

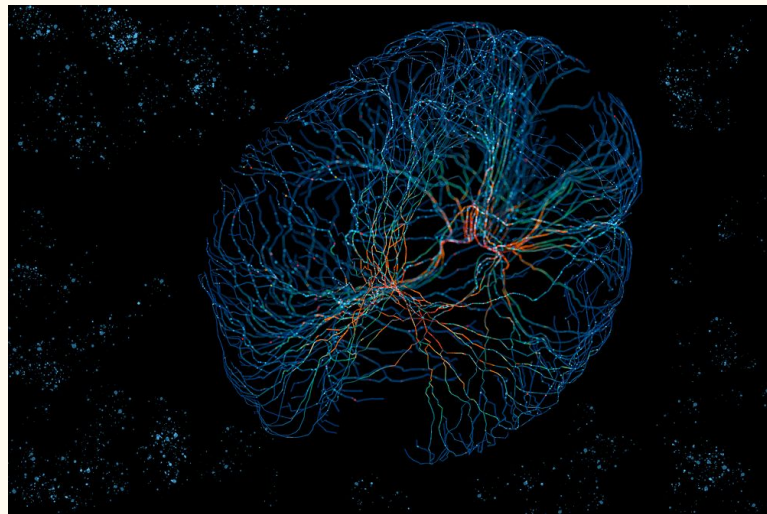
# Parallel Deep Learning Neural Network: Final Presentation

---

Thomas Keller, Zack Humphries, Anuradha Agarwal

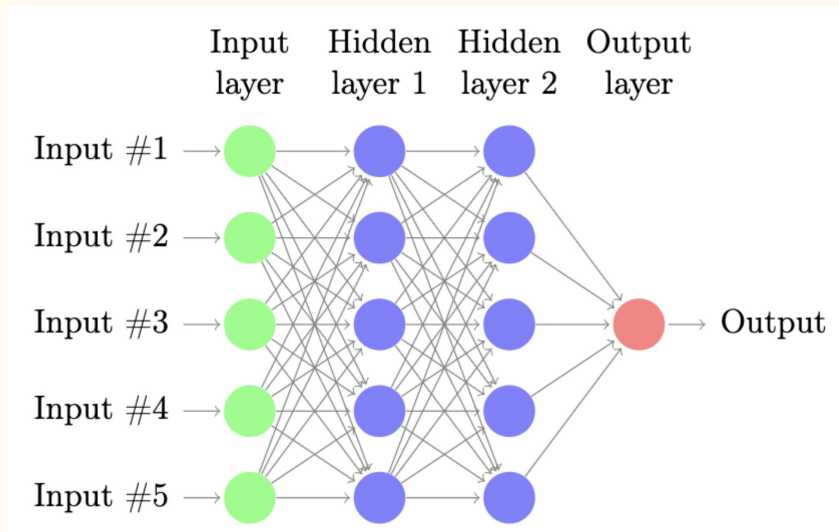
# Outline

- Introduction
- Literature Review
- OpenMP
- CUDA
- Conclusion



# Introduction

- Breast Cancer Data Set
- Deep Neural Network - serial code
  - Binary classification
  - 2 Hidden Layers
  - Output Layer with 1 node
  - 30/4800 inputs
- Parallelizing:
  - OpenMP
  - CUDA



