

# TriForce Updates and Classification with New Data

---

**Matt Zhang** | Apr. 13, 2018



**ILLINOIS**  
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN



# TriForce Updates

Added ability to split data into train-validation-test. If validation split percentage is not specified, the training set will also be the validation set.

TriForce code significantly refactored. Training histories for all tools now stored internally in a single data structure.

Accuracies split by signal/background, and training history stored by batch/epoch are now stored for all tools.

Plotting scripts now part of default analysis tool.

Additional protections added when reading in options file. TriForce warns when options are missing or mis-set. Certain unspecified options are now given default values.

# Required Options

**samplePath** - folders to find samples in

**classPdgID** - pdg IDs of sample classes

**trainRatio** - what percent of the data files to use for the training set

**nEpochs** - end training after this many epochs

**relativeDeltaLossThreshold** - end training early if delta test loss over an epoch falls below this threshold

**relativeDeltaLossNumber** - ...or if it falls below relativeDeltaLossThreshold for this many sets of batches

**batchSize** - number of samples in a batch

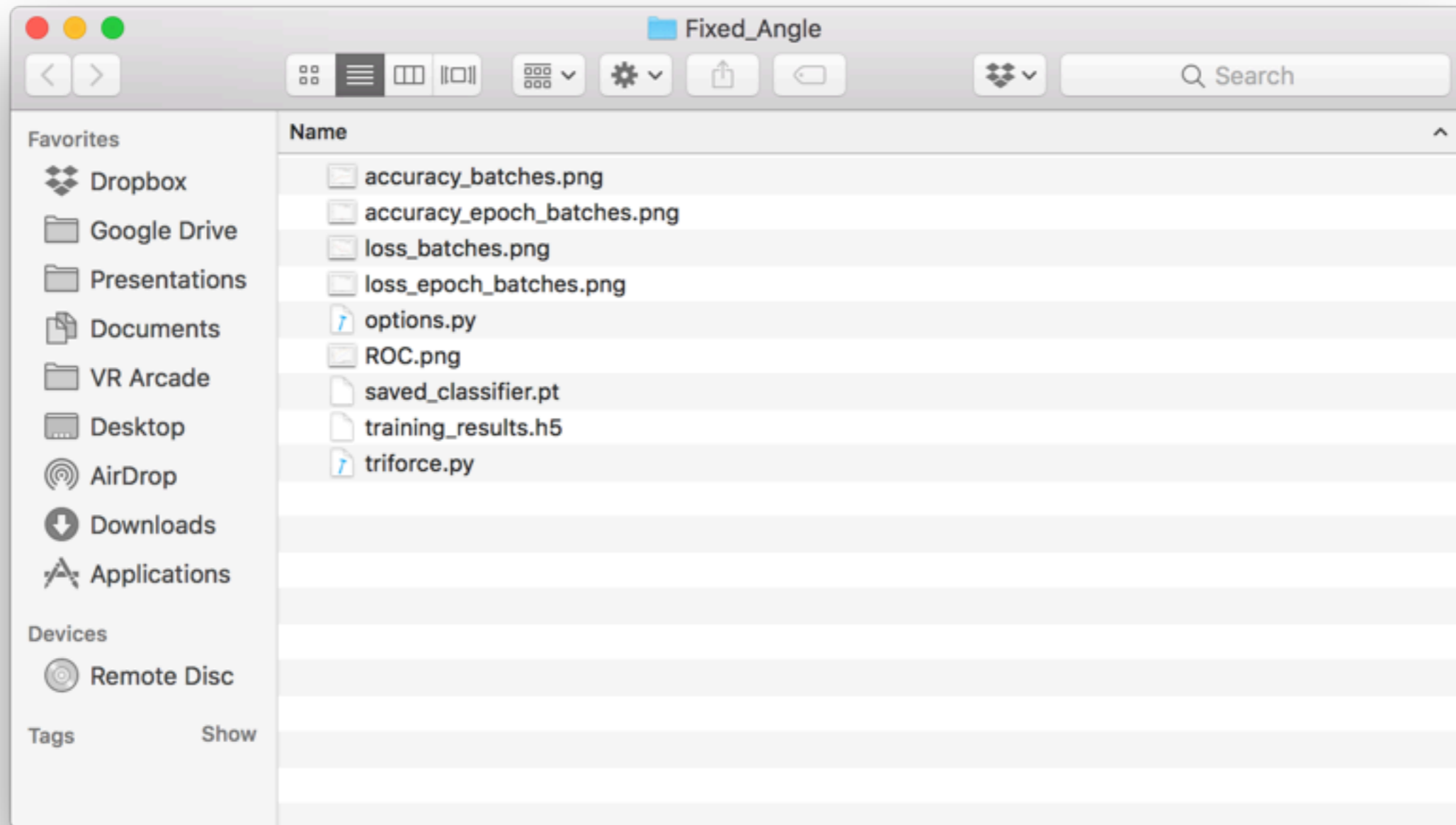
**saveModelEveryNEpochs** - save trained model every N epochs (set to 0 to only save at end of training)

