# Pokémon Image Classification Challenge

Task 1

Team LLMH

# Basics

- The team
  - GALLO Lorenzo 72719
  - LEKBOURI Lina 72697
  - LICHTNER Marc 72690
  - WERCK Hugo 72692

- The work
  - Best public score: 0.93
  - Best private score: 0.93
  - Leaderboard position in the private leaderboard: 2

# Task 1: Multilayer Perceptron (MLP) Classification

# Data Exploration

- Key insights and visualizations (e.g., class distribution, sample images)

- Quantitative analysis of class imbalance
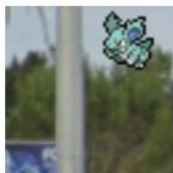
- Dataset characteristics

# Data Exploration

What can we notice?

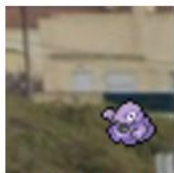- train_labels.csv : Id + labels
- labels + ".png" = images files

| | Id | label |
|---|---|---|
| 0 | 6fc9045b-9983-41e2-be2d-8796ecd97412 | Normal |
| 1 | 874716ce-9048-4e8a-b980-5ed9a5c0110e | Poison |
| 2 | c3613b20-ead8-48e1-8c8d-2f219d8e19d4 | Normal |
| 3 | c7264ebc-ba44-460a-9b2b-df23c04783bc | Normal |
| 4 | a72045db-8fae-458b-993e-23d2aab1a5c6 | Normal |

0a0d5982-d91f-4943-a900-4eb2c0c6b6bc.png
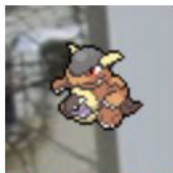
0a7f7b7f-f48a-46d4-9bb9-9dba5e85a70d.png

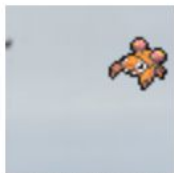0a20dde9-45bd-4145-9ca8-ea1b459b76da.png
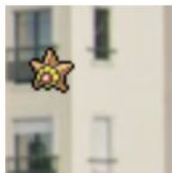
0a66c94e-72a9-490b-80c1-f3e961c77992.png

0a434ba5-fe57-4294-9125-7902b7835cc4.png

0a0494bc-c9e2-43ff-8e12-f81d4bea2cf8.png

0a630ff4-6872-47b6-8f6d-3e69a61ba042.png

0a848df0-558b-47f5-b5a5-6ead7d7970e2.png

# Data Exploration

Global descriptive statistics:

- number of images: 3600
- number of labels: 9
- duplicate ones: 0
- labels and their counts:
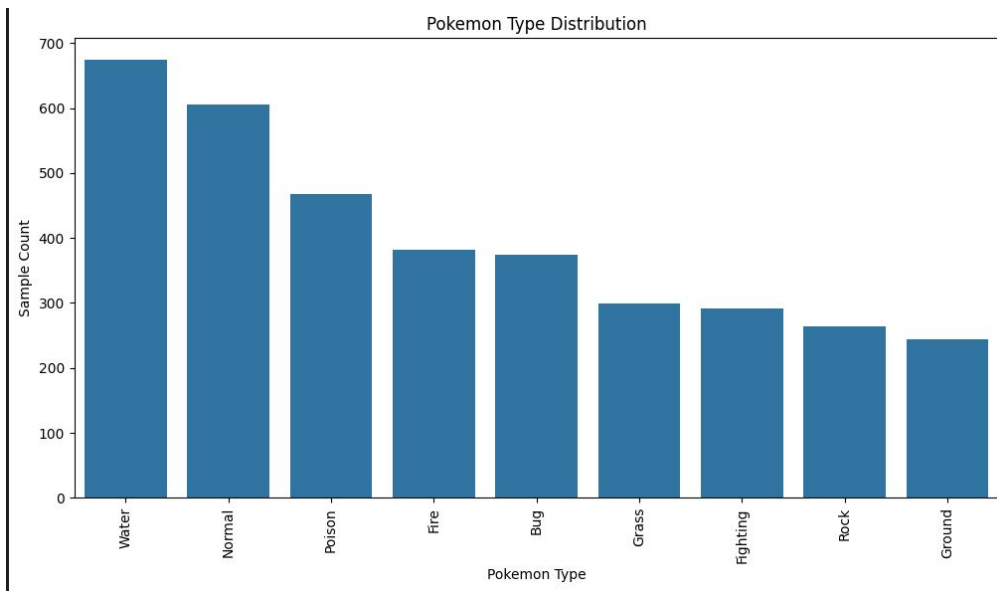
```
Labels and their counts:

label
Water       674
Normal      606
Poison      467
Fire        381
Bug         374
Grass       299
Fighting    291
Rock        264
Ground      244
Name: count, dtype: int64
```

# Data Exploration

- Class distribution visualization

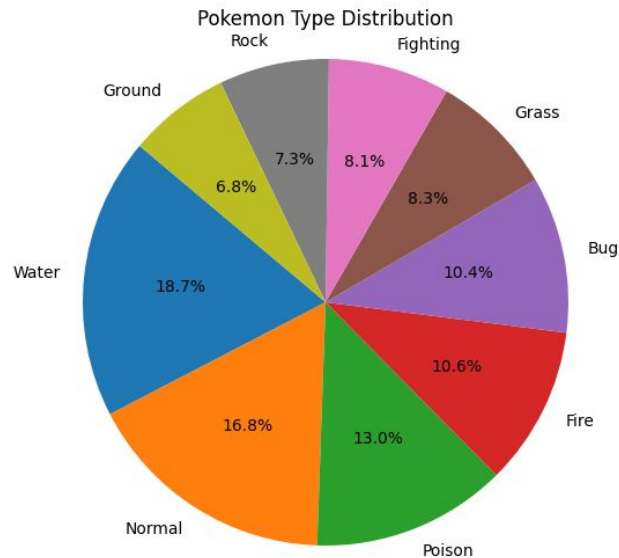  <u>an histogram?</u>

# Data Exploration

- Class distribution visualization

    →a <u>pie chart</u> with percentages

    Best way to see any **imbalance**



Pokemon Type Distribution

# Data Exploration

Images analysis

- size: **64*64 pixels**
- background? It has nothing to do with the pokémon

# Data Exploration

- Overview of some primary types



Primary type: Fire    Primary type: Grass    Primary type: Poison    Primary type: Fire    Primary type: Fighting

# Data Exploration

• How do we balance data ?

  Possible methods:

- undersampling: Decrease the number of instances in the majority class.
- oversampling: Increase the number of instances in the minority class.

→ Choice : **undersampling**, by using the value of the less frequent class.

# Disclaimer

It is the day before the submission of the first task and we achieved to center the Pokemon on the picture. As a result, the accuracy of the model exploded and went from **around 0.3** on validation test to **0.93.**

The result you will see below are the result before the centering.