# Literature Review: Exploring Comparison Code

## Deep Neural Network code

Complete Neural Network

OpenMP

CUDA

Comparison Forward Propagation

Mattia Orlandi

OpenMP

CUDA

Ours

OpenMP

CUDA

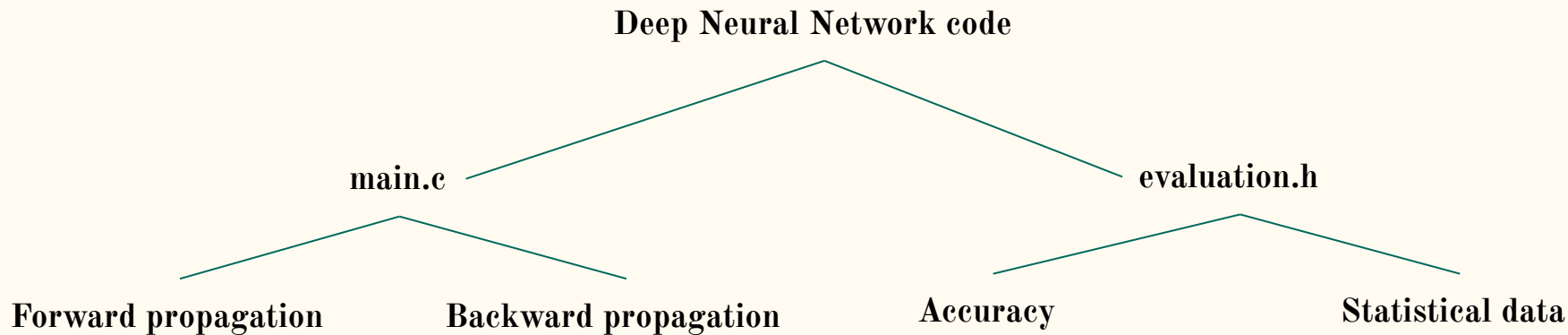
### Legend

Fully Functioning

Partially Functioning (Limited Inputs)

Unavailable

# Serial Code and Accuracy Evaluation



# Accuracy and Confusion matrix

- All metrics are determined by the test data set with 114 test cases.
- 30 inputs, 2 hidden layers
- Epochs: 1500
- Learning rate: 0.001
- 98.25% accuracy

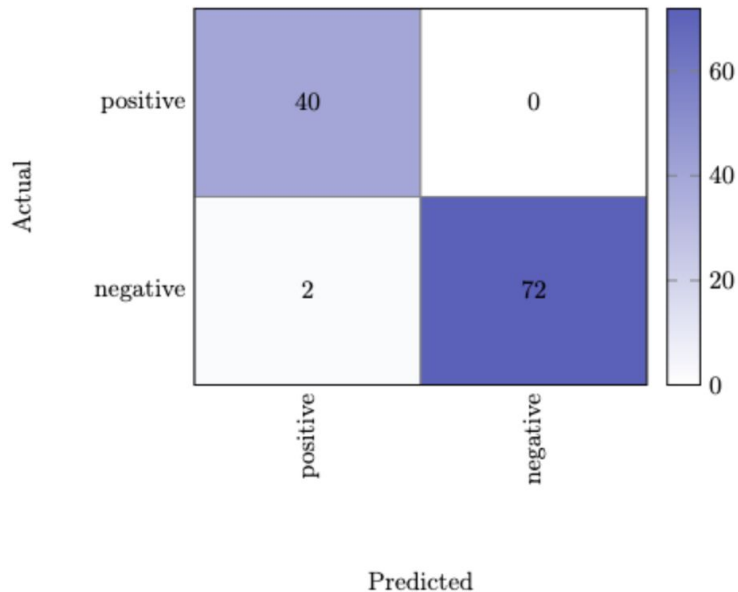


Figure 2: Confusion matrix error metrics

# OpenMP Code: Forward

```
1 // forward pass
2 // compute hidden layer activation
3
4 // hidden layer 1
5 #pragma omp parallel for num_threads(thread_count)
6 for(int j =0; j < numHiddenNodes; j++){
7     double activation = hiddenLayerBias[j];
8
9     for(int k = 0; k < numInputs; k++){
10         activation += trainingInputs[i][k] * hiddenWeights[k][j];
11     }
12
13     hiddenLayer[j] = relu(activation);
14 }
15
16 // hidden layer 2
17 double activation = 0;
18 #pragma omp parallel for reduction(+:activation) num_threads(thread_count)
19 for(int j =0; j < numHiddenNodes2; j++){
20     activation = hiddenLayerBias2[j];
21
22     for(int k = 0; k < numHiddenNodes; k++){
23         activation += hiddenLayer[k] * hiddenWeights2[k][j];
24     }
25
26     hiddenLayer2[j] = relu(activation);
27 }
28
29 // compute output layer activation
30 #pragma omp parallel for reduction(+:activation) num_threads(thread_count)
31 for(int j =0; j < numOutputs; j++){
32     activation = outputLayerBias[j];
33
34     for(int k = 0; k < numHiddenNodes2; k++){
35         activation += hiddenLayer2[k] * outputWeights[k][j];
36     }
37
38     outputLayer[j] = sigmoid(activation);
39 }
```