**outPath** - folder to save outputs in

# Optional Options

- importGPU (False)** - whether to use the GPU packages on the CalTech clusters
- eventsPerFile (10,000)** - how many events are in each file
- nTrainMax (-1)** - maximum number of files in the training set (-1 for no max)
- nValidationMax (-1)** - maximum number of files in the validation set (-1 for no max)
- nTestMax (-1)** - maximum number of files in the test set (-1 for no max)
- validationRatio (-1)** - what percent of the data files to use for the validation set (-1 if validation set = test set)
- nWorkers (0)** - number of workers to use in the data loader
- calculate\_loss\_per\_n\_batches (20)** - calculate loss every N batches
- test\_loss\_eval\_max\_n\_batches (10)** - the maximum number of test batches to use when calculating test loss
- earlyStopping (True)** - set to false to turn off early stopping during training

# Outputs



# Running on New Data

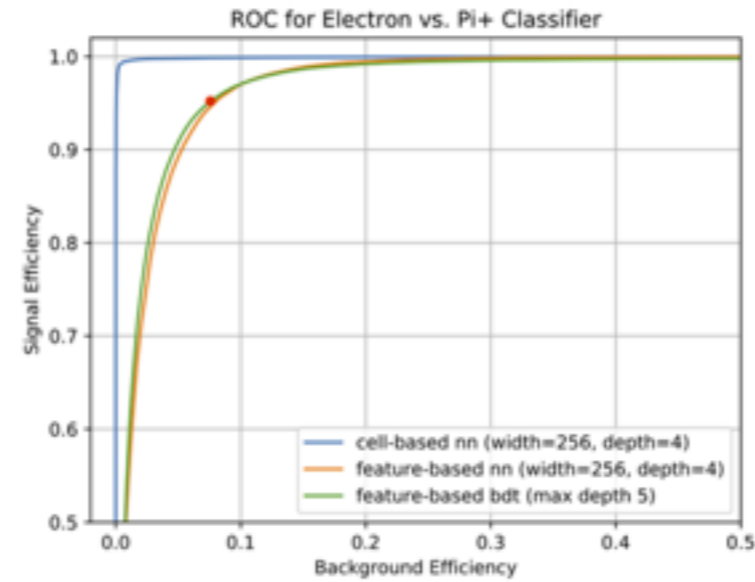
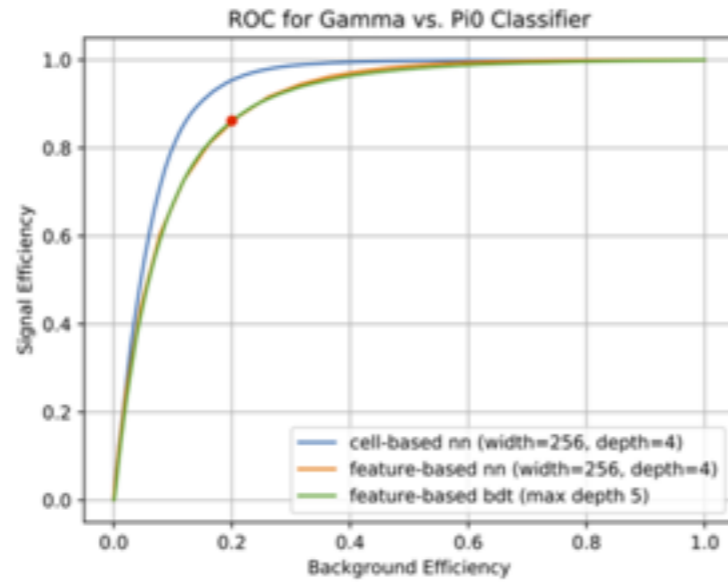
Only ECAL used.

Size set to 51x51x51.

NIPS architecture used (aside from side of input layer) for comparison.

Running on gamma vs. pi0 for both fixed-angle and variable-angle data (using the same architecture).

