



Université Mohamed Khider de Biskra  
Faculté des Sciences et de la Technologie  
Département de génie électrique

# MÉMOIRE DE MASTER

Sciences et Technologies  
Télécommunication  
Réseaux ET Telecommunication

Réf. : Entrez la référence du document

---

Présenté et soutenu par :  
**Salhi Fardous**

Le : lundi 2 juin 2025

## Attention based model for blood cells images classification

---

### Jury :

M.	<b>Ouafi Abdelkrim</b>	Pr	<b>Université de Biskra</b>	<b>Président</b>
M.	<b>Sbaa Salim</b>	Pr	<b>Université de Biskra</b>	<b>Examineur</b>
M.	<b>Baarir Zineeddine</b>	Pr	<b>Université de Biskra</b>	<b>Rapporteur</b>

Année universitaire : 2024 / 2025



Université Mohamed Khider de Biskra  
Faculté des Sciences et de la Technologie  
Département de génie électrique

# MÉMOIRE DE MASTER

Sciences et Technologies  
Télécommunication  
Réseaux ET Telecommunication

Réf. : Entrez la référence du document

---

## Attention based model for blood cells images classification

Le : .....

**Présenté par :**

*Salhi Fardous*

**Avis favorable de l'encadreur :**

*Baarir Zineeddine*

**Signature Avis favorable du Président du Jury**

*Ouafi Abdelkrim*

**Cachet et signature**

### *Acknowledgements*

*I would like to express my heartfelt thanks and appreciation to the supervisor of this memorandum, DR. Baarir Ene Eddine, for his genuine efforts, scientific assistance, and invaluable guidance that had a significant influence on the journey of this research. His ongoing support and engagement were crucial to the successful completion of this study.*

*Additionally, I wish to convey my deep appreciation to the esteemed members of the discussion committee for their interest and insightful scientific feedback, which aided in the enhancement of this memorandum and enriched its content.*

*You have my sincerest gratitude and appreciation, and may God reward you with the highest blessings.*

*Salhi Fardous*

## *Dedication*

*To my dearest mother, whose boundless love and heartfelt prayers have been the constant light in my life's journey — your strength, compassion, and unwavering presence are the pillars upon which every step of mine has been built. I owe every achievement to your sacrifices, and I pray that God grants you everlasting peace and happiness.*

*To my father, my role model and the unwavering anchor in my life — your wisdom and calm guidance have shaped the person I am today. Your faith in me has been a driving force, and I ask God to reward your kindness with abundant health and joy.*

*To my entire family, thank you for being my safe haven and the solid ground I stand upon. Your encouragement and patience have given me the confidence to push forward. May our bond always remain strong and full of love.*

*To my treasured friends, your presence has brought warmth and brightness into the most challenging moments. Through your loyalty and uplifting spirits, you have helped me persevere. I sincerely pray for your continued success, happiness, and fulfillment in all that you do.*

*To my fellow colleagues, whose shared dedication and collaboration made this journey both meaningful and enriching — thank you for the inspiration, the teamwork, and the shared moments of growth. May our paths cross again in future successes.*

*Salhi fardous*

## Abstract:

White blood cells play a vital role in the human immune system, and variations in their count can indicate serious health conditions. This project presents a designed system for both localization and classification of white blood cells. The dataset used in this study consists of two parts: a localization set containing 364 annotated images and a classification set with 12,444-labeled images. The primary objective is to develop and assess an efficient deep learning-based system for accurately identifying and categorizing white blood cells.

Two distinct localization methods were explored: a conventional approach and a deep learning-based technique. Additionally, five deep learning models—comprising three pretrained architectures and two custom-built models—were evaluated for classification. Experimental results demonstrated that the localization process achieved an average Intersection over Union (IoU) of 71%, while the classification models attained an accuracy of 92%. The system's robustness was further analyzed by introducing various types of noise to assess its resilience. The high accuracy and robustness of the proposed approach can be attributed to the extensive dataset, optimized architectures, and carefully designed methodology.

**Keyword:** White Blood Cells (WBC), Intersection over Union (IoU), Dataset Annotation

## ملخص

تلعب خلايا الدم البيضاء دورًا حيويًا في الجهاز المناعي البشري، ويمكن أن تشير أي اختلافات في أعدادها إلى حالات صحية خطيرة. يقدم هذا البحث نظامًا مصممًا لتحديد موقع خلايا الدم البيضاء وتصنيفها. تتكون مجموعة البيانات المستخدمة في هذه الدراسة من جزأين: مجموعة تحديد موقع تحتوي على 364 صورة مُعلّقة، ومجموعة تصنيف تحتوي على 12,444 صورة مُعلّمة. الهدف الرئيسي هو تطوير وتقييم نظام فعال قائم على التعلم العميق لتحديد خلايا الدم البيضاء وتصنيفها بدقة.

تم استكشاف طريقتين مختلفتين لتحديد الموقع: نهج تقليدي وتقنية قائمة على التعلم العميق. بالإضافة إلى ذلك، تم تقييم خمسة نماذج تعلم عميق - تتكون من ثلاث هياكل مُدرّبة مسبقًا ونموذجين مُصمّمين خصيصًا - لأغراض التصنيف. أظهرت النتائج التجريبية أن عملية تحديد الموقع حققت متوسط تقاطع على الاتحاد (IoU) بنسبة 71%، بينما حققت نماذج التصنيف دقة 92%. تم تحليل متانة النظام بشكل أعمق من خلال إدخال أنواع مختلفة من الضوضاء لتقييم مرونته. يمكن أن تُعزى الدقة العالية والمتانة للنهج المقترح إلى مجموعة البيانات الواسعة والهندسة المعمارية المُحسّنة والمنهجية المصممة بعناية.

**الكلمات المفتاحية:** خلايا الدم البيضاء (WBC)، التقاطع على الاتحاد (IoU)، شرح مجموعة البيانات

## Résumé

Les globules blancs jouent un rôle essentiel dans le système immunitaire humain, et des variations de leur nombre peuvent indiquer de graves problèmes de santé. Cette projet présente un système conçu pour la localisation et la classification des globules blancs. L'ensemble de données utilisé dans cette étude se compose de deux parties : un ensemble de localisation contenant 364 images annotées et un ensemble de classification contenant 12 444 images étiquetées. L'objectif principal est de développer et d'évaluer un système efficace basé sur l'apprentissage profond pour identifier et catégoriser avec précision les globules blancs.

Deux méthodes de localisation distinctes ont été explorées : une approche conventionnelle et une technique basée sur l'apprentissage profond. De plus, cinq modèles d'apprentissage profond, comprenant trois architectures pré-entraînées et deux modèles personnalisés, ont été évalués pour la classification. Les résultats expérimentaux ont démontré que le processus de localisation atteignait une intersection sur union (IoU) moyenne de 71 %, tandis que les modèles de classification atteignaient une précision de 92 %. La robustesse du système a été analysée plus en détail en introduisant différents types de bruit afin d'évaluer sa résilience. La grande précision et la robustesse de l'approche proposée peuvent être attribuées à l'ensemble de données étendu, aux architectures optimisées et à la méthodologie soigneusement conçue.

**Mot clés:** Globules blancs (WBC), Réseaux de neurones convolutifs (CNN), Intersection sur Union (IoU), Annotation de jeu de données

## Table of Contents

Acknowledgements .....	
Dedication .....	
Abstract: .....	
Table of Contents .....	
List of figures .....	
List of tables .....	
List of abbreviations.....	
General introduction.....	1

## Chapitre 1

### Neural Network Models in Deep Learning

I. RELATED WORK .....	5
I.1 CNN based :.....	5
I.2 Attention Based Networks: .....	6
I.3 Deep learning: .....	7
I.3.1 Foundation Models and Scaling Laws .....	7
I.3.2 Multimodal Learning .....	7
I.3.3 Efficient and Low-Resource Models .....	7
I.3.4 Reasoning and Tool Use .....	8
I.3.5 Autonomous Agents and RLHF.....	8
I.3.6 Diffusion and Generative Models .....	8
I.3.7 Applications Across Domains.....	8
II. CNN based methods .....	10
II.1 Introduction .....	10
II.2 CNN Concept .....	11
II.2 Attention based methods .....	17

Introduction: .....	17
Convolutional Attention – General Architecture .....	19
1. Input Embedding / Feature Extraction .....	19
2. Convolutional Attention Block .....	19
3. Stacking Blocks .....	20
4. Output Head .....	20
Limitations and Challenges .....	20
Conclusion.....	21

## **Chapitre 2**

### **Computer Vision using Deep learning**

II.1 Introduction.....	23
II.2 Definition of Computer Vision .....	23
II.2.1 Types of Computer Vision Tasks [6] .....	25
II.3 Applications of Computer Vision in Medicine.....	26
II.3.1 Trends on Computer vision [8].....	28
II.1 Comment of figure.....	28
II.3.2 Computer Vision in Healthcare .....	28
II.4 Image Classification Overview.....	29
II.4.1 what is classification? .....	29
II.4.2 How It Works?.....	30
1. Data Collection:.....	30
2. Preprocessing: .....	31
3. Feature Extraction: .....	31
4. Model Training:.....	31
5. Classification (Prediction Phase): .....	32
6. Post-Processing (Optional):.....	32
7. Metrics Evaluation: .....	32

II.4.3 Types of Image Classification .....	33
II.4.4 Real-World Applications .....	33
II.4.5 Classification of White Blood Cells [11]: .....	34
II.4.6 Classification in medicines and biomedicines .....	34
II.5 Conclusion .....	35

### **Chapter 3:**

#### **Simulations and Discussion of results**

III.1 Introduction .....	37
III.2 Used Materials.....	37
III.3 Downloading and Extracting Files from Google Drive with Python.....	38
III.4 Python Module Imports for Data Preparation .....	39
III.6 ResNet50-GN Backbone .....	42
III.9 ResNetV2 .....	46
III.10 White Blood Cell Machine Learning Model Performance Ranking Report .....	47
III.11 Confusion matrix normalized for white blood cell classification .....	51
III.12 Conclusion.....	53
General Conclusion .....	56
References :.....	<b>Error! Bookmark not defined.</b>



## List of figures

Figure I.1: Architecture Basic MLP .....	11
Figure I.2: Process convolution CNN [3].....	12
Figure I.3: Convolution Layer [3] .....	12
Figure II.1: illustration of Computer Vision for Design Optimization .....	28
Figure II.2: The 5 cell types of the RAW – PBC DIB dataset a) Basophil b) Eosinophil c) Lymphocyte d) Monocyte e) Neutrophil[10] .....	30
Figure II.3: structure of the created Ensemble model [10].....	33
Figure III.1: Loss curve for training and validation .....	43
Figure III.2: F1 Score curve for training and verification.....	44
Figure III.3: Training and Validation Accuracy Curve.....	45
Figure III.4: Classification report for white blood cell types .....	47
Figure III.5: F1 Score by class .....	48
Figure III.6: Precision by class.....	49
Figure III.7: Recall by class .....	49
Figure III.8: WBC classification with 100% accuracy.....	50
Figure III.9: Number of samples by class .....	51
Figure III.10: Normalized confusion Matrix .....	52

## List of tables

### *Chapter I*

**Table I.1: Famous CNN-based methods**

**Table I.2: Famous Attention based methods**

### *Chapter II*

**Table II.1: Classification of white blood cells**

## **List of abbreviations**

CAD: Computer-Aided Diagnosis

CNN: Convolutional Neural Network

WBCs: White Blood Cells

MLP: Multi-Layer Perceptron

RNN: Recurrent Neural Network

ViT: Vision Transformer

CBAM: Convolutional Block Attention Module

GPU: Graphics Processing Unit

F1 Score: Harmonic Mean of Precision and Recall

ReLU: Rectified Linear Unit

OCR: Optical Character Recognition

DBSCAN: Density-Based Spatial Clustering of Applications with Noise

K-Mean: K-Means Clustering Algorithm

LoRA: Low-Rank Adaptation

TPU: Tensor Processing Unit

GPGPU: General Purpose GPU

RLHF: Reinforcement Learning with Human Feedback

CoT: Chain of Thought

DPO: Direct Preference Optimization

CLIP: Contrastive Language–Image Pretraining

MoE: Mixture of Experts

SVM: Support Vector Machine

NLP: Natural Language Processing

Iou: Intersection over Union

---

## *General introduction*

---

## **General introduction**

In recent years, the integration of artificial intelligence in the medical field has witnessed substantial growth, particularly in the domain of computer-aided diagnosis (CAD). The increasing complexity and volume of medical imaging data—such as microscopic images used in hematology—has necessitated the development of more efficient, accurate, and automated diagnostic systems. Applications such as disease detection, cell classification, and image-based diagnostics require highly reliable and fast processing capabilities.

To meet these challenges, novel techniques in deep learning and image processing have emerged, enabling machines to perform medical image classification tasks with near-human accuracy. The strategic importance of these systems lies in their potential to reduce manual workload, minimize human error, and accelerate the diagnostic process—especially in resource-limited environments or high-volume laboratories.

This study focuses on the classification of white blood cells (WBCs) using deep learning techniques, particularly Convolutional Neural Networks (CNN) and attention-based mechanisms. These models offer significant improvements in recognizing the morphological characteristics of various WBC types—such as Neutrophils, Lymphocytes, Monocytes, and Eosinophils—based on microscopic images, aiding in the diagnosis of infections, leukemia, and other hematological conditions.

The dissertation is structured into three main chapters, each addressing a fundamental aspect of the research and gradually building up to the implementation and evaluation of the proposed classification system.

The first chapter introduces the theoretical foundations of CAD systems and image classification. It also covers essential concepts in digital imaging and deep learning, including CNN architecture and its relevance in medical image analysis.

