

ML4DQM Updates

Pritam Palit, Shamik Ghosh, Subir Sarkar

Saha Institute of Nuclear Physics, India

11/05/2020

Topics to report

Anomaly detection with Autoencoder for 2017 data

- As a beginner, I first used Autoencoder architecture as given in Francesco's Code for :
chargelInner_PXLayer_1
 - Checked the distribution of Reconstruction Error (mseTop10)
 - Checked the Precision-Recall curve , F1 score and corresponding best Threshold
 - Checked the Confusion Matrix
- Then, I checked the sequence of the consecutive bins and decided to use LSTM Autoencoder
 - Compared the results.

Standard Autoencoder

- The total dataset is split into Train - Validation - Test (72 - 8 - 20)%
 - We Trained only on Good runs and lumis of Train dataset
 - Good runs and lumis selected from Golden json file
 - Test is done on the whole test dataset (including both good and bad)
- Standard DNN based autoencoder with tanh activation in all layers :
input - 10 - 3 - 10 - output
- Optimizer = adam , loss = mseTop10 (Mean of Top 10 highest error)

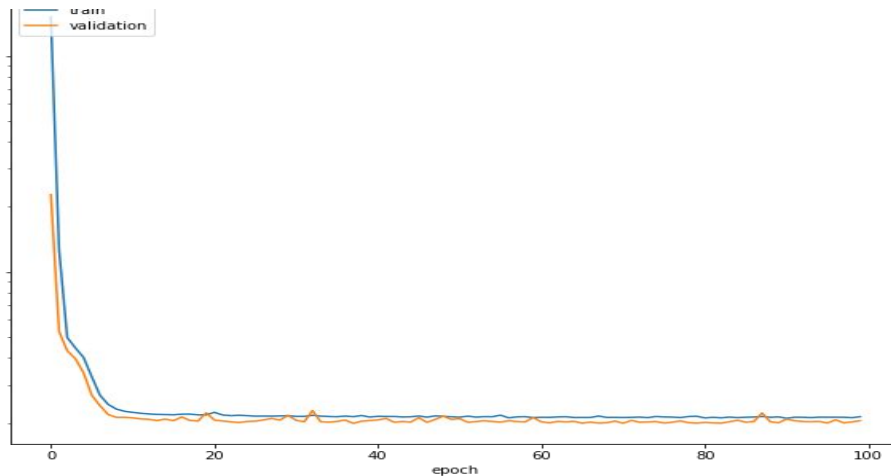
Model: "model_1"

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	(None, 100)	0
dense_1 (Dense)	(None, 10)	1010
dense_2 (Dense)	(None, 3)	33
dense_3 (Dense)	(None, 10)	40
dense_4 (Dense)	(None, 100)	1100

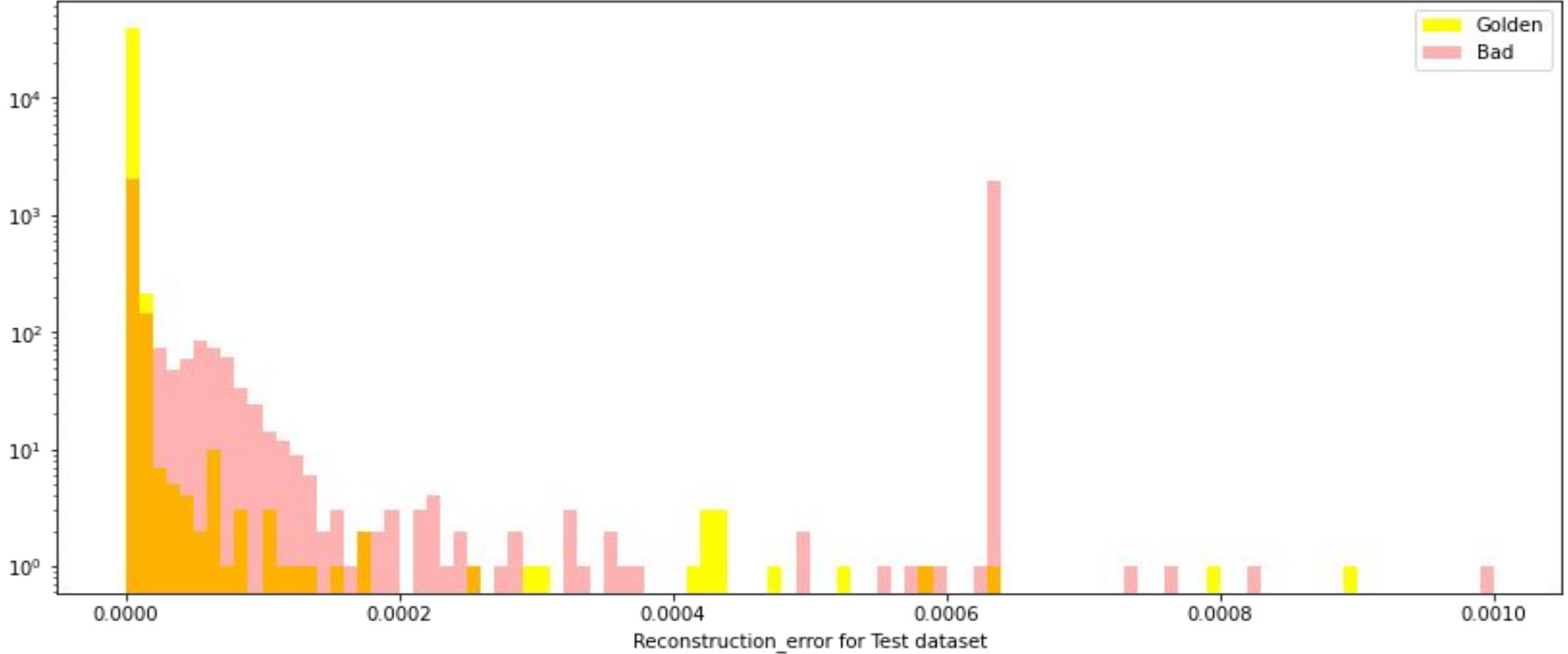
Total params: 2,183

Trainable params: 2,183

Non-trainable params: 0

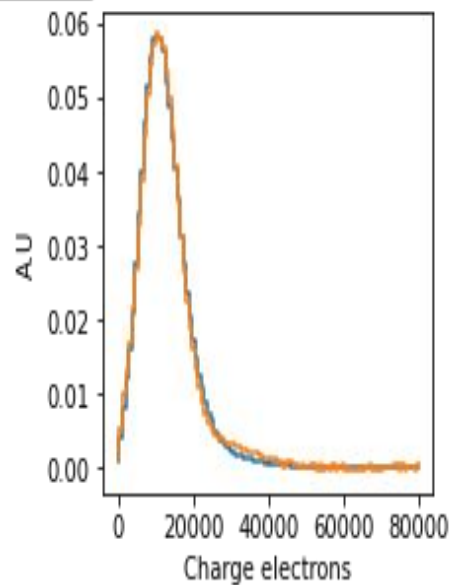


Reconstruction error For Test Data

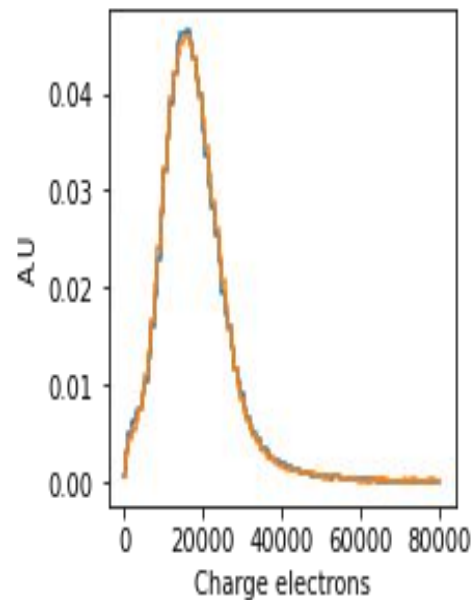


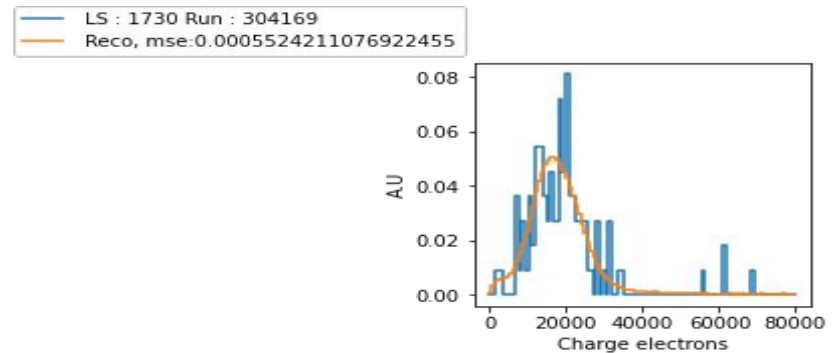
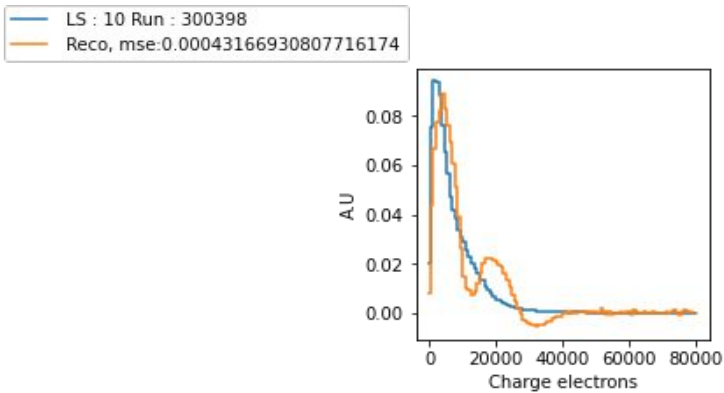
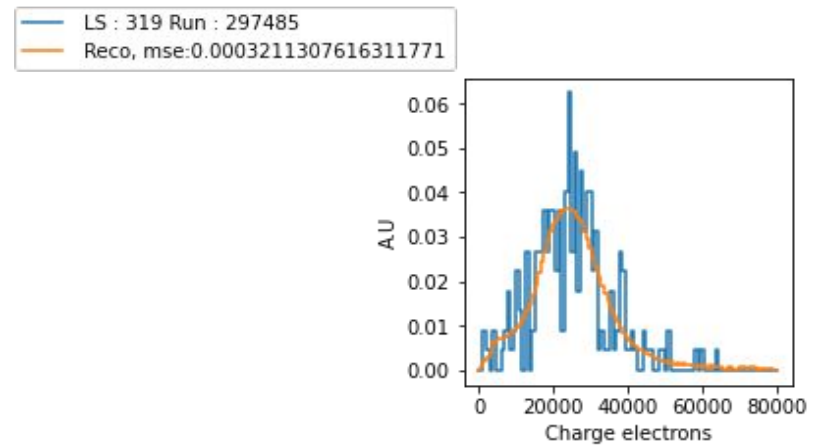
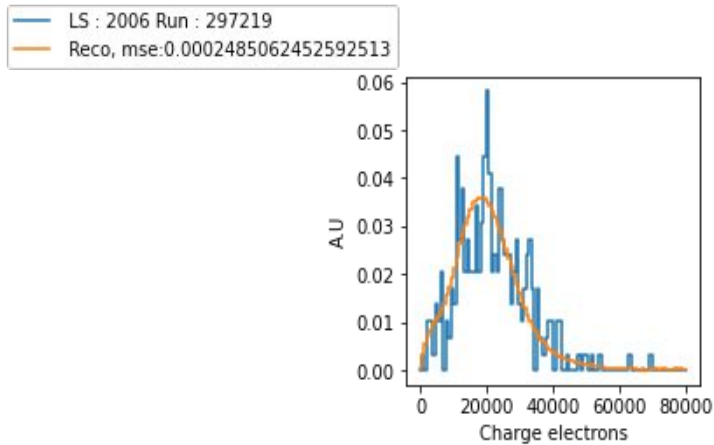
Error < 0.0002

