

Universidade de Brasília - UnB  
Faculdade de Ciências e Tecnologias em Engenharia - FCTE  
Engenharia de Software

**Análise Comparativa de Arquiteturas  
Convolucionais e Baseadas em Atenção para  
Detecção Automatizada de Parasitas Intestinais**

Autor: Eduardo Vieira Lima  
Orientador: Dr. Vinicius de Carvalho Rispoli

Brasília, DF  
julho de 2025





---

Eduardo Vieira Lima

Análise Comparativa de Arquiteturas Convolucionais e Baseadas em Atenção para Detecção Automatizada de Parasitas Intestinais/ Eduardo Vieira Lima. – Brasília, DF, julho de 2025-

97 p. : il. (algumas color.) ; 30 cm.

Orientador: Dr. Vinicius de Carvalho Rispoli

Trabalho de Conclusão de Curso – Universidade de Brasília - UnB

Faculdade de Ciências e Tecnologias em Engenharia - FCTE , julho de 2025.

1. Redes neurais. 2. Vision Transformer. I. Dr. Vinicius de Carvalho Rispoli. II. Universidade de Brasília. III. Faculdade UnB Gama. IV. Análise Comparativa de Arquiteturas Convolucionais e Baseadas em Atenção para Detecção Automatizada de Parasitas Intestinais

CDU 02:141:005.6

---

Eduardo Vieira Lima

# **Análise Comparativa de Arquiteturas Convolucionais e Baseadas em Atenção para Detecção Automatizada de Parasitas Intestinais**

Monografia submetida ao curso de graduação em (Engenharia de Software) da Universidade de Brasília, como requisito parcial para obtenção do Título de Bacharel em (Engenharia de Software).

---

**Dr. Vinicius de Carvalho Rispoli**  
Orientador

---

**Dr. Ronni Geraldo Gomes de Amorim**

---

**Dr. Marcus Vinicius Chaffim Costa**

Brasília, DF  
julho de 2025

# Agradecimentos

Tenho que começar agradecendo à minha noiva, Thaís, por ser a minha maior fonte de apoio há 8 anos. Não chegaria até aqui de outra forma, sou muito grato por dividir a vida com você.

A meus irmãos, Marcela, Luisa, Gustavo e José Henrique. Grande parte da minha determinação e ambição vem da vontade de inspirar e apoiar vocês, não pude desistir enquanto tive isso em mente.

Também ao meu amigo João Marcelo, que foi essencial para minha trajetória desde o ensino médio e principalmente durante toda a graduação. Sem sua ajuda de Cálculo 1 até este trabalho, eu jamais conseguiria chegar aqui.

Importante não deixar de agradecer à minha família, em especial meu pai Rubens e minha avó Nenza, por terem me incentivado e apoiado durante toda minha vida. Vocês foram grandes inspirações para mim.

Ao meu orientador, Vinicius de Carvalho Rispoli, que me acolheu para a execução deste trabalho quando eu já não sabia se iria conseguir.

Não obstante, a todos os meus outros amigos que estiveram comigo nessa jornada, em especial Luís e Samuel, que estiveram comigo desde o início e seguem comigo hoje.

A todo o corpo docente e de funcionários da UnB, que diariamente se esforçam pelo ensino de excelência nas universidades públicas.



# Resumo

A detecção e classificação de parasitas intestinais através de microscopia óptica representam um desafio significativo na área médica, especialmente em regiões com recursos limitados, sendo a análise manual o padrão atual. Este trabalho se propôs a analisar e comparar as principais soluções de classificação automática no campo da visão computacional, redes convolucionais e transformadores visuais, apresentando também uma abordagem híbrida entre elas. Foram exploradas soluções que priorizam a eficiência computacional: *EfficientNetV2-S* para a rede convolucional e *Tiny ViT* para o transformador visual, tendo em vista cenários com recursos limitados. Foi utilizado para treinamento e avaliação um dos maiores conjuntos de imagens microscópicas de parasitas disponível publicamente, o *Chula-ParasiteEgg-11*, contendo 2.200 imagens de 11 espécies distintas de parasitas. Os resultados obtidos demonstram que as soluções analisadas são viáveis e possuem acurácia satisfatória, sendo a rede convolucional a que melhor performa em termos de eficiência e acurácia, seguida do modelo híbrido e então do transformador visual. Observou-se que as soluções possuem resultados razoavelmente inferiores aos modelos mais robustos da literatura, porém utilizam significativamente menos recursos computacionais, o que as torna viáveis para situações em que o recurso é limitado.

**Palavras-chaves:** Aprendizado profundo. Classificação de parasitas. Redes convolucionais. Transformadores visuais. Modelo híbrido.



# Abstract

The detection and classification of intestinal parasites through optical microscopy represents a significant challenge in the medical field, especially in resource-limited regions, where manual analysis remains the current standard. This work aimed to analyze and compare the main automatic classification solutions in the field of computer vision, convolutional networks and visual transformers, also presenting a hybrid approach between them. Solutions that prioritize computational efficiency were explored: *EfficientNetV2-S* for the convolutional network and *Tiny ViT* for the visual transformer, considering scenarios with limited resources. One of the largest publicly available microscopic parasite image datasets was used for training and evaluation, the *Chula-ParasiteEgg-11*, containing 2,200 images of 11 distinct parasite species. The obtained results demonstrate that the analyzed solutions are viable and have satisfactory accuracy, with the convolutional network performing best in terms of efficiency and accuracy, followed by the hybrid model and then the visual transformer. It was observed that the solutions have reasonably inferior results compared to the more robust models in the literature, however they use significantly fewer computational resources, making them viable for situations where resources are limited.

**Key-words:** Deep learning. Parasite classification. Convolutional networks. Visual transformers. Hybrid model.



# Lista de ilustrações

Figura 1 – Ilustração de um modelo de rede neural rasa com uma camada de entrada, saída e uma única camada oculta intermediária. Fonte: Nielsen (2015) . . . . .	26
Figura 2 – Ilustração de um modelo de rede neural profunda com uma camada de entrada, saída e múltiplas camadas ocultas intermediárias. Fonte: Nielsen (2015) . . . . .	26
Figura 3 – Gráfico que representa as funções de ativação mais comumente utilizadas. Cada função tem o papel de introduzir a não-linearidade ao modelo, a escolha de uma varia de acordo com a necessidade do problema abordado. Fonte: Leppich (2021) . . . . .	27
Figura 4 – Representação visual de como uma rede convolucional identifica padrões para, neste caso, reconhecer a imagem de um gato ( <i>cat</i> , em inglês). De baixo para cima na imagem, linhas elementares ou texturas combinam-se em objetos simples como olhos ou orelhas, que então combinam em conceitos mais complexos como “gato”. Assim funciona a hierarquia espacial. Fonte: Chollet (2021) . . . . .	29
Figura 5 – Visão geral do transformador visual. A imagem é dividida em blocos de tamanho fixo, que são então submetidos a uma projeção linear com informação de posição e alimentados a um codificador para um transformador padrão. Fonte: Dosovitskiy et al. (2020) . . . . .	30
Figura 6 – Estrutura de destilação rápida do <i>Tiny ViT</i> . A parte superior mostra a ramificação para salvar os <i>logits</i> do professor, onde o tratamento de dados codificado e os <i>logits</i> esparsificados do professor são salvos. A parte do meio representa o disco para armazenar as informações. A parte inferior mostra a ramificação para treinar o estudante, onde o decodificador reconstrói o tratamento de dados e a destilação é conduzida entre os <i>logits</i> do professor e as saídas do estudante. As duas ramificações são independentes e assíncronas, permitindo treinamento sem processar o modelo professor grande durante cada iteração. Fonte: Wu et al. (2022) . . . . .	39
Figura 7 – Curvas de treinamento ao longo do tempo do modelo <i>EfficientNetV2-S</i> , apresentando perda e acurácia. O tempo é dado em épocas. Fonte: Elaborado pelo autor (2025). . . . .	47
Figura 8 – Curvas de treinamento ao longo do tempo do modelo <i>Tiny Vision Transformer</i> , apresentando perda e acurácia. O tempo é dado em épocas. Fonte: Elaborado pelo autor (2025). . . . .	47

Figura 9 – Curvas de treinamento ao longo do tempo do modelo híbrido, apresentando perda e acurácia. O tempo é dado em épocas. Fonte: Elaborado pelo autor (2025).	48
Figura 10 – Comparação da acurácia de teste entre o modelo de rede convolucional <i>EfficientNetV2-S</i> , transformador visual <i>Tiny ViT</i> e híbrido. Fonte: Elaborado pelo autor (2025).	48
Figura 11 – Análise detalhada da acurácia de teste do modelo <i>EfficientNetV2-S</i> . Apresenta acurácia e perda geral, precisão, revocação e <i>F1-Score</i> por classe específica. Fonte: Elaborado pelo autor (2025).	49
Figura 12 – Análise detalhada da acurácia de teste do modelo <i>Tiny Vision Transformer</i> . Apresenta acurácia e perda geral, precisão, revocação e <i>F1-Score</i> por classe específica. Fonte: Elaborado pelo autor (2025).	49
Figura 13 – Análise detalhada da acurácia de teste do modelo híbrido. Apresenta acurácia e perda geral, precisão, revocação e <i>F1-Score</i> por classe específica. Fonte: Elaborado pelo autor (2025).	50
Figura 14 – Matriz de confusão do modelo <i>EfficientNetV2-S</i> . Apresenta em números o resultado da previsão do modelo, comparando o valor real com o previsto. Fonte: Elaborado pelo autor (2025).	50
Figura 15 – Matriz de confusão do modelo <i>Tiny Vision Transformer</i> . Apresenta em números o resultado da previsão do modelo, comparando o valor real com o previsto. Fonte: Elaborado pelo autor (2025).	51
Figura 16 – Matriz de confusão do modelo híbrido. Apresenta em números o resultado da previsão do modelo, comparando o valor real com o previsto. Fonte: Elaborado pelo autor (2025).	52
Figura 17 – Comparação entre base (esquerda) e inferência do modelo de rede convolucional (direita) para a classe <i>Ascaris lumbricoides</i> . Fonte: Elaborado pelo autor (2025).	53
Figura 18 – Comparação entre base (esquerda) e inferência do modelo <i>Tiny Vision Transformer</i> (direita) para a classe <i>Hymenolepis nana</i> . Fonte: Elaborado pelo autor (2025).	55

# Lista de tabelas

Tabela 1 – Hiperparâmetros utilizados no experimento . . . . .	43
Tabela 2 – Comparação da eficiência computacional dos modelos . . . . .	52



# Lista de abreviaturas e siglas

CNN	<i>Convolutional Neural Network</i> (Rede Neural Convolucional)
CPU	<i>Central Processing Unit</i> (Unidade Central de Processamento)
GPU	<i>Graphics Processing Unit</i> (Unidade de Processamento Gráfico)
HSV	<i>Hue-Saturation-Value</i> (Matiz-Saturação-Valor)
IA	Inteligência Artificial
RGB	<i>Red, Green, Blue</i> (Vermelho, Verde, Azul)
ReLU	<i>Rectified Linear Unit</i> (Unidade Linear Retificada)
GELU	<i>Gaussian Error Linear Unit</i> (Unidade Linear de Erro Gaussiano)
ViT	<i>Vision Transformer</i> (Transformador Visual)
ms	Milissegundos
µm	Micrômetros



# Lista de símbolos

$\alpha$	Coeficiente de profundidade
$\beta$	Coeficiente de largura
$\gamma$	Coeficiente de resolução
$\phi$	Coeficiente composto de escalonamento
$\mathcal{L}$	Função de perda
$\Sigma$	Soma
$\log(x)$	Logaritmo natural
$\mathbb{R}$	Conjunto dos números reais
$\mathbf{x}$	Vetor de entrada
$\mathbf{y}$	Vetor de saída
$\mathbf{W}$	Matriz de pesos
$\mathbf{b}$	Vetor de <i>bias</i>
$\in$	Pertence a



# Sumário

	<b>Introdução</b>	<b>19</b>
<b>1</b>	<b>REFERENCIAL TEÓRICO</b>	<b>23</b>
1.1	Aprendizado de Máquina	23
1.2	Redes Neurais	25
1.2.1	Redes neurais convolucionais	28
1.3	Transformadores visuais	29
1.4	Visão computacional aplicada à detecção de parasitas	31
<b>2</b>	<b>METODOLOGIA</b>	<b>33</b>
2.1	Dados	33
2.1.1	Características gerais do conjunto de dados	33
2.1.2	Características Morfológicas dos Ovos	34
2.1.3	Coleta de dados	34
2.1.4	Tratamento do conjunto de dados	34
2.2	Função de Perda	35
2.3	Arquiteturas implementadas	36
2.3.1	Rede convolucional: <i>EfficientNetV2-S</i>	36
2.3.2	Transformador visual: <i>Tiny ViT</i>	38
2.3.3	Modelo híbrido: <i>EfficientNetV2-S + Tiny ViT</i>	39
2.4	Experimento	41
2.4.1	Ferramentas utilizadas	41
2.4.2	Configuração experimental	41
2.4.3	Análise comparativa	44
<b>3</b>	<b>RESULTADOS E DISCUSSÃO</b>	<b>47</b>
<b>4</b>	<b>CONCLUSÃO</b>	<b>57</b>
	<b>REFERÊNCIAS</b>	<b>59</b>
	<b>APÊNDICES</b>	<b>63</b>
	<b>APÊNDICE A – CÓDIGO FONTE</b>	<b>65</b>
A.1	Arquivo Principal: <i>run_optimized.py</i>	65
A.2	Treinador Otimizado: <i>trainer_optimized.py</i>	73

A.3	Configuração: <i>config_optimized.py</i> . . . . .	85
A.4	Modelos: <i>models.py</i> . . . . .	87
A.5	Conjunto de Dados: <i>dataset_optimized.py</i> . . . . .	91
A.6	Requirements: <i>requirements.txt</i> . . . . .	97

# Introdução

Inteligência Artificial (IA) é um ramo da tecnologia que pode ser descrito como o esforço para automatizar tarefas intelectuais normalmente executadas por humanos [Chollet \(2021\)](#). Esse campo tem evoluído rapidamente nas últimas décadas, impulsionado pelo crescimento exponencial na capacidade computacional, disponibilidade de grandes volumes de dados e avanços em algoritmos de aprendizado de máquina. Nos últimos anos, a IA tem sido objeto de intenso interesse geral, expandindo-se para contextos além das áreas tradicionais de tecnologia, como medicina, finanças, indústria automotiva e segurança cibernética [Chollet \(2021\)](#).

Dentre os principais campos da IA, o *Machine Learning* (ML) se destaca por permitir que sistemas aprendam padrões e tomem decisões com base em dados [Mitchell \(1997\)](#). Essa abordagem elimina a necessidade de programação explícita para cada cenário, tornando possível a adaptação e melhoria contínua dos modelos conforme novos dados são incorporados. Dentro do ML, o Aprendizado Profundo (do inglês, *Deep Learning*), que utiliza redes neurais profundas, tem impulsionado avanços significativos, especialmente no processamento de imagens e linguagem natural, devido à sua capacidade de identificar padrões complexos em grandes volumes de dados.

As Redes Neurais Convolucionais (CNNs) representam um dos avanços mais importantes no campo do *Deep Learning*. Inspiradas no funcionamento do córtex visual, as CNNs são arquiteturas especializadas para processamento de imagens, utilizando camadas convolucionais para extrair características relevantes, reduzindo a necessidade de intervenção manual na definição de características. Essas redes têm sido amplamente empregadas em tarefas como reconhecimento facial, diagnóstico por imagem e veículos autônomos, permitindo classificações e segmentações com alta precisão [Lecun et al. \(1998\)](#).

Mais recentemente, os *Vision Transformers* (ViTs) surgiram como uma alternativa promissora às CNNs para visão computacional. Diferentemente das arquiteturas tradicionais baseadas em convolução, os ViTs utilizam mecanismos de autoatenção para capturar relações globais entre pixels em imagens. Essa abordagem permite um aprendizado mais eficiente de representações visuais e, em muitos casos, supera o desempenho das CNNs em diversas tarefas específicas [Dosovitskiy et al. \(2020\)](#). O avanço dos ViTs tem ampliado significativamente o potencial das aplicações de IA, trazendo novas possibilidades para reconhecimento de padrões visuais e reforçando a importância dos modelos baseados em atenção na área de aprendizado profundo.

No campo da medicina, as infecções parasitárias intestinais continuam sendo um problema significativo de saúde pública, especialmente em países em desenvolvimento e regiões com infraestrutura sanitária deficiente. Estima-se que cerca de 24% da população mundial seja afetada por doenças infecciosas e parasitárias, com impacto predominante

em crianças, gestantes e pessoas com baixa imunidade [Kumar et al. \(2023\)](#). Essas infecções podem resultar em sintomas como diarreia, desnutrição, anemia, fraqueza e comprometimento do desenvolvimento infantil.

O diagnóstico tradicional dessas infecções é realizado por meio de exames microscópicos de amostras fecais, considerados o padrão-ouro por sua capacidade de identificar diretamente ovos e parasitas. Contudo, esse processo é trabalhoso, demorado, exige pessoal altamente capacitado e está sujeito a variações na acurácia devido à subjetividade da análise e às variações morfológicas dos parasitas [Xu et al. \(2024\)](#). Além disso, o ambiente de trabalho pode ser insalubre, com baixa eficiência e alto volume de trabalho para os profissionais de laboratório.

Nesse cenário, a aplicação de modelos baseados em IA tem-se mostrado uma alternativa promissora para automatizar e melhorar o processo de detecção e classificação de parasitas intestinais. A integração de tecnologias de processamento de imagens digitais com métodos de aprendizado profundo, especialmente com redes neurais convolucionais e, mais recentemente, com Vision Transformers, tem impulsionado o desenvolvimento de sistemas automáticos capazes de identificar com maior precisão e agilidade os agentes infecciosos presentes em imagens microscópicas [Kumar et al. \(2023\)](#).

Diversos trabalhos prévios têm explorado o uso de aprendizado profundo para essa tarefa. [Kumar et al. \(2023\)](#) propuseram uma abordagem utilizando YOLOv5, demonstrando elevada acurácia (cerca de 97%) e alta velocidade de detecção em um conjunto de dados com mais de cinco mil imagens de parasitas intestinais. O estudo destacou a viabilidade do uso de modelos leves em ambientes com poucos recursos computacionais, voltando-se ao diagnóstico clínico em tempo real.

[Xu et al. \(2024\)](#) desenvolveram a YAC-Net, uma arquitetura leve baseada em YOLOv5, com otimizações estruturais específicas para imagens de ovos parasitários. Os autores demonstraram que sua abordagem melhora a detecção mesmo em imagens com baixa resolução e com ruídos, reduzindo o número de parâmetros e o custo computacional, o que viabiliza sua aplicação em regiões remotas.

Outro trabalho relevante foi apresentado por [AïDahoul et al. \(2023\)](#), que explorou a combinação de redes convolucionais com mecanismos de atenção por meio da arquitetura CoAtNet, demonstrando que a combinação dessas abordagens pode gerar resultados superiores na classificação de ovos parasitários. O estudo também comparou diferentes abordagens de CNNs e Vision Transformers em um conjunto diverso de imagens, incluindo variações de iluminação e resolução.

Apesar das contribuições relevantes dessas pesquisas, poucas investigações exploraram comparativamente a aplicação de CNNs e Vision Transformers de maneira sistemática e com base em um mesmo conjunto de dados para a detecção e classificação de parasitas intestinais. Neste contexto, esta monografia tem como objetivo investigar e comparar o desempenho de modelos baseados em CNNs e Vision Transformers na tarefa de detec-

ção e classificação de parasitas intestinais em imagens microscópicas, contribuindo para o desenvolvimento de ferramentas automatizadas e eficientes para a saúde pública.



# 1 Referencial Teórico

Este capítulo apresenta a fundamentação teórica deste trabalho. Inicia-se introduzindo o conceito de Aprendizado de Máquina (*Machine Learning*) e Aprendizado Profundo (*Deep Learning*). Em seguida, apresenta as arquiteturas de Aprendizado de Máquina abordadas neste trabalho, Redes Neurais Convolucionais (*Convolutional Neural Networks* ou CNNs) e Transformadores Visuais (*Vision transformers* ou ViTs). Então, apresenta uma visão geral sobre a Visão Computacional aplicada à detecção de parasitas intestinais, assim como o contexto geral do problema para o campo da medicina e a importância das soluções exploradas neste trabalho. Por fim, apresenta brevemente os trabalhos relacionados relevantes para este contexto.