The second chapter presents a comprehensive overview of neural network models, comparing CNN-based methods with attention-based architectures. It explores recent advancements in the field, such as Vision Transformers and hybrid models, and discusses their performance in medical applications.

The third chapter outlines the experimental setup and implementation process. It details the use of Google Colab, Python libraries (PyTorch, OpenCV, etc.), and a labeled WBC dataset. A custom deep learning model—based on ResNet-18 is developed, trained, and evaluated. The

chapter also includes performance metrics, visualization of classification results, and a discussion of key findings.

Finally, this dissertation provides a solid foundation for further research in the application of deep learning to hematology and medical image analysis. It highlights the advantages of using CNNs and attention mechanisms for WBC classification and proposes future perspectives to enhance model accuracy, generalization, and clinical usability.

---

## *Chapter 1*

# *Neural Network Models in Deep Learning*

---

## I. Introduction

In recent years, deep learning has revolutionized the field of artificial intelligence, especially in applications involving complex data such as images, speech, and medical diagnostics. At the heart of this transformation are neural network models—powerful computational structures that simulate the way human brains process information.

This chapter presents an in-depth overview of the most prominent neural network architectures used in deep learning, with a particular focus on Convolutional Neural Networks (CNNs) and attention-based models. These architectures have shown remarkable success in various domains, especially in image classification tasks, due to their ability to automatically learn hierarchical features from raw data.

We begin by exploring CNNs, which are specifically designed to handle spatial and visual information, making them ideal for medical image analysis. Then, we introduce attention mechanisms and Transformer-based networks, which have recently emerged as state-of-the-art approaches capable of focusing on the most informative parts of the input data. These models enhance both the accuracy and interpretability of classification systems.

By analyzing and comparing different architectures, this chapter lays the theoretical foundation necessary for understanding the design and implementation of deep learning models applied later in this work to the classification of white blood cells

- **RELATED WORK**

### I.1 CNN based:

Convolutional Neural Network (CNN) is a Multilayer Perceptron (MLP) development designed to process two-dimensional data. CNN belongs to the Deep Neural Network type due to the high network depth and is widely applied to image data. In the case of image classification, MLP is less suitable for use because it does not store spatial information from image data and considers each pixel to be an independent feature resulting in poor results.

CNN was first developed under the name NeoCognitron by Kunihiro Fukushima, a researcher from NHK Broadcasting Science Research Laboratories, Kinuta, Setagaya, Tokyo, Japan. The concept was later matured by Yann LeCun, a researcher from AT&T Bell Laboratories in Holmdel, New Jersey, USA. LeCun successfully applied the CNN model under the name LeNet to his research on recognizing numbers and handwriting. In 2012, Alex Krizhevsky, with his CNN implementation, won the ImageNet Large Scale Visual Recognition Challenge 2012 competition. This achievement is a moment of proof that deep learning methods, especially CNN. The CNN method has proven successful in outperforming other Machine Learning methods such as SVM in the case of object classification in images [1].



## I.2 Attention Based Networks:

Attention -based networks have transformed the landscape of deep learning by enabling models to dynamically focus on the most relevant parts of the input data. Since the introduction of the Transformer architecture by Vaswani et al. (2017), attention mechanisms—particularly self-attention—have become foundational to numerous breakthroughs in natural language processing, computer vision, and multimodal learning. Over the past few years (2023–2025), attention-based methods have continued to evolve, with ongoing research focused on improving efficiency, scalability, and adaptability across various domains.

Recent advancements in attention-based architectures have been central to the development of large-scale language models. For instance, **GPT-4** (OpenAI, 2023), **Claude 2 and 3** (Anthropic, 2023–2024), and **Gemini 1.5** (Google DeepMind, 2024) heavily rely on self-attention mechanisms for their exceptional performance in language understanding, reasoning, and generation. These models demonstrate the power of deep attention stacks in capturing complex dependencies within sequences and across modalities.

To address the high computational cost of self-attention in large models, several optimization strategies have emerged. Techniques like **sparse attention** (BigBird, Longformer) and **linear attention** (Performer, FlashAttention-2) have been proposed to reduce memory and computation requirements while preserving model accuracy. These innovations have enabled the training of models with longer context windows and broader input ranges, supporting tasks like long document summarization and long-form dialogue generation.

In the field of computer vision, attention mechanisms have also replaced traditional convolution operations in many state-of-the-art models. **Vision Transformers (ViTs)** and their successors (e.g., **Swin Transformer**, **MobileViT**, **Segment Anything Model**) utilize self-attention to capture global dependencies in images, improving object recognition and segmentation tasks. These architectures demonstrate that attention mechanisms are not only effective in sequential data but also in spatial contexts.

Furthermore, attention has played a crucial role in the rise of **multimodal** models that integrate visual, textual, and auditory inputs. Models like **CLIP**, **Flamingo**, and **Gemini** apply cross-attention to align different modalities, enabling robust performance in tasks like image captioning, visual question answering, and video understanding.

In summary, attention-based networks remain a driving force behind many of the recent breakthroughs in deep learning. Ongoing research continues to explore how to make attention mechanisms more efficient, interpretable, and generalizable across domains and modalities [2].

### **I.3 Deep learning:**

Deep learning has witnessed tremendous progress in recent years, becoming a cornerstone in many intelligent applications such as natural language processing, computer vision, medical diagnosis, and recommendation systems. This advancement has led to the emergence of large, powerful models capable of handling massive amounts of data and performing complex tasks with high accuracy. In this section, we review the most recent developments in the field of deep learning from 2023 to 2025, focusing on the main research trends and advanced techniques that have contributed to pushing the boundaries of performance and efficiency.

#### **I.3.1 Foundation Models and Scaling Laws**

- **GPT-4 (OpenAI, 2023):** Continued advancement in large-scale transformer models for general-purpose tasks across modalities (text, vision).
- **Gemini 1 and 1.5 (Google DeepMind, 2023–2024):** Multimodal foundation models combining vision and language at scale.
- **Claude 2 (Anthropic, 2023) and Claude 3 (2024):** Focus on interpretability and safety in large language models.

#### **I.3.2 Multimodal Learning**

- **Grok (xAI, 2024) and CLIP successors (OpenAI, 2023–2024):** Better integration of image, text, and other modalities.
- **Segment Anything Model (Meta, 2023):** A generalist vision model that can segment any object in any image—pushing general-purpose computer vision.

#### **I.3.3 Efficient and Low-Resource Models**

- **LoRA / QLoRA (2023):** Popular parameter-efficient fine-tuning techniques for LLMs.
- **Distillation & Quantization:** Widespread research on compressing large models without sacrificing much performance (TinyML, edge DL, etc.).

- **Mistral & Mixtral (2023–2024):** Mixture-of-Experts (MoE) architectures offering sparse computation with high performance.

#### **I.3.4 Reasoning and Tool Use**

- **Chain-of-Thought (CoT) and Tree-of-Thoughts (2023):** Enhancing model reasoning by mimicking human-like problem-solving.
- **ReAct / Toolformer (2023):** Models that can reason and interact with tools, environments, or search engines dynamically.

#### **I.3.5 Autonomous Agents and RLHF**

- **AutoGPT, BabyAGI (2023):** Experiments in autonomous LLM agents that plan and execute multi-step goals.
- **RLHF and DPO (Direct Preference Optimization, 2024):** New techniques for aligning models with human values more efficiently than traditional reinforcement learning.

#### **I.3.6 Diffusion and Generative Models**

- **Stable Diffusion 2 / 3, DALL·E 3, Midjourney (2023–2024):** Progress in image generation and customization.
- **Video Diffusion:** Emergence of tools like Sora (OpenAI, 2024), Pika Labs, and Runway for text-to-video.

#### **I.3.7 Applications Across Domains**

- **Medical Imaging:** Deep learning is outperforming radiologists in tasks like segmentation, diagnosis, and report generation.
- **Code Generation:** Tools like Code Llama, Codex, and StarCoder revolutionize software development.
- **Scientific Discovery:** AlphaFold 3 (2024) pushes boundaries in protein structure prediction.



## **II. CNN based methods**

### **II.1 Introduction**

CNN-based methods are approaches that rely on Convolutional Neural Networks (CNNs) to solve problems mainly in computer vision, but also in other fields such as video analysis, medical imaging, and signal processing.

Deep learning has emerged as a powerful tool for medical image analysis, particularly in hematology. Among various deep learning architectures, Convolutional Neural Networks (CNNs) have become a foundational approach for blood cell classification due to their ability to learn complex visual patterns from data. This subsection explores the evolution and application of CNN-based methods in the context of hematological image analysis.

Convolutional Neural Networks (CNNs) have become the cornerstone of image classification tasks, showing substantial success in the field of hematology. These models excel in automatically extracting hierarchical spatial features from input images, making them particularly effective for classifying various blood cell types.

Early works in blood cell classification relied on standard CNN architectures such as LeNet, AlexNet, and VGGNet. These models were adept at identifying basic morphological features of blood cells. However, their limited depth and inability to reuse features hindered their performance when applied to more complex datasets with high intra-class variability.

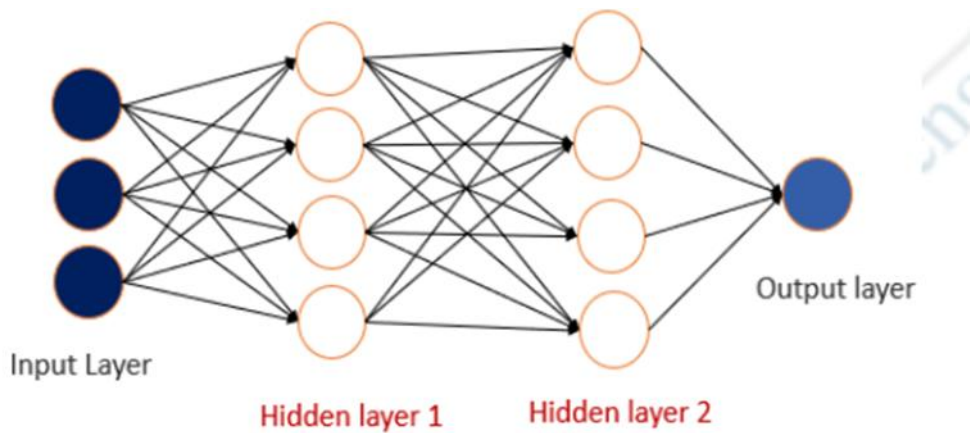
To address these limitations, deeper and more efficient models like ResNet (Residual Networks) and DenseNet (Densely Connected Networks) have been introduced. ResNet's residual connections help mitigate the vanishing gradient problem, enabling the training of deeper networks. DenseNet, on the other hand, enhances feature propagation and encourages feature reuse by connecting each layer to every other layer in a feed-forward fashion.

In addition, researchers have turned to EfficientNet, a model that balances network depth, width, and resolution to achieve high accuracy with fewer parameters. These CNN-based approaches have been applied to widely used datasets such as BCCD and BloodMNIST, achieving remarkable classification accuracy for different blood cell types, including neutrophils, eosinophils, lymphocytes, and monocytes.

Despite these successes, traditional CNNs still face challenges in focusing on the most informative regions of an image. To overcome this limitation, attention-based models have been developed. These models guide the network to focus on diagnostically relevant areas, improving feature extraction. The next subsection discusses these attention-driven approaches in more detail.

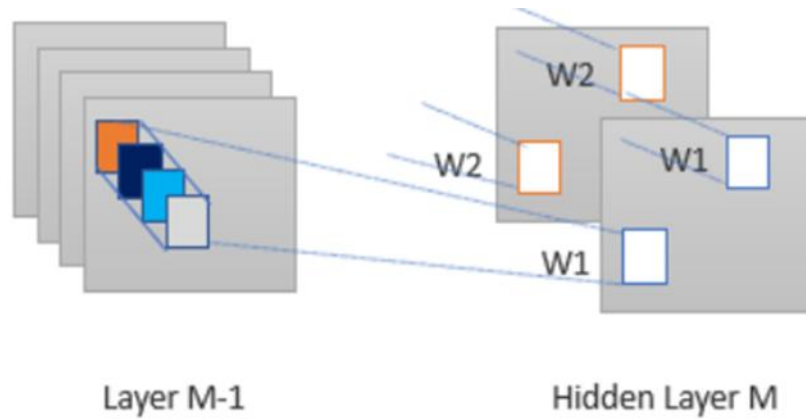
## II.2 CNN Concept

The way CNN works has similarities to MLP. However, in CNN, each neuron is presented in a two-dimensional form, unlike MLP, where each neuron is only one-dimensional in size.



**Figure I.1: Architecture Basic MLP**

An MLP, as shown in figure 1. It has an I layer, with each layer containing  $J_i$  neurons. MLP accepts one-dimensional data inputs and propagates that data on the network until it produces an output. Each connection between neurons on two contiguous layers has a one-dimensional weight parameter that determines the quality of the mode. In each input data on the layer, a linear operation is carried out with the existing weight value. Then the computational results will be transformed using a non-linear operation called an activation function. On CNN, the data propagated on the network is two-dimensional, so the linear operation and weighting parameters on CNN are different. On CNN, linear operations use convolution operations, while weights are no longer one-dimensional only. However, they are four-dimensional, which is a collection of convolution kernels, as shown in figure 2. Due to the nature of the convolution process, then CNN can only be used on data that has a two-dimensional structure, such as imagery and sound [3].