Listing 3: Forward path openMP of main.c

- The first loop computes the activation of the neurons in the first hidden layer.
- Each thread executes the loop body for a subset of the iterations.
- The second loop computes the activation of the neurons in the second hidden layer.
- The third loop computes the activation of the output layer neurons
- In addition, the reduction clause is used to sum up the values of the activation variable across threads.

# OpenMP Code: Backward

```
1 // Backpropagation
2 // Compute change in output weights
3 double deltaOutput[numOutputs];
4 #pragma omp parallel for num_threads(thread_count)
5 for(int j = 0; j < numOutputs; j++){
6     double error = (trainingOutputs[i][j] - outputLayer[j]); // L1
7     deltaOutput[j] = error * dSigmoid(outputLayer[j]);
8 }
9
10 // Compute change in hidden weights (second layer)
11 double deltaHidden2[numHiddenNodes2];
12 #pragma omp parallel for num_threads(thread_count)
13 for(int j = 0; j < numHiddenNodes2; j++){
14     double error = 0.0f;
15     for(int k = 0; k < numOutputs; k++){
16         error += deltaOutput[k] * outputWeights[j][k];
17     }
18     deltaHidden2[j] = error * dRelu(hiddenLayer[j]);
19 }
20
21 // Compute change in hidden weights (first layer)
22 double deltaHidden[numHiddenNodes];
23 #pragma omp parallel for num_threads(thread_count)
24 for(int j = 0; j < numHiddenNodes; j++){
25     double error = 0.0f;
26     for(int k = 0; k < numHiddenNodes2; k++){
27         error += deltaHidden2[k] * hiddenWeights2[j][k];
28     }
29     deltaHidden[j] = error * dRelu(hiddenLayer2[j]);
30 }
31
32 // Apply change in output weights
33 #pragma omp parallel for num_threads(thread_count)
34 for(int j = 0; j < numOutputs; j++){
35     outputLayerBias[j] += deltaOutput[j] * lr;
36     for(int k = 0; k < numHiddenNodes2; k++){
37         outputWeights[k][j] += hiddenLayer2[k] * deltaOutput[j] * lr;
38     }
39 }
40
41 // Apply change in second hidden layer weights
42 #pragma omp parallel for num_threads(thread_count)
43 for(int j = 0; j < numHiddenNodes2; j++){
44     hiddenLayerBias[j] += deltaHidden[j] * lr;
45     for(int k = 0; k < numHiddenNodes; k++){
46         hiddenWeights2[k][j] += hiddenLayer[k] * deltaHidden[j] * lr;
47     }
48 }
49
50 // Apply change in first hidden layer weights
51 #pragma omp parallel for num_threads(thread_count)
52 for(int j = 0; j < numHiddenNodes; j++){
53     hiddenLayerBias[j] += deltaHidden[j] * lr;
54     for(int k = 0; k < numInputs; k++){
55         hiddenWeights[k][j] += trainingInputs[i][k] * deltaHidden[j] * lr;
56     }
57 }
```

