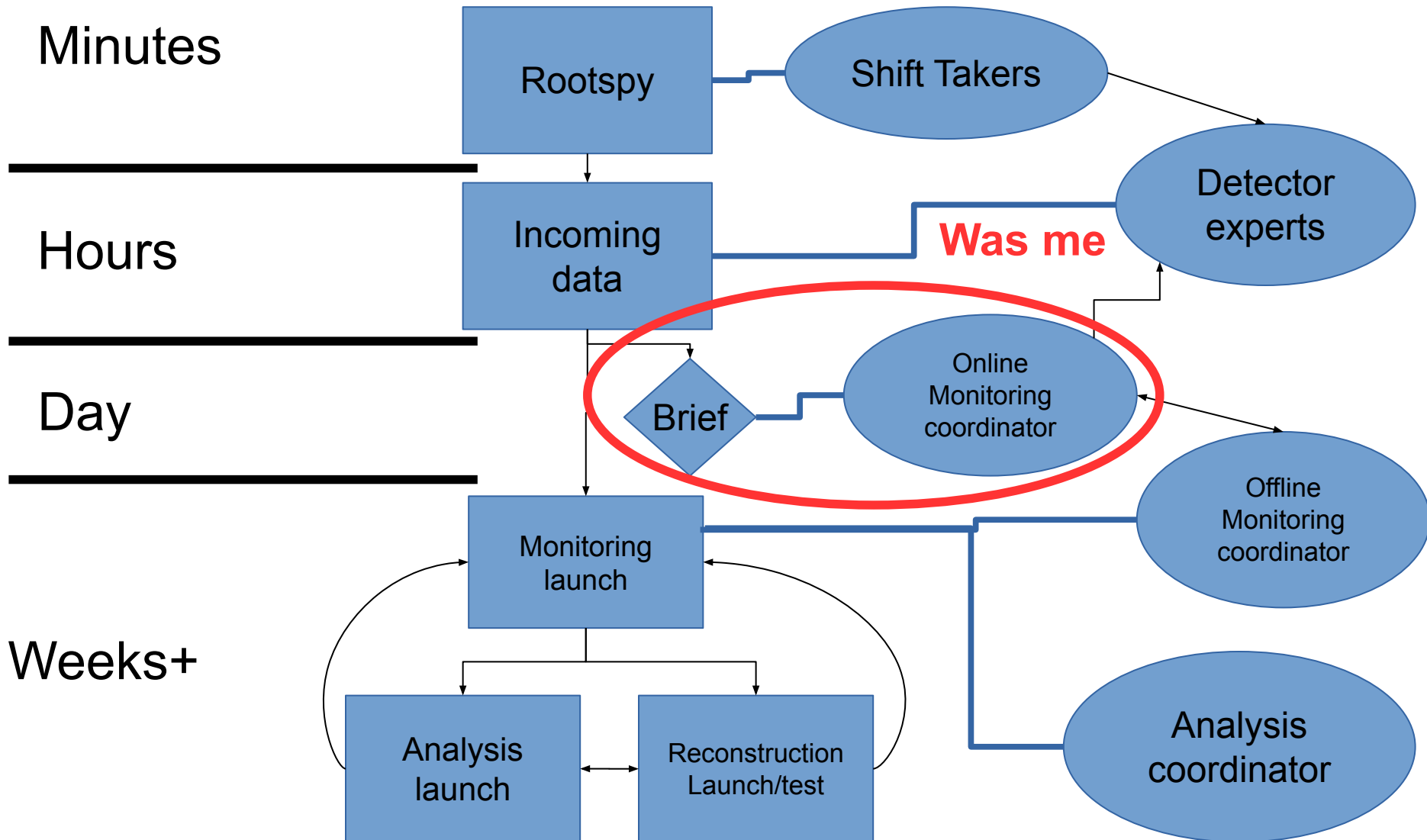




# AI: Data Quality Monitoring

Thomas Britton

# From DAQ to ANA



# The Challenge

- **Every run** produces an initial **22 plots**. More thorough monitoring is performed offline and produces **109 plots**. With a run lasting **~3 hours** every day there are **between ~175 and 875** plots to look at.
  - To preserve sanity I looked at closer to 175 plots, but there is no reason a machine couldn't aid in looking at all of them...
- Often times a single plot being “off” is not an indication of problems. Need to look at all the plots to determine cause and severity
  - Trigger studies: Often look like big problems but are not. Can be hard to catch when shift logs have scant details

# Let's Play

Ladies and gentlemen: the story you are about to hear is true. Only the names have been changed to protect the innocent.

\*Logbook searches not included



## Questions to keep in mind

What is the problem? Is it isolated? Is this due to running conditions or a test?

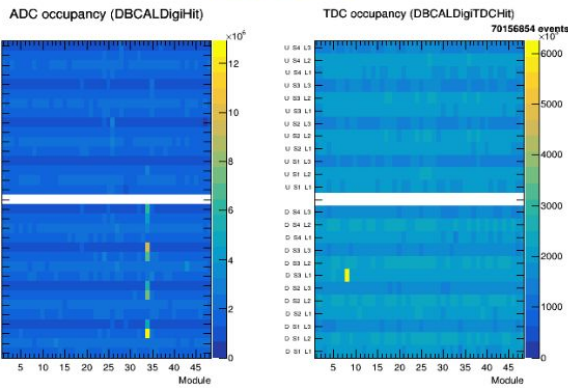
Should we stop taking data/call the shift crew?

