

GPU-Accelerated Backpropagation with C/CUDA

A presentation by Drew Hans

“ *Can a CUDA-enabled graphics processing unit (GPU) be used to reduce the time it takes to train a fully-connected three-layer artificial neural network?* ”

Research Question

Presentation Overview

- Fundamentals of Machine Learning & Artificial Neural Networks
- GPU-Accelerated Computing with CUDA
- My Research Project: Neuronmancer
- Project Challenges, Successes, & Results

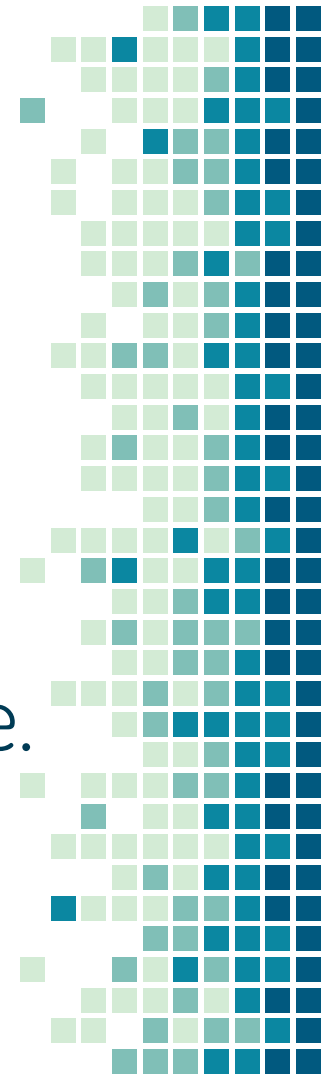
Fundamentals of Machine Learning

Let's cover the basics



Learning

is the process of turning
experience into knowledge.



3 Major Machine Learning Approaches

Supervised Learning

The output you want from the model is known.

Two subgroups:

- Classification
- Regression

Examples:

- Image recognition
- Speech recognition
- Stock price prediction

Unsupervised Learning

The output you want from the model is not known.

One big subgroup:

- Clustering

Examples:

- Netflix recommendations
- Medical imaging
- Gene sequence analysis

Reinforcement Learning

The model has agency, meaning it is able to interact with its environment, and receives "rewards" and "punishments" for its actions.

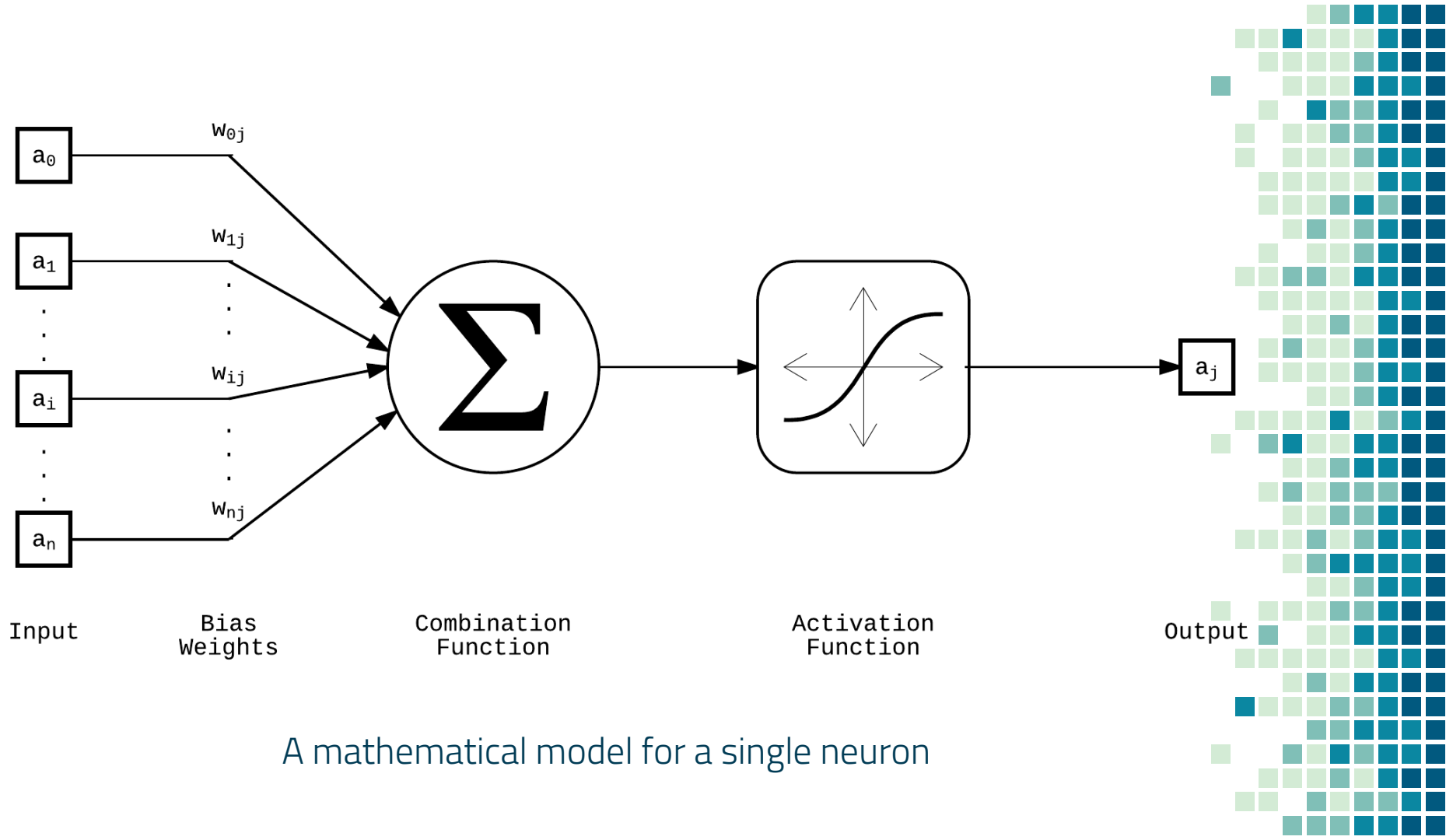
Examples:

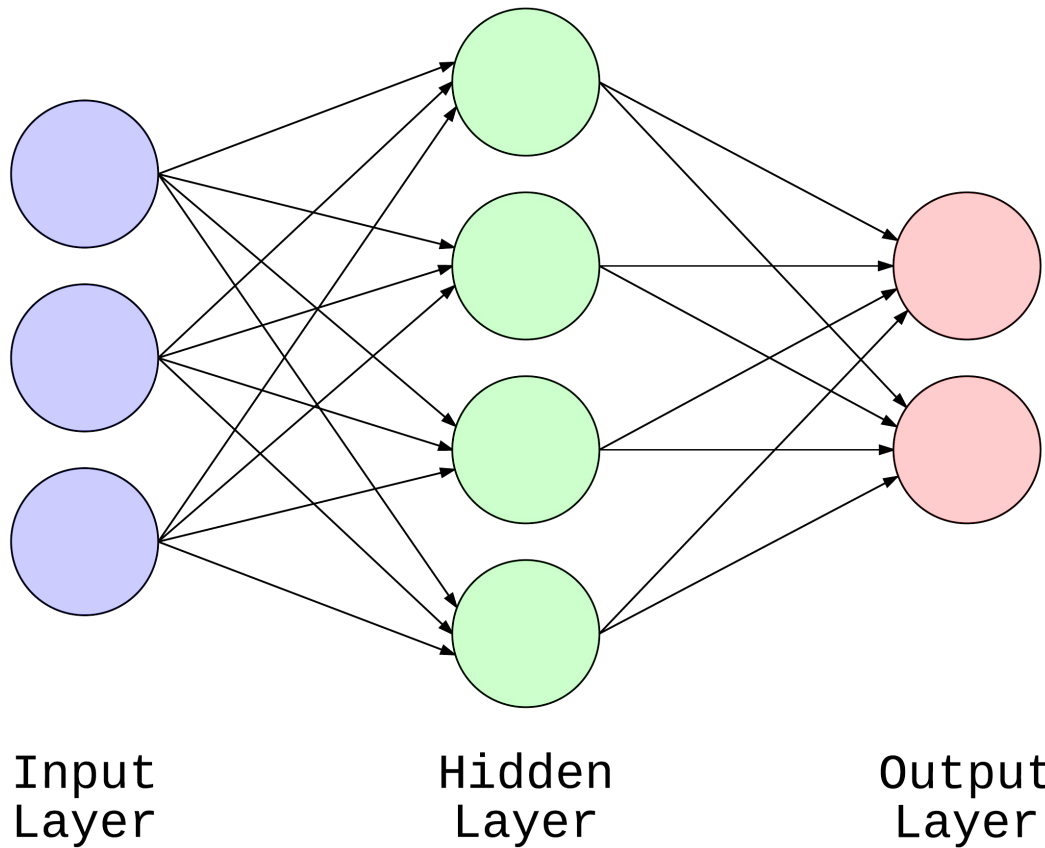
- Playing Checkers
- Playing Go
- Robot Control

