

# LCD and LArIAT Datasets And CaloDNN and LArTPCDNN

Amir Farbin  
(ATLAS/UTA)

LCD Calo Dataset made by M. Pierini (CMS/CERN) + JR Vlimant (CMS/Caltech)  
LArIAT Dataset made by S. Shaksavarani (Neutrinos/UTA) + AF

# Intro

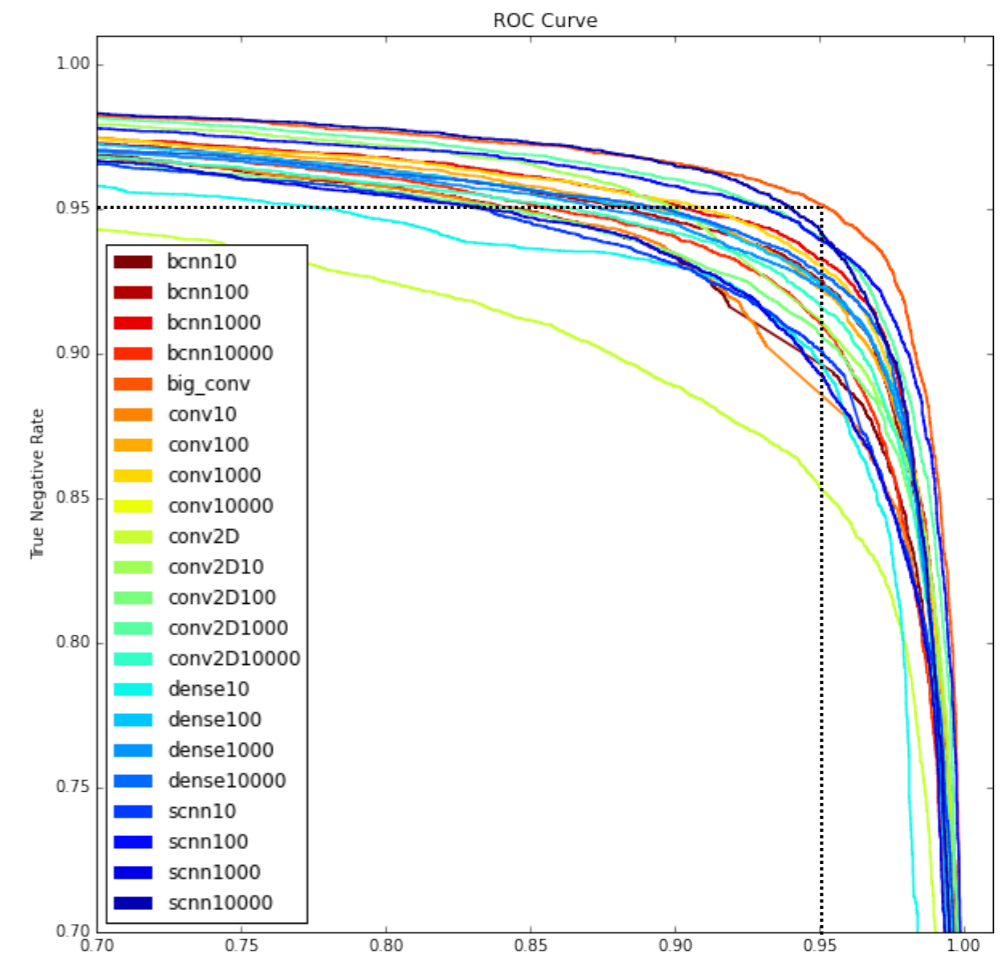
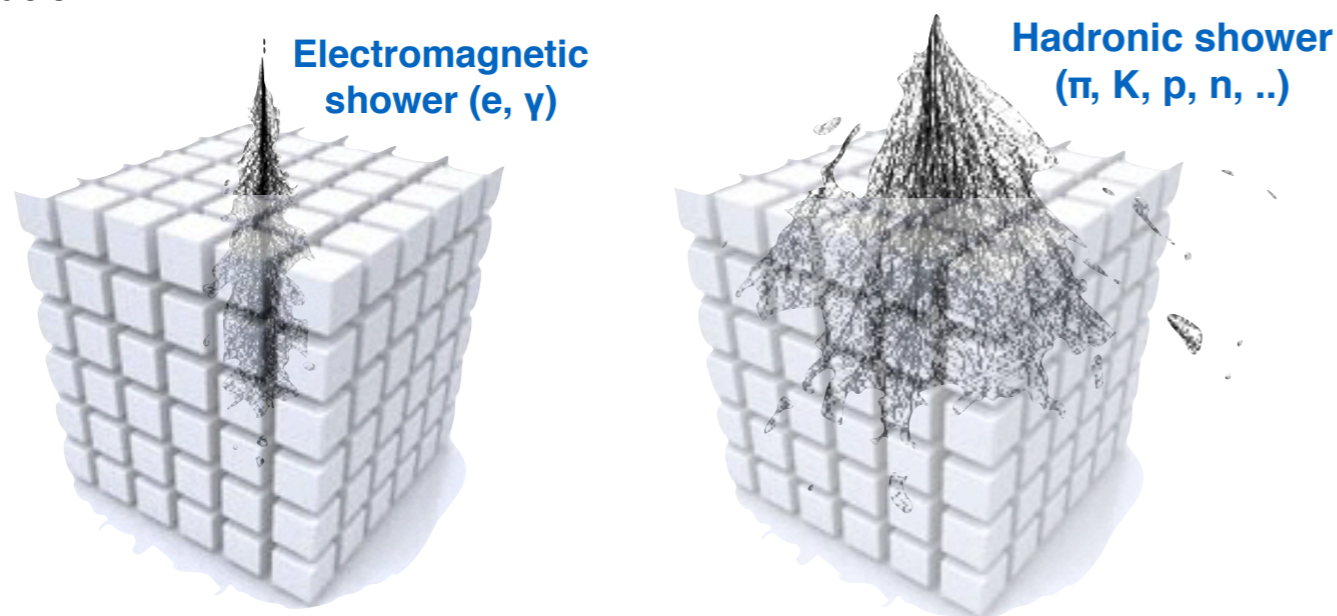
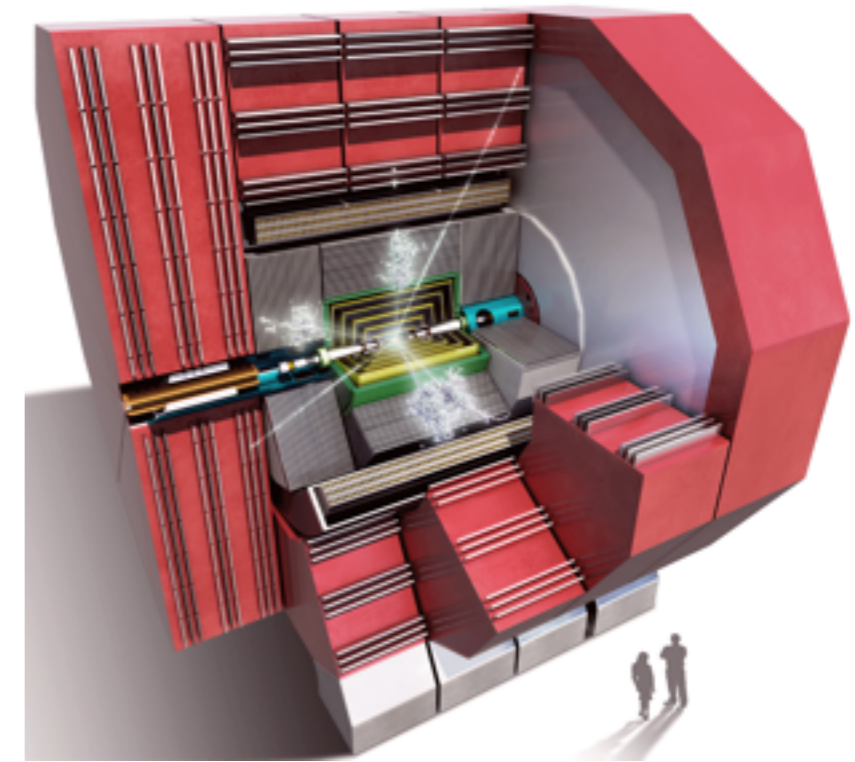
- Reconstruction level DL requires realistic detector simulation... not as easy as 4-vectors or parameterized detectors.
- Experiments are understandably strict about their data. Prohibits:
  - Cross experiment or HEP/ML collaboration
  - Rapid publication of DL R&D (no physics).
- Imaging detectors (Granular Calorimeters, TPCs, Cherenkov, ...) ideally suited for Deep Learning.
- We generated the LCD and LArIAT Datasets to avoid these issues.
  - Dataset and code very similar, so I'll talk about both.
  - Weekly LCD meetings to organize work. Should do for LArIAT.
- Data Science @ LHC (Nov 2015 @ CERN) -> DS@HEP.
  - Experts workshop (July 2015): these datasets were introduced in prim. Goal was to make them public for NIPS... btut we didn't get a workshop and got busy.
  - Goal is to reveal datasets at next workshop. May 8-12 @ FNAL. <https://indico.fnal.gov/conferenceDisplay.py?confId=13497>