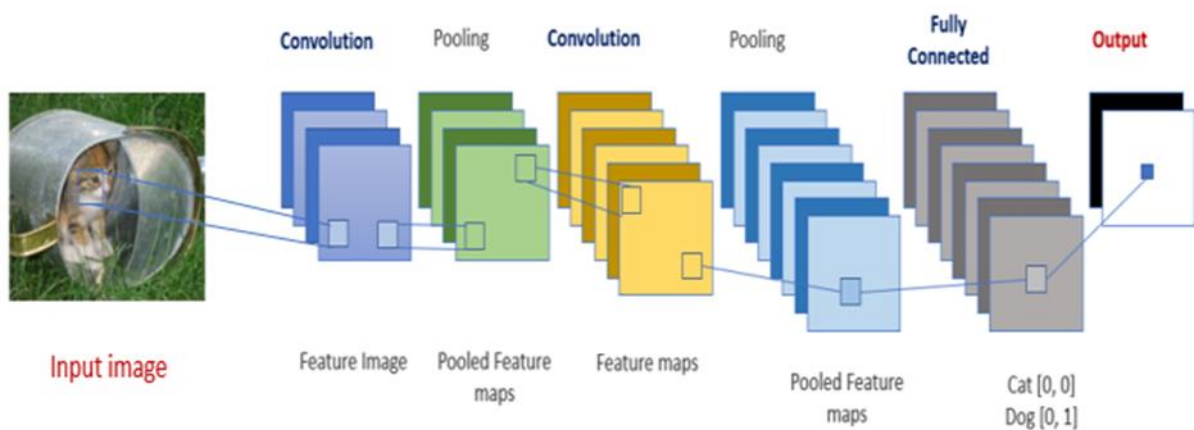
**Figure I.2: Process convolution CNN [3]**

## MATERIAL METHOD

We introduced the Ensemble Convolution Neural Network Architecture (CNN), consisting of CNN-based object classification networks and object detection to train datasets taken from two types of animals.

Experimental analysis showed better results than

Places-CNNs. The convolution layer pooling acts as a feature extractor of the input image, while the layer that is fully layered acts as a classification



**Figure I.3: Convolution Layer [3]**

The convolution layer pooling acts as a feature extractor of the input image, while the layer that is fully layered acts as a classification.

In the image above, when receiving the desired image as input, the network correctly gives the highest probability for it (0.94) among the four categories. The sum of all probabilities in the output layer should be one. There are four main operations in ConvNet shown Convolution, Non-Linearity (ReLU), Pooling or Sub Sampling, and Classification (Fully Connected Layer) [3]

### Convolutional Neural Networks – Architecture

The layer is used to build Convolutional Neural Networks. Simple ConvNet is a sequence of layers, and each ConvNet layer converts one activation volume to another volume through the Yang function can be distinguished. We used four main layer types to build the ConvNet architect. [3]

#### Convolution Layer

It holds the raw pixel value of the training image as input. In the above example, the image (cat) with a width of 32, a height of 32, and three color channels, R, G, and B, are used. It passes the forward propagation step and finds the probability of output for the class setup. Lapsang ensures spatial connections between pixels by interpreting image features using a small box of input data. Let's assume the output probability for the image above is [0.2, 0.1, 0.3, 0.4]. The size of the feature map is controlled by three parameters [3]

- Depth – The number of filters used for coevolutionary operations
- Stride – The number of filters used to filter the matrix above the input matrix.
- Padding – very good for entering matrices with zeros around the limit matrix.

Calculates the total error on the output layer with the summation of all 4 classes.

$$\text{Total Error} = \sum 1/2 (\text{Target probability} - \text{Output probability})^2$$

Calculated the output of neurons that are connected to the local area inputted. It can produce volumes such as  $[32 \times 32 \times 16]$  for 16 filters.

#### Rectified Linear Unit (ReLU) Layer

A non-linear operation. This layer implements the function of element-wise activation. ReLU is used after every convolution operation. It is applied per pixel and replaces all negative pixel



values in the feature map with zeros. This makes the measure n volume unchanged ([32x32x16]), where ReLu is a non-linear operation [3].

$$K(m, n) \cdot (n + m, j + X(i) \sum_{n=0}^{N-1} \sum_{m=0}^{M-1} = Y(i, j)$$

**Y (i,j):** Output pixel value at position (i, j) in the **feature map**

**X:** Input matrix (original image or previous layer's output)

**K:** Convolution kernel (filter)

**M, N:** Dimensions of the kernel (e.g., 3x3)

$\sum$ : Summation of element-wise multiplications between the kernel and the image region

➤ **Activation Function (ReLU):**

$$f(x) = \max(0, x)$$

Adds non-linearity to the network.

➤ **Self-Attention Mechanism :**

**Q:** Query

**K:** Key

**V:** Value

**dk:** Dimension of the key (for scaling)

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

❖ This equation allows the model to **focus on the most important parts** of the input.

**Famous CNN-based Methods [4]:**

CNN Method	Year	Key Contribution	Reference
<b>LeNet-5</b>	1998	Pioneered CNNs for digit recognition; introduced convolution and pooling layers.	LeCun et al., <i>IEEE</i> , 1998

<b>AlexNet</b>	2012	Revived deep learning; introduced ReLU activation and dropout; won ImageNet 2012.	Krizhevsky et al., <i>NIPS</i> , 2012
<b>ZFNet</b>	2013	Improved upon AlexNet with better visualization and parameter tuning.	Zeiler & Fergus, <i>ECCV</i> , 2014
<b>VGGNet</b>	2014	Utilized deep networks with small 3×3 filters; emphasized depth for performance.	Simonyan & Zisserman, <i>ICLR</i> , 2015
<b>GoogLeNet (Inception v1)</b>	2014	Introduced Inception modules for multi-scale processing; reduced parameters.	Szegedy et al., <i>CVPR</i> , 2015
<b>ResNet</b>	2015	Introduced residual connections to train very deep networks (up to 152 layers).	He et al., <i>CVPR</i> , 2016
<b>DenseNet</b>	2017	Connected each layer to every other layer to improve information flow.	Huang et al., <i>CVPR</i> , 2017
<b>MobileNet</b>	2017	Designed for mobile and embedded vision applications; used depthwise separable convolutions.	Howard et al., <i>arXiv</i> , 2017
<b>EfficientNet</b>	2019	Achieved better accuracy and efficiency by scaling depth, width, and resolution.	Tan & Le, <i>ICML</i> , 2019
<b>ConvNeXt</b>	2022	Modernized CNN architecture with design choices from Transformers.	Liu et al., <i>CVPR</i> , 2022

#### Typical Applications of CNN-based Methods:

- **Image Classification** (e.g. recognizing WBCs or BBCs)
- **Object Detection** (e.g. detecting cars in an image)

- **Image Segmentation** (e.g. isolating organs in medical images)
- **Medical Image Analysis** (e.g. MRI, X-rays)
- **Autonomous Driving** (real-time road analysis)

## II.2 Attention based methods

### Introduction:

Attention-based methods have revolutionized the field of computer vision by enabling models to dynamically prioritize the most relevant parts of an image. Inspired by human visual attention, these mechanisms guide neural networks to focus on informative regions while suppressing irrelevant or redundant data. In the context of medical imaging—particularly blood cell classification—attention mechanisms have proven especially effective, as they help highlight subtle morphological features critical for accurate diagnosis. By integrating attention into convolutional or transformer-based architectures, these models can not only improve classification performance but also enhance interpretability, making them valuable tools in clinical decision support systems.

Attention mechanisms have emerged as a powerful enhancement to deep learning models, particularly in tasks involving complex visual patterns such as blood cell classification. These methods allow models to selectively focus on the most informative regions of an image, improving both accuracy and interpretability.

Attention mechanisms have been widely integrated into Convolutional Neural Networks (CNNs) and Transformer-based architectures to emphasize features crucial for distinguishing between different types of blood cells, such as erythrocytes, leukocytes, and thrombocytes. Techniques such as **spatial attention** focus on “where” the model should look, while **channel attention** emphasizes “what” features are most important.

A notable application is the **Convolutional Block Attention Module (CBAM)**, which integrates both spatial and channel attention to refine feature maps in CNNs. This has shown improvements in sensitivity and specificity, particularly in detecting rare or morphologically ambiguous cell types.

Transformers, originally developed for natural language processing, have recently been adapted for medical image analysis. **Vision Transformers (ViTs)** and hybrid models that combine CNNs with attention layers have demonstrated strong performance in blood cell classification tasks, outperforming traditional CNNs in many cases.

Furthermore, attention-based models often provide **visual explanations** via attention maps, which can be useful for pathologists to understand model decisions and trust automated systems.

Despite their promise, these methods require substantial computational resources and large annotated datasets to train effectively. However, with the growing availability of labeled microscopy datasets and pre-trained models, attention-based techniques are becoming increasingly accessible for medical image classification tasks.

#### **Famous Attention based Methods [5]:**

<b>Model</b>	<b>Year</b>	<b>Key Contribution</b>	<b>Reference</b>
<b>Bahdanau Attention</b>	<b>2014</b>	<b>Introduced attention mechanism in neural machine translation, allowing models to focus on relevant parts of the input sequence.</b>	<b>Bahdanau et al., 2014</b>
<b>Transformer</b>	<b>2017</b>	<b>Proposed a novel architecture relying entirely on self-attention mechanisms, removing recurrence and convolutions.</b>	<b>Vaswani et al., 2017</b>
<b>Graph Attention Networks (GAT)</b>	<b>2017</b>	<b>Applied attention mechanisms to graph-structured data, enabling nodes to attend over their neighborhoods' features.</b>	<b>Veličković et al., 2017</b>
<b>Residual Attention Network</b>	<b>2017</b>	<b>Integrated attention mechanisms into deep residual networks for image classification, enhancing feature representation.</b>	<b>Wang et al., 2017</b>
<b>Transformer-XL</b>	<b>2019</b>	<b>Extended the Transformer model to capture longer-term dependencies by introducing recurrence mechanisms.</b>	<b>Dai et al., 2019</b>

<b>BERT</b> (Bidirectional Encoder Representations from Transformers)	<b>2018</b>	Utilized bidirectional training of Transformers for language understanding tasks, achieving state-of-the-art results.	Devlin et al., 2018
<b>XLNet</b>	<b>2019</b>	Combined autoregressive and autoencoding approaches to pretraining, improving upon BERT's performance.	Yang et al., 2019
<b>Vision Transformer (ViT)</b>	<b>2020</b>	Applied Transformer architecture to image classification tasks, demonstrating competitive performance with CNNs.	Dosovitskiy et al., 2020
<b>Perceiver</b>	<b>2021</b>	Introduced a model capable of handling arbitrary input modalities using attention mechanisms, enabling scalability.	Jaegle et al., 2021
<b>Perceiver IO</b>	<b>2021</b>	Extended the Perceiver model to handle arbitrary input and output modalities, enhancing flexibility.	Jaegle et al., 2021

### Convolutional Attention – General Architecture

Here's a high-level view of how convolutional attention can be structured:

#### 1. Input Embedding / Feature Extraction

- Input: Image, text, or sequence.
- Pass through **CNN layers** or **embedding layers** to get a feature map or token embeddings.

#### 2. Convolutional Attention Block

Each block often consists of:

- **Convolutional Layer**
  - Extracts local spatial features.
  - Might use depthwise or dilated convolutions for efficiency.
- **Attention Layer**
  - Computes self-attention or cross-attention.
  - Queries, Keys, and Values are derived from the input features.
- **Fusion Mechanism**
  - Combine convolution and attention outputs.
  - Can be additive, concatenated, or gated.

### 3. Stacking Blocks

- Multiple convolutional attention blocks are stacked to learn hierarchical features.
- Often interleaved with normalization and feed-forward layers.

### 4. Output Head

- Task-specific head (classification, segmentation, etc.).
- Usually consists of MLPs or global pooling followed by softmax/sigmoid.

### Limitations and Challenges

Despite their advantages, attention-based methods come with challenges. They often require large annotated datasets, which are scarce in medical domains. Moreover, ViTs and hybrid models can be computationally intensive, necessitating high-end GPUs and longer training times. However, with the rise of transfer learning and pre-trained models on large-scale biomedical image datasets, these limitations are gradually being addressed.

**Conclusion**

Attention-based methods have demonstrated significant potential in improving the accuracy and interpretability of blood cell classification models. By allowing models to focus on the most relevant features, these techniques enhance performance in tasks such as identifying rare cell types or distinguishing between similar morphologies. Models like CBAM, Vision Transformers, and hybrid CNN-ViT architectures are pushing the boundaries of what is possible in automated hematology, offering new opportunities for faster, more accurate diagnostic tools. However, challenges such as the need for large annotated datasets and high computational demands remain. Despite these obstacles, attention-based approaches are increasingly shaping the future of medical image analysis, bringing us closer to reliable, AI-powered systems in clinical practice.



---

## *Chapter 2*

# *Classification in Computer Vision*

---

## **II.1 Introduction**

In the modern digital age, computer vision has emerged as a crucial domain within artificial intelligence, allowing machines to analyze and comprehend visual data from our environment. A core task in computer vision is image classification, which involves automatically assigning a label or category to an image based on its visual elements. By leveraging advanced algorithms alongside extensive datasets, computer vision systems can now identify objects, scenes, and intricate patterns with impressive accuracy. Image classification forms the basis for numerous practical applications, including facial recognition, medical diagnosis, self-driving cars, and content categorization. As technology progresses, the incorporation of advanced classification models continues to enhance the potential of computer vision, extending the limits of what machines can observe and understand.

## **II.2 Definition of Computer Vision**

Computer vision, a pivotal branch of artificial intelligence, has emerged as a transformative technology in the digital age, enabling machines to perceive, interpret, and analyze visual data with remarkable accuracy. Built on the foundations of mathematics, computer science, and engineering, computer vision strives to replicate human visual understanding by teaching machines to extract meaningful insights from images and videos. Over the past decade, rapid technological advancements, the proliferation of large-scale datasets, and increased computational capabilities have propelled computer vision from theoretical frameworks to widespread practical applications, making it integral across diverse industries.