# Data Exploration

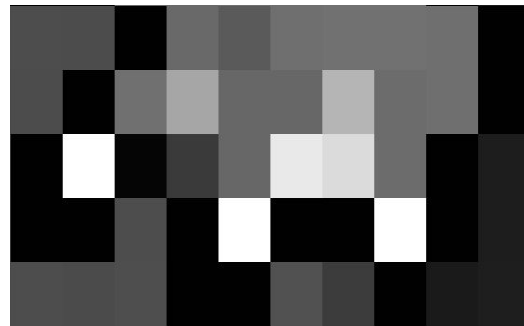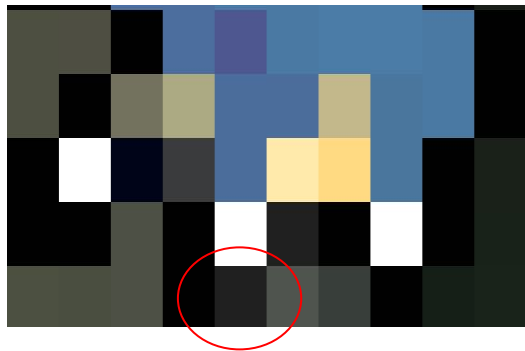- Data preprocessing : removing the background

# Remove the background

- Few examples

# Remove the background

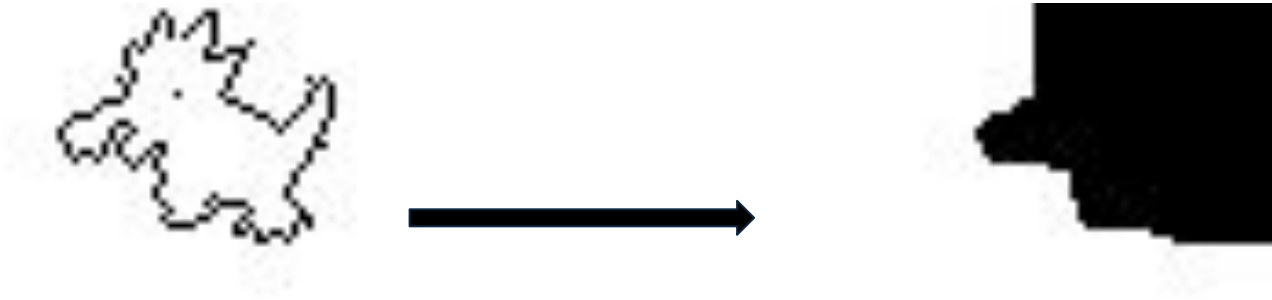**First step**: fix_border: Make every pixels of the border the exact same value

# Remove the background

**Second step**: apply_threshold: Make the border black and everything else white

# Remove the background

**Third step**: erode_and_dilate: Fill the inside of the border using the closing operation which creates the mask as a result

# Remove the background

**Last step**: apply the mask on the original image



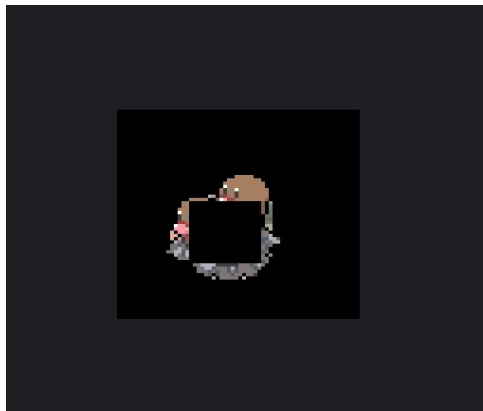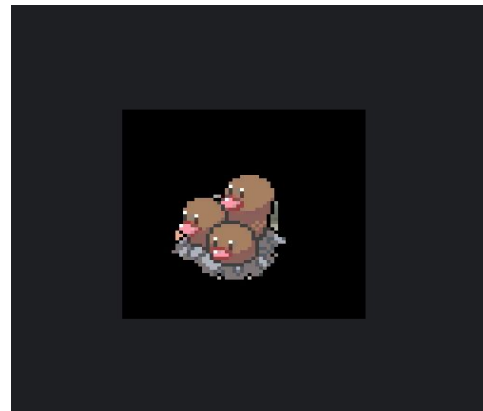AND                                    =

# Remove the background

**Remark**: This filter keep parts of the background because of the number of time we iterate using the closing algorithm. We had to make sure that data from the pokemon don't get remove
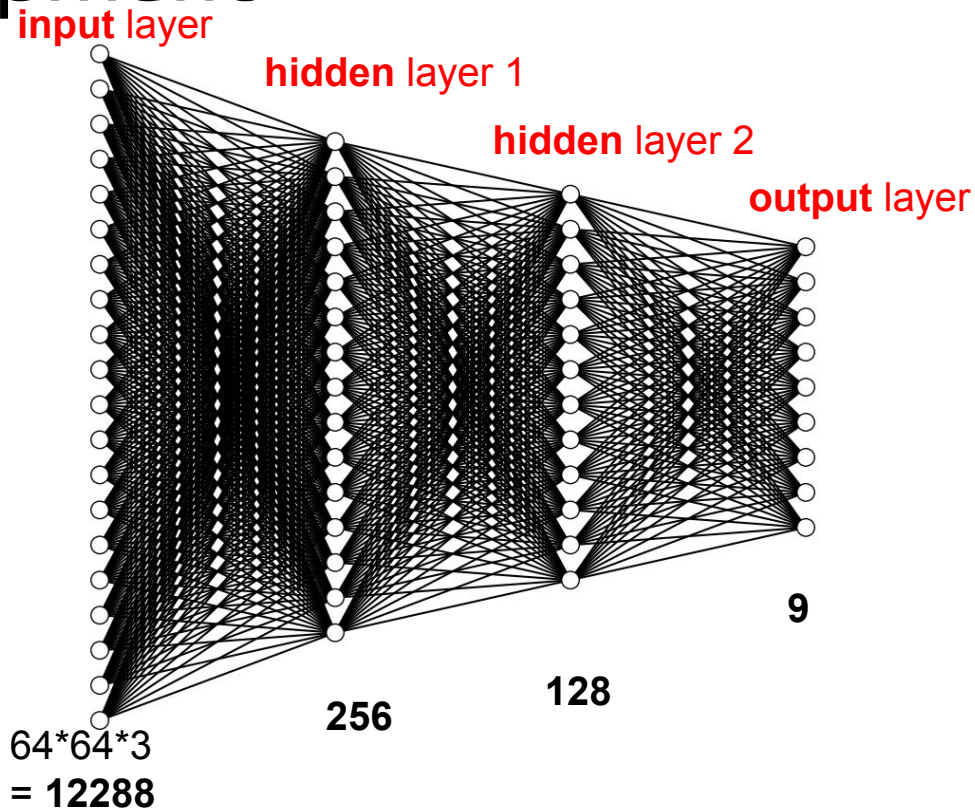


9 iterations



10 iterations

# Model Development

- MLP architecture diagram
- Justification for chosen architecture (number of layers, neurons, activations)
- Ablation study

# Model Development

- MLP architecture diagram:

(Given the big numbers of neurons, we can not represent the real numbers on this slide, so we take random numbers of neurons respecting orders of magnitude.)

**input** layer

**hidden** layer 1

**hidden** layer 2

**output** layer

9

128

256

64*64*3
= **12288**

number of neurons:

# Model Development

- Justification for chosen architecture

**4 layers:**

- input layer: **64*64*3 = 12 288 neurons** (number of pixels in each image by the RGB channel)
- hidden layer 1: 256 neurons
- hidden layer 2: 128 neurons
- output layer: **9 neurons** (number of primary types of pokémon)

# Model Development

- Justification for chosen architecture: trying other architectures.
  ⇒ justification for the **hidden layers**.