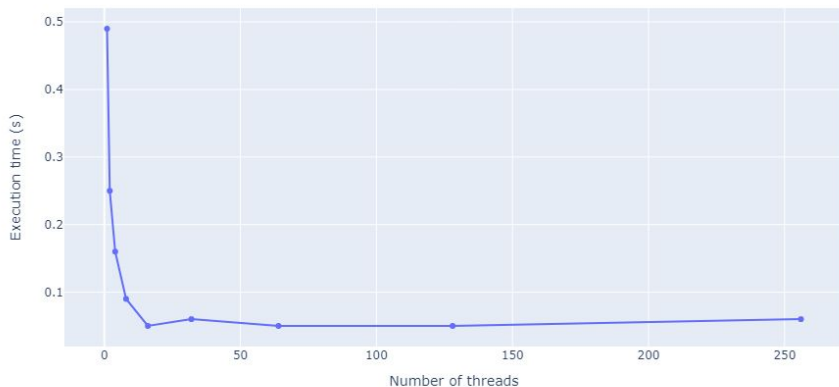
Listing 4: Backward path openMP of main.c

- The code starts by computing the change in output weights using the training inputs and outputs.
- Next, the code computes the change in hidden weights for the second layer and then the first layer.
- After computing the changes in the weights, the code applies them to the network by updating the weights and biases.

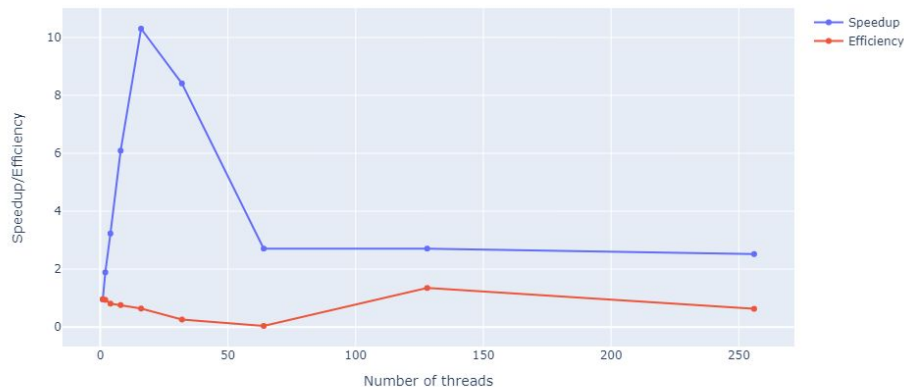


# Time and Speedup (OpenMP)

Execution times for parallelization with OpenMP with only Forward Propagation

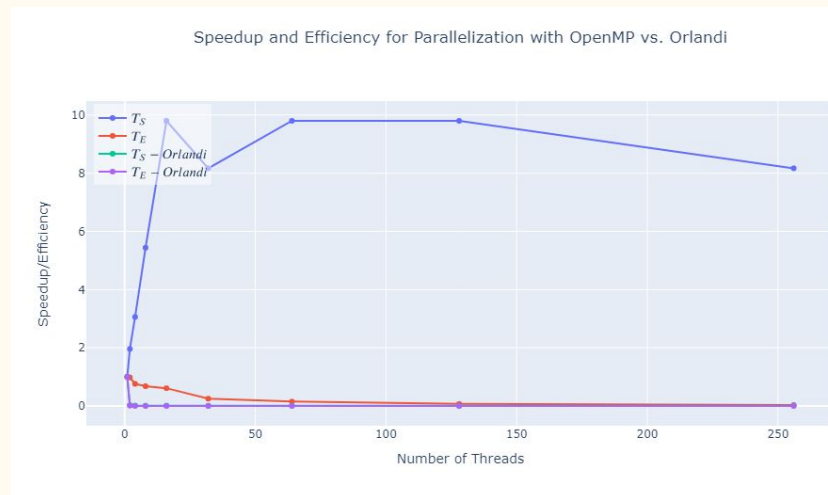
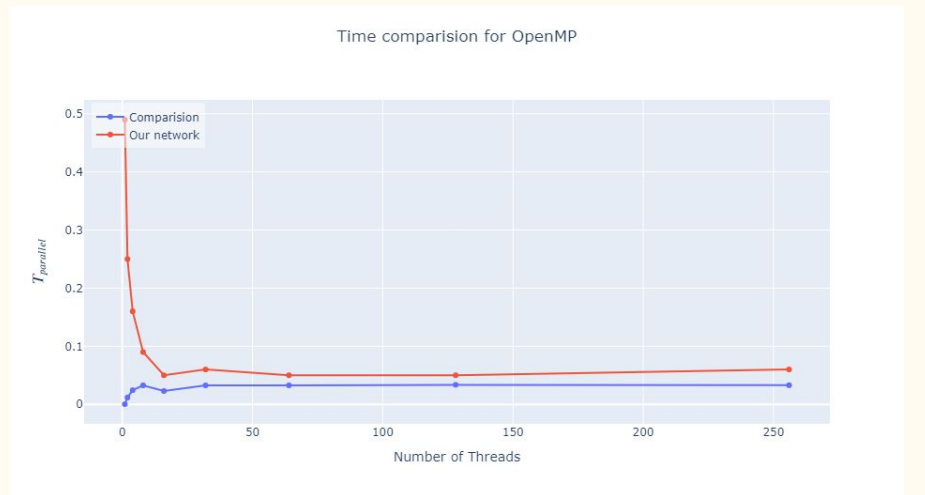