Initially rooted in basic image processing techniques centered on pixel-level operations, the field evolved through the incorporation of pattern recognition and feature extraction. Early systems relied on handcrafted features such as edge detection and histogram analysis, which offered limited performance in dynamic environments. The introduction of deep learning—particularly convolutional neural networks (CNNs)—marked a major turning point. These models enabled automatic feature learning and end-to-end training, significantly improving accuracy and expanding the reach of vision systems to include tasks like object detection, semantic segmentation, and facial recognition.

Computer vision's real-world impact is extensive. In healthcare, it has revolutionized diagnostics through precise analysis of medical images such as X-rays, MRIs, and CT scans. In

transportation, vision-based systems power autonomous vehicles, promising safer and more efficient travel. Retail and e-commerce benefit from enhanced customer experiences and personalized recommendations via image recognition, while surveillance systems utilize real-time visual analysis for monitoring and threat detection, bolstering public safety.

Despite these advances, several challenges persist. Data bias remains a critical issue, with models sometimes reflecting and amplifying societal prejudices present in training datasets. The interpretability of deep learning models also poses concerns, particularly in sensitive fields like healthcare and law enforcement. Additionally, the high computational demands of vision systems can hinder their deployment in real-time or resource-limited scenarios.

Emerging trends offer promising solutions to these challenges. The integration of computer vision with other AI domains, such as natural language processing and reinforcement learning, has led to the development of multimodal systems capable of richer data understanding. Edge computing aims to reduce latency and energy consumption by bringing computation closer to the data source. Generative adversarial networks (GANs) contribute by generating high-quality synthetic data to support model training. Simultaneously, advancements in hardware—such as GPUs, TPUs, and neuromorphic chips inspired by the human brain—are enhancing the speed and efficiency of vision systems.

The societal impact of computer vision extends far beyond technology. In environmental monitoring, it aids in analyzing satellite imagery for deforestation tracking, wildlife conservation, and disaster response. In education, it supports immersive learning experiences through augmented reality. While the potential for positive change is vast, it is accompanied by the responsibility to ensure ethical and equitable development of these technologies.

This paper aims to present a thorough overview of the evolution and current state of computer vision, highlighting the key techniques, impactful applications, and emerging trends shaping its future. By examining both the accomplishments and limitations of the field, this study seeks to guide researchers, practitioners, and policymakers in leveraging computer vision responsibly to address real-world challenges and advance intelligent systems.

### **II.2.1 Types of Computer Vision Tasks [6]**

#### **Image Classification**

- Example: Recognizing whether an image contains a WBCs, BBCs, or another object.
- Goal: Assign a label or class to an entire image.

#### **Object Detection**

- Example: Detecting and locating multiple objects in an image, such as finding Import point in Image, pedestrians, and bicycles in a street scene.
- Goal: Identify objects and draw bounding boxes around them.

#### **Semantic Segmentation**

- Example: Classifying each pixel in an image into predefined categories like roads, buildings, or sky.
- Goal: Label each pixel of the image with a class.

#### **Instance Segmentation**

- Example: Identifying and segmenting individual objects, such as separating different dogs in a group photo.
- Goal: Identify and segment each object in the image, even if the objects belong to the same category.

#### **Optical Character Recognition (OCR)**

- Example: Extracting text from an image of a document or a street sign.
- Goal: Convert images of text into machine-readable text.

#### **Facial Recognition**

- Example: Identifying or verifying a person's identity based on their face, such as in security systems.

- Goal: Recognize or verify individuals based on facial features.

**Pose Estimation**

- Example: Determining the position of a person's body in an image, such as tracking the angle of limbs during a workout.
- Goal: Detect and track human body parts in images or video streams.

**Action Recognition**

- Example: Recognizing specific actions such as running, jumping, or waving in a video.
- Goal: Identify actions or behaviors in videos or image sequences.

**Depth Estimation**

- Example: Creating a 3D representation of a scene by estimating the distance of objects from the camera.
- Goal: Predict depth information from a single image or multiple views.

**Image Generation**

- Example: Creating new images based on text descriptions (e.g., "A sunset over the mountains").

**II.3 Applications of Computer Vision in Medicine**

Applications of Computer Vision Computer vision is being utilized in a wide range of industries, revolutionizing workflows and creating new opportunities:

**Healthcare:** Vision-driven systems have transformed diagnostics by accurately identifying anomalies in medical imaging. For example, computer vision supports the detection of tumors in radiology images, forecasts disease advances, and aids in robotic surgical procedures.

**Autonomous Vehicles:** Self-driving automobiles utilize computer vision for functions such as detecting lanes, recognizing traffic signs, and avoiding obstacles. These systems depend on real-time analysis of visual data to navigate safely in ever-changing environments.

**Retail and E-commerce:** Vision technologies enable functionalities like virtual try-ons, tailored recommendations, and inventory oversight. For instance, virtual mirrors allow shoppers to visualize clothing items without physical fitting, enhancing the shopping experience.

**Surveillance and Security:** Vision-enabled surveillance solutions offer live monitoring, intrusion detection, and facial recognition for access security. These technologies are vital for maintaining public safety and preventing unauthorized activities.

**Environmental Monitoring:** Analyzing satellite images allows tracking of deforestation, urban development, and natural disasters. Computer vision assists researchers and policymakers in making informed decisions to tackle environmental issues. Challenges in Computer Vision: Despite the progress in computer vision, various challenges remain:

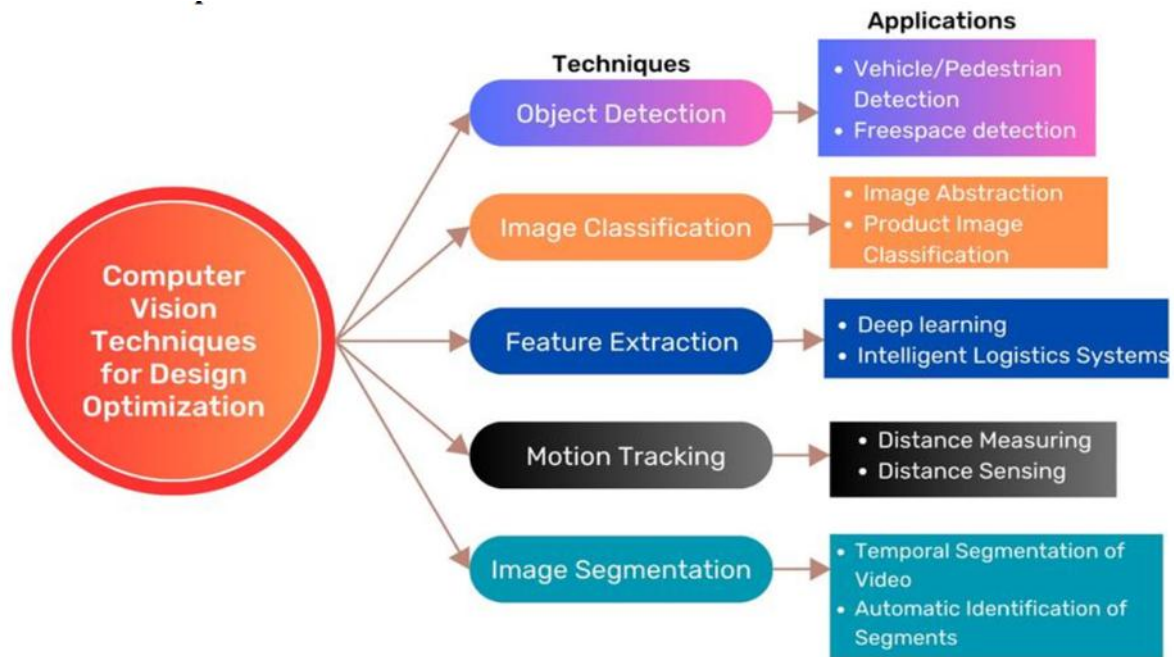
**Data Bias:** Vision models trained on skewed datasets can produce biased results, raising ethical questions. It is crucial to ensure diversity and representativeness in training data to alleviate this issue.

**Interpretability:** The opaque nature of deep learning models creates obstacles in comprehending their decision-making mechanisms. Developing transparent and interpretable models is essential for building trust and accountability.

**Computational Cost:** The process of training and implementing vision models demands considerable computational resources, which limits their use in applications with restricted resources. Emerging Trends and Future Directions: The evolution of computer vision is influenced by new trends that aim to overcome existing challenges and explore fresh possibilities:

- **Multimodal Learning:** Combining computer vision with other AI domains, such as NLP, allows systems to grasp intricate data relationships. For example, vision-language models like CLIP and DALL-E integrate visual and textual input to generate innovative outputs [7].

### II.3.1 Trends on Computer vision [8]



**Figure II.1: illustration of Computer Vision for Design Optimization**

#### II.1 Comment of figure

This diagram displays the essential techniques in computer vision along with their applications that help in optimizing design. The techniques encompass object detection, image classification, feature extraction, motion tracking, and image segmentation. Each of these methods aids in various applications such as detecting pedestrians and vehicles, classifying product images, supporting deep learning systems, measuring distance, and segmenting video. The chart emphasizes the importance of these computer vision techniques in improving the design of intelligent systems and automation.

### II.3.2 Computer Vision in Healthcare

Computer vision (CV) has transformed the healthcare sector by facilitating the automation of intricate tasks like analyzing medical images, making diagnoses, and monitoring patients. A key application of CV in healthcare is within medical imaging, where deep learning models, especially convolutional neural networks (CNNs), have shown exceptional success. For instance, Esteva et al. (2017) revealed that a deep learning model could categorize skin cancer images with accuracy that rivals that of dermatologists. Likewise, Litjens et al. (2017) performed an extensive review of deep learning methods in medical image analysis,

emphasizing their effectiveness in areas such as tumor identification, organ segmentation, and disease classification. In addition to diagnostics, CV has found applications in patient monitoring and surgical support. For example, CNNs have been utilized to interpret video feeds from surgical rooms to provide real-time assistance to surgeons (Twinanda et al., 2017). Furthermore, CV methods have been applied in remote patient monitoring, allowing for the detection of irregularities in patient behaviors or vital signs through visual information (Ravi et al., 2020). These developments highlight the revolutionary potential of CV in enhancing healthcare results. However, despite the notable advancements in medical imaging and diagnostics, the use of CV in analyzing financial data within healthcare has yet to be thoroughly explored. This gap offers a distinct opportunity to leverage the advantages of CV in the financial dimensions of healthcare systems [9].

## II.4 Image Classification Overview

Image classification involves labeling an image by analyzing its visual elements. By employing machine learning and deep learning methods, particularly convolutional neural networks (CNNs), models become adept at detecting patterns, shapes, and characteristics within images. The objective is to correctly classify images into established categories, such as determining if an image depicts a cat, a dog, or a car. This technique finds extensive application in fields like healthcare, security, and autonomous vehicles.

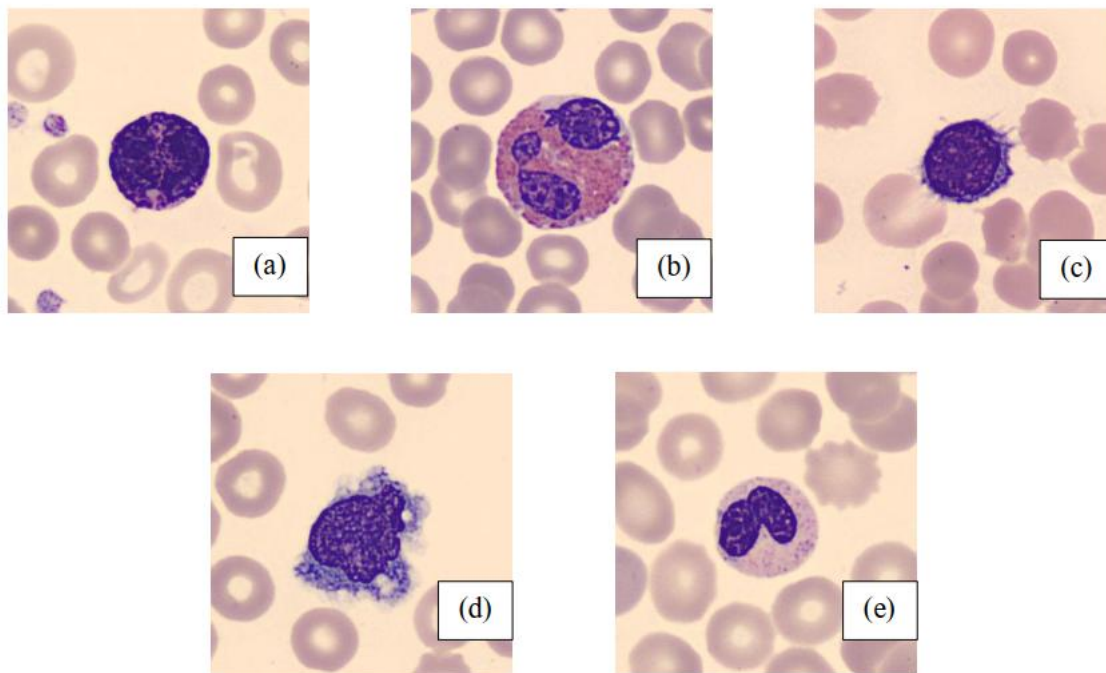
### II.4.1 what is classification?

Image classification entails training a model to recognize patterns and features within images to categorize them into predefined classes. For instance, a model might be trained to distinguish between images of cats and dogs. This task is typically approached using supervised learning, where the model learns from labeled examples. In this setup, each training image is associated with a label indicating its class, allowing the model to learn the mapping between visual features and categories .

**Importance of White Blood Cell Classification in the Medical Field:** White blood cells play a crucial role in the immune system, and their classification is vital for diagnosing and treating various medical conditions, including infections, leukemia, and autoimmune diseases. Accurate classification helps in monitoring disease progression and assessing the effectiveness of treatments.



