# Leveraging Early Exit and Rejection in ML for Efficiency

Florian Valade

Gustave Eiffel University, Paris

September 18, 2024

# Table of Contents

# Introduction to Deep Neural Networks

- Neural networks process data through multiple layers, one after the other.
- Each layer extracts more complex features from the data.
- Deep networks can learn intricate patterns but may be computationally intensive.

# ResNet: Deep Residual Networks

- ResNet (Residual Network) is a popular deep neural network architecture.
- Introduces **skip connections** to ease training of very deep networks.
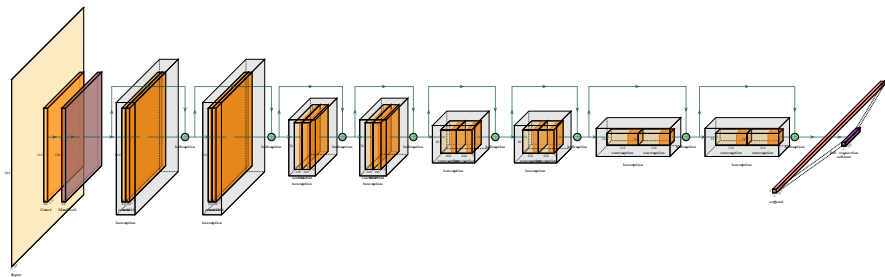- As data advances through layers, features become more refined and complex.



**Figure 1.** Structure of a ResNet with skip connections.

# Motivation for Early Exiting

- Not all data requires processing through all layers.
- **Easy data**: Can be correctly predicted with fewer layers.
- Saves computational resources by not processing all data equally.



(a) Red wine   (b) Volcano

**Figure 2.** Some inputs are easier to classify than others. [1]

Images are from the ImageNet dataset.

# Introducing Early Exit Mechanisms

- Add **early exit** points in the network with auxiliary classifiers.
- Allows the network to make predictions at intermediate layers.
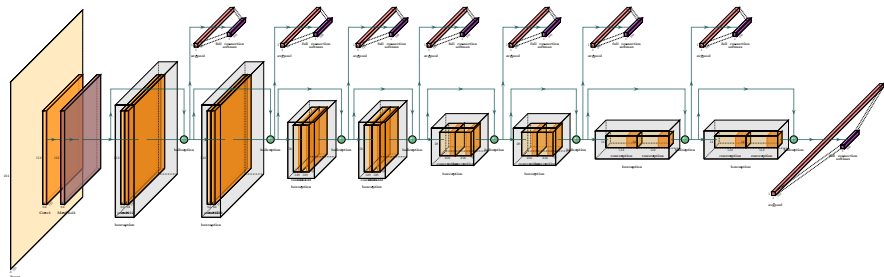- Stops computation early if the prediction is confident.



**Figure 3.** Neural network with early exit points.

# Benefits of Early Exiting

- **Power Savings**: Less computation leads to reduced energy consumption.
- **Faster Inference**: Quick predictions for easy inputs.
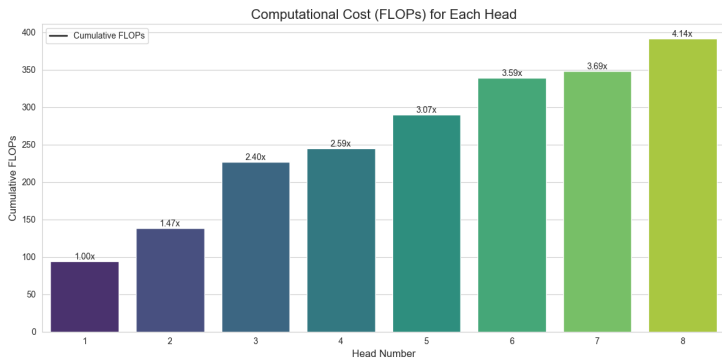- **Efficiency**: Optimizes resource usage.



**Figure 4.** Difference of computation cost per head in ResNet with early exits.

# Applying Early Exit to Large Language Models

- Large Language Models (LLMs) process sequences of words (tokens).
- **Some tokens are easier** to predict than others.
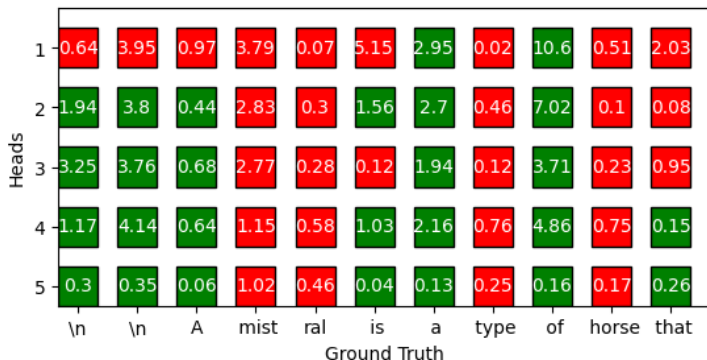- Early exit can reduce computation in language tasks by exiting on easy tokens.



**Figure 5.** A simple LLM architecture diagram. The image shows a simple architecture with Transformer Blocks which are repeated N times.

**Number of Transformer Blocks in Common LLMs**

| | | |
|---|---|---|
| GPT-2 | : | 12 to 48 blocks |
| GPT-3 | : | 96 blocks |
| LLaMA | : | 32 to 80 blocks |

# Early Exit in Language Models: Identifying Token Difficulty

- Early exit helps identify easy and hard tokens in language processing.
- Allows for efficient allocation of computational resources.
- Improves processing speed for simpler parts of the text.



**Figure 6:** Visualization of early exit identifying easy and hard tokens.

# Conclusion

- Early exit optimizes neural network processing by leveraging data difficulty.
- Saves power and improves inference speed without sacrificing accuracy.
- Applicable to various models, enhancing both performance and efficiency.

# Thank You!