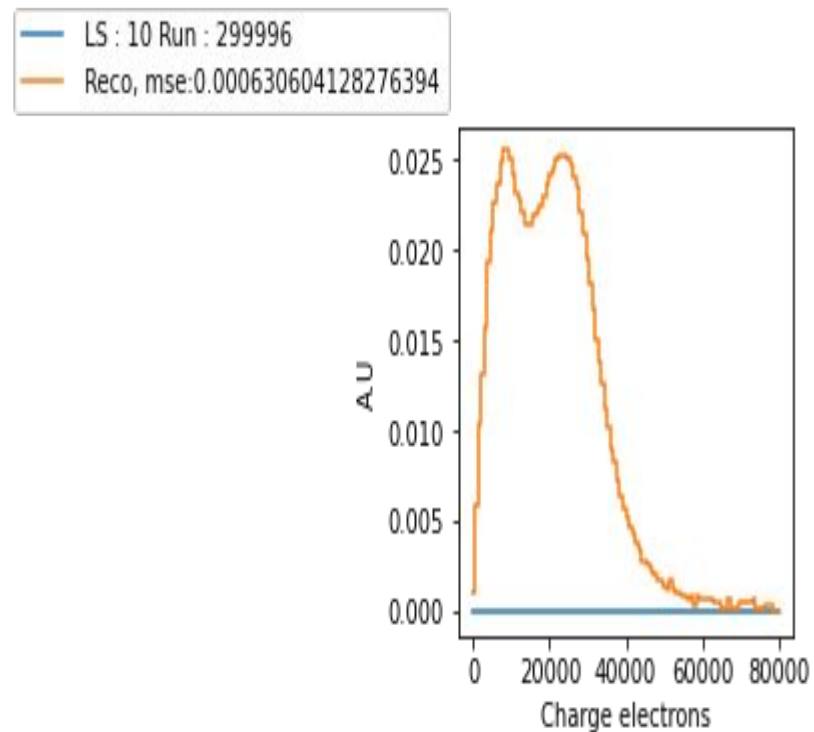
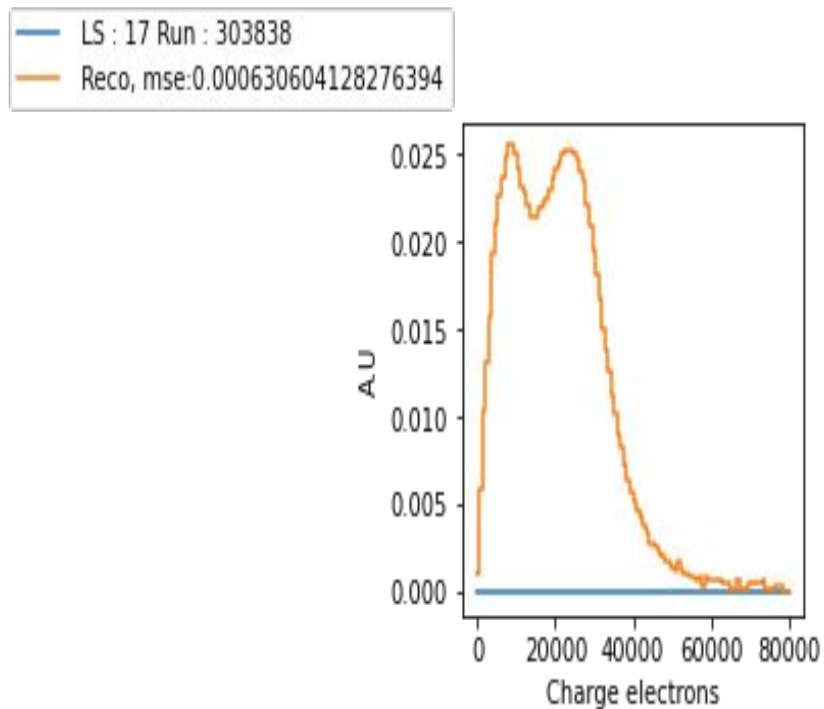
— LS : 449 Run : 301998
— Reco, mse:2.8570818161932563e-06



— LS : 1325 Run : 304366
— Reco, mse:1.276485583392423e-06

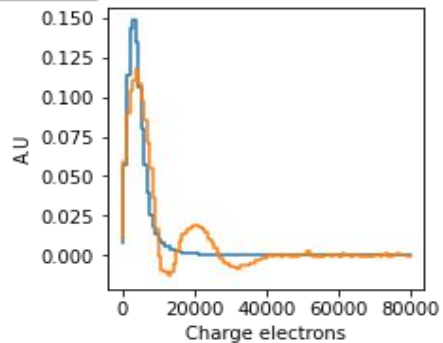


$0.0002 < \text{Error} < 0.0006$ 

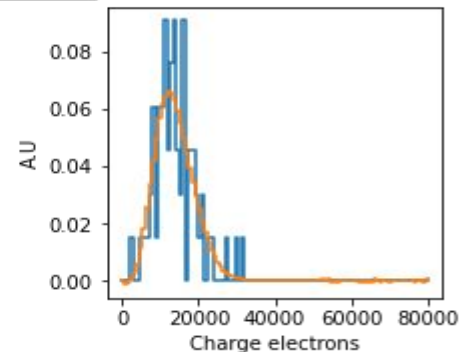
$0.0006 < \text{Error} < 0.0007$ 

$0.0007 < \text{Error} < 0.0009$

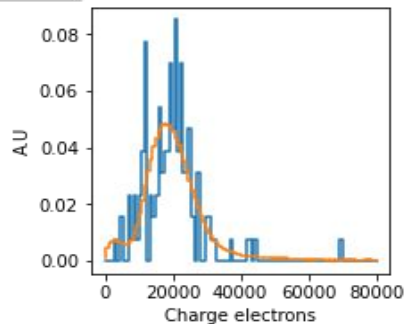
— LS : 115 Run : 301086
— Reco, mse:0.0007505902325838943



— LS : 22 Run : 301393
— Reco, mse:0.0008204580054481982



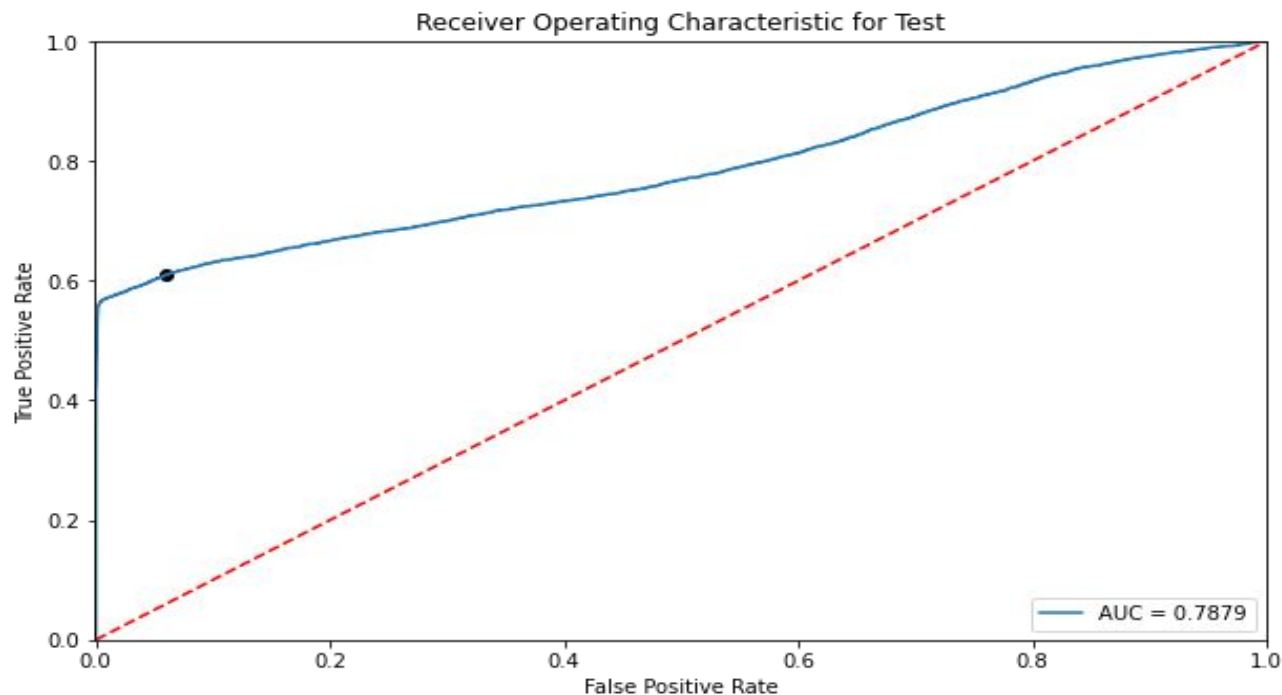
— LS : 175 Run : 299593
— Reco, mse:0.0008987443451908298



ROC Curve

$$\text{true positive rate} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

$$\text{false positive rate} = \frac{\text{false positives}}{\text{false positives} + \text{true negatives}}$$



Here :

Positive = Anomaly

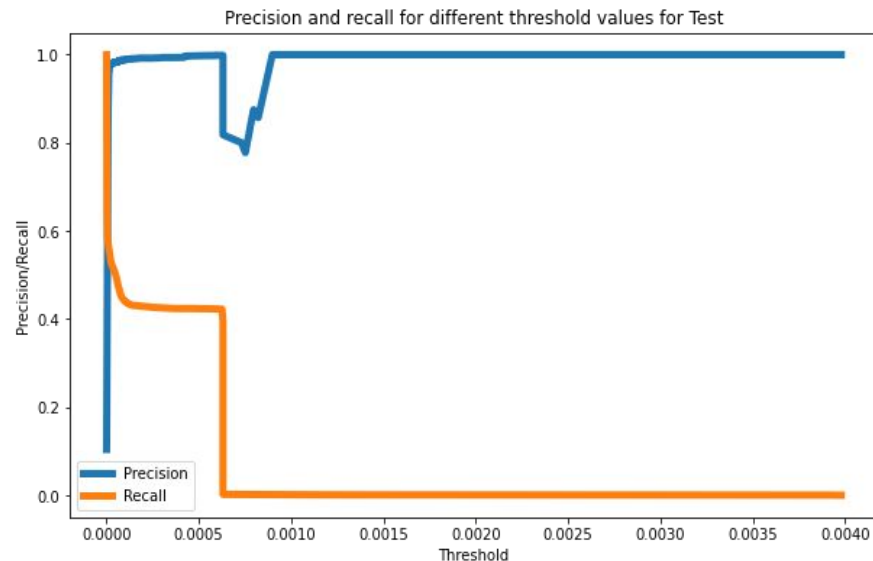
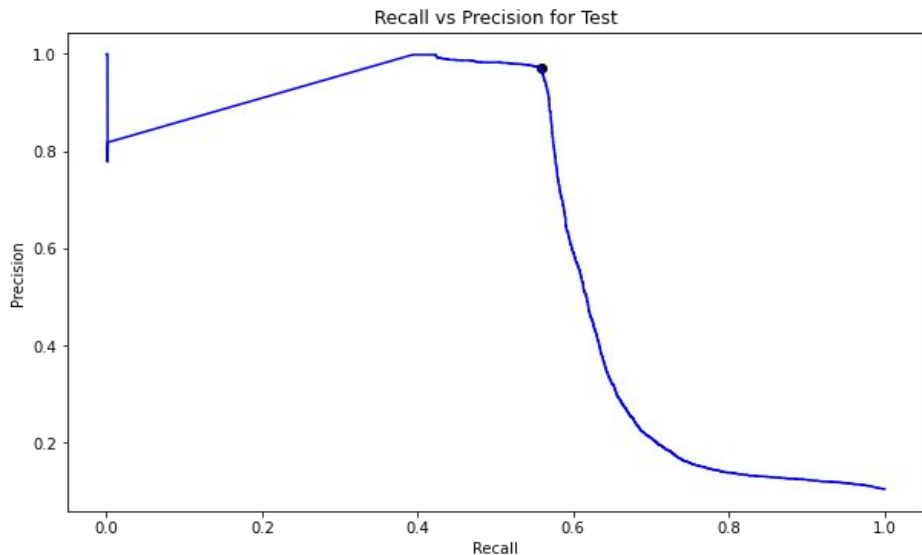
Negative = Good

Precision - Recall (PR) Curve for Imbalanced Dataset

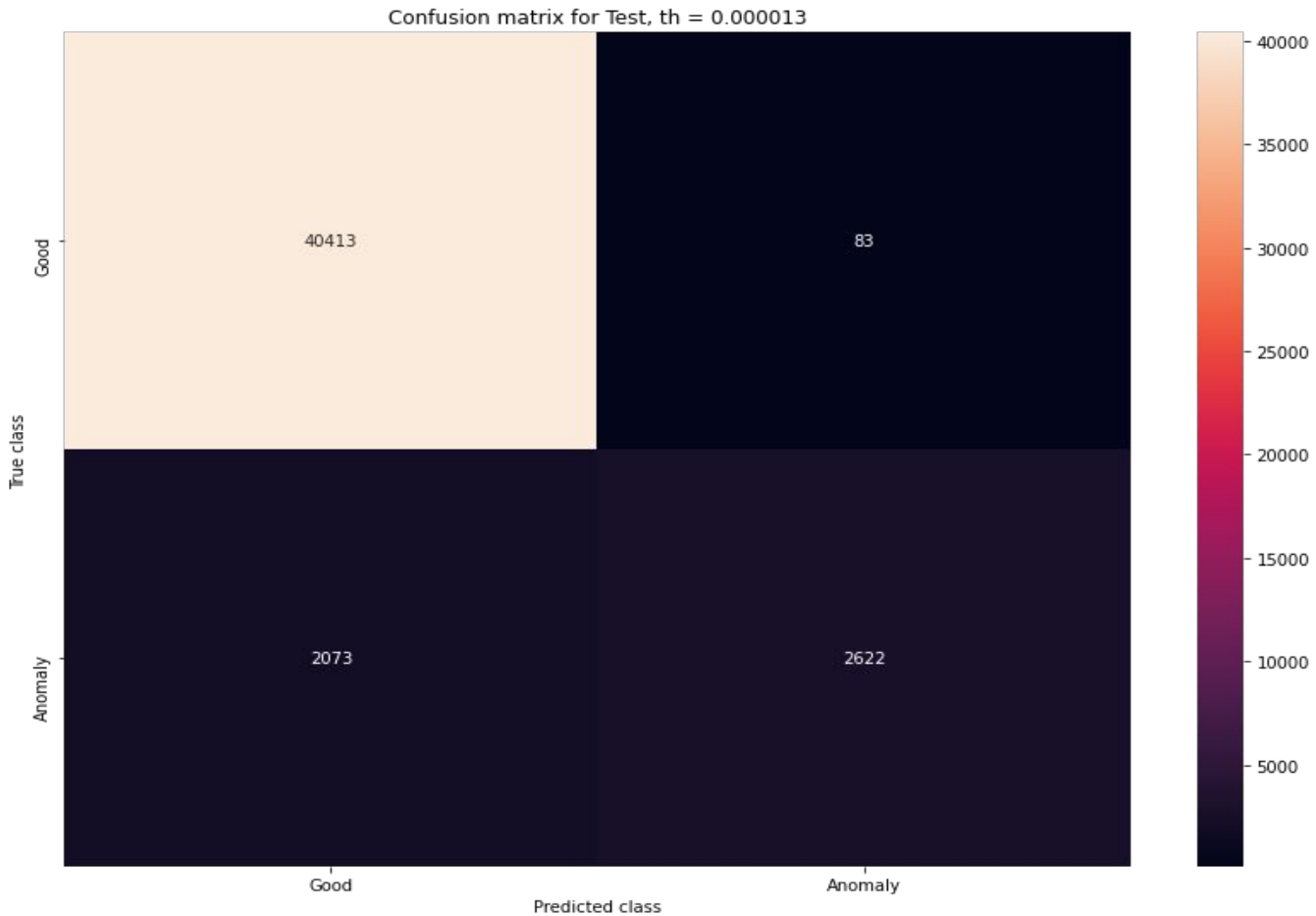
$$\text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