**It's 5:15am the email has just come in...**

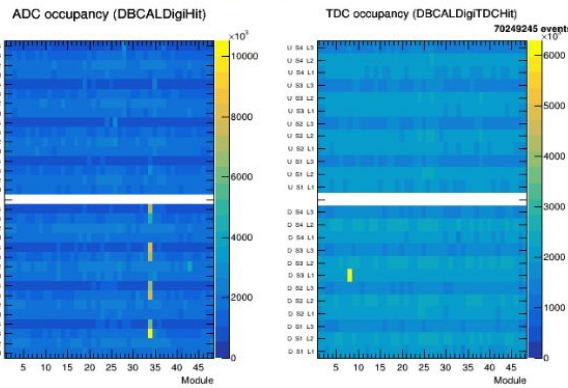
# Let's Play

Anomaly detected! BCAL occupancy changed

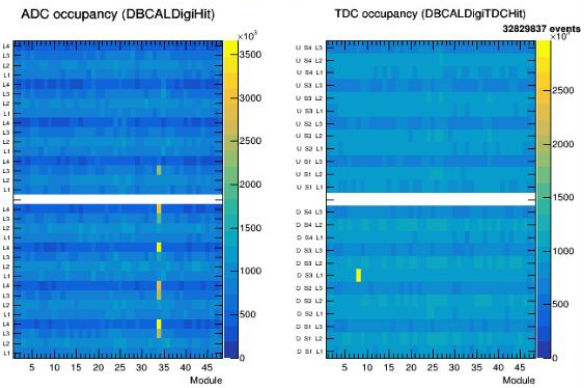
Run 042382 info



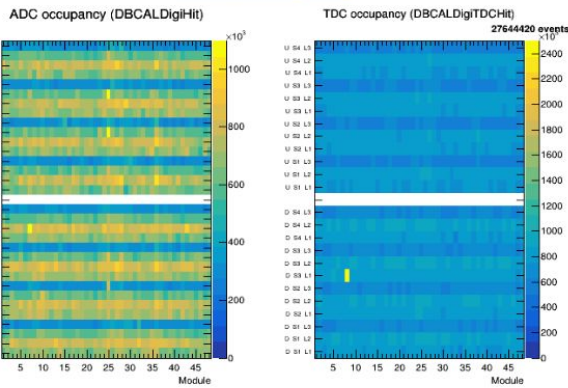
Run 042381 info



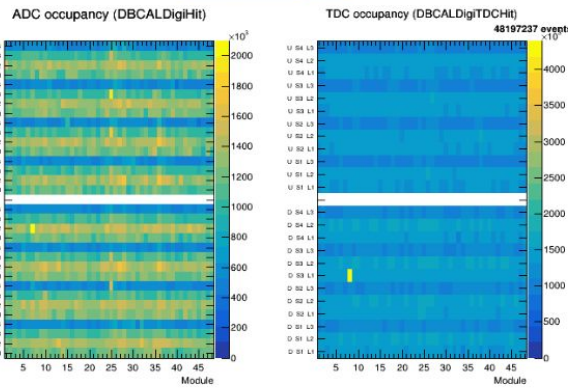
Run 042374 info



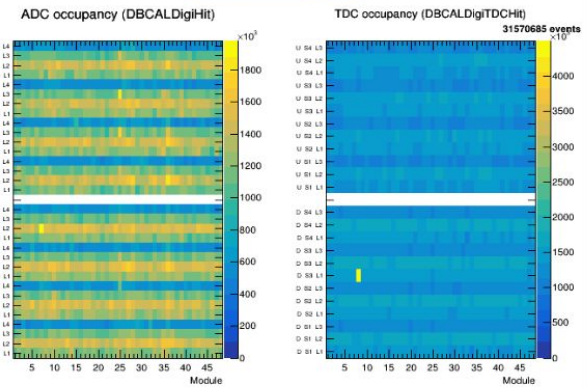
Run 042352 info



Run 042351 info



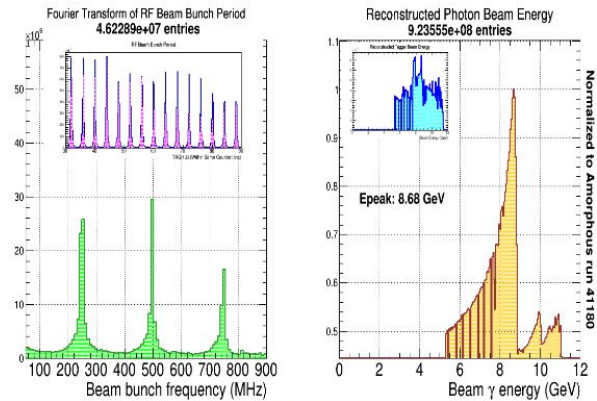
Run 042350 info



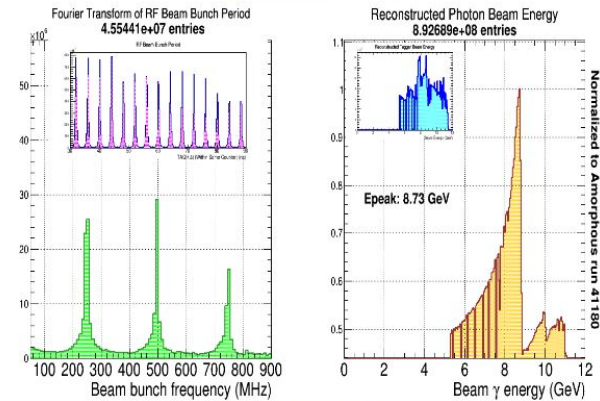
# Let's Play

Nothing looks really off with the beam

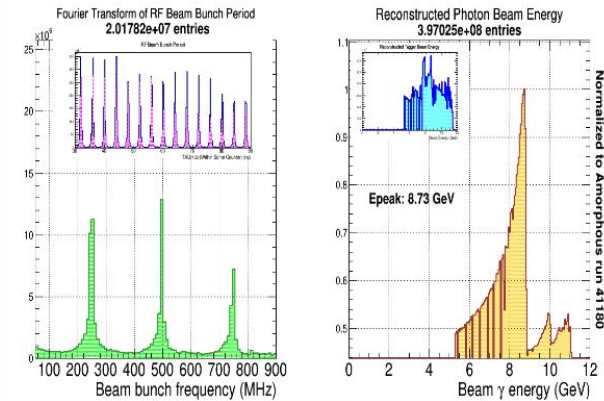
[Run 042382 info](#)