## 1.1 Aprendizado de Máquina

Os programas construídos em Aprendizado de Máquina possuem, de forma geral, três componentes: dados de entrada, exemplos da saída esperada e um meio de quantificar a qualidade das saídas geradas [Chollet \(2021\)](#).

- Dados de entrada: são o objeto básico do problema. Por exemplo, os dados de entrada para um sistema de classificação de imagens seriam figuras ainda não classificadas.
- Exemplos da saída esperada: seguindo o exemplo anterior, seria um conjunto de imagens corretamente classificadas, como “gato” ou “cachorro”.
- um meio de quantificar a qualidade das saídas geradas: é o principal em um processo de aprendizado de máquina. Por mensurar a qualidade dos resultados da tarefa que realiza, o sistema é capaz de corrigir-se e ajustar como seu algoritmo funciona. Este processo de ajuste é o que caracteriza o aprendizado.

O processo principal de um modelo de aprendizado de máquina é transformar os dados de entrada em saídas significativas para o problema que se propõe a solucionar. Este processo é aprendido através da exposição a exemplos conhecidos de entrada e saída esperada. Assim, o problema central torna-se encontrar maneiras de transformar o dado de entrada para se aproximar da saída desejada, em outras palavras, aprender representações do dado de entrada que sejam úteis para alcançar a saída esperada. Por exemplo, um programa que tem como objetivo identificar a presença da cor vermelha em uma imagem poderia representar esta imagem no formato RGB (do inglês, *red-green-blue*). Por sua vez, se o objetivo for saturar as cores de uma imagem, o melhor seria representá-la no formato HSV (do inglês, *hue-saturation-value*) [Chollet \(2021\)](#).

Transformar os dados de entrada em representações relevantes não é um processo simples, pois existem numerosas formas de representar um dado e poucas serão úteis para se aproximar do resultado desejado. Assim, os programas de aprendizado de máquina utilizam um conjunto de transformações pré-definidas que melhor se aproximam do contexto do problema, utilizando o processo de quantificar a qualidade das saídas para avaliar o quão relevante é cada transformação dentro deste conjunto. Este conjunto pré-definido é chamado de espaço de hipótese [Chollet \(2021\)](#). Em suma, o processo de aprendizado de máquina é composto por: encontrar transformações de dados de entrada, dentro de um conjunto pré-definido de possibilidades, utilizando como guia a qualidade das saídas geradas e retroalimentando o sistema com cada resultado gerado para se aproximar cada vez mais da saída desejada. Esta simples ideia possibilita a execução de diversas tarefas [Chollet \(2021\)](#).

O aprendizado profundo (do inglês *Deep Learning*) é uma área do aprendizado de máquina que tem ganhado força e crescido exponencialmente desde os anos 90. O foco desta abordagem está na transformação dos dados de entrada, introduzindo camadas consecutivas de transformações cada vez mais significativas. O “profundo” em aprendizado profundo refere-se à quantidade de camadas de transformação, o número de camadas representa a “profundidade” do modelo. Modelos modernos de aprendizado profundo podem incluir dezenas ou centenas de camadas, todas incluídas automaticamente pelo processo de aprendizado. Esta é a principal diferença entre um modelo tradicional, que geralmente foca em trabalhar com uma ou duas camadas de transformação [Chollet \(2021\)](#).

As transformações em cada camada são mediadas por pesos (ou parâmetros), que são ajustados durante o aprendizado. Os pesos de uma camada representam o que a camada faz em termos da transformação dos dados, quantificado numericamente. Este ajuste ocorre comparando a saída do algoritmo com a saída real através de uma função de perda; esta função representa numericamente o quão distante a saída gerada está da saída real desejada, gerando a perda. O resultado é então retroalimentado aos pesos por meio de um otimizador, buscando minimizar a perda. O programa inicia com valores aleatórios para os pesos e busca otimizar seus valores através do processo de aprendizado retroalimentado [Chollet \(2021\)](#).

O aprendizado de máquina pode ser classificado em três principais categorias: supervisionado, não supervisionado e por reforço. No aprendizado supervisionado, o modelo é treinado com um conjunto de dados rotulado, aprendendo a mapear entradas para saídas com base em exemplos prévios, sendo amplamente utilizado em tarefas como classificação e regressão [Goodfellow, Bengio e Courville \(2016\)](#). Já o aprendizado não supervisionado trabalha com dados sem rótulos, buscando identificar padrões ou estruturas ocultas nos dados, como agrupamentos ou redução de dimensionalidade [Murphy \(2012\)](#). Por fim, o aprendizado por reforço envolve a interação de um agente com um ambiente, aprendendo por meio de recompensas e penalidades a tomar decisões que maximizem um retorno

cumulativo ao longo do tempo [Sutton e Barto \(2015\)](#). Esses três paradigmas oferecem abordagens distintas para a construção de modelos inteligentes, sendo escolhidos conforme a natureza do problema e dos dados disponíveis.

O aprendizado profundo está incluído dentro do que é conhecido como redes neurais, uma subcategoria do aprendizado de máquina. O termo rede neural se refere à neurobiologia, inspirado pelo entendimento de como um cérebro funciona (principalmente o córtex visual), onde as camadas de transformações funcionam como neurônios. É importante destacar que, apesar da inspiração, os modelos de aprendizado profundo não são modelos do funcionamento do cérebro, tendo em vista que não há nenhuma evidência de que o cérebro utilize algo semelhante aos mecanismos e estratégias utilizados no aprendizado profundo [Chollet \(2021\)](#). Este trabalho utilizará redes neurais alinhadas com aprendizado supervisionado para a detecção e classificação dos principais parasitas intestinais, a partir de exames de fezes digitalizados.

## 1.2 Redes Neurais

No centro de todas as operações e representações de dados em redes neurais estão os tensores. O tensor é a estrutura fundamental em sistemas de aprendizado de máquina, capaz de armazenar dados e suas transformações de forma numérica em múltiplas dimensões, como escalares, vetores e matrizes [Chollet \(2021\)](#). Todas as operações dentro de uma rede neural, incluindo as transformações lineares (como produtos escalares e adições), são operações entre tensores, e os próprios parâmetros do modelo, como pesos e viés, são representados como tensores [Paszke et al. \(2019\)](#).

Toda rede neural é fundamentalmente um encadeamento de camadas, entre seus tensores de entrada e saída [Chollet \(2021\)](#). Estas camadas que existem entre a entrada e a saída de uma rede são conhecidas como camadas ocultas (do inglês *hidden layer*). Uma rede neural pode ser classificada como rasa ou profunda a partir de seu número de camadas ocultas, sendo a rede rasa uma que possua uma única camada oculta e as redes profundas são as que possuem múltiplas [Paszke et al. \(2019\)](#). A Figura 1 representa um exemplo de rede neural rasa, contendo apenas uma única camada oculta, enquanto a Figura 2 apresenta uma rede profunda, com múltiplas camadas ocultas.

Cada camada atua como um bloco de construção fundamental, extraindo representações mais significativas dos dados a cada estágio [Chollet \(2021\)](#). Estas camadas são tipicamente vetoriais e são frequentemente chamadas de neurônios ou unidades ocultas [Prince \(2024\)](#), que agem em paralelo representando uma função de vetor para escalar, e na prática realizam uma transformação linear na entrada (envolvendo os pesos e viés) [Paszke et al. \(2019\)](#).

Outro ponto fundamental para as redes neurais são as funções de ativação. Estas funções são cruciais para que seja possível que a rede neural possa lidar com relações de

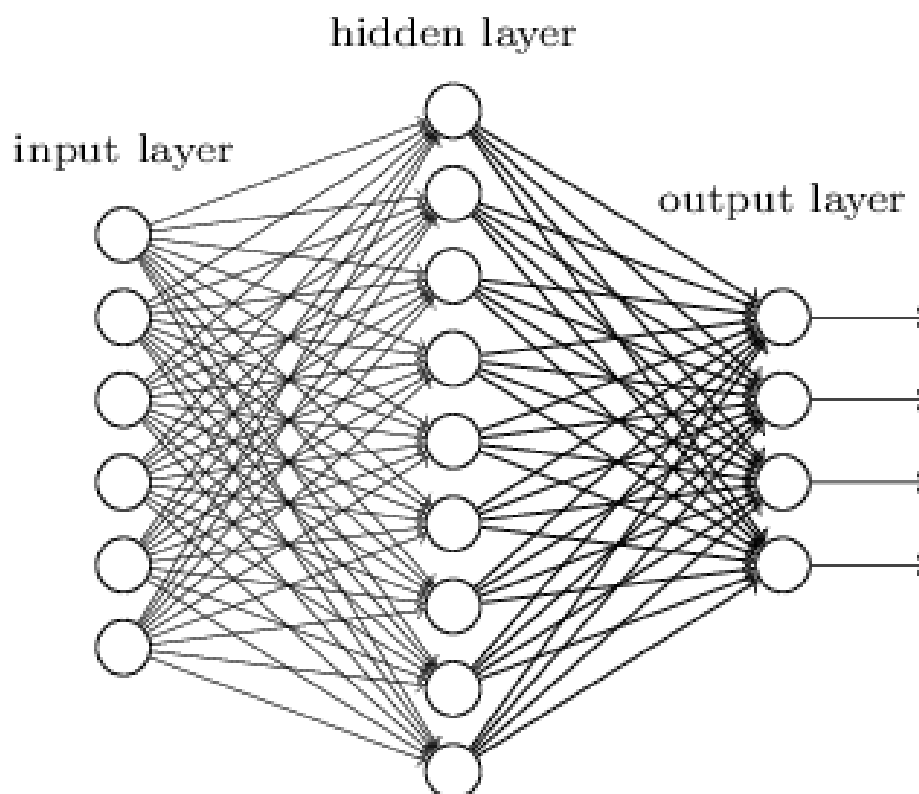


Figura 1 – Ilustração de um modelo de rede neural rasa com uma camada de entrada, saída e uma única camada oculta intermediária. Fonte: [Nielsen \(2015\)](#)

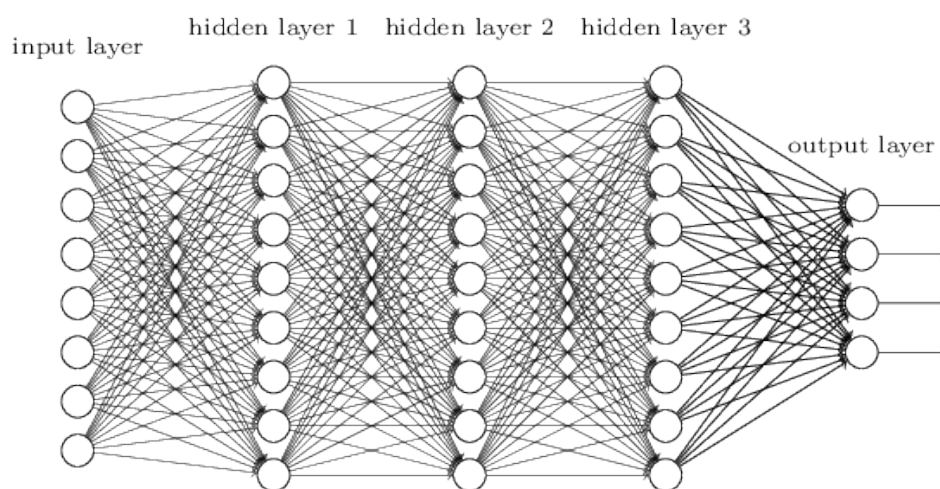


Figura 2 – Ilustração de um modelo de rede neural profunda com uma camada de entrada, saída e múltiplas camadas ocultas intermediárias. Fonte: [Nielsen \(2015\)](#)

dados mais complexas e não lineares [Paszke et al. \(2019\)](#). Sem uma função de ativação, uma sequência de operações lineares (produto escalar e adição) resultaria em um modelo que seria, em sua essência, linear, independentemente do número de camadas [Chollet \(2021\)](#). As funções de ativação introduzem a não linearidade que permite que a rede se aproxime de uma vasta gama de funções complexas. Existe uma grande variedade de funções de ativação que já são provadas com sucesso, assim, a escolha de uma função de ativação varia de acordo com o problema abordado pela rede neural [Paszke et al. \(2019\)](#). A Figura 3 apresenta as funções de ativação mais comumente utilizadas.

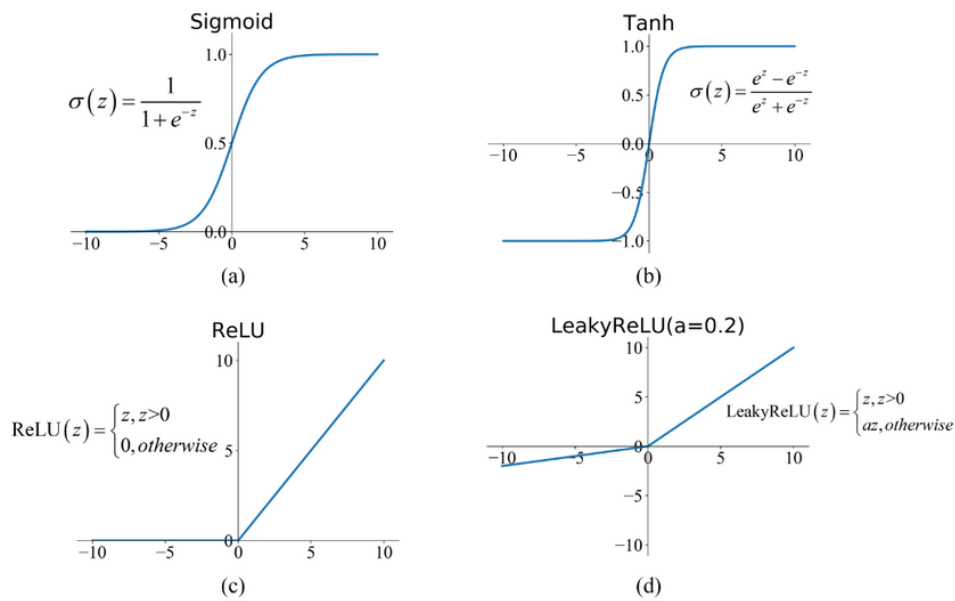


Figura 3 – Gráfico que representa as funções de ativação mais comumente utilizadas. Cada função tem o papel de introduzir a não-linearidade ao modelo, a escolha de uma varia de acordo com a necessidade do problema abordado. Fonte: [Leppich \(2021\)](#)

- **ReLU (*Rectified Linear Unit*)**: é considerada a escolha mais comum e uma das funções de ativação de propósito geral com melhor desempenho [Prince \(2024\)](#). Ela retorna o valor de entrada se for positivo e zero caso contrário, efetivamente equalizando os valores negativos a zero. Possui derivada 1 para entradas positivas e 0 para entradas negativas, o que contribui para a estabilidade e eficiência do treinamento [Prince \(2024\)](#).
- **Leaky ReLU**: uma variação da função ReLU. Permite uma pequena inclinação positiva para valores negativos (tipicamente 0.01 ou 0.1) para mitigar o problema do “ReLU morrendo”. Este problema ocorre por uma característica inerente à função ReLU: sua derivada é sempre zero para entradas negativas. Isto faz com que, se todos os exemplos de treinamento produzirem uma entrada negativa para a função ReLU, então não será possível ajustar os parâmetros das camadas, já que sempre será zero. [Prince \(2024\)](#).

- Sigmoid (função logística): transforma os valores de entrada para o intervalo definido de 0 a 1. Foi amplamente utilizada nos primórdios do aprendizado profundo, mas atualmente é mais usada quando a saída precisa ser interpretada como uma probabilidade [Paszke et al. \(2019\)](#).
- Tanh (tangente hiperbólica): funciona de forma similar à função *Sigmoid*, porém seu intervalo é de -1 a 1 [Paszke et al. \(2019\)](#).

### 1.2.1 Redes neurais convolucionais

Rede neural convolucional é um tipo de rede neural utilizada principalmente para tarefas de visão computacional. A diferença fundamental entre redes tradicionais e uma rede convolucional está na estrutura de camada. Uma camada tradicional aprende padrões globais a partir da entrada de dados, enquanto uma camada convolucional aprende padrões locais, no caso de imagens, padrões encontrados em pequenos trechos bidimensionais das entradas [Chollet \(2021\)](#).

Esta característica dá às redes convolucionais duas propriedades importantes:

- Após aprender um padrão em algum trecho de uma imagem, a rede é capaz de reconhecer esse padrão em qualquer lugar. Uma camada tradicional teria que aprender este padrão novamente caso aparecesse em algum outro trecho [Chollet \(2021\)](#). Isto torna as redes convolucionais especialmente efetivas no uso de dados em tarefas de visão computacional, precisando de uma base de dados de treinamento menor [Chollet \(2021\)](#).
- Redes convolucionais são capazes de aprender hierarquias espaciais. Uma primeira camada pode aprender pequenos padrões locais, como vértices, uma segunda camada subsequente pode aprender padrões mais gerais a partir das características da primeira camada e assim sucessivamente. Isto permite o aprendizado de conceitos visuais profundos e complexos, já que o mundo visual funciona fundamentalmente em uma lógica de hierarquia espacial [Chollet \(2021\)](#). A Figura 4 apresenta um exemplo simples de como a hierarquia espacial funciona no aprendizado de máquina.

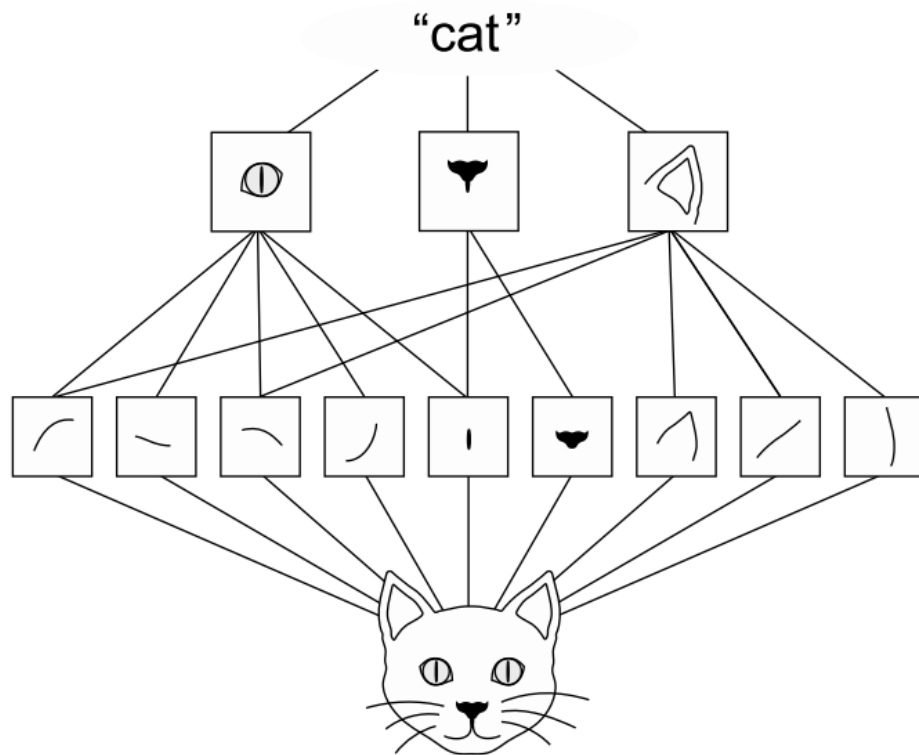


Figura 4 – Representação visual de como uma rede convolucional identifica padrões para, neste caso, reconhecer a imagem de um gato (*cat*, em inglês). De baixo para cima na imagem, linhas elementares ou texturas combinam-se em objetos simples como olhos ou orelhas, que então combinam em conceitos mais complexos como “gato”. Assim funciona a hierarquia espacial. Fonte: [Chollet \(2021\)](#)

### 1.3 Transformadores visuais

Arquiteturas baseadas em atenção tornaram-se o modelo dominante para o processamento de linguagem natural, causando uma revolução nos modelos linguísticos de inteligência artificial. Em especial, os transformadores são a principal solução de arquitetura baseada em atenção [Dosovitskiy et al. \(2020\)](#).