Speedup and efficiency using OpenMP



- The parallel execution times decrease as the number of threads increases, indicating that the parallelization is effective. However, the execution time does not continue to decrease at the same rate as the number of threads increases. The speedup increases as the number of processors increases, but with diminishing returns.

# Forward Propagation Comparison (Orlandi)



- Our DNN takes longer to train compared to Orlandi's; nevertheless, we are able to parallelize the process using OpenMP efficiently. We are able to reach higher speedups, and our code is more efficient comparatively.

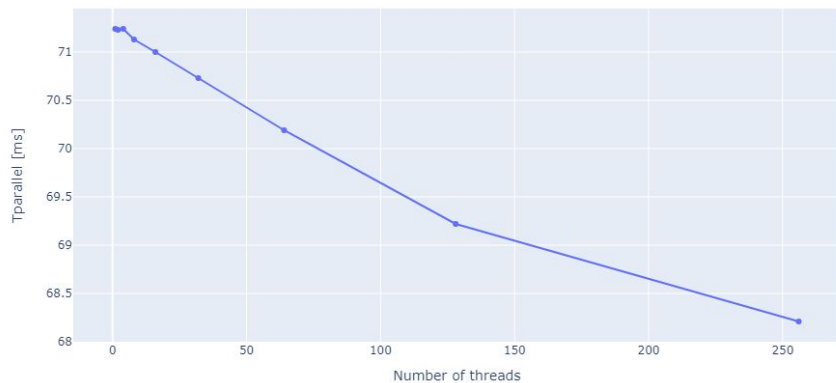
# CUDA Code Forward and Backward

```
1 // Forward propagation kernel
2 __global__ void forward_kernel(double* X, double* W1, double* W2, double* W3, double* b1,
3 double* b2, double* b3, double* hidden1, double* hidden2, double* output) {
4     int tid = blockIdx.x * blockDim.x + threadIdx.x;
5
6     if (tid < HIDDEN_SIZE) {
7         double sum = 0.0;
8         for (int j = 0; j < INPUT_SIZE; j++) {
9             sum += X[j] * W1[j * HIDDEN_SIZE + tid];
10        }
11        hidden1[tid] = relu(sum + b1[tid]);
12    }
13    __syncthreads();
14
15    if (tid < HIDDEN_SIZE) {
16        double sum = 0.0;
17        for (int i = 0; i < HIDDEN_SIZE; i++) {
18            sum += hidden1[i] * W2[i * HIDDEN_SIZE + tid];
19        }
20        hidden2[tid] = relu(sum + b2[tid]);
21    }
22    __syncthreads();
23
24    if (tid == 0) {
25        double sum = 0.0;
26        for (int i = 0; i < HIDDEN_SIZE; i++) {
27            sum += hidden2[i] * W3[i];
28        }
29        *output = sigmoid(sum + b3[0]);
30    }
31 }
32 }
```