# Message

- Everyone is busy, so help is appreciated:
  - Contribute to finalizing data and Nature Scientific Data paper.
  - Collaborate on research.
  - We ask that Dataset paper would be the first, and all work done before DS@HEP WS be collaborative.
- These are large datasets (LCD = 20 GB so far, LArIAT = 20 TB)
  - Distribution and processing require extra thought
  - Code to efficiently read the data should be provided.
- Not clear if we should distribute full running examples... or even collaborative code used for papers.
  - I'll present my packages... open to input and suggestions.
- I feel like I'm often working in a corner may make mistakes.
  - I have lots of questions I have no one to ask.
  - I hope this forum could be a place to share experiences and give advice...

# LCD Calorimeter

- CLIC is a proposed CERN project for a linear accelerator of electrons and positrons to TeV energies (~ LHC for protons)
  - Not a real experiment yet, so we can simulate data and make it public.
  - Simpler geometry than ATLAS...
- The LCD calorimeter is an array of absorber material and silicon sensors comprising the most granular calorimeter design available
  - Data is essentially a 3D image
  - So far several million  $\pi^0$ , Elec, ChPi, Gamma. 10 to 510 GeV. Low energy and Jet samples planned.
  - ECAL (25x25x25) / HCAL (5x5x60) "window". Aux info: Energy, ...
- First studies,  $\pi$  vs  $\gamma$  classification with various DNNs by summer students.
  - Code/results not collected... but should be easy to redo.
  - New version of dataset.
  - Some visualization code exists... Full running example in CaloDNN.
- Many interesting problems: PID Classification, Energy Regression, Shower generative models.



# Join the fun....

Imaging calorimeter data for Machine Learning applications in HEP

---

Josh Bendavid<sup>a</sup> Kaustuv Datta<sup>a,b</sup> Amir Farbin<sup>c</sup> Nikolaus Howe<sup>d,e</sup> Jayesh Mahapatra<sup>d</sup>  
Maurizio Pierini<sup>d</sup> Maria Spiropulu<sup>a</sup> Jean-Roch Vlimant<sup>a</sup>

<sup>a</sup> *California Institute of Technology*

<sup>b</sup> *Reed College*

<sup>c</sup> *University of Texas at Arlington*

<sup>d</sup> *CERN*

<sup>e</sup> *Williams College*

Photon identification and energy measurement with a highly granular calorimeter through Deep Learning

---

Josh Bendavid<sup>a</sup> Kaustuv Datta<sup>a,b</sup> Amir Farbin<sup>c</sup> Nikolaus Howe<sup>d,e</sup> Jayesh Mahapatra<sup>d</sup>  
Maurizio Pierini<sup>d</sup> Maria Spiropulu<sup>a</sup> Jean-Roch Vlimant<sup>a</sup>