  To justify our architecture, we trained our model with **other architectures**, printed the **result on the validation set,** and concluded that the **best results** were given by our chosen one.
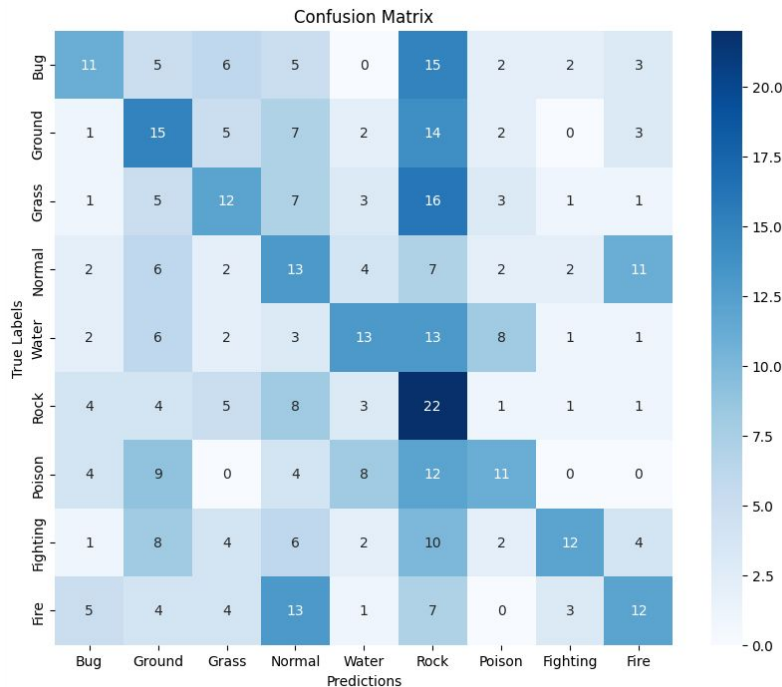
# Model Development

- Justification of the architecture of **hidden layers**.

Configuration 1:
- hidden layer 1: **512**
- hidden layer 2: **256**

→ **a class is more predicted than the other.**



Confusion Matrix

# Model Development

- Justification of the architecture of **hidden layers**.

Configuration 2:
- hidden layer 1: **1024**
- hidden layer 2: **512**
- hidden layer 3: **256**

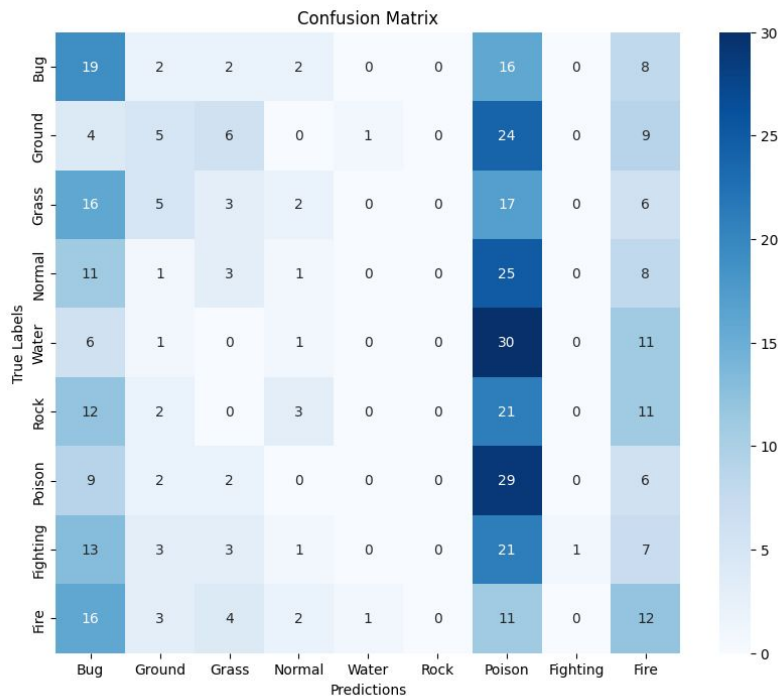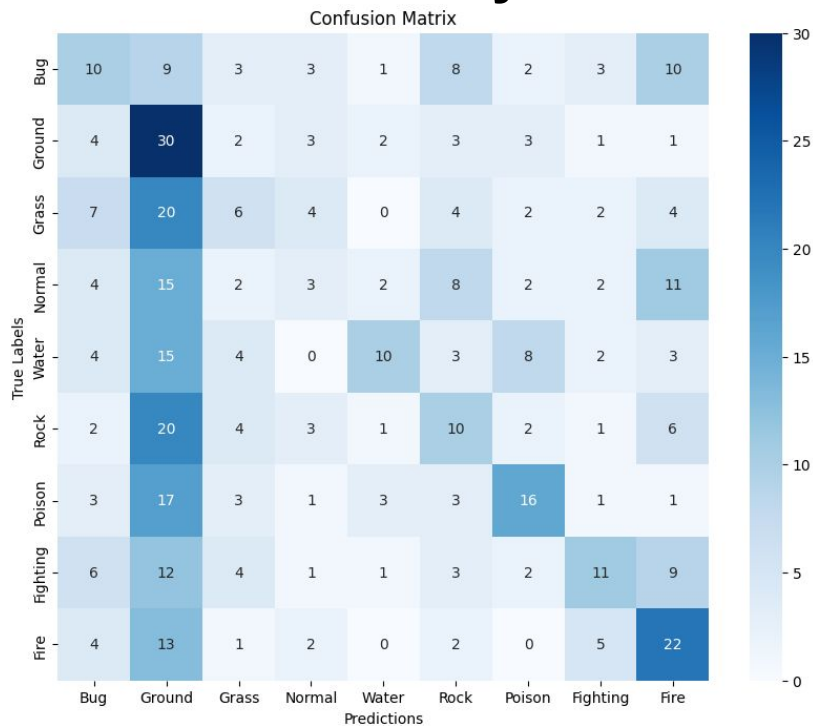→ **2 classes are way more predicted than the other.**



Confusion Matrix

# Model Development

- Justification of the architecture of **hidden layers**.

Configuration 3:
- hidden layer 1: **1024**
- hidden layer 2: **512**

→**a class is more predicted than the other.**
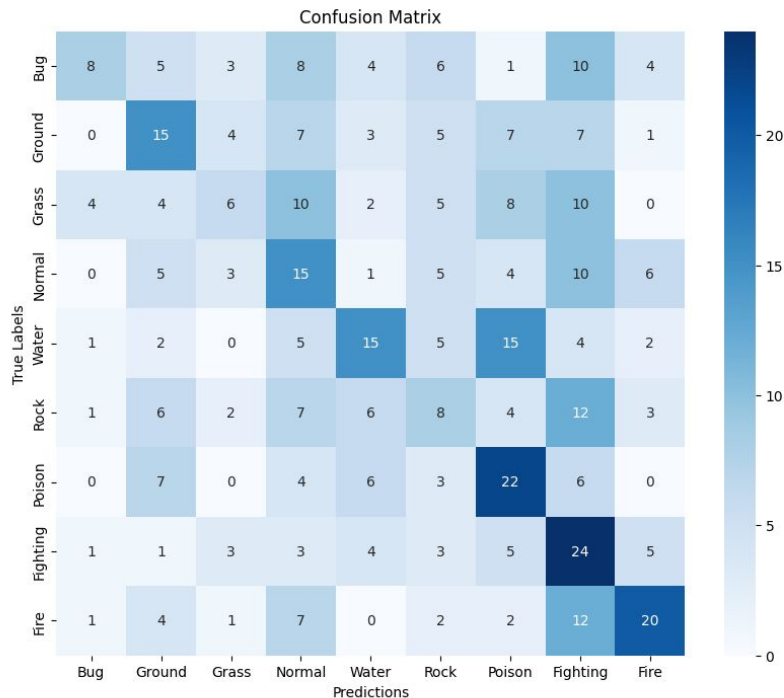


Confusion Matrix

# Model Development

- Justification of the architecture of **hidden layers**.

