the input layer, propagated through hidden layers that add complexity to model nonlinear relationships, and yield predictions in the output layer. Known as universal approximators, MLPs can approximate any continuous function given sufficient neurons in the hidden layers [29, 30].

To illustrate the structure of a neural network, figure 2.2 shows an example of a feed-forward neural network with three input neurons, two hidden layers, each with four neurons, and two output neurons. This architecture could be used, for instance, to classify images of cats and dogs based on three input features, such as height, weight, and width of the animals. The input propagates through the network, with each neuron computing a weighted sum of its inputs followed by an optional nonlinearity. The final output is a prediction, where the class corresponding to the neuron with the highest value is chosen.

However, this simple neural network becomes insufficient for more complex problems, such as image classification, as it requires an increasing number of parameters. For instance, a 224×224 RGB image flattened is very large, making MLPs parameter-heavy and inefficient. The limitations of MLPs were addressed by Convolutional Neural Networks (CNNs), which introduced convolutional and pooling layers to effectively preserve and utilize the spatial information of pixels in two-dimensional images [29].

TODO: Include how the weights are updated with gradients based on the loss function.

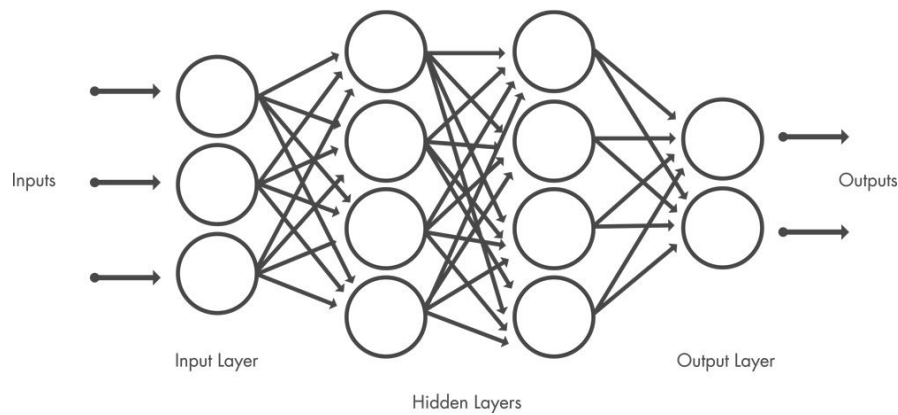


Figure 2.2: Layers of a neural network. Figure from [31]. **TODO: Make this figure.**

2.2.2 Convolutional Neural Networks

Convolutional Neural Networks (CNNs) [32] were introduced to address the limitations of MLPs for image-related tasks. Unlike MLPs, which treat input features as independent, CNNs are designed to recognize patterns in images by applying local filters through convolutional layers, and thereby preserving the two dimensional input of an image, preserving the idea that nearby pixels are related [33, 34, 29].

CNNs consist of several core components, as illustrated in Figure 2.3.

TODO: write this section They have three main layers: convolutional layers, pooling layers, and a fully connected layer. The convolutional layer is the first layer, and serves to extract local features by applying filters to small regions of an image, while pooling layers reduce the spatial dimensions of feature maps, providing invariance to small translations. The CNN can be made up of multiple convolution and pooling layers, but the fully connected layer will be the final layer. Activation functions introduce nonlinearity, allowing the network to capture complex patterns. At the final stage, fully connected or global pooling layers aggregate the extracted features into predictions, enabling tasks such as classification or segmentation.

TODO: What is a convolution? TODO: Inductive bias TODO: Translation equivariance TODO: Stem cells, kernel, stride, convolution.

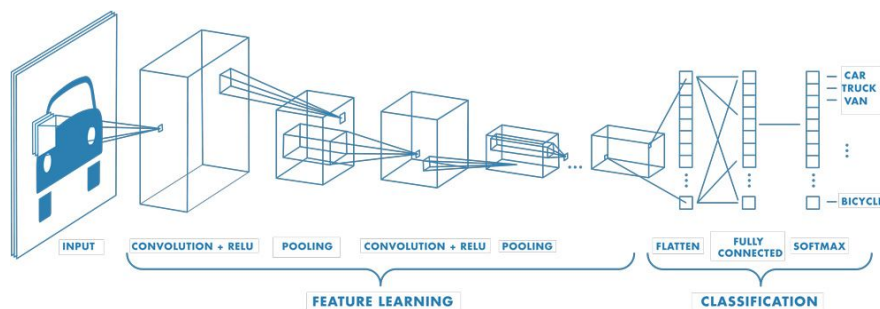


Figure 2.3: Illustration of a convolutional neural network. Figure from [31]. **TODO: Make this figure.**

TODO: Tie to the models used for experiments in this thesis. <https://www.ibm.com/topics/convolutional-neural-networks>

CNNs gained popularity after the introduction of LeNet-5 by LeCun et al. in 1998 [33] that demonstrated the potential of CNNs by recognizing handwritten digits. Later, AlexNet [34] achieved a breakthrough by winning the ImageNet Challenge in 2012, demonstrating the potential of CNNs to handle large-scale image recognition tasks by deepening the architectures and utilizing multiple GPUs for training. Subsequently, the evolution of CNNs has progressed through architectures such as VGGNet [35], GoogLeNet [36], and ResNet [37], which have set the stage for the advancements seen in modern CNNs.

ResNet50 Architecture

ResNet50 is a variant of the Residual Network (ResNet) architecture, developed by Microsoft Research (He et al.) in 2015 [37]. The ResNet architecture was designed to address the vanishing and exploding gradient problem in deep networks. When training a deep neural network, as the number of layers increases, the gradients of the loss function with respect to the weight can potentially vanish or explode during backpropagation [37]. This can lead to slow convergence or unstable updates.

In traditional CNNs, stacked layers learn a direct mapping $\mathcal{H}(x)$ of the input x to the output. The introduction of residual layers, however, allows for the layer to fit a residual mapping, as illustrated in figure 2.4, where the network learns the residual function $\mathcal{F}(x) = \mathcal{H}(x) - x$. Isolating $\mathcal{H}(x)$ yields $\mathcal{H}(x) = \mathcal{F}(x) + x$, meaning that the input x is passed through a skip connection. This architecture reduces the complexity of the optimization process, as the network only has to model the residual component $\mathcal{F}(x)$ rather than the full mapping $\mathcal{H}(x)$ [37].

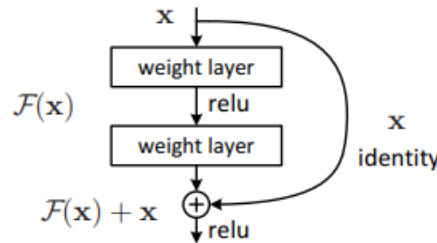


Figure 2. Residual learning: a building block.

Figure 2.4: Residual Learning. Figure from [37]. **TODO: Make this figure.**

The residual network is formed by stacking multiple layers of residual blocks, as the one depicted in figure 2.4. The ResNet architecture can have a varying number of layers, and, as the name implies, the ResNet50 variant has 50 layers. Other architectures include ResNet-34, ResNet-101, and ResNet-152 [38].

The ResNet50 implementation consists of 1 convolutional layer, 4 stages of bottleneck residual blocks with layers [3, 4, 6, 3], global average pooling, and finally a fully connected layer for classification[39], as seen in figure 2.5.

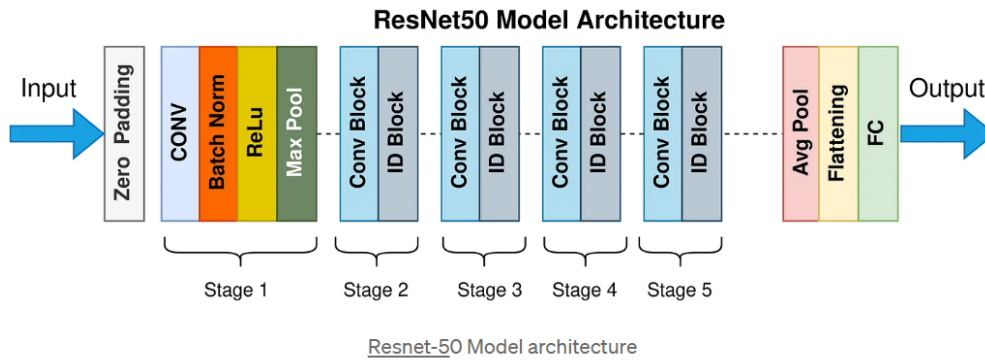


Figure 2.5: ResNet50 Architecture. Figure from <https://towardsdatascience.com/the-annotated-resnet-50-a6c536034758>. **TODO: Make this figure.**

MobileNetV2 Architecture

MobileNetV2 introduced by Sandler et al. [26] is a lightweight CNN model designed primarily to balance model accuracy and computational efficiency, making it suitable for mobile or embedded devices. Building upon the original concepts of MobileNetV1 [40], MobileNetV2 preserves the use of depthwise separable convolutions, a method for reducing the parameters and floating-point operations, while introducing a novel element known as the inverted residual structure.

Inverted Residual Blocks While traditional residual connections described above allow for identity mapping and improved gradient flow, MobileNetV2 employs an inverted residual structure [26]. Instead of mapping from a high-dimensional representation down to a lower-dimensional bottleneck, then reconstructing features at the output, inverted residual blocks begin with a low-dimensional input and expand it to a higher-dimensional space before applying a depthwise convolution. After spatial filtering, the representation is projected back down to a low-dimensional space. This approach, illustrated in Figure 2.6, helps preserve crucial information and maintain a rich feature space without substantially increasing computational cost. The use of a linear bottleneck (i.e., no nonlinear activation in the low-dimensional projection) also helps prevent the destruction of useful information, further improving efficiency and accuracy.

Depthwise Separable Convolution Following MobileNetV1, MobileNetV2 relies on depthwise separable convolutions to factorize the convolution operation into two simpler operations [40]: depthwise convolution and pointwise convolution as shown in figure 2.7. The depthwise convolution applies a single filter to each input channel, and the pointwise convolution (a 1×1 convolution) then recombines the channels to produce the desired output. In comparison, a standard convolution both filters and combines inputs into a new set of outputs. This approach reduces the parameter count and computational load, making the model suitable for devices with limited resources [40, 26].

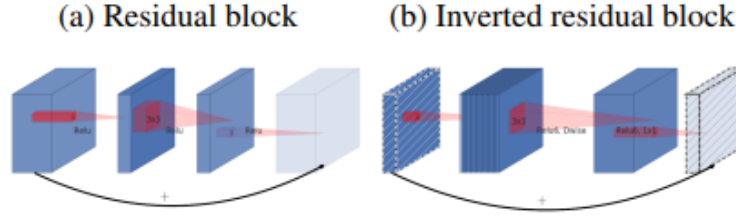


Figure 3: The difference between residual block [8, 30] and inverted residual. Diagonally hatched layers do not use non-linearities. We use thickness of each block to indicate its relative number of channels. Note how classical residuals connects the layers with high number of channels, whereas the inverted residuals connect the bottlenecks. Best viewed in color.

Figure 2.6: Residual and Inverted Residual Blocks. Figure from [26]. **TODO: Make this figure.**

2.2.3 Vision Transformers

TODO: Explain their advantages over CNNs for certain tasks. Mention why they are relevant for handling long-tailed datasets.

Transformers were introduced by Vaswani et al. in 2017 [41], and revolutionized the deep learning field, surpassing Recurrent Neural Networks (RNNs) for Natural Language Processing (NLP) tasks [42, 41]. The design behind transformers is based on self-attention mechanism, allowing the model to weigh the significance of each part of input. As their design lack recurrence or convolutions, transformers use positional embeddings to represent the order of tokens in a sequence [41]. Later, the Vision Transformer (ViT), an adaptation of the transformer architecture for image processing, was introduced by Dosovitskiy et al. (2021) [28]. Here, images are represented as sequences including the class label as a learnable token for classification. The input image is divided into a sequence of patches, which are then flattened and linearly embedded into a vector. The spacial information is preserved by adding positional encodings to the embeddings. Next, the sequence is fed into a transformer encoder identical to that introduced by Vaswani et al., consisting of alternating layers of Multi-head self-attention (MSP) and Multi-Layer Perceptron (MLP) blocks with Layer Norm (LN) applied before every block, and residual connection after every block. The final MLP layer acts as the classification head [28]. The illustration in figure 2.8 shows the architecture of Vision Transformers.

Unlike CNNs, ViTs have much less image-specific inductive bias [28]. While the inductive bias in CNNs help the model generalize to unseen data [43], the ViT architecture does not include built-in assumptions about the structure of the data,

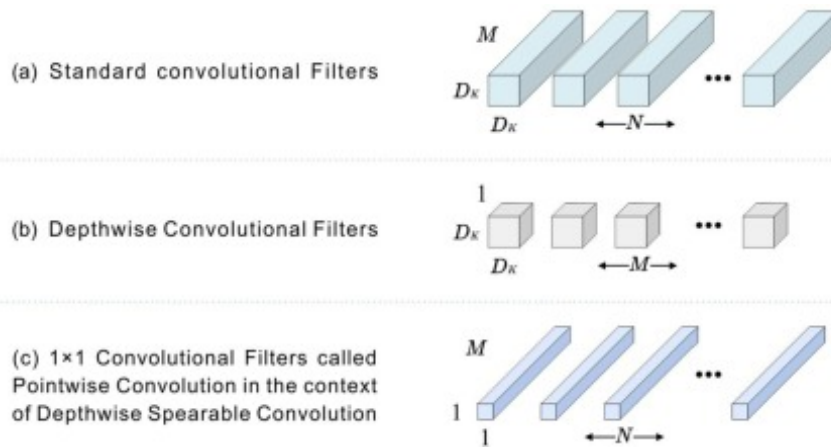


Figure 2.7: Depthwise Separable Convolution. Figure from <https://www.sciencedirect.com/topics/computer-science/depthwise-separable-convolution>. **TODO: Make this figure.**

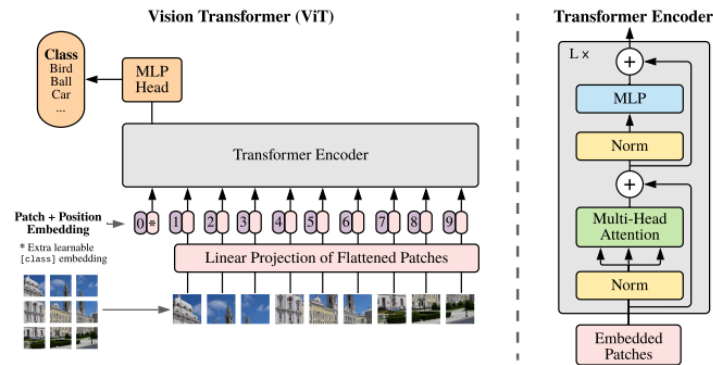


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

Figure 2.8: Vision Transformer architecture. Figure from [28]. **TODO: Make this figure.**

only the MLP layers are local and translationally equivalent [28]. As a result, ViTs require more training data to learn spacial relations compared with CNNs.

TODO: What made the ViT get SOTA? Why were they better than CNNs? Scalability. <https://www.youtube.com/watch?v=QqejVOLNDHA>.

ViT-B/16 Architecture

the ViT-B/16 is a Vision Transformer that leverages the transformer architecture for image classification [28], pre-trained on **TODO: ImageNet-1K or ImageNet-21K?**. One of the main characteristics of the ViT-B/16 is that the input image consists of a fixed 16×16 patch size. It has 12 stacked encoder layers, each computing self-attention with 12 heads [28]. Additionally, each layer includes an MLP, and skip-connections are applied around self-attention and MLP blocks

for better gradient flow. LN is applied before the attention and MLP blocks. Afterwards, the appended class token is extracted, and a fully connected layer maps the token to the class prediction [44].

ConvNeXt Base Architecture

ConvNext, introduced by Liu et al. in 2022 [45], evolutionized the traditional CNN architecture by incorporating elements from ViTs. Starting with the stem cell, the ConvNeXts employ a patchify stem, as seen in figure 2.9. By setting the kernel to 4×4 and the stride to 4, the result is a non-overlapping convolution where no information is shared between them. However, these are later combined in the final layers of the network.

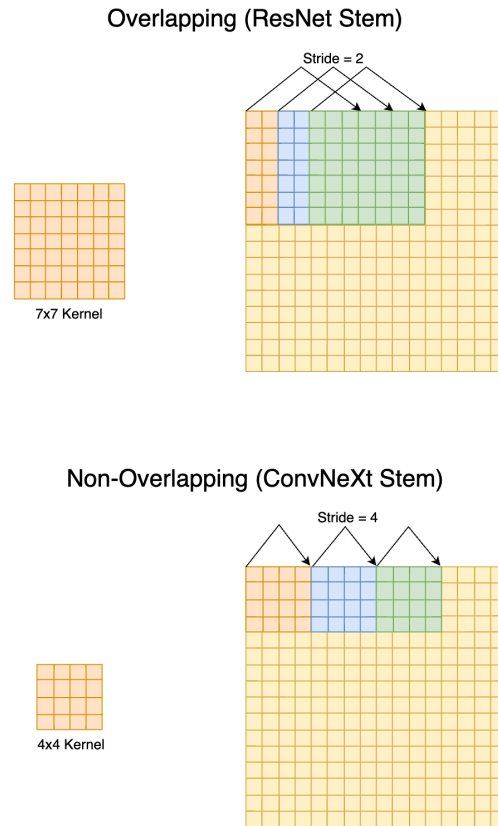


Figure 2.9: ResNet stem vs ConvNeXt stem. Figure from <https://www.kungfu.ai/blog-post/convnext-a-transformer-inspired-cnn-architecture> **TODO: Make this figure.**

Next, inspired by ResNeXt [46], the ConvNeXt use depthwise convolution, already described in the MobileNetV2 section 2.2.2, to increase the FLOPs/accuracy trade-off. Likewise, the ConvNext incorporates inverted bottleneck block, which are an important design in transformers as well as MobileNetV2 [45]. Additionally, the ConvNext architecture uses a *7time7* depthwise convolution in each block on the condition of moving the depthwise convolution layer up. See figure 2.10.