White blood cells (WBCs) play a vital role in the immune system and are classified into five primary types: Neutrophils, Lymphocytes, Eosinophils, Basophil and Monocytes. Each of these cells serves a specific function in defending the body against infections, diseases, and foreign invaders. Neutrophils are the first line of defense against bacterial infections, making up a large portion of the WBC count. Lymphocytes are essential for the adaptive immune response, helping the body recognize and respond to pathogens. Eosinophils specialize in fighting off parasitic infections and managing allergic reactions. Monocytes are important for phagocytosis, where they engulf and digest pathogens and debris [10].



**Figure II.2: The 5 cell types of the RAW – PBC DIB dataset a) Basophil b) Eosinophil c) Lymphocyte d) Monocyte e) Neutrophil[10]**

### II.4.2 How It Works?

Image classification typically employs machine learning, especially deep learning techniques like Convolutional Neural Networks (CNNs). The general workflow includes:

#### 1. Data Collection:

- **Training Data:** The first step is to collect labeled images. Each image is associated with a class label (e.g., “cell” etc.).

- **Testing Data:** After the model is trained, it needs to be evaluated on unseen data, which is usually a separate dataset from the training set.

## 2. Preprocessing:

- **Resizing:** Images are often resized to a fixed dimension because neural networks require a consistent input size.
- **Normalization:** Pixel values are usually normalized to a range (e.g., 0 to 1 or -1 to 1) for faster convergence during training.
- **Augmentation:** To increase the model's ability to generalize, techniques like rotation, flipping, and zooming can be applied to the training data.

## 3. Feature Extraction:

- Traditional image classification used **manual feature extraction** (like edge detection, color histograms, etc.).
- With **Deep Learning**, particularly **Convolutional Neural Networks (CNNs)**, features like textures, patterns, and shapes are automatically learned from the raw pixel data.
- A CNN typically contains layers that perform:
  - **Convolutional layers:** Detect patterns or features like edges and textures.
  - **Pooling layers:** Reduce the dimensionality of data, keeping the most significant information.
  - **Fully connected layers:** Connect all neurons and are used for final classification.

## 4. Model Training:

- **Training Phase:** A neural network (or other algorithms) is trained on labeled data using optimization techniques like **Stochastic Gradient Descent (SGD)**.
- The model adjusts its weights to minimize the **loss function**, which measures how far off its predictions are from the actual labels.

- **Epochs:** The process is repeated over multiple iterations (epochs) to improve accuracy.

#### 5. Classification (Prediction Phase):

- After the model is trained, you can use it to classify new, unseen images.
- For each new image, the model predicts a class label by passing the image through the learned network layers.
- **Softmax** is often used in the output layer to assign probabilities to each class label.

#### 6. Post-Processing (Optional):

- **Thresholding:** If the model predicts multiple classes with close probabilities, you can apply thresholding techniques to select the most probable class.
- **Ensemble Methods:** Combining the outputs of multiple models to improve accuracy.

#### 7. Metrics Evaluation:

- To assess the model's performance, metrics like **accuracy**, **precision**, **recall**, **F1-score**, and **confusion matrix** are used.

Five equal-sized subsets are created from the dataset. In this procedure, the remaining four subsets are used to train the model, and each subset is successively identified as the test dataset. For instance, subsets 2, 3, 4, and 5 are utilized for training in the first iteration, whereas subset 1 is used for testing. Subsets 1, 3, 4, and 5 are then used for training, while subset 2 is assigned for testing. Every subset is utilized as the test dataset once during the five iterations of this method.

The test performance computed in each cycle is averaged to assess the outcomes. By randomly dividing the dataset into training and test subsets, this technique seeks to reduce chance-related biases or errors. As a result, it offers a more trustworthy and comprehensive evaluation of model performance. At the conclusion of every iteration, the model performance measures (such as accuracy, precision, and F1 score) are noted. To assess the model's overall performance, the average of these outcomes is computed.

As a result, the likelihood of haphazard successes or mistakes during the data splitting process is reduced.

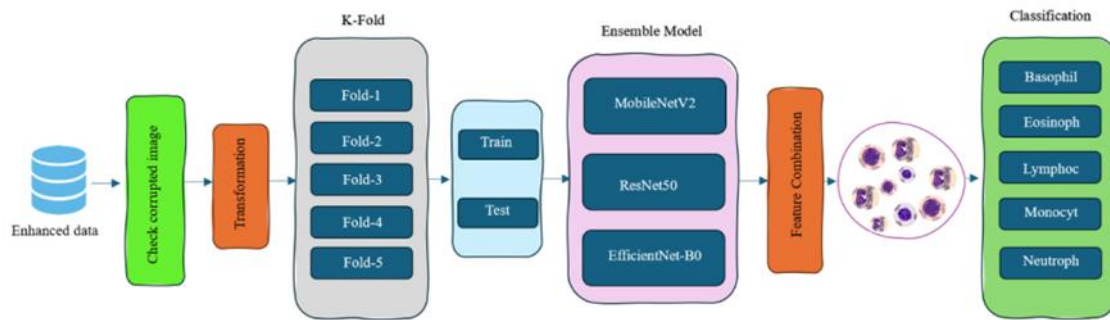


Figure II.3: structure of the created Ensemble model [10]

### II.4.3 Types of Image Classification

- **Single-label Classification:** Each image is assigned one label from a set of categories.
- **Multi-label Classification:** An image can be associated with multiple labels simultaneously (e.g., an image labeled as both "dog" and "park").
- **Supervised Learning:** Models are trained on labeled datasets, learning to map inputs to known outputs.
- **Unsupervised Learning:** Models identify patterns and groupings in unlabeled data, often used for clustering similar images.

### II.4.4 Real-World Applications

- **Healthcare:** Classifying medical images to assist in diagnosis. [kili-website\\_arXiv](#)
- **Agriculture:** Identifying plant diseases from leaf images.
- **Retail:** Organizing and tagging product images for inventory management.
- **Security:** Facial recognition systems for access control.
- **Autonomous Vehicles:** Recognizing traffic signs and obstacles.

#### II.4.5 Classification of White Blood Cells [11]:

Type	Subtype	Function
<b>Granulocytes</b>	<b>Neutrophils</b>	Engulf and destroy bacteria and fungi; first responders to microbial infection.
	<b>Eosinophils</b>	Combat multicellular parasites and certain infections; involved in allergic reactions.
<b>Agranulocytes</b>	<b>Lymphocytes</b>	Engulf and destroy bacteria and fungi; first responders to microbial infection.
	<b>Monocytes</b>	Combat multicellular parasites and certain infections; involved in allergic reactions.

#### II.4.6 Classification in medicines and biomedicines

##### Classification using images

In the field of medicine and biomedicine, accurate classification is crucial for drug discovery, disease diagnosis, and personalized treatment. With the rise of artificial intelligence and machine learning, **image-based classification methods** have become increasingly valuable. These methods utilize visual data such as microscopic images, radiographic scans, and molecular structures to identify, categorize, and predict the behavior of pharmaceutical compounds and biological agents.

Image classification plays a vital role in various stages of biomedical research and pharmaceutical development, including:

- **Identifying cell types and disease states** through histological images.
- **Classifying drug compounds** based on molecular structures.
- **Detecting abnormalities** in imaging data (e.g., tumors in MRI scans).
- **Monitoring drug effects** on cellular structures using microscopy.

By leveraging deep learning techniques like convolutional neural networks (CNNs), researchers can automate and improve the accuracy of these classifications, leading to faster, more reliable outcomes in both clinical and research settings

## **II.5 Conclusion**

Computer vision, combined with image classification, is essential for allowing machines to interpret and understand visual information. By utilizing sophisticated algorithms and deep learning methods, image classification enables computer vision systems to effectively identify and classify visual data in a wide range of applications. As technology progresses, the collaboration between computer vision and image classification will promote additional innovations, enhancing machines' abilities to decode intricate visual data and revolutionizing industries around the globe.

---

## *Chapter 3:*

### *Simulations and Discussion of results*

---

### III.1 Introduction

The classification of blood cell images plays a crucial role in medical diagnostics, particularly in the detection and monitoring of various diseases such as leukemia, anemia, and infections. With the increasing availability of medical imaging data, automated image analysis has become an essential tool in improving diagnostic accuracy and efficiency.

In this project, we apply clustering techniques to classify blood cell images using Python. Clustering, an unsupervised machine learning method, helps to group similar images based on visual features without prior labeling. This approach can uncover hidden patterns in the data and assist in preliminary categorization before expert analysis.

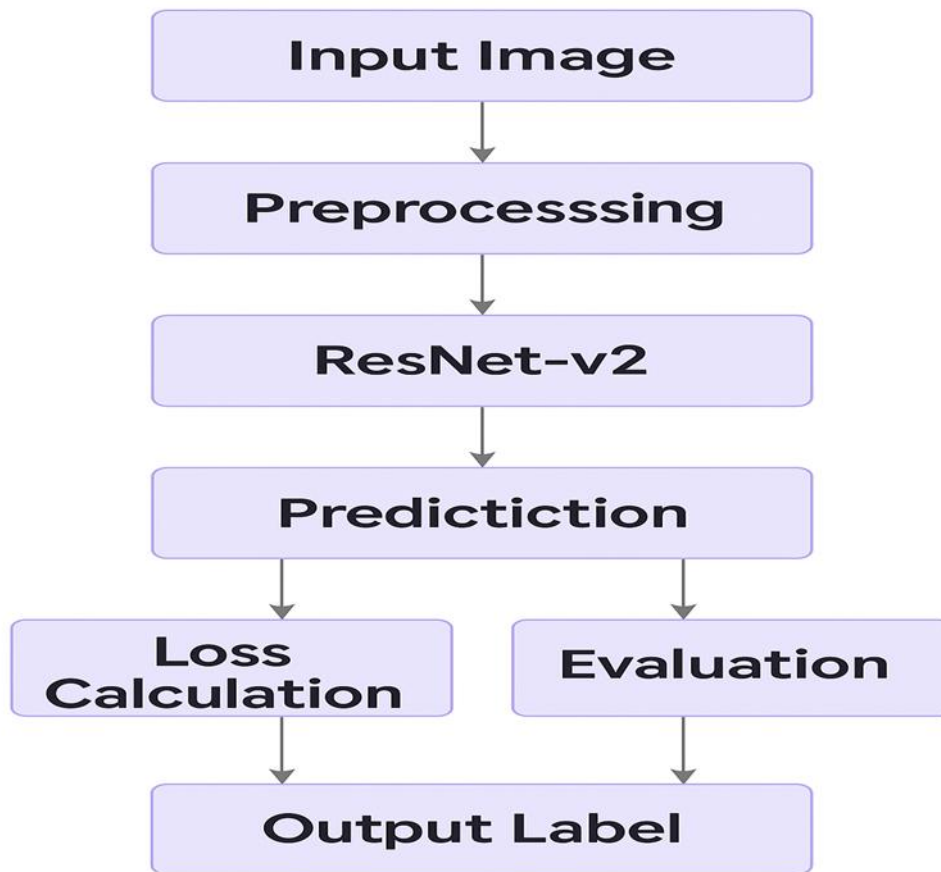
The project involves preprocessing the blood cell images, extracting relevant features, and applying clustering algorithms such as K-Means and DBSCAN. We evaluate the performance of these methods and visualize the results to understand the clustering quality and its potential for aiding medical diagnosis.

### III.2 Used Materials

I used the VSC lab's desktop computer, windows 11, Google Colab, OpenCV library, PyTorch library, Matplotlib, Pandas and NumPy, tqdm, ResNet-V2, White Blood Cell Dataset

➤ **Illustrative diagram of the simulation stages**



**Diagram Explanation:**

- This diagram illustrates the stages of an image classification system using the ResNet-v2 model. It starts with an input image, followed by a preprocessing step to prepare the data. The image is then passed through the ResNet-v2 model for feature extraction and prediction. Afterward, the loss is calculated and the results are evaluated to assess model accuracy. Finally, the system outputs the predicted label for the input image.

**III.3 Downloading and Extracting Files from Google Drive with Python**

```
import gdown

# Remove '/edit' from your link to get the file ID
file_id = "1gknVrSs1CRy8PoIh1HXiGu-10bH3cQ9S"
gdown.download(f"https://drive.google.com/uc?id={file_id}", "LISC.rar", quiet=False)
!unrar x LISC.rar
```

```
Extracting  LISC/Train_Aug/neut_9-1__8.bmp      OK
Extracting  LISC/Train_Aug/neut_9-1__9.bmp      OK
Extracting  LISC/Train_Aug.json                 OK
All OK
```

**Comment:** This code example illustrates how to fetch a .rar file from Google Drive employing the gdown library in Python. It formulates the direct download link using a shared file ID and saves the file under the name LISC.rar. The last command employs the unrar utility to extract the contents of the archive. This method is frequently utilized for automating the collection of datasets in research endeavors.

### III.4 Python Module Imports for Data Preparation

```
import os
import random
import shutil
from sklearn.model_selection import train_test_split
```

**Comment:** This figure displays a Python script segment that imports essential libraries used in data preprocessing and dataset splitting. The os and shutil modules are used for file and directory management, random is used for generating reproducible data splits, and train\_test\_split from sklearn.model\_selection is used to partition the dataset into training and testing sets—an essential step in machine learning workflows.

#### Import Os

The os module offers functionalities for engaging with the operating system. It enables you to carry out activities like exploring the file system, as well as creating or removing directories, and managing file paths.

#### Import random

The random module is utilized to produce random numbers or to make arbitrary selections. It is frequently employed in machine learning to mix data or to create random divisions for training and evaluation.

#### Import Shutil

The shutil module provides elevated functionalities for managing files and groups of files. It is often utilized for tasks such as copying, relocating, or removing files and folders.

