# Milestone 1
# Team 12

| Name | ID | Sec |
|------|-----|-----|
| **Adham Ehab Saleh** | 2100679 | 1 |
| **Ahmed Mohamed Ramadan** | 2100323 | 1 |
| **Ahmed Salah Eldin Abdelrahman** | 2100505 | 1 |
| **Ahmed Yasser Hosney** | 2101101 | 1 |
| **Omar Magdy Abdelsattar** | 2100273 | 2 |

**Delivered to:** Dr. Hossam El-Din Hassan

Eng. Abdallah Mohamed Mahmoud

Eng. Dina Zakaria Mahmaud

# Table of Contents

# Table of figures

# Problem Definition (Milestone 1: Library Validation & XOR)

## Background and Context

The fundamental objective of this project is to develop a foundational neural network library from scratch using only Python and NumPy. Before applying this library to complex tasks like image reconstruction, it is essential to validate the correctness of the core mathematical algorithms, specifically forward propagation, backward propagation, and parameter optimization. To achieve this, the library must be tested against a benchmark problem that requires non-linear decision boundaries.

## Problem Statement

The specific challenge for Part 1 is to construct and train a Multi-Layer Perceptron (MLP) capable of solving the XOR (Exclusive OR) problem. The XOR function is a classic binary classification problem where the target output is true (1) if and only if the inputs differ, and false (0) otherwise.

Unlike logical AND or OR gates, the XOR function is **not linearly separable**, meaning a single linear decision boundary (like a perceptron) cannot solve it. This makes it an ideal candidate to verify the library's ability to learn non-linear representations using hidden layers and activation functions.

The dataset consists of the four-standard truth-table inputs:

- **Inputs:** [(0,0), (0,1), (1,0), (1,1)]

- **Targets:** [0, 1, 1, 0]

## Project Objectives (Milestone 1)

The primary goal of this phase is to demonstrate that the custom-built library functions correctly by converging on a solution for the XOR problem.

Specific objectives include:

- **Library Architecture:** Implementing modular classes for Dense layers, activations (ReLU, Sigmoid, Tanh), and Loss functions (MSE).

- **Gradient Verification:** Validating the backpropagation algorithm by performing **numerical gradient checking** to ensure the analytical gradients match the numerical approximations.

- **Model Implementation:** Constructing a specific MLP architecture (e.g., 2-4-1) using Tanh and Sigmoid activations.

- **Optimization:** successfully minimizing the Mean Squared Error (MSE) using Stochastic Gradient Descent (SGD) to achieve correct classification for all 4 input cases.

# Methodology and Library Design

## System Architecture

To meet the requirement of a modular, extensible framework, the library was designed using an Object-Oriented Programming (OOP) approach. The core architecture relies on an abstract base class structure where specific components (Layers, Activations, Optimizers) inherit common interfaces. This ensures that the network model can treat all components uniformly during the forward and backward passes.

The library structure consists of the following key modules:

- **Layer (Base Class):** Defines the abstract methods forward() and backward() that all derived classes must implement.

- **Dense (Fully Connected Layer):** Inherits from Layer. It initializes and stores the trainable parameters—Weights ($W$) and Biases ($b$)—and computes the linear transformation $Z = X.W + b$. It also stores the gradients ($\frac{\partial L}{\partial W}, \frac{\partial L}{\partial b}$) computed during backpropagation.

- **Network (Container):** Acts as the orchestrator. It maintains a list of layers, manages the flow of data through the network, and executes the training loop.

# Mathematical Implementation

## Activation Functions

Activation functions are implemented as subclasses of the Layer class to maintain modularity. Each activation class implements:

1. **Forward:** Applies the non-linear function $f(x)$ (e.g., Sigmoid, Tanh, ReLU).

2. **Backward:** Computes the derivative $f'(x)$ and applies the chain rule to propagate the error gradient to the previous layer.

## Loss Function

The training performance is measured using the Mean Squared Error (MSE), as required for regression and binary classification tasks in this project.

The loss is calculated as:

$$L = \frac{1}{N}\sum(Y_{true} - Y_{pred})^2$$

The derivative with respect to the prediction $Y_{pred}$, which initializes the backpropagation process, is implemented as $2 \times (\frac{Y_{pred} - Y_{true}}{N})$

## Optimization (SGD)

Parameter updates are handled by the Stochastic Gradient Descent (SGD) optimizer. The optimizer iterates through every trainable layer in the network and updates the weights and biases using the gradients computed during the backward pass:

$$W_{new} = W_{old} - \eta \frac{\partial L}{\partial W}$$

where $\eta$ represents the learning rate.

# Validation Strategy

## Numerical Gradient Checking

Before training, the mathematical correctness of the backpropagation implementation was verified using Gradient Checking. This technique compares the analytical gradient (computed by the chain rule in the code) against a numerical approximation derived from the limit definition of a derivative.