Artificial Neural Networks

A type of machine learning model





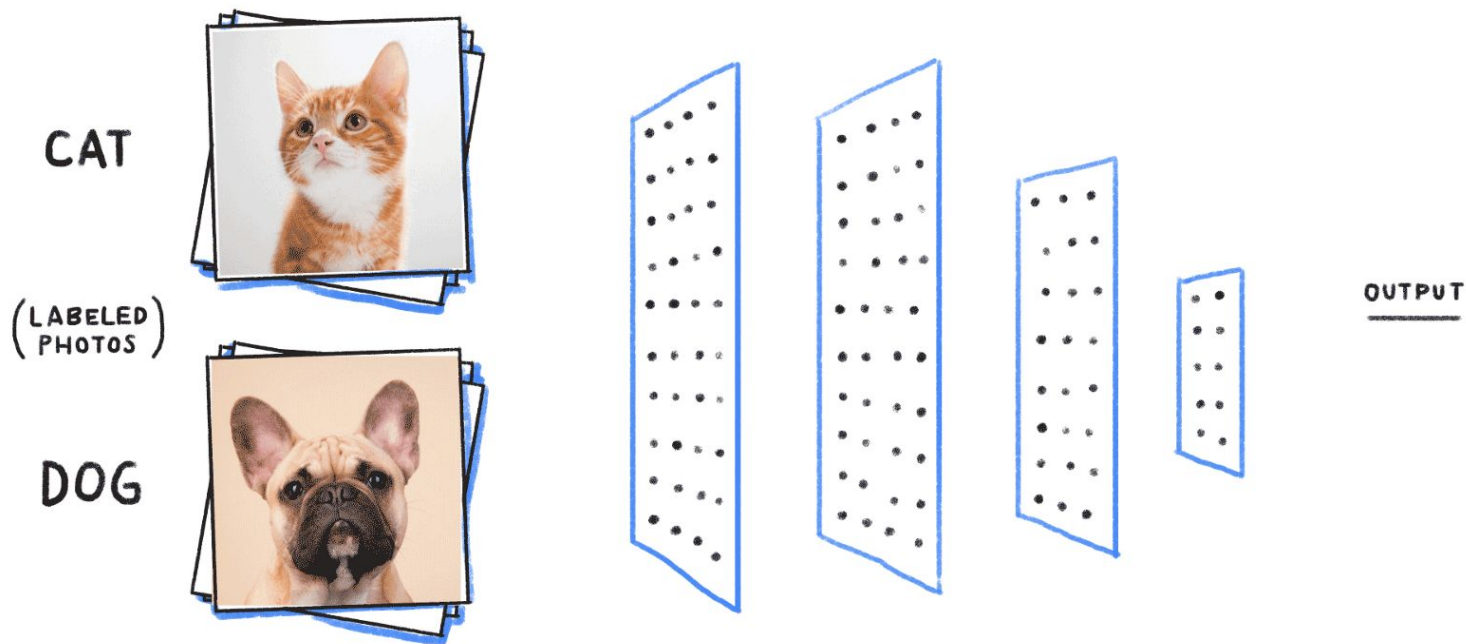


Input
Layer

Hidden
Layer

Output
Layer

Visualization of a simple 3-layer feedforward artificial neural network



Visualization of a 4-layer neural network that identifies cats!