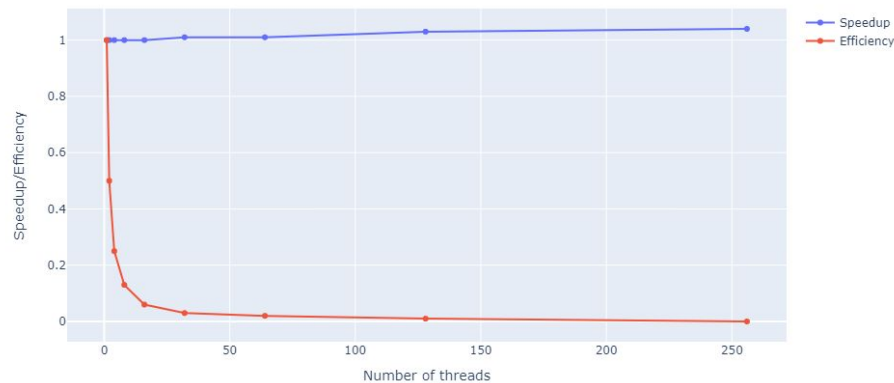
```
1 __global__ void backward_kernel(double* X, double* W1, double* W2, double* W3, double* b1,
2 double* b2, double* b3, double* hidden1, double* hidden2, double* output, double target)
3 ) {
4     int tid = blockIdx.x * blockDim.x + threadIdx.x;
5
6     if (tid == 0) {
7         // calculate output weight
8         double d_output = (*output - target) * dSigmoid(*output);
9
10        // calculate hidden 2 weight
11        double d_hidden2[HIDDEN_SIZE];
12        for (int i = 0; i < HIDDEN_SIZE; i++) {
13            d_hidden2[i] = dRelu(hidden2[i]) * W3[i] * d_output;
14        }
15
16        // calculate hidden 1 weight
17        double d_hidden1[HIDDEN_SIZE];
18        for (int i = 0; i < HIDDEN_SIZE; i++) {
19            double sum = 0.0;
20            for (int j = 0; j < HIDDEN_SIZE; j++) {
21                sum += W3[j * HIDDEN_SIZE + i] * d_hidden2[j];
22            }
23            d_hidden1[i] = dRelu(hidden1[i]) * sum;
24        }
25
26        // update hidden weights 1 & bias
27        for (int i = 0; i < HIDDEN_SIZE; i++) {
28            for (int j = 0; j < INPUT_SIZE; j++) {
29                W1[j * HIDDEN_SIZE + i] -= LEARNING_RATE * X[j] * d_hidden1[i];
30            }
31            b1[i] -= LEARNING_RATE * d_hidden1[i];
32        }
33
34        // update hidden weights 2 & bias
35        for (int i = 0; i < HIDDEN_SIZE; i++) {
36            for (int j = 0; j < HIDDEN_SIZE; j++) {
37                W2[j * HIDDEN_SIZE + i] -= LEARNING_RATE * hidden1[i] * d_hidden2[j];
38            }
39            b2[i] -= LEARNING_RATE * d_hidden2[i];
40        }
41
42        // update output weights & bias
43        for (int i = 0; i < HIDDEN_SIZE; i++) {
44            W3[i] -= LEARNING_RATE * hidden2[i] * d_output;
45        }
46        b3[0] -= LEARNING_RATE * d_output;
47    }
48 }
```

# Time and Speedup (CUDA)

Execution times for parallelization with Cuda



Speedup and efficiency for parallelization with CUDA



- The execution time remains almost constant, around 71 seconds, with only a slight improvement as the number of processors increases. This indicates that the parallelization is not effective in reducing the overall execution time.
- Note that there our CUDA code could not handle up to 4800 inputs (so limited to 570). Plus, we were unable to compare it to Orlandi's code so these times should be taken with “a grain of salt”

# Conclusion

- Forward Propagation
- Backward Propagation
- OpenMP - speedup
  - 8 and 16 Threads
- CUDA - no speedup (570 inputs)
  - Matrix size
  - Lack of knowledge
  - Time