# Paper Results

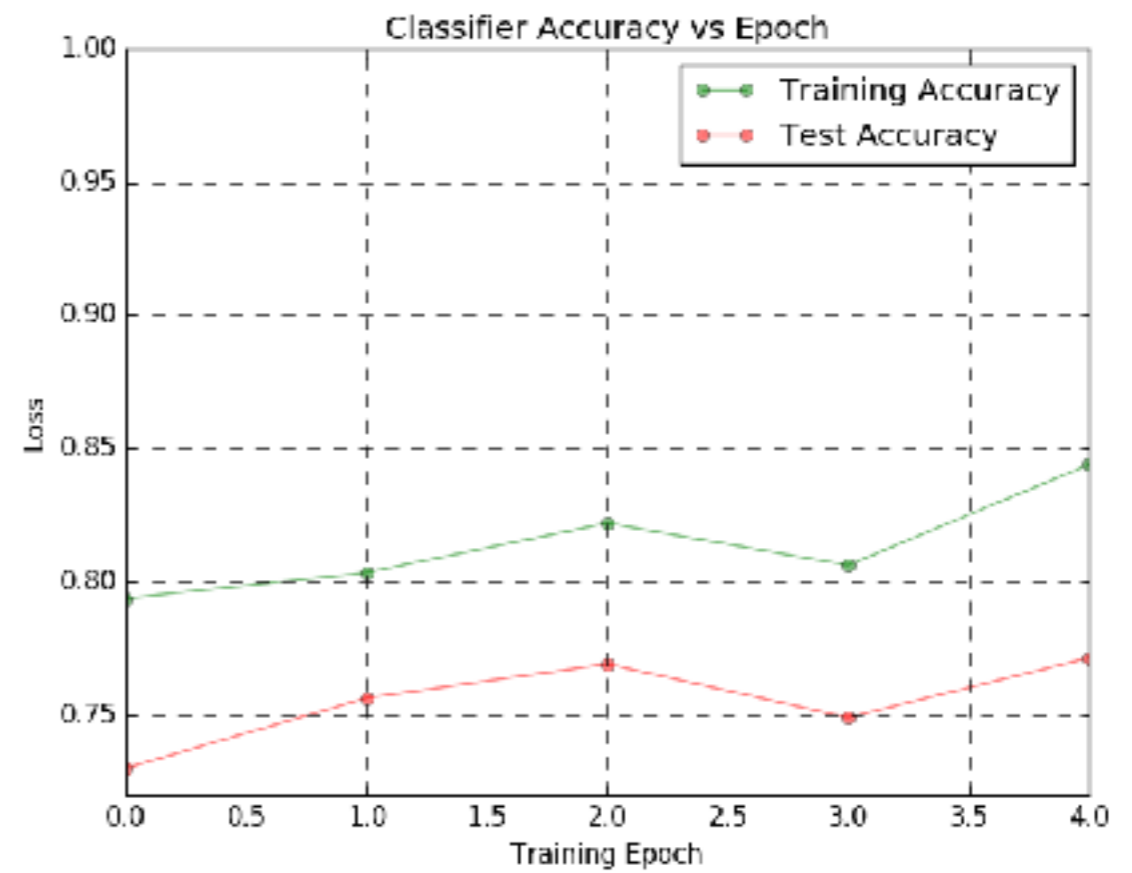
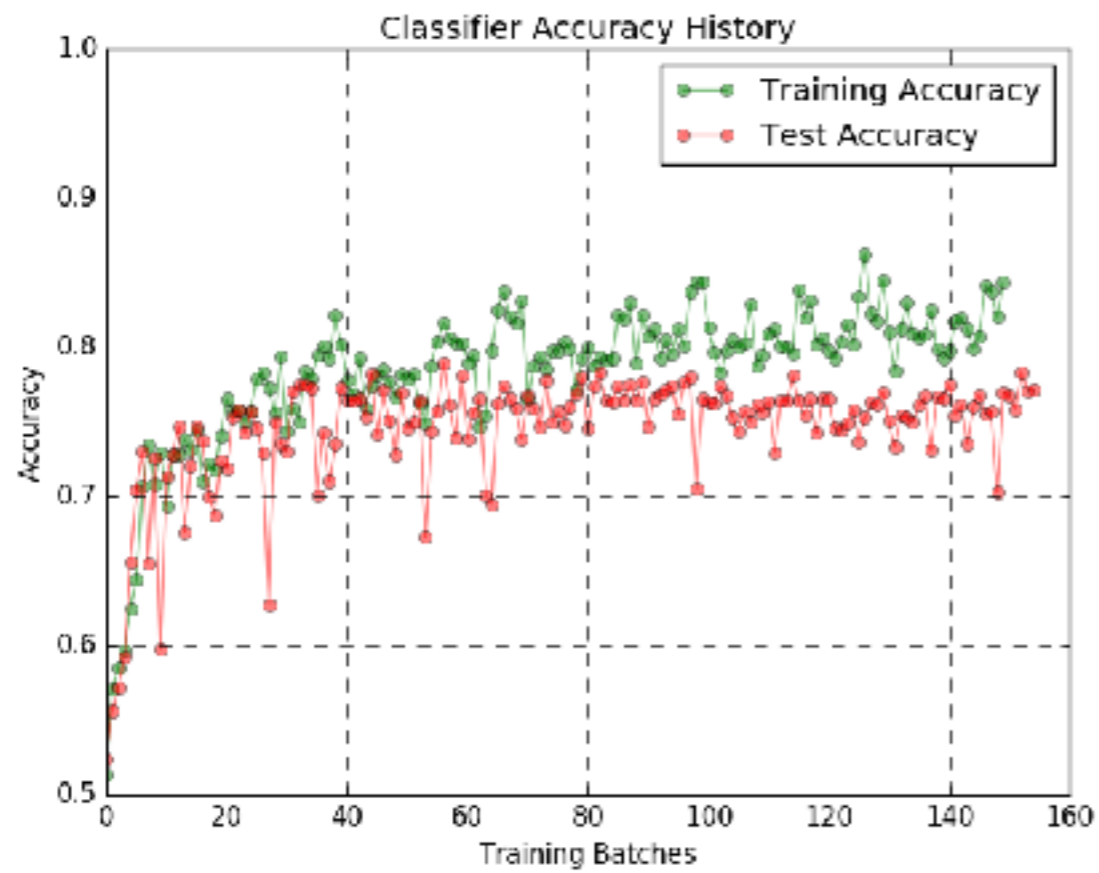