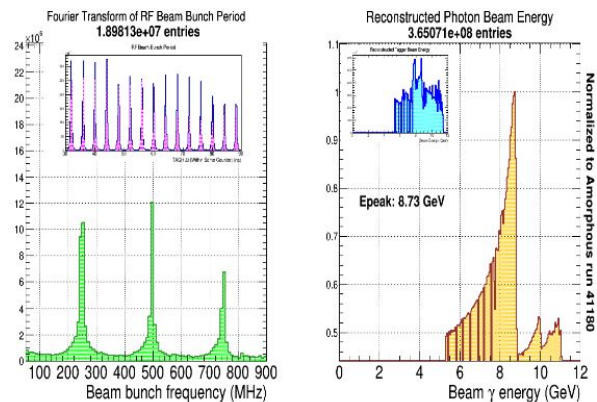
[Run 042381 info](#)



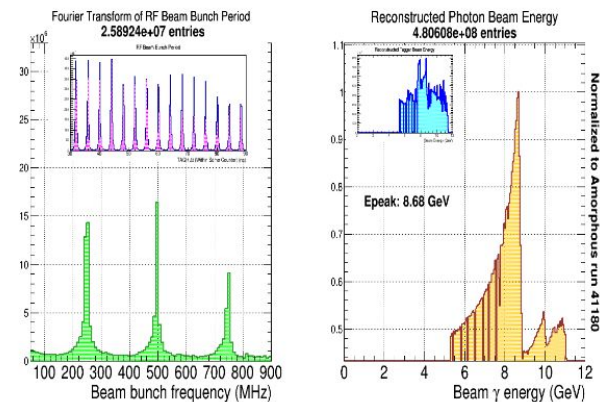
[Run 042374 info](#)



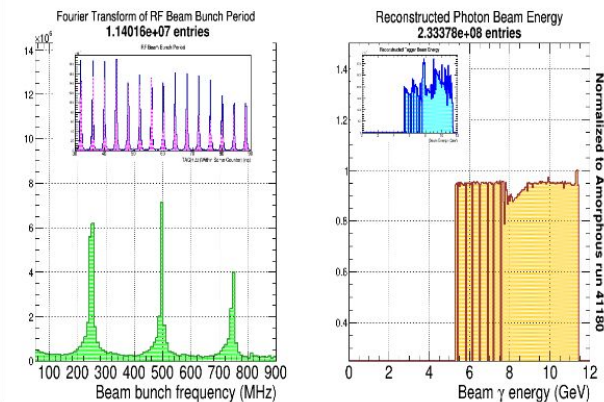
[Run 042352 info](#)



[Run 042351 info](#)



[Run 042350 info](#)

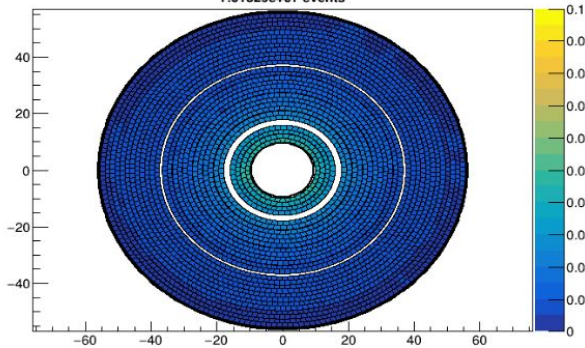


# Let's Play

CDC "dims"....hmmmm. Did you remember 42351 from the BCAL looked normal? Probably uncorrelated.....

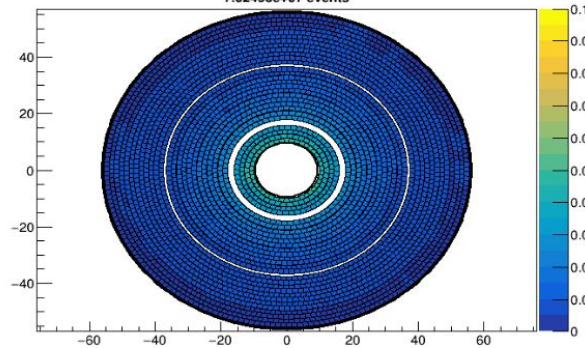
[Run 042382 info](#)

CDC Occupancy  
7.01829e+07 events



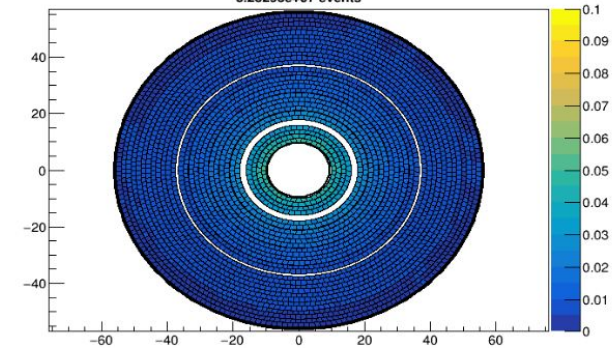
[Run 042381 info](#)

CDC Occupancy  
7.02493e+07 events



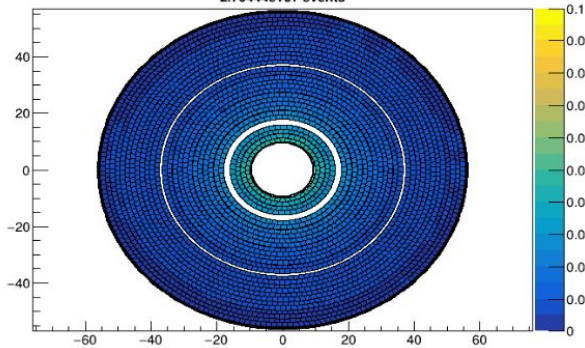
[Run 042374 info](#)