Configuration 4:
- hidden layer 1: **256**
- hidden layer 2: **128**
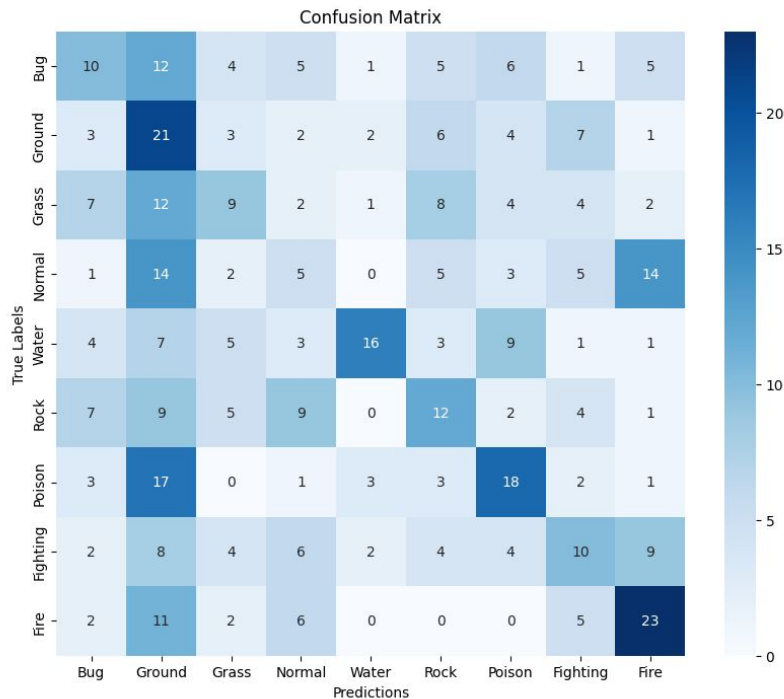- hidden layer 3: **64**

→ **less good on the test data.**



Confusion Matrix

# Model Development

- Justification of the architecture of **hidden layers**.

Final configuration:
- hidden layer 1: **256**
- hidden layer 2: **128**

**→ the one we choose, because it is also better on the test data.**



Confusion Matrix

# Model Development

- That is why we choose this hidden architecture.

**4 layers:**

- input layer: **64*64*3 neurons** (number of pixels in each image by the RGB channel
- hidden layer 1: **256** neurons
- hidden layer 2: **128** neurons
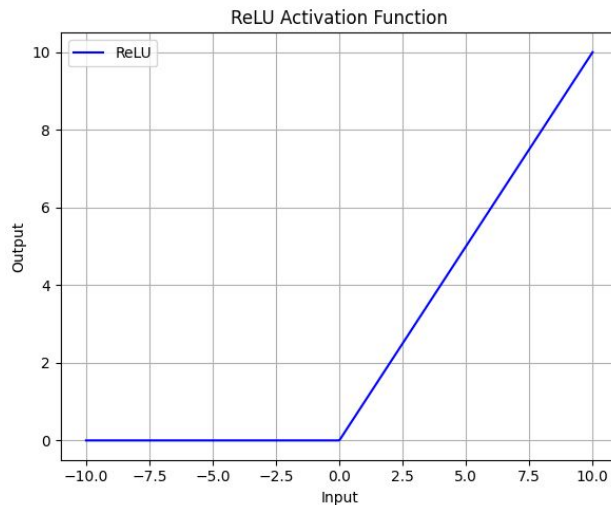- output layer: **9 neurons** (number of primary types of pokémon)

# Model Development

- Justification for activations

**Activation functions:**

- **ReLU** (from input to h1 and from h1 to h2): $\text{ReLU}(x) = \max(0, x)$
- **Softmax** (from h2 to output): $\text{Softmax}(x_i) = \dfrac{e^{x_i}}{\sum_j e^{x_j}}$
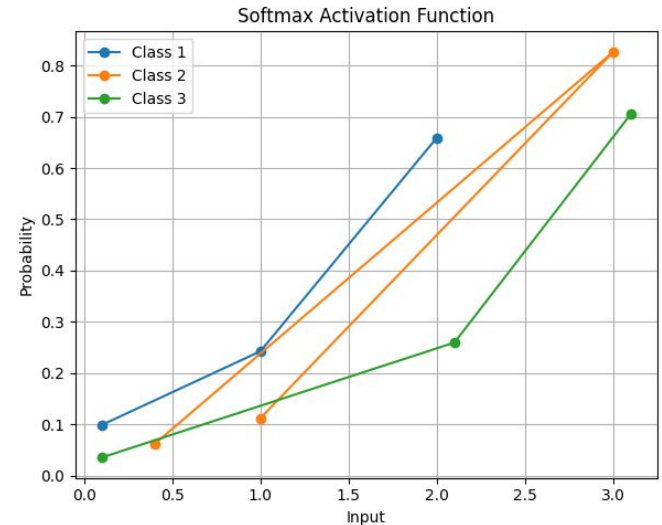
# Model Development



ReLU Activation Function

- Justification for activations

  - **ReLU activation function, why?**

    - Efficient and computationally cheap;

    - Avoids vanishing gradient issues common in other activations;

    - Improve generalization.

# Model Development

- Justification for activations

  - **Softmax activation function, why?**

    - Converts raw scores into probabilities, making the model interpretable;

    - Ensures sum of outputs equals 1.



Softmax Activation Function

# Model Development

- Ablation study: understanding the impact of every component in the performance

Experiments Conducted:

- **removing ReLU** → Model struggles to learn complex patterns, lower accuracy;
- **replacing ReLU with Sigmoid** → Slower convergence, possible vanishing gradient issues;
- **removing Softmax** → Outputs unbounded raw scores instead of interpretable probabilities;
- **reducing hidden layers** →Model underfits, struggles with feature extraction.

# Model Development

- Ablation study: conclusions on the experiments conducted

  - ReLU significantly improves **training efficiency and performance**;
  - Softmax is essential for **classification tasks**;
  - Deeper architectures capture **more complex representations**.

# Training Efficiency

- GPU usage

- Strategies employed for efficient training (early stopping, batch sizes, etc.)
- Training time

# Training Efficiency

**GPU usage**:

By using the GPU, it took 1min59 to compute 200 epochs whereas it took 4min54 with the CPU. This is what we used to take advantage of the GPU if available:
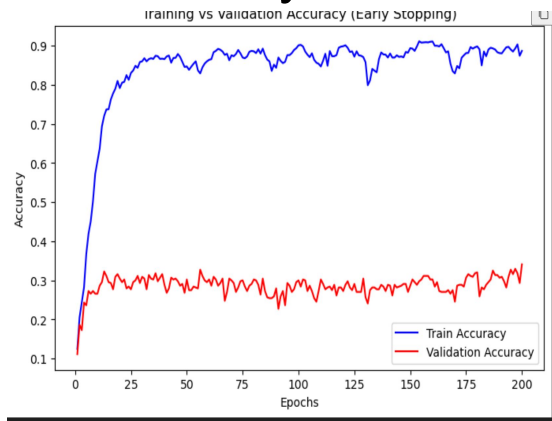
```python
# Check if GPU is available
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print(f"Using device: {device}")
```
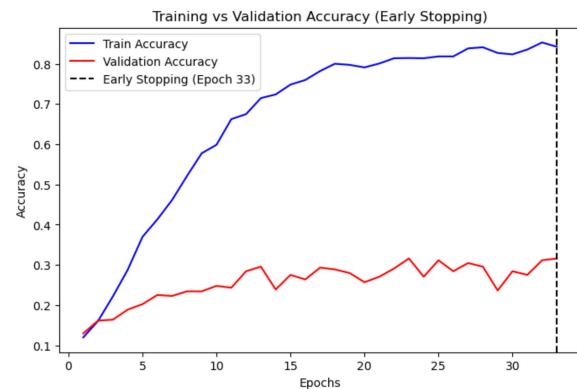
# Training Efficiency

- Strategies employed for efficient training (early stopping, batch sizes, etc.)