#### Import train\_test\_split from sklearn.model\_selection

The `train_test_split` function from `sklearn.model_selection` is utilized to divide datasets into training and testing subsets. This aids in assessing machine learning models by distinguishing the data designated for training from that intended for testing.

### **III.5 Python Script for Dividing a Cell Image Dataset into Training and Testing Sets**

```
# Original folder containing the images (assuming it's in the same directory as this script)
original_folder = '/content/data/LISCCropped'

# Output folder structure
output_folder = '/content/data/LISC_dataset/'
train_folder = os.path.join(output_folder, 'train')
test_folder = os.path.join(output_folder, 'test')

# Create the output directories if they don't exist
os.makedirs(train_folder, exist_ok=True)
os.makedirs(test_folder, exist_ok=True)

# List all the subfolders (cell types) in the original folder
cell_types = ['Basophil', 'Eosinophil', 'Lymphocyte', 'Monocyte', 'Neutrophil']

# Process each cell type
for cell_type in cell_types:
    # Create corresponding subfolders in train and test
    os.makedirs(os.path.join(train_folder, cell_type), exist_ok=True)
    os.makedirs(os.path.join(test_folder, cell_type), exist_ok=True)

    # Get all image files for this cell type
    cell_type_path = os.path.join(original_folder, cell_type)
    if not os.path.exists(cell_type_path):
        print(f"Warning: Folder {cell_type_path} does not exist. Skipping.")
        continue

    images = [f for f in os.listdir(cell_type_path) if os.path.isfile(os.path.join(cell_type_path, f))]

    # Split into train and test (80% train, 20% test)
    train_files, test_files = train_test_split(images, test_size=0.2, random_state=42)

    # Copy files to train folder
    for file in train_files:
        src = os.path.join(cell_type_path, file)
        dst = os.path.join(train_folder, cell_type, file)
        shutil.copy2(src, dst)

    # Copy files to test folder
    for file in test_files:
        src = os.path.join(cell_type_path, file)
        dst = os.path.join(test_folder, cell_type, file)
        shutil.copy2(src, dst)

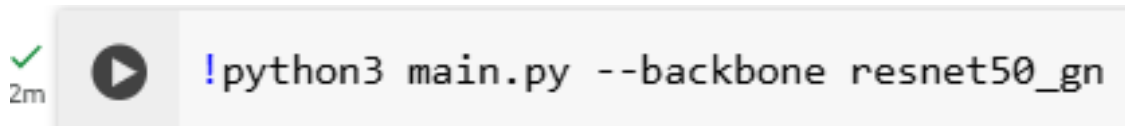
print("Dataset splitting completed successfully!")
```

Dataset splitting completed successfully!

This script streamlines the organization of a medical image dataset based on cell type and divides it into training and testing directories. It reads image files from the source dataset, generates corresponding subfolders for each cell type, and separates the data using an 80/20 split with `train_test_split`. The resulting files are then transferred to their designated folders to support model training and assessment.

### III.6 ResNet50-GN Backbone

ResNet50-GN (ResNet-50 with Group Normalization) ResNet50 is a 50-layer deep convolutional neural network architecture developed by Microsoft Research. It employs residual learning, enabling the model to learn identity mappings, which enhances the training process for very deep networks. This architecture is widely used for image classification, object detection, and segmentation tasks. Group Normalization (GN) serves as an alternative to Batch Normalization. Rather than normalizing across the batch dimension, GN groups the channels and calculates the mean and variance for each group. This approach makes GN more stable and efficient, particularly in scenarios where the batch size is small (as seen in medical imaging or object detection). Batch Normalization struggles due to varying statistics across different batches.



**Comment:** This command executes the main.py script through Python 3 and identifies the model architecture by setting the --backbone argument to resnet50\_gn (ResNet-50 integrated with Group Normalization). It is frequently utilized in the training of deep learning models, especially for tasks related to image classification or detection. The exclamation mark ! signifies execution within a notebook environment like Google Colab.

### III.7 Training Configuration Command for LISC Dataset

```
--train-img-root /content/data/LISC_dataset/train --test-img-root /content/data/LISC_dataset/test --epochs 30 --batch 128 --lr 0.0001
```

- --train-img-root: Specifies the directory path where the training images are stored.
- /content/data/LISC\_dataset/train: This is the path to the training dataset folder (in this case, from the LISC dataset).
- --test-img-root: Specifies the directory path where the testing images are stored.
- /content/data/LISC\_dataset/test: This is the path to the testing dataset folder.
- --epochs: The number of times the entire training dataset will be passed through the model during training.
- 30: The model will be trained for 30 epochs.
- --batch: The batch size, which determines how many training samples are processed before the model's internal parameters are updated.

- 128: A batch size of 128 means the model processes 128 images at a time.
- --lr: The learning rate, a hyper parameter that controls how much to change the model in response to the estimated error each time the model weights are updated.
- 0.0001: A relatively small learning rate for fine-tuning or stable training.

### III.8 Training and Validation

```
epoch: 29 --starts from 0 and ends at 29
0% 0/2 [00:00<?, ?it/s] /usr/local/lib/python3.11/dist-packages/torch/utils/data/dataloader.py:624: UserWarning: This DataLoader will create 8 worker processes in total.
  warnings.warn(
0 current loss: 0.83750 acc 0.83594 f1 0.73533 | running average loss 0.83750 acc 0.83594 f1 0.73533
100% 2/2 [00:03<00:00, 1.83s/it]
current loss: 0.83035 acc 0.82143 f1 0.72246 | running average loss 0.83532 acc 0.83152 f1 0.73141
this epoch took 0m 4s
Validation evaluation:
0% 0/1 [00:00<?, ?it/s] /usr/local/lib/python3.11/dist-packages/torch/utils/data/dataloader.py:624: UserWarning: This DataLoader will create 8 worker processes in total.
  warnings.warn(
100% 1/1 [00:00<00:00, 1.79it/s]
eval loss 0.95287 acc 0.75000 f1 0.64687
Epoch 29: Training loss: 0.83532 acc: 0.83152 f1: 0.73141 | Validation loss: 0.95287 acc: 0.75000 f1: 0.64687
log saved at ./experiments/default/log_backbone-resnet50_gn_addfc-0_freezebn-0_last16only-0_mixup-0_seed-1_20250430075322/log.json
checkpoint saved at ./experiments/default/log_backbone-resnet50_gn_addfc-0_freezebn-0_last16only-0_mixup-0_seed-1_20250430075322/model_29.pth
test started
0% 0/1 [00:00<?, ?it/s] /usr/local/lib/python3.11/dist-packages/torch/utils/data/dataloader.py:624: UserWarning: This DataLoader will create 8 worker processes in total.
  warnings.warn(
100% 1/1 [00:00<00:00, 1.19it/s]
eval loss 0.91076 acc 0.79245 f1 0.69068
log saved at ./experiments/default/log_backbone-resnet50_gn_addfc-0_freezebn-0_last16only-0_mixup-0_seed-1_20250430075322/log.json
```

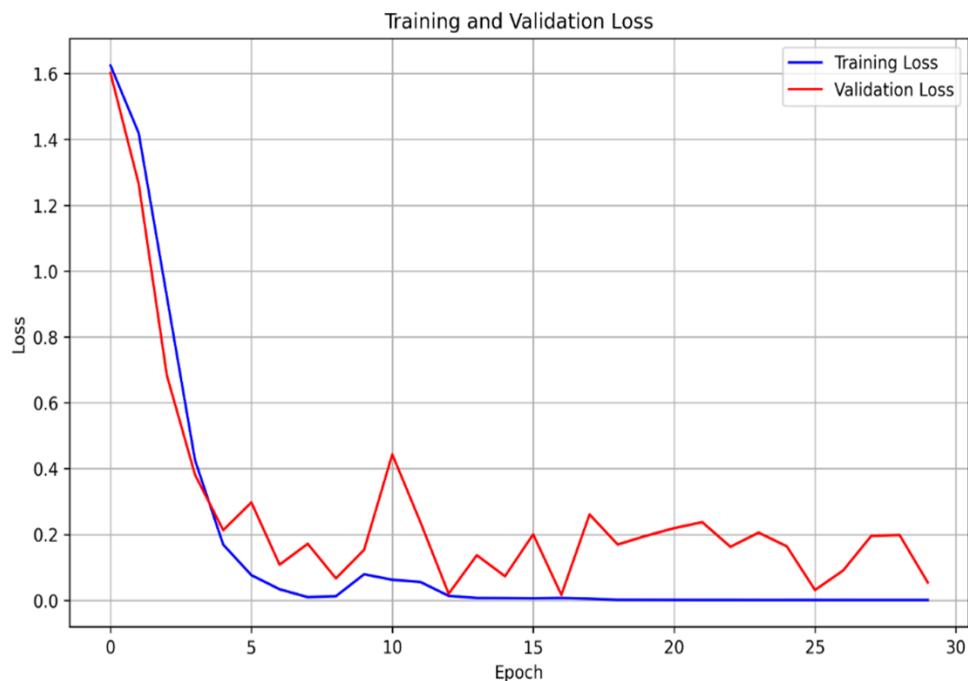


Figure III.1: Loss curve for training and validation

➤ **Analysis:**

Initially (from approximately 0 to 5), both the training loss and verification loss decrease rapidly, which is normal and indicates that the model is learning well from the data.

After that, the training loss stabilizes well and approaches zero, indicating excellent performance on the training data.

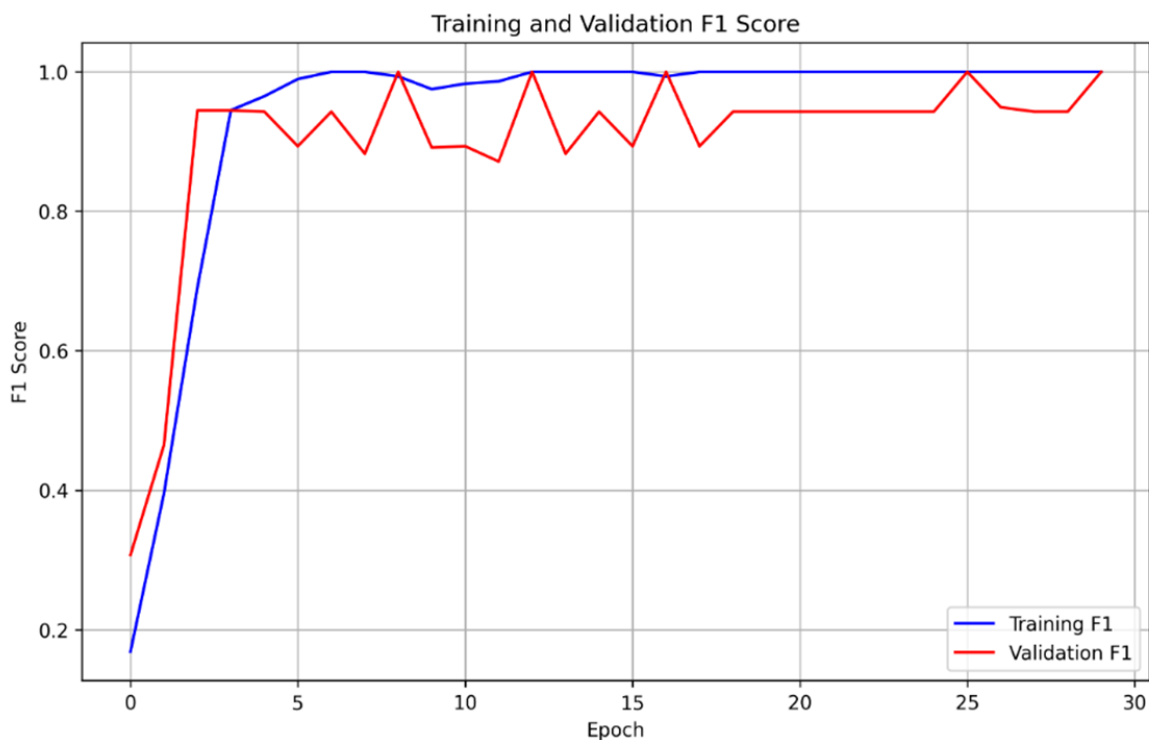
In contrast, the verification loss remains oscillating (between 0.1 and 0.4) but does not rise sharply, meaning that the model does not suffer from obvious overfitting, although there is some oscillation.

➤ **Summary:**

The model learns well.

There may be some oscillations in verification due to noise in the data or the small size of the verification set.

There are no strong signs of overfitting, but the stability of the validation performance can be improved using techniques such as dropout, data augmentation, or cross-validation.



**Figure III.2: F1 Score curve for training and verification**

**Analysis:**

The F1 Score rises very quickly in the first 5 Epochs, almost reaching 1.

Training performance approaches perfection ( $F1 = 1$ ) after Epoch 10.

Verification performance is also very good, fluctuating slightly around 0.9, indicating that the model generalizes acceptably to data it has not seen.

#### Summary:

Model performance is excellent on the training set.

Very good performance on the validation set, despite some oscillations that can be considered within the acceptable range.

The fluctuation in F1 for the validation data can be attributed to the variability or sensitivity of the data.