- PR curve for anomaly detection in a imbalanced dataset
 - Less interested in true negative (as well as fpr of ROC curve)
- `fscore = (2 * precision * recall) / (precision + recall)`
- **Best Threshold=0.000013 => F-Score=0.709**



Confusion Matrix



Want to reduce both F.P & F.N

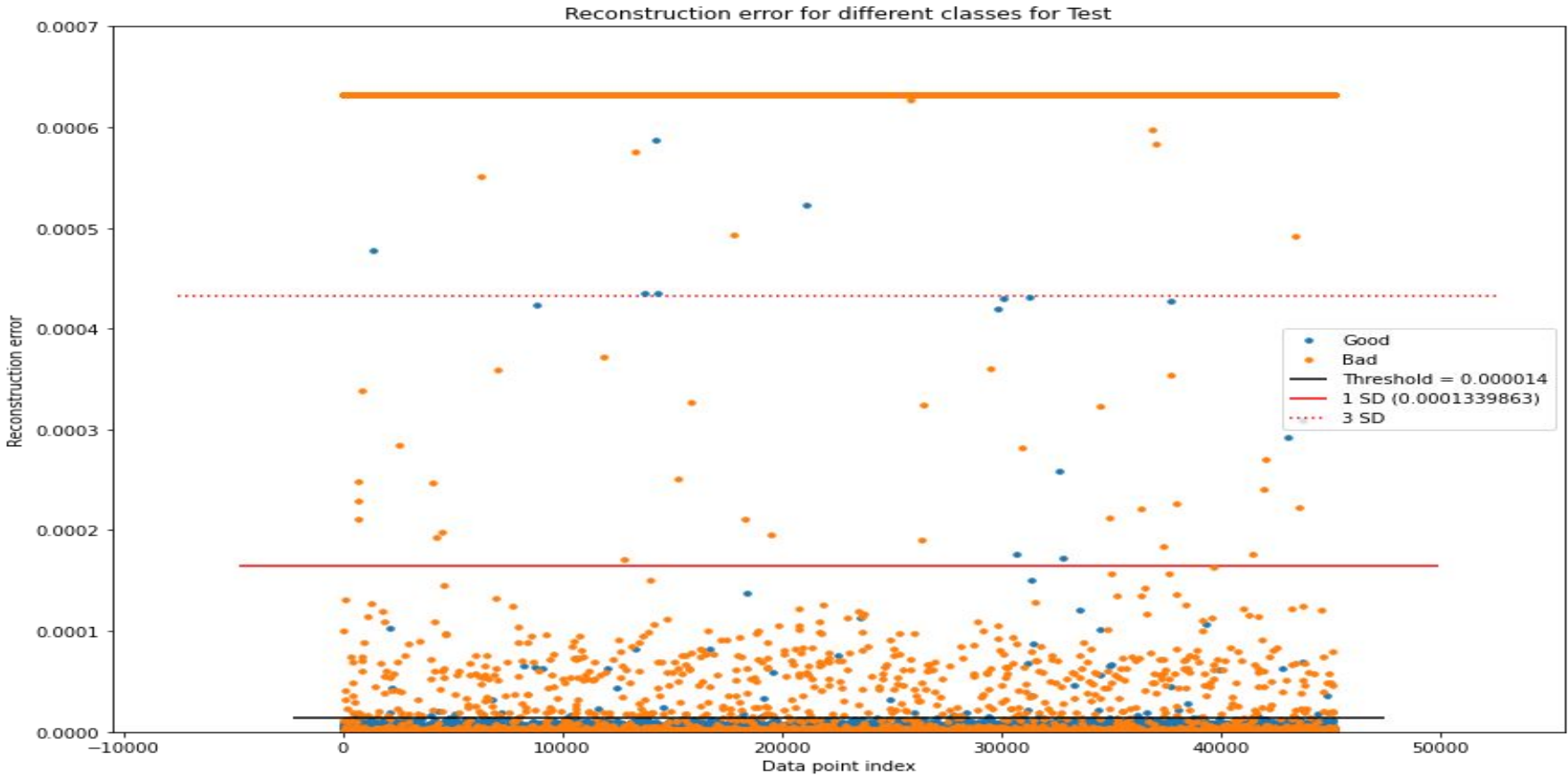
F.P = 0.21 % of N

F.N. = 44.2 % of P

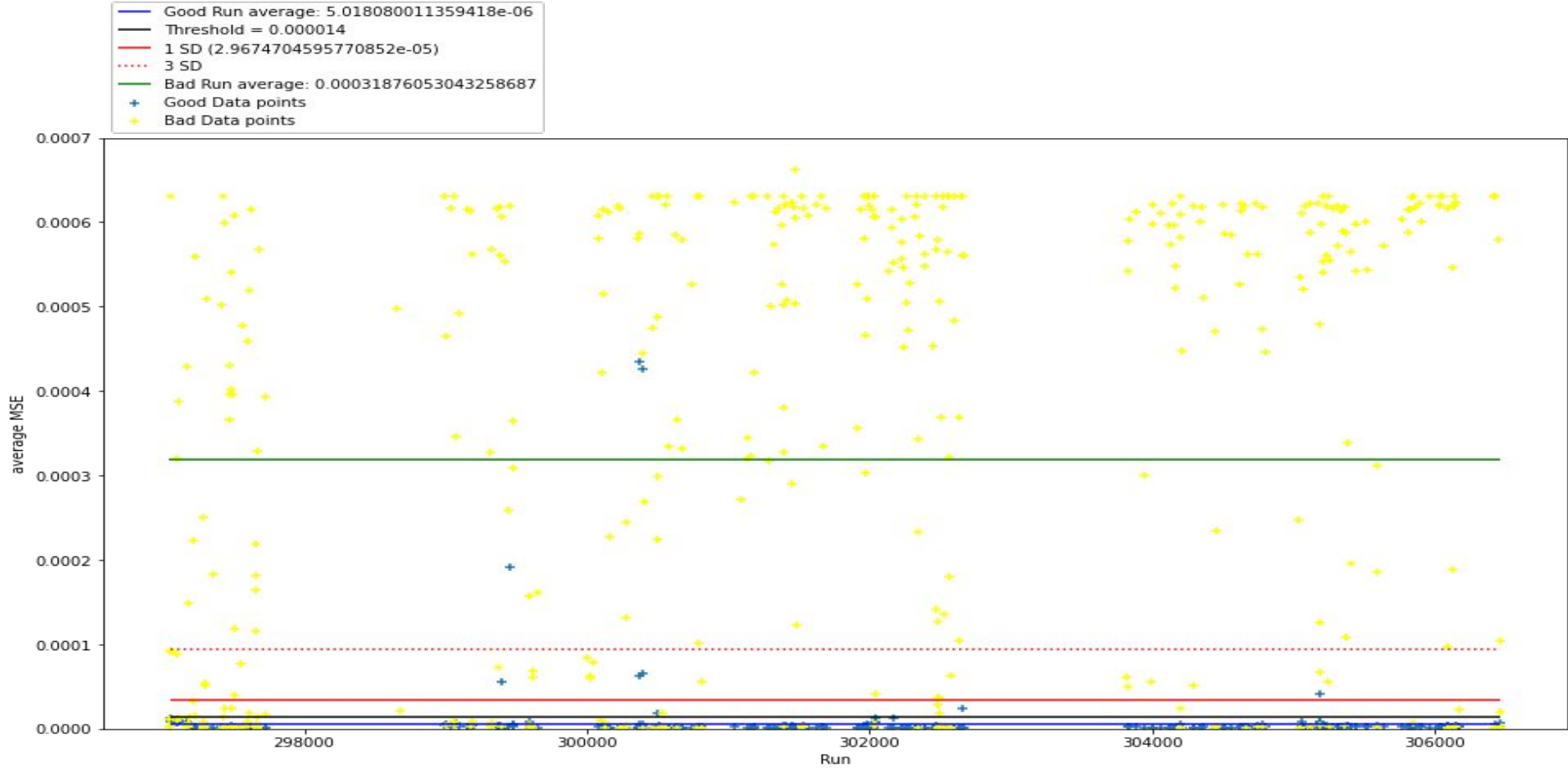
Suggestions?

Should we call it Positive or Bad lumi at all?

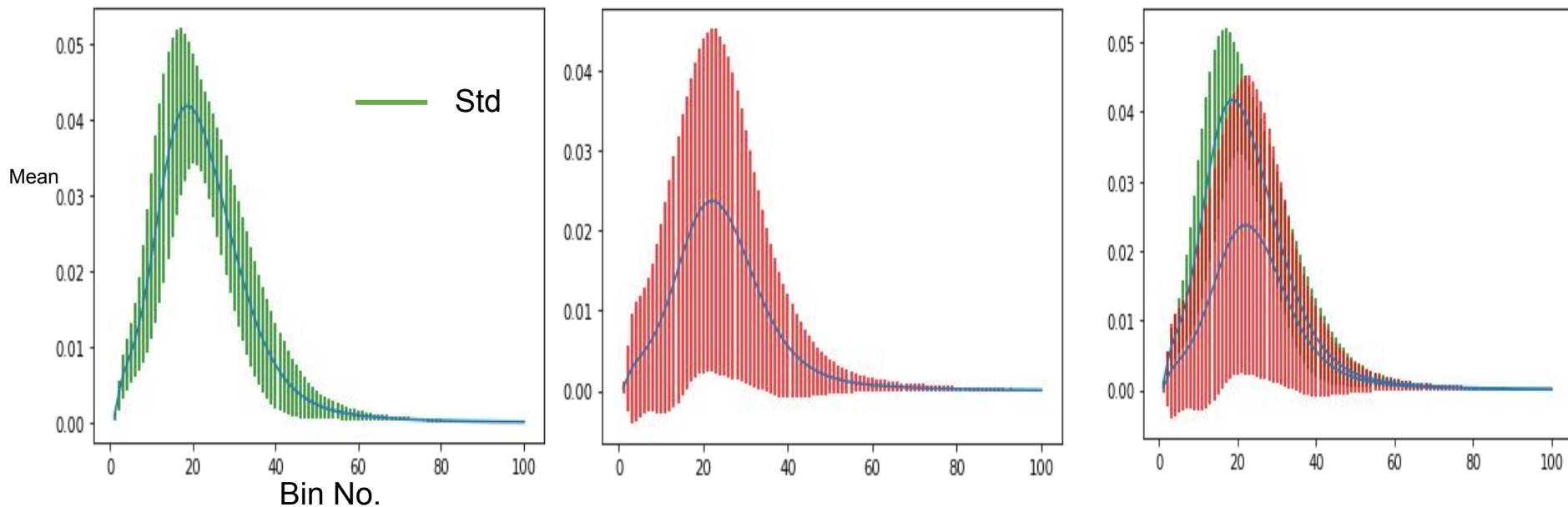
Reconstruction error for test dataset (Good & Bad)



Reconstruction error Global Trend over runs



Mean and Std of **Good** and **Bad** data over all lumis



Good : Possibility of some correlation over the sequence of bins

Bad : random fluctuation of std => Poor correlation

=> We try to investigate the sequential information of bins by LSTM AE

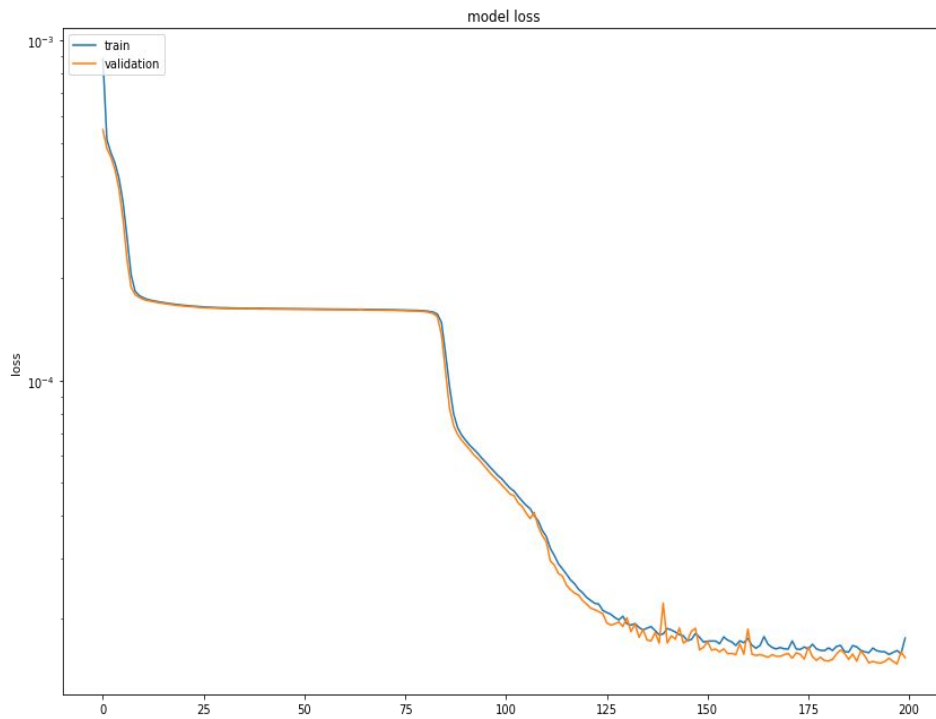
LSTM Autoencoder

- Using individual bin as a step in a sequence, we tried to encode histogram as sequence in LSTM (timestep = 100, no. of features = 1)
- Tried to keep no. of parameters as close possible to the standard AE to compare