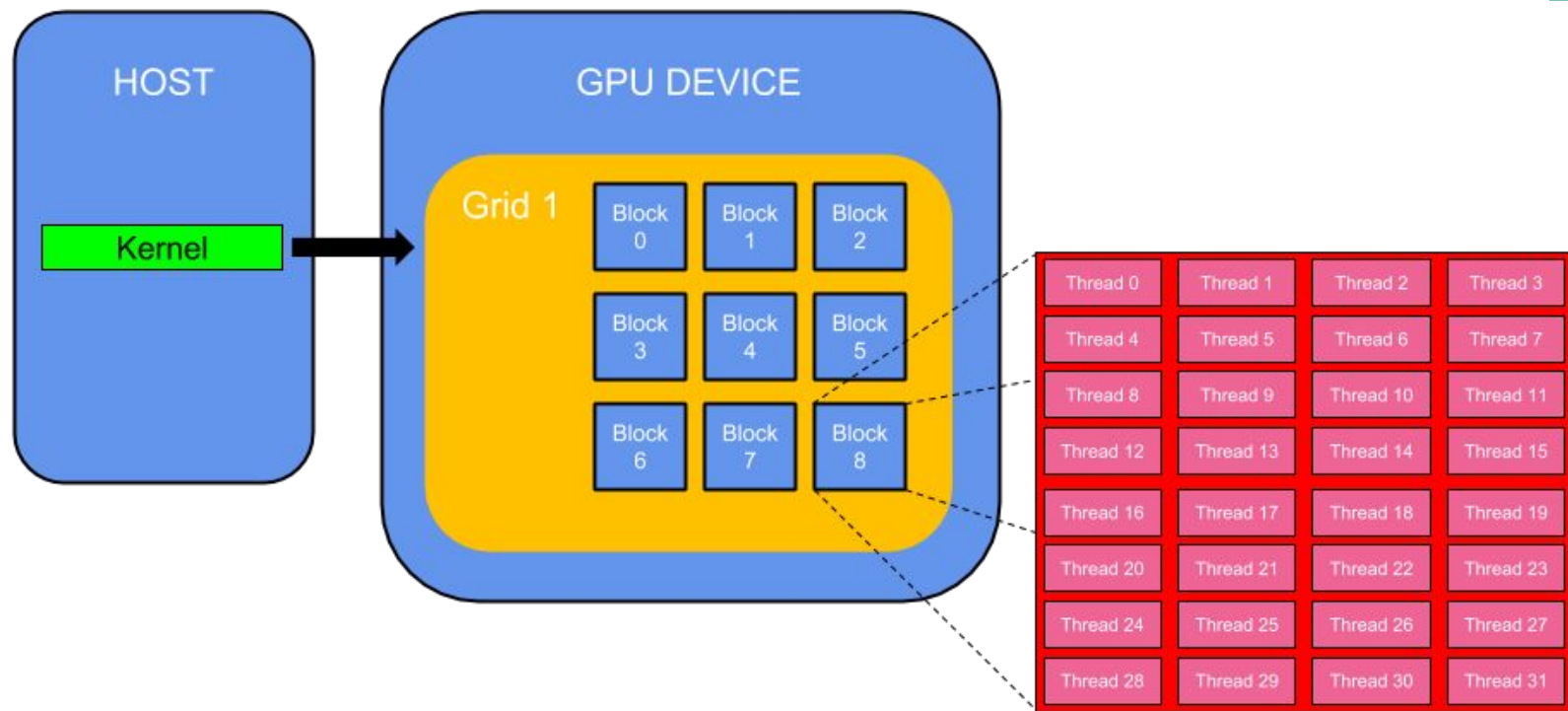
GPU-Accelerated Computing with CUDA

A parallel programming model





This is a CUDA-enabled GPU device - the GTX 1080



Visualization of a simple 3-layer feedforward artificial neural network

```
int main(void) {  
  
    const int size = 5;  
  
    const int a[size] = { 1, 2, 3, 4, 5 };  
    const int b[size] = { 10, 20, 30, 40, 50 };  
    int c[size] = { 0 };  
  
    // Add vectors sequentially  
    for (int i = 0; i < size; i++) {  
        c[i] = a[i] + b[i];  
    }  
  
} // end main function
```

A simple C program for adding array elements sequentially.

```

int main(void) {
    const int size = 5;

    const int host_a[size] = { 1, 2, 3, 4, 5 };
    const int host_b[size] = { 10, 20, 30, 40, 50 };
    int host_c[size] = { 0 };

    int* dev_a = 0;
    int* dev_b = 0;
    int* dev_c = 0;

    cudaSetDevice(0); // run kernels on GPU 0

    // allocate GPU memory for device copies of a, b, c
    cudaMalloc((void**)&dev_a, size * sizeof(int));
    cudaMalloc((void**)&dev_b, size * sizeof(int));
    cudaMalloc((void**)&dev_c, size * sizeof(int));

    // copy a and b values from host memory into allocated GPU memory
    cudaMemcpy(dev_a, host_a, size * sizeof(int), cudaMemcpyHostToDevice);
    cudaMemcpy(dev_b, host_b, size * sizeof(int), cudaMemcpyHostToDevice);

    // add array values in parallel (launch 1 block with "size" threads)
    kernel_add<<< 1, size >>>(dev_c, dev_a, dev_b);

    // copy c values from GPU memory into host memory
    cudaMemcpy(host_c, dev_c, size * sizeof(int), cudaMemcpyDeviceToHost);

    // free GPU memory once we're done with it
    cudaFree(dev_a);
    cudaFree(dev_b);
    cudaFree(dev_c);
} // end main function

```

```

__global__
void kernel_add(int* dev_c,
               const int* dev_a,
               const int* dev_b ) {
    int i = threadIdx.x;
    dev_c[i] = dev_a[i] + dev_b[i];
} // end kernel_add cuda kernel

```

A simple C / CUDA
program for adding
array elements
sequentially.