# References

- W. Nick Street Dr. William H. Wolberg and Olvi L. Mangasarian, UCI machine learning repository: Breast cancer wisconsin (diagnostic) data set, University of California, Irvine, School of Information and Computer Sciences, 1995.
- Orlandi, Neural network implementation with opemmp and cuda (2020).
- [https://github.com/nihil21/parallel\\_nn/](https://github.com/nihil21/parallel_nn/)
- <https://github.com/thomkell/DeepNeuralNetwork>

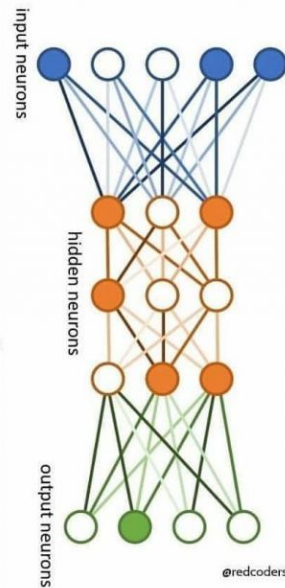
# Questions?

---

**THIS IS A NEURAL NETWORK.**

**IT MAKES MISTAKES.  
IT LEARNS FROM THEM.**

**BE LIKE A NEURAL NETWORK.**



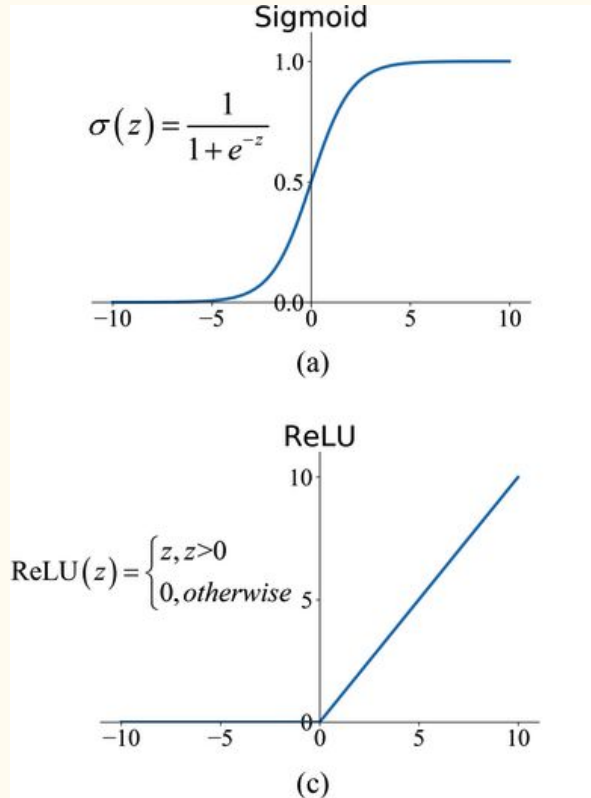
# Problem Statement and Dataset

- Both the forward and backward propagation processes involve matrix multiplication (dot product)
  - Although the results of the previous layer are the inputs for the next layer, the multiplication itself can be parallelized
- We plan to use a publicly available dataset with 569 instances of breast cancer data to predict whether certain attributes of a mass are benign or malignant
  - In order to see the effectiveness of our parallelization, we replicated the data from 30 inputs to 4800 inputs





# Backup Slide: Activation Functions



# Literature Review: Exploring Comparison Code

- Mattia Orlandi's "Neural Network Implementation with OpenMP and CUDA"
  - Explores the parallelization of only the forward propagation
    - Inner-Loop: Parallelization of the matrix multiplication
      - What we focused on
    - Outer-Loop: Parallelization of each epoch (individual forward propagation)
      - Unrealistic as there would be data discrepancies when including a back propagation
      - All of the weights and biases need to be updated through gradient calculations (gradient descent)
- We created a fully-functioning parallelized neural network with OpenMP, a parallelized forward propagation with OpenMP to compare to Orlandi's version
- We were unable to run Orlandi's CUDA code on Tuckoo as his code does not recognize Tuckoo's internal architecture
  - Able to create a partially functioning parallelized neural network with CUDA and a partially-functioning parallelized forward propagation with CUDA