Model	$\gamma$ vs. $\pi^0$				$e$ vs. $\pi$			
	acc.	AUC	$\Delta\epsilon_{sig}$	$\Delta R_{bkg}$	acc.	AUC	$\Delta\epsilon_{sig}$	$\Delta R_{bkg}$
BDT	83.1%	89.8%	-	-	93.8%	98.0%	-	-
DNN (features)	82.8%	90.2%	0.9%	0.95	93.6%	98.0%	-0.1%	0.95
DNN (cells)	87.2%	93.5%	9.4%	1.63	99.4%	99.9%	4.9%	151

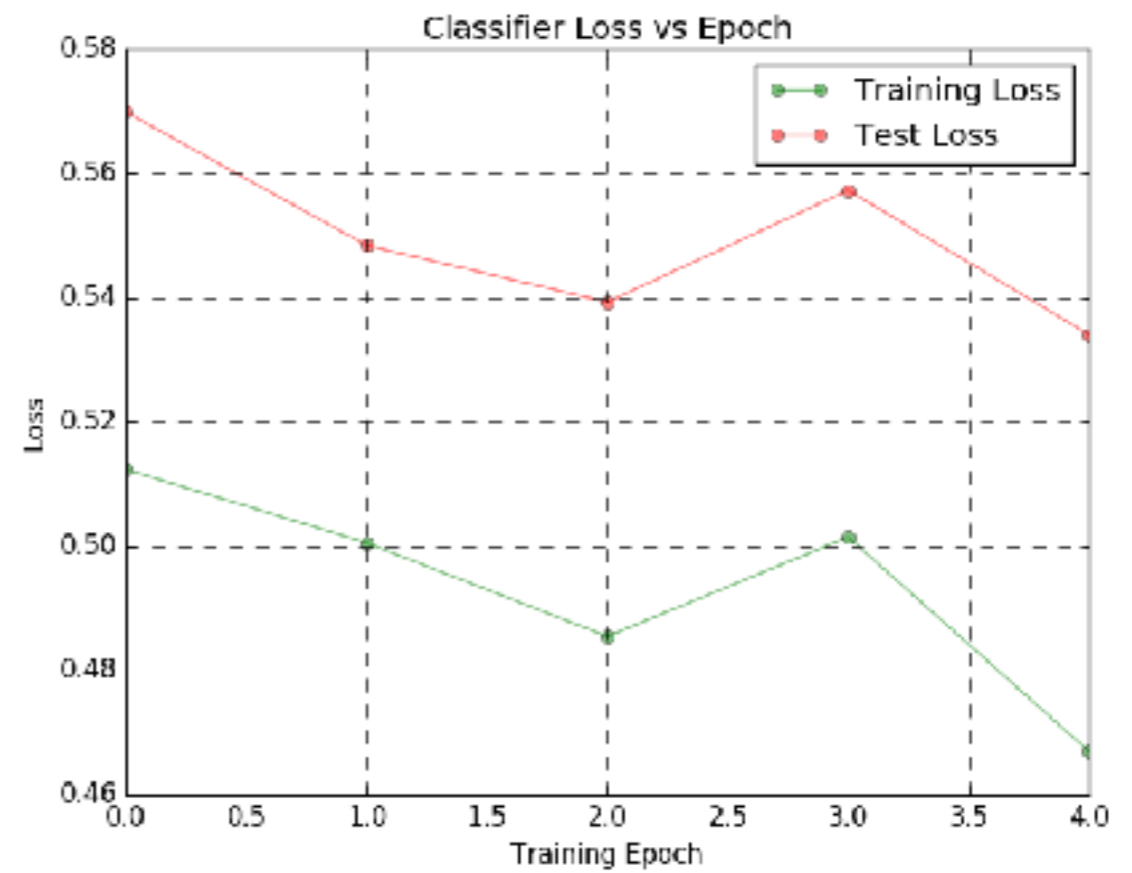
# Fixed Angle



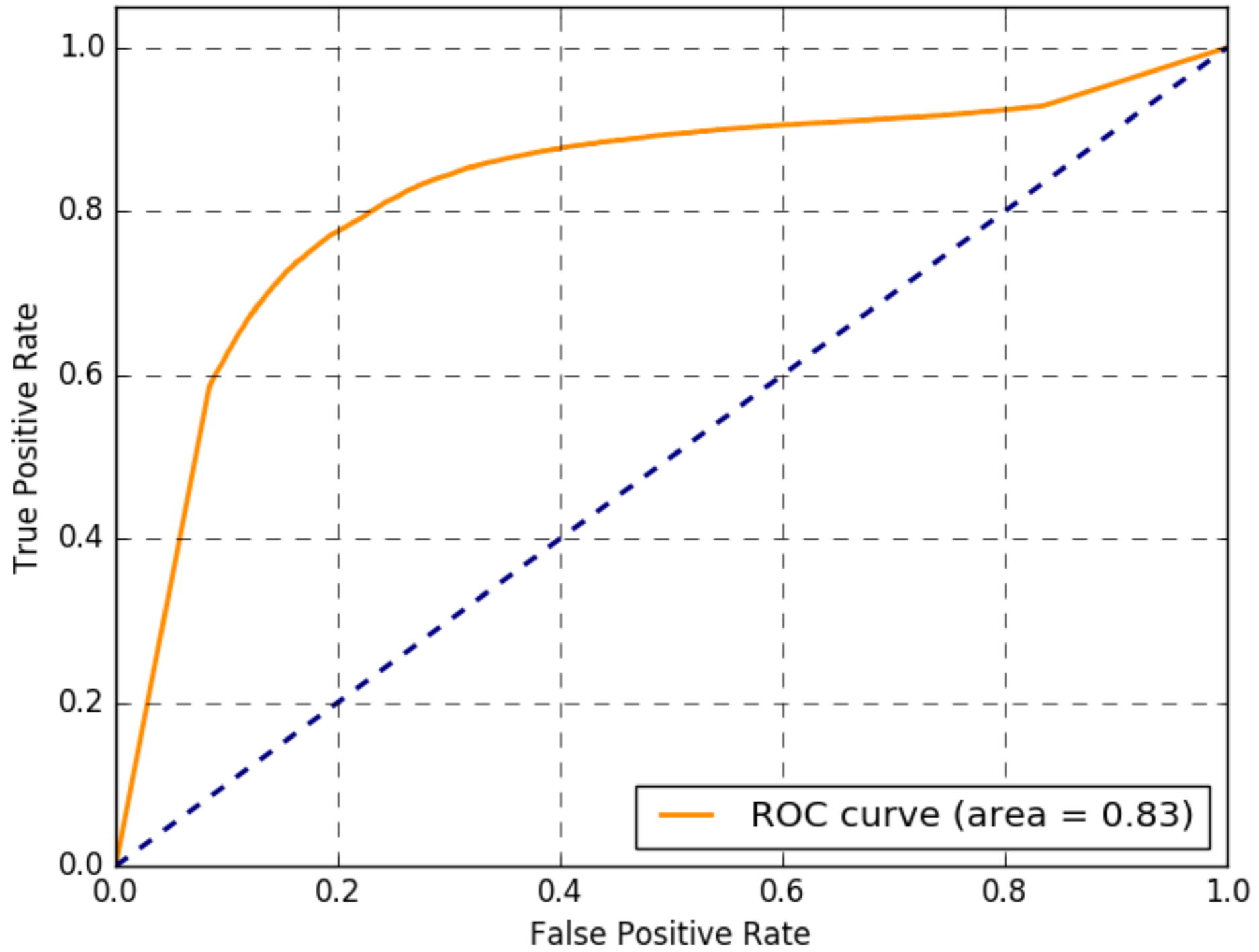
# Accuracy



# Loss



ROC Curve for Classification



# Takeaway

Net has a harder time training with the larger data window.