CDC Occupancy  
3.28298e+07 events



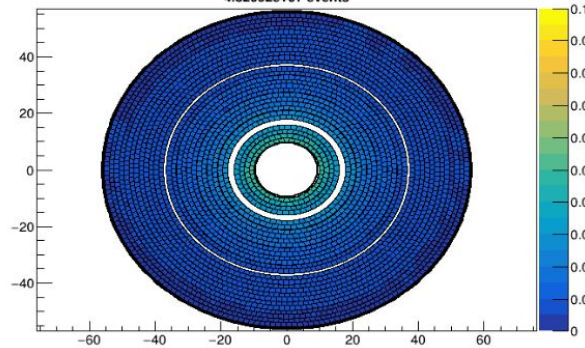
[Run 042352 info](#)

CDC Occupancy  
2.76444e+07 events



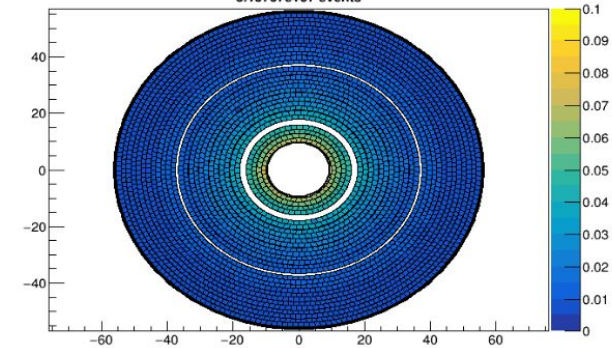
[Run 042351 info](#)

CDC Occupancy  
4.82092e+07 events



[Run 042350 info](#)

CDC Occupancy  
3.15707e+07 events

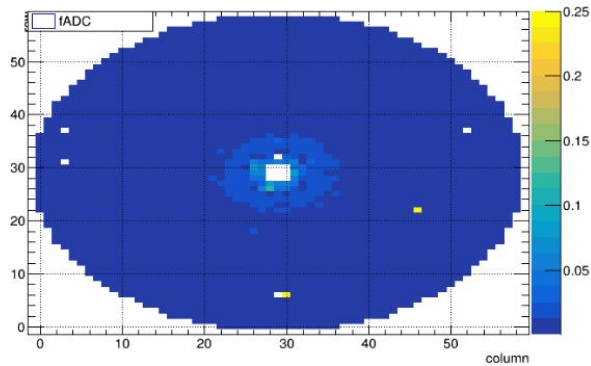


# Let's Play

FCAL looks normal. Note a few new dead modules....

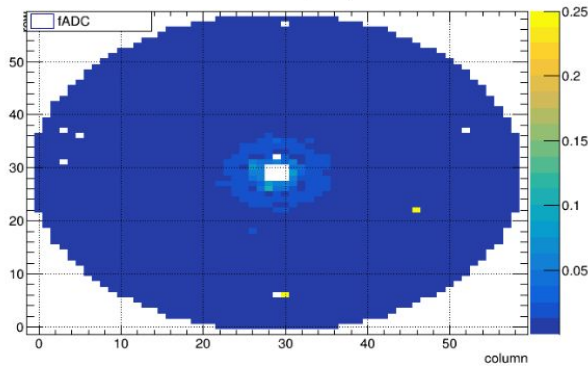
[Run 042382 info](#)

FCAL Occupancy



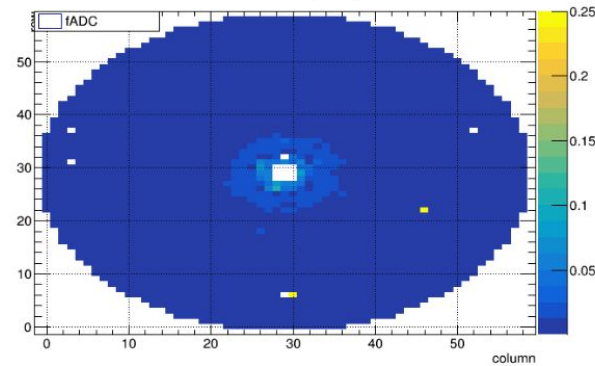
[Run 042381 info](#)

FCAL Occupancy



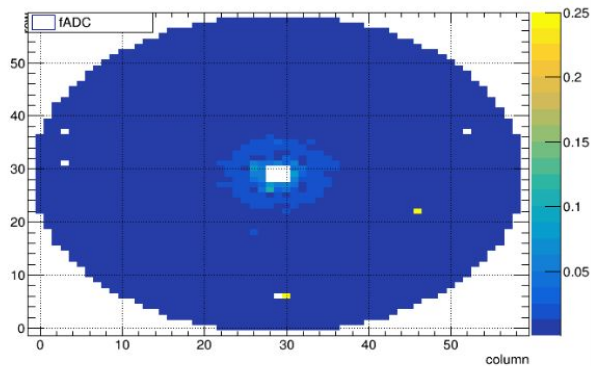
[Run 042374 info](#)

FCAL Occupancy



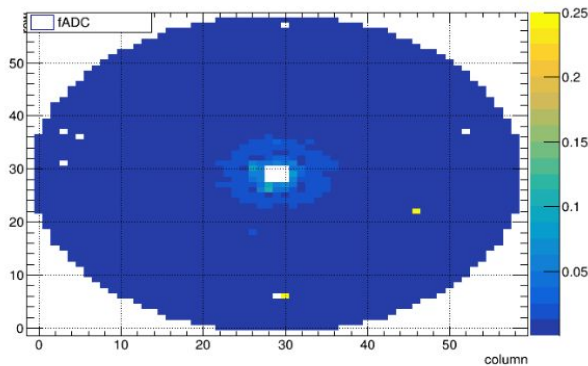
[Run 042352 info](#)