# Serial Code

## Forward Propagation

```
1 // forward pass
2 // compute hidden layer activation
3
4 // hidden layer 1
5 for(int j =0; j < numHiddenNodes; j++){
6     double activation = hiddenLayerBias[j];
7
8     for(int k = 0; k < numInputs; k++){
9         activation += trainingInputs[i][k] * hiddenWeights[k][j];
10    }
11
12    hiddenLayer[j] = relu(activation);
13 }
14
15 // hidden layer 2
16 for(int j =0; j < numHiddenNodes2; j++){
17     double activation = hiddenLayerBias2[j];
18
19     for(int k = 0; k < numHiddenNodes; k++){
20         activation += hiddenLayer[k] * hiddenWeights2[k][j];
21    }
22
23    hiddenLayer2[j] = relu(activation);
24 }
25
26 // compute output layer activation
27 for(int j =0; j < numOutputs; j++){
28     double activation = outputLayerBias[j];
29
30     for(int k = 0; k < numHiddenNodes2; k++){
31         activation += hiddenLayer2[k] * outputWeights[k][j];
32    }
33
34    outputLayer[j] = sigmoid(activation);
35 }
```

## Backward Propagation

```
1 // Backpropagation
2 // Compute change in output weights
3 double deltaOutput[numOutputs];
4 for(int j = 0; j < numOutputs; j++){
5     double error = (trainingOutputs[i][j] - outputLayer[j]); // L1
6     deltaOutput[j] = error * dSigmoid(outputLayer[j]);
7 }
8
9 // Compute change in hidden weights (second layer)
10 double deltaHidden2[numHiddenNodes2];
11 for(int j = 0; j < numHiddenNodes2; j++){
12     double error = 0.0f;
13     for(int k = 0; k < numOutputs; k++){
14         error += deltaOutput[k] * outputWeights[j][k];
15     }
16     deltaHidden2[j] = error * dRelu(hiddenLayer[j]);
17 }
18
19 // Compute change in hidden weights (first layer)
20 double deltaHidden[numHiddenNodes];
21 for(int j = 0; j < numHiddenNodes; j++){
22     double error = 0.0f;
23     for(int k = 0; k < numHiddenNodes2; k++){
24         error += deltaHidden2[k] * hiddenWeights2[j][k];
25     }
26     deltaHidden[j] = error * dRelu(hiddenLayer2[j]);
27 }
28
29 // Apply change in output weights
30 for(int j = 0; j < numOutputs; j++){
31     outputLayerBias[j] += deltaOutput[j] * lr;
32     for(int k = 0; k < numHiddenNodes2; k++){
33         outputWeights[k][j] += hiddenLayer2[k] * deltaOutput[j] * lr;
34     }
35 }
36
37 // Apply change in second hidden layer weights
38 for(int j = 0; j < numHiddenNodes2; j++){
39     hiddenLayerBias[j] += deltaHidden[j] * lr;
40     for(int k = 0; k < numHiddenNodes; k++){
41         hiddenWeights2[k][j] += hiddenLayer[k] * deltaHidden[j] * lr;
42     }
43 }
44
45 // Apply change in first hidden layer weights
46 for(int j = 0; j < numHiddenNodes; j++){
47     hiddenLayerBias[j] += deltaHidden[j] * lr;
48     for(int k = 0; k < numInputs; k++){
49         hiddenWeights[k][j] += trainingInputs[i][k] * deltaHidden[j] * lr;
50     }
51 }
```

# Discussion/Challenges

- As a showcase of the parallelization of neural networks, we have expanded upon Orlandi's forward propagation parallelization and have implemented a fully functioning parallelized neural network.
- Finally, as with any software project, future iterations could benefit from additional testing and optimization to ensure optimal performance and accuracy.
- Overall, we are proud of the performance of our neural network model and have a deeper understanding of how neural networks and parallelization work.