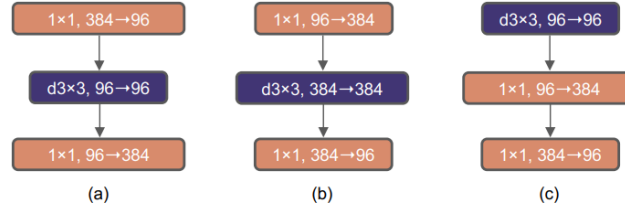


Figure 3. **Block modifications and resulted specifications.** (a) is a ResNeXt block; in (b) we create an inverted bottleneck block and in (c) the position of the spatial depthwise conv layer is moved up.

Figure 2.10: Block modifications. Figure from [45] **TODO: Make this figure.**

Instead of using the commonly used activation function in ResNets, the Rectified Linear Unit (ReLU), the ConvNeXt use Gaussian Error Linear Unit (GeLU). This activation function has shown to perform better than ReLU for transformers [45]. Additionally, ConvNext draws inspiration from transformers by removing most instances of normalization layers, only leaving them after depthwise convolutions and further replacing the Batch Normalization with Layer Normalization. Finally, ConvNeXt introduces separate downsampling layers between blocks rather than implementing downsampling within the blocks themselves.

2.2.4 Transfer Learning

In this thesis, transfer learning is used to

Transfer learning involves training a model on data from a source domain $TS = P(y|X)$ and then transferring it to a target domain $TT = P(y|X)$ – typically with less data available.

2.3 Class Re-balancing Methods for Long-Tailed Learning

Following the paper *Deep Long-Tailed Learning: A Survey* [1], the existing deep long-tailed learning methods are grouped into three main categories based on their technical approach: class re-balancing, information augmentation, and module improvement. These categories are further divided onto sub-categories: re-sampling, class-sensitive learning, logit adjustment, transfer learning, data augmentation, representation learning, classifier desing, decoupled training, and ensemble learning as shown on figure 2.11. This thesis will focus on class re-balancing long-tailed learning. The backgrounds of the methods used in this thesis are described in this section.

TODO: Introduce the cost/loss

TODO: consider if it makes sense to mention all the methods in the survey paper, or just mention the class re-balancing. TODO: Remove the figure of long-tailed categories. Potentially move it to related work.

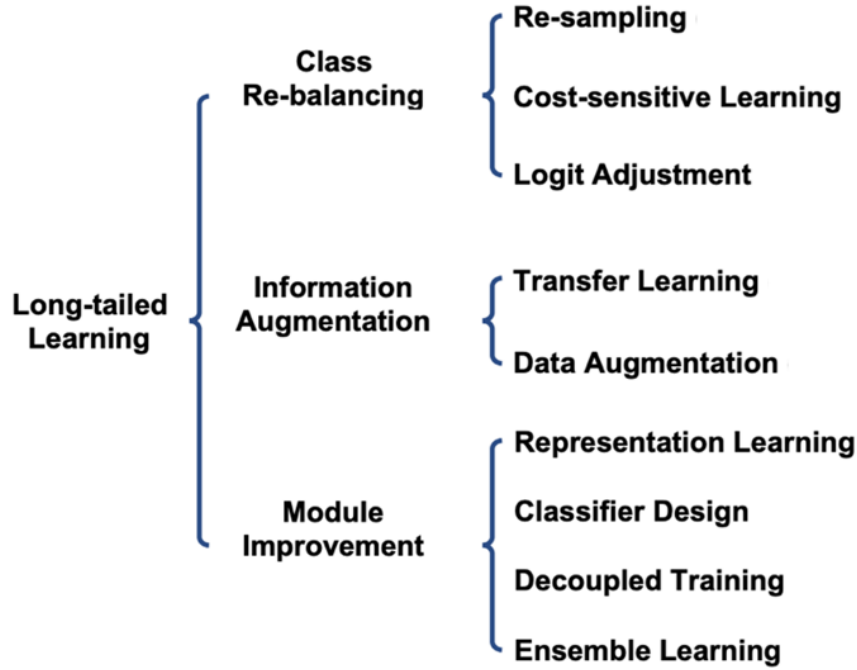


Figure 2.11: Long-tailed categories as described by *Zhang et al.* [1].

The class re-balancing method aims to re-balance the effect of the imbalanced training dataset, and has three main sub-categories: re-sampling, class-sensitive learning, and logit adjustment, according to Zhang et al. [1].

TODO: Include little about re-sampling, and focus on class-sensitive learning and an little on logit adjustment.

2.3.1 Data Modification

The traditional way to sample when training deep networks is using mini-batch gradient descent with random sampling. This means that each sample has an equal probability of being sampled. When sampling from an imbalanced dataset, samples from head classes naturally occur more often, and thus have higher chance of being sampled than samples from tail classes, making the resulting deep models biased towards head classes. Re-sampling is a method that addresses this problem by adjusting the number of samples per class in each sample batch for model training.

TODO: Rewrite this section to closer resemble the related work.

2.3.2 Loss Modification

TODO: make the notation consistent.

Traditional training methods using the standard loss function, cross-entropy loss, can lead the model to be biased towards head classes, as the loss ignores class imbalance and thus generate an uneven amount of gradients for different classes

[1], as described later in this section. Consequently, tail classes are often misclassified. To address this imbalance, class-sensitive learning methods modify the loss function **TODO: or how losses are calculated** to pay more attention to minority classes, thus improving overall performance.

The three categories of class re-balancing strategies explored in this thesis are *re-weighting*, *re-margining*, and *logit adjustment*. Re-weighting involves adjusting the contribution of each class to the loss by multiplying them with different weights. By multiplying weights to the loss function, the model allocates a higher penalty to misclassification of underrepresented classes, thus re-balancing the training process [1]. Re-margining, on the other hand, modifies the decision boundaries by introducing class-dependent margins, ensuring that minority classes have more room to establish discriminative features in the embedding space [1, 16].

In what follows, the theory behind several representative loss functions commonly used in class-sensitive learning is presented, beginning with an introduction to loss functions followed by the baseline Softmax Cross-Entropy loss. Subsequently, the theory behind multiple re-weighting schemes including Weighted Softmax Cross-Entropy, Focal Loss, Class-Balanced Loss, Balanced Softmax Loss, and Equalization Loss, as well as the re-margining method LDAM Loss. Each of these methods seeks to improve the imbalance in training through modifications that ultimately improve the classification performance on tail classes.

Table 2.1 summarizes the loss functions employed for class-sensitive learning used in this thesis, outlining their formulations and categorizations into re-weighting or re-margining strategies.

Table 2.1: Summary of losses.

Loss Name	Formulation	Type
Softmax loss [47]	$\mathcal{L}_{ce} = -\log(p_y)$	-
Focal loss [48]	$\mathcal{L}_{fl} = -(1 - p_y)^\gamma \log(p_y)$	re-weighting
Weighted Softmax loss [1]	$\mathcal{L}_{wce} = -\frac{1}{\pi_y} \log(p_y)$	re-weighting
Class-balanced loss [25]	$\mathcal{L}_{cb} = -\frac{1-\gamma}{1-\gamma^{n_y}} \log(p_y)$	re-weighting
Balanced Softmax loss [49]	$\mathcal{L}_{bs} = -\log\left(\frac{\pi_y \exp(z_y)}{\sum_j \pi_j \exp(z_j)}\right)$	logit adjustment
Equalization loss [50]	$\mathcal{L}_{eq} = -\log\left(\frac{\exp(z_y)}{\sum_j \omega_j \exp(z_j)}\right)$	re-margining
LDAM loss [16]	$\mathcal{L}_{ldam} = -\log\left(\frac{\exp(z_y - \Delta_y)}{\sum_j \exp(z_j - \Delta_j)}\right)$	re-margining

Introduction to Loss Functions

Loss functions are a fundamental component in machine learning and deep learning [29, 51]. They serve as a measure of how far a model’s predictions are from the actual values, where higher loss indicates greater prediction error, whereas a lower loss implies more accurate predictions. Thus, the goal of training a model is to minimize the loss.

The loss function plays a crucial role, as it guides the model to adjust its param-

eters (weights and biases) through an optimization algorithm [51]. The algorithm calculates the gradient of the loss function with respect to each parameter to determine how the parameter affect the loss. Afterwards, the model update its parameter to reduce the loss, following:

$$w_{t+1} = w_t - \eta \frac{\partial L}{\partial w_t} \quad (2.3)$$

Here, w are the parameters, η is the learning rate, and L is the loss function. This is an iterative process, and continues until the loss reaches a minimum or the model converges.

From equation (2.3), the choice of loss function is critical as it influences the model's interpretation of prediction errors and parameter adjustment. For image classification tasks, the standard loss function is the cross-entropy loss combined with the softmax activation function [29] described in the following.

TODO: This is the SGD. Include the Adam Optimizer. The Adam optimizer introduced by Kingma et al. [52] Adam is a replacement optimization algorithm for stochastic gradient descent for training deep learning models. Adam combines the best properties of the AdaGrad and RMSProp algorithms to provide an optimization algorithm that can handle sparse gradients on noisy problems. Adam is relatively easy to configure where the default configuration parameters do well on most problems.

Softmax Cross-Entropy Loss The *Softmax Cross-Entropy (CE) loss* is a fundamental building block in training deep classifiers and is widely regarded as the baseline in classification tasks [1, 53, 47]. The softmax loss, as it is commonly known as, is the combination of the cross-entropy loss and the softmax activation function [29, 51]. First, the need of the softmax activation function will be decribed.

The *Softmax* function transforms the raw output scores (logits) of the final layer of a neural network into a probability distribution over K classes. For an input $\mathbf{z} = [z_1, z_2, \dots, z_K]$, the Softmax function for class i is defined as:

$$P(y = i \mid \mathbf{z}) = \frac{\exp(z_i)}{\sum_{j=1}^K \exp(z_j)} \quad (2.4)$$

Here, $\exp(z_i)$ ensures that all values are positive, and dividing by the sum normalizes the probabilities so that they sum to 1. This normalization is crucial for classification, as it allows the network's outputs to represent the likelihood of each class.

The *Cross-Entropy loss* measures the difference between the predicted probability distribution \mathbf{P} (produced by Softmax) and the true distribution \mathbf{y} (the one-hot encoded ground truth). It is defined as:

$$\mathcal{L}_{\text{CE}} = - \sum_{i=1}^K y_i \log(P(y = i \mid \mathbf{z})) \quad (2.5)$$

For a single example where the true class is c , this simplifies to:

$$\mathcal{L}_{\text{CE}} = -\log(p_y) \quad (2.6)$$

This equation penalizes incorrect predictions by heavily weighting the log of the predicted probability for the true class. The loss is minimized when the predicted probability $P(y = c \mid \mathbf{z})$ approaches 1, indicating high confidence in the correct class.

While this approach provides a robust and stable training objective, it treats all classes equally and does not account for class imbalance.

Bias The standard Softmax Loss (eq. (2.6)) estimates class probabilities based on the assumption that the training and validation labels follow the same distribution, $p(y = j)$, in other words how common is the label j in the data [49]. When training with an imbalanced dataset the estimated probabilities ϕ can become biased when applied to data with a different class distribution, e.g. a balanced test set, which can lead to errors in predictions.

Re-Weighting

Weighted Softmax Cross-Entropy Loss The *Weighted Softmax Cross-Entropy (WCE) loss* modifies the standard CE loss (equation (2.6)) by introducing a weighting factor w . For multiclass classification, the CE loss it is multiplied by the inverse class frequencies $1/\pi_y$ [1, 48], given by:

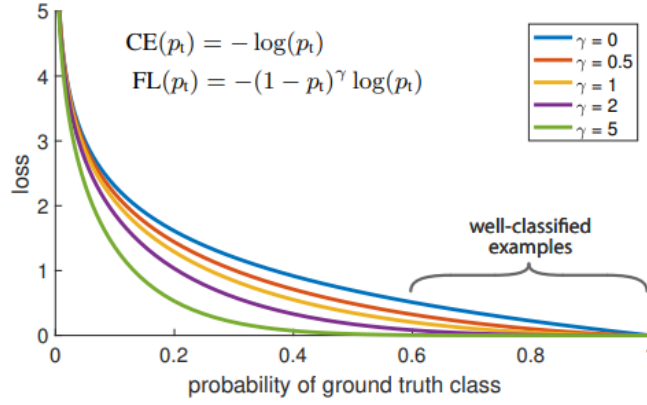
$$\mathcal{L}_{\text{WCE}} = -w \log(p_y) = -\frac{1}{\pi_y} \log(p_y) \quad (2.7)$$

This weighted CE loss ensures that minority classes contribute more to the overall loss, addressing the imbalance during training.

Focal Loss *Focal Loss* [48] addresses class imbalance by improving model performance on hard-to-classify examples by focusing on wrongly classified examples. The technique for this is called down-weighting, and uses the prediction probabilities to dynamically adjust the contribution of each sample to the loss. Well-classified examples with high probabilities p_y are down-weighted by adding a modulating factor $(1 - p_y)^\gamma$ to the cross-entropy loss (eq. (2.6)). Here, $\gamma \geq 0$ is a tunable focusing parameter. The Focal Loss modifies the CE loss by applying the inverse prediction probability as follows:

$$L_{\text{fl}} = -(1 - p_y)^\gamma \log(p_y) \quad (2.8)$$

This factor increases the weight of misclassified examples, ensuring the model prioritizes learning from challenging samples. Figure 2.12 shows how the focal loss for different values of γ . When $\gamma = 0$ the modulating factor equals 1, and the focal loss becomes the standard CE loss. The blue line in figure 2.12 represents the standard cross-entropy loss for $\gamma = 0$.

Figure 2.12: Figure from [48]. TODO: Make this figure

Class-Balanced Loss *Class-Balanced (CB) Loss* [25] introduces a re-weighting strategy based on the effective number of samples per class to re-balance the loss. Instead of simply using raw class frequencies, CB Loss estimates how much additional information new samples provide, acknowledging an information overlap among data, and as the number of samples increases, the marginal benefit of extracted features from new data diminishes. As illustrated in figure 2.13, the feature space is the map of all possible data, and each sample occupies a part of the space. Collecting more samples of a class means that there is a probability that their features overlap, meaning that additional samples diminishes new information. In feature space, the probability of newly sampled data with volume 1 is overlapping p , which means that the probability of adding new information to the feature space is $(1 - p)$.

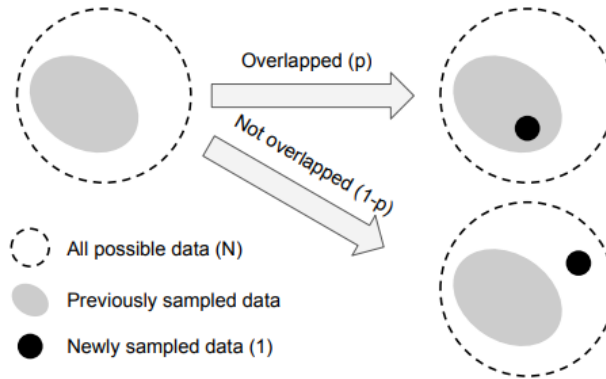


Figure 2. Giving the set of all possible data with volume N and the set of previously sampled data, a new sample with volume 1 has the probability of p being overlapped with previous data and the probability of $1 - p$ not being overlapped.

Figure 2.13: Figure from [25]. TODO: Make this figure

Cui et al. [25] introduced the effective number of samples as the expected volume of samples, denoted E_n , where $n \in \mathbb{Z}_{>0}$ is the number of samples. The effective

number of samples is defined as:

$$E_n = \frac{1 - \beta^n}{1 - \beta} \quad (2.9)$$

where $\beta = (N - 1)/N$. The hyperparameter $\beta \in [0, 1)$ determines how fast E_n grows.

As a result, the CB loss introduces a class-balanced re-weighting term, inversely proportional to the effective number of classes. Adding the CB term to the standard softmax cross-entropy given in equation (2.6) yield the following:

$$L_{cb} = -\frac{1 - \beta}{1 - \beta^{n_y}} \log(p_y) \quad (2.10)$$

where n_y is the number of samples in the ground truth class y . Applying the inverse of the effective number ensures that minority classes contribute more to the total loss, as the effective number in equation (2.9) grows large for minority classes (small n_y) and small for majority classes (large n_y). When $\beta = 0$ the loss is equivalent to the CE loss. Contrary, when $\beta \rightarrow 1$ the re-weighting effect grows.

Equalization Loss *Equalization (EQ) Loss* [50] aims to mitigate the over-suppression of tail classes, which occurs when these underrepresented classes serve predominantly as negative examples for the more frequent classes. The idea is to ignore the gradients for rare classes so that they are not excessively penalized when they appear as negatives, preventing their learned representations from being overshadowed. An analysis of gradient impact on classes is seen in figure 2.14.

The EQ loss was designed for object recognition, but the authors decided to adopt the principles into image classification, officially called the *Softmax Equalization (SEQ) Loss*, defined as:

$$\mathcal{L}_{SEQL} = -\sum_{j=1}^C y_j \log(\tilde{p}_j) \quad (2.11)$$

where

$$\tilde{p}_j = \frac{e^{z_j}}{\sum_{k=1}^C \tilde{w}_k e^{z_k}} \quad (2.12)$$

and the weight term is given by:

$$\tilde{w}_k = (1 - \beta T_\lambda(f_k))(1 - y_k) \quad (2.13)$$

Here, y_j is the ground truth class, z_j is the logit for class j , f_k is the frequency of class k , $T_\lambda(\cdot)$ is the threshold on class frequency. As there is no background category for image classification tasks, β is introduced to randomly maintain the gradient of negative samplers. β is a random variable with probability γ of taking value 1 and probability $1 - \gamma$ of taking value 0.

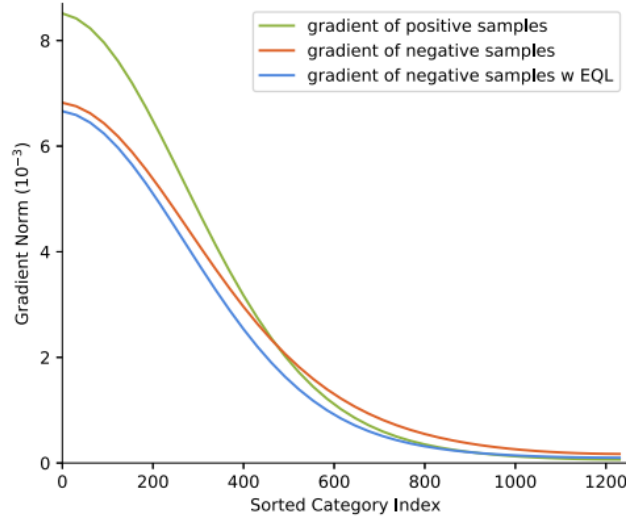


Figure 1: The overall gradient analysis on positive and negative samples. We collect the average L_2 norm of gradient of weights in the last classifier layer. Categories' indices are sorted by their instance counts. Note that for one category, proposals of all the other categories and the background are negative samples for it.