Neuronmancer (v3)

My research project

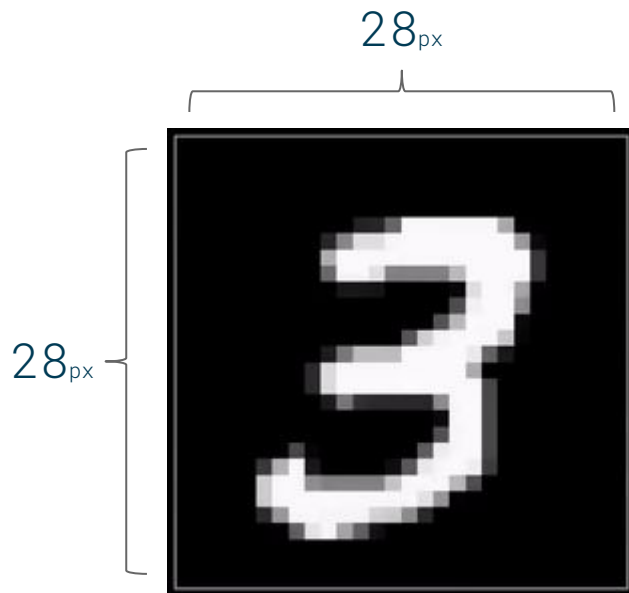


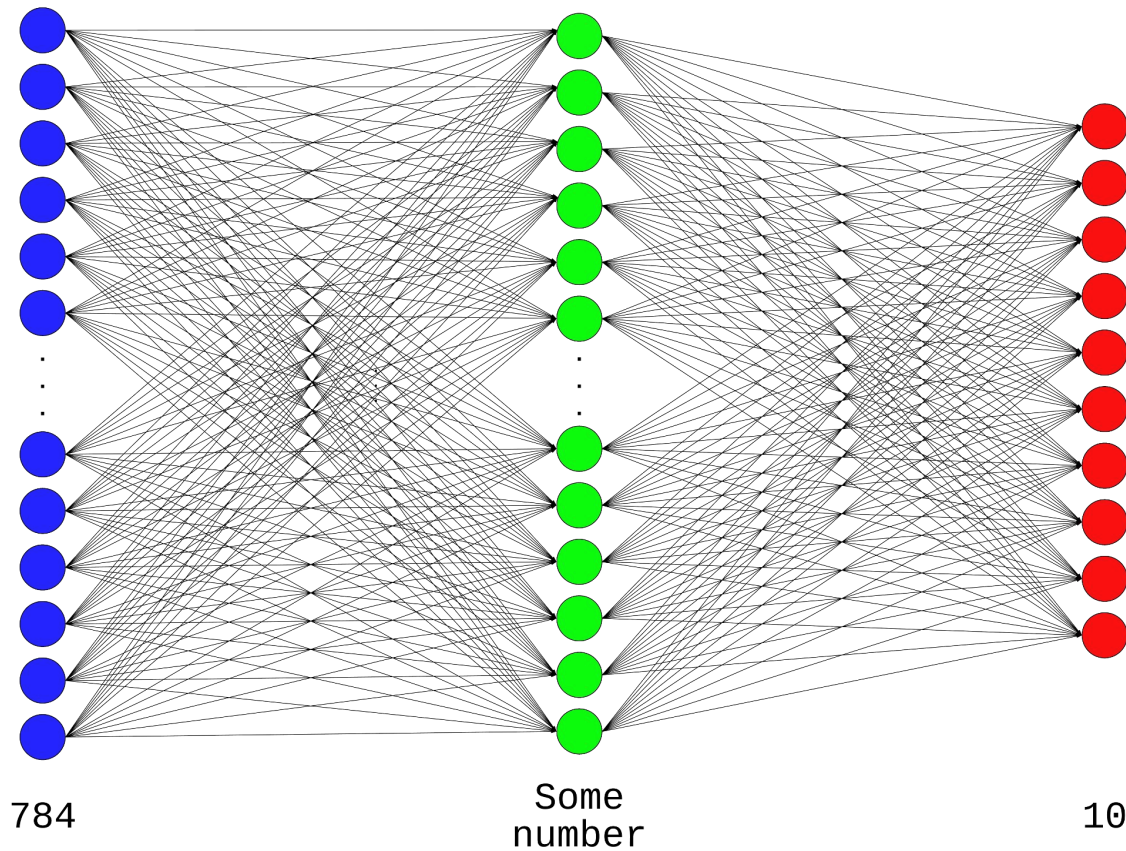
(8, 9) (0, 0) (2, 2) (6, 6) (7, 7) (8, 8) (3, 3) (9, 9)
 (0, 0) (4, 4) (6, 6) (7, 7) (4, 4) (6, 6) (8, 8) (0, 0)
 (7, 7) (8, 8) (3, 3) (1, 1) (5, 5) (7, 7) (1, 1) (7, 7)
 (1, 1) (1, 1) (6, 6) (3, 3) (0, 0) (2, 2) (9, 9) (3, 3)
 (1, 1) (1, 1) (0, 0) (4, 4) (9, 9) (2, 2) (0, 0) (0, 0)
 (2, 2) (0, 0) (2, 2) (7, 7) (1, 1) (8, 8) (6, 6) (4, 4)
 (1, 1) (6, 6) (3, 3) (4, 4) (5, 5) (9, 9) (1, 1) (3, 3)
 (3, 3) (8, 8) (5, 5) (4, 4) (2, 7) (7, 7) (4, 4) (2, 2)
 (8, 8) (5, 5) (8, 8) (6, 6) (7, 7) (3, 3) (4, 4) (6, 6)
 (1, 1) (9, 9) (9, 9) (6, 6) (0, 0) (3, 3) (7, 7) (2, 2)