The numerical gradient was approximated using the central difference formula:

$$\frac{\partial L}{\partial W} \approx \frac{L(W + \epsilon) - L(W - \epsilon)}{2\epsilon}$$

A small $\epsilon$ (e.g., $10^{-7}$) was used to perturb the weights. The relative difference between the analytical and numerical gradients was checked to ensure it fell below a strict threshold (e.g., $10^{-7}$), confirming the backpropagation logic is correct.

## XOR Benchmark Architecture

To validate the library's ability to learn non-linear boundaries, the XOR problem was selected as the primary unit test.

The specific architecture constructed for this test is a Multi-Layer Perceptron (MLP) with the following topology:

- **Input Layer:** 2 Neurons (representing the two binary inputs).

- **Hidden Layer:** 4 Neurons with **Tanh** activation.

- **Output Layer:** 1 Neuron with **Sigmoid** activation (outputting a probability between 0 and 1).

This "2-4-1" architecture provides sufficient capacity to map the non-linearly separable XOR inputs to the correct Boolean outputs.

# Results

First, we needed to validate the gradient logic, so in the notebook there is a section to validate it, and its output is as follows

```
Generating random data for testing...
--- Gradient Check for Dense Layer ---
    Analytical Gradient Norm: 0.958773
    Numerical Gradient Norm:  0.958773
    Relative Error:           1.62e-13
  ✅ PASS: Backpropagation is correct.
--------------------------------------------------
```

*Figure 1.Gradient Logic Validation.*

Then, we trained the network for 10000 epochs

```
Epoch 1000/10000    error=0.002678
Epoch 2000/10000    error=0.001048
Epoch 3000/10000    error=0.000640
Epoch 4000/10000    error=0.000458
Epoch 5000/10000    error=0.000356
Epoch 6000/10000    error=0.000290
Epoch 7000/10000    error=0.000245
Epoch 8000/10000    error=0.000212
Epoch 9000/10000    error=0.000186
Epoch 10000/10000    error=0.000166
```

*Figure 2.Training of The Network.*

Then, we tested the network

```
XOR Predictions after training:
Input: [0 0]  Predicted: 0.0051  True: 0
Input: [0 1]  Predicted: 0.9864  True: 1
Input: [1 0]  Predicted: 0.9864  True: 1
Input: [1 1]  Predicted: 0.0164  True: 0
```