Figure 2.14: Gradient analysis. [50].

Logit Adjustment

Logit adjustment is a class re-balancing technique that aims to optimize the class imbalance by adjusting the model outputs, i.e. prediction logits, typically based on prior class probabilities [54, 49]. **TODO: Write about the label distribution shifts between training and test data.**

The output model outputs are computed by the following equation:

$$\eta_j = \theta_j^T f(x) \quad (2.14)$$

where θ_j are the weights for the last layer for class j , and $f(x)$ is the feature extractor function [49].

Balanced Softmax Loss The *Balanced Softmax (BS)* Loss modifies the standard softmax loss (eq. (2.6)) by directly incorporating class priors π_y into the logits before computing probabilities [49]. Unlike approaches that operate solely in the loss space, BS Loss integrates class frequency adjustments at the probability computation stage, effectively neutralizing the bias introduced by imbalanced class distributions.

In *Balanced Meta-Softmax for Long-Tailed Visual Recognition*, the authors represent Softmax regression as a multinomial distribution, ϕ , dependent on η , defined as:

$$\phi_j = \frac{e^{\eta_j}}{\sum_{i=1}^k e^{\eta_i}} \quad (2.15)$$

where k are the number of classes.

The Balanced Softmax uses the output logits η (eq. (2.14)) to parametrize two separate probabilities: ϕ for testing, and $\hat{\phi}$ for training. This separation is used when training with an imbalanced dataset to eliminate the discrepancy between training and testing distributions [49].

Assume ϕ to be the desired conditional probability of the balanced dataset, with the form $\phi_j = p(y = j | x) = \frac{p(x|y=j)}{p(x)} \cdot \frac{1}{k}$, and $\hat{\phi}$ to be the desired conditional probability of the imbalanced training set, with the form $\hat{\phi}_j = \hat{p}(y = j | x) = \frac{p(x|y=j)}{\hat{p}(x)} \cdot \frac{n_j}{\sum_{i=1}^k n_i}$. If ϕ is expressed by the standard Softmax function of the model output η , then $\hat{\phi}$ can be expressed as:

$$\hat{\phi}_j = \frac{n_j e^{\eta_j}}{\sum_{i=1}^k n_i e^{\eta_i}}. \quad (2.16)$$

where n_j is the number of samples in class j .

The authors proved that this method can accomodate the label distribution shifts between training and test sets [49]. ultimately, the Balanced Softmax loss function is becomes:

$$\mathcal{L}_{bs} = -\log \left(\frac{\pi_y \exp(z_y)}{\sum_j \pi_j \exp(z_j)} \right) \quad (2.17)$$

where z_y represents the logit for the true class, π_y is the class prior (e.g. normalized frequency of samples from class y), and the term in the denominator normalizes the probabilities across all classes while accounting for priors.

Margin Modification (Re-Margining)

Re-margining aims to solve class imbalance by improving classification by introducing margin between classes, encouraging higher margins between rare classes. The margin γ_i of the i -th class is the minimum distance of the data in the i -th class to the decision boundary. Figure 2.15 illustrates how the margin affects the decision boundary for binary classification. As the decision boundary separates the regions of input space corresponding to different class predictions, by increasing the margin, the model becomes more confident in classification.

The margin is defined as:

$$\gamma(x, y) = f(x)_y - \max_{j \neq y} f(x)_j \quad (2.18)$$

which is the distance between the prediction of the sample (x, y) and the decision boundary.

The training margin for class j is defined as:

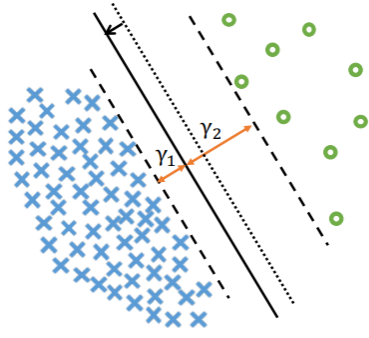


Figure 1: For binary classification with a linearly separable classifier, the margin γ_i of the i -th class is defined to be the minimum distance of the data in the i -th class to the decision boundary. We show that the test error with the uniform label distribution is bounded by a quantity that scales in $\frac{1}{\gamma_1 \sqrt{n_1}} + \frac{1}{\gamma_2 \sqrt{n_2}}$. As illustrated here, fixing the direction of the decision boundary leads to a fixed $\gamma_1 + \gamma_2$, but the trade-off between γ_1, γ_2 can be optimized by shifting the decision boundary. As derived in Section 3.1, the optimal trade-off is $\gamma_i \propto n_i^{-1/4}$ where n_i is the sample size of the i -th class.

Figure 2.15: Margins [16]. **TODO: make this figure.**

$$\gamma_j = f(x)_y - \max_{i \in S_j} \gamma(x_i, y_i) \quad (2.19)$$

The classical generalization bound is defined as

$$\text{imbalanced test error} \lesssim \frac{1}{\gamma_{\min}} \sqrt{\frac{C(\mathcal{F})}{n}} \quad (2.20)$$

In a balanced setting:

$$\frac{1}{2} \text{err}[\text{class1}] + \frac{1}{2} \text{err}[\text{class2}] \lesssim \left(\frac{1}{\gamma_1 \sqrt{n_1}} + \frac{1}{\gamma_2 \sqrt{n_2}} \right) \sqrt{C(\mathcal{F})} \quad (2.21)$$

This boundary can be shifted by taking the minimum: $\min \gamma_1 \sqrt{n_1} + \frac{1}{\gamma_2 \sqrt{n_2}}$, such that $\gamma_1 + \gamma_2 = C$. This yields the optimal solution:

$$\gamma_i \propto n_i^{-1/4} \quad (2.22)$$

with the standard solution being $\gamma_1 = \gamma_2$ [16].

LDAM Loss The *Label-Distribution-Aware Margin (LDAM)* Loss [16] represents a class-sensitive approach based on re-margining. Instead of altering the loss directly, LDAM modifies the margin applied to each class' decision boundary, ensuring that the model is more confident in classification for minority classes. The goal is to apply the optimal boundary in equation (2.22) by increasing the distance between the largest and second largest logit for each class and introducing that to the cross-entropy loss (eq. (2.6)).

The class-dependent margin for multiple classes is defined as:

$$\gamma_j = \frac{C}{n_j^{1/4}} \quad (2.23)$$

where C is a tunable hyper-parameter, and n_j is the number of samples in class j . Applying the margin through the loss function, inspired by the Hinge Loss **TODO: reference**, will encourage the model to respect the margins in equation (2.23):

$$\mathcal{L}_{ldam-hg} = \max \left(\max_{j \neq y} \{z_j\} - z_y + \Delta_y, 0 \right) \quad (2.24)$$

where

$$\Delta_j = \frac{C}{n_j^{1/4}} \quad \text{for } j \in \{1, \dots, k\}. \quad (2.25)$$

and z_y is the model output (logit) for the true class y , and $\max_{j \neq y} \{z_j\}$ is the highest score among incorrect classes. Adding the margin Δ_j ensures that the true class logit is separated from the highest incorrect score by the at least Δ_j . However, since the function in equation (2.24) is non-differentiable at certain points, a modifications to the cross-entropy loss is used to address this. Here, the class-dependent margins are introduced directly into the logits as follows:

$$L_{ldam} = -\log \left(\frac{\exp(z_y - \Delta_y)}{\sum_j \exp(z_j - \Delta_j)} \right) \quad (2.26)$$

Chapter 3

Methodology

This chapter describes the methods and approaches used in the experiments. This includes the selection of dataset, models, loss functions, and evaluation strategies.

3.1 Overview of Approach

Image classification involves predicting the correct category of an image from a set of predefined classes. This project focuses on addressing the challenge of class imbalance, where some classes have many samples while others have very few.

The CIFAR-100 dataset [21] was chosen because it is small enough to work with efficiently and is commonly used in other research. To create a long-tailed version of the dataset, its class distribution was adjusted to match the distribution of the ImageNet-LT dataset used in [1]. This adjustment ensures that the experiments simulate real-world class imbalance.

Several loss functions designed to handle class imbalance were tested in this project, including Softmax Cross-Entropy (CE), Focal Loss (FL), Class-Balanced (CB) Loss, Balanced Softmax (BS) Loss, Equalization (EQ) Loss, and LDAM Loss. These loss functions modify how the models learn, making it easier for them to focus on classes with fewer samples.

The experiments were performed using four model architectures: ResNet50, MobileNetV2, ConvNeXt Base, and ViT-B/16. All models were pretrained on ImageNet, which provided a strong starting point for fine-tuning on the CIFAR-100 dataset. These models were chosen because they are well-suited for image classification tasks and represent different design approaches.

Each model was trained on both a balanced training set and a long-tailed training set. The performance of the models was evaluated using a balanced test set, a long-tailed test set, and the long-tailed test set divided into head, middle, and tail classes. This division made it possible to analyze the performance for classes with different numbers of samples. The main evaluation metric used was top-1 accuracy, which measures how often the model’s top prediction is correct.

In what follows, each step of the methodology will be explained in detail.

3.2 Dataset Preparation and Specifications

A description of this section here.

Following the dataset structure used in *Deep Long-Tailed Learning: A Survey*, the CIFAR100 dataset was modified to create a long-tailed training set and a balanced test set.

3.2.1 Benchmark Dataset Selection

There are a number of benchmark datasets for long-tailed image classification tasks, including ImageNet-LT, Places365-Lt, CIFAR-100-LT, and iNaturalist 2018 **TODO: references**. The previous three are sampled from ImageNet, Places365, and CIFAR-100, following the Pareto distribution, as described in section 2.1, while the iNaturalist 2018 is a natural long-tailed dataset. Table 3.1 summarizes the four datasets and their data specifications.

Table 3.1: Summary of long-tailed benchmarks for image classification.

Dataset	# Classes	# Training Data	# Test Data
ImageNet-LT	1,000	115,846	50,000
CIFAR100-LT	100	50,000	10,000
Places-LT	365	62,500	36,500
iNaturalist 2018	8,142	437,513	24,426

With its 100 classes and only 60,000 samples, the CIFAR-100 dataset was chosen as a benchmark for the experiments conducted in this thesis based on its manageable size, compared to the other benchmarks, and widespread use as a standard in image classification research **TODO: reference**. It’s comparatively small size allows for rapid experimentation and comparisons, whereas a larger benchmark requires more computational effort and time for experiments, thus it is well-suited for multiple configurations and evaluations. However, CIFAR-100 has some limitations compared to larger dataset, like ImageNet-LT or iNaturalist. For example, CIFAR-100 images are of relatively low resolution (32×32 pixels) compared to ImageNet (224×224 pixels), and may limit the ability of models to capture fine-grained details. Likewise, the smaller number of classes and samples can result in less diversity compared to datasets like iNaturalist 2018, which contains nearly nine times as many samples in the training data as CIFAR-100-LT. Despite these drawbacks, CIFAR-100-LT remains a practical choice for this thesis due to its computational efficiency.

TODO: Describe the type of images in CIFAR-100 and include examples.

3.2.2 Data Characteristics: Class Distribution

Understanding the class distribution in long-tailed datasets is essential for evaluating the impact of the class imbalance on model performance. This section analyzes the class distribution of the benchmark ImageNet-LT used in [1], which

serves as a reference for creating a similar distribution across the CIFAR-100-LT training, evaluation, and test datasets.

The ImageNet-LT dataset, introduced by Liu et al. [6], has been widely used in empirical studies, including *Deep Long-Tailed Learning: A Survey* [1]. The class distributions of its training, validation, and test sets are shown in Figures 3.1, 3.2, and 3.3, respectively. The ImageNet-LT is constructed from the original ImageNet-2012 following the Pareto-distribution (Eq. (2.2)) with power value $\alpha = 6$ [6], resulting in an imbalance ratio of 256, meaning that the most frequent class has 256 times more samples than the least frequent class. As suggested in figure 3.1, the most frequent class has 1280 images, and the least frequent class 5 images [6]. The power value of 6 is relatively large, resulting in the steep imbalance and very few samples in the tail, as seen in figure 3.1. The underrepresentation in the tail classes poses a significant challenge for models to learn from tail classes. Both the validation and test sets are balanced, as shown in figures 3.2 and 3.3, respectively. There are 20 samples per class in the validation set, and 50 samples per class in the test set. **TODO: remove the figures, and just leave it in writing.**

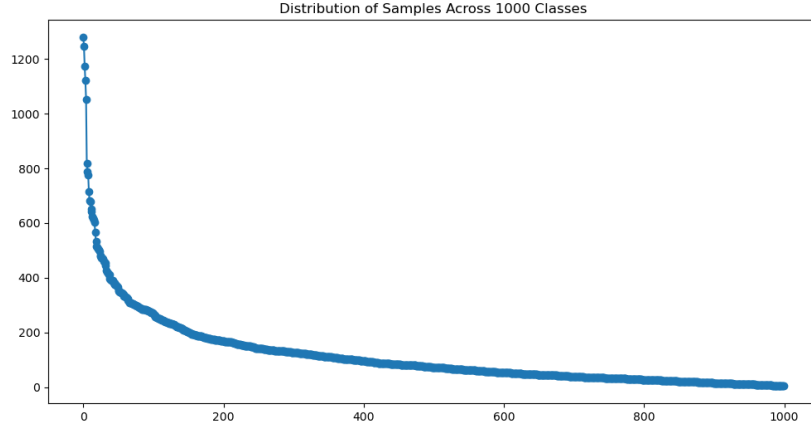


Figure 3.1: The class distribution of the training images for the ImageNet-LT dataset shows a long-tailed distribution.

However, investigating the CIFAR-100-LT dataset used in *Learning Imbalanced Datasets with Label-Distribution-Aware Margin Loss* by Cao et al. [16] shows that the training dataset follows an exponential decay, and not a Pareto distribution. The imbalance factor in their studies is 100, meaning that the most frequent class has 100 times more samples than the least frequent class, following equation (2.1). This means that the imbalance is not as steep as with the Pareto distribution, resulting in a middle section of classes with fewer samples than the head classes, but more samples than the tail classes, see figure 3.4. The following section will describe the preparation of CIFAR-100 for the experiments in this thesis.

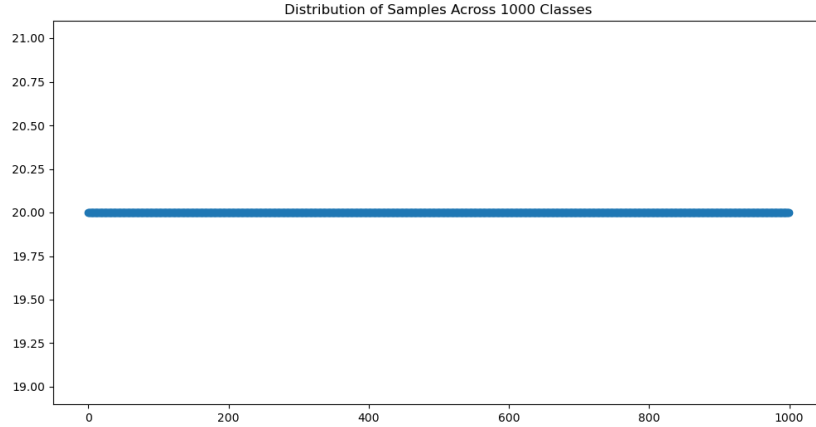


Figure 3.2: The class distribution of the validation images for the ImageNet-LT dataset shows that there are 20 samples of each class.

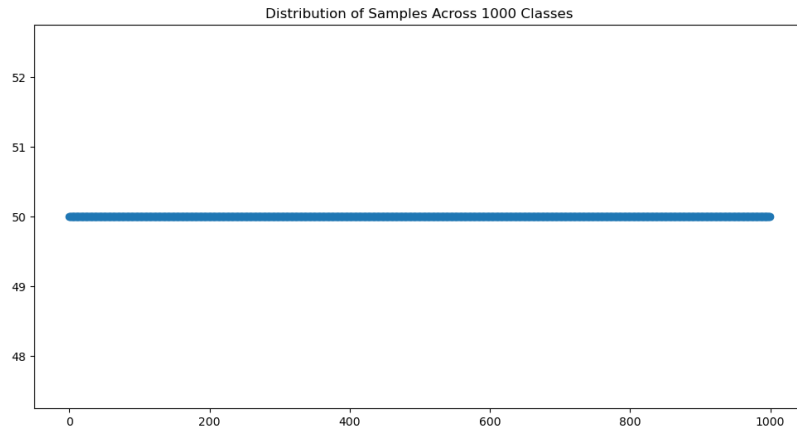


Figure 3.3: The class distribution of the test images for the ImageNet-LT dataset shows that there are 50 samples of each class.

3.2.3 CIFAR-100-LT

The experiments conducted in this thesis are trained and evaluated on the CIFAR-100 dataset. This was downloaded using the PyTorch `torchvision.datasets.CIFAR100` [55]. The training and test datasets were preprocessed by converting the images to tensors using the `ToTensor` transformation and saved as `.pth` files for efficient loading during experiments.

The class distributions used in experiments in this thesis: Balanced training set, balanced test set, balanced validation set, long-tailed training set, long-tailed test set, long-tailed test set split into head, middle and tail. The class distributions can be seen in table 3.2.

To address the needs of the experiments in this thesis, the original CIFAR-100 dataset was modified to create a new split from the training data. Specifically, the

Table 3.2: Class Distributions Used in Experiments

Dataset Type	Samples
Balanced Training Set	45,000
Balanced Validation Set	10,000
Balanced Test Set	5,000
Long-Tailed Training Set	TODO: -
Long-Tailed Test Set	TODO: -
Head	TODO: -
Middle	TODO: -
Tail	TODO: -

training data was split into 450 samples per class for training and 50 samples per class for testing, mirroring the test data for the ImageNet. The original test set of the CIFAR-100 dataset was retained as the validation set for evaluation during training. All datasets were saved locally to maintain reproducibility.