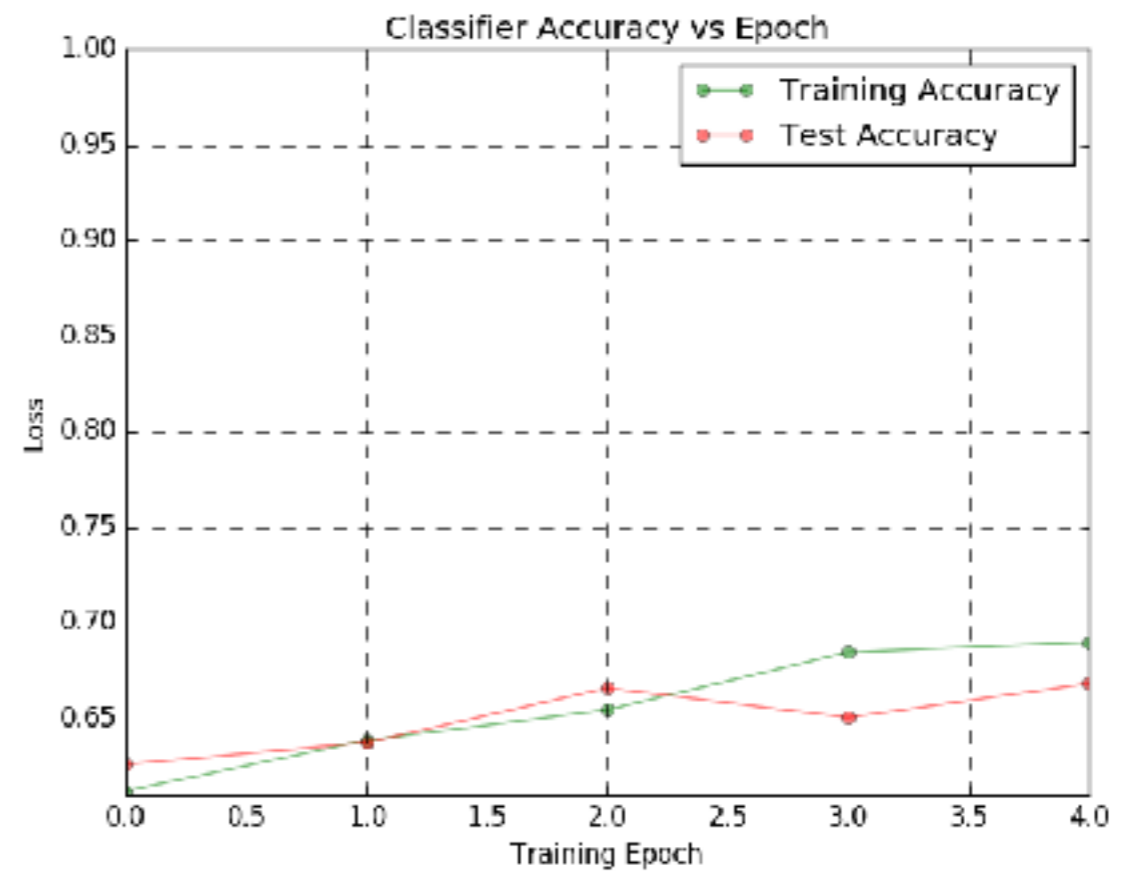
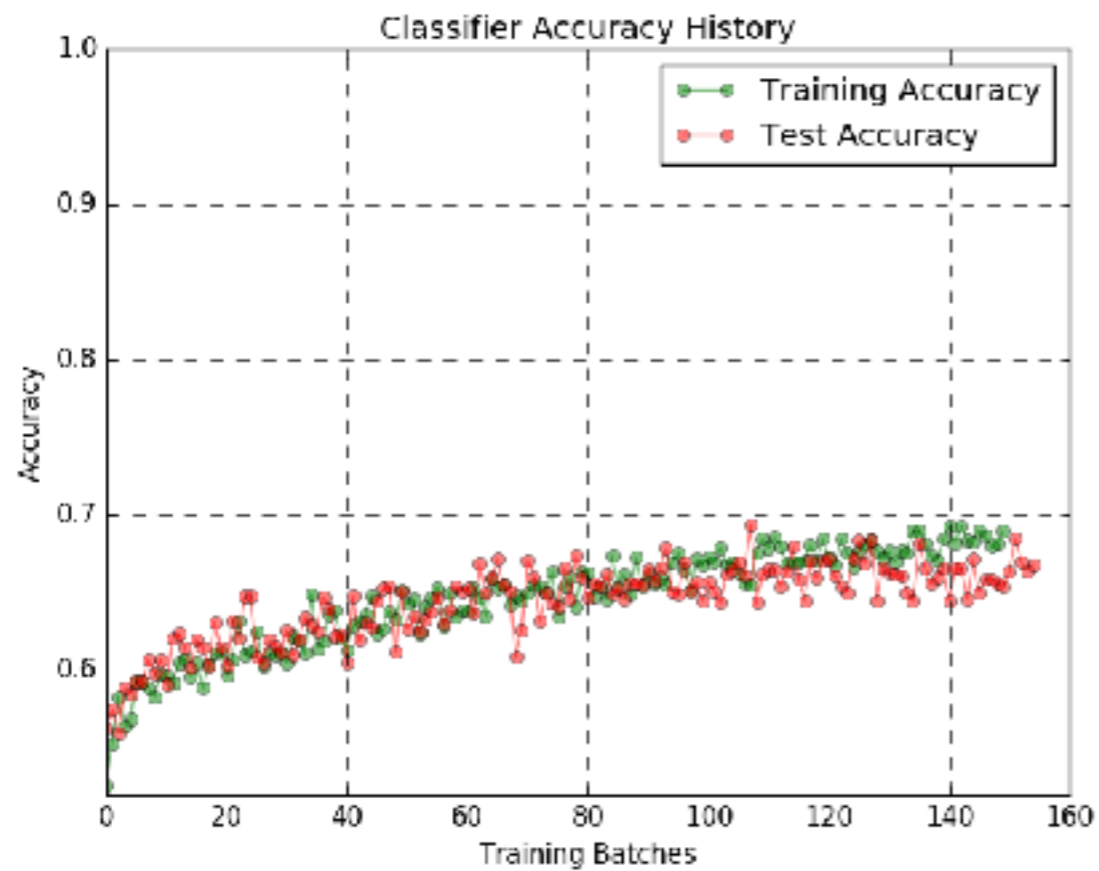
Results are not too bad, but number of neurons in hidden layers should probably be increased.

