# ECE 408
# Final Project Report

## Group Name

*ConvolutionallyUnsampledDiscreteAlgorithm*

## Members

Zain Paya

Aditya Bhargava

Arvind Kamal

# Contents

# MILESTONE 1: Getting Started

## 1.1 RUN THE BASELINE FORWARD PASS

```
* Running python m1.1.py
New Inference
Loading fashion-mnist data... done
Loading model... done
EvalMetric: {'accuracy': 0.8673}
```

*Fig 1.1: Output for CPU run on m1.1.py*

## 1.2 RUN THE BASELINE GPU IMPLEMENTATION

```
* Running /usr/bin/time python m1.2.py
New Inference
Loading fashion-mnist data... done
Loading model...[00:35:03] src/operator/././cudnn_algoreg-inl.h:112: Running performance tests
to find the best convolution algorithm, this can take a while... (setting env variable MXNET_CU
DNN_AUTOTUNE_DEFAULT to 0 to disable)
done
EvalMetric: {'accuracy': 0.8673}
1.38user 0.68system 0:01.89elapsed 108%CPU (0avgtext+0avgdata 904628maxresident)k
0inputs+0outputs (0major+155169minor)pagefaults 0swaps
```

*Fig 1.2: Output for GPU run on m1.2.py with time outputs (real, user, sys)*

## 1.3 GENERATE A NVPROF PROFILE

- Top 2 time consuming kernels: *implicit_convolve_sgemm, pooling_fw_4d_kernel*



*Fig 1.3: Output table for running times of all kernels via nvprof*

# MILESTONE 2:  New CPU Layer in MXNet

## 2.1 ADD CPU FORWARD IMPLEMENTATION



*Fig 2.1: Output for CPU Forward Convolution on m2.1.py (default)*



*Fig 2.2: Output for CPU Forward Convolution on m2.1.py with time outputs (elapsed, user, system)*



*Fig 2.3: Comparative output for ece-408-low of size 10000*

*Fig 2.4: Comparative output for ece-408-high of size 10000*

# TEAM MEMBERS' ROLES

- **Zain Paya**
  - Implemented Forward Convolution with an extra layer for Milestone 2

- **Aditya Bhargava**
  - Implemented Forward Convolution with an extra layer for Milestone 2

- **Arvind Kamal**
  - Implemented Forward Convolution with an extra layer for Milestone 2