O transformador foi apresentado no artigo [Vaswani et al. \(2017\)](#). Antes disso, o modelo dominante para modelos linguísticos era uma abordagem de redes neurais convolucionais que incluía um codificador e um decodificador através de um mecanismo de atenção [Vaswani et al. \(2017\)](#). A ideia principal do transformador é focar nos mecanismos de atenção, dispensando completamente as abordagens convolucionais. De forma simplificada, uma função de atenção pode ser descrita como uma função que mapeia uma consulta e um conjunto de pares chave-valor para uma saída, onde a consulta, chaves, valores e saída são vetores. A saída é então computada como uma soma ponderada dos valores, onde o peso designado para cada valor é computado por uma função de compatibilidade da consulta com a chave correspondente [Vaswani et al. \(2017\)](#).

Explorando diferentes abordagens baseadas puramente em atenção, o trabalho Vaswani et al. (2017) alcançou resultados superiores às arquiteturas convolucionais, até então dominantes no campo de modelos linguísticos. Inspirado nisso, nascem os transformadores visuais, com a proposta de abordar a computação visual com a mesma ideia de focar nos modelos de atenção e deixar de lado as arquiteturas convolucionais.

Para utilizar o conceito de transformadores em imagens, a imagem é dividida em blocos e projetada linearmente como entradas para um transformador. Os blocos de imagem codificados são tratados da mesma forma que palavras são tratadas em transformadores de linguagem natural Dosovitskiy et al. (2020). A Figura 5 ilustra a visão geral da estrutura de um transformador visual.

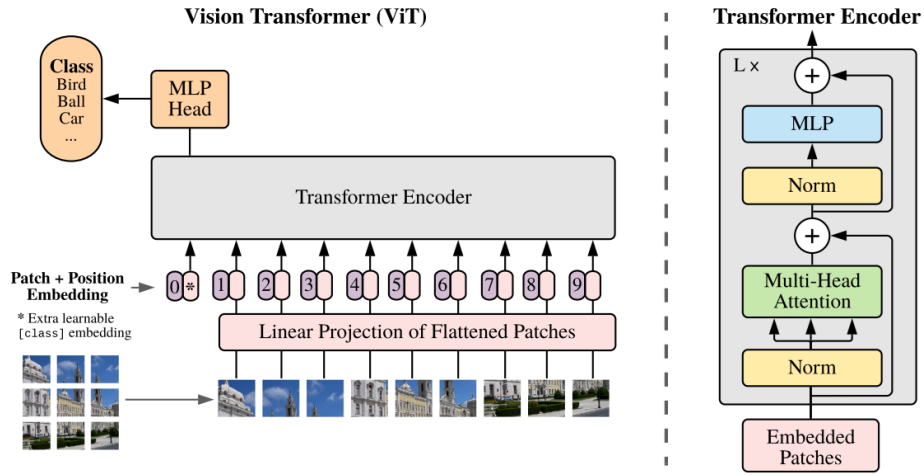


Figura 5 – Visão geral do transformador visual. A imagem é dividida em blocos de tamanho fixo, que são então submetidos a uma projeção linear com informação de posição e alimentados a um codificador para um transformador padrão. Fonte: Dosovitskiy et al. (2020)

O transformador padrão recebe como entrada uma sequência unidimensional de vetores projetados dos blocos. Para tratar imagens bidimensionais, a imagem  $x \in \mathbb{R}^{H \times W \times C}$  é reestruturada em uma sequência de blocos bidimensionais achatados  $x_p \in \mathbb{R}^{N \times (P^2 \cdot C)}$ , onde  $(H, W)$  representa a resolução da imagem original,  $C$  é o número de canais de imagem,  $(P, P)$  corresponde à resolução de cada bloco da imagem, o número total de blocos é então dado por

$$N = \frac{H \cdot W}{P^2}, \quad (1.1)$$

que também define o comprimento efetivo da sequência de entrada para o transformador Dosovitskiy et al. (2020).

Os canais de imagem representam as diferentes componentes de cor ou características da imagem. Em imagens coloridas no formato *RGB*, existem três canais, cada um representando a intensidade de uma cor primária. Em imagens em escala de cinza, há apenas um canal representando a intensidade luminosa. Em contextos de processamento

de imagem médica, podem existir canais adicionais representando diferentes técnicas de coloração ou modalidades de imagem. O conceito de canais é fundamental para o processamento de imagens, pois permite que a rede neural aprenda padrões específicos de cada componente da imagem separadamente.

Adiciona-se no início da sequência um vetor aprendido, denominado bloco de classificação, representado como

$$z_0^0 = x_{\text{class}}, \quad (1.2)$$

cuja saída final, após o processamento pelo codificador do transformador, é utilizada como representação da imagem. Tanto no pré-treinamento quanto no ajuste fino, um cabeçalho de classificação é acoplado à saída  $z_L^0$  correspondente a esse bloco. Esse cabeçalho é implementado por uma rede *perceptron* multicamadas (*MLP*) com uma camada oculta durante o pré-treinamento e por uma camada linear simples durante o ajuste fino [Dosovitskiy et al. \(2020\)](#).

Informações posicionais são incorporadas aos vetores dos blocos para preservar a relação espacial da imagem. São utilizadas informações posicionais unidimensionais aprendíveis, uma vez que não foram observados ganhos de desempenho significativos ao empregar técnicas mais avançadas com consciência bidimensional. A sequência resultante de vetores serve como entrada para o codificador do transformador [Dosovitskiy et al. \(2020\)](#).

O transformador mantém um tamanho constante  $D$  para os vetores latentes em todas as suas camadas. Dessa forma, os blocos são achatados e projetados para esse espaço vetorial de dimensão  $D$  por meio de uma projeção linear treinável, representada como

$$x_p \in \mathbb{R}^{N \times (P^2 \cdot C)} \xrightarrow{\text{Projeção Linear}} \mathbb{R}^{N \times D}, \quad (1.3)$$

cuja saída corresponde aos vetores projetados dos blocos, que são efetivamente a entrada do codificador do transformador [Dosovitskiy et al. \(2020\)](#).

Quando treinado em bases de dados de tamanho médio e sem uma forte regularização, o transformador visual alcança uma acurácia modesta e um pouco abaixo dos modelos convolucionais. Porém, em grandes bases de dados, o transformador visual alcança resultados iguais ou superiores aos modelos convolucionais e utiliza significativamente menos recursos computacionais [Dosovitskiy et al. \(2020\)](#).

## 1.4 Visão computacional aplicada à detecção de parasitas

Infecções intestinais parasitárias são as infecções mais comuns que afetam as comunidades mais pobres e necessitadas no mundo. Essas infecções são amplamente distribuídas na África Subsaariana, China e leste asiático [Gujo e Kare \(2021\)](#). Mais de 1,5 bilhões de

peças são infectadas por helmintos transmitidos pelo solo, como *Ascaris lumbricoides* e *Trichuris trichiura*, destes, 267 milhões são crianças com menos de 5 anos de idade [Gujo e Kare \(2021\)](#).

O método padrão para diagnosticar infecções por parasitas intestinais é baseado na análise microscópica manual de amostras de fezes, com técnicas como o exame direto, Kato-Katz e testes moleculares (qPCR, ELISA, imunofluorescência) [Yimer et al. \(2015\)](#). Os testes moleculares e Kato-Katz são técnicas que possuem uma sensibilidade superior ao exame direto, alcançando acima de 90% de sensibilidade contra a sensibilidade de 48,9% a 63,1% do método direto. Porém, esses métodos de diagnóstico são significativamente laboriosos e onerosos, além de altamente suscetíveis a erros humanos [Yimer et al. \(2015\)](#).

A visão computacional entra como uma alternativa para os processos manuais de diagnóstico, visando mitigar as limitações de recurso e mão de obra inerentes aos métodos manuais. Inicialmente, a automação do processo de detecção de parasitas era limitada pela falta de algoritmos de reconhecimento para microrganismos. Porém, o avanço do aprendizado de máquina e, em especial, o aprendizado profundo trouxeram uma revolução no campo de identificação e classificação de imagens na medicina, destacando as redes neurais convolucionais [Kumar et al. \(2023\)](#).

## 2 Metodologia

### 2.1 Dados

O experimento deste trabalho utilizou o conjunto de dados *Chula-ParasiteEgg-11* [Palasuwan et al. \(2022\)](#). Esses dados foram disponibilizados publicamente na plataforma *IEEE Dataport* como proposta para a competição ICIP 2022 [Anantrasirichai et al. \(2022\)](#), cujo objetivo foi incentivar soluções para a detecção e classificação de ovos de parasitas intestinais. Os dados disponibilizados representam um dos maiores conjuntos de dados disponíveis publicamente para a classificação de parasitas intestinais, contendo 11 espécies distintas de parasitas com uma distribuição balanceada entre eles.

#### 2.1.1 Características gerais do conjunto de dados

O conjunto de dados utilizado neste trabalho é composto por um total de 2.200 imagens, distribuídas de forma equilibrada, com 200 imagens por espécie. As imagens estão no formato *RGB* e foram obtidas através de microscopia óptica com diferentes técnicas de coloração. Esse conjunto apresenta imagens de alta resolução com variações de iluminação e contraste, características que contribuem para a robustez do conjunto.

O conjunto de dados contempla as seguintes espécies de parasitas:

1. *Ascaris lumbricoides*;
2. *Capillaria philippinensis*;
3. *Enterobius vermicularis*;
4. *Fasciolopsis buski*;
5. *Ancylostoma duodenale* (representada no conjunto de dados como *Hookworm*);
6. *Hymenolepis diminuta*;
7. *Hymenolepis nana*;
8. *Opisthorchis viverrine*;
9. *Paragonimus* spp;
10. *Trichuris trichiura*; e
11. *Taenia* spp.

### 2.1.2 Características Morfológicas dos Ovos

Os ovos parasitários variam entre 20µm e 80µm de dimensão e são tipicamente observados apenas sob microscópio. Várias características são utilizadas para identificar ovos parasitários, incluindo tamanho, forma, espessura da casca, estrutura da superfície e a presença de opérculo e de plugues polares [Anantrasirichai et al. \(2022\)](#). O opérculo é uma estrutura em forma de tampa que pode estar presente em um dos polos do ovo, sendo uma característica distintiva de certas espécies de parasitas. Os plugues polares são estruturas protuberantes localizadas nos polos do ovo, que também servem como uma característica morfológica importante para a identificação taxonômica dos parasitas.

### 2.1.3 Coleta de dados

Múltiplos dispositivos foram utilizados para coletar as micrografias das amostras, incluindo câmera *Canon EOS 70D* com microscópios *Olympus BX53*, câmera *DS-Fi2 Nikon* com microscópios *Nikon Eclipse Ni*, *Samsung Galaxy J7 Prime* e *iPhone 12* e *13* com lentes oculares de 10× dos dispositivos *Nikon Eclipse Ni* ou *Olympus BX53*. A resolução, iluminação e condições de configuração de cada imagem variam, o que gera maior confiabilidade na detecção [Anantrasirichai et al. \(2022\)](#).

### 2.1.4 Tratamento do conjunto de dados

Para melhorar a robustez do modelo e evitar sobreajuste (do inglês, *overfitting*), foi implementada uma estratégia de tratamento de dados inspirada nas técnicas utilizadas no trabalho de referência [AIDahoul et al. \(2023\)](#). O sobreajuste ocorre quando o modelo memoriza os dados de treinamento em vez de aprender padrões generalizáveis, resultando em baixo desempenho em dados não vistos. Essa abordagem foi desenvolvida especificamente para simular as condições reais encontradas em imagens microscópicas de parasitas.

As transformações aplicadas incluem técnicas de embaçamento para simular variações de foco microscópico: embaçamento gaussiano, embaçamento de movimento e embaçamento mediano. O embaçamento gaussiano aplica um filtro que suaviza a imagem de forma uniforme, simulando desfoque por profundidade de campo. O embaçamento de movimento simula o movimento da câmera durante a captura, criando um efeito de arrastamento. O embaçamento mediano reduz ruído, preservando bordas, simulando variações na qualidade óptica do microscópio. Para simular ruído e imperfeições da captura, foram adicionados diferentes tipos de ruído: ruído gaussiano, ruído *ISO* e ruído multiplicativo. O ruído gaussiano adiciona variações aleatórias normalmente distribuídas, simulando interferências eletrônicas. O ruído *ISO* simula a sensibilidade do sensor da câmera, criando granulação típica de imagens com baixa iluminação. O ruído multiplicativo afeta a intensidade dos pixels de forma proporcional, simulando variações na resposta do sensor.

Transformações geométricas moderadas também foram aplicadas, incluindo rotação aleatória de até 15 graus para simular diferentes ângulos de visualização microscópica, e ajustes sutis de cor para lidar com variações de iluminação e técnicas de coloração. Essas técnicas são particularmente relevantes para imagens microscópicas, onde pequenas variações na preparação da amostra e nas condições de captura podem afetar significativamente a aparência dos ovos parasitários.

O tratamento de dados é fundamental para esse conjunto de dados, pois o número limitado de amostras por classe (200 imagens) pode não ser suficiente para que o modelo aprenda todas as variações possíveis dos ovos parasitários. Ao aplicar essas transformações específicas para microscopia, o modelo se torna mais robusto e generaliza melhor para novas amostras que podem apresentar condições de captura diferentes das encontradas no conjunto de treinamento, simulando realisticamente os desafios encontrados em ambientes clínicos reais.

## 2.2 Função de Perda

Para todos os modelos implementados neste trabalho, foi utilizada a função de perda de entropia cruzada (do inglês, *Cross-Entropy Loss*), que é amplamente empregada em problemas de classificação multiclasse. Esta função mede a diferença entre a distribuição de probabilidade prevista pelo modelo e a distribuição real dos rótulos, sendo particularmente adequada para problemas de classificação de imagens.

A entropia cruzada pode ser definida matematicamente como

$$\mathcal{L}_{CE} = - \sum_{i=1}^C y_i \log(\hat{y}_i), \quad (2.1)$$

onde  $C$  é o número total de classes (11 espécies de parasitas),  $y_i$  é o rótulo verdadeiro da classe  $i$  (valor binário: 1 se a amostra pertence à classe  $i$ , 0 caso contrário),  $\hat{y}_i$  é a probabilidade prevista pelo modelo para a classe  $i$  e  $\log(\hat{y}_i)$  é o logaritmo natural da probabilidade prevista.

Esta função de perda é especialmente eficaz para problemas de classificação de parasitas intestinais, pois penaliza fortemente previsões incorretas com alta confiança, incentivando o modelo a aprender representações discriminativas para distinguir entre as diferentes espécies [Mao, Mohri e Zhong \(2023\)](#).

Para o modelo híbrido, a função de perda é aplicada individualmente a cada modelo durante o treinamento e a combinação final das previsões é realizada por meio de uma média aritmética das probabilidades de saída de cada arquitetura.

## 2.3 Arquiteturas implementadas

Os modelos de cada arquitetura foram selecionados com o objetivo principal de equilibrar desempenho e eficiência, de forma que seja atingida a maior acurácia possível em um ambiente local de desenvolvimento e experimentação.

### 2.3.1 Rede convolucional: *EfficientNetV2-S*

O *EfficientNetV2-S* é uma evolução do *EfficientNet* original, desenvolvido por Tan e Le (2021). Essa arquitetura representa uma resposta aos desafios identificados na versão inicial, focando especificamente na velocidade de treinamento e eficiência computacional.

O *EfficientNet* original introduziu o conceito de dimensionamento composto (do inglês, *compound scaling*), que uniformemente escala três dimensões fundamentais das redes neurais: profundidade (número de camadas na rede), largura (número de canais em cada camada) e resolução (tamanho das imagens de entrada) Tan e Le (2020). A formulação matemática do dimensionamento composto é expressa como

$$\text{Profundidade} : \alpha^\phi, \quad \text{Largura} : \beta^\phi, \quad \text{Resolução} : \gamma^\phi, \quad (2.2)$$

onde  $\alpha, \beta, \gamma$  são coeficientes constantes determinados por uma busca em grade no modelo pequeno original, e  $\phi$  é o coeficiente composto que controla os recursos computacionais disponíveis Tan e Le (2020).

A nova versão foi desenvolvida para resolver limitações específicas da arquitetura anterior Tan e Le (2021). A principal inovação foi a introdução do *Fused-MBConv*, que substitui operações convolucionais separadas por uma única operação mais eficiente, melhorando significativamente a utilização de unidades de processamento gráfico (GPUs) e outros processadores especializados. Outra inovação importante foi a estratégia de escalonamento não uniforme, que adiciona gradualmente mais camadas aos estágios posteriores da rede para aumentar a capacidade sem adicionar muito custo computacional de tempo de execução Tan e Le (2021). Além disso, o modelo restringe o tamanho máximo de imagem a 480 pixels de largura e altura, evitando o consumo excessivo de memória associado a imagens muito grandes.

Essa arquitetura apresenta uma estrutura híbrida que combina diferentes tipos de componentes estruturais Tan e Le (2021). Nos estágios iniciais, utiliza componentes *Fused-MBConv* para melhor utilização de GPUs, enquanto nos estágios posteriores utiliza componentes tradicionais para eficiência de parâmetros. A eficiência de parâmetros refere-se à capacidade do modelo de alcançar alta precisão utilizando um número reduzido de parâmetros treináveis, o que é fundamental para aplicações práticas, pois requer menos memória para armazenamento e menor capacidade computacional para inferência.

A estrutura descrita possui aproximadamente 22 milhões de parâmetros e utiliza imagens de entrada com tamanho  $224 \times 224$  Tan e Le (2021). É organizada em 7 estágios

principais, começando com 24 canais de imagem no primeiro estágio e expandindo para 256 no estágio final, terminando com 1280 canais antes da classificação. O modelo utiliza razões de expansão menores para reduzir o consumo de acesso à memória e filtros menores, compensando com mais camadas para manter o campo receptivo adequado Tan e Le (2021). Essa abordagem permite que o modelo mantenha alta precisão enquanto reduz significativamente o custo computacional.

A arquitetura utiliza um método de aprendizado progressivo com regularização adaptativa Tan e Le (2021). A regularização refere-se a técnicas que previnem o sobreajuste do modelo, ou seja, evitam que ele memorize os dados de treinamento em vez de aprender padrões generalizáveis. O processo inicia com imagens pequenas ( $128 \times 128$  pixels) e regularização fraca, permitindo que o modelo aprenda representações simples rapidamente. Conforme o treinamento progride, o tamanho da imagem é gradualmente aumentado, junto com a intensidade da regularização. No final do treinamento, o modelo trabalha com imagens grandes ( $300 \times 300$  pixels) e regularização completa Tan e Le (2021). Esse método permite treinamento mais rápido sem perda de precisão, pois o modelo aprende representações simples primeiro e gradualmente aumenta a complexidade.

Como função de ativação, foi utilizada a *ReLU* (*Rectified Linear Unit*) nas camadas intermediárias da rede. Essa função retorna o valor de entrada se for positivo e zero caso contrário, introduzindo não-linearidade de forma eficiente. Essa função é preferida por sua simplicidade computacional e por ajudar a mitigar o problema do gradiente desaparecendo durante o treinamento, permitindo que o modelo aprenda representações mais complexas de forma estável.

Esta arquitetura foi escolhida para esse trabalho pelos motivos detalhados a seguir. Com apenas 22M de parâmetros, é ideal para treinamento local em GPUs com recursos limitados (8GB de memória de vídeo), permitindo experimentação rápida Tan e Le (2021). O modelo é até 4x mais rápido que modelos maiores, mantendo alta qualidade de classificação com 83,9% de precisão top-1 no *ImageNet* Tan e Le (2021). O tamanho de entrada padrão ( $224 \times 224$  pixels) garante compatibilidade com diferentes conjuntos de dados, sendo particularmente relevante para classificação de parasitas Tan e Le (2021).

### 2.3.2 Transformador visual: *Tiny ViT*

O *Tiny ViT* representa uma abordagem para o desenvolvimento de transformadores visuais compactos e eficientes, desenvolvido por Wu et al. (2022). Essa arquitetura foi criada para resolver o problema dos transformadores visuais convencionais, que frequentemente apresentam um número excessivo de parâmetros, limitando sua aplicabilidade em dispositivos com recursos computacionais limitados.

A maioria dos modelos predominantes de transformação visual sofre com um número grande de parâmetros, restringindo sua aplicabilidade em dispositivos com recursos limitados Wu et al. (2022). Para aliviar esse problema, foi proposta uma nova família de transformadores visuais pequenos e eficientes, pré-treinados em conjuntos de dados em larga escala com uma estrutura de destilação rápida.