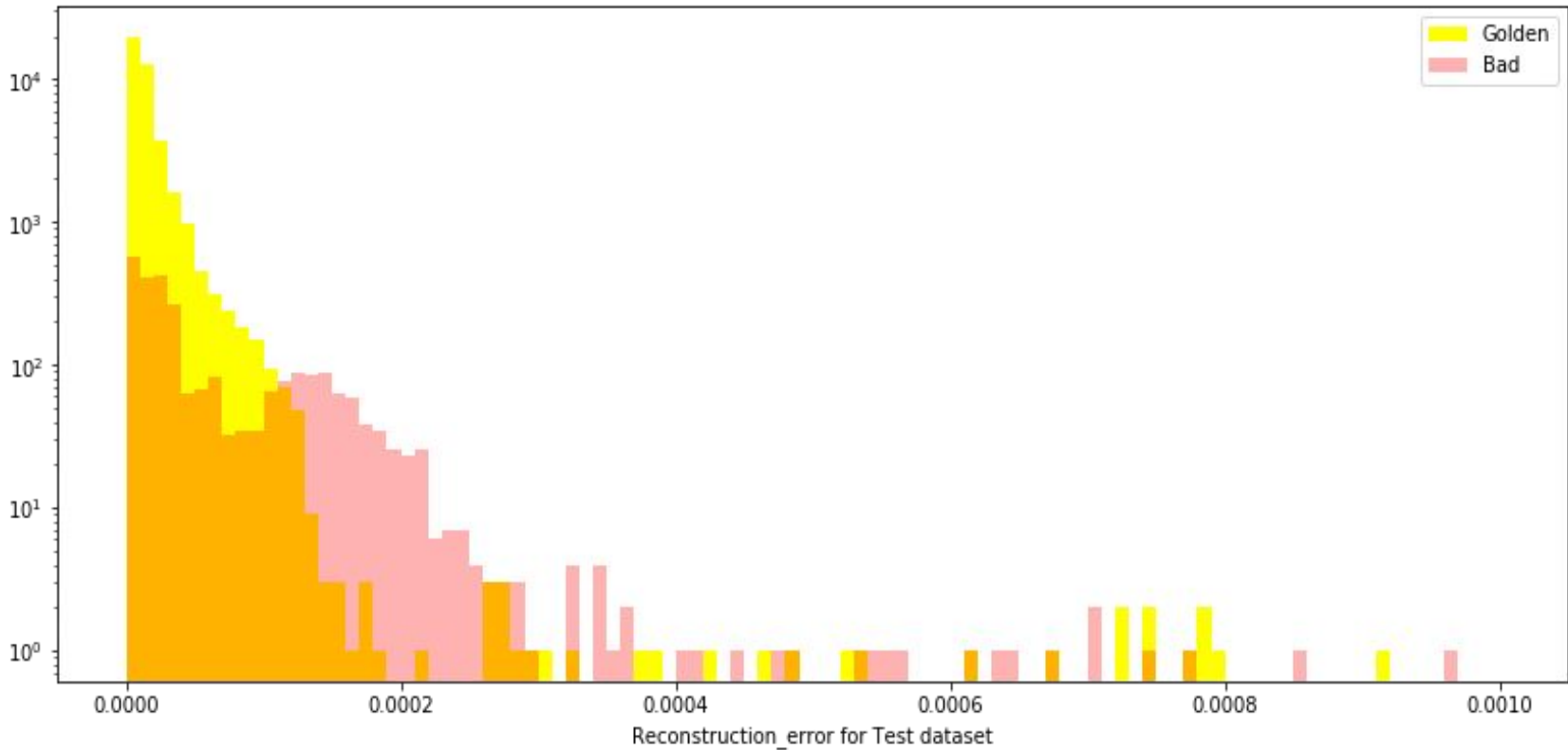
```
(145591, 100, 1)
Model: "sequential_1"
```

Layer (type)	Output Shape	Param #
lstm_1 (LSTM)	(None, 100, 10)	480
lstm_2 (LSTM)	(None, 3)	168
repeat_vector_1 (RepeatVecto	(None, 100, 3)	0
lstm_3 (LSTM)	(None, 100, 3)	84
lstm_4 (LSTM)	(None, 100, 10)	560
time_distributed_1 (TimeDist	(None, 100, 1)	11

```
Total params: 1,303
Trainable params: 1,303
Non-trainable params: 0
```

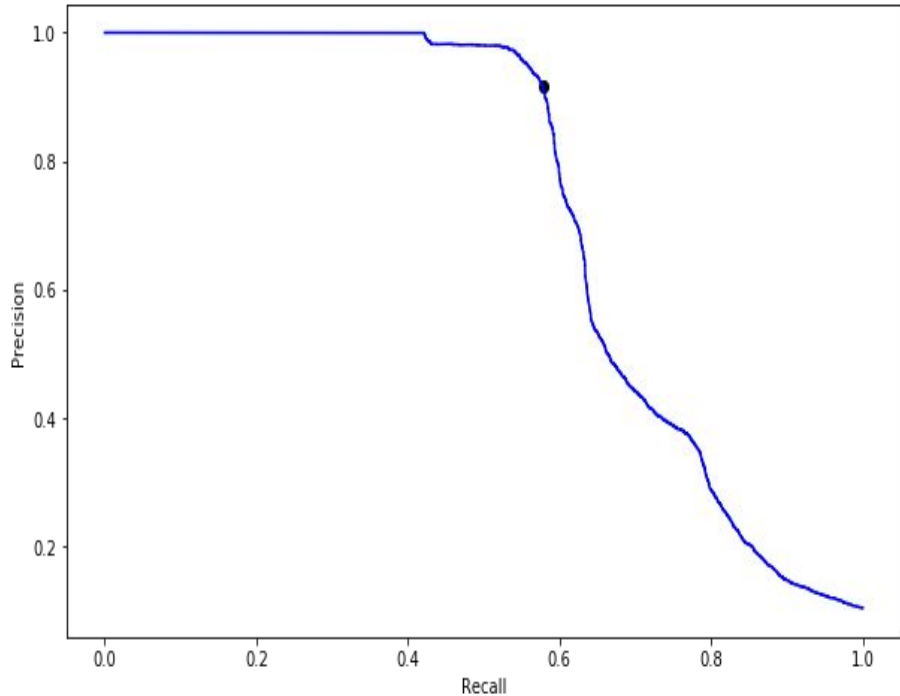


Reconstruction error (mseTop10)

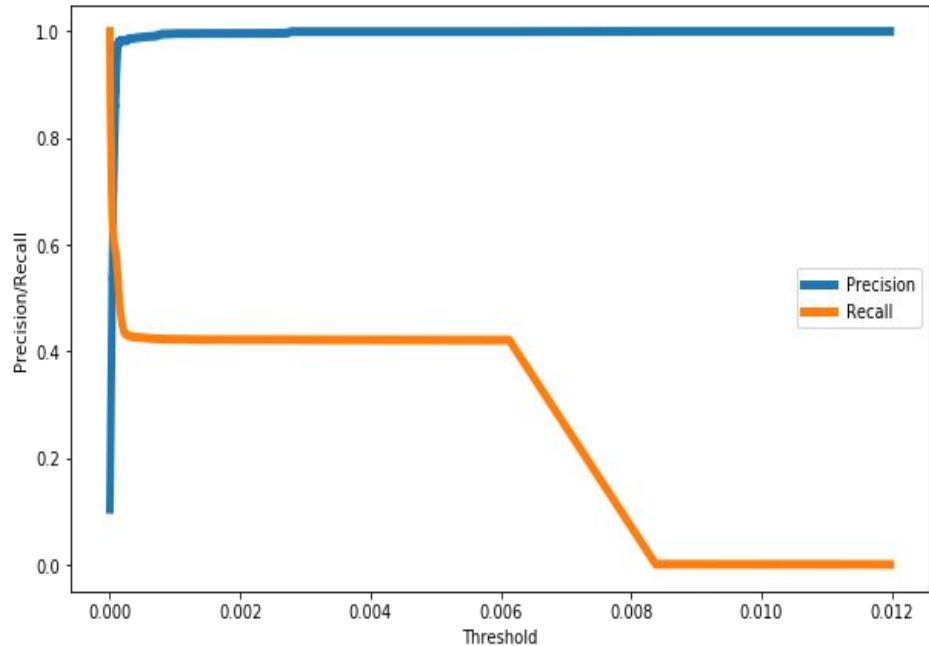


17 PR Curve

Recall vs Precision for Test

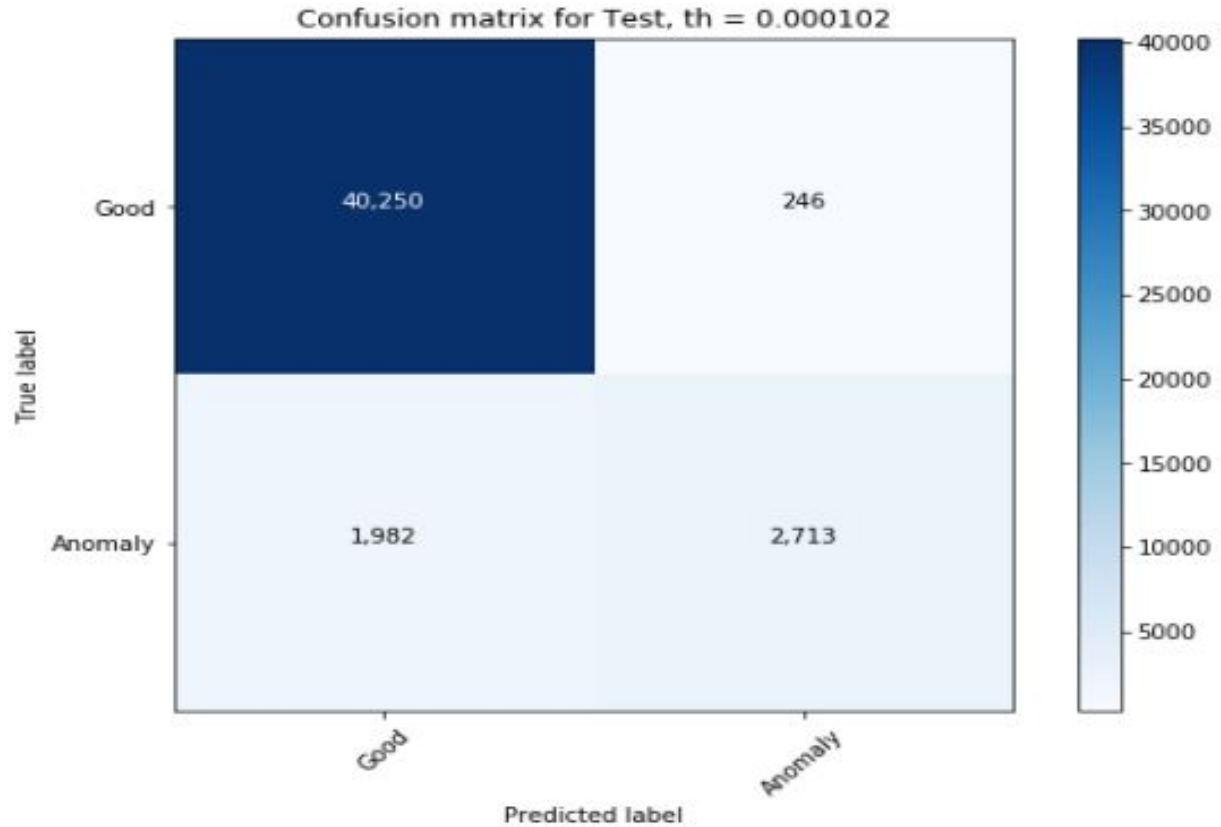


Precision and recall for different threshold values for Test



Best Threshold = 0.000102, F1 Score = 0.709 (Same as Standard AE)