Long-tailed Training Set: To simulate real-world scenarios with class imbalances, the training dataset was modified to introduce an exponential imbalance across the 100 classes. Following the procedure of [16], the imbalance was created using exponential decay. Here, the number of samples per class decreases exponentially, controlled by the imbalance factor, following equation (2.1). For this thesis, an imbalance factor of 100 was applied. This means that the most frequent class contains a 100 times more samples than the least frequent class.

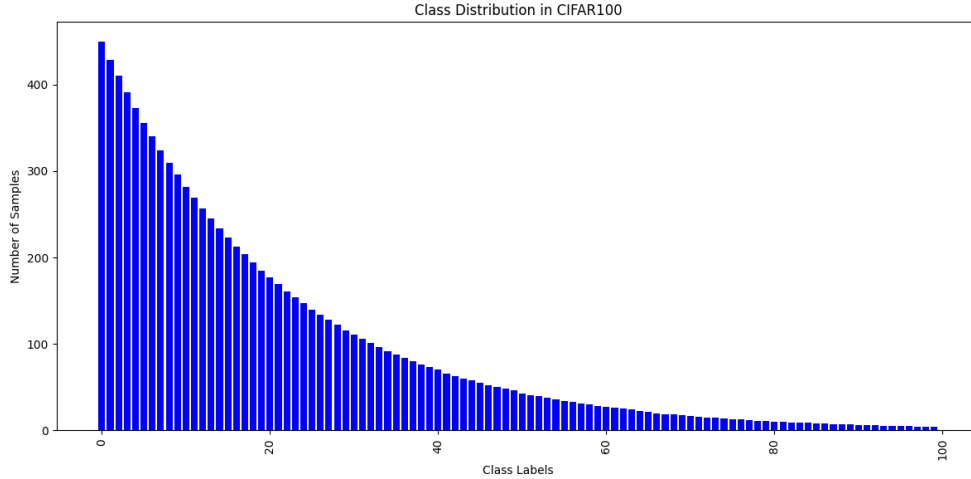


Figure 3.4: The class distribution of the CIFAR-100 long-tailed training set with 450 samples in the most frequent class.

The resulting class distribution varied from the most frequent class having 450 samples to the least frequent class having only 4 samples, as shown in figure 3.4. This imbalance ensured no class was left with zero samples, maintaining the integrity of all classes for training. The dataset was saved locally to maintain reproducibility.

Long-Tailed Test Set: To evaluate the performance of the model under similar conditions to the imbalanced training set, an imbalanced test set was created from the previously split test dataset. The imbalance in the test set mirrors the exponential distribution used for the training data, with the same imbalance factor of 0.01. The class distribution in the test set follows the same order of classes (from most to least frequent) as the imbalanced training set. No class has fewer than one sample. The class distribution can be seen in figure 3.5. The dataset was saved locally to maintain reproducibility.

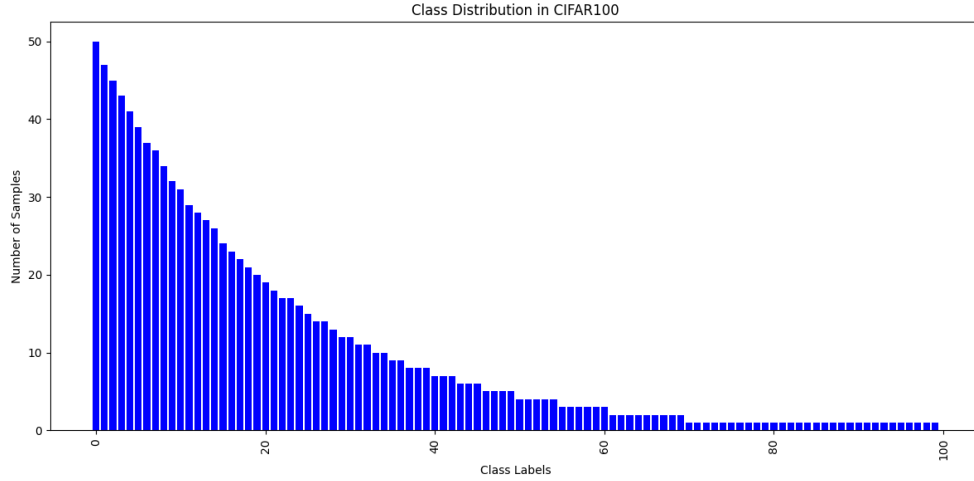


Figure 3.5: The class distribution of the CIFAR-100 long-tailed test set with 50 samples in the most frequent class.

To further analyze performance, the test set was divided into three subsets: head, middle, and tail classes. The head classes comprised the top 33% most frequent classes, the middle classes included the next 33%, and the tail classes represented the bottom 33%. See table 3.2. This categorization allows for a detailed evaluation of the methods across different classes.

3.2.4 Data Augmentation

Data augmentation and normalization were applied to the CIFAR-100 dataset to improve the generalization capability of the models. For the CIFAR-100 dataset, the augmentations were applied using the CIFAR-100 RGB statistics, with mean values [0.4914, 0.4822, 0.4465], and standard deviations [0.2023, 0.1994, 0.2010] [22]. These values are used to normalize the images after applying the transformations. However, the models used in this study are all pretrained on ImageNet, and choosing the ImageNet RGB statistics could potentially influence the training, as the pretrained weights are optimized for inputs with ImageNet statistics. Future work could explore the impact of using ImageNet statistics instead.

Images were resized to 224x224 pixels to match the input requirements of the ImageNet-pretrained models. Random cropping with a padding of 4 was applied to introduce spatial variations, while horizontal flipping with a 50% probability

further augmented the data. After these transformations, images were normalized using the CIFAR-100-specific mean and standard deviation values, ensuring consistency with the dataset.

For validation, a simpler approach was used, where images were resized to 224x224 pixels and normalized using the same mean and standard deviation values as the training set.

3.3 Long-tailed Learning Techniques

This section outlines the methods employed to address the challenges posed by a long-tailed data distribution. Specifically, it discusses the choice of model architecture and class re-balancing methods, including their strengths, limitations, and rationale based on prior research and feasible implementations.

3.3.1 Model Selection

A list of subjects to include in this section:

- Mention the model architectures chosen for training, and describe why they are appropriate for deep long-tailed learning with reference to the background section.
- Discuss the strengths and limitations of these models in addressing the challenges posed by imbalanced data.
- Describe how they were pretrained (ImageNet-1K, ImageNet-21K) and what that means for the training on CIFAR100.
- Write about their implementation. The implementation is the same for all models.
- Include the training loss for both balanced and imbalanced data.
- Compare to the results in their original paper.

This section will discuss the choice of model architectures addressing the challenges of long-tailed datasets. The models evaluated in this thesis are four common state-of-the-art (SOTA) models: ResNet-50, MobileNetV2, ViT-B/16, and ConvNeXt-Base. These SOTA models were selected to represent a diverse range of architectures for solving image classification tasks and to pose a balance between computational efficiency and state-of-the-art performance. The models were modified and trained using transfer learning with pre-trained weights from ImageNet [56]. Model complexity can be viewed in table 3.3.

For all models, the fully connected layer is modified to accommodate for the CIFAR-100 dataset's 100 classes.

Transfer learning: <https://www.kaggle.com/code/riteshsinha/transfer-learning-using-res>

Fine-tuning: <https://medium.com/@caterine/fine-tuning-a-model-using-pytorch-part-2-e>

https://www.tensorflow.org/tutorials/images/transfer_learning

Table 3.3: Models. **TODO: caption**

Model	Type	Pre-trained	Parameters (millions)	top-1 accuracy
MobileNetV2 [26]	CNN	ImageNet-1K	-	72.154 [57]
ResNet50 [37]	CNN	ImageNet-1K	-	80.858 [39]
ConvNeXt-Base [27]	CNN	ImageNet-1K	88.591.464	84.062 [58]
ViT-B/16 [28]	ViT	ImageNet-21K (ft on ImageNet-1K)	-	-

ResNet-50 Resnet-50 [37] is a convolutional neural network build with residual connections that mitigates the vanishing gradient problem, as discussed in section 2.2.2. Introduced by He et al. in 2015, the ResNet architecure achieved unprecedented results at the time and won ILSVRC 2015 in image classification tasks [59].

With its 50 layers, the ResNet-50 is capable of strong feature extraction, making it a solid choice for image classification on long-tailed datasets. Furthermore, ResNet-50 benefits from transfer learning, ensuring rich hierarchical features learned from a large-scale dataset, enabling better generalization on the imbalanced CIFAR-100 dataset [60, 61, 62]. Transfer learning was initialized with weights pre-trained on ImageNet-1K, and the corresponding top-1 accuracy can be viewed in table 3.3.

ResNet-50 has more parameters than its sibling architectures ResNet-16 and -34, and fewer than those of ResNet-101 and -152. The fewer layers and parameters, the model is faster to train and requires less memory [37]. This architecture often performs well on small datasets. The ResNet-50 generally provides a good trade-off between accuracy and speed. The larger ResNet architectures require more computational power, but are often more accurate on larger datasets [63].

The ResNet-50 is a good starting point for transfer learning due to its balance of performance and efficiency. To lower computational resources and time yet still maintaining accuracy, the ResNet-50 architecture is chosen as a reference point for the CNN design.

ResNets are often used in scenarios that benefit from transfer learning.

TODO: Write this and find references

MobileNetV2 The MobileNetV2 model is a successor to the ResNet architecture, adjusted to run on mobile devices. It requires fewer computational resources due to the introduction of inverted residual blocks and depthwise separable convolutions, as described in section 2.2.2. This architecture was chosen to represent an example of a modern and lightweight CNN architecture that has achieved state-of-the-art performance [26]. Its lightweight design makes it particularly attractive for applications where computational resources are limited. By including MobileNetV2 among the evaluated architectures, the performance is compared across varying complexity levels, providing insights into how efficiency-oriented designs fare against more complex models. This comparison is especially relevant if the application domain involves real-time processing or deployment on mobile or embedded devices.

ViT-B/16 The introduction of ViTs in 2021 [28] suddenly changed the status of CNNs as the state-of-the-art image classification architecture [45]. However, the CIFAR-100 images are only 32x32 pixels. The convolution and stride, as discussed in the previous section.

The implementation of the ViT-B/16 was done through the timm library [64], trained on ImageNet-21K and fine-tuned on ImageNet-1K, as suggested for ViTs [28]. The PyTorch library has a version of the ViT-B/16 trained on ImageNet-1K [44]. For consistency and direct comparison, the implementation of ViT-B/16 pre-trained on ImageNet-1k would be preferred, however, since this is a vision transformer, the model benefits from pretraining on a larger dataset and then fine-tuned due to the lack of inductive bias [28, 65], as discussed in section 2.2.3.

ConvNeXt Base In the paper *A ConvNet for the 2020s* by Liu et al. (2022) [45] the authors presents the ConvNeXt model architecture, a CNN with elements of ViTs (see section 2.2.3), which achieved state-of-the-art performance in image classification tasks **TODO: was it image classification?**.

Rationale for Model Selection These models were chosen to represent a spectrum of design paradigms; from the traditional ResNet to lightweight architectures and transformers. By evaluating models of varying complexity, this study aims to identify architectures best suited for long-tailed learning under different resource constraints and data distributions.

3.3.2 Selection of Loss Function

A list of subjects to include in this section:

- Describe the different loss functions and why they are appropriate for deep long-tailed learning with reference to the background section.
- Rationale for each loss function’s inclusion, focusing on its expected benefits for imbalanced classes and how it addresses the bias toward majority classes.
- Write pseudo code of the implementations

Softmax Loss Softmax Cross-Entropy is the standard loss function for classification tasks. It serves as a baseline in this study. However, it tends to favor head classes in imbalanced datasets because these classes dominate the gradient updates, as discussed in section 2.3.2. CE loss provides a reference point to assess the improvements achieved by other loss functions.

Weighted Softmax Loss

Focal Loss

Class-Balanced Loss

Balanced Softmax Loss The implementation details for the balanced softmax loss followed **TODO: just reference the paper.**

Equalization Loss

LDAM Loss

3.3.3 Optimizer

TODO: Influence of optimizer. Find SGD references. [54, 66].

3.4 Evaluation Metrics

The top-1 accuracy is a widely used benchmark evaluation metric in image classification tasks [1] **TODO: find more references.** It measures the proportion of samples where the model’s top prediction matches the ground truth label.

Due to the metric being well-suited for tasks with a single label corresponding to each input, the top-1 accuracy is used as the evaluation metric to measure the model performance on both balanced and long-tailed datasets in this thesis. Its role as a benchmark ensures compatibility with prior research, allowing direct comparisons of the results.

Performance under class imbalance is analyzed by calculating top-1 accuracy for subsets of the long-tailed test set, divided into head, middle, and tail classes. This separation of classes provides an insight into the model’s ability to handle the challenges of long-tailed data, such as preserving high performance on tail classes while maintaining accuracy on head classes.

3.4.1 Baselines

Baselines are crucial when evaluating the performance of deep learning methods **TODO: reference?**, especially in tasks involving class imbalance. In this thesis, the Softmax Cross-Entropy Loss, which is widely recognized for its simplicity and effectiveness in balanced image classification (see section 2.3.2) **TODO: references**, serves as the primary baseline for both balanced and long-tailed training. This ensures a consistent benchmark for understanding the performance of the proposed long-tailed methods while also facilitating comparability with prior research.

Additional to the softmax baseline, all methods across all models are trained on the balanced CIFAR-100 dataset to establish performance baselines. This allows for direct comparisons of the results when trained on the CIFAR-100-LT data with the same models and loss functions. This setup represents the optimal training conditions where all classes are equally represented, serving as a benchmark for the best achievable performance across all test sets.

3.5 Model Training and Evaluation

All models are trained on all loss functions with both the balanced and long-tailed versions of the CIFAR-100 dataset.

Chapter 4

Experimental Setup

This chapter focuses on the on the implementation details of the experiments conducted in this thesis. Here, the specifics of the training configurations are described.

Dataset Specifications The experiments in this thesis are conducted with CIFAR-100 [21] benchmark. The CIFAR-100 is split into a training set of 45,000 samples and a test set of 5,000 samples. Furthermore, the CIFAR-100-LT was created from the training set with an imbalance factor of 100 with the number of samples per class ranging from 450 in the most frequent class to just 4 in the least frequent class. The test set mirrors this distribution and maintains the same imbalance factor. See section 3.2.

To further analyze model performance across different class frequencies, the imbalanced test set is divided into three subsets: the head test set, containing the top one-third most frequent classes, the middle test set, which includes the middle one-third of classes, and the tail test set. The validation set is the original CIFAR-100 test set with 10,000 samples distributed equally across the 100 classes.

Training Configuration All experiments are implemented using PyTorch [67], with models ResNet50 [37], MobileNetV2 [26], and ConvNeXt Base [27] initialized using pretrained weights on the ImageNet-1K dataset, and model ViT-B/16 [28] with ImageNet-21K. Each model is modified to adapt to the CIFAR-100 dataset by replacing the final classification layer with a 100-class fully connected layer.

Training is performed over 90 epochs with a batch size of 128. The Adam optimizer [52] is employed with an initial learning rate of 0.001, which is reduced by a factor of 0.1 every 30 epochs using a StepLR scheduler. The training is conducted on four NVIDIA TITAN X (Pascal, 12 GB). See appendix ?? for further specifications.

Evaluation The metric used for performance evaluation is the top-1 accuracy.

Reproducibility To ensure reproducibility, a fixed random seed of 42 is used across all experiments, and GPU computations enforced to be deterministic. Configuration files in YAML format are used to document hyperparameters, dataset

settings, and model configurations, allowing for exact replication of experiments. Datasets and model checkpoints are saved to facilitate reuse and further analysis.

Chapter 5

Results and Discussion

This chapter presents the experimental results, comparing model performances across head, middle, and tail classes, and discusses the impact of different model architectures combined with different loss functions. First, the main findings are presented, followed by a detailed analysis of the results across experiments, and lastly a summary and discussion of the results.

TODO: Compare the different types of losses against each other. See [54].

5.1 Main Findings

Across all models, Balanced Softmax Loss demonstrated the highest performance on tail classes for models trained on long-tailed datasets while maintaining consistent performance on head, middle, and overall long-tailed test sets. The highest top-1 accuracy for tail classes was achieved by the ResNet50 architecture (Acc1: 0.6053), closely followed by the ConvNeXt Base architecture (Acc1: 0.5789). However, this improved tail-class performance comes at the cost of head-class accuracy, where ConvNeXt Base outperforms ResNet50 with a top-1 accuracy of 0.8685 compared to 0.8270. Overall, ConvNeXt Base demonstrates better performance across all classes (Acc1: 0.8230) compared to ResNet50 (Acc1: 0.7916). See Tables 5.4 and 5.8 for reference. These findings, however, require further statistical analysis to confirm their significance and that the observed differences are not simply due to random variability.

Class-Balanced Loss consistently underperformed, warranting further investigation into its implementation. Similarly, the ViT-B/16 architecture demonstrated the lowest overall accuracy when trained on both balanced and long-tailed datasets (Acc1: 59.06 %, see Table 5.5), despite having the highest reported benchmark performance trained with balanced CIFAR-100 (Acc1: 93.95 %) among all model architectures investigated in this thesis [68]. This discrepancy suggests that the ViT-B/16 configuration used for the experiments in this thesis may not be well-suited for smaller-scale datasets such as CIFAR-100 without careful optimization as discussed in sections 2.2.3 and 3.3.1 **TODO: write these sections.**

5.2 Overall Results

This section presents the overall results of all experiments conducted in this thesis, commenting on the best and worst performance of each loss design on a given model, but not directly comparing loss designs or models. It provides an overview of all findings to understand how each combination of architecture and loss function behaves.

TODO: Include that all models are trained with a balanced CIFAR-100 to provide a baseline for the class re-balancing strategies.

5.2.1 MobileNetV2

Results from Balanced Training

Table 5.1: Top-1 accuracy results for MobileNetV2 on the balanced dataset across all loss functions.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax	0.7978	0.8059	0.8069	0.7870	0.8684
Focal loss	0.8014	0.8011	0.7998	0.7870	0.8947
Weighted Softmax loss	0.7978	0.8059	0.8069	0.7870	0.8684
Class-balanced loss	0.7978	0.8059	0.8069	0.7870	0.8684
Balanced Softmax loss	0.8034	0.8030	0.8069	0.7574	0.9211
Equalization loss	0.7994	0.8040	0.8057	0.7692	0.9211
LDAM loss	0.7828	0.7821	0.7808	0.7574	0.9211

From Table 5.1, Balanced Softmax Loss achieves the best accuracy on the balanced test dataset (Acc1: 0.8034) and the highest accuracy on tail classes (Acc1: 0.9211). Not surprisingly, Softmax, Weighted Softmax, and Class-Balanced Loss yield the same accuracies across all test datasets due to shared cross-entropy-based designs that become indistinguishable when trained on balanced data according to equations (2.6), (2.7), and (2.10) described in section 2.3. Overall, all methods achieve comparable results to the softmax baseline, while outperforming the best published result for MobileNetV2 trained with CIFAR-100, as seen in table 5.9.