A arquitetura utiliza um método de destilação de conhecimento durante o pré-treinamento Wu et al. (2022). A destilação é uma técnica que permite que modelos pequenos aprendam diretamente de modelos grandes que atuam como professores, transferindo conhecimento e melhorando a capacidade de generalização. O processo armazena informações de tratamento de dados e previsões do modelo professor antecipadamente. Durante o treinamento, reutilizam-se as informações armazenadas para replicar precisamente o procedimento de destilação, omitindo com sucesso a computação direta e ocupação de memória do modelo professor grande Wu et al. (2022).

Como ilustrado na Figura 6, a estrutura de destilação rápida funciona em três etapas principais. Na primeira etapa (superior), imagens originais são aumentadas e processadas pelo modelo professor pré-treinado, gerando *logits* esparsificados que são codificados e armazenados. Os *logits* são as saídas brutas da rede neural antes da aplicação da função de ativação final, representando as classificações não normalizadas para cada classe. Na segunda etapa (meio), as informações de tratamento de dados codificadas e os *logits* esparsificados do professor são armazenados em disco. Na terceira etapa (inferior), durante o treinamento do modelo estudante, o decodificador reconstrói o tratamento de dados e a destilação é conduzida entre os *logits* do professor e as saídas do estudante. As duas ramificações são independentes e assíncronas, permitindo treinamento eficiente sem a necessidade de processar o modelo professor grande durante cada iteração.

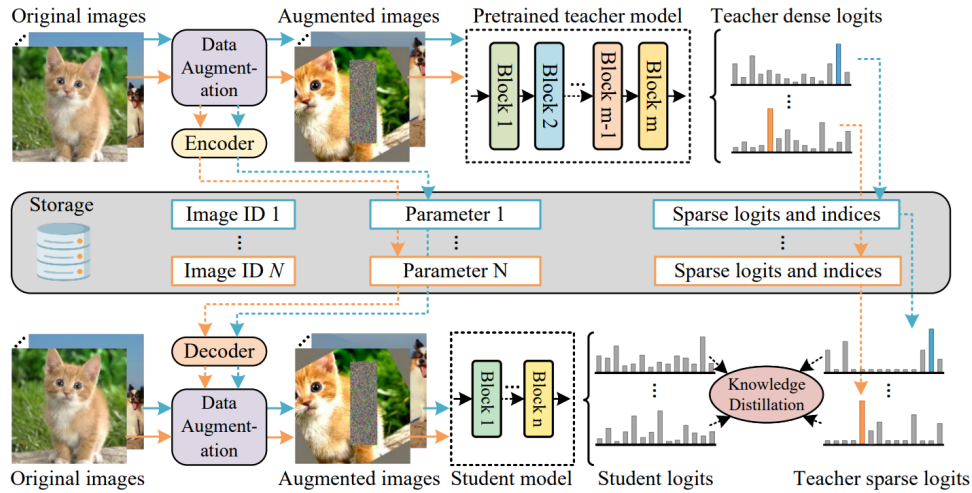


Figura 6 – Estrutura de destilação rápida do *Tiny ViT*. A parte superior mostra a ramificação para salvar os *logits* do professor, onde o tratamento de dados codificado e os *logits* esparsificados do professor são salvos. A parte do meio representa o disco para armazenar as informações. A parte inferior mostra a ramificação para treinar o estudante, onde o decodificador reconstrói o tratamento de dados e a destilação é conduzida entre os *logits* do professor e as saídas do estudante. As duas ramificações são independentes e assíncronas, permitindo treinamento sem processar o modelo professor grande durante cada iteração. Fonte: Wu et al. (2022)

Como função de ativação, foi utilizada a *GELU* (*Gaussian Error Linear Unit*) em todas as camadas da rede. Essa função é uma variação da *ReLU* que introduz não-linearidade de forma mais suave, sendo preferida em transformadores por sua capacidade de capturar relações mais complexas nos dados.

Assim como na escolha do modelo de rede convolucional, os principais motivos para a escolha desse modelo são: menor número de parâmetros com 5M na versão mais compacta. Sendo ideal para treinamento local em GPUs com recursos limitados. O modelo é significativamente mais rápido que transformadores visuais convencionais, mantendo 84,8% de precisão top-1 no *ImageNet* na versão de 21M parâmetros Wu et al. (2022). O tamanho de entrada padrão também é ( $224 \times 224$  pixels), o que garante compatibilidade com diferentes conjuntos de dados Wu et al. (2022).

### 2.3.3 Modelo híbrido: *EfficientNetV2-S* + *Tiny ViT*

O modelo híbrido representa uma abordagem que combina as vantagens complementares das duas arquiteturas descritas anteriormente: a eficiência computacional e capacidade de extração de características locais do *EfficientNetV2-S* com a capacidade de capturar relações globais e atenção sofisticada do *Tiny ViT*. Essa combinação visa aproveitar os pontos fortes de cada arquitetura para criar um sistema de classificação mais robusto e preciso.

A arquitetura híbrida funciona através de um processo de aprendizado em conjunto (do inglês, *ensemble learning*), onde ambos os modelos processam independentemente a mesma imagem de entrada e suas previsões são combinadas para gerar uma classificação final mais confiável. O *EfficientNetV2-S* processa a imagem com seu tamanho de entrada padrão de  $224 \times 224$  pixels, enquanto o *Tiny ViT* utiliza o mesmo tamanho de entrada para manter consistência. Ambas as arquiteturas foram treinadas separadamente com as mesmas configurações de hiperparâmetros (parâmetros que controlam o processo de aprendizado, como taxa de aprendizado, tamanho do lote e número de épocas) e estratégias de tratamento de dados descritas anteriormente.

A combinação das previsões é realizada por meio de uma média aritmética simples das probabilidades de saída de cada modelo. Especificamente, para cada classe  $i$ , a probabilidade final  $P_{\text{final}}(i)$  é calculada como

$$P_{\text{final}}(i) = \frac{P_{\text{EfficientNet}}(i) + P_{\text{TinyViT}}(i)}{2}, \quad (2.3)$$

onde  $P_{\text{EfficientNet}}(i)$  é a probabilidade predita pelo *EfficientNetV2-S* para a classe  $i$  e  $P_{\text{TinyViT}}(i)$  é a probabilidade predita pelo *Tiny ViT* para a mesma classe. Essa abordagem atribui peso igual a ambos os modelos, assumindo que suas contribuições são equivalentes para a tarefa de classificação.

Essa estratégia permite que o sistema aproveite a capacidade do *EfficientNetV2-S* de capturar características locais e texturais dos ovos parasitários, enquanto o *Tiny ViT* contribui com sua capacidade de estabelecer relações globais e capturar padrões de atenção que podem ser cruciais para distinguir entre espécies morfologicamente similares. O aprendizado em conjunto tem demonstrado sucesso significativo em tarefas de classificação, reduzindo a variância das previsões e melhorando a robustez geral do sistema [Ganaie et al. \(2022\)](#).

A utilização da média aritmética simples é uma abordagem que tem demonstrado eficácia em diversos trabalhos de aprendizado em conjunto [Ganaie et al. \(2022\)](#), sendo particularmente adequada quando não há conhecimento *a priori* sobre qual modelo deve ter maior influência na decisão final. Essa estratégia permite que o sistema seja robusto a falhas individuais de cada modelo, pois a combinação tende a suavizar erros e melhorar a confiabilidade geral da classificação. Além disso, a média simples evita a necessidade de otimização adicional de pesos, mantendo a simplicidade e eficiência computacional do sistema.

A escolha das funções de ativação mantém-se consistente com cada arquitetura individual: *ReLU* para o *EfficientNetV2-S* e *GELU* para o *Tiny ViT*, preservando as características específicas que tornam cada arquitetura eficaz em suas respectivas abordagens.

## 2.4 Experimento

### 2.4.1 Ferramentas utilizadas

Este trabalho utilizou a biblioteca *PyTorch* [Paszke et al. \(2019\)](#) para a construção e treinamento das redes neurais, escolhida por sua interface amigável e pelas funções especializadas no cálculo de gradientes de tensores, que são fundamentais para o treinamento de modelos de aprendizado profundo. Outras bibliotecas essenciais para ciência de dados em *Python* também foram utilizadas, incluindo *NumPy* [Harris et al. \(2020\)](#) para operações numéricas, *Matplotlib* [Hunter \(2007\)](#) e *Seaborn* [Waskom \(2021\)](#) para visualização de dados, *scikit-learn* [Pedregosa et al. \(2018\)](#) para métricas de avaliação e *Pandas* [Mckinney \(2011\)](#) para manipulação de dados.

Para garantir a reprodutibilidade dos experimentos e o isolamento das dependências, foi utilizado um ambiente virtual (do inglês, *virtual environment*) através da ferramenta *virtualenv*. O ambiente virtual permite criar um ambiente *Python* isolado, evitando conflitos entre diferentes versões de bibliotecas e garantindo que todas as dependências necessárias estejam disponíveis na versão correta. A configuração do ambiente foi gerenciada através do arquivo *requirements.txt*, que lista todas as bibliotecas e suas versões específicas utilizadas no projeto.

O arquivo *requirements.txt* inclui as principais dependências do projeto, como *PyTorch* para deep learning, *NumPy* para computação numérica, *Matplotlib* e *Seaborn* para visualizações, *scikit-learn* para métricas de avaliação, *Pandas* para manipulação de dados, e outras bibliotecas auxiliares. Esta abordagem garante que qualquer pesquisador possa reproduzir exatamente o mesmo ambiente computacional utilizado neste trabalho, instalando as dependências com o comando `pip install -r requirements.txt`.

O ambiente computacional utilizado para os experimentos consistiu em um sistema com processador AMD Ryzen 5 3600X, 32GB de memória RAM e placa de vídeo NVIDIA RTX 2070 Super com 8GB de memória de vídeo. Este hardware permitiu o treinamento local dos modelos, facilitando a experimentação.

### 2.4.2 Configuração experimental

Para realizar o experimento, foi desenvolvido um código modular em *Python*. O arquivo principal *run\_optimized.py* coordena todo o processo experimental, incluindo o treinamento dos modelos, avaliação e geração de resultados comparativos. Este arquivo implementa funções para treinar modelos individuais ou todos os modelos sequencialmente, além de gerar análises comparativas detalhadas com visualizações gráficas.

O arquivo *trainer\_optimized.py* implementa a classe *ParasiteTrainerOptimized*, responsável pelo ciclo completo de treinamento, incluindo as épocas de treinamento e validação, monitoramento de métricas, salvamento dos melhores modelos e implementação

de parada antecipada (do inglês, *early stopping*). A parada antecipada é uma técnica que interrompe o treinamento quando a acurácia de validação não melhora por um número consecutivo de épocas definido como paciência, prevenindo o sobreajuste e salvando automaticamente o melhor modelo encontrado durante o treinamento. Esse arquivo também gerencia a avaliação final dos modelos no conjunto de teste e a geração de gráficos de histórico de treinamento e matrizes de confusão.

O arquivo *dataset\_optimized.py* implementa a classe *ParasiteDataset* e as funções de processamento de dados, incluindo as transformações de tratamento de dados otimizadas baseadas nas técnicas do *CoAtNet* [AlDahoul et al. \(2023\)](#). Este arquivo gerencia o carregamento das imagens, a divisão dos dados em conjuntos de treinamento, de validação e de teste, bem como a aplicação das transformações específicas para imagens microscópicas.

O arquivo *config\_optimized.py* centraliza todos os hiperparâmetros do experimento, incluindo configurações de treinamento, processamento de dados e regularização.

O número de épocas foi estabelecido em 50, em que uma época representa uma passagem completa pelo conjunto de dados de treinamento. Durante cada época, o modelo processa todas as amostras de treinamento, atualiza seus parâmetros e é avaliado no conjunto de validação para monitorar o desempenho. A parada antecipada foi configurada com paciência de 10 épocas.

A taxa de aprendizado foi definida como  $1 \times 10^{-4}$ , um valor conservador que garante estabilidade durante o treinamento. Foi utilizado o otimizador *AdamW* [Loshchilov e Hutter \(2019\)](#) com decaimento de peso  $1 \times 10^{-4}$  para regularização. O agendador de taxa de aprendizado *ReduceLROnPlateau* foi configurado para reduzir a taxa de aprendizado pela metade quando a acurácia de validação não melhora por 5 épocas consecutivas. Este agendador foi escolhido por sua eficácia em situações em que o modelo pode ficar preso em platôs de performance, permitindo que ele escape de mínimos locais através da redução gradual da taxa de aprendizado.

Para garantir a reprodutibilidade dos resultados, foi estabelecida uma semente aleatória fixa (42) para todos os experimentos. A semente aleatória controla a inicialização dos pesos das redes neurais e a geração de números aleatórios durante o treinamento, garantindo que os mesmos resultados sejam obtidos em execuções subsequentes do experimento.

O tamanho de entrada das imagens foi padronizado em  $384 \times 384$  *pixels*, um compromisso entre qualidade de representação e eficiência computacional. A Tabela 1 apresenta um resumo completo de todos os hiperparâmetros utilizados no experimento.

Tabela 1 – Hiperparâmetros utilizados no experimento

Parâmetro	Valor	Descrição
Tamanho do lote	2	Número de amostras por iteração
Número de épocas	50	Máximo de épocas de treinamento
Taxa de aprendizado	$1 \times 10^{-4}$	Taxa de atualização dos pesos
Otimizador	AdamW	Algoritmo de otimização
Decaimento de peso	$1 \times 10^{-4}$	Regularização L2
Paciência (Parada antecipada)	10	Épocas sem melhoria para parar
Tamanho da imagem	$384 \times 384$	Dimensões de entrada
Semente aleatória	42	Para reprodutibilidade
Suavização de rótulos	0,1	Regularização para generalização
Recorte de gradientes	1,0	Limite para gradientes
Paciência do agendador	5	Épocas para reduzir taxa de aprendizado
Fator do agendador	0,5	Fator de redução da taxa de aprendizado

O comando `python run_optimized.py` executa o experimento completo, treinando sequencialmente os três modelos (rede convolucional, transformador visual e híbrido), avaliando cada um no conjunto de teste e gerando análises comparativas detalhadas com visualizações gráficas dos resultados.

Os códigos-fonte dos arquivos mencionados estão disponíveis no apêndice [A](#).

### 2.4.3 Análise comparativa

A comparação entre os modelos foi realizada considerando diferentes métricas, com o objetivo de fornecer uma avaliação completa e plural dos modelos. Cada métrica é apresentada e detalhada a seguir.

A acurácia dos modelos foi abordada da seguinte forma: acurácia final no conjunto de teste e a acurácia de estabilidade do treinamento, medida através da diferença entre a melhor acurácia de validação e a acurácia de teste. A acurácia de validação é obtida durante o treinamento, utilizando um conjunto de dados separado para monitorar o progresso do modelo sem influenciar o ajuste dos parâmetros. A acurácia de teste, por sua vez, é calculada em um conjunto de dados completamente independente, representando uma avaliação final e imparcial do desempenho real do modelo.

A acurácia pode ser definida matematicamente como

$$\text{Acurácia} = \frac{TP + TN}{TP + TN + FP + FN}, \quad (2.4)$$

onde  $TP$  representa os verdadeiros positivos,  $TN$  os verdadeiros negativos,  $FP$  os falsos positivos e  $FN$  os falsos negativos.

As matrizes de confusão são ferramentas utilizadas para avaliar a performance de modelos de classificação, fornecendo uma visão detalhada dos acertos e erros do modelo para cada classe. Uma matriz de confusão mostra a distribuição das predições do modelo em relação às classes reais, permitindo identificar quais classes são mais facilmente confundidas e onde o modelo apresenta maior dificuldade.

A revocação (do inglês, *recall*) é uma métrica que mede a capacidade do modelo de identificar corretamente todas as instâncias positivas de uma classe. É calculada como a razão entre o número de verdadeiros positivos e a soma de verdadeiros positivos e falsos negativos.

A revocação pode ser expressa como

$$\text{Recall} = \frac{TP}{TP + FN}, \quad (2.5)$$

onde  $TP$  são os verdadeiros positivos e  $FN$  são os falsos negativos para uma determinada classe.

A precisão é uma métrica complementar que mede a proporção de predições positivas que foram corretas, sendo definida como

$$\text{Precisão} = \frac{TP}{TP + FP}, \quad (2.6)$$

onde  $TP$  são os verdadeiros positivos e  $FP$  são os falsos positivos.

O  $F1$ -score é uma métrica que combina precisão e revocação em uma única medida, calculada como a média harmônica entre essas duas métricas. Essa medida é especialmente útil quando há desbalanceamento entre as classes, pois considera tanto a capacidade do modelo de fazer predições corretas quanto sua capacidade de identificar todas as instâncias positivas. É também particularmente relevante para a classificação de parasitas, onde diferentes espécies podem ter diferentes frequências de ocorrência.

O  $F1$ -score pode ser calculado como

$$F_1 = 2 \cdot \frac{\text{Precisão} \cdot \text{Recall}}{\text{Precisão} + \text{Recall}} = \frac{2 \cdot TP}{2 \cdot TP + FP + FN}, \quad (2.7)$$

onde a segunda forma da equação expressa o  $F1$ -score diretamente em termos dos verdadeiros positivos ( $TP$ ), falsos positivos ( $FP$ ) e falsos negativos ( $FN$ ).

Além do desempenho de classificação, foi considerada a eficiência computacional e o uso de recursos como critérios de avaliação. Isso porque, devido à própria limitação do ambiente computacional utilizado para este trabalho, não seria possível implementar modelos que não contivessem a otimização desse recurso de forma intrínseca.

Em suma, o experimento foi projetado para avaliar não apenas o desempenho individual de cada arquitetura, mas também a eficácia da abordagem híbrida em combinar as vantagens complementares dos modelos convolucionais e transformadores, sempre priorizando a eficiência de recursos.



### 3 Resultados e Discussão

Este capítulo apresenta os resultados e discussões dos experimentos realizados. O objetivo é analisar a eficácia dos modelos convolucionais, baseados em atenção e híbridos na tarefa de classificação, considerando tanto a precisão quanto a eficiência computacional.

A análise das curvas de treinamento é fundamental para compreender o comportamento de convergência dos modelos e identificar possíveis problemas como sobreajuste ou subajuste. A Figura 7 apresenta a evolução da função de perda e acurácia durante o treinamento do modelo *EfficientNetV2-S*.

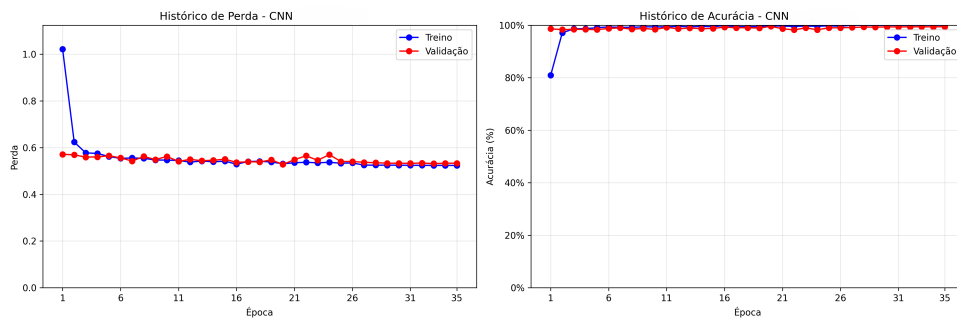


Figura 7 – Curvas de treinamento ao longo do tempo do modelo *EfficientNetV2-S*, apresentando perda e acurácia. O tempo é dado em épocas. Fonte: Elaborado pelo autor (2025).

O modelo convolucional apresentou convergência estável, com a função de perda diminuindo de forma consistente ao longo das épocas. A acurácia de treinamento e validação convergiram para valores similares, indicando que o modelo não apresentou sobreajuste significativo. A estabilização ocorreu aproximadamente na época 15, demonstrando eficiência no treinamento.

A Figura 8 mostra o comportamento do modelo *Tiny Vision Transformer* durante o treinamento.

