

# Title of the Final Project

First Author  
`firstrauthor@i1.org`

Second Author  
`secondauthor@i2.org`

Third Author  
`thirdauthor@i2.org`

Fourth Author  
`fourthauthor@i2.org`

## Abstract

*Brief summary of the work developed as well as the main results and contributions.*

## 1. Introduction

Here, the problem to be solved is described (what do we want to do?), the motivation (why is it relevant to do it?), and the goals (what specific objectives are we going to address in order to solve the problem?).

We use the L<sup>A</sup>T<sub>E</sub>X format of the CVPR conference. This document can be written either in English or Spanish.

What follows is a tentative approximation of the sections the document should have. If students consider that they need others, as well as subdividing the sections into different subsections, they can do it without problem.

Students can take inspiration from the cs231n final projects website: <http://cs231n.stanford.edu/project.html>, where different projects are presented and developed.

This whole final report can have from 6 to 8 pages (no more and no less).

## 2. Background

This section presents the fundamental concepts necessary to understand the work.

## 3. Related Works

It presents what has been done in the field previously, and what the best methods are currently. It is very important, in general, not only in this section, to adequately document the relevant literature. To do this, you must use the .bib file, in the way I show here: [1-3]

## 4. Methods

### 4.1. Modelos de Clasificación

En nuestro problema tenemos dos etapas de clasificación de imágenes. Primero, una clasificación de hojas en tipos de plantas, y luego una clasificación de las enfermedades presentes en las hojas de cada tipo de planta.

Estos son dos problemas bastante diferentes en cuanto a la complejidad de la tarea, principalmente porque en la primera etapa las diferencias entre clases son más notorias (muy diferentes tipos de plantas), mientras que en la segunda etapa las diferencias son más sutiles (diferentes enfermedades que pueden afectar a la misma planta, y que a veces pueden parecerse mucho entre sí). Además, nos enfrentamos al problema de la diferencia de tamaños de las imágenes de las hojas y la de sus enfermedades.

Con esto en mente, tiene sentido estudiar una amplia gama de modelos de clasificación, y elegir los más adecuados para cada etapa del problema.

He dividido el estudio inicial de los modelos en dos partes: arquitecturas clásicas y arquitecturas modernas basadas en Transformers. En las tablas 1 y 2 se presentan las comparaciones detalladas de las diferentes arquitecturas estudiadas.

#### 4.1.1 Modelos clásicos

En la tabla 1 se presentan las comparaciones detalladas de las diferentes arquitecturas clásicas estudiadas: GoogLeNet (Inception), ResNet, EfficientNetV1 y EfficientNetV2. Aquí se describen las características principales de cada modelo, sus ventajas y desventajas, y por qué se eligieron para este proyecto.

A lo largo del tiempo se han desarrollado técnicas interesantes que han llevado al avance de estos modelos. A continuación los describiré brevemente:

1. **GoogLeNet (Inception)**: Introdujo los módulos Inception que permiten capturar características a múltiples escalas mediante convoluciones paralelas de diferentes tamaños. Esto mejora la eficien-