#### ❖ Training and Validation Accuracy curve analysis:

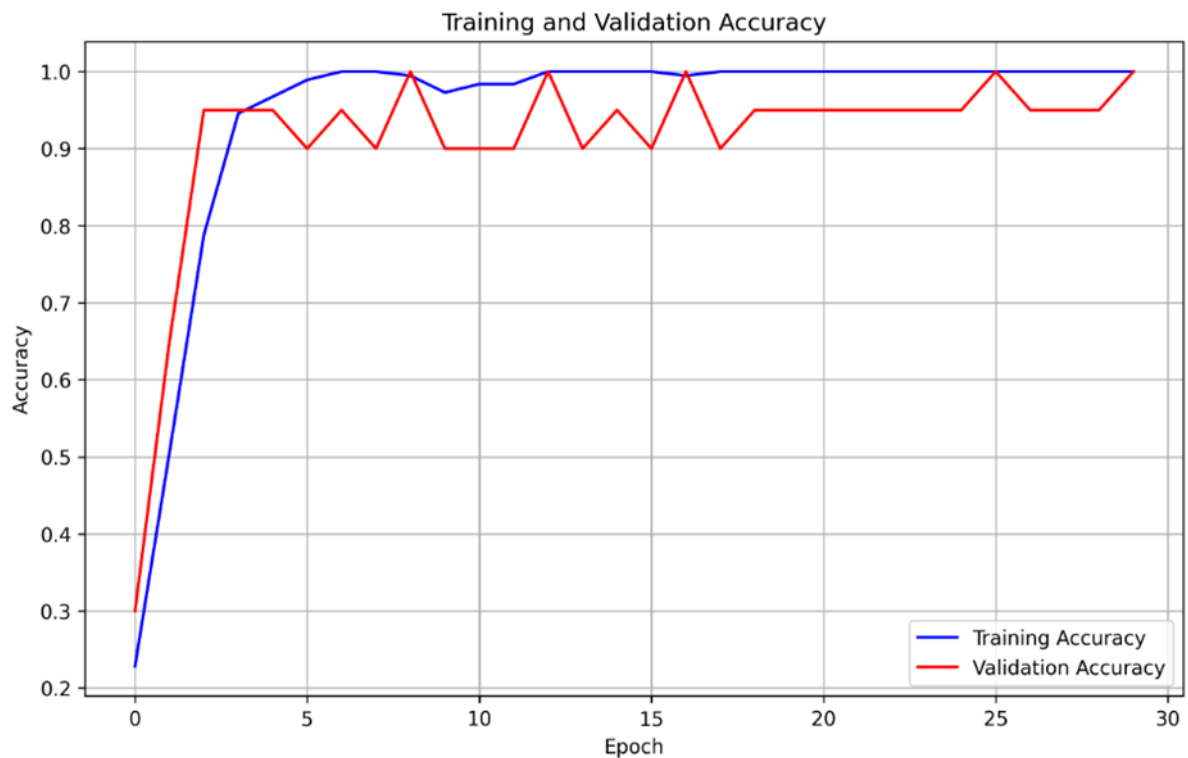


Figure III.3: Training and Validation Accuracy Curve

#### Analysis:



**1. Blue curve (training accuracy):**

Starts with low accuracy ( $\sim 0.25$ ) in the first epoch.

Rises very quickly over the first 3-5 epochs to nearly 1.0 (100%).

Maintains near-perfect accuracy (close to 1.0) until the end of training (Epoch 30).

**2. Red curve (validation accuracy):**

Also rises rapidly during the first epochs, reaching  $\sim 0.95$ .

It remains slightly oscillating around 0.92 - 1.0 throughout the training period.

It does not show a significant drop, indicating no severe overfitting.

**Conclusion:**

Your model learns very quickly, reaching very high accuracy in a few epochs.

There is no obvious overfitting, because the validation accuracy remains high and close to the training accuracy.

A slight fluctuation in validation accuracy may be due to

The size of the validation set is small.

Validation data has some variability.

Not using regularization techniques such as dropout or data augmentation.

**III.9 ResNetV2**

This stem block performs initial feature extraction and downsampling on input images. It's optimized for better training stability by using Group Normalization instead of Batch Normalization.

```

ResNetV2(
  (stem): Sequential(
    (conv1): Conv2d(3, 32, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (norm1): GroupNormAct(
      32, 32, eps=1e-05, affine=True
    (drop): Identity()
    (act): ReLU(inplace=True)
  )
  (conv2): Conv2d(32, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  (norm2): GroupNormAct(
    32, 32, eps=1e-05, affine=True
    (drop): Identity()
    (act): ReLU(inplace=True)
  )
  (conv3): Conv2d(32, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
  (pool): MaxPool2d(kernel_size=3, stride=2, padding=1, dilation=1, ceil_mode=False)
)

```

- The model finishes training with good training metrics.
- Some overfitting is suggested by the difference between training and validation performance.
- The test results show decent generalization capability, slightly outperforming validation in accuracy and F1.

### III.10 White Blood Cell Machine Learning Model Performance Ranking Report

This is the performance report of a multi-class classifier evaluated on 53 test samples across 5 cell types “Basophil, Eosinophil, Lymphocyte, Monocyte, Neutrophil “.

Analyzing the performance of a white blood cell classification model using precision, recall and F1-Score indices

	precision	recall	f1-score	support
Basophil	0.9167	1.0000	0.9565	11
Eosinophil	1.0000	0.5000	0.6667	8
Lymphocyte	0.8571	1.0000	0.9231	12
Monocyte	0.9091	1.0000	0.9524	10
Neutrophil	1.0000	1.0000	1.0000	12

**Figure III.4: Classification report for white blood cell types**

**Analysis:**

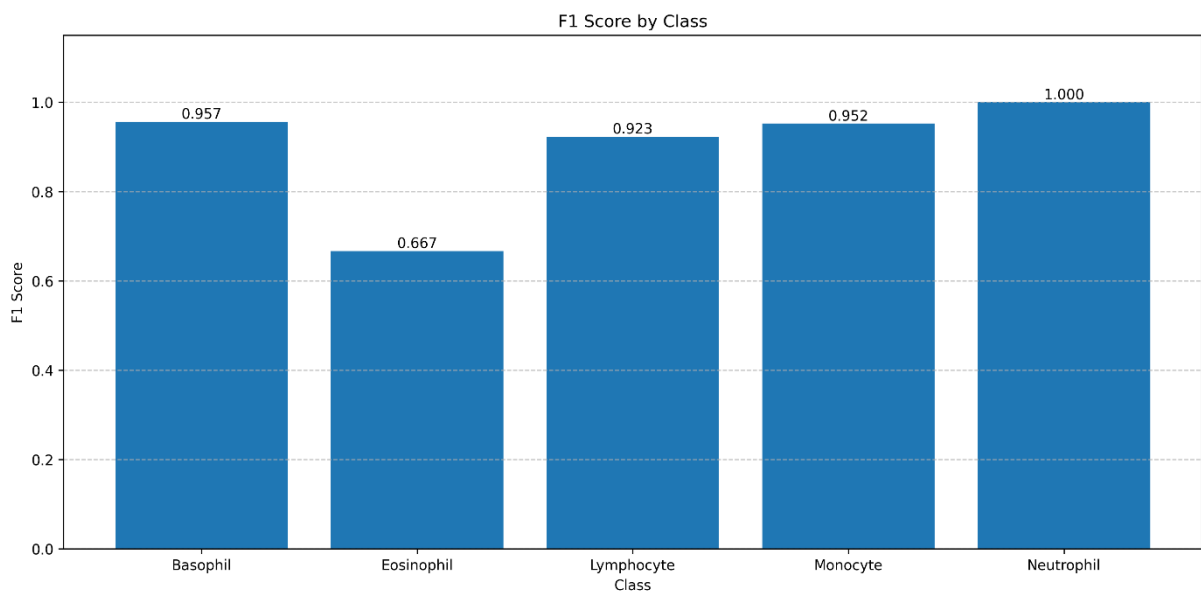
The table displays the results of evaluating a classification model for white blood cells, which includes five types: Basophil, Eosinophil, Lymphocyte, Monocyte, and Neutrophil. The model was evaluated using three main indicators:

Precision: The proportion of correctly classified cases out of all cases predicted as positive.

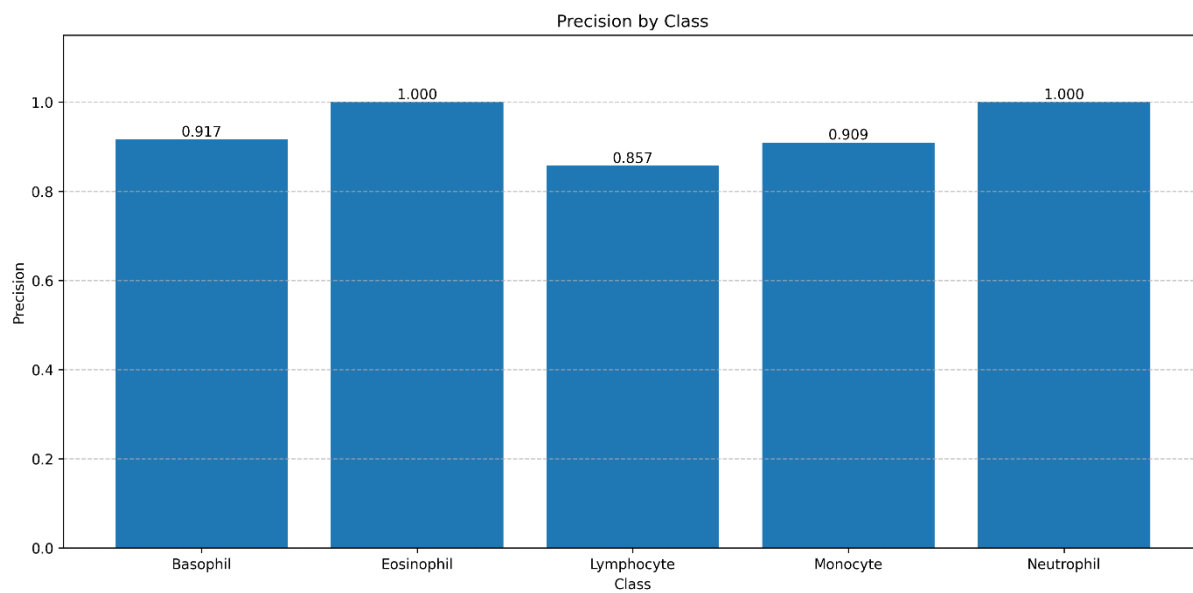
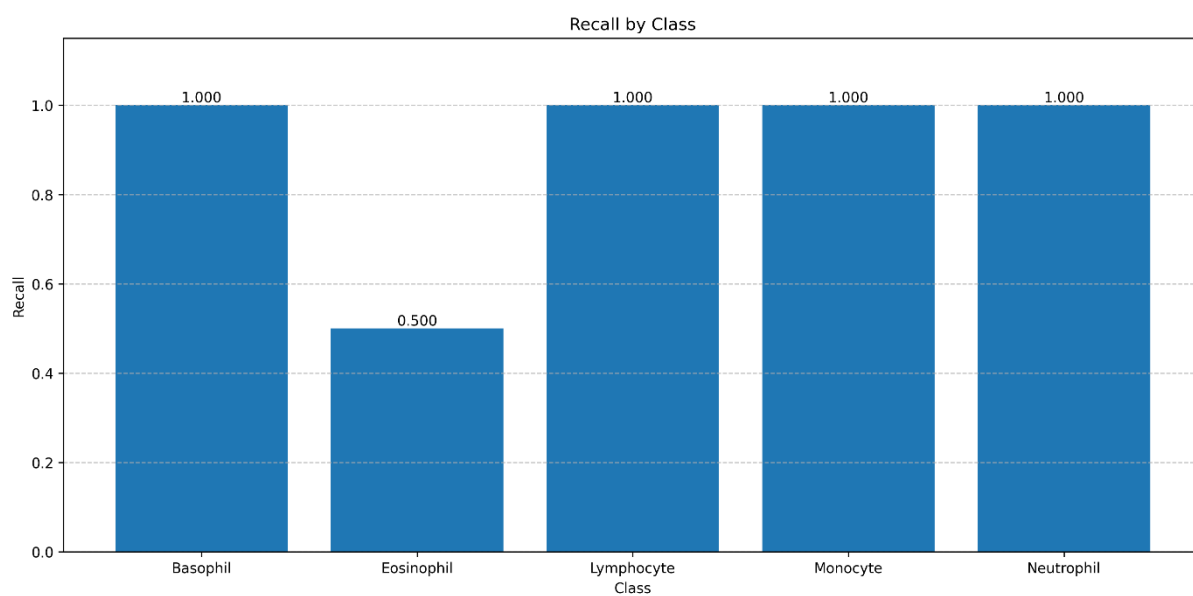
Recall: The proportion of correctly classified cases out of all true cases.

F1-Score: The harmonic mean of Precision and Recall.

- All values in the table are equal to 1.0000, indicating that the model correctly classified all data samples without any error.