Will run hyperparameter scan on Blue Waters.

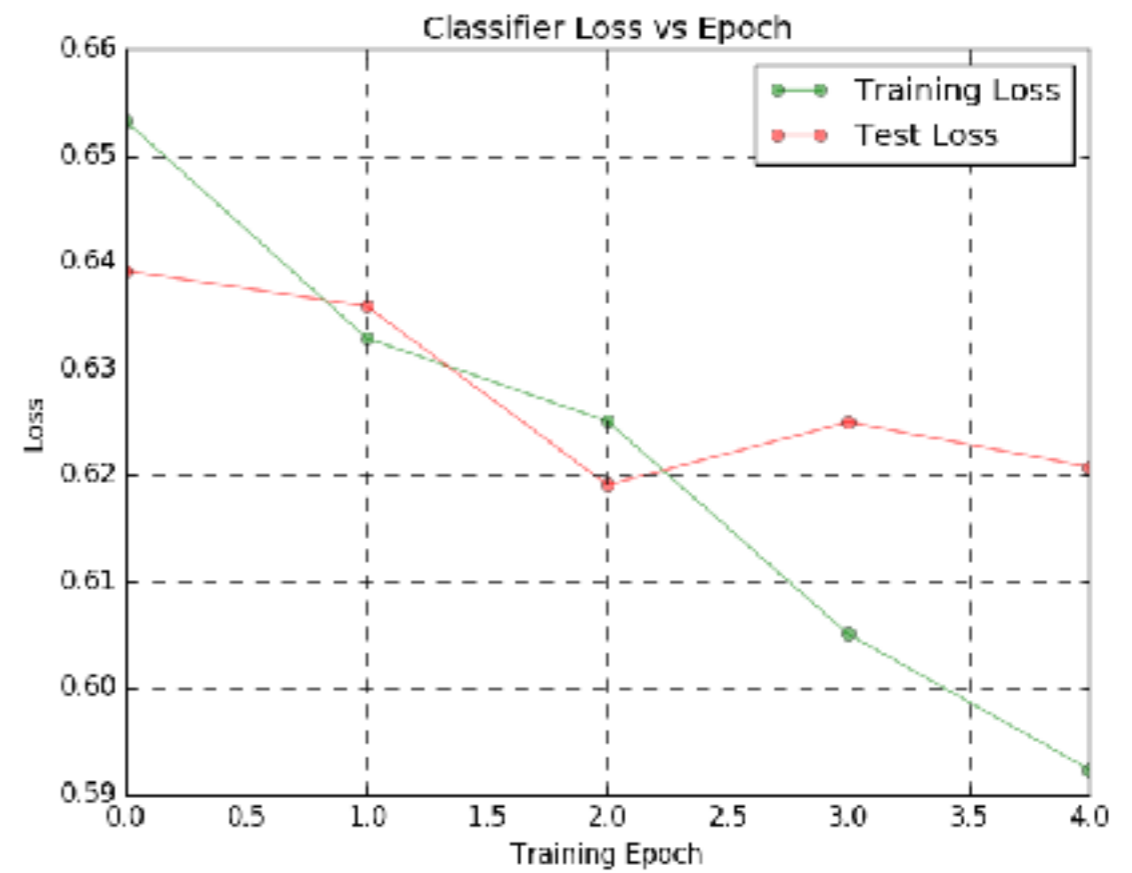
Will also train using GoogLeNet.