Característica	GoogLeNet (Inception)	ResNet	EfficientNet (V1)	EfficientNetV2
Primera Publicación	2014 (Szegedy et al., Google)	2015 (He et al.)	2019 (Tan & Le, Google)	2021 (Tan & Le, Google)
Idea Principal	Módulos Inception con convoluciones paralelas de múltiples escalas	Conexiones residuales (skip) para permitir el entrenamiento de redes muy profundas	Escalado compuesto: escalar conjuntamente profundidad, anchura y resolución usando un coeficiente fijo	Escalado compuesto mejorado + entrenamiento más rápido: arquitectura más inteligente + regularización adaptativa
Bloque Básico	Módulo Inception: convoluciones paralelas $1 \times 1$ , $3 \times 3$ , $5 \times 5$ + max pooling, concatenadas	<ul style="list-style-type: none"> <li>BasicBlock (<math>2 \times 3 \times 3</math> convs) para ResNet18/34</li> <li>Bottleneck (<math>1 \times 1 \rightarrow 3 \times 3 \rightarrow 1 \times 1</math> convs) para ResNet50+</li> </ul>	MBConv (Mobile Inverted Bottleneck): <ul style="list-style-type: none"> <li><math>1 \times 1</math> exp <math>\rightarrow 3 \times 3</math> depthwise</li> <li><math>\rightarrow</math> SE <math>\rightarrow 1 \times 1</math> proj</li> </ul>	Fused-MBConv (etapas iniciales): conv $3 \times 3$ regular en lugar de $1 \times 1$ +depthwise MBConv (etapas posteriores) SE solo en bloques posteriores
Conexiones Skip	$\times$ No (pero versiones posteriores como Inception-ResNet sí)	✓ Si (residual aditivo)	✓ Sí (MBConv usa residual cuando strides=1)	✓ Si
Atención	$\times$ No	$\times$ No	✓ Squeeze-and-Excitation (SE) en cada bloque	✓ SE, pero solo en etapas posteriores (para reducir overhead)
Método de Escalado	Manual: aumento de módulos Inception apilados	Aumento manual de profundidad ( $18 \rightarrow 34 \rightarrow 50 \rightarrow 101 \rightarrow 152$ )	Escalado compuesto uniforme: Profundidad $\times \alpha$ , Anchura $\times \beta$ , Res $\times \gamma$	Escalado compuesto no uniforme + aprendizaje progresivo (imágenes pequeñas $\rightarrow$ grandes durante entrenamiento) Menor que V1 para la misma precisión: S: 384, M: 480, L: 480–512
Resolución de Entrada	Fija ( $224 \times 224$ )	Fija (usualmente $224 \times 224$ ), pero flexible	Escala con el tamaño del modelo: B0: $224 \rightarrow$ B7: 600	
Tamaños Típicos	<ul style="list-style-type: none"> <li>GoogLeNet (6.8M params)</li> <li>Inception-v3 (23.8M)</li> <li>Inception-v4 (42.7M)</li> </ul>	<ul style="list-style-type: none"> <li>ResNet18 (11M params)</li> <li>ResNet34 (21M)</li> <li>ResNet50 (25M)</li> <li>ResNet101 (44M)</li> </ul>	<ul style="list-style-type: none"> <li>B0 (5.3M)</li> <li>B3 (12M)</li> <li>B5 (30M)</li> <li>B7 (66M)</li> </ul>	<ul style="list-style-type: none"> <li>S (21M)</li> <li>M (54M)</li> <li>L (120M)</li> </ul>
Precisión (Imagenet)	<ul style="list-style-type: none"> <li>GoogLeNet: <math>\sim 69.8\%</math></li> <li>Inception-v3: <math>\sim 78.8\%</math></li> </ul>	<ul style="list-style-type: none"> <li>ResNet50: <math>\sim 76\%</math></li> <li>ResNet101: <math>\sim 78\%</math></li> </ul>	<ul style="list-style-type: none"> <li>B0: <math>\sim 77\%</math></li> <li>B3: <math>\sim 82\%</math></li> <li>B7: <math>\sim 84.4\%</math></li> </ul>	<ul style="list-style-type: none"> <li>S: <math>\sim 83.9\%</math></li> <li>M: <math>\sim 85.2\%</math></li> <li>L: <math>\sim 85.7\%</math></li> </ul>
Velocidad de Entrenamiento	Moderada	Moderada	Lenta (debido a entradas de alta resolución y profundidad)	Mucho más rápida (hasta $11\times$ vs V1 con precisión similar)
Eficiencia de Inferencia	Buena; diseño eficiente en parámetros para su época	Buena para modelos pequeños (18/34); más pesada para 101+	Muy eficiente en parámetros, pero alta resolución afecta latencia	Mejor relación FLOPs/precisión; optimizado para TPU/GPU
Mejor Para	<ul style="list-style-type: none"> <li>Comprensión histórica de CNNs</li> <li>Extracción de características multi-escala</li> </ul>	<ul style="list-style-type: none"> <li>Modelos base</li> <li>Transfer learning con datos limitados</li> <li>Prototipado rápido</li> </ul>	<ul style="list-style-type: none"> <li>Despliegue con recursos limitados</li> <li>Cuando se necesita alta precisión con menos parámetros</li> </ul>	<ul style="list-style-type: none"> <li>Sistemas de producción que necesitan velocidad + precisión</li> <li>Entrenamiento desde cero o fine-tuning</li> </ul>
fastai / PyTorch	✓ Nativo en torchvision (googlenet, inception_v3)	✓ Nativo en torchvision y fastai	✓ Vía timm (ej., efficientnet_b0)	✓ Vía timm (ej., efficientnetv2_s)
Debilidad Principal	Arquitectura obsoleta; menor precisión que modelos modernos	Eficiencia de parámetros subóptima; obsoleto vs modelos modernos	Entrenamiento lento; SE añade overhead computacional; alta resolución = mucha memoria	Ecosistema ligeramente menos maduro; no está en torchvision

Table 1. Comparación de las Arquitecturas GoogLeNet (Inception), ResNet, EfficientNetV1 y EfficientNetV2

cia computacional y la capacidad de extracción de características.