MNIST DATASET SAMPLES

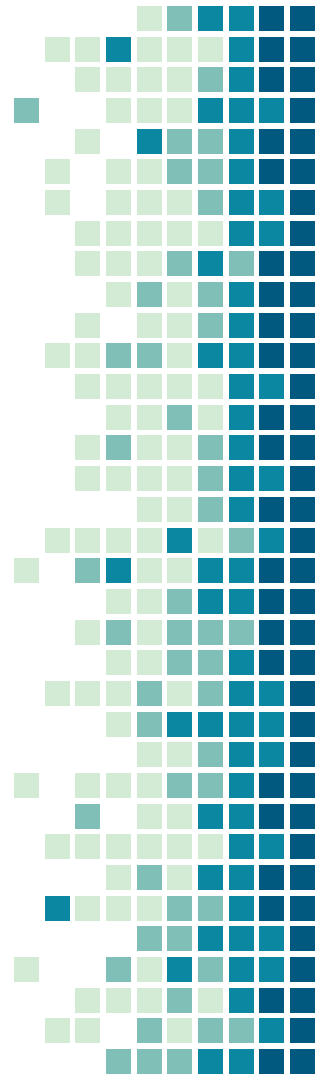
$28 \times 28 = 784$ pixels
in a single image.

10 possible classifications
for any given image.