FCAL Occupancy



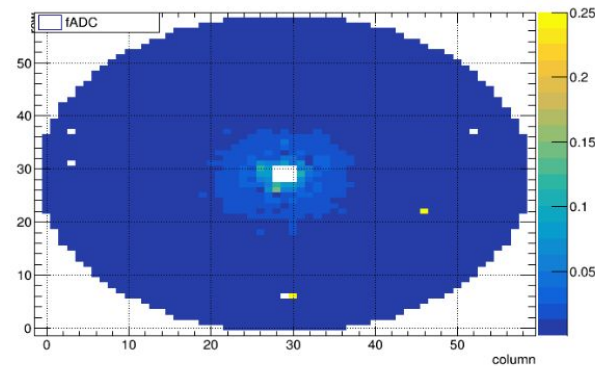
[Run 042351 info](#)

FCAL Occupancy



[Run 042350 info](#)

FCAL Occupancy



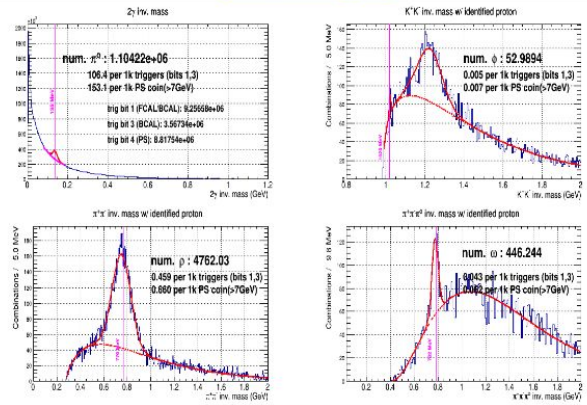


# Let's Play

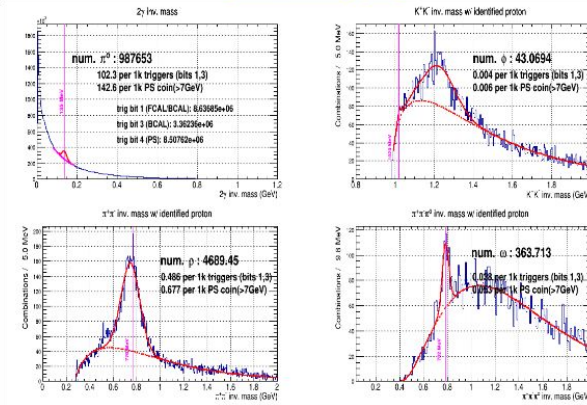
13 plots later....