2. **ResNET (Capas Residuales):** Introdujo las conexiones residuales (skip connections) que facilitan el entrenamiento de redes muy profundas al mitigar el problema del desvanecimiento del gradiente. Esto permitió construir redes con cientos de capas.

3. **EfficientNetV1 (Escalado Compuesto):** Pro-

puso un método sistemático para escalar redes neuronales ajustando conjuntamente la profundidad, anchura y resolución de entrada mediante un coeficiente fijo. Esto llevó a modelos más eficientes en términos de precisión y tamaño.

4. **EfficientNetv2 (Progesive Learning):** Mejoró la velocidad de entrenamiento mediante una arquitectura más inteligente y una regularización adaptativa.

Característica	ConvNeXt	Swin Transformer	MaxViT
Primera Publicación	2022 (Liu et al., Meta/FAIR)	2021 (Liu et al., Microsoft)	2022 (Tu et al., Google)
Idea Principal	Modernización de ResNet con técnicas inspiradas en Transformers, demostrando que CNNs puras pueden competir	Transformer jerárquico con atención en ventanas locales y desplazamiento para conexión entre ventanas	Combinación de convoluciones + atención local (block) + atención global diluida (grid) en cada bloque
Tipo de Arquitectura	CNN pura (sin atención)	Transformer puro (sin convoluciones en backbone)	Híbrido CNN + Transformer
Bloque Básico	Inverted bottleneck con depthwise conv $7 \times 7$ , LayerNorm, GELU	Swin Transformer Block: Window-MSA → MLP → Shifted-Window-MSA → MLP	MBConv → Block Attention → Grid Attention
Mecanismo de Atención	✗ No usa atención (solo convoluciones)	✓ Multi-head Self-Attention en ventanas $7 \times 7$ , con shift en capas alternas	✓ Block Attention (local en ventanas) + Grid Attention (global diluida)
Receptive Field Global	Solo en capas profundas (por apilamiento)	Gradual mediante shifted windows	✓ Desde las primeras capas (via Grid Attention)
Complejidad Computacional	$O(n)$ lineal respecto a resolución	$O(n)$ lineal (atención local en ventanas fijas)	$O(n)$ lineal (block + grid attention)
Representación Jerárquica	✓ Sí (4 etapas con downsampling)	✓ Sí (patch merging entre etapas)	✓ Sí (estructura multi-escala)
Resolución de Entrada	$224 \times 224$ (escalable)	$224 \times 224$ (escalable a alta resolución)	$224 \times 224$ a $512 \times 512$
Tamaños Típicos	<ul style="list-style-type: none"> <li>ConvNeXt-T (29M)</li> <li>ConvNeXt-S (50M)</li> <li>ConvNeXt-B (89M)</li> <li>ConvNeXt-L (198M)</li> </ul>	<ul style="list-style-type: none"> <li>Swin-T (29M)</li> <li>Swin-S (50M)</li> <li>Swin-B (88M)</li> <li>Swin-L (197M)</li> </ul>	<ul style="list-style-type: none"> <li>MaxViT-T (31M)</li> <li>MaxViT-S (69M)</li> <li>MaxViT-B (120M)</li> <li>MaxViT-L (212M)</li> </ul>
Precisión (ImageNet-1K)	<ul style="list-style-type: none"> <li>T: ~82.1%</li> <li>B: ~83.8%</li> <li>L: ~84.3%</li> <li>XL (21K): ~87.8%</li> </ul>	<ul style="list-style-type: none"> <li>T: ~81.3%</li> <li>B: ~83.5%</li> <li>L: ~86.4%</li> <li>L (21K): ~87.3%</li> </ul>	<ul style="list-style-type: none"> <li>T: ~83.6%</li> <li>B: ~85.0%</li> <li>L: ~86.4%</li> <li>L (21K): ~89.5%</li> </ul>
Velocidad de Entrenamiento	Rápido (operaciones de convolución optimizadas)	Moderado (overhead de atención)	Moderado-lento (múltiples tipos de atención)
Eficiencia de Inferencia	Excelente (convoluciones muy optimizadas en GPU)	Buena (atención local eficiente)	Buena (pero más complejo que los otros)
Backbone Universal	✓ Sí (detección, segmentación)	✓ Sí (muy popular para downstream tasks)	✓ Sí (clasificación, detección, segmentación)
Mejor Para	<ul style="list-style-type: none"> <li>Cuando se prefiere simplicidad de CNNs</li> <li>Inferrencia eficiente</li> <li>Transfer learning</li> </ul>	<ul style="list-style-type: none"> <li>Backbone para tareas densas</li> <li>Estado del arte en detección/segmentación</li> <li>Preentrenamiento a gran escala</li> </ul>	<ul style="list-style-type: none"> <li>Máxima precisión</li> <li>Cuando se necesita contexto global desde el inicio</li> <li>Tareas que requieren relaciones de largo alcance</li> </ul>
Disponibilidad	✓ torchvision + timm	✓ torchvision + timm	✓ timm (maxvit_*)
Debilidad Principal	Menor capacidad de modelar dependencias globales que Transformers	Conexión entre ventanas limitada por shifted windows	Mayor complejidad de implementación; menos maduro que Swin/ConvNeXt