Confusion Matrix



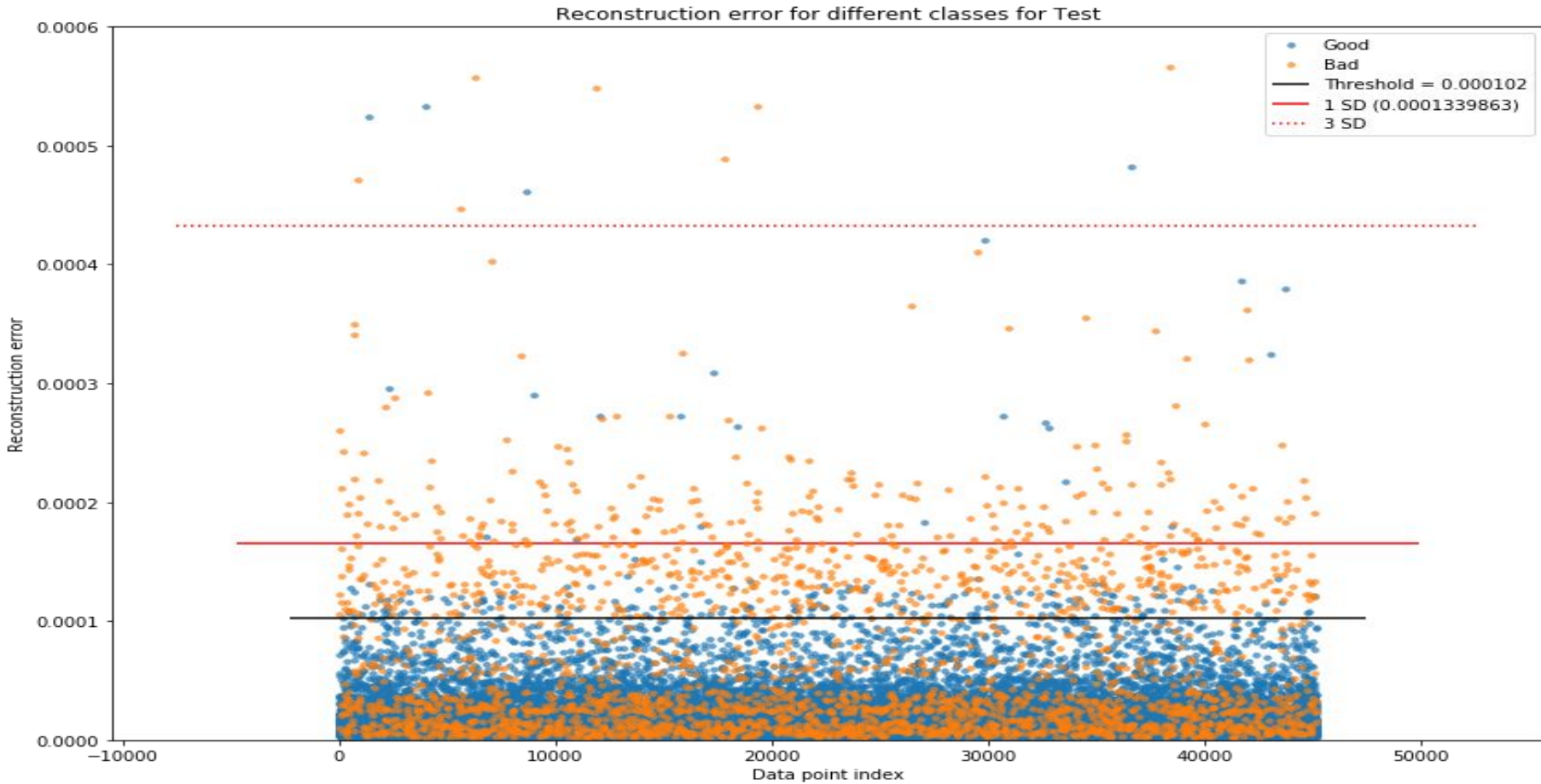
FP = 0.61% of N

FN = 42.2% of P

Comparable with standard AE

Till now, no further improvement wrt standard AE

Reconstruction error for test dataset (Good & Bad)



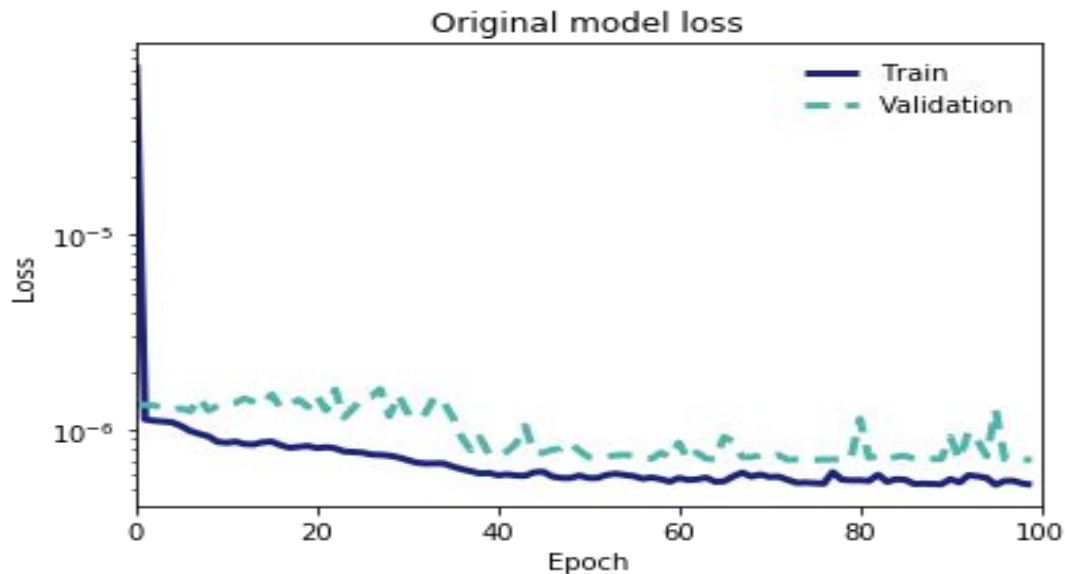
Future plan & Acknowledgements

- We further want to investigate the models and the parameters more.
- On behalf of Tracker, we want to start looking into Hist1D Normalized HitResiduals_TIB_layer* for 2017 data.
- Reference for LSTM & PR Curve:
 - <https://towardsdatascience.com/step-by-step-understanding-lstm-autoencoder-layers-ffab055b6352>
 - <https://machinelearningmastery.com/roc-curves-and-precision-recall-curves-for-classification-in-python/>
- Thanks Francesco for sharing the code snippets! It helped me a lot as a beginner.

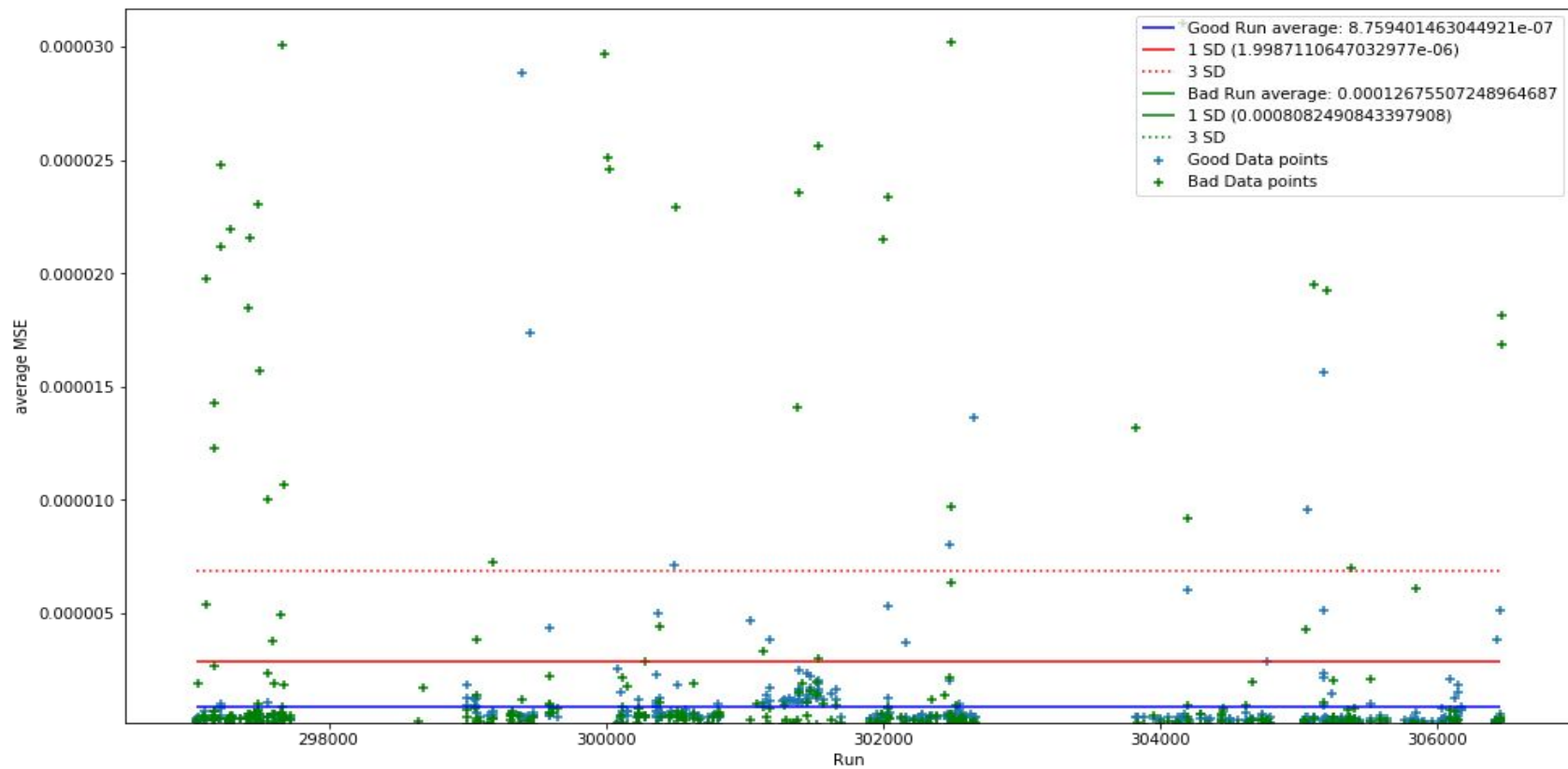
Thank you!

Standard Autoencoder

- Standard DNN based autoencoder with tanh activation in intermediate layers final layer. input - 10 - 3 - 10 - output
- Optimizer = adam , loss = mseTop10
Good runs and lumis selected from Golden json file



GlobalMSETrend

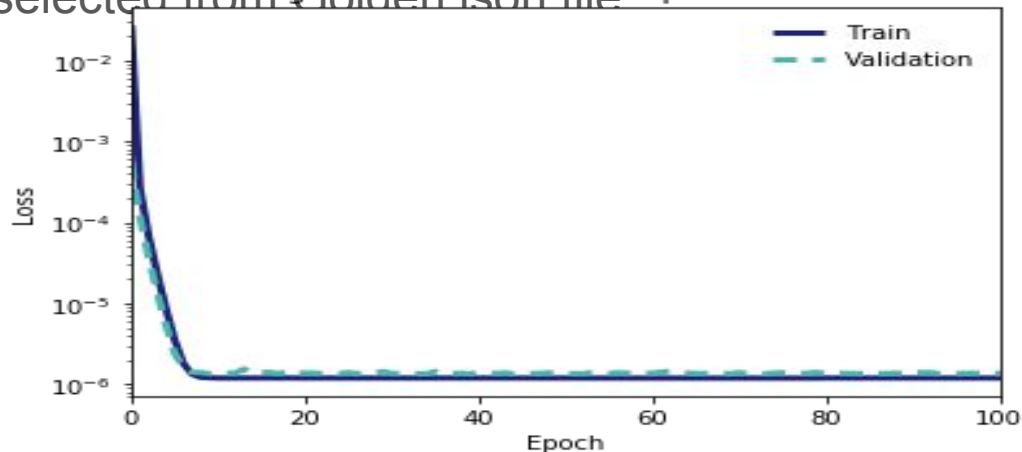


Deep autoencoder

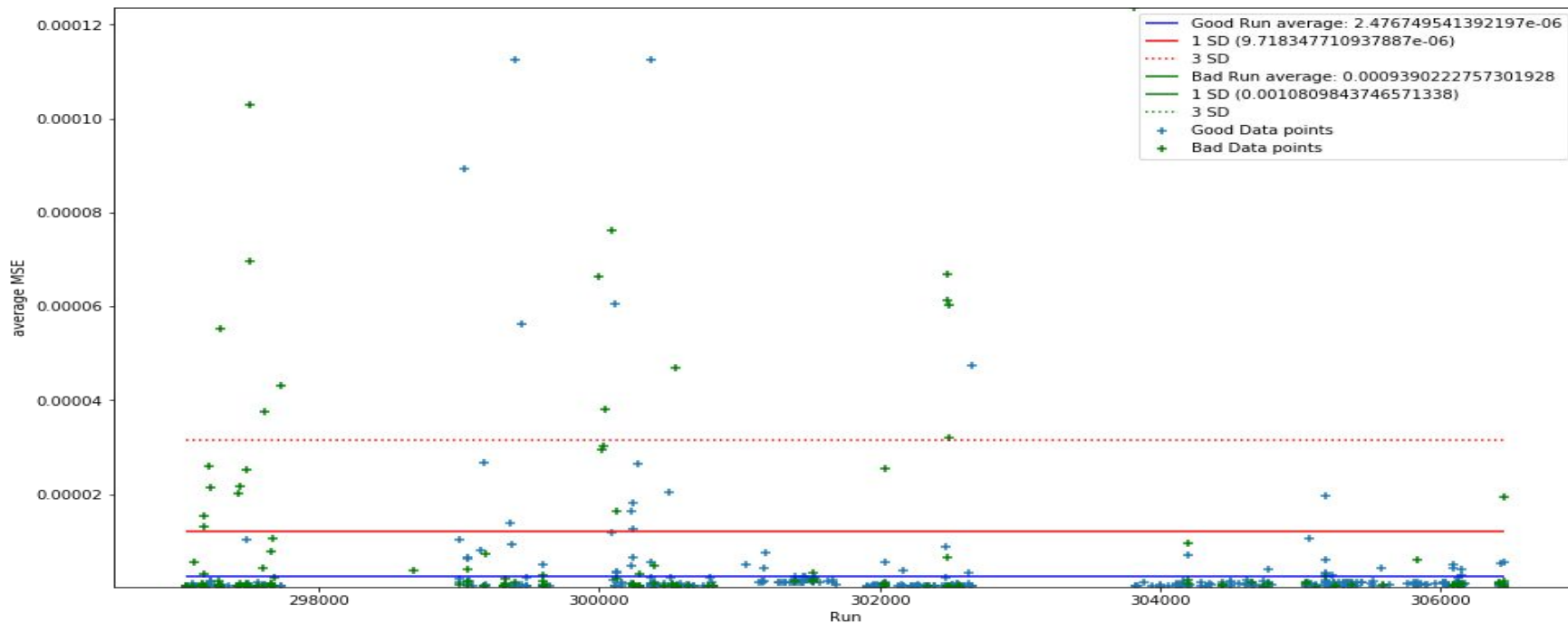
- Deeper DNN based autoencoder with relu activation in intermediate layers and tanh in the final layer. (input - 20 - 10 - 20 - output)