Physics signals still in tolerance

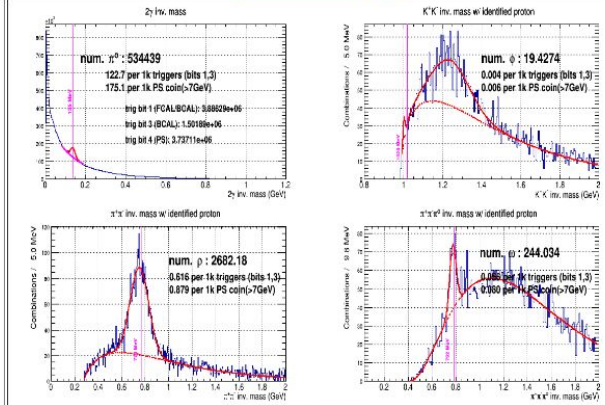
**Run 042382 info**



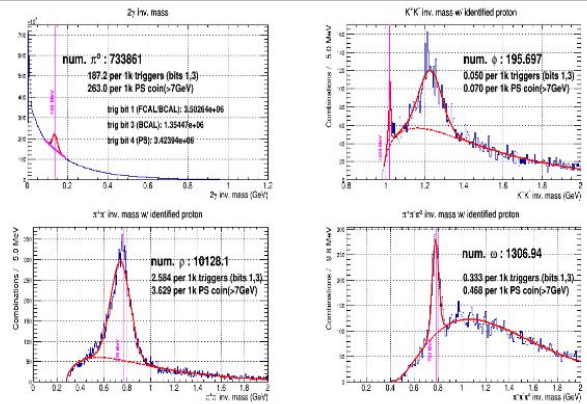
**Run 042381 info**



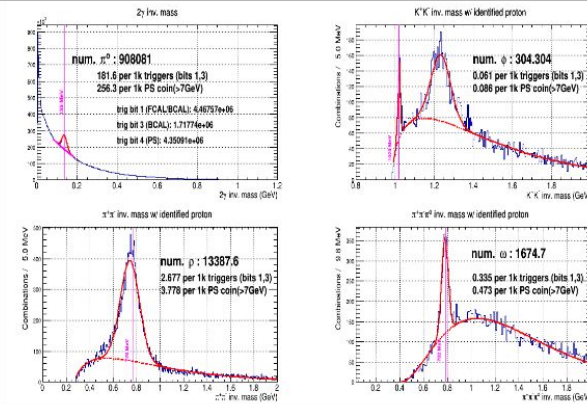
**Run 042374 info**



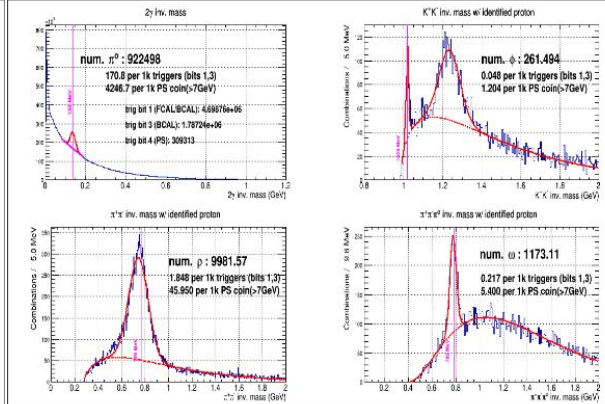
**Run 042352 info**



**Run 042351 info**



**Run 042350 info**

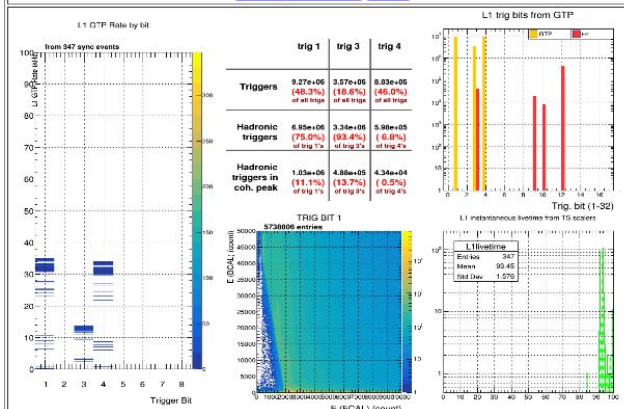


# Let's Play

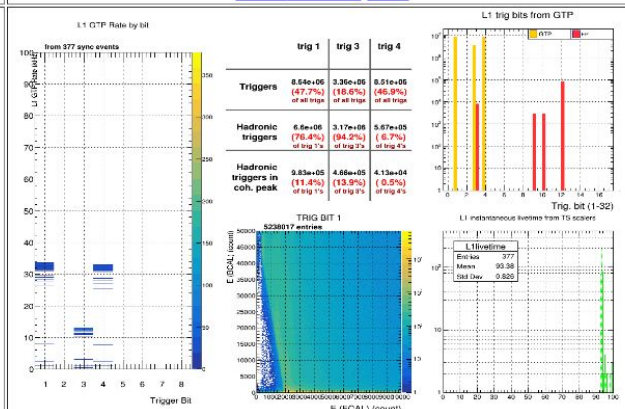
5 plots later...