Visualization of neural network models created by my program



MY DATA GATHERING PROCESS



```

root@ubuntu: ~/Desktop/neuromancer
Note, these parameters can be changed in the main.h File.

what would you like to do?
Enter 0 to quit, 1 to train on host, 2 to train on GPU device, or 3 to evaluate:
-3

--- beginning evaluation! ---
== sample 10000 of 10000 complete
== sample 10000 of 10000 complete
== sample 10000 of 10000 complete
== sample 10000 of 10000 complete
== sample 10000 of 10000 complete
== sample 10000 of 10000 complete
== sample 10000 of 10000 complete
== sample 10000 of 10000 complete
== sample 10000 of 10000 complete
== sample 10000 of 10000 complete
--- evaluation complete ---

  PREDICTED + PREDICTED + PREDICTED + PREDICTED + PREDICTED + PREDICTED + PREDICTED + PREDICTED + PREDICTED + PREDICTED +
  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| ACTUAL 0 | 33 | 15 | 0 | 343 | 0 | 44 | 0 | 10 | 0 | 513 |
| ACTUAL 1 | 2 | 0 | 0 | 827 | 0 | 140 | 2 | 23 | 0 | 342 |
| ACTUAL 2 | 25 | 11 | 0 | 145 | 0 | 173 | 0 | 11 | 4 | 401 |
| ACTUAL 3 | 10 | 0 | 0 | 767 | 0 | 44 | 0 | 14 | 16 | 119 |
| ACTUAL 4 | 14 | 1 | 0 | 534 | 0 | 23 | 0 | 211 | 5 | 190 |
| ACTUAL 5 | 7 | 0 | 0 | 512 | 0 | 42 | 0 | 0 | 15 | 107 |
| ACTUAL 6 | 18 | 7 | 0 | 440 | 0 | 336 | 0 | 24 | 2 | 322 |
| ACTUAL 7 | 4 | 5 | 0 | 293 | 0 | 17 | 0 | 104 | 3 | 413 |
| ACTUAL 8 | 3 | 0 | 0 | 529 | 0 | 22 | 0 | 113 | 4 | 103 |
| ACTUAL 9 | 1 | 1 | 0 | 528 | 0 | 17 | 0 | 291 | 1 | 176 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
Accuracy = 13.22%
MisclassificationRate = 86.78%

Evaluation completed in 1.00 seconds!

what would you like to do?
Enter 0 to quit, 1 to train on host, 2 to train on GPU device, or 3 to evaluate:
-3

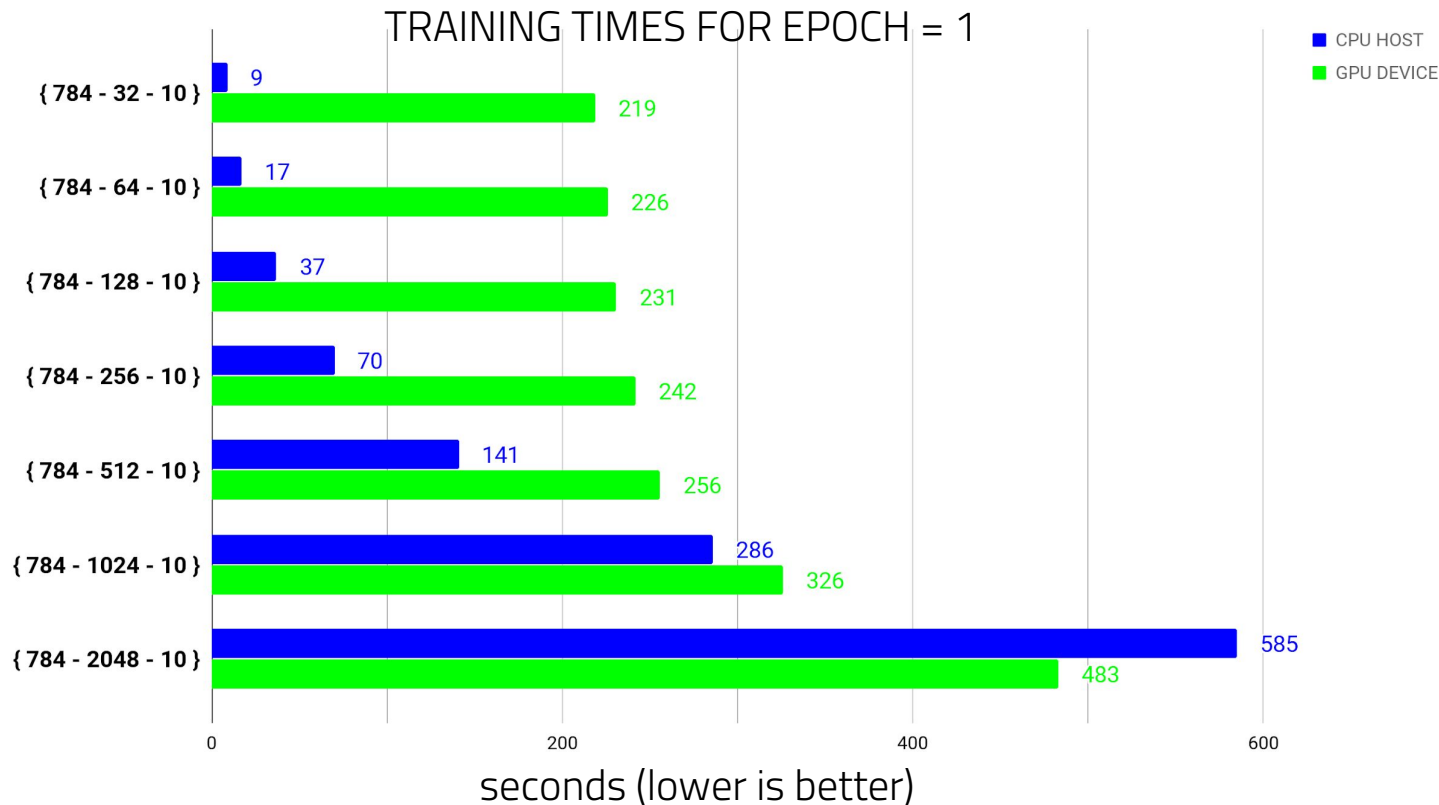
--- beginning training on host ---
--- starting epoch 1 of 1 ---
== sample 10000 of 10000 complete

```

RESULTS (1 of 2)

Intel® Core™ i7-4770K CPU (3.50GHz)

NVIDIA GeForce GTX 670 (1344 CUDA Cores)