Optimizer = adam , loss = mseTop10

Good runs and lumis selected from Golden-Isom file

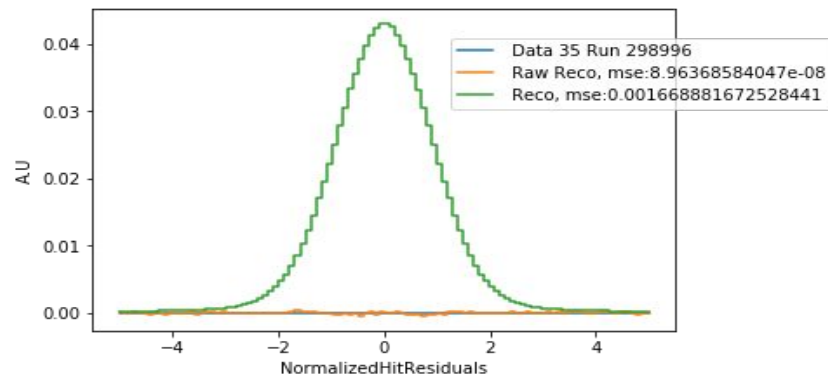
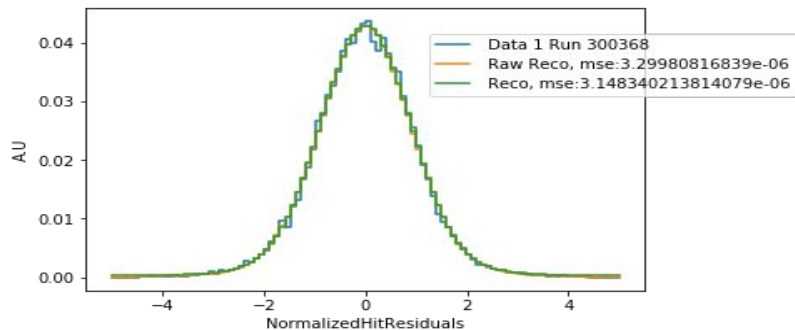
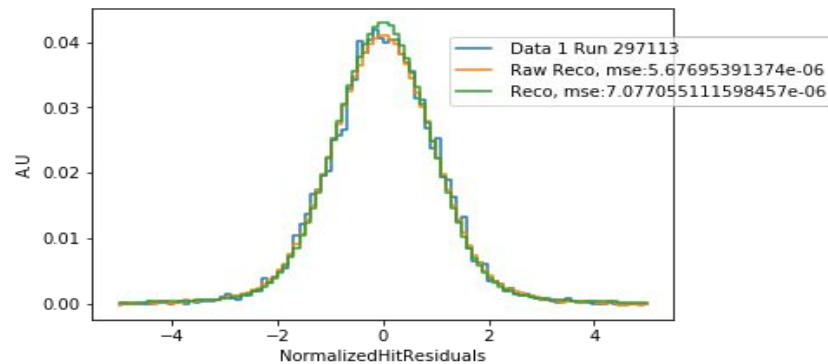
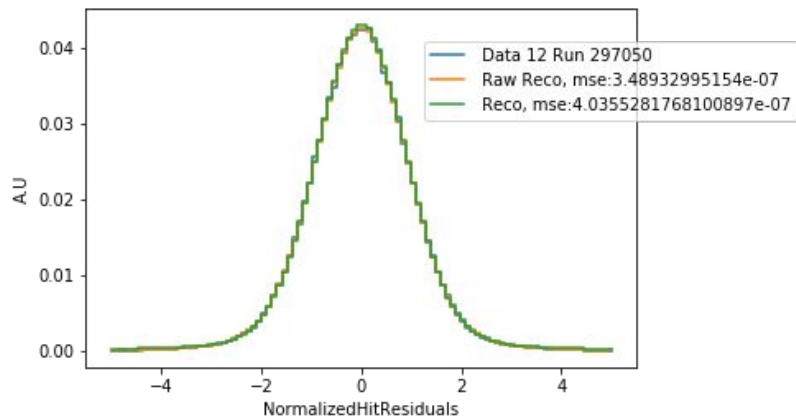


GlobalMSETrend

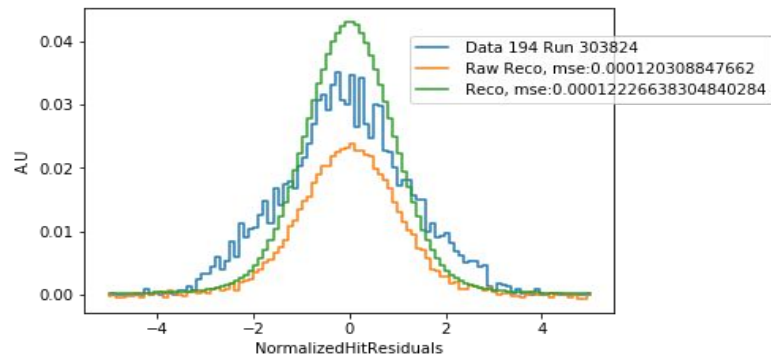
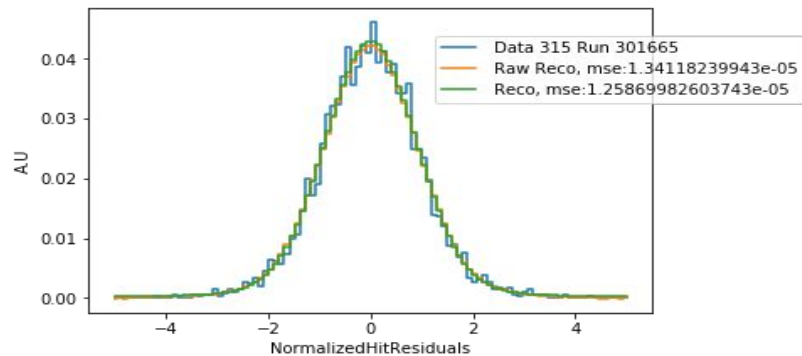
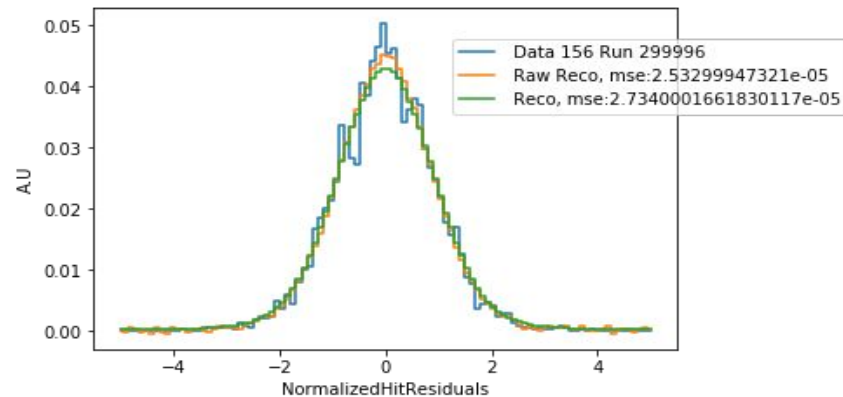
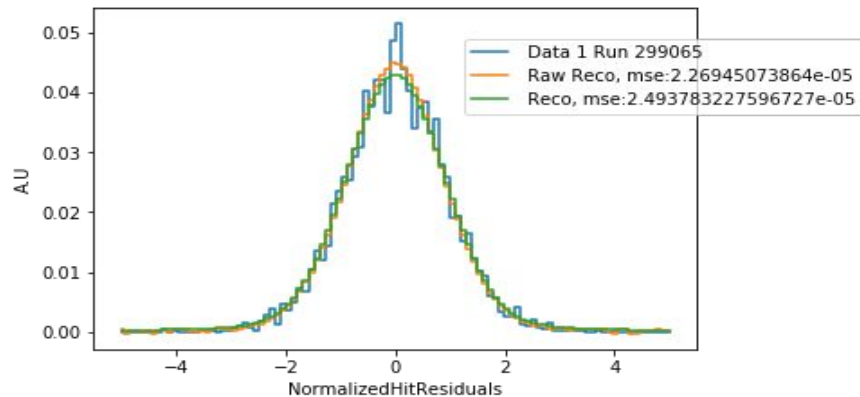


Deep AE works better in global mse trend
Good and bad lumis are not well separated.

Predictions for good lumis [Raw reco = first AE, Reco = Deep AE]



Predictions for bad lumis [Raw reco = first AE, Reco = Deep AE]



Using LSTM autoencoder for chargeInner_PXLayer_1

Using individual bin as a step in a sequence, we tried to encode histogram as sequence in LSTM.

Model: "sequential_2"

Layer (type)	Output Shape	Param #
layer_0 (LSTM)	(None, 10)	480
layer_1 (RepeatVector)	(None, 100, 10)	0
layer_2 (LSTM)	(None, 100, 10)	840
layer_3 (TimeDistributed)	(None, 100, 1)	11