Table 2. Comparación de Arquitecturas Modernas: ConvNeXt, Swin Transformer y MaxViT

5. **MBConv (Mobile Inverted Bottleneck)**: Utilizado en EfficientNet, este bloque combina convoluciones depthwise separables con una expansión y proyección de canales, lo que reduce significativamente el número de parámetros y operaciones computacionales. Se basa en que las capas más finas (bottlenecks) ya tienen la información necesaria, y las convoluciones depthwise son más eficientes.

6. **Fused-MBConv**: Introducido en EfficientNetV2, este bloque combina convoluciones regulares con depthwise en las primeras etapas de la red para mejorar la velocidad de entrenamiento sin sacrificar la precisión.

7. **Squeeze-and-Excitation (SE)**: Introdujo bloques de atención que recalibran dinámicamente las características de los canales, mejorando la capacidad del modelo para enfocarse en características relevantes.

Estas son algunas de las técnicas clave que han impulsado el desarrollo de modelos clásicos de clasificación de imágenes. Cada una de ellas ha contribuido a mejorar la precisión, eficiencia y capacidad de generalización de las redes neuronales convolucionales.

#### 4.1.2 Modelos modernos

El desarrollo de los Transformers ha revolucionado el campo del aprendizaje profundo, y su aplicación

a la visión por computadora ha llevado a la creación de modelos modernos que superan a las arquitecturas clásicas en muchas tareas. En la tabla 2 se presentan las comparaciones detalladas de las diferentes arquitecturas modernas estudiadas: ConvNeXt, Swin Transformer y MaxViT. Aquí se describen las características principales de cada modelo, sus ventajas y desventajas, y por qué se eligieron para este proyecto.

A continuación se describen brevemente las técnicas clave de estos modelos modernos:

1. **Convolución Depthwise separable**
2. **LayerNorm:**
3. **GELU**
4. **Swin (patch + transformers):** Introduce un enfoque jerárquico con atención en ventanas locales y desplazamiento para conectar ventanas, lo que permite manejar imágenes de alta resolución de manera eficiente.
5. **MaxViT:** Combina convoluciones con atención local (block) y atención global diluida (grid) en cada bloque, permitiendo capturar tanto características locales como relaciones de largo alcance desde las primeras capas.
6. **Atención en Ventanas Desplazadas:** Utilizada en Swin Transformer, esta técnica permite que la atención se aplique en ventanas locales, pero con un desplazamiento entre capas para conectar diferentes regiones de la imagen.
7. **Atención Global Diluida:** Empleada en MaxViT, esta técnica permite capturar información global sin el alto costo computacional de la atención completa, utilizando una estrategia de muestreo diluido.

## 5. Experiments

The data used, the experimental validation protocol, the metrics used, the experiments carried out, the results obtained, and their discussion are presented here.

### 5.1. Dataset

## 6. Conclusions

Section that presents, briefly and as a summary, the main conclusions of the work carried out. It also usually includes future possible works. That is, what are the most promising lines to continue with this work, as well as possible proposals for improvement. IMPORTANT: these are the scientific conclusions reached in the project; not your personal conclusions about the work you have done!

## References

- [1] Stéphane Lathuilière, Pablo Mesejo, Xavier Alameda-Pineda, and Radu Horaud. A comprehensive analysis of deep regression. *IEEE transactions on pattern analysis and machine intelligence*, 42(9):2065–2081, 2019. 1
- [2] Pablo Mesejo, Daniel Pizarro, Armand Abergel, Olivier Rouquette, Sylvain Beorchia, Laurent Poincloux, and Adrien Bartoli. Computer-aided classification of gastrointestinal lesions in regular colonoscopy. *IEEE transactions on medical imaging*, 35(9):2051–2063, 2016. 1
- [3] Víctor A Vargas-Pérez, Pablo Mesejo, Manuel Chica, and Oscar Cordón. Deep reinforcement learning in agent-based simulations for optimal media planning. *Information Fusion*, 91:644–664, 2023. 1