<sup>a</sup> *California Institute of Technology*

<sup>b</sup> *Reed College*

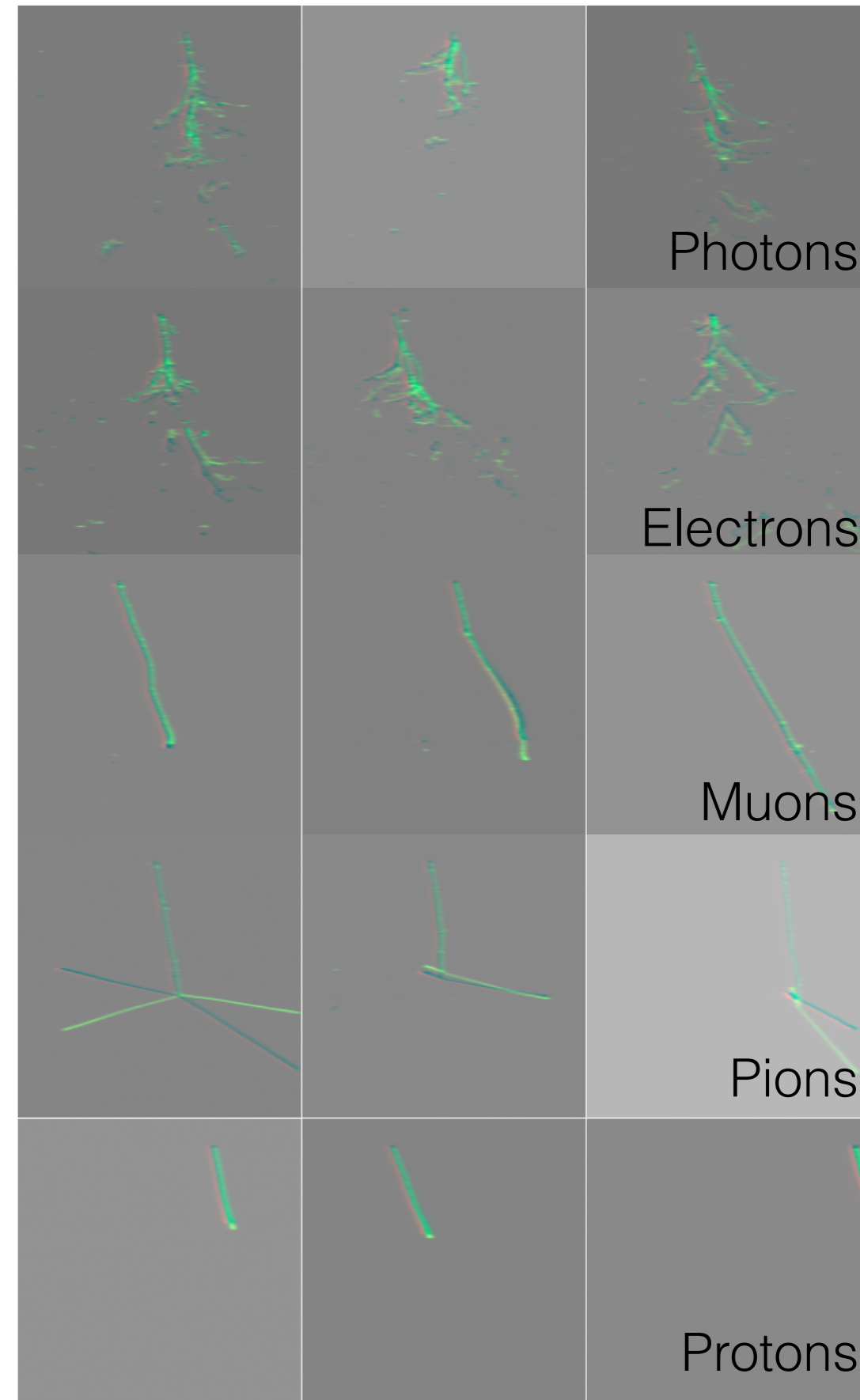
<sup>c</sup> *University of Texas at Arlington*

<sup>d</sup> *CERN*

<sup>e</sup> *Williams College*

# LArIAT Data

- LArIAT is a small LArTPC detector: 2 wire planes with 240 wires each, 4096 samples.
- 1 M each of: antielectron, kaonPlus, nue\_CC, nutaubar\_CC pionMinus, antimuon, nue\_NC, nutaubar\_NC, pionPlus, antiproton, muon, numubar\_CC, nutau\_CC, electron, numubar\_NC nutau\_NC, proton, nuebar\_CC, numu\_CC, photon, kaonMinus, nuebar\_NC, numu\_NC, pion\_0
- Data: Sim done.
  - Raw ADC readout: 2 x 4096 x 240 (essentially no noise)
  - Geant4 charge deposits. SparseTensor allows creating 3D images of any resolution. (Needs reprocessing of data-prep steps)
  - Aux info: type of interaction, energy, ...
- Studies:
  - Preliminary studies very promising.
  - Subsequent work (P. Sadowski + ?) showed impressive classification performance using siamese inception model trained for 1 week.
  - A bit of work on energy regression... not as straightforward.
  - Progress stalled...
- Interesting problems: PID classification, Energy Regression, Compression/ Noise suppression, 2x 2D -> 3D (DNN tomography)