TODO: This paragraph could potentially be moved to an overall discussion section, as it includes a discussion of other model behavior. Although LDAM Loss demonstrates lower overall accuracy, its high tail-class accuracy matches that of Balanced Softmax and Equalization Loss, suggesting that tail-class performance may have reached a saturation point due to the limited number of samples in tail classes. This notion is further supported by identical performance across multiple long-tailed subsets with other backbones (see tables 5.3 and 5.7), suggesting potential limitations in class-specific data quality or quantity. These limitations could be investigated through testing on a larger dataset.

Results from Long-Tailed Training

Table 5.2: Top-1 accuracy results for MobileNetV2 on the long-tailed dataset across all loss functions.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax	0.5282	0.7735	0.8341	0.5917	0.2368
Focal loss	0.5200	0.7745	0.8389	0.5917	0.1579
Weighted Softmax loss	0.5016	0.7231	0.7808	0.5503	0.2105
Class-balanced loss	0.1936	0.0913	0.0521	0.2485	0.2632
Balanced Softmax loss	0.5796	0.7650	0.8069	0.6331	0.4211
Equalization loss	0.5310	0.7650	0.8235	0.5917	0.2368
LDAM loss	0.4264	0.5899	0.6137	0.5444	0.2632

From Table 5.2, Balanced Softmax Loss outperforms other losses in terms of balanced test accuracy (0.5796) and especially tail classes (0.4211). Although Softmax and Focal Loss surpass Balanced Softmax on the long-tailed dataset and head classes, the improvements Balanced Softmax provides for tail and middle classes are notable. Class-Balanced Loss shows an overwhelming underperformance compared with the balanced training in table 5.1, and now appears severely limited, possibly due to implementation issues that must be investigated **TODO: investigate**. The difference in performance patterns compared to the softmax baseline shows that loss functions designed for imbalance can provide more nuanced performance gains.

5.2.2 ResNet50

Results from Balanced Training

Table 5.3: Evaluation results for ResNet50 trained on the custom balanced dataset, showing Acc1.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.8324	0.8421	0.8448	0.8047	0.9474
Focal loss	0.8310	0.8344	0.8341	0.8166	0.9211
Weighted Softmax loss	0.8324	0.8421	0.8448	0.8047	0.9474
Class-balanced loss	0.8324	0.8421	0.8448	0.8047	0.9474
Balanced Softmax loss	0.8310	0.8430	0.8460	0.8107	0.9211
Equalization loss	0.8292	0.8373	0.8412	0.7929	0.9474
LDAM loss	0.7990	0.7983	0.8069	0.7337	0.8947

From Table 5.3, Softmax, Weighted Softmax, and Class-Balanced Loss yield identical performance when trained on balanced data due to their cross-entropy term

in equations (2.10), (2.7), and (2.10). Balanced Softmax Loss demonstrates strong performance on the long-tailed test set (0.8430) and head classes (0.8460), while still performing comparably on tail classes. Overall, the performance of the class-rebalancing methods compare to that of the softmax baseline. The best performing design yield an accuracy of 83.24%, not far behind the best published result for ResNet50 with an accuracy of 86.90% as seen in table 5.9.

Results from Long-Tailed Training

Table 5.4: Evaluation results for ResNet50 trained on the long-tailed dataset, showing Acc1.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.5522	0.7954	0.8531	0.6391	0.2105
Focal loss	0.5456	0.7935	0.8483	0.6272	0.3158
Weighted Softmax loss	0.4976	0.7336	0.7915	0.5562	0.2368
Class-balanced loss	0.2052	0.1836	0.1445	0.3787	0.1842
Balanced Softmax loss	0.5908	0.7916	0.8270	0.6568	0.6053
Equalization loss	0.5452	0.7897	0.8389	0.6450	0.3421
LDAM loss	0.3742	0.5937	0.6469	0.4438	0.0789

From Table 5.4, Balanced Softmax Loss significantly improves tail-class accuracy (0.6053), outperforming other losses, while also outperforming on balanced and middle sets. Softmax Loss achieves the best performance on the long-tailed dataset and head classes, indicating that while Balanced Softmax improves minority-class accuracy, Softmax remains strong for majority classes. Class-Balanced Loss and LDAM Loss both show underwhelming results, pointing toward the necessity of further investigation for Class-Balanced Loss **TODO: investigate** or combining LDAM with DRW [16].

5.2.3 ViT-B/16

Results from Balanced Training

From Table 5.5, LDAM Loss achieves the overall best performance for ViT-B/16 trained on balanced data. However, the overall performance remain far below the benchmark results for Vision Transformers on CIFAR-100, see table 5.9. This suggests that further adjustments are needed, possibly including extended training time, data augmentation, or different optimizers, as ViT architecture differs from CNNs, as describe in section 2.2.3.

Results from Long-Tailed Training

From Table 5.6, performance remains low across all loss functions. Balanced Softmax Loss slightly improves results on balanced and middle categories, while

Table 5.5: Evaluation results for ViT-B/16 trained on the custom balanced dataset, showing Acc1.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.5620	0.5671	0.5521	0.6036	0.7368
Focal loss	0.5516	0.5538	0.5438	0.5680	0.7105
Weighted Softmax loss	0.5620	0.5671	0.5521	0.6036	0.7368
Class-balanced loss	0.5620	0.5671	0.5521	0.6036	0.7368
Balanced Softmax loss	0.5628	0.5642	0.5640	0.5325	0.7105
Equalization loss	0.5634	0.5519	0.5462	0.5503	0.6842
LDAM loss	0.5906	0.6013	0.5924	0.6095	0.7632

Table 5.6: Evaluation results for ViT-B/16 trained on the long-tailed dataset, showing Acc1.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.2254	0.4367	0.5071	0.1775	0.0263
Focal loss	0.2210	0.4206	0.4834	0.1953	0.0263
Weighted Softmax loss	0.1284	0.1760	0.1919	0.1302	0.0263
Class-balanced loss	0.0558	0.0076	0.0000	0.0237	0.1053
Balanced Softmax loss	0.2460	0.4244	0.4822	0.2130	0.0789
Equalization loss	0.2168	0.4215	0.4893	0.1716	0.0263
LDAM loss	0.1570	0.2750	0.3140	0.1361	0.0263

Softmax Loss achieves the highest accuracy on the long-tailed dataset and head classes. The best accuracy on tail classes (0.1053) is achieved by Class-Balanced Loss, though it fails everywhere else. This pattern suggests that the ViT-B/16 architecture may not be sufficient for robust performance on smaller-scale, long-tailed datasets and may require major adjustments in the training pipeline, see section 2.2.3.

5.2.4 ConvNeXt Base

Results from balanced Training Dataset

From Table 5.7, Balanced Softmax Loss achieves the best balanced test accuracy (0.8364) while the Softmax architectures tie for top performance on the long-tailed dataset, head classes, and tail classes. The saturation of tail-class performance is likely due to the quality or quantity of tail classes. Balanced Softmax’s advantage on balanced data may stem from its logit-adjustment mechanism, which can enhance generalization even when classes are balanced [49] **TODO: investigate**. Although the best published result for a ConvNext architecture trained with CIFAR-100 is the Conv2NeXt with an accuracy of 83.82% (see table 5.9), the best performance of on the balanced test set is in alignment, yielding an accuracy of 83.64%.

Table 5.7: Evaluation results for ConvNeXt Base trained on the custom balanced dataset, showing Acc1.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.8332	0.8535	0.8566	0.8166	0.9474
Focal loss	0.8314	0.8487	0.8507	0.8284	0.8947
Weighted Softmax loss	0.8332	0.8535	0.8566	0.8166	0.9474
Class-balanced loss	0.8332	0.8535	0.8566	0.8166	0.9474
Balanced Softmax loss	0.8364	0.8344	0.8365	0.7988	0.9474
Equalization loss	0.8318	0.8468	0.8448	0.8343	0.9474
LDAM loss	0.8316	0.8373	0.8412	0.8047	0.8947

Results from Long-Tailed Training Dataset

Table 5.8: Evaluation results for ConvNeXt Base trained on the long-tailed dataset, showing Acc1.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.5972	0.8316	0.8898	0.6568	0.3158
Focal loss	0.5938	0.8145	0.8685	0.6568	0.3158
Weighted Softmax loss	0.4090	0.6356	0.6848	0.4911	0.1842
Class-balanced loss	0.0142	0.0019	0.0000	0.0000	0.0526
Balanced Softmax loss	0.6460	0.8230	0.8685	0.6509	0.5789
Equalization loss	0.5956	0.8278	0.8768	0.6923	0.3421
LDAM loss	0.3770	0.5956	0.6445	0.4260	0.2632

From Table 5.8, Balanced Softmax Loss outperforms on the balanced test set (0.6460) and tail classes (0.5789) by a considerable margin. Softmax achieves the best performance on the long-tailed dataset and head classes, while Equalization Loss leads on middle classes. The poor performance of Class-Balanced Loss stands out, calling for further investigation into its implementation **TODO: investigate**.

5.3 Benchmarks

Table 5.9 compares the best results from the models trained on the balanced CIFAR-100 dataset in this thesis with published benchmarks using the same or similar architectures. This comparison serves as a reference to ensure the results align with expectations and to identify any potential discrepancies.

Notably, there are no published benchmarks for ConvNeXt Base trained on CIFAR-100. The closest benchmark available is from Feng et al. [69], where Conv2NeXt achieved a top-1 accuracy of 83.82%, which is highly similar to the result obtained here.

For MobileNetV2, the results in this thesis surpass the published benchmarks. This could be attributed to improved loss function designs. On the other hand,

ResNet50 exhibits a slightly lower accuracy compared to benchmarks, which may be explained by differences in optimization strategies, such as the use of Adam optimizer instead of SGD [54].

The most significant deviation is observed with ViT-B/16, where the accuracy falls far short of the benchmark. This discrepancy suggests that the configuration for ViT-B/16 may not be optimized for training with small-scale datasets, see section 2.2.3.

Overall, these results indicate that the experimental setup is effective for most models but highlights potential limitations for transformer-based architectures such as ViT-B/16. Future investigation is necessary to understand and address these discrepancies.

Table 5.9: Comparison of model performance on balanced CIFAR-100 with published benchmarks.

Model	Result (Acc1)	Benchmark Result (Acc1)
MobileNetV2	80.34%	73.20% [70]
ResNet50	83.24%	86.90% [71]
ViT-B/16	59.06%	93.95% [68]
ConvNeXt Base	83.64%	83.82% [69]

TODO: The original ViT paper has a CIFAR-100 benchmark as well [28].

5.4 Performance Comparisons

The models are compared by taking the mean of the results of the loss functions for each model. Excluding the generally underperforming Class-Balanced Loss, the mean and standard deviation across all loss designs for each model is plotted in figure 5.1. **TODO: Include a table of the values.** Here, the performance of the models MobileNetV2, ResNet50, ViT-B/16, and ConvNeXt Base is shown across the evaluation categories: balanced, long-tailed, head, middle, and tail. Each model is slightly offset along the x-axis to avoid overlap. Each point represents the mean accuracy of a model for a specific category and the bars represent the standard deviation. Again, the mean and standard deviation excludes those of the Class-Balanced Loss. The error bars indicate the variability in accuracy across different loss functions for each model and category. A longer error bar demonstrate a less consistent performance, dependent on the chosen loss function, while shorter error bars indicate consistent performance adhering to the specific model. From figure 5.1, the shorter errorbars for MobileNetV2 indicate that its architecture may be more robust to the type of re-weighting strategy used than other architectures. Contrary, the ResNet50V2 performance on tail classes shows greater sensitivity to the chosen loss design. ViT-B/16 shows a consistently lower mean accuracy, indicating that transformer-based approaches may require additional adaptations for smaller-scale or imbalanced datasets, as discussed in 3.3.1. ConvNeXt Base maintains strong and stable performance across several categories, benefiting from Balanced Softmax Loss as seen in table 5.8.

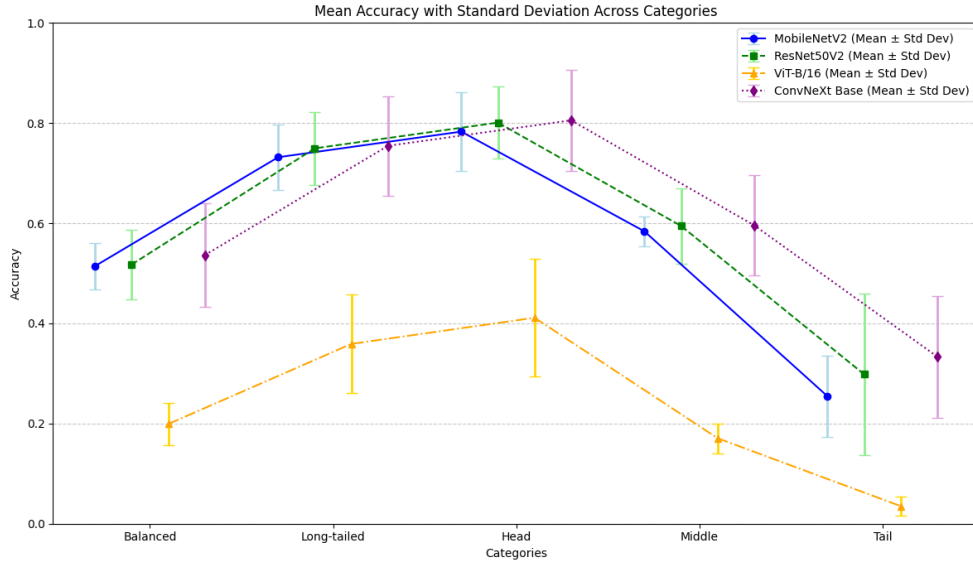


Figure 5.1: Mean Accuracy with Standard Deviation Across Categories for MobileNetV2, ResNet50, ViT-B/16, and ConvNeXt Base trained with CIFAR-100-LT. These values are excluding those of the Class-Balanced Loss.

The Balanced Softmax Loss emerges as the overall best performing loss design across all models with the least variance in performance across test categories, as figures 5.2, 5.3, 5.5, and 5.4 illustrate.

TODO: Include variance across test categories in either figures or tables.

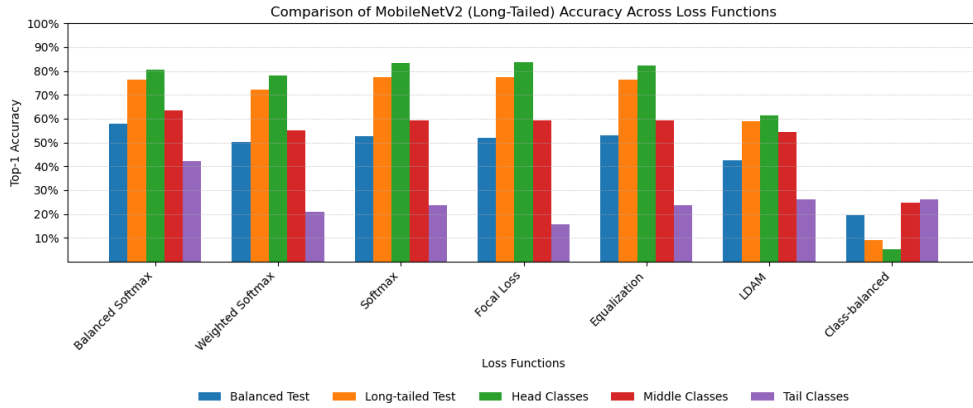


Figure 5.2: MobileNetV2 top 1 accuracy across test categories (Balanced, Long-tailed, Head, Middle, Tail) for different loss functions.

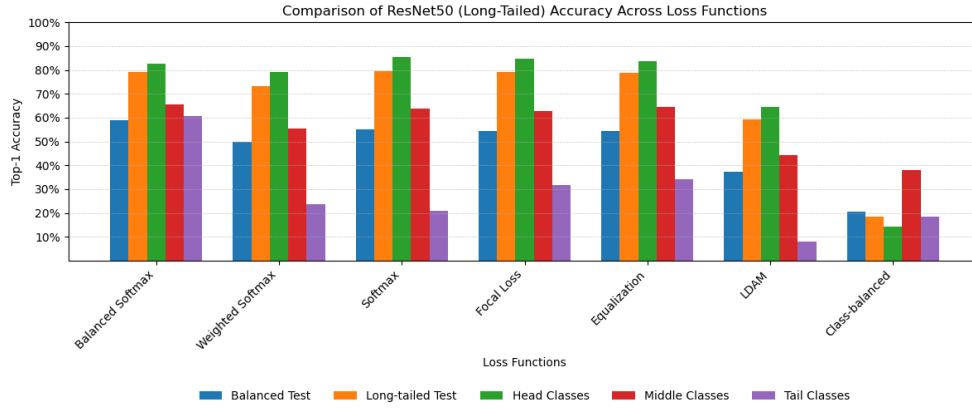


Figure 5.3: ResNet50 top 1 accuracy across test categories (Balanced, Long-tailed, Head, Middle, Tail) for different loss functions.

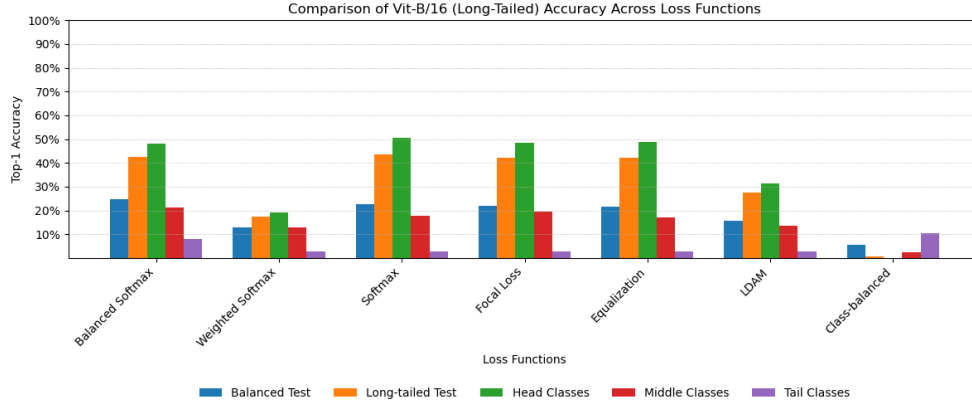


Figure 5.4: ViT-B/16 top 1 accuracy across test categories (Balanced, Long-tailed, Head, Middle, Tail) for different loss functions.