  - **Early stopping**: Avoid useless computations by stopping the training when the model don't improve anymore

  - **Batch size**: In order to select the batch sizes, we followed the "linear scaling rule". We chose a batch size of 32 and a learning rate of 0.001.

# Early stopping

We used early stopping because the model was quickly converging to a maximum around 25 epochs and won't improve it's accuracy on the validation dataset anymore



Without early stopping
accuracy of 0.3087

With early stopping
accuracy of 0.2977
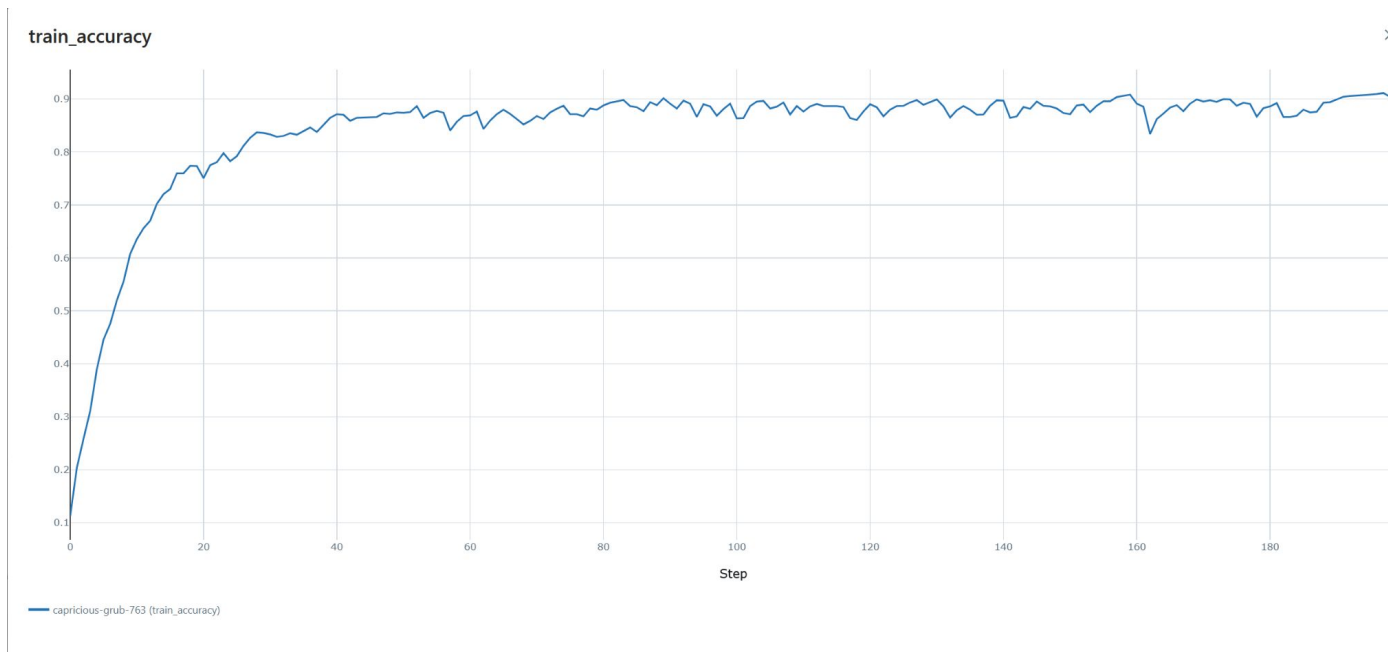
# Training Efficiency

Training time

By using early stopping and the gpu, the training took **24 seconds.**

# Performance Evaluation

- Justification of the metrics used (macro F1 score, accuracy, confusion matrix)
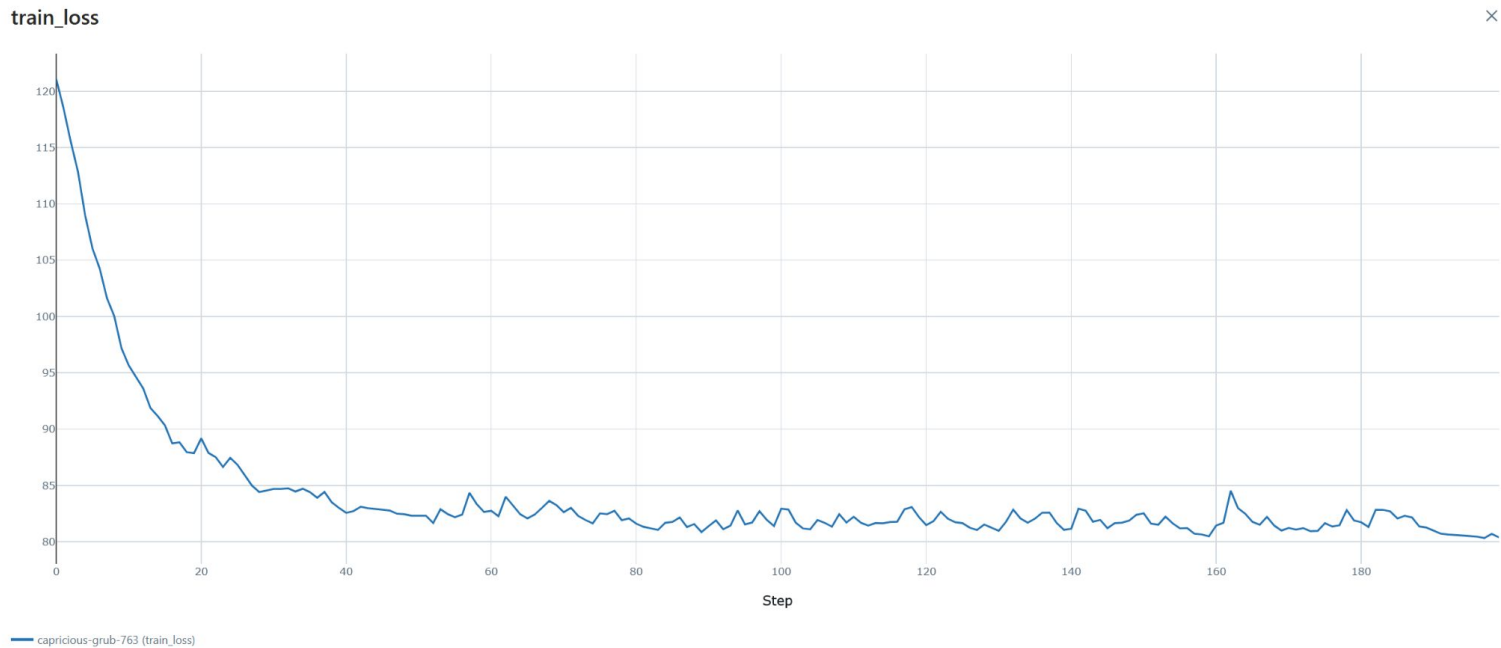- Interpretation of results