# Variable Angle

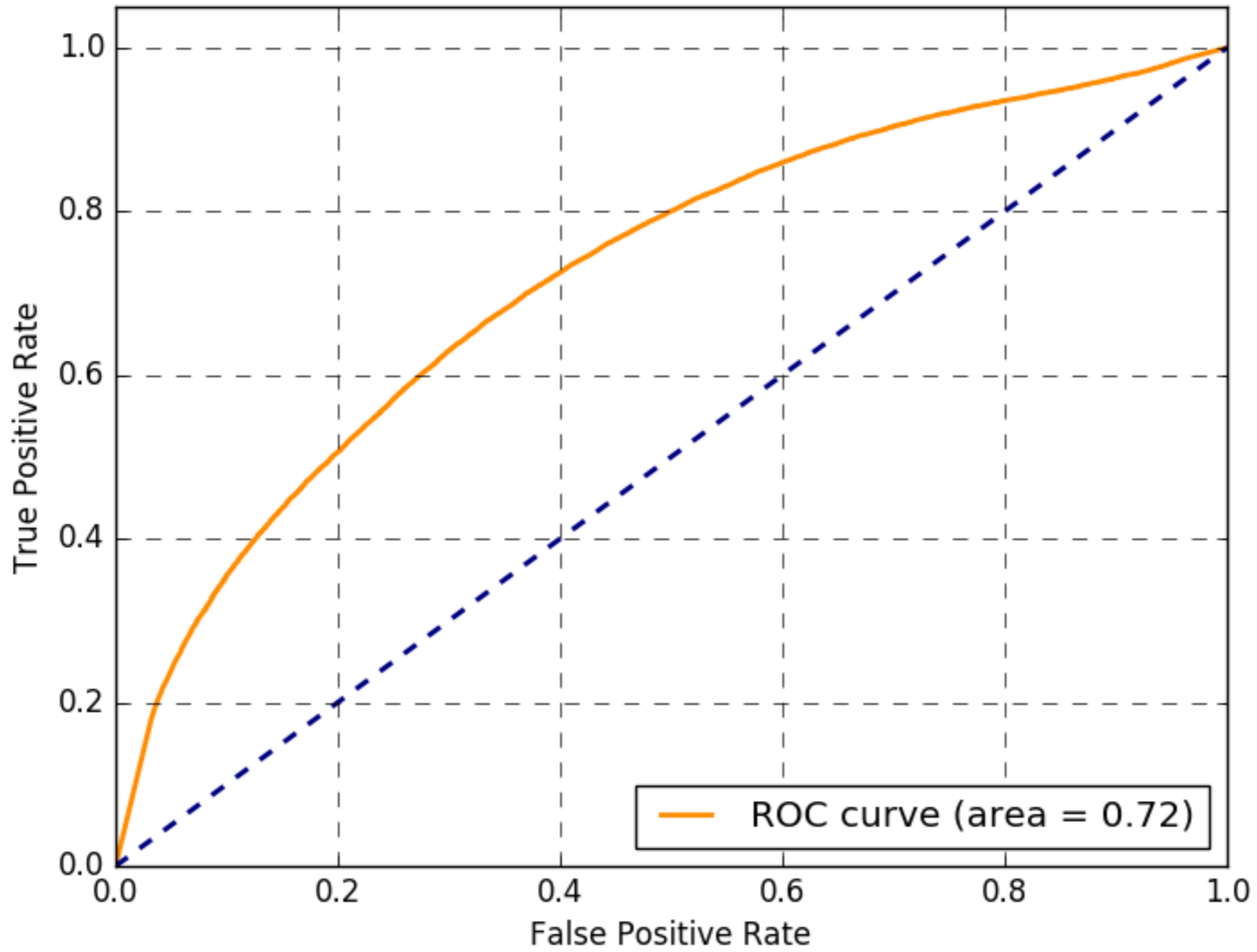
# Accuracy



# Loss



ROC Curve for Classification





# Takeaway

I am incredibly surprised that we were able to do this well on the variable-angle data.

Simply adding a single extra input for the particle eta may be enough to push up the accuracy. Will try this.