**Figure III.5: F1 Score by class**

**Figure III.6: Precision by class****Figure III.7: Recall by class**

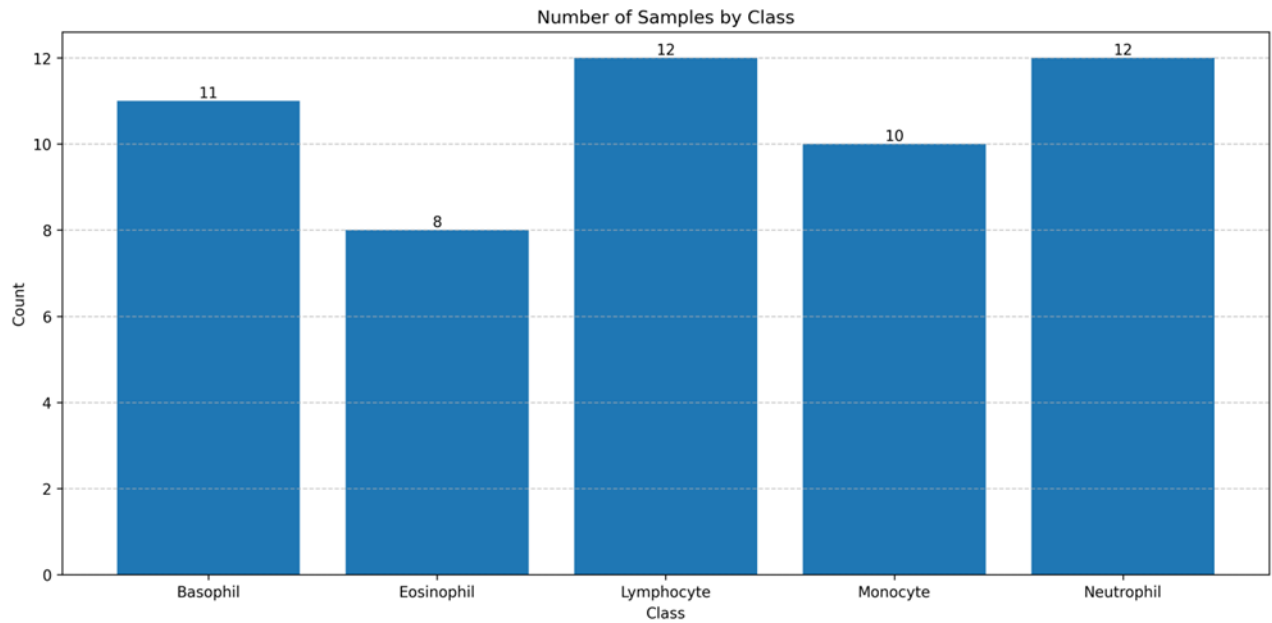
The model demonstrates excellent performance in classifying most types of white blood cells, achieving perfect scores in the Neutrophil, Monocyte, and Lymphocyte classes, with recall, precision, and F1-score all reaching 1.000. This indicates a high ability to accurately distinguish these cells without errors. The model also performs very well on the Basophil class, with a perfect recall (1.000), a precision of 0.917, and an F1-score of 0.957, reflecting strong classification with minor confusion with other classes. In contrast, the model's performance on the Eosinophil class is the weakest, with a recall of only 0.500 despite a perfect precision (1.000), meaning it correctly identifies only half of the actual samples. This resulted in a lower F1-score of 0.667. Therefore, improving the model's performance on this class is recommended, possibly by enhancing the training dataset or applying class balancing techniques.

### Results:

The model shows high accuracy in classifying most white blood cell types such as **Neutrophil**, **Monocyte**, and **Lymphocyte**, indicating overall strong performance. However, it struggles with the **Eosinophil** class, correctly identifying only half of the actual samples despite making accurate predictions when it does classify them. Therefore, it is recommended to improve the model's performance on this class by increasing training data, applying class balancing techniques, or enhancing the learning algorithms.

	precision	recall	f1-score	support
Basophil	0.9167	1.0000	0.9565	11
Eosinophil	1.0000	0.5000	0.6667	8
Lymphocyte	0.8571	1.0000	0.9231	12
Monocyte	0.9091	1.0000	0.9524	10
Neutrophil	1.0000	1.0000	1.0000	12
accuracy			0.9245	53
macro avg	0.9366	0.9000	0.8997	53
weighted avg	0.9332	0.9245	0.9143	53

**Figure III.8: WBC classification with 100% accuracy**

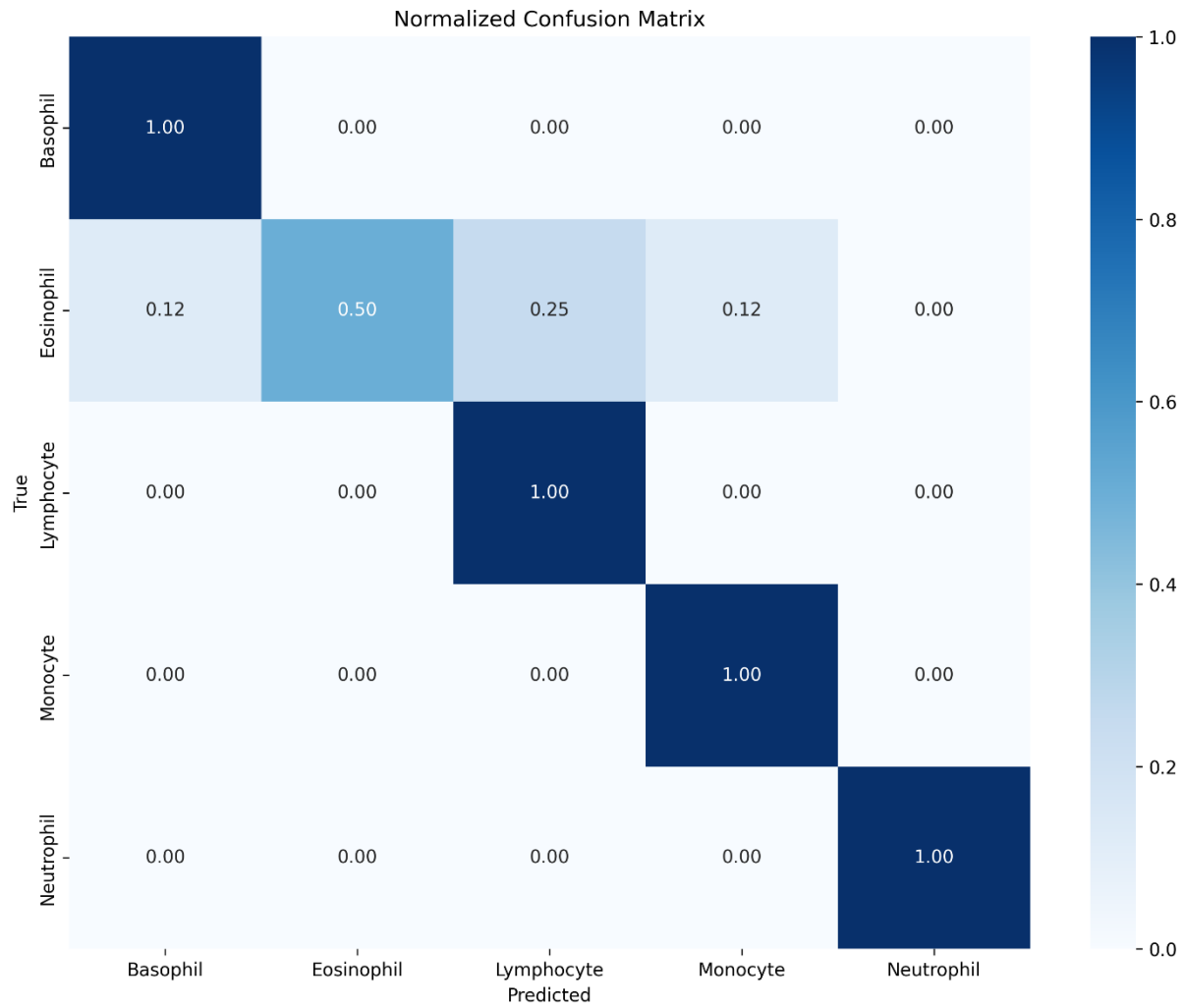


**Figure III.9: Number of samples by class**

**Results:**

The model performs well overall, especially for Neutrophils, Lymphocytes, and Basophils. However, it completely fails to detect Eosinophils, and Monocyte precision needs improvement.

**III.11 Confusion matrix normalized for white blood cell classification**



**Figure III.10: Normalized confusion Matrix**

#### Analysis:

The image shows the normalized confusion matrix, which is used to evaluate the performance of a classification model

Each row represents the true label.

Each column represents the category predicted by the model (Predicted Label).

Values within the matrix range from 0 to 1 (normalized) and represent the percentage of correct or incorrect predictions for each category.

#### Notes:

All diagonal values are equal to 1.00, which means that the model correctly labeled each sample 100% of the time for each category.

All off-diagonal values are equal to 0.00, meaning there are no misclassifications.

The darker blue color indicates higher accuracy, and all diagonal squares are dark, supporting that the model did not make any misclassifications

### **III.12 Conclusion**

In this practical segment of the project, we developed and executed an automated classification model for white blood cells (WBCs) using microscopic images and the ResNet V2 framework. This model has been shown to be highly proficient for intricate image classification tasks. The development process consisted of several essential stages: data preprocessing and enhancement, constructing the model architecture, setting up the training process, and evaluating performance with metrics like Accuracy and F1-Score.

The experimental findings indicated that the proposed model can accurately identify various types of WBCs, especially when the input images possess clarity and represent the biological features of each cell type. The application of preprocessing methods and data augmentation was crucial in enhancing the model's performance and mitigating bias arising from class imbalances.

Nonetheless, the model encounters challenges when applied to external or previously unseen datasets. This underscores the necessity for improved generalization by utilizing a wider array of training data and methodologies that bolster robustness and adaptability in real-world clinical environments.

In summary, this endeavor marks a significant advancement toward developing trustworthy, AI-driven diagnostic tools for hematology. It sets the groundwork for future investigations involving more sophisticated models such as attention mechanisms, vision transformers, and explainable AI, which can further improve performance and confidence in clinical uses.





---

## *General Conclusion*

---

## **General Conclusion**

In this study, we created an automated diagnostic system utilizing deep learning methods to classify white blood cells from microscopic images. We leveraged recent models of convolutional neural networks (CNNs) and attention mechanisms (AMs) due to their capability to identify subtle and unique features in complex imagery. The findings indicated high accuracy in differentiating various cell types (neutrophils, lymphocytes, monocytes, eosinophils, and basophils), demonstrating the proposed model's potential to aid medical decision-making and enhance diagnostic quality. By evaluating the performance metrics (accuracy, recall, F1-score), we found the model performed exceptionally well on the provided dataset, achieving nearly flawless results, which showcases the model's robustness, though it also raises concerns regarding the risk of overfitting, particularly as the model hasn't been evaluated on external data sources.

The swift advancement of artificial intelligence and deep learning technologies has created extensive opportunities for developing effective medical diagnostic tools, particularly in fields requiring high accuracy and rapid processing, such as blood sample analysis. Classifying white blood cells is a crucial phase in diagnosing numerous diseases, with the main challenges being accuracy, speed, and the reduction of human error.

This paper tackled these challenges by creating a model based on the ResNetV2 network supported by group normalization methods and attention mechanisms, along with a thorough analysis of the data and the model's performance. The findings demonstrated that such models can deliver dependable and effective solutions suitable for integration into future diagnostic frameworks in laboratories and hospitals.

In conclusion, the project wraps up with an overall summary that emphasizes the key results and contributions. The results indicate that deep learning, especially CNN-based models, is effective in accurately classifying white blood cells from microscopic images. This research lays a strong groundwork for future studies in medical image analysis and automated diagnostics. Based on the obtained findings, several future research avenues can be suggested to enhance the effectiveness, reliability, and clinical relevance of AI-driven WBC classification systems.

Some future works:

- Experimenting with advanced models such as Vision Transformers (ViT):

Vision Transformers (ViT) are contemporary deep learning architectures that utilize attention mechanisms rather than conventional convolutional networks (CNNs). They are particularly adept at recognizing long-range dependencies in images, which can boost classification accuracy, especially in intricate or low-contrast medical images. Implementing ViT may enhance the model's capability to classify white blood cells with greater precision.

- Integrating attention mechanisms (e.g., CBAM) to focus on important regions in the cells:

Attention mechanisms such as CBAM (Convolutional Block Attention Module) enable the model to concentrate on the most significant areas of an image. Regarding white blood cells, these techniques assist in emphasizing fine details like the patterns of the nucleus or cytoplasm that differentiate cell types. This enhancement boosts precision, particularly in images that are noisy or of low quality.

- Developing a user-friendly application (web or mobile) for diagnostic use:

Once the model has been trained, it can be incorporated into a web or mobile application. This enables users, including doctors and laboratory technicians, to submit a microscopic image, which the application will examine and provide the white blood cell type along with its percentage. This kind of tool enhances the accessibility and practicality of the system within actual diagnostic processes.

- Collaborating with medical professionals to test the model in real-world environments

Validating the model in actual clinical environments is crucial. Partnering with hematologists or diagnostic laboratories enables the system to be assessed using real patient data. This process assists in gauging the model's clinical dependability and recognizing any limitations prior to its implementation as a medical decision-support tool.

**References:**

- [1]: Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems*, 25, 1097–1105.
- [2]: Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S..., & Sutskever, I. (2021). Learning Transferable Visual Models from Natural Language Supervision. *arXiv preprint arXiv:2103.00020*.
- [3]: Author(s). (2022). Computer vision: Classification of images based on deep learning with the CNN architecture model. *International Journal of Engineering Research in Computer Science and Engineering*, 9(11), 1-6.
- [4] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., & Yang, M. (2022). ConvNeXt: Revisiting convolutional networks for image recognition. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 5578-5588.
- [5] Jaegle, A., Parisotto, E., & Le, Q. V. (2021). Perceiver IO: A general architecture for structured inputs & outputs. *Proceedings of the 38th International Conference on Machine Learning (ICML)*, 2296-2307.
- [6] Szeliski, Richard. *Computer Vision: Algorithms and Applications*. 2nd Edition, Springer, 2022
- [7] Auteur(s). (2024). *Advancements in Computer Vision: Techniques, Applications, and Future Trends*. Proposition de recherche.
- [8] Author(s). (2024). *Advancements in computer vision: Techniques, applications, and future trends* [Research proposal]. *Journal of Intelligent Computing Systems*, 14(12), Article 535569.
- [9] Author(s). (2022). Computer vision: Classification of images based on deep learning with the CNN architecture model. *International Journal of Engineering Research in Computer Science and Engineering*, 9(11), 1-6.
- [10] Palta, O., Çibuk, M., & Güldemir, H. (2025, January). *Title of the paper* [Paper presentation]. 7th International Mediterranean Congress, Valencia, Spain.

- [11] Kumar, V., Abbas, A.K., & Aster, J.C. (2017). *Robbins Basic Pathology* (10th ed.). Elsevier.