# Technical Challenges

- Data comes as many h5 files, each containing  $O(1000)$  events, organized into directories by particle type.
- Needs to be read, mixed, “labeled”, and normalized.... can be time consuming.
- Doesn't fit in memory...
- Very difficult to keep the GPU fed with data. GPU utilization often  $< 10\%$ , rarely  $> 50\%$ .
- Keras python generator mechanism:
  - Allows reading on the fly and parallel read
  - Found 2 problems: (Am I crazy?)
    - Multiprocessing requires the generators to be `thread_safe`, which means putting in a locking mechanism which only allows one process to read the data at a time. So  $> 2$  processes not useful.
    - Easy to mess up and have parallel generator instances deliver overlapping data.
  - LCD data is  $\sim x10$  slower with naive Keras generator vs preloading in memory.
- I wrote a standalone parallel generator: `DLKit/ThreadedGenerator`:
  - Python Global Interpreter Lock (GIL) allows only one thread to run at a time... so must use multiprocessing.
  - Current implementation: Filler process sends requests (file/block) via multiprocessing queues to workers processes that deliver data to corresponding threads via pipes that feed the generator via thread queues.
  - Bottle neck is the process to thread pipe... data needs to be serialized. Working on share memory solution...
  - Data can be premixed. Premix:  $\sim 2x$  slower than data in memory. Mix as you go:  $\sim 4x$  slower than data in memory.
  - System resources become problem when running many trainings in same system. Working on framework upgrade to simultaneously train several models with same data.

# DLKit

- Thin layer on top of Keras.
- My personal DNN framework. I imagine many of you would write something similar...
- Handles book keeping for comparing large number of training sessions (e.g. for hyper parameter scan or optimization)
- Tools necessary to setup HEP problems.
- I have several HEP problems setup using this package:
  - EventClassificationDNN, MEDNN, CaloDNN, LArTPCDNN, ...
- Hyperas or Spearmin integration demonstrated, but needs work.
- Keras / MPI Integration also in the works.
- Already ran on BlueWaters and Titan.
- <https://bitbucket.org/anomalousai/dlkit/src>

## Source

🔗 master ▾ | 📄 ▾ | [DLKit](#) / [DLTools](#) /

---

⬆️ ..

---

📄 [Callbacks.py](#)

---

📄 [GPUQueuesNJobs.sh](#)

---

📄 [LoadModel.py](#)

---

📄 [ModelWrapper.py](#)

---

📄 [Permutator.py](#)

---

📄 [Printh5File.py](#)

---

📄 [README.md](#)

---

📄 [ScanAnalysis.py](#)

---

📄 [SparseTensorDataSet.py](#)

---

📄 [TarResults.sh](#)

---

📄 [ThreadedGenerator.py](#)

---

📄 [\\_\\_init\\_\\_.py](#)













---

📄 [clean.sh](#)



# CaloDNN/LArTPCDNN

- Instantiates generators for efficiently reading or premixing data.
- Provides out-of-the-box running realistic (not toy) models.
- Orchestrates running large HP scans.
  - Makes tables...
  - Jupyter notebook analysis in works.
- Generates standard plots.
- <https://github.com/UTA-HEP-Computing/CaloDNN>
- Polishing up package for public...
- Gearing up for a big BlueWaters run...
  - Large HP Scan (not optimization)
  - “Regularization”: training time.