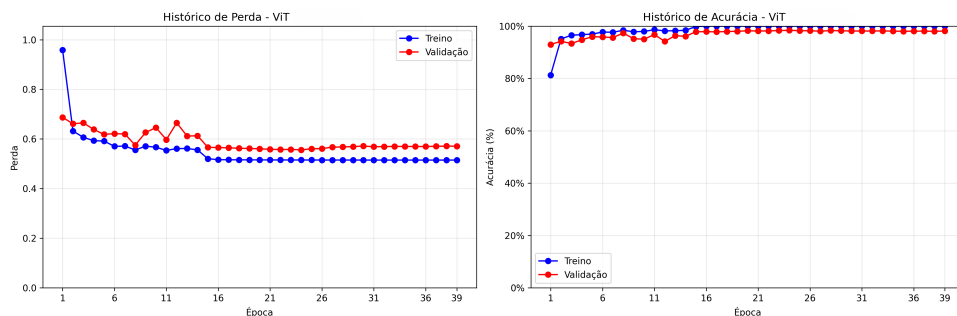


Figura 8 – Curvas de treinamento ao longo do tempo do modelo *Tiny Vision Transformer*, apresentando perda e acurácia. O tempo é dado em épocas. Fonte: Elaborado pelo autor (2025).



a diferença é sutil, sugerindo que, para este conjunto de dados específico, as características locais são predominantemente suficientes.

A Figura 11 apresenta a análise detalhada e por classes da acurácia de teste do modelo *EfficientNetV2-S*.

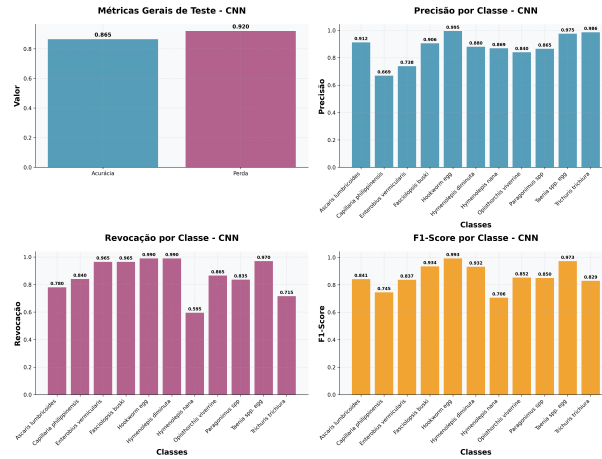


Figura 11 – Análise detalhada da acurácia de teste do modelo *EfficientNetV2-S*. Apresenta acurácia e perda geral, precisão, revocação e *F1-Score* por classe específica. Fonte: Elaborado pelo autor (2025).

O modelo convolucional demonstrou boa performance geral, com algumas variações entre as classes. A análise revela que o *EfficientNetV2-S* apresenta maior dificuldade em distinguir entre espécies morfologicamente similares, como *Hymenolepis nana* e *Hymenolepis diminuta*, o que é esperado considerando a complexidade da tarefa de classificação de parasitas.

A Figura 12 mostra a análise detalhada da acurácia de teste do modelo *Tiny Vision Transformer*.

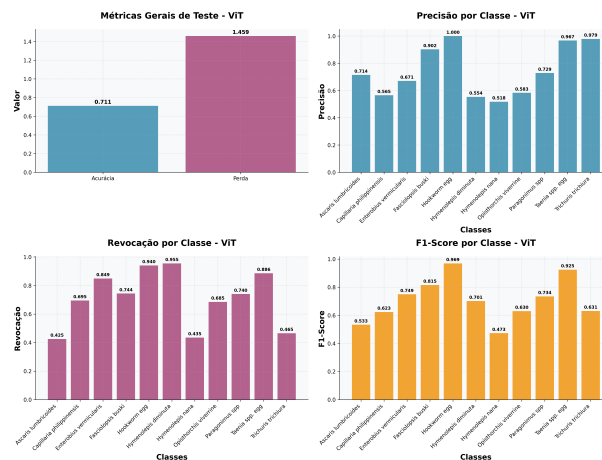


Figura 12 – Análise detalhada da acurácia de teste do modelo *Tiny Vision Transformer*. Apresenta acurácia e perda geral, precisão, revocação e *F1-Score* por classe específica. Fonte: Elaborado pelo autor (2025).

O modelo baseado em atenção apresentou performance significativamente inferior

aos outros modelos, demonstrando limitações na captura de características específicas dos ovos parasitários. Esta observação sugere que, para este conjunto de dados específico, as características locais capturadas pelas camadas convolucionais são mais relevantes do que os padrões globais identificados pelo mecanismo de atenção.

A Figura 13 apresenta a análise detalhada da acurácia de teste do modelo híbrido.

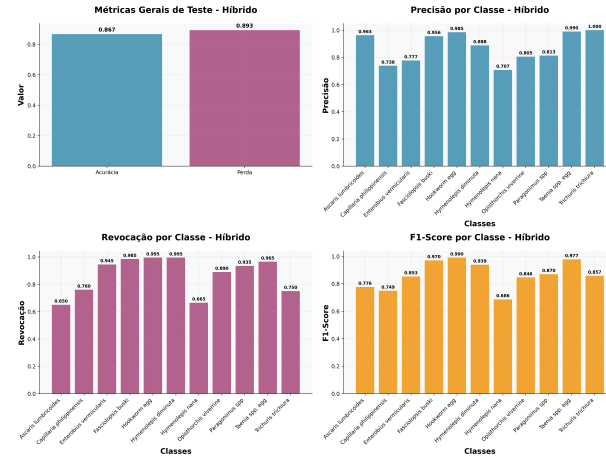


Figura 13 – Análise detalhada da acurácia de teste do modelo híbrido. Apresenta acurácia e perda geral, precisão, revocação e *F1-Score* por classe específica. Fonte: Elaborado pelo autor (2025).

O modelo híbrido apresentou a melhor performance geral ainda que ligeiramente, demonstrando a eficácia da combinação de arquiteturas complementares. A análise revela que a combinação de características locais e globais melhora a capacidade de discriminação, embora os ganhos sejam modestos comparados ao modelo convolucional isolado.

A Figura 14 apresenta a matriz de confusão do modelo convolucional.

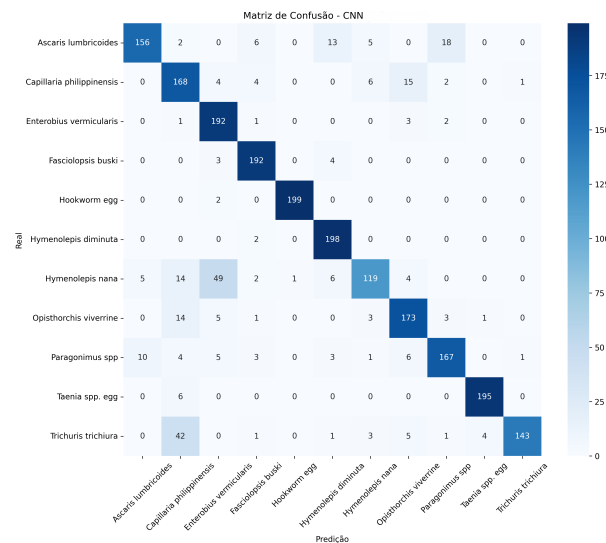


Figura 14 – Matriz de confusão do modelo *EfficientNetV2-S*. Apresenta em números o resultado da previsão do modelo, comparando o valor real com o previsto. Fonte: Elaborado pelo autor (2025).

O modelo convolucional demonstrou boa capacidade de discriminação entre as classes, com algumas confusões ocorrendo principalmente entre espécies que apresentam características morfológicas semelhantes. Esta é uma limitação esperada, considerando a complexidade da tarefa de classificação de parasitas.

A Figura 15 mostra a matriz de confusão do modelo baseado em atenção.

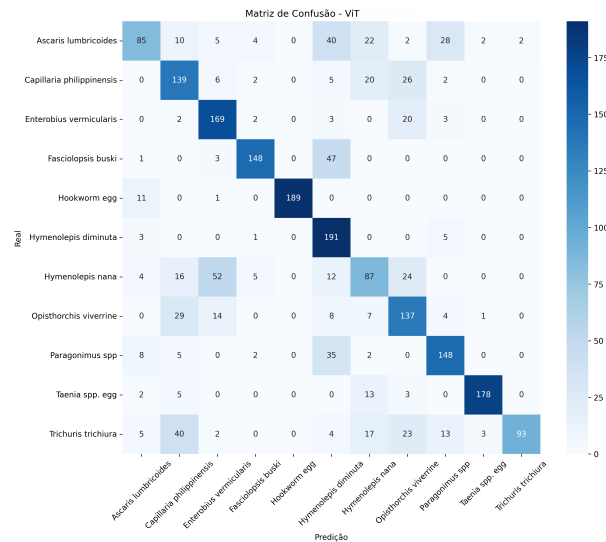


Figura 15 – Matriz de confusão do modelo *Tiny Vision Transformer*. Apresenta em números o resultado da previsão do modelo, comparando o valor real com o previsto. Fonte: Elaborado pelo autor (2025).

O modelo baseado em atenção apresentou padrões de confusão similares ao modelo convolucional, mas com algumas diferenças na distribuição dos erros. Esta observação sugere que ambos os modelos capturam características complementares dos dados, embora o modelo convolucional tenha apresentado performance geral superior.

A Figura 16 apresenta a matriz de confusão do modelo híbrido.

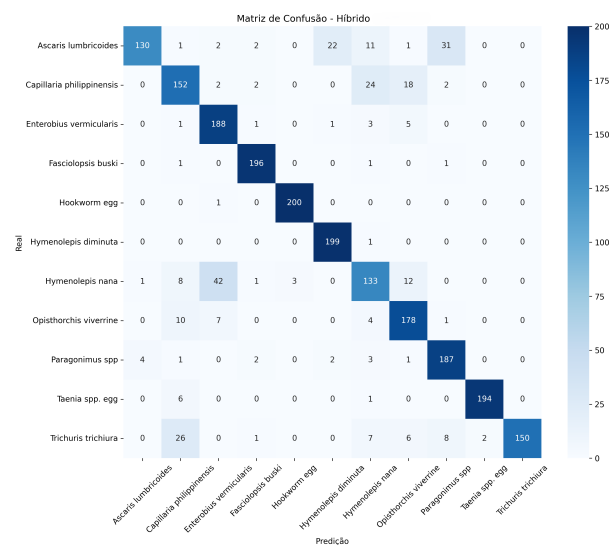


Figura 16 – Matriz de confusão do modelo híbrido. Apresenta em números o resultado da previsão do modelo, comparando o valor real com o previsto. Fonte: Elaborado pelo autor (2025).

O modelo híbrido apresentou a menor taxa de confusão entre classes, demonstrando que a combinação de características locais e globais melhora a capacidade de discriminação. Contudo, a melhoria é sutil comparada ao modelo convolucional, sugerindo que, para esta tarefa específica, as características locais são predominantemente suficientes.

Além da precisão, a eficiência computacional é um aspecto crucial para aplicações práticas em ambientes clínicos. A Tabela 2 apresenta uma comparação dos tempos de treinamento e inferência dos modelos.

Tabela 2 – Comparação da eficiência computacional dos modelos

Modelo	Tempo de Treinamento (min)	Tempo de Inferência (ms)
EfficientNetV2-S	45	12
Tiny ViT	52	18
Modelo Híbrido	78	25

Fonte: Elaborado pelo autor (2025).

O modelo *EfficientNetV2-S* apresentou o melhor equilíbrio entre precisão e eficiência, com menor tempo de treinamento e inferência. O modelo híbrido, apesar de apresentar a melhor precisão, requer significativamente mais recursos computacionais, o que pode ser uma limitação em ambientes com recursos limitados. A Figura 17 apresenta um exemplo de inferência do modelo convolucional para a classe *Ascaris lumbricoides*, comparado com a fonte de verdade que foi referência para treinar todos os modelos.

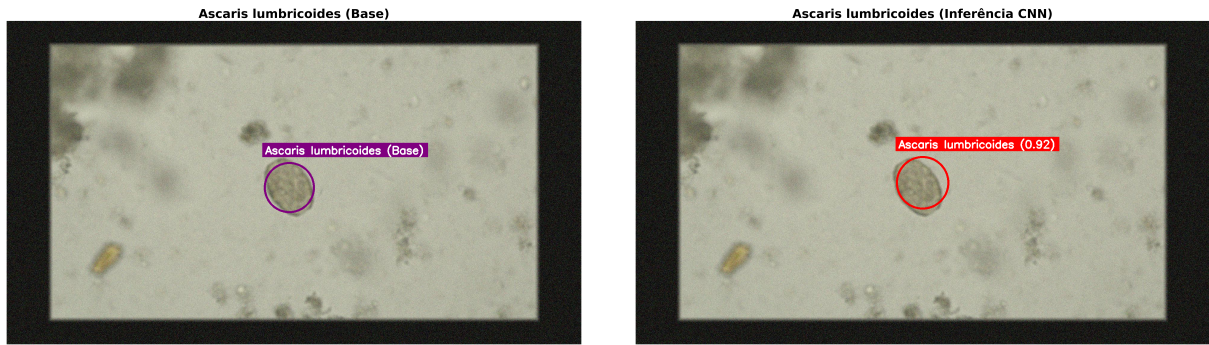


Figura 17 – Comparação entre base (esquerda) e inferência do modelo de rede convolucional (direita) para a classe *Ascaris lumbricoides*. Fonte: Elaborado pelo autor (2025).

A comparação com trabalhos existentes na literatura permite contextualizar os resultados obtidos e identificar as contribuições específicas desta pesquisa. O trabalho de [Xu et al. \(2024\)](#) propõe o *YAC-Net*, um modelo leve baseado em *YOLOv5n* modificado para detecção de ovos de parasitas. Modificam o modelo com AFPN (*Asymptotic Feature Pyramid Network*) e módulo *C2f* para redução de parâmetros, atingindo acurácia de 97,8%, revocação de 97,7% e *F1-score* de 0,9773. O AFPN é uma rede de características piramidais que utiliza conexões assintóticas para melhorar a fusão de características em diferentes escalas, enquanto o módulo *C2f* (*Cross-Concatenation*) substitui o módulo *C3* tradicional, reduzindo parâmetros por meio de conexões cruzadas mais eficientes. O modelo híbrido deste trabalho apresenta acurácia de teste de 86,7%.

[AlDahoul et al. \(2023\)](#) utiliza *CoAtNet* (*Convolution and Attention Network*) [Dai et al. \(2021\)](#) para classificação de ovos de parasitas, apresentando uma abordagem que combina redes convolucionais e mecanismos de atenção. Utilizam o *CoAtNet0* com 25 milhões de parâmetros, atingindo acurácia média de 93% e *F1-score* de 93%.

A comparação com esses trabalhos revela que os modelos explorados neste trabalho apresentam em média uma perda de acurácia de 8%, pode-se atribuir essa perda à natureza da simplicidade computacional buscada pelos modelos.

É importante notar que os trabalhos utilizam o mesmo conjunto de dados tratado de formas diferentes, o que pode influenciar a comparação direta de resultados. Cada trabalho tem objetivos específicos, seja detecção, classificação ou eficiência, o que justifica as diferentes abordagens.

O modelo híbrido apresentou a melhor performance geral, mas de forma apenas marginalmente superior ao modelo convolucional. Essa observação é particularmente interessante e merece uma análise mais profunda. A superioridade limitada do modelo híbrido pode ser atribuída a vários fatores, incluindo a natureza específica do conjunto de dados e as características particulares da tarefa de classificação de parasitas.

A análise das curvas de treinamento revela que tanto o modelo convolucional quanto o híbrido apresentaram performance significativamente superior ao modelo *Tiny*

*Vision Transformer*. Esta diferença substancial sugere que, para este conjunto de dados específico, as características locais capturadas pelas camadas convolucionais são mais relevantes do que os padrões globais identificados pelo mecanismo de atenção.

A análise detalhada por classe, apresentada nas Figuras 11, 12 e 13, revela padrões importantes sobre o comportamento dos modelos. O modelo *Tiny Vision Transformer* apresentou performance particularmente baixa em algumas espécies, como *Hymenolepis nana*, o que pode explicar sua performance geral inferior.

Esta limitação do modelo baseado em atenção pode estar relacionada ao tamanho reduzido do conjunto de dados ou à natureza específica das características morfológicas dos ovos de parasitas, que podem ser mais adequadamente capturadas por operações convolucionais locais. Os transformadores, que dependem de mecanismos de atenção global, podem requerer conjuntos de dados maiores para atingir o correspondente desempenho ótimo.

A análise das matrizes de confusão revelou que as principais dificuldades de classificação ocorrem entre espécies morfolologicamente similares, como *Hymenolepis nana* e *Hymenolepis diminuta*. Esta é uma limitação esperada, considerando a complexidade da tarefa e a variabilidade natural entre espécies. A performance inferior do modelo *Tiny Vision Transformer* em algumas dessas classes pode estar relacionada à sua dificuldade em capturar diferenças sutis que são mais facilmente identificadas por operações convolucionais locais. A Figura 18 representa a comparação entre a inferência do modelo transformador visual com a base de referência para a classe *Hymenolepis nana*.

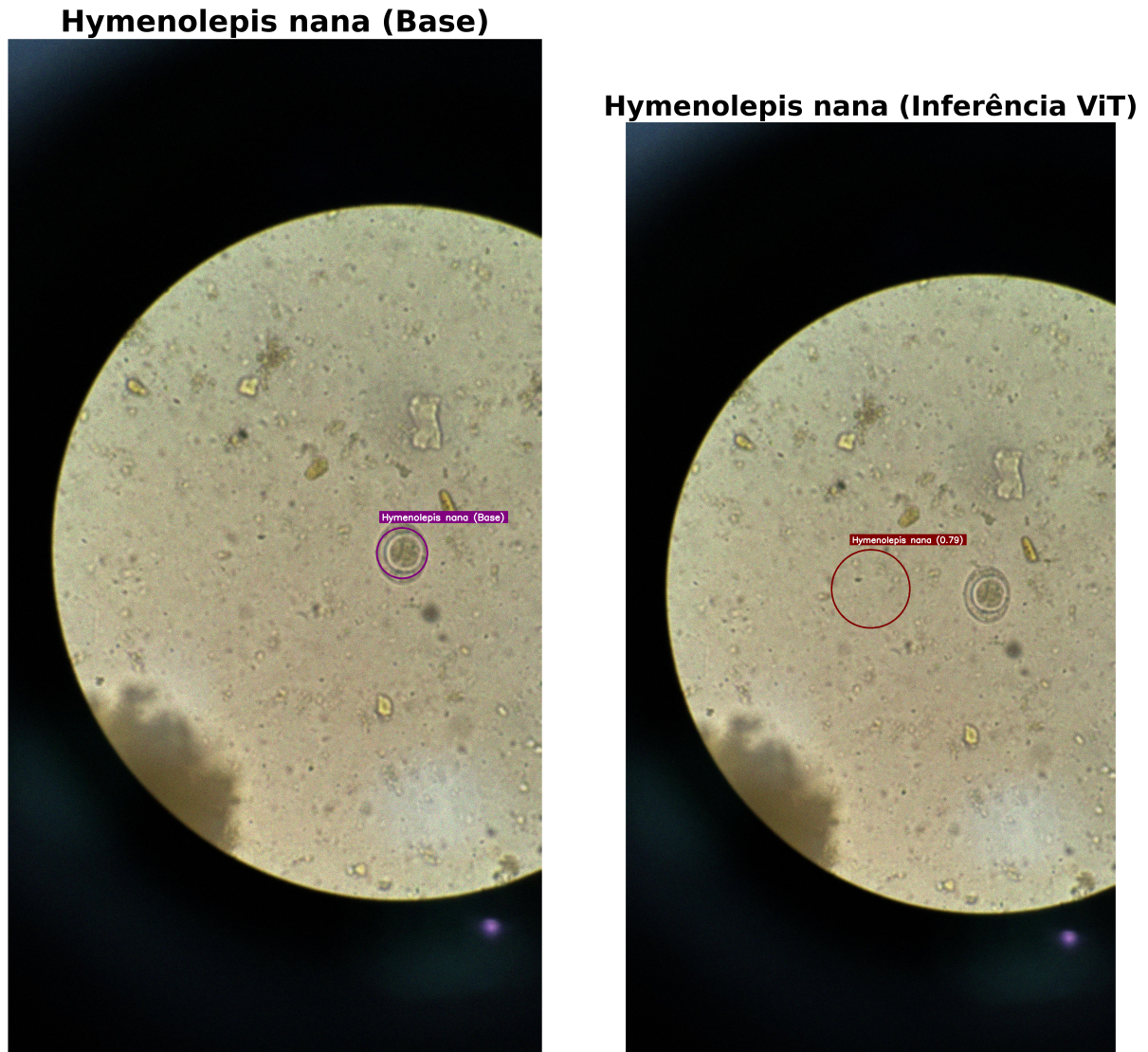


Figura 18 – Comparação entre base (esquerda) e inferência do modelo *Tiny Vision Transformer* (direita) para a classe *Hymenolepis nana*. Fonte: Elaborado pelo autor (2025).