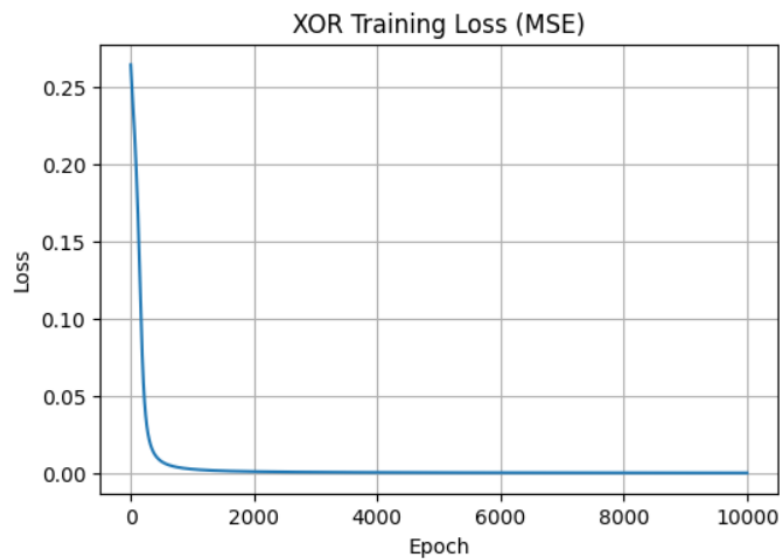*Figure 3.Prediction After Training.*



*Figure 4.Loss Visualization for Model.*
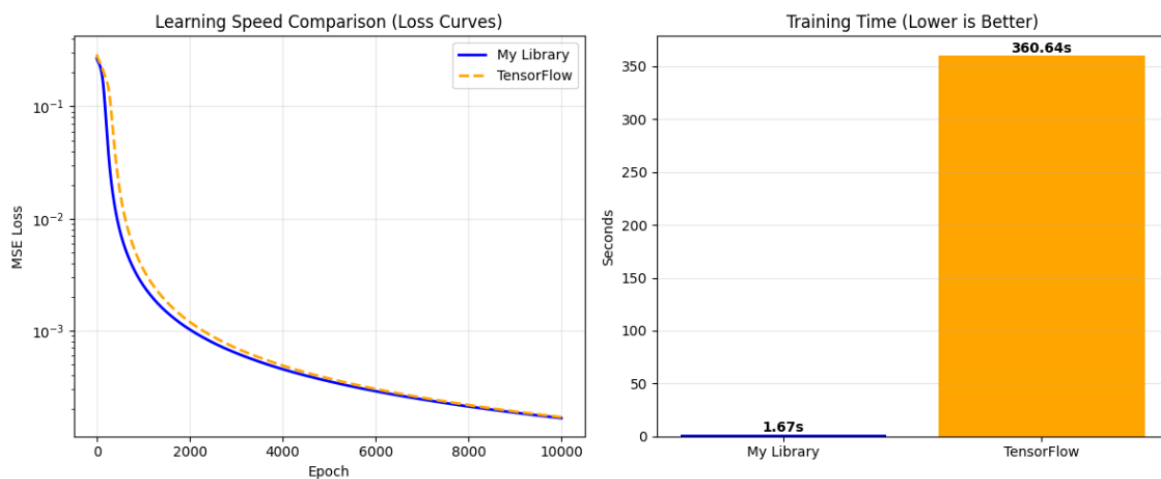
We then compared this custom library with TensorFlow



*Figure 5.Comparison between Custom Model and TensorFlow.*

# Appendices

## Appendix A: activations code

```python
# -*- coding: utf-8 -*-

"""activations.py

Activation layers for neural networks.

Each activation function is implemented as a Layer with forward and backward
methods.
"""

import numpy as np

from .layers import Layer

class ReLU(Layer):
    """

    Rectified Linear Unit activation function.

    f(x) = max(0, x)

    f'(x) = 1 if x > 0, else 0
    """

    def __init__(self):
        super().__init__()

    def forward(self, input_data):
        """

        Apply ReLU activation: output = max(0, input)
        """

        self.input = input_data

        self.output = np.maximum(0, input_data)

        return self.output
```

```python
    def backward(self, output_error, learning_rate):
        """

        Backpropagate through ReLU.

        Gradient is 1 where input > 0, otherwise 0.

        """

        # Element-wise multiplication of output_error with derivative

        input_error = output_error * (self.input > 0)

        return input_error

class Sigmoid(Layer):

    """

    Sigmoid activation function.

    f(x) = 1 / (1 + exp(-x))

    f'(x) = f(x) * (1 - f(x))

    """

    def __init__(self):

        super().__init__()

    def forward(self, input_data):

        """

        Apply sigmoid activation: output = 1 / (1 + exp(-input))

        """

        self.input = input_data

        self.output = 1 / (1 + np.exp(-input_data))

        return self.output
```

```python
    def backward(self, output_error, learning_rate):
        """

        Backpropagate through Sigmoid.

        Gradient is sigmoid(x) * (1 - sigmoid(x))

        """

        sigmoid_derivative = self.output * (1 - self.output)

        input_error = output_error * sigmoid_derivative

        return input_error

class Tanh(Layer):
    """

    Hyperbolic tangent activation function.

    f(x) = tanh(x) = (exp(x) - exp(-x)) / (exp(x) + exp(-x))

    f'(x) = 1 - tanh^2(x)

    """

    def __init__(self):

        super().__init__()

    def forward(self, input_data):
        """

        Apply tanh activation: output = tanh(input)

        """

        self.input = input_data

        self.output = np.tanh(input_data)

        return self.output
```

```python
    def backward(self, output_error, learning_rate):
        """

        Backpropagate through Tanh.

        Gradient is 1 - tanh^2(x)

        """

        tanh_derivative = 1 - self.output ** 2

        input_error = output_error * tanh_derivative

        return input_error

class Softmax(Layer):
    """

    Softmax activation function (typically used for multi-class
classification).

    f(x_i) = exp(x_i) / sum(exp(x_j) for all j)

    Note: Softmax backward pass is complex when used with cross-entropy loss.

    This implementation handles the general case.

    """

    def __init__(self):
        super().__init__()

    def forward(self, input_data):
        """

        Apply softmax activation: output_i = exp(x_i) / sum(exp(x_j))

        Uses numerical stability trick: subtract max before exp

        """

        self.input = input_data

        # Numerical stability: subtract max value

        exp_values = np.exp(input_data - np.max(input_data))

        self.output = exp_values / np.sum(exp_values)

        return self.output
```

```python
def backward(self, output_error, learning_rate):
    """
    Backpropagate through Softmax.
    For a single sample, the Jacobian is:
    J_ij = s_i * (δ_ij - s_j) where s is softmax output, δ is Kronecker
    delta
    """
    n = self.output.shape[0]
    # Compute Jacobian matrix of softmax
    # J[i,j] = output[i] * (1 if i==j else 0) - output[i] * output[j]
    jacobian = np.diagflat(self.output) - np.outer(self.output,
    self.output)
    # Multiply Jacobian with output_error
    input_error = np.dot(jacobian, output_error)
    return input_error
```