# Performance Evaluation

• Train accuracy

   • Measures the percentage of correctly classified training samples, indicating how well the model fits the training data.
   • The curve shows an increasing trend, reaching over 80%. This suggests the model is learning effectively from the training data.



train_accuracy

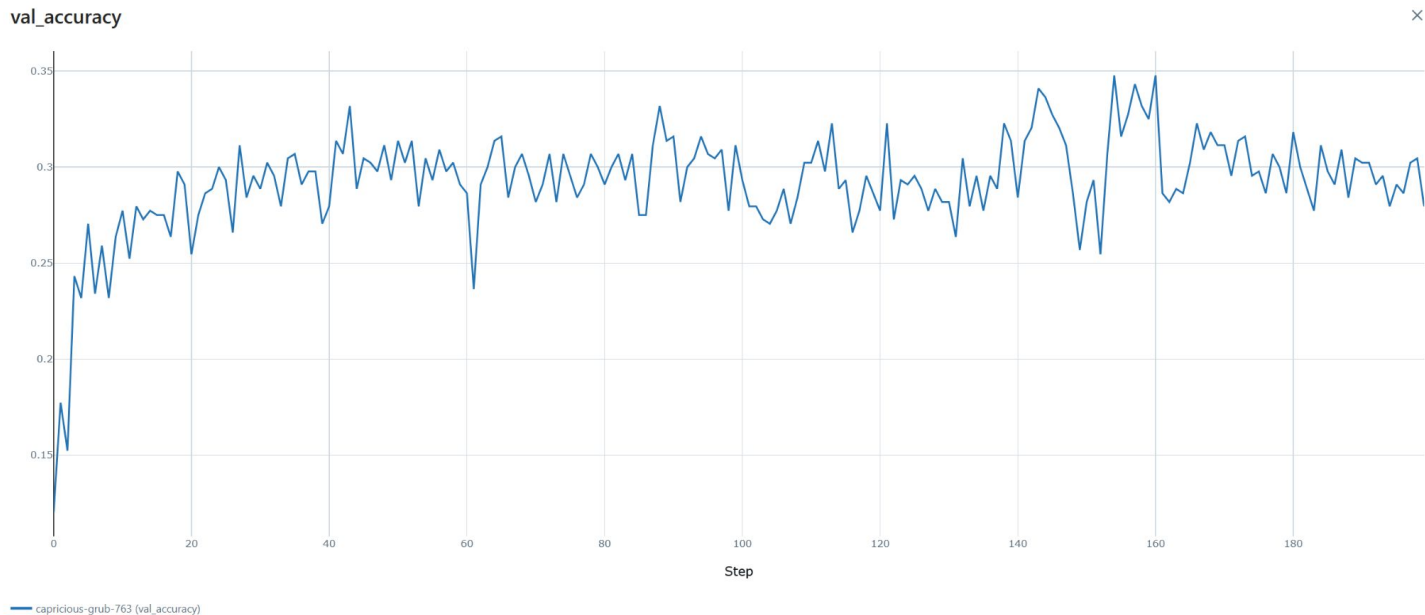capricious-grub-763 (train_accuracy)

# Performance Evaluation

- Train loss
    - Represents the error during training, helping to track convergence and detect overfitting.
    - The loss decreases rapidly, indicating successful learning.

**train_loss** ✕



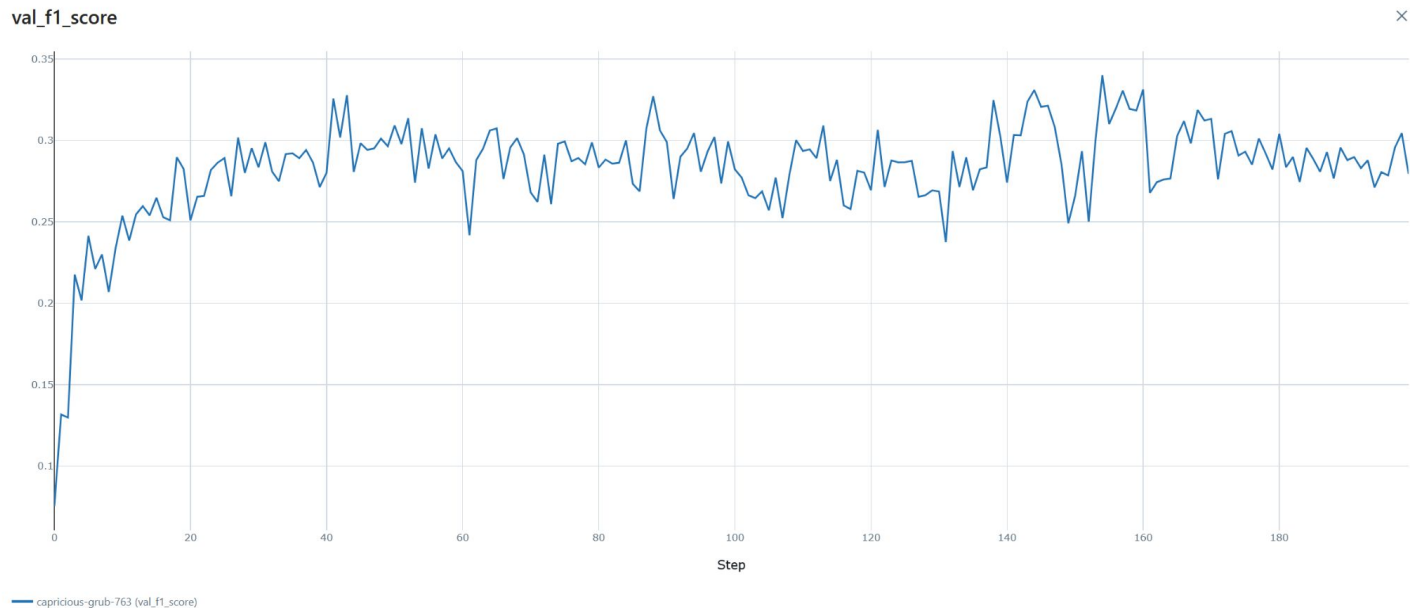capricious-grub-763 (train_loss)

# Performance Evaluation

- Validation accuracy
    - Evaluates performance on unseen data, reflecting generalization ability.
    - The curve fluctuates significantly around 30%, which suggests the model struggles to generalize.

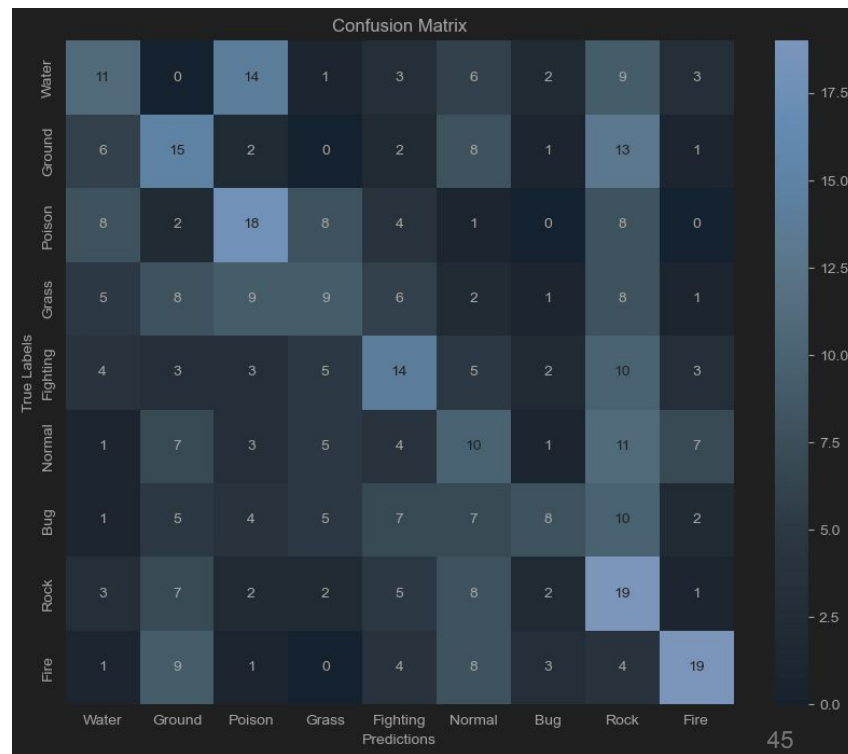# Performance Evaluation

- Macro F1 score

  - Balances precision and recall, crucial for handling class imbalances.