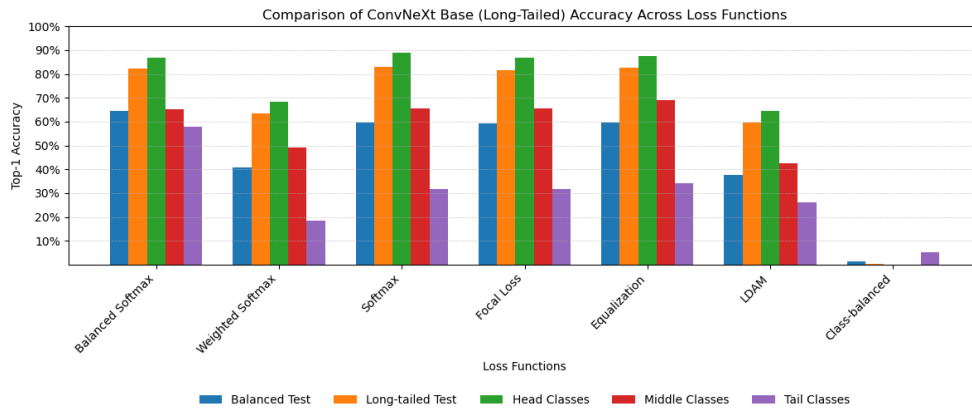


Figure 5.5: ConvNeXt Base top 1 accuracy across test categories (Balanced, Long-tailed, Head, Middle, Tail) for different loss functions.

This result is further underlined from figure 5.6, where the mean and standard deviation of the loss functions across ResNet50, MobileNetV2, ViT-B/16, and ConvNeXt is calculated and plotted. Here, the Balanced Softmax Loss shows a better average performance on tail classes, while also displaying the smoothest variation across evaluation categories **TODO: calculate variance**, ignoring the Class-Balanced loss. However, the standard deviation suggests that the performance of the Balanced Softmax Loss on tail classes depend on the model architecture, as earlier discussed. Table 5.10 display the exact values of the mean and standard deviations shown in figure 5.6.

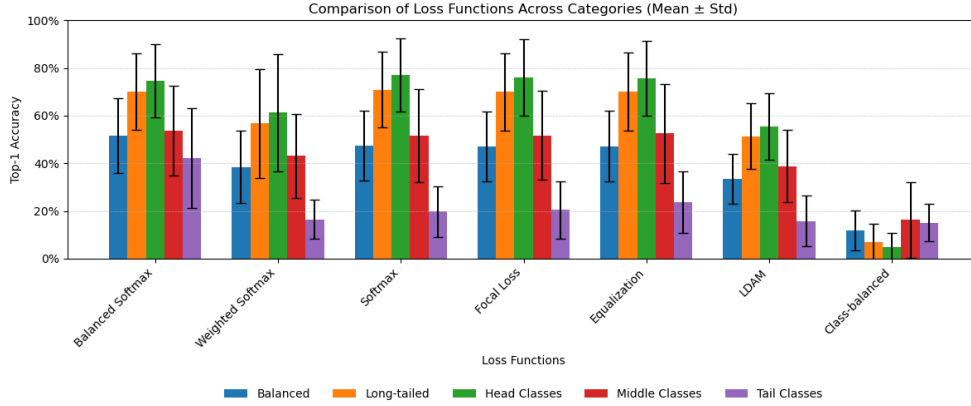


Figure 5.6: Performance trends of different loss functions across evaluation categories. Error bars indicate the standard deviation of accuracy across models, highlighting variability in performance for each loss function.

Table 5.10: Performance Summary Across Categories (Mean \pm Std)

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Balanced Softmax	0.5156 \pm 0.1577	0.7010 \pm 0.1610	0.7462 \pm 0.1540	0.5385 \pm 0.1881	0.4211 \pm 0.2097
Weighted Softmax	0.3841 \pm 0.1522	0.5671 \pm 0.2290	0.6123 \pm 0.2462	0.4320 \pm 0.1761	0.1645 \pm 0.0819
Softmax	0.4758 \pm 0.1466	0.7093 \pm 0.1587	0.7710 \pm 0.1537	0.5163 \pm 0.1970	0.1974 \pm 0.1061
Focal Loss	0.4701 \pm 0.1462	0.7008 \pm 0.1624	0.7598 \pm 0.1599	0.5178 \pm 0.1876	0.2039 \pm 0.1211
Equalization	0.4721 \pm 0.1494	0.7010 \pm 0.1629	0.7571 \pm 0.1558	0.5252 \pm 0.2072	0.2368 \pm 0.1289
LDAM	0.3504 \pm 0.1224	0.5187 \pm 0.1443	0.5548 \pm 0.1792	0.3947 \pm 0.1886	0.1587 \pm 0.1078
Class-balanced	0.1172 \pm 0.0814	0.0711 \pm 0.0841	0.0490 \pm 0.0624	0.1628 \pm 0.1477	0.1513 \pm 0.0946

5.5 Summary and Discussion

The overall best performance on tail classes was achieved by the Balanced Softmax Loss, which proved to be consistent across all models while also achieving competing results on other test categories. Furthermore, these findings highlight the impact of model architecture combined with loss designs, as the standard deviation demonstrates when comparing figures 5.1 and 5.6. Here, the longer errorbars suggests the dependency of loss design and model architecture, respectively. Further investigation of model and loss design combinations could lead to greater tail-class performance.

Although the CB Loss’ overall superior performance on tail classes, this performance accuracy often came with trade-offs in performance on head classes. For example, while ResNet50 achieved the highest top 1 accuracy of 60.53% on tail classes using Balanced Softmax, compared that of ConvNeXt Base (57.89%), its head-class accuracy (82.70%) lagged behind ConvNeXt Base (86.85%). This trade-off highlights the challenge of optimizing for both head and tail classes simultaneously and indicates that Balanced Softmax prioritizes tail-class adjustments at the expense of head-class performance. These results, however, should be challenged by multiple runs to check for statistical significance.

Class-Balanced Loss, in contrast, consistently underperformed across all models and datasets. This discrepancy between its intended purpose and observed results suggests possible issues with its implementation **TODO: investigate**. Further investigation is necessary to understand these limitations and determine whether its performance can be improved. **TODO: Make a table with original findings for comparison**.

Likewise, the ViT-B/16 architecture underperformed significantly across all loss functions, categories, and training datasets, with the best accuracy of 59.06% on a balanced training dataset compared to its benchmark of 93.95% **TODO: reference**. These results indicate that the default ViT-B/16 architecture may not be well-suited for long-tailed datasets without further optimization [54, 66]. These could include a carefully chosen optimizer, longer training, or larger dataset. Comparatively, the CNN architectures displayed overall better performance with the same setup. **TODO: here’s another reference** <https://arxiv.org/pdf/2205.01580v1>.

TODO: Include discussion of the choice of adam optimizer over the SGD.

TODO: Include a discussion of the implementation of the EQ Loss.

The performance of the LDAM could potentially improve by including the deferred re-weighting scheme (DRW), also introduced by Cao et al. [16]. According to related work [54], the DRW significantly improves performance on long-tailed datasets in contrast to relying solely on LDAM for effective class re-balancing. **TODO: reference benchmark performance**.

The decision to derive test data from the training set, rather than splitting the validation set to match the class distribution to that of the ImageNet-LT (see section 3.2), could also influence the results, since the training data would include more samples. Likewise, the performance of the combinations of model architectures and loss designs could benefit from training with a larger dataset, like ImageNet and ImageNet-LT. Furthermore, these larger long-tailed benchmarks follow the Pareto distribution rather than an exponential decay, mirroring a natural occurring long-tailed dataset, as discussed in sections 2.1 and 3.2.

While Balanced Softmax stood out as the most robust loss function overall with the smallest variance across categories **TODO: make this**, statistical validation is required to confirm the significance of these findings by running multiple training sessions.

Overall, these findings emphasize the importance of aligning loss functions with both dataset characteristics and model architectures to address the challenges of

deep long-tailed learning. They also highlight the need for a nuanced approach that balances improvements in tail-class accuracy with minimal trade-offs in head-class and overall performance.

TODO: Adam optimization: Figure 1: Mean and standard deviation over 5 runs of per-class weight norms for a ResNet-32 under momentum and Adam optimisers. We use long-tailed (“LT”) versions of CIFAR-10 and CIFAR-100, and sort classes in descending order of frequency; the first class is 100 times more likely to appear than the last class. Both optimisers yield solutions with comparable balanced error. However, the weight norms have incompatible trends: under momentum, the norms are strongly correlated with class frequency, while with Adam, the norms are anti-correlated or independent of the class frequency. Consequently, weight normalisation under Adam is ineffective for combatting class imbalance [54]. adam optimizer: <https://medium.com/@weidagang/demystifying-the-adam-optimizer-in-machine-Useful-for-large-datasets> [52]

Chapter 6

Conclusion and Future Work

The findings from the experiments conducted in this thesis underlines the difficulty of achieving strong performances with deep long-tailed learning. Specifically, the results emphasizes the importance of thoughtfully combining model architectures, loss design, and choice of configurations. Among the evaluated approaches, the Balanced Softmax Loss consistently achieved the best performance on tail classes, while maintaining a comparable performance on other categories across multiple CNN backbones. In contrast, Class-Balanced Loss underperformed across all models and categories, suggesting a fault in implementation **TODO: investigate**. Likewise, ViT-B/16’s performance gap from both CNN-based architectures and its own benchmark imply that vision transformer architectures may require extensive adaptations to reach their full potential.

6.1 Revisiting the Goals of the Thesis

TODO: Add the goal from the introduction section and answer them below for the sake of the reader. In revisiting the goals of this thesis, it has become clear that the investigation of the deep long-tailed learning technique *Class Re-Balancing* has yielded valuable insight. The findings included that performance on tail classes often comes with a trade-off in performance on head classes. Similarly, adequate performances on all classes was often attributed to the performance on head classes and not inherently on tail classes, although the Balanced Softmax Loss convincingly overcame this challenge across CNN-based architectures. This aligns with the goal of investigating the ability of deep long-tailed learning techniques to maintain overall accuracy.

Secondly, the evaluations of different model architectures, namely the CNN-based and Vision Transformers, in combination with different class re-balancing designs has revealed that model design has an effect on the performance of long-tailed learning, and that the MobileNetV2 architecture exhibits greater robustness against choice of loss design than other architectures evaluated in this thesis.

Finally, these findings contribute to a more comprehensive understanding of what to consider when training a deep network with long-tailed data, and, in doing so, serves as a guide to deep learning practitioners seeking more informed choices

when encountered with a long-tailed problem.

6.2 Future Work

Future investigation would benefit from evaluating experiments with larger datasets, such as ImageNet or iNaturalist to gain more insight into the performance of models and long-tailed techniques. Additionally, training ViT-B/16 on larger datasets while carefully considering its configurations could prove its full potential. Furthermore, evaluating the LDAM including the DRW and, potentially, combined with other loss functions on CIFAR-100-LT could reveal under which conditions tail-class performance is maximized. Finally, other techniques, such as re-sampling in combination with class re-balancing, should be explored.

Bibliography

- [1] Yifan Zhang et al. “Deep long-tailed learning: A survey”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).
- [2] Chongsheng Zhang et al. *A Systematic Review on Long-Tailed Learning*. 2024. arXiv: 2408.00483 [cs.LG]. URL: <https://arxiv.org/abs/2408.00483>.
- [3] N. V. Chawla et al. “SMOTE: Synthetic Minority Over-sampling Technique”. In: *Journal of Artificial Intelligence Research* 16 (June 2002), pp. 321–357. ISSN: 1076-9757. DOI: 10.1613/jair.953. URL: <http://dx.doi.org/10.1613/jair.953>.
- [4] Hui Han, Wen-Yuan Wang, and Bing-Huan Mao. “Borderline-SMOTE: a new over-sampling method in imbalanced data sets learning”. In: *Proceedings of the 2005 International Conference on Advances in Intelligent Computing - Volume Part I*. ICIC’05. Hefei, China: Springer-Verlag, 2005, pp. 878–887. ISBN: 3540282262. DOI: 10.1007/11538059_91. URL: https://doi.org/10.1007/11538059_91.
- [5] Xi Yin et al. *Feature Transfer Learning for Deep Face Recognition with Under-Represented Data*. 2019. arXiv: 1803.09014 [cs.CV]. URL: <https://arxiv.org/abs/1803.09014>.
- [6] Ziwei Liu et al. *Large-Scale Long-Tailed Recognition in an Open World*. 2019. arXiv: 1904.05160 [cs.CV]. URL: <https://arxiv.org/abs/1904.05160>.
- [7] Luis Perez and Jason Wang. *The Effectiveness of Data Augmentation in Image Classification using Deep Learning*. 2017. arXiv: 1712.04621 [cs.CV]. URL: <https://arxiv.org/abs/1712.04621>.
- [8] Connor Shorten and Taghi M. Khoshgoftaar. “A survey on image data augmentation for deep learning”. In: *Journal of Big Data* 6.1 (2019), pp. 1–48.
- [9] Hongyi Zhang et al. *mixup: Beyond Empirical Risk Minimization*. 2018. arXiv: 1710.09412 [cs.LG]. URL: <https://arxiv.org/abs/1710.09412>.
- [10] Sangdoo Yun et al. *CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features*. 2019. arXiv: 1905.04899 [cs.CV]. URL: <https://arxiv.org/abs/1905.04899>.
- [11] Zhengzhuo Xu, Zenghao Chai, and Chun Yuan. *Towards Calibrated Model for Long-Tailed Visual Recognition from Prior Perspective*. 2021. arXiv: 2111.03874 [cs.CV]. URL: <https://arxiv.org/abs/2111.03874>.
- [12] Jianfeng Wang et al. *RSG: A Simple but Effective Module for Learning Imbalanced Datasets*. 2021. arXiv: 2106.09859 [cs.CV]. URL: <https://arxiv.org/abs/2106.09859>.

-
- [13] Bingyi Kang et al. *Decoupling Representation and Classifier for Long-Tailed Recognition*. 2020. arXiv: 1910.09217 [cs.CV]. URL: <https://arxiv.org/abs/1910.09217>.
 - [14] Boyan Zhou et al. *BBN: Bilateral-Branch Network with Cumulative Learning for Long-Tailed Visual Recognition*. 2020. arXiv: 1912.02413 [cs.CV]. URL: <https://arxiv.org/abs/1912.02413>.
 - [15] Xudong Wang et al. *Long-tailed Recognition by Routing Diverse Distribution-Aware Experts*. 2022. arXiv: 2010.01809 [cs.CV]. URL: <https://arxiv.org/abs/2010.01809>.
 - [16] Kaidi Cao et al. *Learning Imbalanced Datasets with Label-Distribution-Aware Margin Loss*. 2019. arXiv: 1906.07413 [cs.LG]. URL: <https://arxiv.org/abs/1906.07413>.
 - [17] MEJ Newman. “Power laws, Pareto distributions and Zipf’s law”. In: *Contemporary Physics* 46.5 (Sept. 2005), pp. 323–351. ISSN: 1366-5812. DOI: 10.1080/00107510500052444. URL: <http://dx.doi.org/10.1080/00107510500052444>.
 - [18] Grant Van Horn et al. *The iNaturalist Species Classification and Detection Dataset*. 2018. arXiv: 1707.06642 [cs.CV]. URL: <https://arxiv.org/abs/1707.06642>.
 - [19] Jia Deng et al. “ImageNet: A large-scale hierarchical image database”. In: *2009 IEEE Conference on Computer Vision and Pattern Recognition*. 2009, pp. 248–255. DOI: 10.1109/CVPR.2009.5206848.
 - [20] Charika de Alvis and Suranga Seneviratne. *A Survey of Deep Long-Tail Classification Advancements*. 2024. arXiv: 2404.15593 [cs.LG]. URL: <https://arxiv.org/abs/2404.15593>.
 - [21] Alex Krizhevsky and Geoffrey Hinton. *Learning Multiple Layers of Features from Tiny Images*. Technical Report. University of Toronto, 2009.
 - [22] Di Kai. *LDAM-DRW: Learning from Long-Tailed Data*. GitHub repository, [Online]. Available: <https://github.com/kaidic/LDAM-DRW/tree/master>. Accessed: 2024-09-18.
 - [23] LG AI Research. *[ICML 2022] Part 1: Long-Tail Distribution Learning*. <https://www.lgresearch.ai/blog/view/?seq=257&page=1&pageSize=12>. Accessed: 2024-11-26.
 - [24] Grant Van Horn and Pietro Perona. *The Devil is in the Tails: Fine-grained Classification in the Wild*. 2017. arXiv: 1709.01450 [cs.CV]. URL: <https://arxiv.org/abs/1709.01450>.
 - [25] Yin Cui et al. *Class-Balanced Loss Based on Effective Number of Samples*. 2019. arXiv: 1901.05555 [cs.CV]. URL: <https://arxiv.org/abs/1901.05555>.
 - [26] Mark Sandler et al. “Mobilenetv2: Inverted residuals and linear bottlenecks”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018, pp. 4510–4520.
 - [27] Agastya Todi et al. “ConvNext: A Contemporary Architecture for Convolutional Neural Networks for Image Classification”. In: *2023 3rd International Conference on Innovative Sustainable Computational Technologies (CISCT)*. IEEE. 2023, pp. 1–6.