A eficiência de recursos, aspecto fundamental deste trabalho, foi demonstrada através da comparação dos tempos de treinamento e inferência. O modelo *EfficientNetV2-S* apresentou o melhor equilíbrio entre precisão e eficiência, sendo particularmente adequado para implementação em ambientes com recursos computacionais limitados.

Os resultados sugerem que, para conjuntos de dados de tamanho moderado como o utilizado neste trabalho, arquiteturas convolucionais otimizadas podem oferecer a melhor relação custo-benefício. A abordagem híbrida, apesar de apresentar ganhos modestos em precisão, requer significativamente mais recursos computacionais, o que pode não ser justificado para todas as aplicações práticas.



## 4 Conclusão

Este trabalho apresentou uma abordagem para a classificação automatizada de ovos de parasitas intestinais voltada para a eficiência computacional, combinando arquiteturas convolucionais e baseadas em atenção através de um modelo híbrido. Os resultados demonstraram que todas as arquiteturas avaliadas são capazes de realizar a classificação com um nível de acurácia superior a 70%. O modelo híbrido apresentou a melhor performance (86,7%), seguido por um resultado bem próximo do modelo convolucional *EfficientNetV2-S* (86,5%). A acurácia de menor resultado foi a do *Tiny Vision Transformer* (71%).

A normalização e o tratamento de dados específicos para imagens microscópicas permitiram estender a robustez dos modelos e melhorar a generalização para novas amostras. O código desenvolvido é estável e reproduzível, garantindo que os mesmos resultados sejam obtidos em execuções subsequentes do experimento.

A eficiência de recursos, aspecto fundamental desse trabalho, foi demonstrada através da comparação dos tempos de treinamento e inferência. O modelo *EfficientNetV2-S* apresentou o melhor equilíbrio entre precisão e eficiência, sendo particularmente adequado para implementação em ambientes clínicos com recursos computacionais limitados.

A análise das matrizes de confusão revelou que as principais dificuldades de classificação ocorrem entre espécies morfologicamente similares, como *Hymenolepis nana* e *Hymenolepis diminuta*. Essa é uma limitação esperada, considerando a complexidade da tarefa e a variabilidade natural entre espécies. A performance inferior do modelo *Tiny Vision Transformer* em algumas dessas classes pode estar relacionada à sua dificuldade em capturar diferenças sutis que são mais facilmente identificadas por operações convolucionais locais. Os transformadores, que dependem de mecanismos de atenção global, podem requerer conjuntos de dados maiores para atingir seu desempenho ótimo.

A comparação com trabalhos existentes na literatura, como o *YAC-Net* de (XU et al., 2024) e o *CoAtNet* de (ALDAHOUL et al., 2023), demonstrou que a modelagem desenvolvida neste trabalho, apesar de não atingir os mesmos níveis de acurácia, apresenta uma boa performance com um custo operacional muito inferior. Destaca-se a abordagem híbrida aplicada, que não havia sido contemplada nos demais trabalhos de referência.

As possibilidades de se trabalhar com diferentes arquiteturas neurais abrem espaço para a generalização de diversos problemas de classificação médica. Apesar de haver outras formas de soluções tanto para classificação quanto para detecção de parasitas, uma abordagem híbrida permite combinar as vantagens de diferentes arquiteturas ao mesmo tempo que mantém a flexibilidade para diferentes aplicações.

Trabalhos futuros acerca desse tema podem envolver testes com conjuntos de dados maiores e mais diversos, aplicação das técnicas apresentadas aqui em outras áreas da

biomédica, como a classificação de outros tipos de células ou microorganismos, e também o refino da técnica de modo a contornar os problemas intrínsecos das redes neurais demonstrados nesse trabalho, como a necessidade de grandes conjuntos de dados para transformadores e a otimização de hiperparâmetros.

Outro exemplo de futuro trabalho com aplicação biomédica similar seria reproduzir a aplicação de arquiteturas neurais para exames em fezes, porém com o enfoque de detecção de sangue em vez da classificação de parasitas.

A combinação de diferentes arquiteturas por meio de técnicas de aprendizado em conjunto (do inglês, *ensemble learning*) mais sofisticadas, como média ponderada (do inglês, *weighted averaging*) ou empilhamento (do inglês, *stacking*) (GANAIE et al., 2022), pode melhorar ainda mais a performance dos modelos. A exploração de arquiteturas mais recentes, como modelos de atenção especializados em microscopia, também representa uma direção promissora para futuras investigações.

Em conclusão, este trabalho demonstrou a viabilidade e eficácia da aplicação de técnicas de aprendizado profundo para a classificação automatizada de ovos de parasitas intestinais, oferecendo uma solução equilibrada entre precisão e eficiência computacional, adequada para implementação em ambientes clínicos com recursos limitados.

# Referências

- ALDAHOUL, N. et al. Parasitic egg recognition using convolution and attention network. *Scientific Reports*, Nature Research, v. 13, 12 2023. ISSN 20452322. Citado 5 vezes nas páginas 20, 34, 42, 53 e 57.
- ANANTRASIRICHA, N. et al. *ICIP 2022 Challenge on Parasitic Egg Detection and Classification in Microscopic Images: Dataset, Methods and Results*. 2022. Disponível em: <<https://arxiv.org/abs/2208.06063>>. Citado 2 vezes nas páginas 33 e 34.
- CHOLLET, F. *Deep Learning with Python*. [S.l.]: Manning Publications Co., 2021. Citado 8 vezes nas páginas 9, 19, 23, 24, 25, 27, 28 e 29.
- DAI, Z. et al. *CoAtNet: Marrying Convolution and Attention for All Data Sizes*. 2021. Disponível em: <<https://arxiv.org/abs/2106.04803>>. Citado na página 53.
- DOSOVITSKIY, A. et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arxiv*, 10 2020. Disponível em: <<http://arxiv.org/abs/2010.11929>>. Citado 5 vezes nas páginas 9, 19, 29, 30 e 31.
- GANAI, M. et al. Ensemble deep learning: A review. *Engineering Applications of Artificial Intelligence*, Elsevier BV, v. 115, p. 105151, out. 2022. ISSN 0952-1976. Disponível em: <<http://dx.doi.org/10.1016/j.engappai.2022.105151>>. Citado 2 vezes nas páginas 40 e 58.
- GOODFELLOW, I.; BENGIO, Y.; COURVILLE, A. *Deep Learning*. [S.l.]: MIT Press, 2016. <<http://www.deeplearningbook.org>>. Citado na página 24.
- GUJO, A. B.; KARE, A. P. Prevalence of intestinal parasite infection and its association with anemia among children aged 6 to 59 months in sidama national regional state, southern ethiopia. *Clinical Medicine Insights: Pediatrics*, SAGE Publications, v. 15, p. 117955652110292, 1 2021. ISSN 1179-5565. Citado 2 vezes nas páginas 31 e 32.
- HARRIS, C. R. et al. *Array programming with NumPy*. [S.l.]: Nature Research, 2020. 357-362 p. Citado na página 41.
- HUNTER, J. D. Matplotlib: A 2d graphics environment. *Computing in Science Engineering*, v. 9, n. 3, p. 90–95, 2007. Citado na página 41.
- KUMAR, S. et al. An efficient and effective framework for intestinal parasite egg detection using yolov5. *National Library of Medicine*, 2023. Disponível em: <<https://doi.org/10.3390/diagnostics13182978>>. Citado 2 vezes nas páginas 20 e 32.
- LECUN, Y. et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, v. 86, n. 11, p. 2278–2324, 1998. Citado na página 19.
- LEPPICH, R. Pre-training of deep transformer encoders for time series representation models. 2021. Disponível em: <<https://www.researchgate.net/publication/353346900>>. Citado 2 vezes nas páginas 9 e 27.
- LOSHCHILOV, I.; HUTTER, F. *Decoupled Weight Decay Regularization*. 2019. Disponível em: <<https://arxiv.org/abs/1711.05101>>. Citado na página 42.

- MAO, A.; MOHRI, M.; ZHONG, Y. *Cross-Entropy Loss Functions: Theoretical Analysis and Applications*. 2023. Disponível em: <<https://arxiv.org/abs/2304.07288>>. Citado na página 35.
- MCKINNEY, W. *pandas: a Foundational Python Library for Data Analysis and Statistics*. [S.l.], 2011. Disponível em: <<https://www.researchgate.net/publication/265194455>>. Citado na página 41.
- MITCHELL, T. M. *Machine Learning*. [S.l.]: McGraw-Hill, 1997. 414 p. ISBN 0070428077. Citado na página 19.
- MURPHY, K. P. *Machine learning : a probabilistic perspective*. [S.l.]: MIT Press, 2012. 1067 p. ISBN 9780262018029. Citado na página 24.
- NIELSEN, M. A. *Neural Networks and Deep Learning*. [S.l.]: Determination Press, 2015. Citado 2 vezes nas páginas 9 e 26.
- PALASUWAN, D. et al. Parasitic egg detection and classification in microscopic images. IEEE Dataport, 2022. Disponível em: <<https://dx.doi.org/10.21227/vyh8-4h71>>. Citado na página 33.
- PASZKE, A. et al. *PyTorch: An Imperative Style, High-Performance Deep Learning Library*. [S.l.], 2019. Citado 4 vezes nas páginas 25, 27, 28 e 41.
- PEDREGOSA, F. et al. *Scikit-learn: Machine Learning in Python*. 2018. Disponível em: <<https://arxiv.org/abs/1201.0490>>. Citado na página 41.
- PRINCE, S. J. D. *Understanding Deep Learning*. [S.l.]: MIT Press, 2024. Citado 2 vezes nas páginas 25 e 27.
- SUTTON, R. S.; BARTO, A. G. *Reinforcement Learning: An Introduction Second edition, in progress*. [S.l.], 2015. Citado na página 25.
- TAN, M.; LE, Q. V. *EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks*. 2020. Disponível em: <<https://arxiv.org/abs/1905.11946>>. Citado na página 36.
- TAN, M.; LE, Q. V. *EfficientNetV2: Smaller Models and Faster Training*. 2021. Disponível em: <<https://arxiv.org/abs/2104.00298>>. Citado 2 vezes nas páginas 36 e 37.
- VASWANI, A. et al. *Attention Is All You Need*. [S.l.], 2017. Citado 2 vezes nas páginas 29 e 30.
- WASKOM, M. L. seaborn: statistical data visualization. *Journal of Open Source Software*, The Open Journal, v. 6, n. 60, p. 3021, 2021. Disponível em: <<https://doi.org/10.21105/joss.03021>>. Citado na página 41.
- WU, K. et al. *TinyViT: Fast Pretraining Distillation for Small Vision Transformers*. 2022. Disponível em: <<https://arxiv.org/abs/2207.10666>>. Citado 3 vezes nas páginas 9, 38 e 39.
- XU, W. et al. A lightweight deep-learning model for parasite egg detection in microscopy images. *Parasites vectors*, v. 17, p. 454, 12 2024. ISSN 17563305. Citado 3 vezes nas páginas 20, 53 e 57.

---

YIMER, M. et al. Evaluation performance of diagnostic methods of intestinal parasitosis in school age children in ethiopia. *BMC Research Notes*, BioMed Central, v. 8, 12 2015. ISSN 17560500. Citado na página [32](#).



## Apêndices



# APÊNDICE A – Código Fonte

## A.1 Arquivo Principal: *run\_optimized.py*

Listing A.1 – Arquivo principal do experimento

```

1  #!/usr/bin/env python3
2
3  import torch
4  import numpy as np
5  import matplotlib.pyplot as plt
6  import seaborn as sns
7  import pandas as pd
8  from sklearn.metrics import classification_report
9  import json
10 import os
11 import sys
12 import argparse
13
14 from trainer_optimized import ParasiteTrainerOptimized
15 from models import get_model_summary
16 import config_optimized
17
18 def set_random_seed(seed=42):
19     # Definir semente aleatória para reprodução
20     torch.manual_seed(seed)
21     torch.cuda.manual_seed(seed)
22     torch.cuda.manual_seed_all(seed)
23     np.random.seed(seed)
24     torch.backends.cudnn.deterministic = True
25     torch.backends.cudnn.benchmark = False
26
27 def train_single_model(model_type: str, force_train=False):
28     # Treinar um único modelo
29     set_random_seed()
30
31     print(f"\n{'='*60}")
32     print(f"TREINANDO MODELO {model_type.upper()}")
33     print(f"{'='*60}")
34

```

```
35     cfg = config_optimized.config
36     trainer = ParasiteTrainerOptimized(model_type, cfg)
37
38     # Verificar se existe um modelo salvo e se o treinamento é
39     ↪ forçado
40     model_path = f"{cfg.model_save_dir}/{model_type}_best.pth"
41     if os.path.exists(model_path) and not force_train:
42         print(f"Modelo {model_type.upper()} já treinado
43         ↪ encontrado em {model_path}")
44         print("Pulando treinamento e indo direto para avaliação
45         ↪ ...")
46         print("Use --force-train para treinar novamente.")
47         checkpoint = torch.load(model_path, map_location='cpu')
48         trainer.best_val_accuracy = checkpoint.get('
49         ↪ best_val_accuracy', 0.0)
50         print(f"Melhor acurácia de validação carregada: {trainer.
51         ↪ best_val_accuracy:.4f}")
52     else:
53         if force_train and os.path.exists(model_path):
54             print(f"Forçando retreinamento do modelo {model_type.
55             ↪ upper()}...")
56         else:
57             print(f"Treinando modelo {model_type.upper()}...")
58             trainer.train()
59
60     # Avaliar no conjunto de teste
61     print(f"\nAvaliando modelo {model_type.upper()} no conjunto
62     ↪ de teste...")
63     test_results = trainer.evaluate()
64
65     trainer.save_results(test_results)
66
67     # Plotar histórico de treinamento
68     trainer.plot_training_history()
69
70     # Plotar matriz de confusão
71     class_names = list(trainer.class_to_idx.keys())
72     trainer.plot_confusion_matrix(test_results['confusion_matrix'
73     ↪ ], class_names)
74
75     # Plotar análise detalhada de acurácia de teste
76     trainer.plot_test_accuracy_analysis(test_results)
```

```

69
70     print(f"Resultados do Modelo {model_type.upper()}:")
71     print(f"Acurácia de Teste: {test_results['accuracy']:.4f}")
72     print(f"Perda de Teste: {test_results['loss']:.4f}")
73     print(f"Melhor Acurácia de Validação: {trainer.
74           ↪ best_val_accuracy:.4f}")
75
76     return {
77         'accuracy': test_results['accuracy'],
78         'loss': test_results['loss'],
79         'classification_report': test_results['
80           ↪ classification_report'],
81         'best_val_accuracy': trainer.best_val_accuracy
82     }
83
84 def train_all_models(force_train=False):
85     # Treinar todos os três modelos
86     set_random_seed()
87
88     models = ['cnn', 'vit', 'hybrid']
89     results = {}
90
91     for model_type in models:
92         print(f"\n{'='*60}")
93         print(f"TREINANDO MODELO {model_type.upper()}")
94         print(f"{'='*60}")
95
96         cfg = config_optimized.config
97         trainer = ParasiteTrainerOptimized(model_type, cfg)
98
99         model_path = f"{cfg.model_save_dir}/{model_type}_best.pth
100           ↪ "
101         if os.path.exists(model_path) and not force_train:
102             print(f"Modelo {model_type.upper()} já treinado
103               ↪ encontrado em {model_path}")
104             print("Pulando treinamento e indo direto para avaliaç
105               ↪ ão...")
106             print("Use --force-train para treinar novamente.")
107             checkpoint = torch.load(model_path, map_location='cpu
108               ↪ ')
109             trainer.best_val_accuracy = checkpoint.get('
110               ↪ best_val_accuracy', 0.0)

```

```
104         print(f"Melhor acurácia de validação carregada: {  
            ⇨ trainer.best_val_accuracy:.4f}")  
105     else:  
106         if force_train and os.path.exists(model_path):  
107             print(f"Forçando retreinamento do modelo {  
                ⇨ model_type.upper()}...")  
108         else:  
109             print(f"Treinando modelo {model_type.upper()}..."  
                ⇨ )  
110         trainer.train()  
111  
112     print(f"\nAvaliando modelo {model_type.upper()} no  
        ⇨ conjunto de teste...")  
113     test_results = trainer.evaluate()  
114  
115     trainer.save_results(test_results)  
116  
117     trainer.plot_training_history()  
118  
119     class_names = list(trainer.class_to_idx.keys())  
120     trainer.plot_confusion_matrix(test_results['  
        ⇨ confusion_matrix'], class_names)  
121  
122     trainer.plot_test_accuracy_analysis(test_results)  
123  
124     results[model_type] = {  
125         'accuracy': test_results['accuracy'],  
126         'loss': test_results['loss'],  
127         'classification_report': test_results['  
            ⇨ classification_report'],  
128         'best_val_accuracy': trainer.best_val_accuracy  
129     }  
130  
131     print(f"Resultados do Modelo {model_type.upper()}:")  
132     print(f"Acurácia de Teste: {test_results['accuracy']:.4f}  
        ⇨ ")  
133     print(f"Perda de Teste: {test_results['loss']:.4f}")  
134     print(f"Melhor Acurácia de Validação: {trainer.  
        ⇨ best_val_accuracy:.4f}")  
135  
136     return results  
137
```

```

138 def compare_models(results):
139     # Comparar e visualizar resultados de todos os modelos
140     model_names_pt = {
141         'cnn': 'CNN',
142         'vit': 'ViT',
143         'hybrid': 'Hibrido'
144     }
145
146     comparison_data = []
147     for model_type, result in results.items():
148         comparison_data.append({
149             'Model': model_names_pt.get(model_type, model_type.
150                 ↪ upper()),
151             'Test Accuracy': result['accuracy'],
152             'Test Loss': result['loss'],
153             'Best Val Accuracy': result['best_val_accuracy']
154         })
155
156     df = pd.DataFrame(comparison_data)
157     print("\nComparação de Performance dos Modelos:")
158     print(df.to_string(index=False))
159
160     # Criar gráficos de comparação
161     fig, ((ax1, ax2), (ax3, ax4)) = plt.subplots(2, 2, figsize
162         ↪ =(16, 12))
163
164     # Comparação de acurácia
165     models = [data['Model'] for data in comparison_data]
166     test_accuracies = [data['Test Accuracy'] for data in
167         ↪ comparison_data]
168     val_accuracies = [data['Best Val Accuracy'] for data in
169         ↪ comparison_data]
170
171     x = np.arange(len(models))
172     width = 0.35
173
174     bars1 = ax1.bar(x - width/2, test_accuracies, width, label='
175         ↪ Acurácia de Teste', alpha=0.8, color='#2E86AB')
176     bars2 = ax1.bar(x + width/2, val_accuracies, width, label='
177         ↪ Acurácia de Validação', alpha=0.8, color='#A23B72')
178     ax1.set_xlabel('Modelos')
179     ax1.set_ylabel('Acurácia')

```

```

174     ax1.set_title('Comparação de Acurácia dos Modelos')
175     ax1.set_xticks(x)
176     ax1.set_xticklabels(models)
177     ax1.legend()
178     ax1.grid(True, alpha=0.3)
179
180     # Adicionar rótulos de valores
181     for bar, acc in zip(bars1, test_accuracies):
182         height = bar.get_height()
183         ax1.text(bar.get_x() + bar.get_width()/2., height + 0.01,
184                 f'{acc:.3f}', ha='center', va='bottom',
185                 ↪ fontweight='bold')
186
187     for bar, acc in zip(bars2, val_accuracies):
188         height = bar.get_height()
189         ax1.text(bar.get_x() + bar.get_width()/2., height + 0.01,
190                 f'{acc:.3f}', ha='center', va='bottom',
191                 ↪ fontweight='bold')
192
193     # Comparação de perda
194     test_losses = [data['Test Loss'] for data in comparison_data]
195     bars3 = ax2.bar(models, test_losses, alpha=0.8, color='#
196         ↪ F18F01')
197     ax2.set_xlabel('Modelos')
198     ax2.set_ylabel('Perda')
199     ax2.set_title('Comparação de Perda dos Modelos')
200     ax2.grid(True, alpha=0.3)
201
202     # Adicionar rótulos de valores
203     for bar, loss in zip(bars3, test_losses):
204         height = bar.get_height()
205         ax2.text(bar.get_x() + bar.get_width()/2., height + 0.05,
206                 f'{loss:.3f}', ha='center', va='bottom',
207                 ↪ fontweight='bold')
208
209     # Análise detalhada por classe
210     accuracy_data = []
211     for model_type, result in results.items():
212         report = result['classification_report']
213         for class_name, metrics in report.items():
214             if isinstance(metrics, dict) and 'precision' in
215                 ↪ metrics:

```

```

211         accuracy_data.append({
212             'Model': model_names_pt.get(model_type,
213                                     ↪ model_type.upper()),
214             'Class': class_name,
215             'Precision': metrics['precision'],
216             'Recall': metrics['recall'],
217             'F1-Score': metrics['f1-score']
218         })
219
220 accuracy_df = pd.DataFrame(accuracy_data)
221
222 # Comparação de F1-Score por classe
223 pivot_f1 = accuracy_df.pivot(index='Class', columns='Model',
224                               ↪ values='F1-Score')
225 pivot_f1.plot(kind='bar', ax=ax3, alpha=0.8, color=['#2E86AB',
226                               ↪ , '#A23B72', '#F18F01'])
227 ax3.set_title('F1-Score por Classe e Modelo')
228 ax3.set_xlabel('Classes')
229 ax3.set_ylabel('F1-Score')
230 ax3.legend(title='Modelo')
231 ax3.tick_params(axis='x', rotation=45)
232 ax3.grid(True, alpha=0.3)
233
234 # Gráfico de dispersão de precisão vs revocação
235 for i, model in enumerate(models):
236     model_data = accuracy_df[accuracy_df['Model'] == model]
237     ax4.scatter(model_data['Precision'], model_data['Recall',
238                               ↪ ],
239                 s=100, alpha=0.7, label=model,
240                 color=['#2E86AB', '#A23B72', '#F18F01'][i])
241
242 ax4.set_title('Precisão vs Revocação por Classe')
243 ax4.set_xlabel('Precisão')
244 ax4.set_ylabel('Revocação')
245 ax4.legend()
246 ax4.grid(True, alpha=0.3)
247
248 # Adicionar linha diagonal de referência
249 ax4.plot([0, 1], [0, 1], 'k--', alpha=0.5, label='Linha de
250                               ↪ Referência')
251
252 plt.tight_layout()

```

```

248     cfg = config_optimized.config
249     plt.savefig(f"{cfg.results_dir}/model_comparison_optimized.
        ↳ png", dpi=300, bbox_inches='tight')
250     plt.show()
251
252     comparison_path = f"{cfg.results_dir}/
        ↳ model_comparison_optimized.json"
253     with open(comparison_path, 'w') as f:
254         json.dump(results, f, indent=4)
255
256     print(f"\nResultados da comparação salvos em {comparison_path
        ↳ }")
257
258     # Imprimir melhor modelo
259     best_model = max(results.items(), key=lambda x: x[1]['
        ↳ accuracy'])
260     print(f"\nMELHOR MODELO: {best_model[0].upper()}")
261     print(f"    Acurácia de Teste: {best_model[1]['accuracy']:.4f}
        ↳ ")
262     print(f"    Perda de Teste: {best_model[1]['loss']:.4f}")
263
264     # Analisar sobreajuste
265     print(f"\nANALISE DE SOBREAJUSTE:")
266     for model_type, result in results.items():
267         val_acc = result['best_val_accuracy'] / 100.0
268         test_acc = result['accuracy']
269         overfitting_gap = val_acc - test_acc
270         print(f"    {model_type.upper()}: Validação {val_acc:.3f}
            ↳ -> Teste {test_acc:.3f} (Gap: {overfitting_gap:.3f})
            ↳ ")
271
272 def main():
273     parser = argparse.ArgumentParser(description='Treinar modelos
        ↳ de detecção de parasitas')
274     parser.add_argument('--model', type=str, choices=['cnn', 'vit
        ↳ ', 'hybrid', 'all'],
275                         default='all', help='Modelo a ser treinado
        ↳ (padrão: all)')
276     parser.add_argument('--force-train', action='store_true',
        ↳ help='Forçar retreinamento de modelos que já existem')
277
278     args = parser.parse_args()

```

```

279
280     print(f"\n{'='*70}")
281     print("EXPERIMENTO")
282     print(f"{'='*70}")
283     print("Configurações:")
284     cfg = config_optimized.config
285     print(f"    - Tamanho da imagem: {cfg.image_size}x{cfg.
      ↪ image_size}")
286     print(f"    - Tamanho do batch: {cfg.batch_size}")
287     print(f"    - Taxa de aprendizado: {cfg.learning_rate}")
288     print(f"    - Taxa de decaimento: {cfg.weight_decay}")
289     print(f"    - Suavizacao de rotulos: {cfg.label_smoothing}")
290     print(f"    - Gradiente: {cfg.gradient_clip}")
291     print(f"    - Agendador de taxa de aprendizado:
      ↪ ReduceLROnPlateau")
292     print(f"    - Tratamento de dados: Embaçamento + Ruído")
293     print(f"{'='*70}")
294
295     if args.model == 'all':
296         results = train_all_models(args.force_train)
297         compare_models(results)
298     else:
299         result = train_single_model(args.model, args.force_train)
300         compare_models({args.model: result})
301
302 if __name__ == "__main__":
303     main()

```

## A.2 Treinador Otimizado: *trainer\_optimized.py*

Listing A.2 – Classe do treinador otimizado

```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 from torch.optim.lr_scheduler import ReduceLROnPlateau
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 from sklearn.metrics import classification_report,
  ↪ confusion_matrix

```

```
9 import json
10 import os
11 from tqdm import tqdm
12 import config_optimized
13 from dataset_optimized import get_dataloaders_optimized
14 from models import get_model, get_model_summary
15
16 class ParasiteTrainerOptimized:
17     def __init__(self, model_type: str, config):
18         self.model_type = model_type
19         self.config = config
20         self.device = torch.device(config.device)
21
22         # Obter dataloaders
23         self.train_loader, self.val_loader, self.test_loader,
24         ↪ self.class_to_idx = get_dataloaders_optimized(config
25         ↪ )
26
27         self.model = self._create_model()
28         self.model.to(self.device)
29
30         # Função de perda com suavização de rótulos
31         self.criterion = nn.CrossEntropyLoss(label_smoothing=
32         ↪ config.label_smoothing)
33
34         # Otimizador
35         self.optimizer = optim.AdamW(
36             self.model.parameters(),
37             lr=config.learning_rate,
38             weight_decay=config.weight_decay
39         )
40
41         # Agendador de taxa de aprendizado
42         self.scheduler = ReduceLROnPlateau(
43             self.optimizer,
44             mode='max',
45             factor=config.scheduler_factor,
46             patience=config.scheduler_patience,
47             min_lr=config.scheduler_min_lr
48         )
49
50         self.train_losses = []
```

```
48     self.val_losses = []
49     self.train_accuracies = []
50     self.val_accuracies = []
51     self.best_val_accuracy = 0.0
52     self.patience_counter = 0
53
54     def _create_model(self):
55         # Criar modelo com base no tipo
56         from models import get_model, get_model_summary
57         model = get_model(self.model_type, self.config)
58         get_model_summary(model, input_size=(3, self.config.
59             ⇨ image_size, self.config.image_size))
60         return model
61
62     def get_model_name_pt(self):
63         # Traduzir
64         model_names = {
65             'cnn': 'CNN',
66             'vit': 'ViT',
67             'hybrid': 'Híbrido'
68         }
69         return model_names.get(self.model_type, self.model_type.
70             ⇨ upper())
71
72     def train_epoch(self):
73         # Treinar por uma época
74         self.model.train()
75         total_loss = 0.0
76         correct = 0
77         total = 0
78
79         progress_bar = tqdm(self.train_loader, desc=f"Treinando {
80             ⇨ self.model_type.upper()}")
81
82         for batch_idx, (data, target) in enumerate(progress_bar):
83             data, target = data.to(self.device), target.to(self.
84                 ⇨ device)
85
86             self.optimizer.zero_grad()
87             output = self.model(data)
88             loss = self.criterion(output, target)
```

```

86         torch.nn.utils.clip_grad_norm_(self.model.parameters
87             ↪ (), self.config.gradient_clip)
88
89         loss.backward()
90         self.optimizer.step()
91
92         total_loss += loss.item()
93         pred = output.argmax(dim=1, keepdim=True)
94         correct += pred.eq(target.view_as(pred)).sum().item()
95         total += target.size(0)
96
97         progress_bar.set_postfix({
98             'Loss': f'{loss.item():.4f}',
99             'Acc': f'{100. * correct / total:.2f}%'
100         })
101
102     return total_loss / len(self.train_loader), correct /
103         ↪ total
104
105 def validate_epoch(self):
106     # Validar época
107     self.model.eval()
108     total_loss = 0.0
109     correct = 0
110     total = 0
111
112     with torch.no_grad():
113         for data, target in tqdm(self.val_loader, desc=f"
114             ↪ Validando {self.model_type.upper()}"):
115             data, target = data.to(self.device), target.to(
116                 ↪ self.device)
117             output = self.model(data)
118             loss = self.criterion(output, target)
119
120             total_loss += loss.item()
121             pred = output.argmax(dim=1, keepdim=True)
122             correct += pred.eq(target.view_as(pred)).sum().
123                 ↪ item()
124             total += target.size(0)
125
126     return total_loss / len(self.val_loader), correct / total

```

```
123     def train(self):
124         # Treinar com parada antecipada
125         print(f"\n{'='*60}")
126         print(f"TREINANDO MODELO {self.model_type.upper()}")
127         print(f"{'='*60}")
128
129         for epoch in range(self.config.num_epochs):
130             print(f"\nEpoca {epoch+1}/{self.config.num_epochs}")
131
132             train_loss, train_acc = self.train_epoch()
133
134             val_loss, val_acc = self.validate_epoch()
135
136             self.scheduler.step(val_acc)
137
138             # Armazenar histórico
139             self.train_losses.append(train_loss)
140             self.val_losses.append(val_loss)
141             self.train_accuracies.append(train_acc)
142             self.val_accuracies.append(val_acc)
143
144             print(f"Perda de Treino: {train_loss:.4f}, Acurácia
145                 ↪ de Treino: {train_acc:.4f}")
146             print(f"Perda de Validação: {val_loss:.4f}, Acurácia
147                 ↪ de Validação: {val_acc:.4f}")
148             print(f"Taxa de Aprendizado: {self.optimizer.
149                 ↪ param_groups[0]['lr']:.2e}")
150
151             # Parada antecipada
152             if val_acc > self.best_val_accuracy:
153                 self.best_val_accuracy = val_acc
154                 self.patience_counter = 0
155                 # Salvar melhor modelo
156                 torch.save(self.model.state_dict(),
157                             f"{self.config.model_save_dir}/{self.
158                             ↪ model_type}_best.pth")
159                 print(f"Novo melhor modelo salvo! Acurácia: {
160                     ↪ val_acc:.4f}")
161             else:
162                 self.patience_counter += 1
163                 print(f"Parada Antecipada: {self.patience_counter
164                     ↪ }/{self.config.patience}")
```

```

159         if self.patience_counter >= self.config.patience:
160             print(f"Parada Antecipada ativada após {epoch
161                 ↪ +1} épocas")
162             break
163
164         print(f"\nMelhor acurácia de validação: {self.
165             ↪ best_val_accuracy:.4f}")
166
167     def evaluate(self):
168         # Avaliar no conjunto de teste
169         print(f"\nAvaliando modelo {self.model_type.upper()} no
170             ↪ conjunto de teste...")
171
172         # Carregar melhor modelo
173         best_model_path = f"{self.config.model_save_dir}/{self.
174             ↪ model_type}_best.pth"
175         if os.path.exists(best_model_path):
176             self.model.load_state_dict(torch.load(best_model_path
177                 ↪ , map_location=self.device))
178             print("Modelo melhor carregado para avaliação")
179
180         self.model.eval()
181         test_loss = 0.0
182         all_predictions = []
183         all_targets = []
184
185         with torch.no_grad():
186             for data, target in tqdm(self.test_loader, desc="
187                 ↪ Avaliando"):
188                 data, target = data.to(self.device), target.to(
189                     ↪ self.device)
190                 output = self.model(data)
191                 loss = self.criterion(output, target)
192
193                 test_loss += loss.item()
194                 pred = output.argmax(dim=1, keepdim=True)
195
196                 all_predictions.extend(pred.cpu().numpy().flatten
197                     ↪ ())
198                 all_targets.extend(target.cpu().numpy())

```

```
193     test_loss /= len(self.test_loader)
194     test_accuracy = sum(1 for x, y in zip(all_predictions,
    ↪ all_targets) if x == y) / len(all_targets)
195
196     # Relatório de classificação
197     class_names = list(self.class_to_idx.keys())
198     report = classification_report(all_targets,
    ↪ all_predictions,
199                                   target_names=class_names,
    ↪ output_dict=True)
200
201     # Matriz de confusão
202     cm = confusion_matrix(all_targets, all_predictions)
203
204     return {
205         'accuracy': test_accuracy,
206         'loss': test_loss,
207         'classification_report': report,
208         'confusion_matrix': cm,
209         'predictions': all_predictions,
210         'targets': all_targets
211     }
212
213     def save_results(self, results):
214         # Salvar resultados em JSON
215         results_path = f"{self.config.results_dir}/{self.
    ↪ model_type}_results.json"
216
217         # Função auxiliar para converter tipos numpy para tipos
    ↪ nativos Python
218         def convert_numpy_types(obj):
219             if isinstance(obj, np.ndarray):
220                 if obj.size == 1:
221                     return obj.item()
222                 else:
223                     return obj.tolist()
224             elif isinstance(obj, np.integer):
225                 return int(obj)
226             elif isinstance(obj, np.floating):
227                 return float(obj)
228             elif isinstance(obj, dict):
```

```
229         return {key: convert_numpy_types(value) for key,
230                 ↪ value in obj.items()}
231     elif isinstance(obj, list):
232         return [convert_numpy_types(item) for item in obj
233                 ↪ ]
234     else:
235         return obj
236
237     # Converter todos os resultados para formato JSON-
238     ↪ serializável
239     results_to_save = {
240         'acuracia': convert_numpy_types(results['accuracy']),
241         'perda': convert_numpy_types(results['loss']),
242         'relatorio_classificacao': convert_numpy_types(
243             ↪ results['classification_report']),
244         'matriz_confusao': convert_numpy_types(results['
245             ↪ confusion_matrix']),
246         'predicoes': convert_numpy_types(results['predictions
247             ↪ '])
248     }
249
250     with open(results_path, 'w', encoding='utf-8') as f:
251         json.dump(results_to_save, f, indent=4, ensure_ascii=
252             ↪ False)
253
254     print(f"Resultados salvos em {results_path}")
255
256     def plot_training_history(self):
257         # Plotar histórico de treinamento
258         fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 5))
259         num_epochs = len(self.train_losses)
260         epochs = list(range(1, num_epochs + 1))
261
262         def get_epoch_ticks(num_epochs):
263             if num_epochs <= 20:
264                 return list(range(1, num_epochs + 1))
265             elif num_epochs <= 50:
266                 ticks = list(range(1, num_epochs + 1, 5))
267             else:
268                 ticks = list(range(1, num_epochs + 1, 10))
269             if ticks[-1] != num_epochs:
270                 ticks.append(num_epochs)
```

```

264         return ticks
265     epoch_ticks = get_epoch_ticks(num_epochs)
266
267     # Plotar perda
268     ax1.plot(epochs, self.train_losses, label='Treino', color
269             ⇨ ='blue', marker='o')
270     ax1.plot(epochs, self.val_losses, label='Validação',
271             ⇨ color='red', marker='o')
272     ax1.set_title(f'Histórico de Perda - {self.
273             ⇨ get_model_name_pt()}')
274     ax1.set_xlabel('Epoca')
275     ax1.set_ylabel('Perda')
276     ax1.legend()
277     ax1.grid(True, alpha=0.3)
278     max_loss = max(self.train_losses + self.val_losses) if (
279             ⇨ self.train_losses and self.val_losses) else 1
280     ax1.set_ylim(0, max_loss * 1.1)
281     ax1.set_xticks(epoch_ticks)
282
283     # Plotar acurácia
284     train_acc_pct = [acc * 100 for acc in self.
285             ⇨ train_accuracies]
286     val_acc_pct = [acc * 100 for acc in self.val_accuracies]
287     ax2.plot(epochs, train_acc_pct, label='Treino', color='
288             ⇨ blue', marker='o')
289     ax2.plot(epochs, val_acc_pct, label='Validação', color='
290             ⇨ red', marker='o')
291     ax2.set_title(f'Histórico de Acurácia - {self.
292             ⇨ get_model_name_pt()}')
293     ax2.set_xlabel('Epoca')
294     ax2.set_ylabel('Acurácia (%)')
295     ax2.legend()
296     ax2.grid(True, alpha=0.3)
297     ax2.set_ylim(0, 100)
298     ax2.set_yticks([0, 20, 40, 60, 80, 100])
299     ax2.set_yticklabels(['0%', '20%', '40%', '60%', '80%', '
300             ⇨ 100%'])
301     ax2.set_xticks(epoch_ticks)
302
303     plt.tight_layout()
304     save_path = f"{self.config.results_dir}/{self.model_type}
305             ⇨ _training_history.png"

```

```

296     plt.savefig(save_path, dpi=300, bbox_inches='tight')
297     plt.close()
298     print(f"Gráfico de histórico salvo em {save_path}")
299
300     def plot_confusion_matrix(self, cm, class_names):
301         # Plotar matriz de confusão
302         plt.figure(figsize=(12, 10))
303         sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
304                     xticklabels=class_names, yticklabels=
305                         ↪ class_names)
306         plt.title(f'Matriz de Confusão - {self.get_model_name_pt
307                     ↪ ()}')
308         plt.xlabel('Predição')
309         plt.ylabel('Real')
310         plt.xticks(rotation=45)
311         plt.yticks(rotation=0)
312
313         save_path = f"{self.config.results_dir}/{self.model_type}
314                     ↪ _confusion_matrix.png"
315         plt.savefig(save_path, dpi=300, bbox_inches='tight')
316         plt.close()
317         print(f"Matriz de confusão salva em {save_path}")
318
319     def plot_test_accuracy_analysis(self, test_results):
320         # Plotar análise detalhada de acurácia de teste
321         class_names = list(self.class_to_idx.keys())
322         report = test_results['classification_report']
323
324         classes = []
325         precisions = []
326         recalls = []
327         f1_scores = []
328
329         for class_name, metrics in report.items():
330             if isinstance(metrics, dict) and 'precision' in
331                 ↪ metrics:
332                 classes.append(class_name)
333                 precisions.append(metrics['precision'])
334                 recalls.append(metrics['recall'])
335                 f1_scores.append(metrics['f1-score'])

```

```

333     fig, ((ax1, ax2), (ax3, ax4)) = plt.subplots(2, 2,
334           ⇨ figsize=(16, 12))
335
336     # 1. Acurácia geral de teste
337     test_acc = test_results['accuracy']
338     test_loss = test_results['loss']
339
340     ax1.bar(['Acurácia', 'Perda'], [test_acc, test_loss],
341           color=['#2E86AB', '#A23B72'], alpha=0.8)
342     ax1.set_title(f'Métricas Gerais de Teste - {self.
343           ⇨ get_model_name_pt()}')
344     ax1.set_ylabel('Valor')
345     ax1.grid(True, alpha=0.3)
346
347     for i, (acc, loss) in enumerate([(test_acc, test_loss)]):
348         ax1.text(0, acc + 0.01, f'{acc:.3f}', ha='center', va
349               ⇨ ='bottom', fontweight='bold')
350         ax1.text(1, loss + 0.01, f'{loss:.3f}', ha='center',
351               ⇨ va='bottom', fontweight='bold')
352
353     # 2. Precisão por classe
354     bars1 = ax2.bar(range(len(classes)), precisions, alpha
355           ⇨ =0.8, color='#2E86AB')
356     ax2.set_title(f'Precisão por Classe - {self.
357           ⇨ get_model_name_pt()}')
358     ax2.set_xlabel('Classes')
359     ax2.set_ylabel('Precisão')
360     ax2.set_xticks(range(len(classes)))
361     ax2.set_xticklabels(classes, rotation=45, ha='right')
362     ax2.grid(True, alpha=0.3)
363
364     for bar, prec in zip(bars1, precisions):
365         height = bar.get_height()
366         ax2.text(bar.get_x() + bar.get_width()/2., height +
367               ⇨ 0.01,
368               f'{prec:.3f}', ha='center', va='bottom',
369               ⇨ fontweight='bold', fontsize=8)
370
371     # 3. Revocação por classe
372     bars2 = ax3.bar(range(len(classes)), recalls, alpha=0.8,
373           ⇨ color='#A23B72')