## Appendix B: layers code

```python
# -*- coding: utf-8 -*-
"""layers.ipynb

Automatically generated by Colab.

Original file is located at

    https://colab.research.google.com/drive/1WpydwzFuapT959b5hGNqf3LVipBNLL6D
"""

import numpy as np

class Layer:

    def __init__(self):

        self.input = None

        self.output = None

    def forward(self, input_data):

        raise NotImplementedError

    def backward(self, output_error, learning_rate):

        raise NotImplementedError

class Dense(Layer):

    def __init__(self, input_size, output_size):

        super().__init__()

        limit = np.sqrt(6.0 / (input_size + output_size))

        self.weights = np.random.uniform(-limit, limit, size=(output_size,

        input_size))

        self.biases = np.zeros((output_size,))

        self.input = None

        self.grad_weights = None

        self.grad_biases = None
```

```python
def forward(self, input_data):

    self.input = input_data

    self.output = np.dot(self.weights, input_data) + self.biases

    return self.output

def backward(self, output_error, learning_rate):

    assert output_error.shape == (self.weights.shape[0],), \

        f"Expected output_error shape {(self.weights.shape[0],)}, " \

        f"got {output_error.shape}"

    assert self.input is not None, \

        "Input not stored. Did you call forward() first?"

    input_error = np.dot(self.weights.T, output_error)

    weights_error = np.outer(output_error, self.input)

    biases_error = output_error

    self.grad_weights = weights_error

    self.grad_biases = biases_error

    return input_error
```

## Appendix C: losses code

```python
"""losses.py

Loss functions for neural network training.

Each loss function is implemented with both the loss calculation and its derivative.

"""

import numpy as np

def mse_loss(y_true, y_pred):
    """

    Mean Squared Error (MSE) loss function.

    Formula: MSE = (1/n) * Σ(y_pred - y_true)^2

    Args:

        y_true: Ground truth values (target)

        y_pred: Predicted values from the network

    Returns:

        Scalar loss value

    """

    return np.mean((y_pred - y_true) ** 2)

def mse_loss_prime(y_true, y_pred):
    """

    Derivative of Mean Squared Error loss with respect to predictions.

    Formula: dMSE/dy_pred = (2/n) * (y_pred - y_true)

    Args:

        y_true: Ground truth values (target)

        y_pred: Predicted values from the network

    Returns:

        Gradient of the loss with respect to predictions

    """

    return 2 * (y_pred - y_true) / y_true.size
```

## Appendix D: network code

```python
# lib/network.py

import numpy as np

from .optimizer import SGD

class Network:

    def __init__(self):

        self.layers = []

        self.loss = None

        self.loss_prime = None

    def add(self, layer):

        self.layers.append(layer)

    def use_loss(self, loss, loss_prime):

        self.loss = loss

        self.loss_prime = loss_prime

    def predict(self, input_data):

        samples = len(input_data)

        result = []

        for i in range(samples):

            output = input_data[i]

            for layer in self.layers:

                output = layer.forward(output)

            result.append(output)

        return result

    def train(self, x_train, y_train, epochs, learning_rate):

        # Initialize optimizer

        optimizer = SGD(learning_rate=learning_rate)

        samples = len(x_train)

        for i in range(epochs):
```

```python
        total_error = 0

        for j in range(samples):

            # Forward pass

            output = x_train[j]

            for layer in self.layers:

                output = layer.forward(output)

            # Track loss

            total_error += self.loss(y_train[j], output)

            # Backward pass

            error = self.loss_prime(y_train[j], output)

            for layer in reversed(self.layers):

                # Compute gradients without updating weights immediately

                error = layer.backward(error, learning_rate=None)

            # Optimizer step: Update weights and biases

            for layer in self.layers:

                optimizer.update(layer)

        # Log average error every 1000 epochs

        total_error /= samples

        if (i + 1) % 1000 == 0:

            print(f"Epoch {i + 1}/{epochs}    error={total_error:.6f}")
```

## Appendix E: optimizer code

```python
# lib/optimizer.py

import numpy as np

class SGD:
    """

    Stochastic Gradient Descent (SGD) optimizer.

    Updates weights and biases using:

        W = W - learning_rate * dW

        b = b - learning_rate * db

    """

    def __init__(self, learning_rate=0.01):

        self.learning_rate = learning_rate

    def update(self, layer):
        """

        Perform one gradient descent step on a single layer.

        Assumes layer has:

            layer.weights

            layer.biases

            layer.grad_weights

            layer.grad_biases

        """

        if hasattr(layer, 'weights') and layer.grad_weights is not None:

            layer.weights -= self.learning_rate * layer.grad_weights

        if hasattr(layer, 'biases') and layer.grad_biases is not None:

            layer.biases -= self.learning_rate * layer.grad_biases
```