Increased dead time at 42351

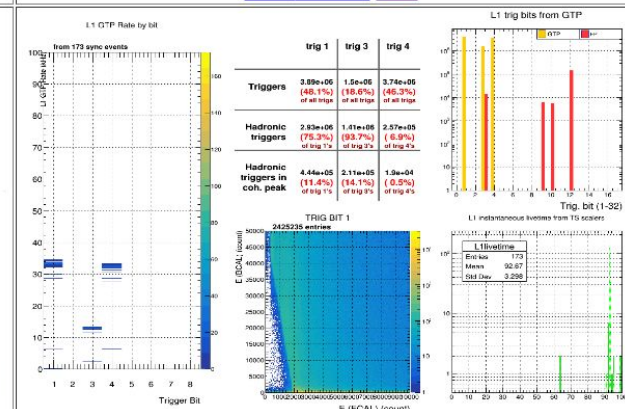
**Run 042382 info**



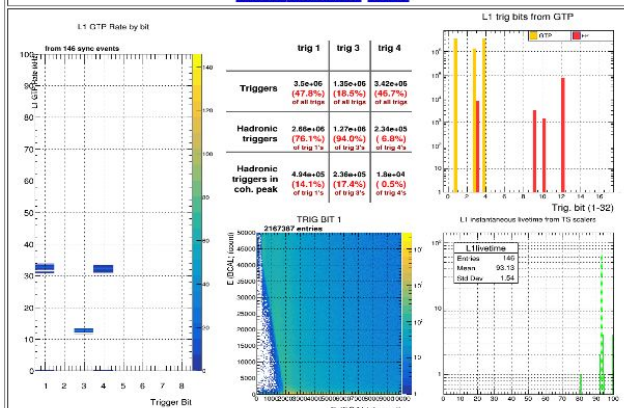
**Run 042381 info**



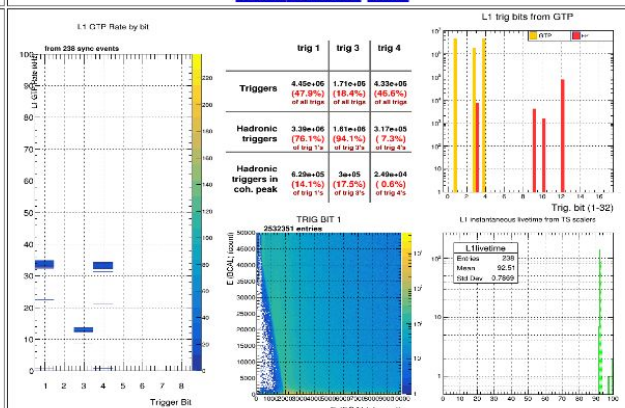
**Run 042374 info**



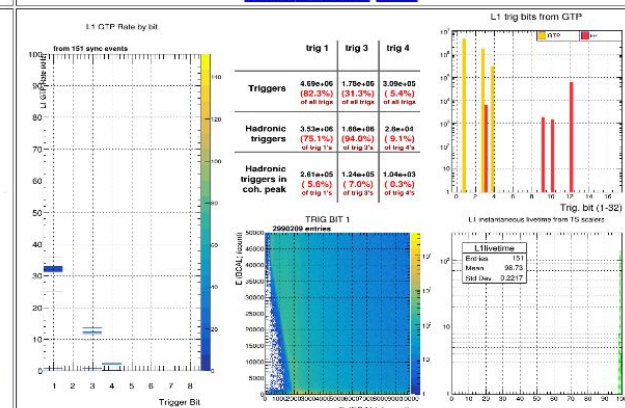
**Run 042352 info**



**Run 042351 info**

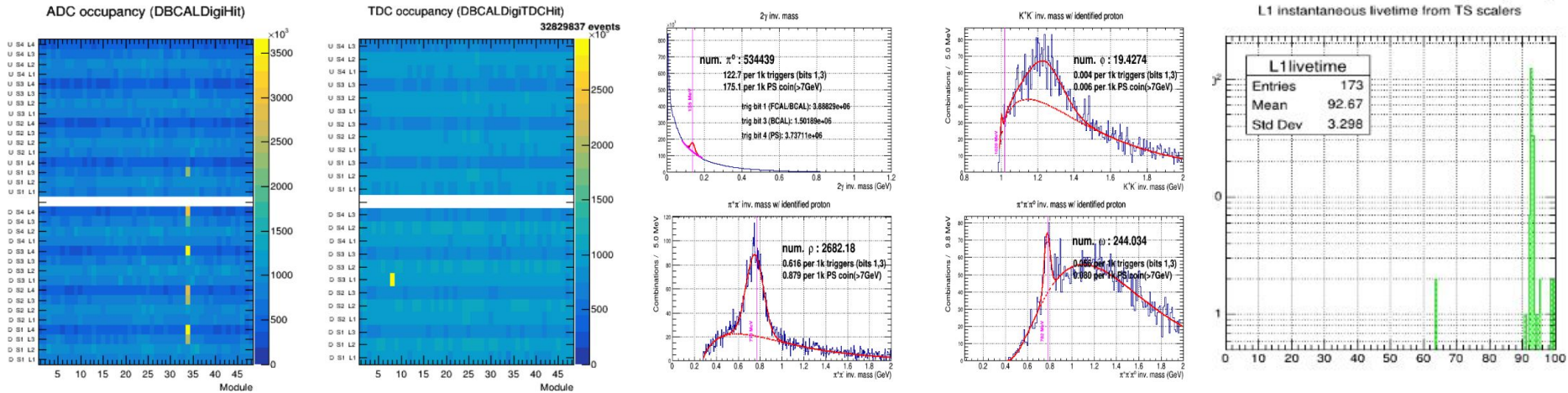


**Run 042350 info**



# Conclusion

Likely, several BCAL modules have become noisy contributing to an increased deadtime. The reason for occurrence is unknown, but it likely happened during run 42351. It is not significantly impacting physics data. Detector experts should be notified and further investigation should be performed. An opportunistic access to the hall may be needed



## Epilogue:

In roughly 24 hours the noisy components were found and replaced. Normal operations resumed

# AI

- As demonstrated, what is involved is a lot of pattern recognition and paying attention to the correlations between plots coupled with domain knowledge
  - When plot A looks like X plot B should look like Y but it looks like Z and by looking at plot C probably Q is wrong...
- This is precisely where AI can shine
  - Large, multi-dimensional datasets full of correlations

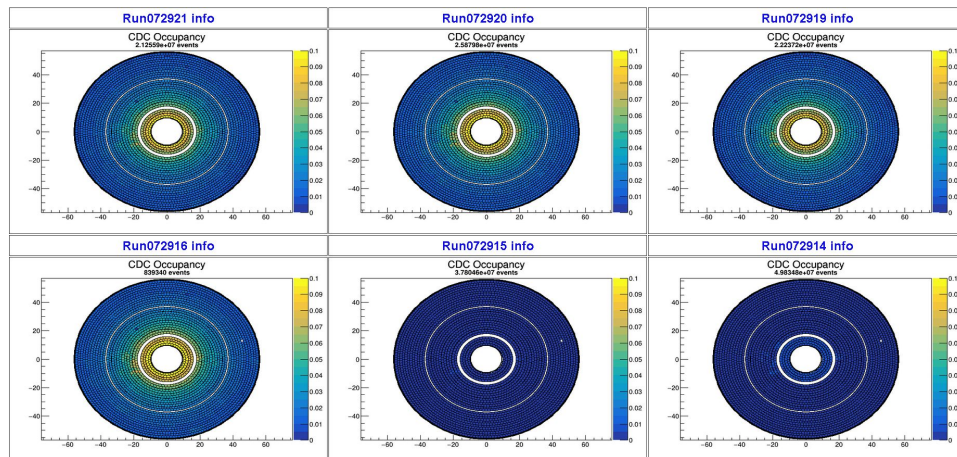
# Considerations

- Many considerations in developing and deploying AI models
  - Data collection:
    - How will you get your training data?
    - How will test data be given?
  - Labeling (in the case of supervised learning):
    - How will you label all of your data?
    - What will that cost?
  - Model development:
    - What requirements does the model need to fulfill? (e.g. inference time, error rates etc)
  - Sociological:
    - How will people interact with the system?
    - Are there any barriers to deployment from the human side?

# Data Collection

- GlueX already had monitoring plots from all the detectors from all the runs in image format
- Using these is slightly wasteful in terms of data needed for inference
- Get to use all of the developments in machine vision
  - Convolutional Neural Nets