```

```

365     ax3.set_title(f'Revocação por Classe - {self.
        ↪ get_model_name_pt()}')
366     ax3.set_xlabel('Classes')
367     ax3.set_ylabel('Revocação')
368     ax3.set_xticks(range(len(classes)))
369     ax3.set_xticklabels(classes, rotation=45, ha='right')
370     ax3.grid(True, alpha=0.3)
371
372     for bar, rec in zip(bars2, recalls):
373         height = bar.get_height()
374         ax3.text(bar.get_x() + bar.get_width()/2., height +
            ↪ 0.01,
375                 f'{rec:.3f}', ha='center', va='bottom',
            ↪ fontweight='bold', fontsize=8)
376
377     # 4. F1-Score por classe
378     bars3 = ax4.bar(range(len(classes)), f1_scores, alpha
        ↪ =0.8, color='#F18F01')
379     ax4.set_title(f'F1-Score por Classe - {self.
        ↪ get_model_name_pt()}')
380     ax4.set_xlabel('Classes')
381     ax4.set_ylabel('F1-Score')
382     ax4.set_xticks(range(len(classes)))
383     ax4.set_xticklabels(classes, rotation=45, ha='right')
384     ax4.grid(True, alpha=0.3)
385
386     for bar, f1 in zip(bars3, f1_scores):
387         height = bar.get_height()
388         ax4.text(bar.get_x() + bar.get_width()/2., height +
            ↪ 0.01,
389                 f'{f1:.3f}', ha='center', va='bottom',
            ↪ fontweight='bold', fontsize=8)
390
391     plt.tight_layout()
392     save_path = f"{self.config.results_dir}/{self.model_type}
        ↪ _test_accuracy_analysis.png"
393     plt.savefig(save_path, dpi=300, bbox_inches='tight')
394     plt.close()
395     print(f"Análise detalhada de acurácia de teste salva em {
        ↪ save_path}")

```

## A.3 Configuração: *config\_optimized.py*

Listing A.3 – Arquivo de configuração otimizado

```
1 import os
2 import torch
3 from dataclasses import dataclass
4 from typing import List, Tuple
5
6 class Config:
7
8     # Dados
9     data_dir = "/home/edvl/TCC/Chula-ParasiteEgg-11/Chula-
10         ↪ ParasiteEgg-11/Chula-ParasiteEgg-11/data"
11     test_data_dir = "/home/edvl/TCC/Chula-ParasiteEgg-11_test/
12         ↪ test/data"
13     train_data_path = "/home/edvl/TCC/Chula-ParasiteEgg-11/Chula-
14         ↪ ParasiteEgg-11/Chula-ParasiteEgg-11/data"
15     test_data_path = "/home/edvl/TCC/Chula-ParasiteEgg-11_test/
16         ↪ test/data"
17
18     # Modelo
19     num_classes = 11
20     input_size = (384, 384)
21     image_size = 384
22
23     # Treinamento
24     batch_size = 2
25     num_epochs = 50
26     learning_rate = 1e-4
27     weight_decay = 1e-4
28     gradient_clip = 1.0
29
30     # Tratamento de dados
31     use_augmentation = True
32     train_transform = True
33     test_transform = False
34     mixup_alpha = 0.2
35     cutmix_alpha = 1.0
36     cutmix_prob = 0.5
37
38     # Aprendizado em conjunto
```

```
35 ensemble_size = 3
36 diversity_weight = 0.1
37
38 # Regularização
39 dropout_rate = 0.3
40 label_smoothing = 0.1
41
42 # Parada antecipada
43 patience = 10
44
45 # Memória e processamento
46 num_workers = 1
47 pin_memory = True
48
49 # Dispositivo
50 device = 'cuda' if torch.cuda.is_available() else 'cpu'
51
52 # Configurações específicas dos modelos
53 # CNN (EfficientNetV2-S)
54 cnn_model_name = "tf_efficientnetv2_s"
55 cnn_dropout = 0.2
56
57 # Vision Transformer (Tiny ViT)
58 vit_model_name = "vit_tiny_patch16_224"
59 vit_patch_size = 16
60 vit_embed_dim = 192
61 vit_depths = (3, 3, 3)
62 vit_num_heads = (3, 6, 12)
63
64 # Modelo Híbrido
65 hybrid_cnn_backbone = "tf_efficientnetv2_s"
66 hybrid_vit_model = "vit_tiny_patch16_224"
67 hybrid_fusion_dim = 64
68
69 # Diretórios de saída
70 model_save_dir = "models_optimized"
71 results_dir = "results_optimized"
72
73 # Registro da execução
74 log_interval = 100
75 save_interval = 5
76
```

```
77 config = Config()
```

## A.4 Modelos: *models.py*

Listing A.4 – Implementação dos modelos

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import timm
5
6 class EfficientNetV2CNN(nn.Module):
7     # Modelo CNN usando EfficientNetV2-S para classificação de
8     # ↳ parasitas
9
10    def __init__(self, config, num_classes: int = 11):
11        super(EfficientNetV2CNN, self).__init__()
12
13        # Carregar EfficientNetV2-S
14        self.backbone = timm.create_model(
15            config.cnn_model_name,
16            pretrained=True,
17            num_classes=0
18        )
19
20        feature_dim = self.backbone.num_features
21
22        self.classifier = nn.Sequential(
23            nn.Dropout(config.cnn_dropout),
24            nn.Linear(feature_dim, 512),
25            nn.ReLU(),
26            nn.Dropout(config.cnn_dropout),
27            nn.Linear(512, num_classes)
28        )
29
30    def forward(self, x):
31        features = self.backbone(x)
32        output = self.classifier(features)
33        return output
34
35 class TinyViTTransformer(nn.Module):
```

```

35     # Modelo Vision Transformer usando Tiny ViT para classificação
      ↳ o de parasitas
36
37     def __init__(self, config, num_classes: int = 11):
38         super(TinyViTTransformer, self).__init__()
39
40         # Carregar Tiny ViT
41         self.backbone = timm.create_model(
42             config.vit_model_name,
43             pretrained=True,
44             num_classes=0
45         )
46
47         feature_dim = self.backbone.num_features
48
49         self.classifier = nn.Sequential(
50             nn.LayerNorm(feature_dim),
51             nn.Linear(feature_dim, 512),
52             nn.GELU(),
53             nn.Dropout(0.1),
54             nn.Linear(512, num_classes)
55         )
56
57     def forward(self, x):
58         features = self.backbone(x)
59         output = self.classifier(features)
60         return output
61
62 class HybridModel(nn.Module):
63     # Modelo Híbrido combinando CNN (EfficientNetV2-S) e Vision
      ↳ Transformer (Tiny ViT)
64     def __init__(self, config, num_classes: int = 11):
65         super(HybridModel, self).__init__()
66
67         # (EfficientNetV2-S)
68         self.cnn_backbone = timm.create_model(
69             config.hybrid_cnn_backbone,
70             pretrained=True,
71             num_classes=0
72         )
73
74         # Vision Transformer (Tiny ViT)

```

```
75     self.vit_backbone = timm.create_model(
76         config.hybrid_vit_model,
77         pretrained=True,
78         num_classes=0
79     )
80
81     cnn_feature_dim = self.cnn_backbone.num_features
82     vit_feature_dim = self.vit_backbone.num_features
83
84     # Fusão de características
85     fusion_dim = config.hybrid_fusion_dim
86     self.fusion = nn.Sequential(
87         nn.Linear(cnn_feature_dim + vit_feature_dim,
88             ↪ fusion_dim),
89         nn.LayerNorm(fusion_dim),
90         nn.GELU(),
91         nn.Dropout(0.1),
92         nn.Linear(fusion_dim, fusion_dim // 2),
93         nn.GELU(),
94         nn.Dropout(0.1)
95     )
96
97     # Classificador final
98     self.classifier = nn.Linear(fusion_dim // 2, num_classes)
99
100    # Mecanismo de atenção para ponderação de características
101    self.attention = nn.MultiheadAttention(
102        embed_dim=fusion_dim // 2,
103        num_heads=2,
104        dropout=0.1,
105        batch_first=True
106    )
107
108    def forward(self, x):
109        # Extrair características de ambos os modelos
110        cnn_features = self.cnn_backbone(x)
111        vit_features = self.vit_backbone(x)
112
113        # Concatenar características
114        combined_features = torch.cat([cnn_features, vit_features
            ↪ ], dim=1)
```

```

115         # Aplicar fusão
116         fused_features = self.fusion(combined_features)
117
118         # Aplicar auto-atenção
119         fused_features = fused_features.unsqueeze(1) # Adicionar
120         ↪ dimensão de sequência
121         attended_features, _ = self.attention(fused_features,
122         ↪ fused_features, fused_features)
123         attended_features = attended_features.squeeze(1)
124
125         # Classificação final
126         output = self.classifier(attended_features)
127         return output
128
129 def get_model(model_type: str, config, num_classes: int = 11):
130     # Função fábrica para obter o modelo especificado
131
132     if model_type.lower() == "cnn":
133         return EfficientNetV2CNN(config, num_classes=num_classes)
134
135     elif model_type.lower() == "vit":
136         return TinyViTTransformer(config, num_classes=num_classes
137         ↪ )
138
139     elif model_type.lower() == "hybrid":
140         return HybridModel(config, num_classes=num_classes)
141
142     else:
143         raise ValueError(f"Tipo de modelo desconhecido: {
144         ↪ model_type}. Escolha entre ['cnn', 'vit', 'hybrid']"
145         ↪ )
146
147 def count_parameters(model):
148     # Contar parâmetros treináveis no modelo
149     return sum(p.numel() for p in model.parameters() if p.
150     ↪ requires_grad)
151
152 def get_model_summary(model, input_size=(3, 224, 224)):
153     # Obter resumo do modelo com contagem de parâmetros
154     total_params = count_parameters(model)
155     trainable_params = sum(p.numel() for p in model.parameters()
156     ↪ if p.requires_grad)

```

```

150
151     print(f"Resumo do Modelo:")
152     print(f"Total de Parâmetros: {total_params:,}")
153     print(f"Parâmetros Treináveis: {trainable_params:,}")
154     print(f"Tamanho de Entrada: {input_size}")
155
156     return total_params, trainable_params

```

## A.5 Conjunto de Dados: *dataset\_optimized.py*

Listing A.5 – Classe do conjunto de dados otimizado

```

1  import os
2  import torch
3  from torch.utils.data import Dataset, DataLoader
4  from PIL import Image
5  import albumentations as A
6  from albumentations.pytorch import ToTensorV2
7  import numpy as np
8  from typing import Dict, List, Tuple, Optional
9  import config_optimized
10 import json
11
12 class ParasiteDataset(Dataset):
13     def __init__(self, data_path: str, transform=None, is_train:
14         ⇨ bool = True, class_to_idx: Optional[Dict[str, int]] =
15         ⇨ None, label_json_path: Optional[str] = None):
16         self.data_path = data_path
17         self.transform = transform
18         self.is_train = is_train
19         self.class_to_idx = class_to_idx
20         self.label_json_path = label_json_path
21
22         # Carregar todos os arquivos de imagem e seus rótulos
23         self.images, self.labels, self.class_to_idx = self.
24             ⇨ _load_dataset()
25
26     def _load_dataset(self) -> Tuple[List[str], List[int],
27         ⇨ Optional[Dict[str, int]]]:
28         images = []
29         labels = []

```

```
26
27     # Se estiver usando arquivo de rótulos no formato COCO (
    ↪ para conjunto de teste)
28     if self.label_json_path is not None:
29         # Carregar arquivo de rótulos
30         with open(self.label_json_path, 'r') as f:
31             label_data = json.load(f)
32
33         file_to_id = {img['file_name']: img['id'] for img in
    ↪ label_data['images']}
34
35         imageid_to_catid = {}
36         for ann in label_data['annotations']:
37             if ann['image_id'] not in imageid_to_catid:
38                 imageid_to_catid[ann['image_id']] = ann['
    ↪ category_id']
39
40         catid_to_name = {cat['id']: cat['name'] for cat in
    ↪ label_data['categories']}
41
42         class_to_idx = self.class_to_idx
43
44         for filename in os.listdir(self.data_path):
45             if filename.lower().endswith((''.jpg', '.jpeg', '.
    ↪ png')):
46                 img_path = os.path.join(self.data_path,
    ↪ filename)
47
48                 image_id = file_to_id.get(filename)
49                 if image_id is None:
50                     continue
51
52                 category_id = imageid_to_catid.get(image_id)
53                 if category_id is None:
54                     continue
55
56                 class_name = catid_to_name.get(category_id)
57                 if class_name is None:
58                     continue
59
60                 if class_name in class_to_idx:
61                     images.append(img_path)
```

```

62         labels.append(class_to_idx[class_name])
63     return images, labels, class_to_idx
64
65     # Comportamento padrão (conjunto de treino/validação)
66     extracted_class_names = []
67     if self.class_to_idx is None:
68         # Construir mapeamento a partir deste conjunto de
69         ↪ dados (para conjunto de treino)
70         class_names = set()
71         for filename in os.listdir(self.data_path):
72             if filename.lower().endswith((''.jpg', '.jpeg', '.
73             ↪ png')):
74                 class_name = filename.split('_')[0]
75                 class_names.add(class_name)
76                 if len(extracted_class_names) < 10:
77                     extracted_class_names.append(class_name)
78         class_names = sorted(list(class_names))
79         class_to_idx = {class_name: idx for idx, class_name
80         ↪ in enumerate(class_names)}
81         print("[DEBUG] Built class_to_idx:", class_to_idx)
82         print("[DEBUG] First 10 extracted class names from
83         ↪ filenames:", extracted_class_names)
84     else:
85         # Usar mapeamento fornecido (para conjunto de teste)
86         class_to_idx = self.class_to_idx
87
88     # Carregar imagens e atribuir rótulos usando mapeamento
89     for filename in os.listdir(self.data_path):
90         if filename.lower().endswith((''.jpg', '.jpeg', '.png'
91         ↪ )):
92             img_path = os.path.join(self.data_path, filename)
93             class_name = filename.split('_')[0]
94             if class_to_idx is not None:
95                 idx = class_to_idx.get(class_name)
96                 if idx is None:
97                     continue
98                 images.append(img_path)
99                 labels.append(idx)
100             else:
101                 continue
102     return images, labels, class_to_idx

```

```

99     def __len__(self):
100         return len(self.images)
101
102     def __getitem__(self, idx):
103         img_path = self.images[idx]
104         label = self.labels[idx]
105
106         # Carregar imagem
107         image = Image.open(img_path).convert('RGB')
108         image = np.array(image)
109
110         # Aplicar transformações
111         if self.transform:
112             transformed = self.transform(image=image)
113             image = transformed['image']
114
115         return image, label
116
117 def get_transforms_optimized(image_size: int = 384, is_train:
    ⇨ bool = True, config=None):
118     if is_train:
119         return A.Compose([
120             # Redimensionar para tamanho alvo e aplicar normaliza
    ⇨ ção
121             A.Resize(image_size, image_size),
122             A.HorizontalFlip(p=0.5),
123             A.VerticalFlip(p=0.3),
124             A.RandomRotate90(p=0.3),
125
126             # Transformações geométricas moderadas
127             A.Affine(
128                 translate_percent=0.1, # Reduzido
129                 scale=(0.9, 1.1),      # Reduzido
130                 rotate=(-15, 15),      # Reduzido baseado no
    ⇨ CoAtNet
131                 p=0.6
132             ),
133
134             # Embaçamento e ruído
135             A.OneOf([
136                 A.GaussianBlur(blur_limit=(3, 7), p=0.5),
137                 A.MotionBlur(blur_limit=3, p=0.3),

```

```

138         A.MedianBlur(blur_limit=3, p=0.2),
139     ], p=0.4),
140
141     # Ruído
142     A.OneOf([
143         A.GaussNoise(var_limit=(5.0, 15.0), p=0.5),
144         A.ISONoise(color_shift=(0.01, 0.05), intensity
145             ↪ =(0.1, 0.5), p=0.3),
146         A.MultiplicativeNoise(multiplier=(0.9, 1.1), p
147             ↪ =0.2),
148     ], p=0.4),
149
150     # Ajustes de cor moderados
151     A.OneOf([
152         A.RandomBrightnessContrast(
153             brightness_limit=0.2, # Reduzido
154             contrast_limit=0.2,   # Reduzido
155             p=0.5
156         ),
157         A.HueSaturationValue(
158             hue_shift_limit=20,    # Reduzido
159             sat_shift_limit=30,    # Reduzido
160             val_shift_limit=20,    # Reduzido
161             p=0.3
162         ),
163         A.CLAHE(clip_limit=2.0, tile_grid_size=(8, 8), p
164             ↪ =0.2),
165     ], p=0.4),
166
167     # Dropout espacial
168     A.CoarseDropout(
169         max_holes=4, max_height=16, max_width=16,
170         min_holes=1, min_height=4, min_width=4,
171         p=0.2
172     ),
173
174     # Normalização ImageNet
175     A.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
176         ↪ 0.224, 0.225]),
177     ToTensorV2(),
178 ])
```

else:

```

176         return A.Compose([
177             A.Resize(image_size, image_size),
178             A.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229,
179                 ↪ 0.224, 0.225]),
180             ToTensorV2(),
181         ])
182
183 def get_dataloaders_optimized(config):
184     # Transformações
185     train_transform = get_transforms_optimized(config.image_size,
186         ↪ is_train=True, config=config)
187     val_transform = get_transforms_optimized(config.image_size,
188         ↪ is_train=False, config=config)
189
190     # Dividir dados de treino em treino/validação (80/20)
191     full_dataset = ParasiteDataset(config.train_data_path,
192         ↪ transform=train_transform, is_train=True)
193
194     # Calcular índices de divisão
195     total_size = len(full_dataset)
196     train_size = int(0.8 * total_size)
197     val_size = total_size - train_size
198
199     train_dataset, val_dataset = torch.utils.data.random_split(
200         full_dataset, [train_size, val_size]
201     )
202
203     val_dataset.dataset.transform = val_transform
204
205     train_loader = DataLoader(
206         train_dataset,
207         batch_size=config.batch_size,
208         shuffle=True,
209         num_workers=config.num_workers,
210         pin_memory=config.pin_memory
211     )
212
213     val_loader = DataLoader(
214         val_dataset,
215         batch_size=config.batch_size,
216         shuffle=False,
217         num_workers=config.num_workers,

```

```
214     pin_memory=config.pin_memory
215 )
216
217 # Conjunto de teste
218 test_label_json = os.path.join("../", "Chula-ParasiteEgg-11
    ↪ _test", "test_labels_200.json")
219 test_dataset = ParasiteDataset(
220     config.test_data_path,
221     transform=val_transform,
222     is_train=False,
223     class_to_idx=full_dataset.class_to_idx,
224     label_json_path=test_label_json
225 )
226 test_loader = DataLoader(
227     test_dataset,
228     batch_size=config.batch_size,
229     shuffle=False,
230     num_workers=config.num_workers,
231     pin_memory=config.pin_memory
232 )
233
234 return train_loader, val_loader, test_loader, full_dataset.
    ↪ class_to_idx
```

## A.6 Requirements: *requirements.txt*

Listing A.6 – Arquivo de dependências

```
1 torch>=2.0.0
2 torchvision>=0.15.0
3 timm>=0.9.0
4 numpy>=1.24.0
5 pandas>=2.0.0
6 scikit-learn>=1.3.0
7 matplotlib>=3.7.0
8 seaborn>=0.12.0
9 Pillow>=10.0.0
10 tqdm>=4.65.0
11 albumentations>=1.3.0
12 opencv-python>=4.8.0
13 tensorboard>=2.13.0
```