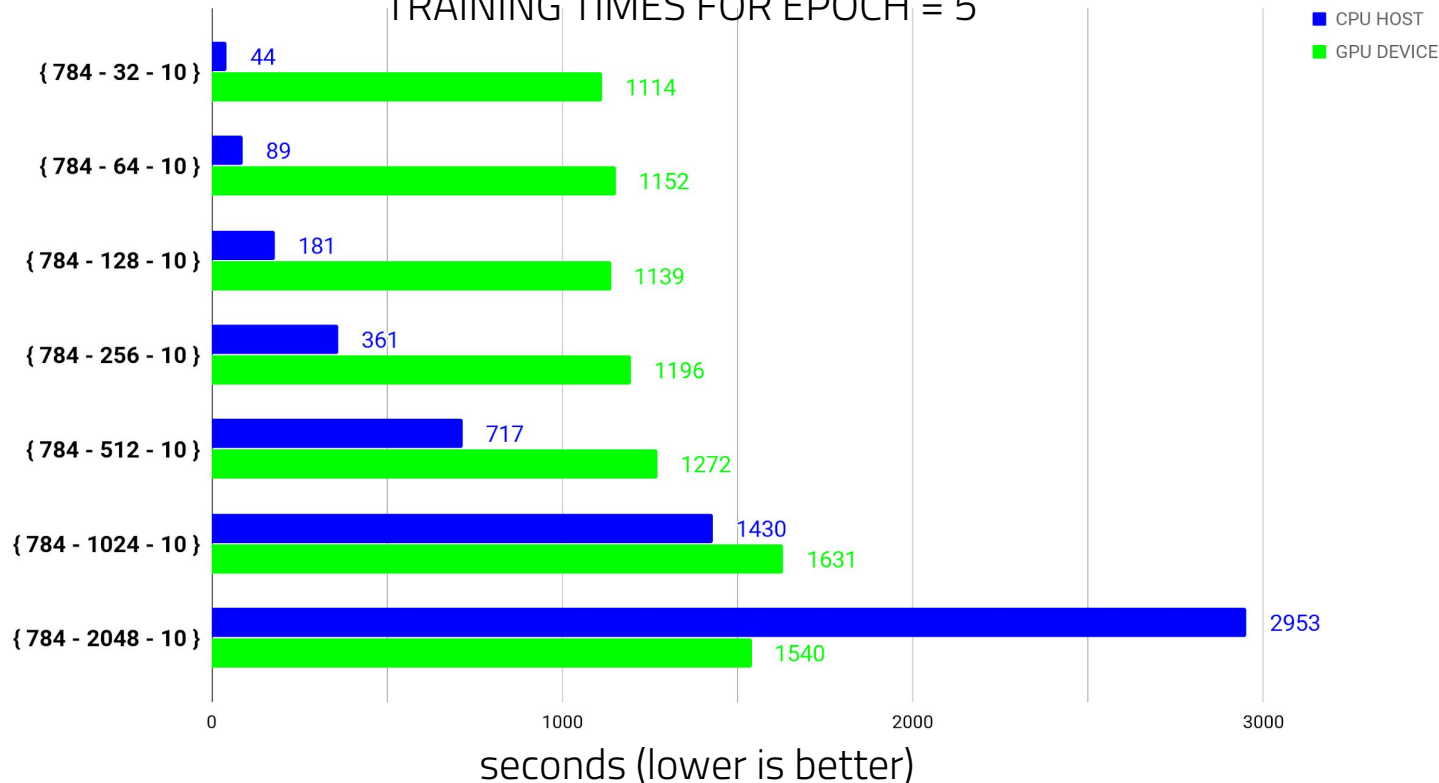


RESULTS (2 of 2)

Intel[®] Core[™] i7-4770K CPU (3.50GHz)

NVIDIA GeForce GTX 670 (1344 CUDA Cores)

TRAINING TIMES FOR EPOCH = 5



CHALLENGES

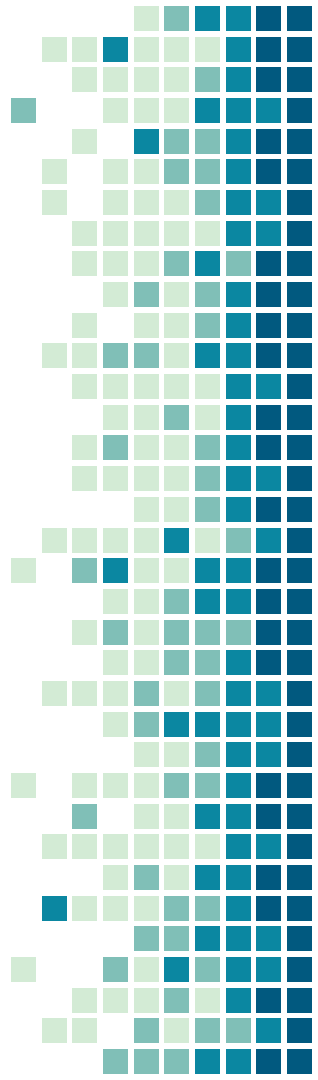
- Working with Visual Studio 2015
- Learning how to write Makefiles for GNU Make
- Learning how CUDA-enabled GPUs work
- Learning how to think like a C programmer
- ... oh, and those two awesome times I spent an entire weekend rewriting Neuronmancer



Acknowledgments

Special thanks to all my friends who encouraged me when I felt like giving up.

Also, to my professors for giving me the freedom to explore a topic I found interesting.




THANKS!

Any questions?

You can find me at:

 github.com/DrewHans555

 drewhans555@gmail.com

// end of presentation