- [28] Alexey Dosovitskiy et al. *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*. 2021. arXiv: 2010.11929 [cs.CV]. URL: <https://arxiv.org/abs/2010.11929>.
- [29] Aston Zhang et al. *Dive into Deep Learning*. <https://D2L.ai>. Cambridge University Press, 2023.
- [30] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. “Multilayer feedforward networks are universal approximators”. In: *Neural Networks* 2.5 (1989), pp. 359–366. ISSN: 0893-6080. DOI: [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8). URL: <https://www.sciencedirect.com/science/article/pii/0893608089900208>.
- [31] *What Is a Convolutional Neural Network?* MathWorks, [Online]. Available: <https://www.mathworks.com/discovery/convolutional-neural-network.html>. Accessed: 28-Nov-2024.
- [32] Yann Lecun et al. “Comparison of learning algorithms for handwritten digit recognition”. In: Jan. 1995.
- [33] Y. LeCun et al. “Backpropagation Applied to Handwritten Zip Code Recognition”. In: *Neural Computation* 1.4 (1989), pp. 541–551. DOI: 10.1162/neco.1989.1.4.541.
- [34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. “ImageNet Classification with Deep Convolutional Neural Networks”. In: *Advances in Neural Information Processing Systems*. Ed. by F. Pereira et al. Vol. 25. Curran Associates, Inc., 2012. URL: https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf.
- [35] Karen Simonyan and Andrew Zisserman. *Very Deep Convolutional Networks for Large-Scale Image Recognition*. 2015. arXiv: 1409.1556 [cs.CV]. URL: <https://arxiv.org/abs/1409.1556>.
- [36] Christian Szegedy et al. *Going Deeper with Convolutions*. 2014. arXiv: 1409.4842 [cs.CV]. URL: <https://arxiv.org/abs/1409.4842>.
- [37] Kaiming He et al. *Deep Residual Learning for Image Recognition*. 2015. arXiv: 1512.03385 [cs.CV]. URL: <https://arxiv.org/abs/1512.03385>.
- [38] Kaiming He et al. *Identity Mappings in Deep Residual Networks*. 2016. arXiv: 1603.05027 [cs.CV]. URL: <https://arxiv.org/abs/1603.05027>.
- [39] TorchVision Contributors. *ResNet Implementation in TorchVision*. <https://github.com/pytorch/vision/blob/main/torchvision/models/resnet.py>. Accessed: 2024-12-18. 2023.
- [40] Andrew G. Howard et al. *MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications*. 2017. arXiv: 1704.04861 [cs.CV]. URL: <https://arxiv.org/abs/1704.04861>.
- [41] Ashish Vaswani et al. *Attention Is All You Need*. 2023. arXiv: 1706.03762 [cs.CL]. URL: <https://arxiv.org/abs/1706.03762>.
- [42] V7 Labs. *Vision Transformer Guide*. <https://www.v7labs.com/blog/vision-transformer-guide>. Accessed: 2024-12-18. 2024.

-
- [43] Christopher Kim. *What is the Inductive Bias in Machine Learning?* <https://medium.com/@chkim345/what-is-the-inductive-bias-in-machine-learning-212a5f53e9aa>. Accessed: 2024-12-18. 2020.
 - [44] PyTorch Team. *Vision Transformer (ViT-B-16)*. https://pytorch.org/vision/main/models/generated/torchvision.models.vit_b_16.html. Accessed: 2024-12-04. 2024.
 - [45] Zhuang Liu et al. *A ConvNet for the 2020s*. 2022. arXiv: 2201.03545 [cs.CV]. URL: <https://arxiv.org/abs/2201.03545>.
 - [46] Saining Xie et al. *Aggregated Residual Transformations for Deep Neural Networks*. 2017. arXiv: 1611.05431 [cs.CV]. URL: <https://arxiv.org/abs/1611.05431>.
 - [47] *torch.nn.CrossEntropyLoss — PyTorch documentation*. <https://pytorch.org/docs/stable/generated/torch.nn.CrossEntropyLoss.html>. Accessed: 2024-11-20.
 - [48] Tsung-Yi Lin et al. *Focal Loss for Dense Object Detection*. 2018. arXiv: 1708.02002 [cs.CV]. URL: <https://arxiv.org/abs/1708.02002>.
 - [49] Jiawei Ren et al. *Balanced Meta-Softmax for Long-Tailed Visual Recognition*. 2020. arXiv: 2007.10740 [cs.LG]. URL: <https://arxiv.org/abs/2007.10740>.
 - [50] Jingru Tan et al. *Equalization Loss for Long-Tailed Object Recognition*. 2020. arXiv: 2003.05176 [cs.CV]. URL: <https://arxiv.org/abs/2003.05176>.
 - [51] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.
 - [52] Diederik P. Kingma and Jimmy Ba. *Adam: A Method for Stochastic Optimization*. 2017. arXiv: 1412.6980 [cs.LG]. URL: <https://arxiv.org/abs/1412.6980>.
 - [53] X. L. Chaitanya Asawa. *CS231n: Convolutional Neural Networks for Visual Recognition*. <http://cs231n.github.io/>. Stanford, [Online]. 2024.
 - [54] Aditya Krishna Menon et al. *Long-tail learning via logit adjustment*. 2021. arXiv: 2007.07314 [cs.LG]. URL: <https://arxiv.org/abs/2007.07314>.
 - [55] PyTorch Contributors. *torchvision.datasets.CIFAR100*. Accessed: 2024-12-13. 2024. URL: <https://pytorch.org/vision/0.17/generated/torchvision.datasets.CIFAR100.html>.
 - [56] Olga Russakovsky et al. “ImageNet Large Scale Visual Recognition Challenge”. In: *International Journal of Computer Vision (IJCV)* 115.3 (2015), pp. 211–252. DOI: 10.1007/s11263-015-0816-y.
 - [57] PyTorch Team. *MobileNetV2 — TorchVision 0.15.0 Documentation*. Accessed: 2024-12-22. 2024. URL: https://pytorch.org/vision/main/models/generated/torchvision.models.mobilenet_v2.html#torchvision.models.mobilenet_v2.
 - [58] TorchVision Contributors. *ConvNeXt Implementation in TorchVision*. <https://github.com/pytorch/vision/blob/main/torchvision/models/convnext.py>. Accessed: 2024-12-18. 2023.
 - [59] Sumit Saha. *Review: ResNet — Winner of ILSVRC 2015 (Image Classification, Localization, Detection)*. Accessed: 2024-12-22. 2018. URL: <https://towardsdatascience.com/review-resnet-winner-of-ilsvrc-2015-image-classification-localization-detection-e39402bfa5d8>.

- [60] Le Liu et al. “A Transfer Learning Method Based on ResNet Model”. In: *Data Mining and Big Data*. Ed. by Ying Tan et al. Singapore: Springer Singapore, 2021, pp. 250–260.
- [61] Mohammad Razavi, Samira Mavaddati, and Hamidreza Koohi. “ResNet deep models and transfer learning technique for classification and quality detection of rice cultivars”. In: *Expert Systems with Applications* 247 (2024), p. 123276. ISSN: 0957-4174. DOI: <https://doi.org/10.1016/j.eswa.2024.123276>. URL: <https://www.sciencedirect.com/science/article/pii/S0957417424001416>.
- [62] Kenneth Chan. *A Guide to Transfer Learning with Keras Using ResNet50*. Accessed: 2024-12-23. 2019. URL: <https://medium.com/@kenneth.ca95/a-guide-to-transfer-learning-with-keras-using-resnet50-a81a4a28084b>.
- [63] Piyush Nagpal, Shivani Atul Bhinge, and Ajitkumar Shitole. “A Comparative Analysis of ResNet Architectures”. In: *2022 International Conference on Smart Generation Computing, Communication and Networking (SMART GENCON)*. 2022, pp. 1–8. DOI: 10.1109/SMARTGENCON56628.2022.10083966.
- [64] Hugging Face. *ViT Base Patch16-224 by Google*. <https://huggingface.co/google/vit-base-patch16-224>. Accessed: 2024-12-04. 2024.
- [65] Alexander Kolesnikov et al. *Big Transfer (BiT): General Visual Representation Learning*. 2020. arXiv: 1912.11370 [cs.CV]. URL: <https://arxiv.org/abs/1912.11370>.
- [66] Ilya Loshchilov and Frank Hutter. *Fixing Weight Decay Regularization in Adam*. 2018. URL: <https://openreview.net/forum?id=rk6qdGgCZ>.
- [67] Adam Paszke et al. *PyTorch: An Imperative Style, High-Performance Deep Learning Library*. 2019. arXiv: 1912.01703 [cs.LG]. URL: <https://arxiv.org/abs/1912.01703>.
- [68] Ching-Hsun Tseng et al. “Perturbed Gradients Updating within Unit Space for Deep Learning”. In: *2022 International Joint Conference on Neural Networks (IJCNN)*. IEEE, July 2022, pp. 01–08. DOI: 10.1109/ijcnn55064.2022.9892245. URL: <http://dx.doi.org/10.1109/IJCNN55064.2022.9892245>.
- [69] Jianwei Feng et al. “Conv2NeXt: Reconsidering Conv NeXt Network Design for Image Recognition”. In: *2022 International Conference on Computers and Artificial Intelligence Technologies (CAIT)*. 2022, pp. 53–60. DOI: 10.1109/CAIT56099.2022.10072172.
- [70] Geon Park et al. *BiTAT: Neural Network Binarization with Task-dependent Aggregated Transformation*. 2022. arXiv: 2207.01394 [cs.CV]. URL: <https://arxiv.org/abs/2207.01394>.
- [71] Ross Wightman, Hugo Touvron, and Hervé Jégou. *ResNet strikes back: An improved training procedure in timm*. 2021. arXiv: 2110.00476 [cs.CV]. URL: <https://arxiv.org/abs/2110.00476>.
- [72] PyTorch Team. *ConvNeXt Base Model*. https://pytorch.org/vision/main/models/generated/torchvision.models.convnext_base.html. Accessed: 2024-12-04. 2024.

Appendix A

Results

Tables of the results from training.

A.1 MobileNetV2

MobileNetV2 trained on custom balanced dataset and imbalanced dataset on different loss functions.

Table A.1 show the top 1 accuracies for MobileNetV2 on various loss functions. Table A.2 show the loss, top 1 accuracy, and F1 score. Table A.3 show the top 1 accuracies for MobileNetV2 on various loss functions. Table A.4 show the loss, top 1 accuracy, and F1 score.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax	0.7978	0.8059	0.8069	0.7870	0.8684
Focal loss	0.8014	0.8011	0.7998 4	0.7870	0.8947
Weighted Softmax loss	0.7978	0.8059	0.8069	0.7870	0.8684
Class-balanced loss	0.7978	0.8059	0.8069	0.7870	0.8684
Balanced Softmax loss	0.8034	0.8030	0.8069	0.7574	0.9211
Equalization loss	0.7994	0.8040	0.8057	0.7692	0.9211
LDAM loss	0.7828	0.7821	0.7808	0.7574	0.9211

Table A.1: Evaluation results for MobileNetV2 trained on the custom balanced dataset, showing Acc1.

Loss Function	Balanced			Long-tailed			Head			Middle			Tail		
	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1
Softmax	1.1455	0.7978	0.7967	1.1415	0.8059	0.8208	1.1208	0.8069	0.8587	1.3060	0.7870	0.8368	0.8690	0.8684	0.8684
Focal Loss	0.6765	0.8014	0.8001	0.7063	0.8011	0.8175	0.7005	0.7998	0.8531	0.8028	0.7870	0.8293	0.4055	0.8947	0.8860
Weighted Softmax	1.1455	0.7978	0.7967	1.1415	0.8059	0.8208	1.1208	0.8069	0.8587	1.3060	0.7870	0.8368	0.8690	0.8684	0.8684
Class-balanced	1.1455	0.7978	0.7967	1.1415	0.8059	0.8208	1.1208	0.8069	0.8587	1.3060	0.7870	0.8368	0.8690	0.8684	0.8684
Balanced Softmax	1.1289	0.8034	0.8011	1.1848	0.8030	0.8145	1.1469	0.8069	0.8553	1.4407	0.7574	0.8123	0.8872	0.9211	0.9298
Equalization	0.9992	0.7994	0.7983	1.0385	0.8040	0.8192	1.0118	0.8057	0.8564	1.2539	0.7692	0.8213	0.6035	0.9211	0.9211
LDAM	13.8126	0.7828	0.7817	13.5566	0.7821	0.7955	13.6884	0.7808	0.8325	14.7496	0.7574	0.8016	5.3231	0.9211	0.9035

Table A.2: Evaluation results for MobileNetV2 trained on the custom balanced dataset, showing Loss, Acc1, and F1 scores for each dataset split.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax	0.5282	0.7735	0.8341	0.5917	0.2368
Focal loss	0.5200	0.7745	0.8389	0.5917	0.1579
Weighted Softmax loss	0.5016	0.7231	0.7808	0.5503	0.2105
Class-balanced loss	0.1936	0.0913	0.0521	0.2485	0.2632
Balanced Softmax loss	0.5796	0.7650	0.8069	0.6331	0.4211
Equalization loss	0.5310	0.7650	0.8235	0.5917	0.2368
LDAM loss	0.4264	0.5899	0.6137	0.5444	0.2632

Table A.3: Evaluation results for MobileNetV2 trained on the long-tailed dataset showing Acc1.

Loss Function	Balanced			Long-tailed			Head			Middle			Tail		
	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1
Softmax	3.2503	0.5282	0.4884	1.2212	0.7735	0.7578	0.8136	0.8341	0.8492	2.4604	0.5917	0.6444	4.7629	0.2368	0.2544
Focal Loss	2.3526	0.5200	0.4818	0.8022	0.7745	0.7602	0.5177	0.8389	0.8528	1.5864	0.5917	0.6625	3.6343	0.1579	0.1667
Weighted Softmax	3.1412	0.5016	0.4690	1.2817	0.7231	0.7104	0.8786	0.7808	0.8015	2.3365	0.5503	0.6213	5.1836	0.2105	0.1912
Class-balanced	4.3308	0.1936	0.1751	4.2181	0.0913	0.0854	4.4197	0.0521	0.0795	2.8788	0.2485	0.2767	6.3575	0.2632	0.2368
Balanced Softmax	3.1185	0.5796	0.5572	1.1630	0.7650	0.7685	0.7989	0.8069	0.8422	2.1612	0.6331	0.6872	4.8108	0.4211	0.4123
Equalization	3.0593	0.5310	0.4911	1.1563	0.7650	0.7499	0.7241	0.8235	0.8398	2.4487	0.5917	0.6524	4.9284	0.2368	0.2544
LDAM	21.4896	0.4264	0.3980	7.9893	0.5899	0.5909	5.6756	0.6137	0.6581	10.3379	0.5444	0.6121	49.1197	0.2632	0.2895

Table A.4: Evaluation results for MobileNetV2 trained on the long-tailed dataset, showing Loss, Acc1, and F1 scores for each dataset split.

A.2 ResNet50V2

ResNet50V2 trained on custom balanced dataset and imbalanced dataset on different loss functions.

Table A.5 show the top 1 accuracies for ResNet50V2 on various loss functions. Table A.6 show the loss, top 1 accuracy, and F1 score.

Table A.7 show the top 1 accuracies for ResNet50V2 on various loss functions. Table A.8 show the loss, top 1 accuracy, and F1 score.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.8324	0.8421	0.8448	0.8047	0.9474
Focal loss	0.8310	0.8344	0.8341	0.8166	0.9211
Weighted Softmax loss	0.8324	0.8421	0.8448	0.8047	0.9474
Class-balanced loss	0.8324	0.8421	0.8448	0.8047	0.9474
Balanced Softmax loss	0.8310	0.8430	0.8460	0.8107	0.9211
Equalization loss	0.8292	0.8373	0.8412	0.7929	0.9474
LDAM loss	0.7990	0.7983	0.8069	0.7337	0.8947

Table A.5: Evaluation results for ResNet50V2 trained on the custom balanced dataset, showing Acc1.

Loss Function	Balanced			Long-tailed			Head			Middle			Tail		
	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1
Softmax	0.9823	0.8324	0.8310	0.9874	0.8421	0.8520	0.9917	0.8448	0.8860	1.0934	0.8047	0.8467	0.4205	0.9474	0.9386
Focal Loss	0.5627	0.8310	0.8300	0.5578	0.8344	0.8474	0.5555	0.8341	0.8788	0.6294	0.8166	0.8607	0.2920	0.9211	0.9123
Weighted Softmax	0.9823	0.8324	0.8310	0.9874	0.8421	0.8520	0.9917	0.8448	0.8860	1.0934	0.8047	0.8467	0.4205	0.9474	0.9386
Class-balanced	0.9823	0.8324	0.8310	0.9874	0.8421	0.8520	0.9917	0.8448	0.8860	1.0934	0.8047	0.8467	0.4205	0.9474	0.9386
Balanced Softmax	1.0198	0.8310	0.8301	0.9689	0.8430	0.8549	0.9601	0.8460	0.8893	1.1309	0.8107	0.8539	0.4440	0.9211	0.9123
Equalization	0.8795	0.8292	0.8279	0.9079	0.8373	0.8495	0.8888	0.8412	0.8877	1.1374	0.7929	0.8453	0.2495	0.9474	0.9386
LDAM	9.8339	0.7990	0.7979	10.1092	0.7983	0.8119	9.8723	0.8069	0.8596	12.5229	0.7337	0.7823	4.6362	0.8947	0.8772

Table A.6: Evaluation results for ResNet50V2 trained on the custom balanced dataset, showing Loss, Acc1, and F1 scores for each dataset split.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.5522	0.7954	0.8531	0.6391	0.2105
Focal loss	0.5456	0.7935	0.8483	0.6272	0.3158
Weighted Softmax loss	0.4976	0.7336	0.7915	0.5562	0.2368
Class-balanced loss	0.2052	0.1836	0.1445	0.3787	0.1842
Balanced Softmax loss	0.5908	0.7916	0.8270	0.6568	0.6053
Equalization loss	0.5452	0.7897	0.8389	0.6450	0.3421
LDAM loss	0.3742	0.5937	0.6469	0.4438	0.0789

Table A.7: Evaluation results for ResNet50V2 trained on the long-tailed dataset, showing Acc1.