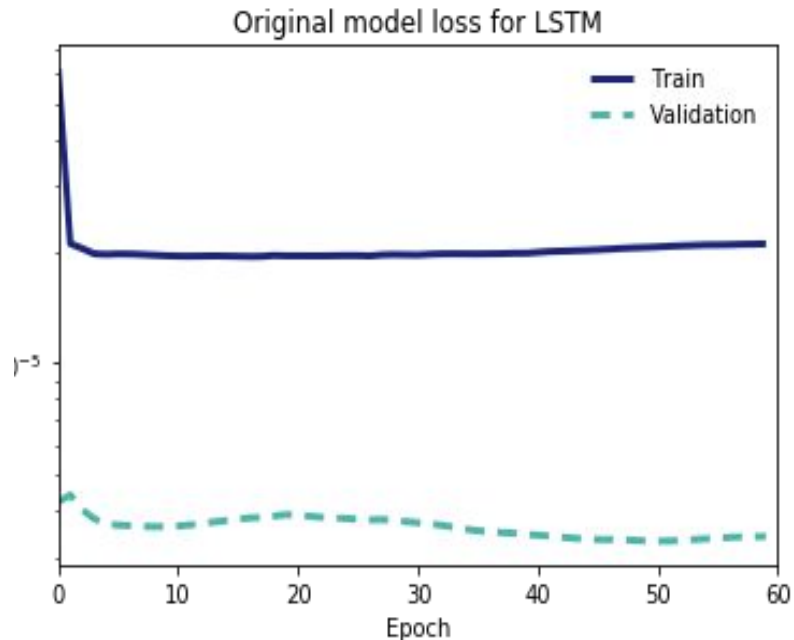
Total params: 1,331

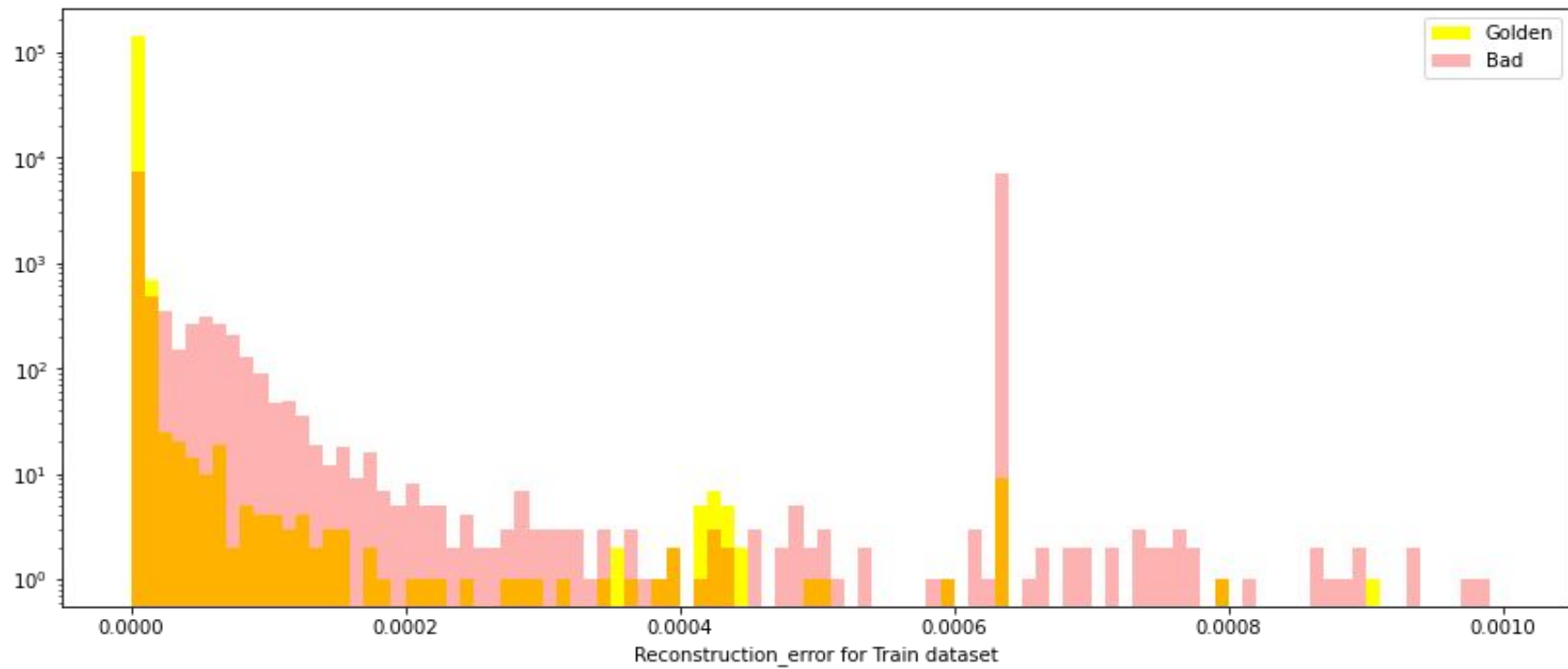
Trainable params: 1,331

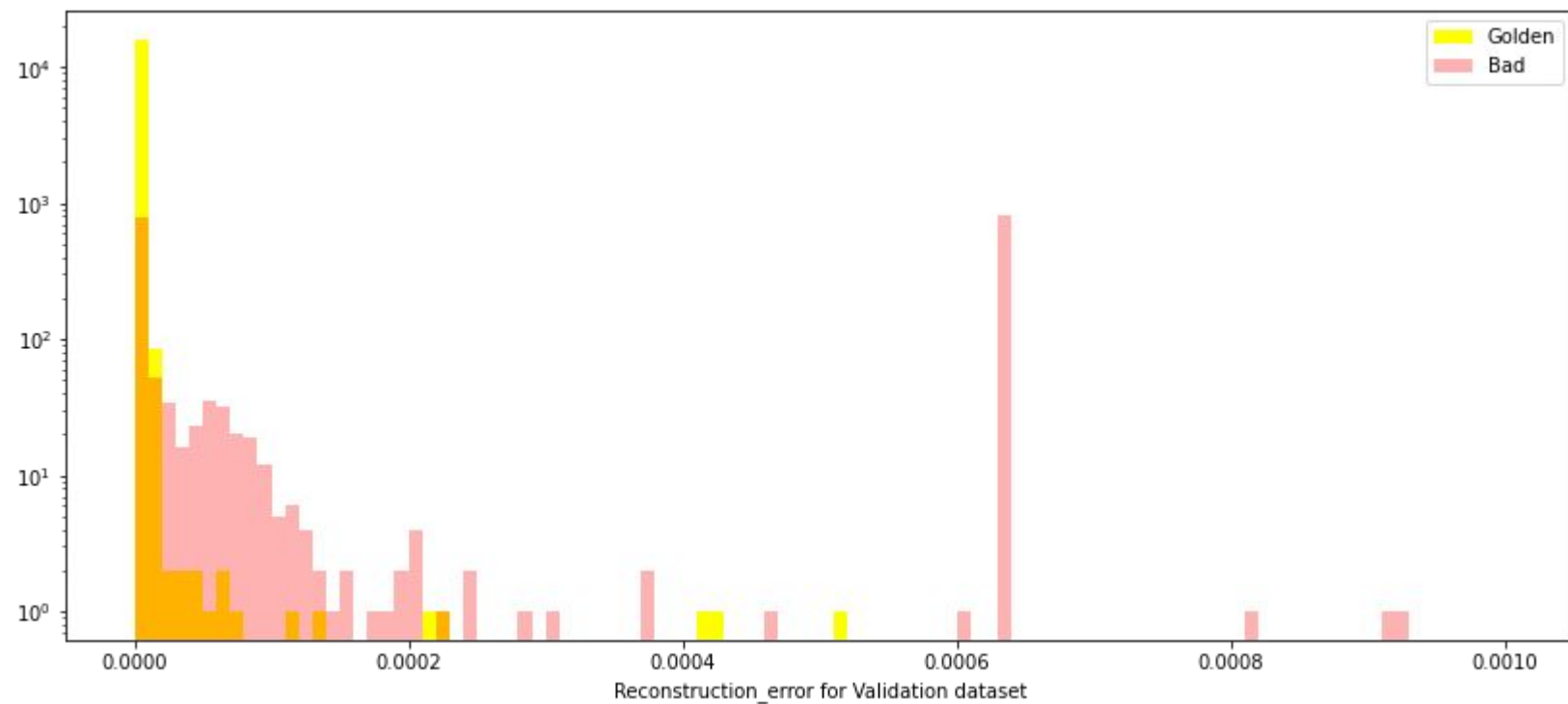
Non-trainable params: 0

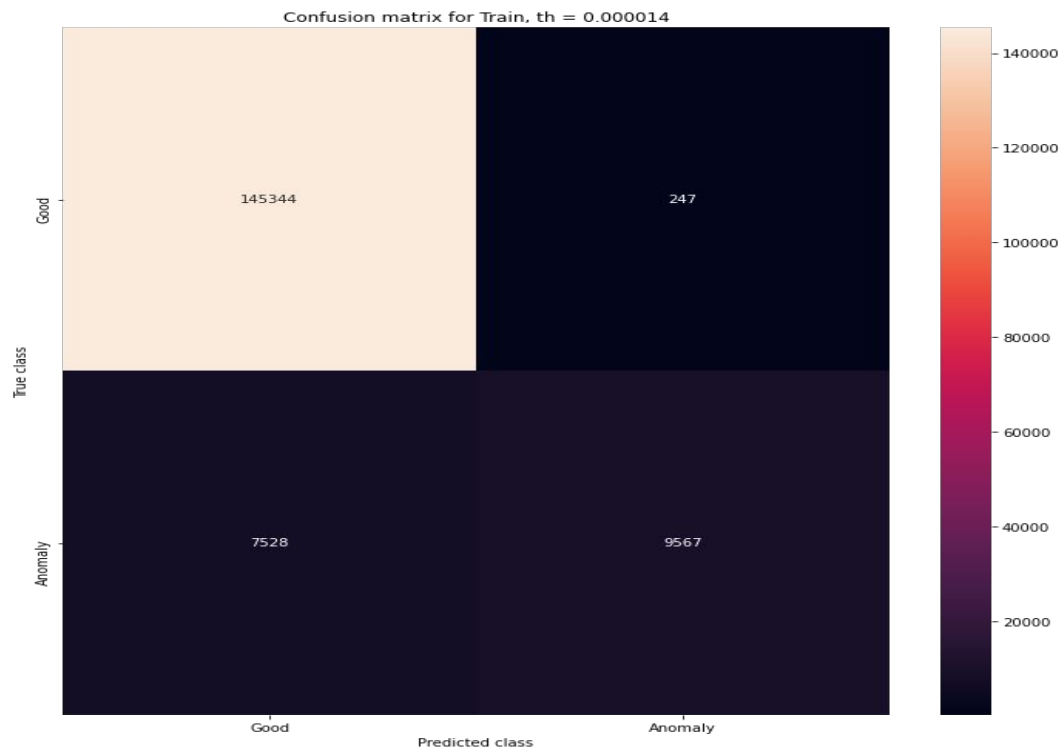
Yet to arrive at optimised network architecture.

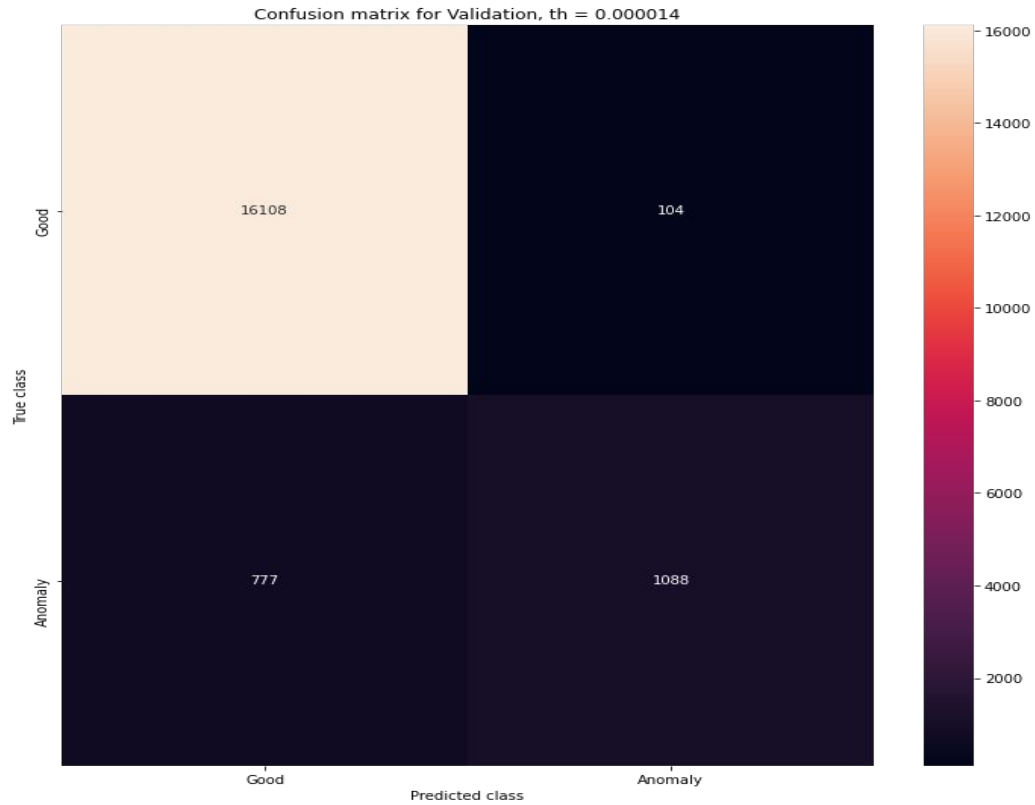
Work in progress.