  - The high variance suggests unstable predictions across classes.

# Performance Evaluation

- Confusion matrix

  - provides a detailed evaluation of the model's performance beyond a single accuracy score.

  - Strong performance on Poison (18), Fighting (14), Rock (19), and Fire (19), indicating the model learns distinct features for these types.

  - Frequent misclassifications:

    - Water → Poison (14x) and Ground → Rock (13x) suggest visual similarities.

    - Grass → Poison and Fighting ↔ Bug/Rock indicate shared features.

    - Normal is scattered, showing weak representation.

  - Improvements: Better feature extraction, class balancing, and data augmentation could enhance performance.
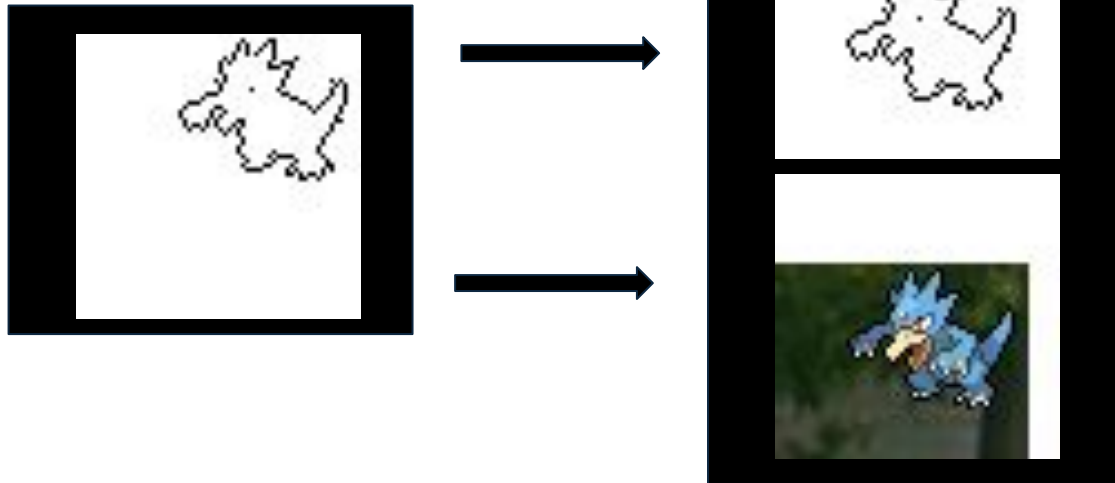
# Breakthrough: centering the Pokemons

We achieve to **center** the Pokemon in addition to removing the background.
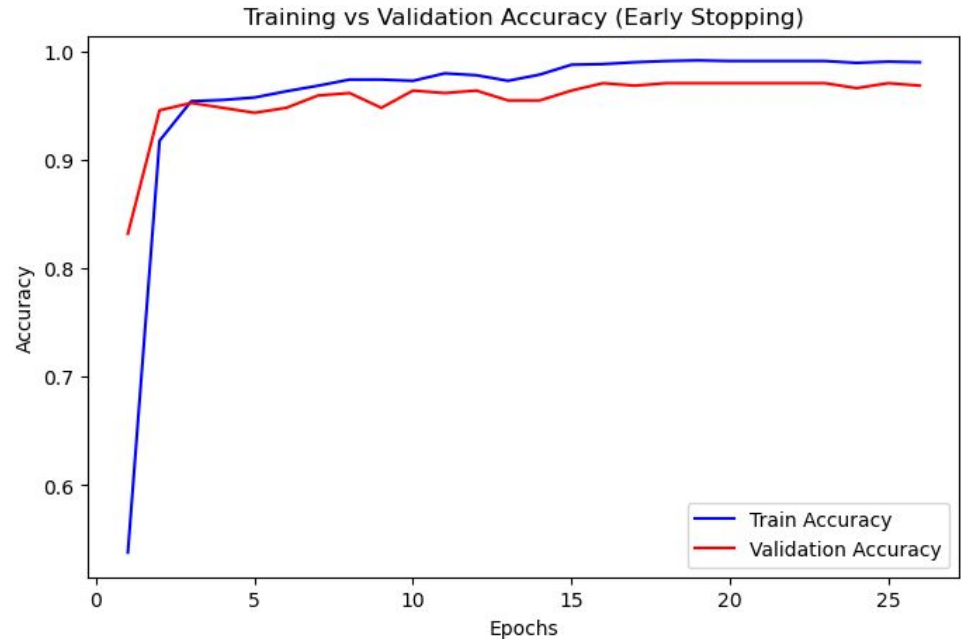
# Centering the Pokemon

We used the thresholded image to locate and get the dimension of the Pokemon easily. Basically we "draw" a rectangle around the Pokemon.With this we know exactly what part we wanted to move on the picture. Once we have this, the process is the same to remove the background
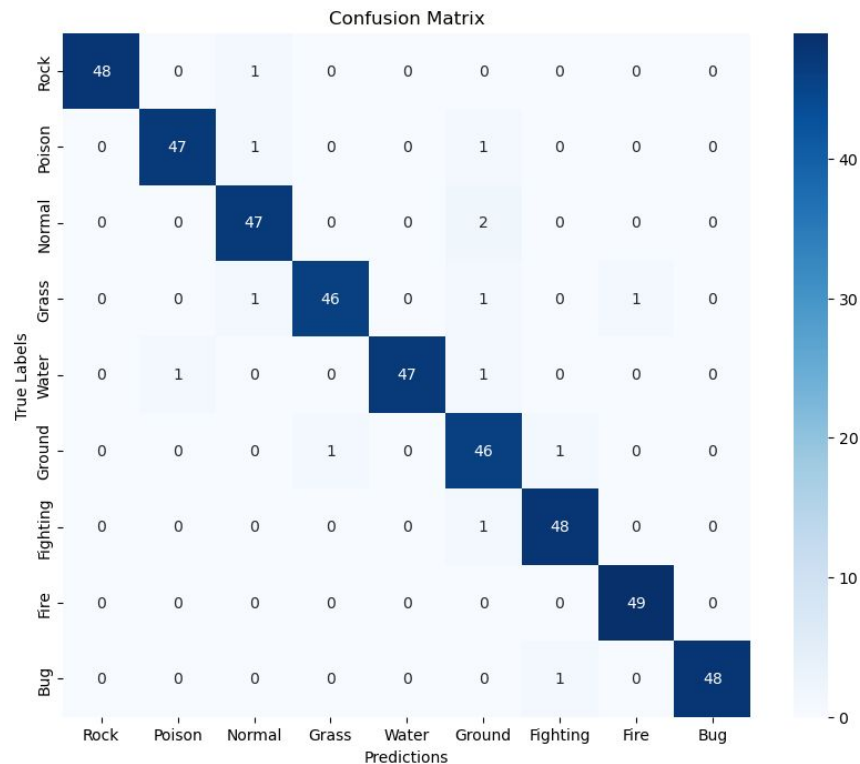
# Training efficiency

As the validation accuracy converged even more aggressively, the training stopped even earlier than before.