Loss Function	Balanced			Long-tailed			Head			Middle			Tail		
	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1
Softmax	3.0907	0.5522	0.5138	1.0524	0.7954	0.7798	0.6888	0.8531	0.8654	2.0330	0.6391	0.6996	4.7658	0.2105	0.2018
Focal Loss	2.0718	0.5456	0.5089	0.6284	0.7935	0.7789	0.3983	0.8483	0.8583	1.3258	0.6272	0.7054	2.6364	0.3158	0.3158
Weighted Softmax	3.7904	0.4976	0.4591	1.3481	0.7336	0.7198	0.8630	0.7915	0.8098	2.2625	0.5562	0.6209	6.9808	0.2368	0.2456
Class-balanced	4.5887	0.2052	0.1928	3.7422	0.1836	0.1932	3.7880	0.1445	0.2045	2.4884	0.3787	0.4138	8.3052	0.1842	0.1737
Balanced Softmax	3.1081	0.5908	0.5654	1.0452	0.7916	0.7895	0.6873	0.8270	0.8602	2.3422	0.6568	0.7135	3.2275	0.6053	0.5965
Equalization	3.0166	0.5452	0.5071	1.0315	0.7897	0.7756	0.7418	0.8389	0.8511	1.8754	0.6450	0.7061	3.6342	0.3421	0.3509
LDAM	22.7933	0.3742	0.3337	8.2056	0.5937	0.5784	5.3320	0.6469	0.6680	12.3074	0.4438	0.5450	53.4080	0.0789	0.0789

Table A.8: Evaluation results for ResNet50V2 trained on the long-tailed dataset, showing Loss, Acc1, and F1 scores for each dataset split.

A.3 ViT-B/16

ViT-B/16 trained on custom balanced dataset and imbalanced dataset on different loss functions.

Table A.9 show the top 1 accuracies for ViT-B/16 on various loss functions. Table A.10 show the loss, top 1 accuracy, and F1 score.

Table A.11 show the top 1 accuracies for ViT-B/16 on various loss functions. Table A.12 show the loss, top 1 accuracy, and F1 score.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.5620	0.5671	0.5521	0.6036	0.7368
Focal loss	0.5516	0.5538	0.5438	0.5680	0.7105
Weighted Softmax loss	0.5620	0.5671	0.5521	0.6036	0.7368
Class-balanced loss	0.5620	0.5671	0.5521	0.6036	0.7368
Balanced Softmax loss	0.5628	0.5642	0.5640	0.5325	0.7105
Equalization loss	0.5634	0.5519	0.5462	0.5503	0.6842
LDAM loss	0.5906	0.6013	0.5924	0.6095	0.7632

Table A.9: Evaluation results for ViT-B/16 trained on the custom balanced dataset, showing Acc1.

Loss Function	Balanced			Long-tailed			Head			Middle			Tail		
	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1
Softmax	4.6431	0.5620	0.5593	4.6089	0.5671	0.5951	4.6420	0.5521	0.6367	4.7263	0.6036	0.6648	3.3521	0.7368	0.7281
Focal Loss	2.3473	0.5516	0.5488	2.3562	0.5538	0.5869	2.4288	0.5438	0.6324	2.2330	0.5680	0.6355	1.2929	0.7105	0.6930
Weighted Softmax	4.6431	0.5620	0.5593	4.6089	0.5671	0.5951	4.6420	0.5521	0.6367	4.7263	0.6036	0.6648	3.3521	0.7368	0.7281
Class-balanced	4.6431	0.5620	0.5593	4.6089	0.5671	0.5951	4.6420	0.5521	0.6367	4.7263	0.6036	0.6648	3.3521	0.7368	0.7281
Balanced Softmax	4.7131	0.5628	0.5592	4.6809	0.5642	0.5929	4.7739	0.5640	0.6471	4.8161	0.5325	0.5998	2.0138	0.7105	0.7105
Equalization	4.2603	0.5634	0.5614	4.4906	0.5519	0.5884	4.6109	0.5462	0.6410	4.3952	0.5503	0.6014	2.0079	0.6842	0.6754
LDAM	48.2745	0.5906	0.5926	47.4149	0.6013	0.6348	49.6692	0.5924	0.6790	42.0117	0.6095	0.6780	21.3751	0.7632	0.7281

Table A.10: Evaluation results for ViT-B/16 trained on the custom balanced dataset, showing Loss, Acc1, and F1 scores for each dataset split.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.2254	0.4367	0.5071	0.1775	0.0263
Focal loss	0.2210	0.4206	0.4834	0.1953	0.0263
Weighted Softmax loss	0.1284	0.1760	0.1919	0.1302	0.0263
Class-balanced loss	0.0558	0.0076	0.0000	0.0237	0.1053
Balanced Softmax loss	0.2460	0.4244	0.4822	0.2130	0.0789
Equalization loss	0.2168	0.4215	0.4893	0.1716	0.0263
LDAM loss	0.1570	0.2750	0.3140	0.1361	0.0263

Table A.11: Evaluation results for ViT-B/16 trained on the long-tailed dataset, showing Acc1.

Loss Function	Balanced			Long-tailed			Head			Middle			Tail		
	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1
Softmax	13.5272	0.2254	0.1871	6.7999	0.4367	0.4216	5.3024	0.5071	0.5248	11.3663	0.1775	0.2303	19.7511	0.0263	0.0263
Focal Loss	7.5701	0.2210	0.1850	3.6474	0.4206	0.4016	2.8064	0.4834	0.4914	6.1246	0.1953	0.2666	11.3091	0.0263	0.0263
Weighted Softmax	6.5391	0.1284	0.1144	3.9782	0.1760	0.1902	3.4559	0.1919	0.2357	4.7288	0.1302	0.1541	11.0975	0.0263	0.0351
Class-balanced	4.9938	0.0558	0.0368	5.8487	0.0076	0.0028	6.2065	0.0000	0.0000	4.7503	0.0237	0.0292	4.0694	0.1053	0.0746
Balanced Softmax	13.3583	0.2460	0.2123	6.7016	0.4244	0.4175	5.2929	0.4822	0.5121	11.3472	0.2130	0.2710	17.3287	0.0789	0.0877
Equalization	13.4511	0.2168	0.1786	6.7202	0.4215	0.4062	5.2340	0.4893	0.5051	11.6755	0.1716	0.2353	17.4650	0.0263	0.0263
LDAM	43.5990	0.1570	0.1321	17.1295	0.2750	0.2769	11.5903	0.3140	0.3466	30.5370	0.1361	0.1791	80.9656	0.0263	0.0351

Table A.12: Evaluation results for ViT-B/16 trained on the long-tailed dataset, showing Loss, Acc1, and F1 scores for each dataset split.

A.4 ConvNeXt Base

ConvNeXt Base trained on custom balanced dataset and imbalanced dataset on different loss functions.

Table A.13 show the top 1 accuracies for ConvNeXt Base on various loss functions. Table A.14 show the loss, top 1 accuracy, and F1 score.

Table A.15 show the top 1 accuracies for ConvNeXt Base on various loss functions. Table A.16 show the loss, top 1 accuracy, and F1 score.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.8332	0.8535	0.8566	0.8166	0.9474
Focal loss	0.8314	0.8487	0.8507	0.8284	0.8947
Weighted Softmax loss	0.8332	0.8535	0.8566	0.8166	0.9474
Class-balanced loss	0.8332	0.8535	0.8566	0.8166	0.9474
Balanced Softmax loss	0.8364	0.8344	0.8365	0.7988	0.9474
Equalization loss	0.8318	0.8468	0.8448	0.8343	0.9474
LDAM loss	0.8316	0.8373	0.8412	0.8047	0.8947

Table A.13: Evaluation results for ConvNeXt Base trained on the custom balanced dataset, showing Acc1.

Loss Function	Balanced			Long-tailed			Head			Middle			Tail		
	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1
Softmax	0.9904	0.8332	0.8323	0.9594	0.8535	0.8661	0.9571	0.8566	0.9010	1.1028	0.8166	0.8603	0.3731	0.9474	0.9386
Focal Loss	0.5686	0.8314	0.8301	0.5597	0.8487	0.8608	0.5640	0.8507	0.8975	0.6046	0.8284	0.8730	0.2629	0.8947	0.8947
Weighted Softmax	0.9904	0.8332	0.8323	0.9594	0.8535	0.8661	0.9571	0.8566	0.9010	1.1028	0.8166	0.8603	0.3731	0.9474	0.9386
Class-balanced	0.9904	0.8332	0.8323	0.9594	0.8535	0.8661	0.9571	0.8566	0.9010	1.1028	0.8166	0.8603	0.3731	0.9474	0.9386
Balanced Softmax	1.0008	0.8364	0.8350	0.9829	0.8344	0.8478	0.9720	0.8365	0.8853	1.1509	0.7988	0.8418	0.4780	0.9474	0.9386
Equalization	0.9124	0.8318	0.8302	0.9030	0.8468	0.8594	0.8550	0.8448	0.8899	1.2187	0.8343	0.8779	0.4981	0.9474	0.9474
LDAM	12.2036	0.8316	0.8308	11.0787	0.8373	0.8485	10.9948	0.8412	0.8882	12.5744	0.8047	0.8592	6.2892	0.8947	0.8947

Table A.14: Evaluation results for ConvNeXt Base trained on the custom balanced dataset, showing Loss, Acc1, and F1 scores for each dataset split.

Text.

Loss Function	Balanced	Long-tailed	Head	Middle	Tail
Softmax loss	0.5972	0.8316	0.8898	0.6568	0.3158
Focal loss	0.5938	0.8145	0.8685	0.6568	0.3158
Weighted Softmax loss	0.4090	0.6356	0.6848	0.4911	0.1842
Class-balanced loss	0.0142	0.0019	0.0000	0.0000	0.0526
Balanced Softmax loss	0.6460	0.8230	0.8685	0.6509	0.5789
Equalization loss	0.5956	0.8278	0.8768	0.6923	0.3421
LDAM loss	0.3770	0.5956	0.6445	0.4260	0.2632

Table A.15: Evaluation results for ConvNeXt Basetrained on the long-tailed dataset, showing Acc1.

Loss Function	Balanced			Long-tailed			Head			Middle			Tail		
	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1	Loss	Acc1	F1
Softmax	2.7006	0.5972	0.5645	0.9552	0.8316	0.8202	0.5867	0.8898	0.9013	2.1181	0.6568	0.7177	3.9670	0.3158	0.3158
Focal Loss	1.8210	0.5938	0.5615	0.6024	0.8145	0.8002	0.3485	0.8685	0.8791	1.3247	0.6568	0.7197	3.0291	0.3158	0.3070
Weighted Softmax	4.5284	0.4090	0.3763	2.0092	0.6356	0.6266	1.5444	0.6848	0.7054	3.1309	0.4911	0.5827	6.0533	0.1842	0.1930
Class-balanced	5.0105	0.0142	0.0016	6.1643	0.0019	0.0000	6.5523	0.0000	0.0000	5.0450	0.0000	0.0000	3.4200	0.0526	0.0164
Balanced Softmax	2.6574	0.6460	0.6273	0.9120	0.8230	0.8237	0.5457	0.8685	0.8952	2.1945	0.6509	0.7250	3.3453	0.5789	0.5965
Equalization	2.5527	0.5956	0.5586	0.9349	0.8278	0.8139	0.6192	0.8768	0.8907	1.9792	0.6923	0.7375	3.2293	0.3421	0.3404
LDAM	39.0426	0.3770	0.3448	12.8480	0.5956	0.5812	8.3491	0.6445	0.6597	25.5813	0.4260	0.5105	72.0081	0.2632	0.2719

Table A.16: Evaluation results for ConvNeXt Base trained on the long-tailed dataset, showing Loss, Acc1, and F1 scores for each dataset split.

Appendix B

Experimental Setup Details

B.1 Dataset Specifications

Table B.1 provides an overview of the dataset splits used in this thesis. The training set consists of 45,000 samples, with 450 samples per class, generated by splitting the original CIFAR-100 training set into training and test sets. The test set derived from the split consists of 5,000 samples, with 50 samples per class. The validation set, unaltered from the original CIFAR-100 test set, contains 10,000 samples with an equal distribution of 100 samples per class.

To simulate real-world class imbalance scenarios, an exponential imbalance was introduced into the training and test sets. The imbalance factor (`imb_factor`) was set to 0.01, resulting in a significant reduction of samples for the least frequent classes. Table B.2 provides a summary of the imbalance characteristics.

Additionally, the imbalanced test set was further divided into three subsets based on class frequencies in the training data:

- **Head Test Set:** Includes the top one-third most frequent classes.
- **Middle Test Set:** Includes the middle one-third of classes.
- **Tail Test Set:** Includes the bottom one-third least frequent classes.

The sample specifics are presented in table B.1.

Table B.1: Dataset Specifications

Dataset Component	Total Samples	Samples per Class
Training Set	45,000	450
Validation Set	10,000	100
Test Set	5,000	50
Head Test Set	TODO: investigate	Variable
Middle Test Set	TODO: investigate	Variable
Tail Test Set	TODO: investigate	Variable

Table B.2: Imbalance Specifications

Aspect	Details
Imbalance Type	Exponential
Imbalance Factor	0.01
Training Set Distribution	Most frequent class: 450 samples Least frequent class: 4 samples
Test Set Distribution	Mirrors training distribution Ensures no class has fewer than 1 sample

B.2 Data Preprocessing

Table B.3 summarizes the preprocessing steps applied to the training, validation, and test datasets. **TODO: reference CIFAR100 statistics.**

Table B.3: Data Preprocessing Steps

Dataset	Preprocessing Steps
Training	<ul style="list-style-type: none"> • Resize to 224×224 pixels • Random crop to 224×224 with 4 pixels of padding • Random horizontal flip • Normalize using CIFAR-100 statistics: Mean = [0.4914, 0.4822, 0.4465], Std = [0.2023, 0.1994, 0.2010]
Validation/Test	<ul style="list-style-type: none"> • Resize to 224×224 pixels • Normalize using CIFAR-100 statistics: Mean = [0.4914, 0.4822, 0.4465], Std = [0.2023, 0.1994, 0.2010]

B.3 Model Architecture Settings

Four different architectures were used: MobileNetV2, ResNet50V2, ViT-B/16, and ConvNeXt Base. All models were initialized with pretrained weights from ImageNet to leverage transfer learning. The modifications ensured that the architectures were adapted to the CIFAR-100 dataset while retaining the general features learned during pretraining.

MobileNetV2 The MobileNetV2 architecture is pretrained on ImageNet-1K dataset, and the classification layer was replaced with a 100-class fully connected layer. **TODO: reference**

ResNet50V2 The ResNet50V2 architecture is pretrained on ImageNet-1K dataset, and the fc-layer is replaced with a 100-class fully connected layer. **TODO: reference**

ViT-B/16 The ViT-B/16 architecture is pretrained on the ImageNet-21K dataset and fine-tuned on the ImageNet-1K dataset [64]. The head is replaced with a 100-class fully connected layer.

ConvNeXt Base The ConvNeXt Base architecture is pretrained on the ImageNet-1K dataset [72], and the final layer is replaced with a 100-class fully connected layer.

The specifications of the model architectures can be seen in table B.4 **TODO: make this table prettier.**

Table B.4: Model Architecture Settings

Model Name	Pretrained Weights	Modifications for CIFAR-100
MobileNetV2	MobileNet_V2_Weights. IMAGENET1K_V1	Replaced classification layer with a 100-class fully connected layer
ResNet50V2	ResNet50_Weights. IMAGENET1K_V2	Replaced the <code>fc</code> layer with a 100-class fully connected layer
ViT-B/16	<code>timm vit_base_patch16_224</code> pretrained	Replaced the <code>head</code> with a 100-class fully connected layer
ConvNeXt Base	ConvNeXt_Base_Weights. DEFAULT	Replaced the final layer of the classifier with a 100-class fully connected layer

B.4 Training Configurations

This section outlines the key hyperparameters and settings used during the training and evaluation of the models.

- **Optimizer:** Adam
- **Learning Rate:** Initial value of 0.001.
- **Learning Rate Scheduler:** StepLR with a step size of 30 epochs and a decay factor of 0.1.
- **Batch Size:** 128.
- **Number of Epochs:** 90.
- **Weight Decay:** Default value of $1\text{e-}4$.
- **Class Weights:** Dynamically computed based on the training dataset and passed to the loss function. **TODO: Reference to Methodology section.**
- **Number of GPUs:** 4. See section B.6 for specifications.
- **Device Setup:** Utilized `torch.nn.DataParallel` for multi-GPU training.
- **Checkpoint Criteria:** The best model is saved based on the highest Top-1 validation accuracy.

B.5 Evaluation Metrics

The primary metric used to evaluate model performance during validation and testing is the top-1 accuracy. Alongside, the F1 score was calculated (macro F1 for balanced datasets, and weighted F1 for imbalanced dataset) but never used for evaluation. The F1 scores for all experiments can be seen in appendix A.

B.6 Hardware and Software Configurations

The experiments were conducted on a high-performance computing system with the following hardware and software configurations:

B.6.1 Hardware

- **GPUs:** 4 NVIDIA TITAN X (Pascal), each with 12 GB memory.
- **RAM:** 125 GiB
- **Swap Space:** 63 GiB
- **CUDA Version:** 12.4
- **Driver Version:** 550.90.07

B.6.2 Software

- **Operating System:** Ubuntu 22.04.4 LTS (Jammy Jellyfish)
- **Python Version:** 3.11.8
- **Deep Learning Frameworks and Libraries Versions:**
 - PyTorch (≥ 1.7)
 - Torchvision ($\geq 0.8.0$)
 - Tensorboard (≥ 1.14)

B.7 Reproducibility Considerations

To ensure that the experiments conducted in this thesis are reproducible, the following measures were implemented:

- **Random Seed:**
 - A fixed random seed of 42 was used for all experiments to ensure consistent initialization across runs.
 - Randomness was controlled for:

- * Python's `random` library.
- * NumPy (`np.random.seed`).
- * PyTorch (`torch.manual_seed` and `torch.cuda.manual_seed`).
- `torch.nn.deterministic` was set to `True` to enforce deterministic behavior in GPU computations.
- **Configuration Management:**
 - All hyperparameters, dataset settings, and model configurations were defined in YAML configuration files.
 - This allows for the exact replication of experiments by reusing the configuration files.
- **Saved Artifacts:**
 - Datasets were saved, making them accessible for evaluation or reuse in future experiments.
 - Model checkpoints were saved after achieving the best validation accuracy.

B.8 Implementation Faults

TODO: Revisit this section later on.

- The wrong version of the ViT-B/16 architecture might have implemented. The version implemented is the `vit-base-patch16-224` [64] via `timm` which was pretrained on the ImageNet-21K and later fine-tuned on ImageNet-1K, whereas the other models are pretrained solely on the ImageNet-1K. The version on ViT-B/16 that is pretrained on the ImageNet-1K can be implemented as the `torchvision` version `vit_b_16` [44].
- Class-Balanced Loss.