 <a href="#">Analysis.py</a>
 <a href="#">ClassificationArguments.py</a>
 <a href="#">ClassificationExperiment.py</a>
 <a href="#">ClassificationScanConfig.py</a>
 <a href="#">LCDDData.py</a>
 <a href="#">Models.py</a>
 <a href="#">README.md</a>
 <a href="#">ScanJob.py</a>
 <a href="#">ScanJob.sh</a>
 <a href="#">SubmitMerge.sh</a>
 <a href="#">__init__.py</a>
 <a href="#">requirements.txt</a>

Last login: Tue Feb 28 08:47:35 2017 from 192.168.1.13

afarbin@thecount:~\$ cd LCD/DLKit/

afarbin@thecount:~/LCD/DLKit\$ source setup.sh

(Keras) afarbin@thecount:~/LCD/DLKit\$ python -m CaloDNN.ClassificationExperiment --help

```
usage: ClassificationExperiment.py [-h] [-C CONFIG] [-L LOADMODEL]
                                   [--gpu GPUID] [--cpu] [--NoTrain]
                                   [--NoAnalysis] [--Test] [-s HYPERPARAMSET]
                                   [--nopremix] [--preload] [-r RUNNINGTIME]
```

#### optional arguments:

```
-h, --help            show this help message and exit
-C CONFIG, --config CONFIG
                       Use specified configuration file.
-L LOADMODEL, --LoadModel LOADMODEL
                       Loads a model from specified directory.
--gpu GPUID           Use specified GPU.
--cpu                 Use CPU.
--NoTrain             Do not run training.
--NoAnalysis          Do not run analysis.
--Test               Run in test mode (reduced examples and epochs).
-s HYPERPARAMSET, --hyperparamset HYPERPARAMSET
                       Use specified (by index) hyperparameter set.
--nopremix            Do not use the premixed inputfile. Mix on the fly.
--preload             Preload the data into memory. Caution: requires lots
                       of memory.
-r RUNNINGTIME, --runningtime RUNNINGTIME
                       End training after specified number of seconds.
```

(Keras) afarbin@thecount:~/LCD/DLKit\$ █

# ScanConfig.py

```
6
7 # Input for Premixed Generator
8 InputFile="/data/afarbin/LCD/LCD-Merged-All.h5"
9 # Input for Mixing Generator
10 FileSearch="/data/afarbin/LCD/*/*.h5"
11
12 # Generation Model
13 Config={
14     "GenerationModel":"'Load'",
15     "MaxEvents":int(3.e6),
16     "NTestSamples":100000,
17     "NClasses":4,
18
19     "Epochs":1000,
20     "BatchSize":1024,
21
22     # Configures the parallel data generator that read the input.
23     # These have been optimized by hand. Your system may have
24     # more optimal configuration.
25     "n_threads":4, # Number of workers
26     "multiplier":2, # Read N batches worth of data in each worker
27
28     # How weights are initialized
29     "WeightInitialization":"'normal'",
30
31     # Normalization determined by hand.
32     "ECAL":True,
33     "ECALNorm":150.,
34
35     # Normalization needs to be determined by hand.
36     "HCAL":True,
37     "HCALNorm":150.,
--
```

```

38
39 # Set the ECAL/HCAL Width/Depth for the Dense model.
40 # Note that ECAL/HCAL Width/Depth are changed to "Width" and "Depth",
41 # if these parameters are set.
42 "HCALWidth":32,
43 "HCALDepth":2,
44 "ECALWidth":32,
45 "ECALDepth":2,
46
47 # No specific reason to pick these. Needs study.
48 # Note that the optimizer name should be the class name (https://keras.io/optimizers/)
49 "loss":"'categorical_crossentropy'",
50
51 # Specify the optimizer class name as True (see: https://keras.io/optimizers/)
52 # and parameters (using constructor keywords as parameter name).
53 # Note if parameter is not specified, default values are used.
54 "optimizer":"'SGD'",
55 #"lr":0.01,
56 #"decay":0.001,
57
58 # Parameter monitored by Callbacks
59 "monitor":"'val_loss'",
60
61 # Active Callbacks
62 # Specify the Callback class name as True (see: https://keras.io/callbacks/)
63 # and parameters (using constructor keywords as parameter name,
64 # with classname added).
65 "ModelCheckpoint":True,
66 "Model_Chekpoin_save_best_only":False,
67
68 # Configure Running time callback
69 # Set RunningTime to a value to stop training after N seconds.
70 "RunningTime": 3600,
71 }

```

```

72
73 # Parameters to scan and their scan points.
74 Params={ "Width": [32,64,128,256,512],
75          "Depth": range(1,5),
76          "lr": [0.1,0.01,0.001],
77          "decay": [0.1,0.01,0.001],
78          }
79

```