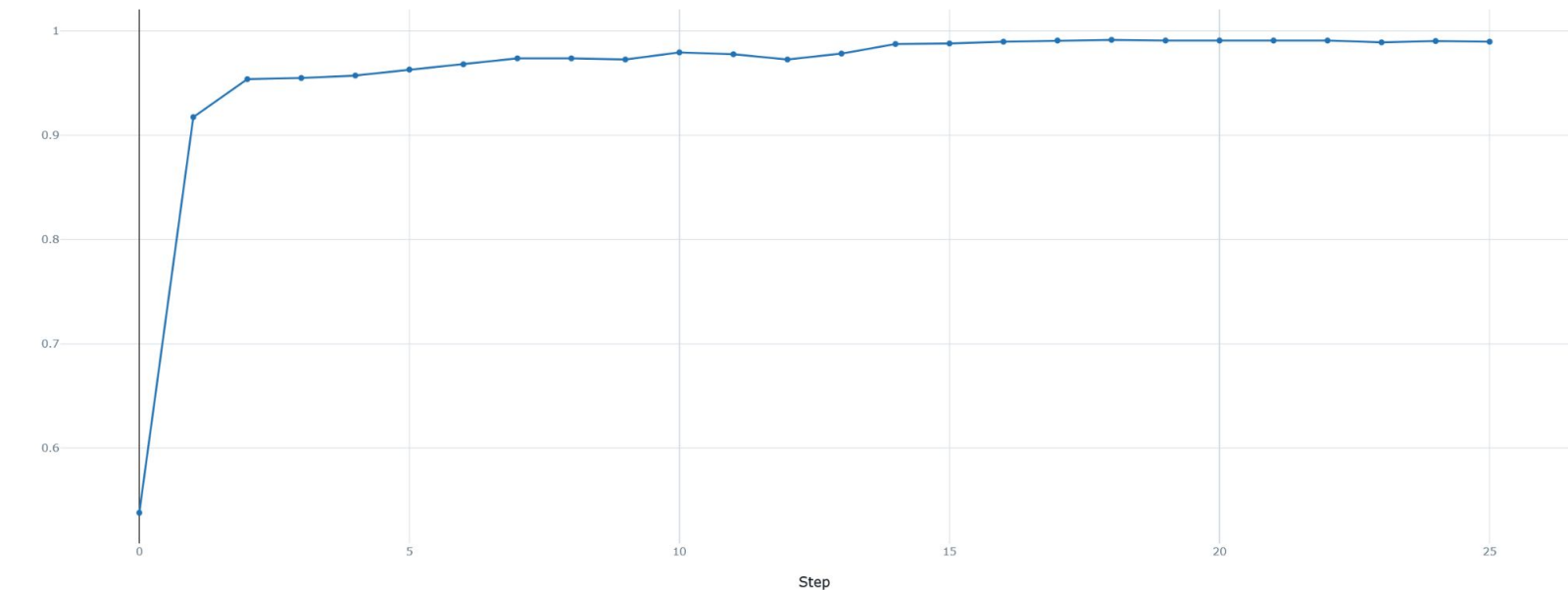The training took only 26 epochs which represent 18sec using gpu and early stopping



Training vs Validation Accuracy (Early Stopping)

# Performance evaluation

The new **confusion matrix**:



Confusion Matrix

# Performance Evaluation



train_accuracy

luxuriant-mouse-243 (train_accuracy)

# Performance Evaluation



train_loss
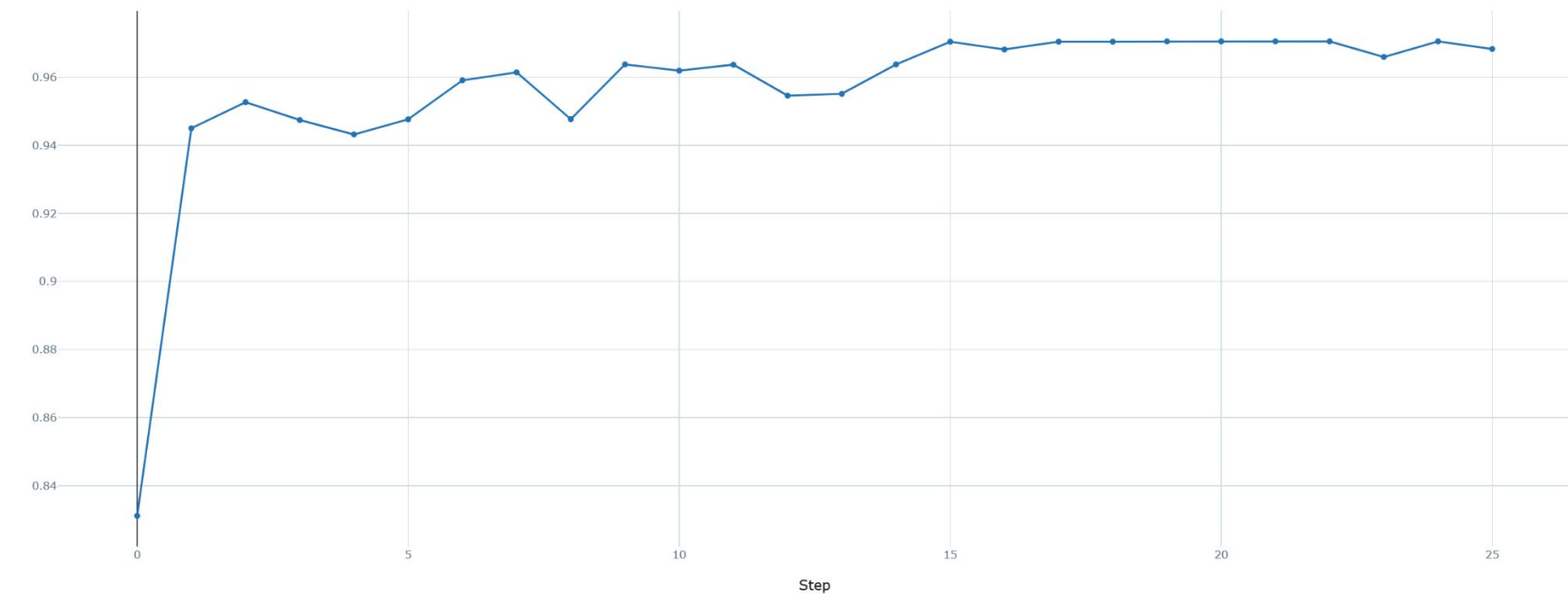
luxuriant-mouse-243 (train_loss)

# Performance Evaluation



val_accuracy

luxuriant-mouse-243 (val_accuracy)

# Performance Evaluation

# Lessons learned and insights gained

- Even though by centering the Pokemon, the MLP was actually pretty accurate, we can still do better using a **CNN** model.
- **Preprocessing** correctly the data is the key to get **good results**.

# What went wrong

- Installing MLFlow locally didn't work, so we had to do a free trial on Databricks.
- It takes a long time for MLFlow to work in all environments.

# What went great

- Data exploration;
- Teamwork;
- Comprehension of how MLP work;
- Image preprocessing.