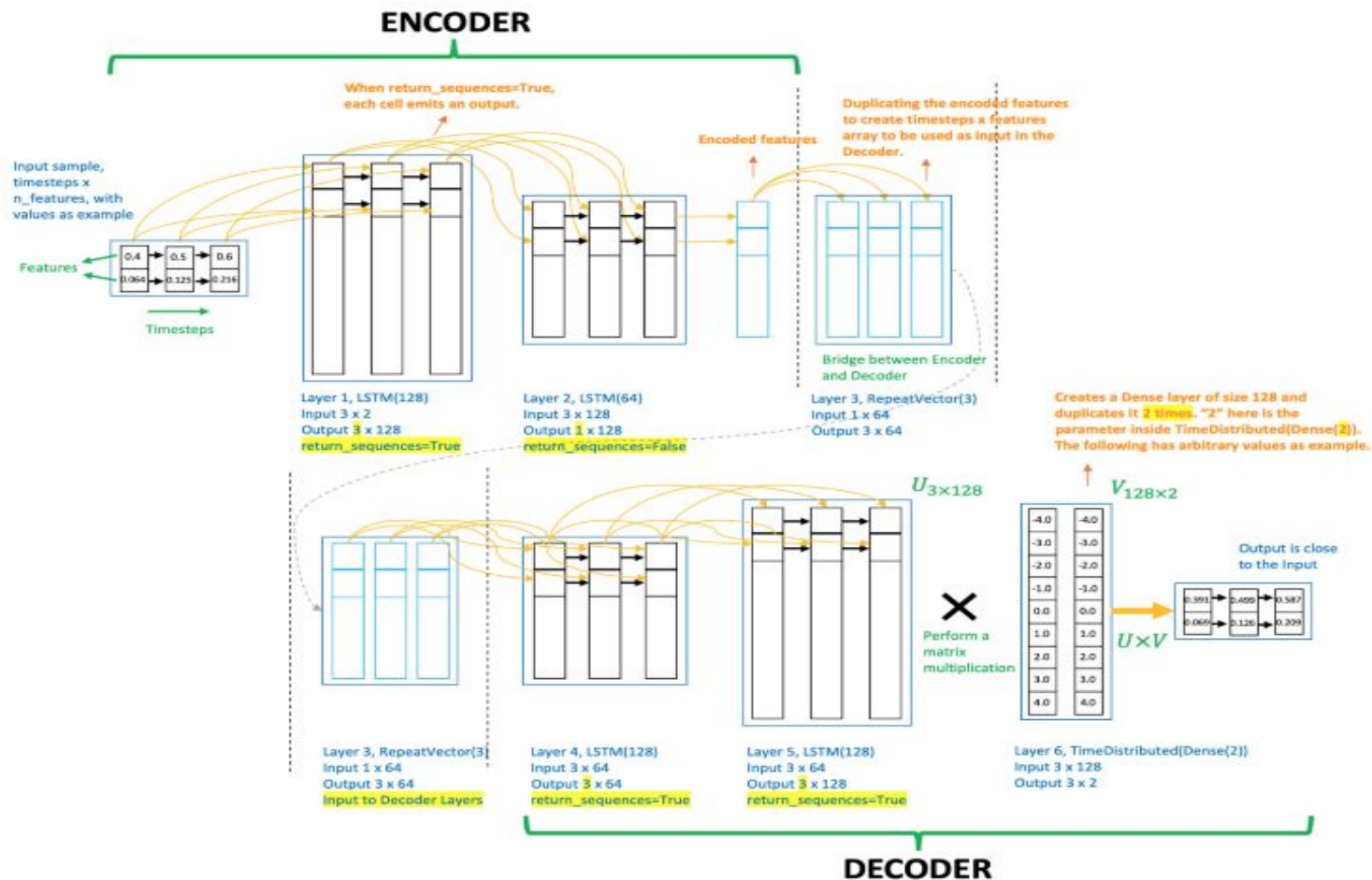
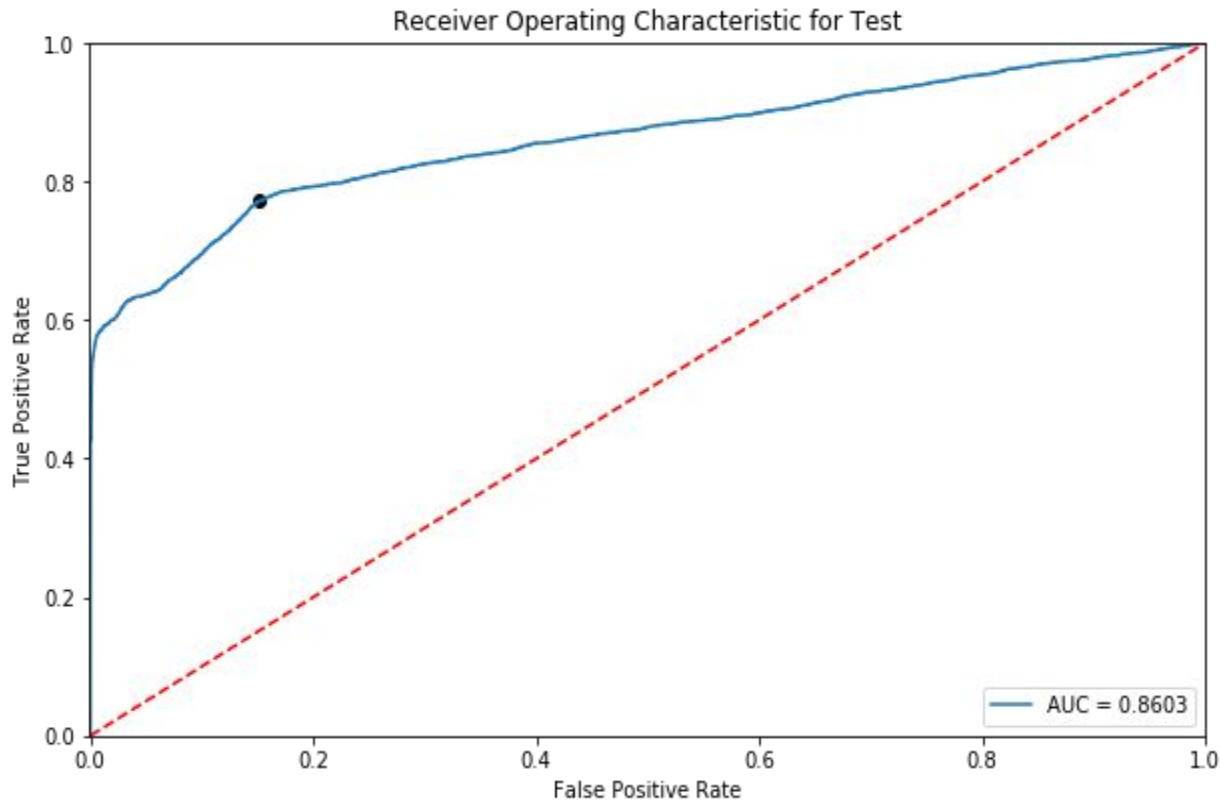
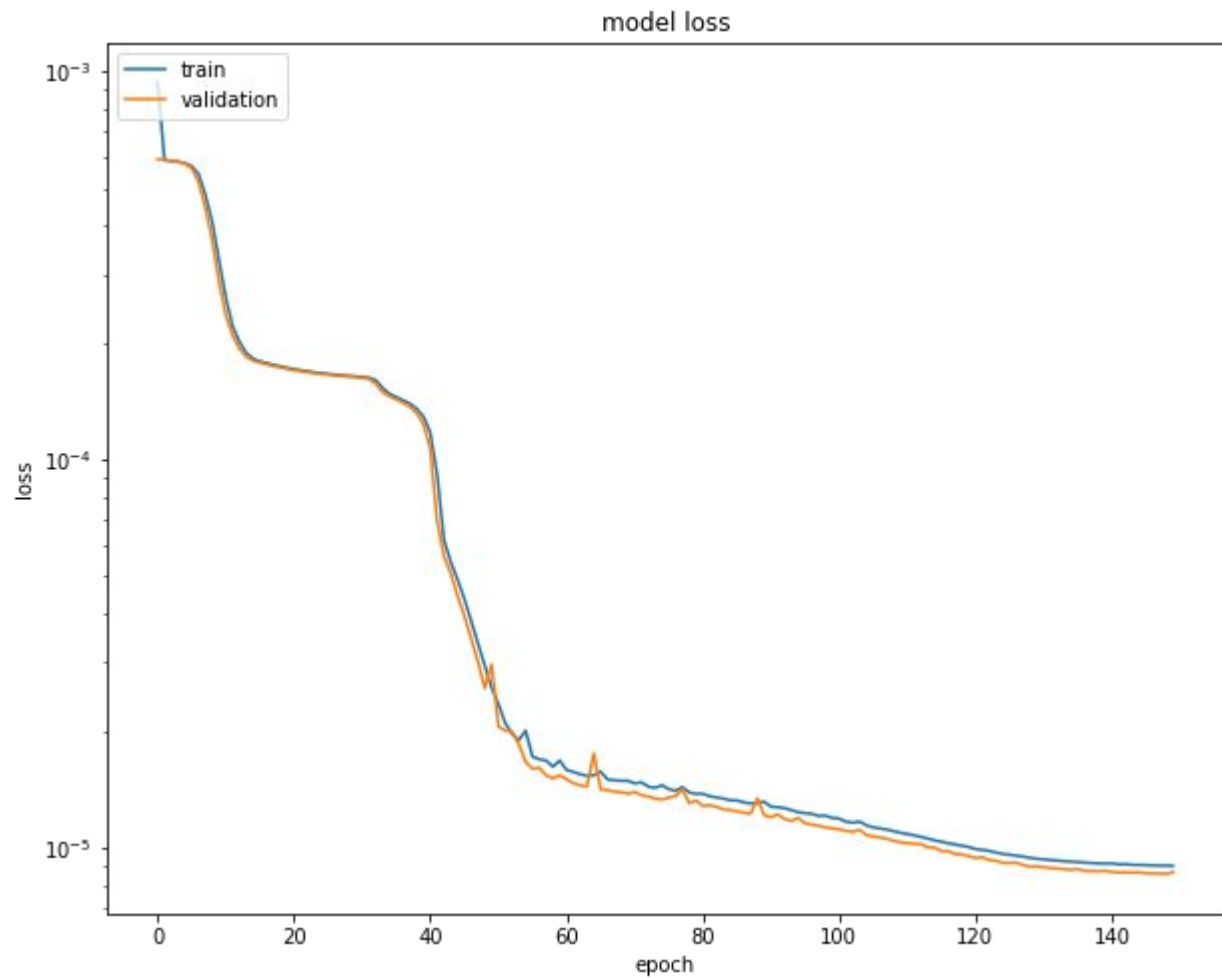
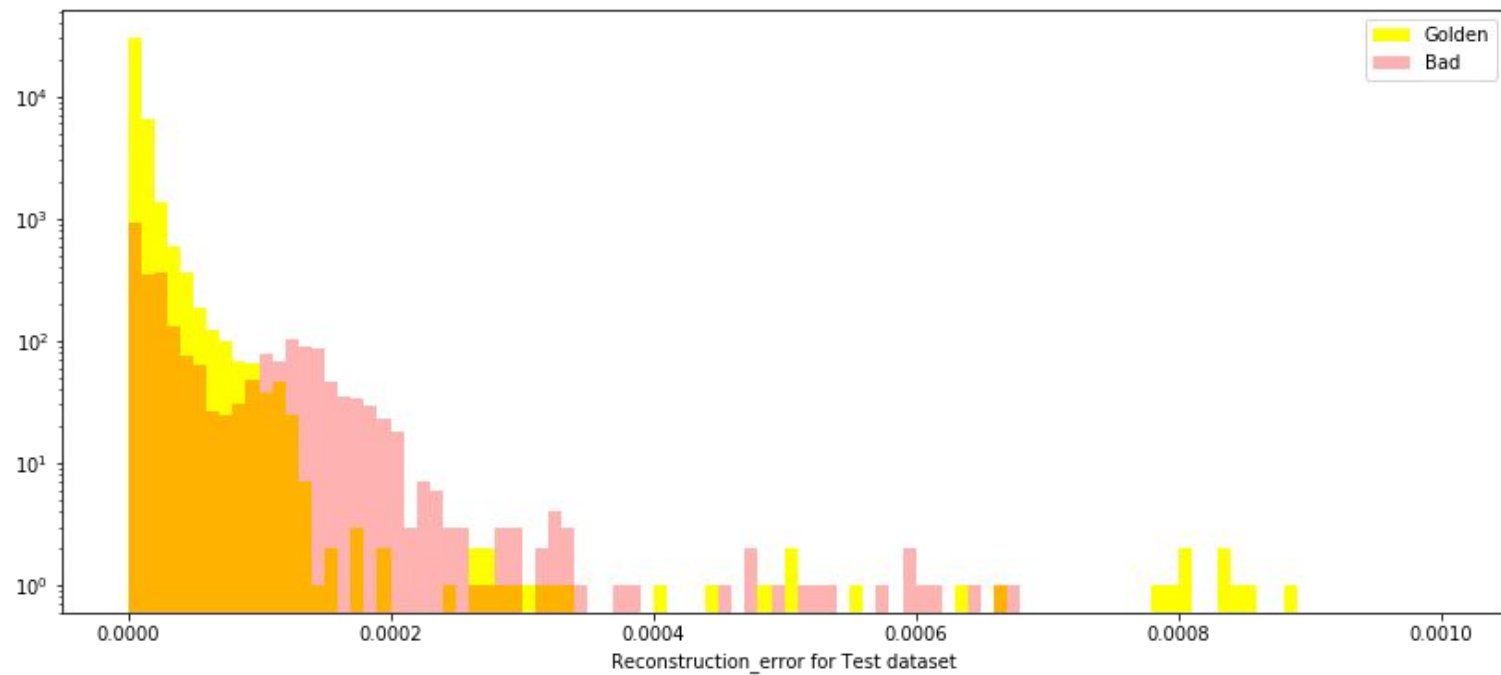


Figure 2.3. LSTM Autoencoder Flow Diagram.

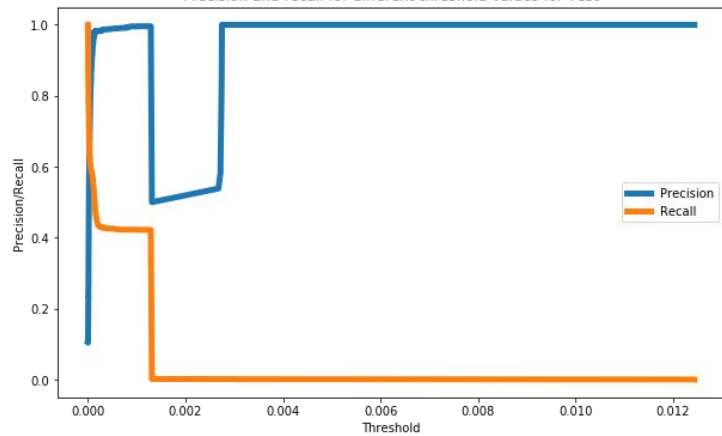
LSTM ROC



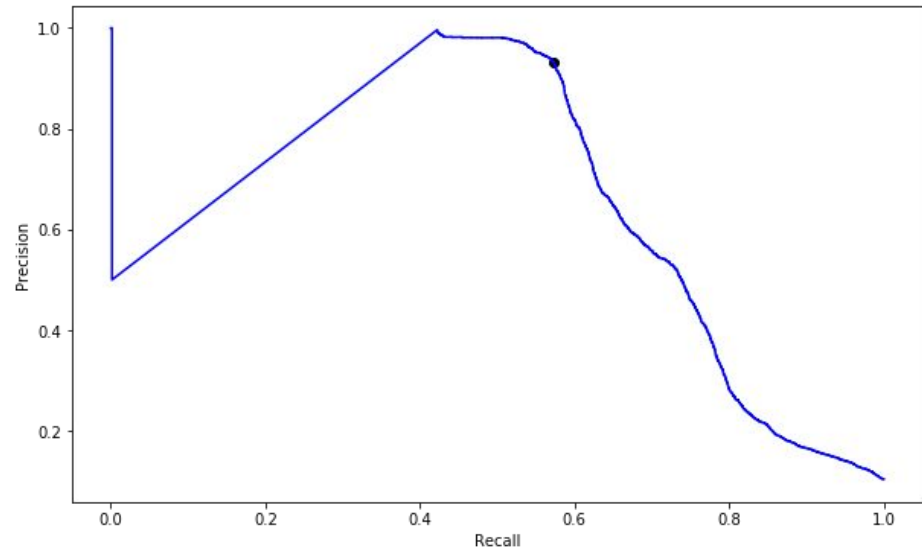




Precision and recall for different threshold values for Test



Recall vs Precision for Test



Confusion matrix for Test, th = 0.000094