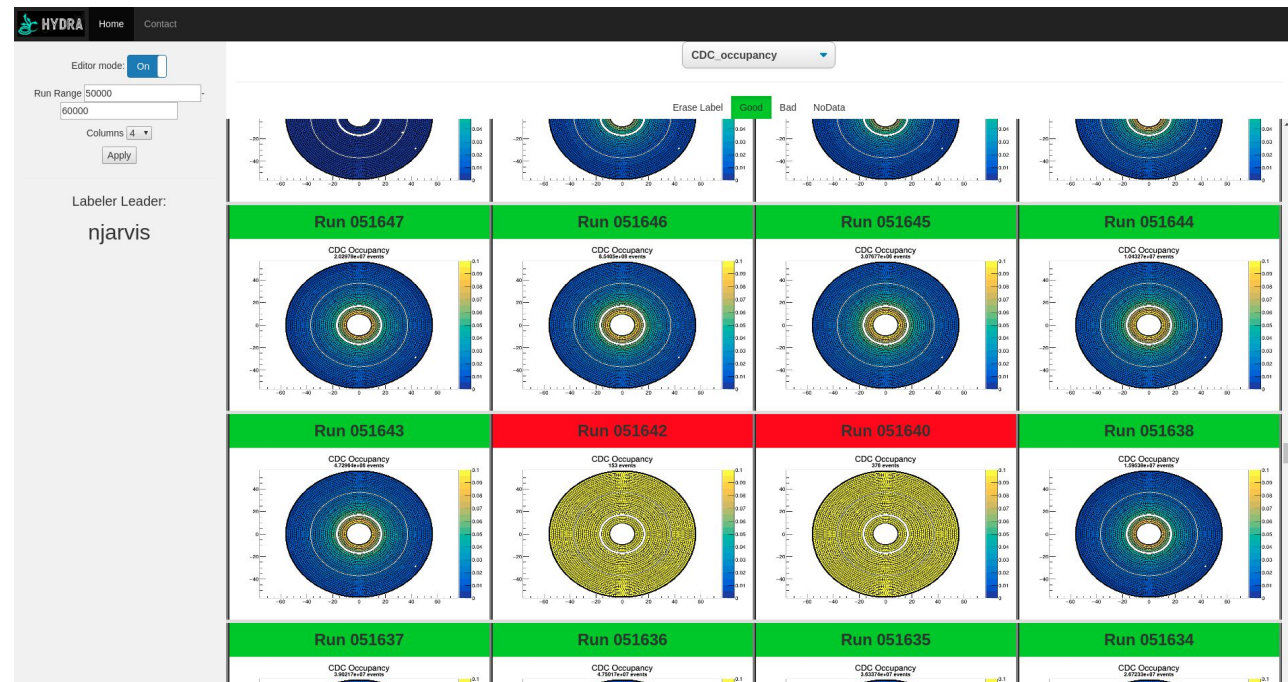
- Cuts down on startup time



# Labeling

- Needed an easy way to label vast numbers of images with a low barrier of entry
- Settled on a webpage

**Able to label hundreds of plots in minutes**



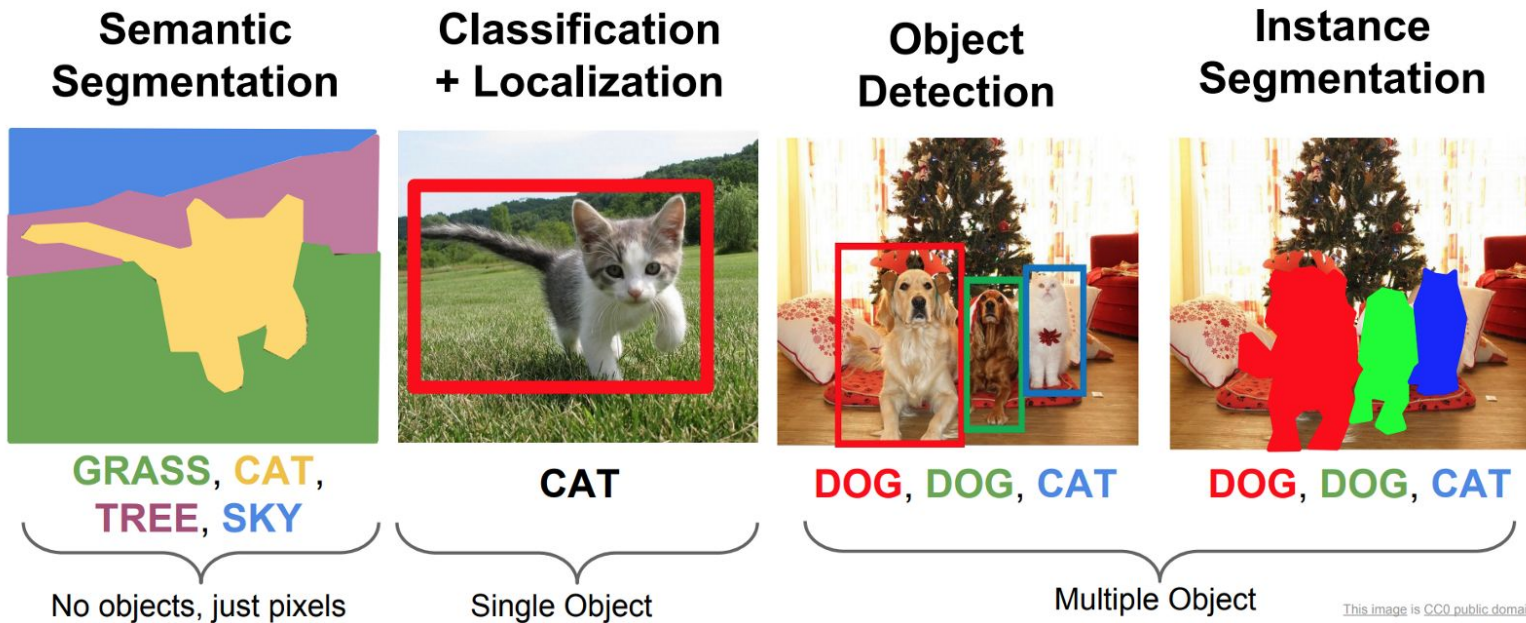
# A Challenge

- Labels are horribly unbalanced (Think ~5% bad).
  - This is ultimately expected as typically detectors are designed to work..
- This can cause problems as the machine may not have enough examples to adequately optimize the weights to classify the lesser labels
  - “It is more advantageous to get really good at identifying cats over dogs or boats as the distribution appears to be heavily skewed to cats”
- Label balancing is the process of repeating examples of lesser labels randomly to make the number of examples equal the biggest label
- This can lead to unnecessarily long training times as it increases, artificially, sample size
  - Strategic undersampling? It doesn't take many all white images for it to perfectly classify the image as NoData. Maybe curation?
  - Random undersampling of the larger categories to match the smaller categories?
  - Weight specific examples more than others?



# Model Development

- Because I decided to use the already available images the algorithm choice was pretty clear
  - **Convolutional neural nets**



This image is CC0 public domain

# CNN (math)

$$G[m,n] = (f * h)[m,n] = \sum_j \sum_k h[j,k]f[m-j,n-k]$$

## Kernel Convolution Example

Input Image

10	10	10	10	10	10
10	10	10	10	10	10
10	10	10	10	10	10
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0

Kernel

1	2	1
0	0	0
-1	-2	-1

=

Feature Map


\* stride of 1

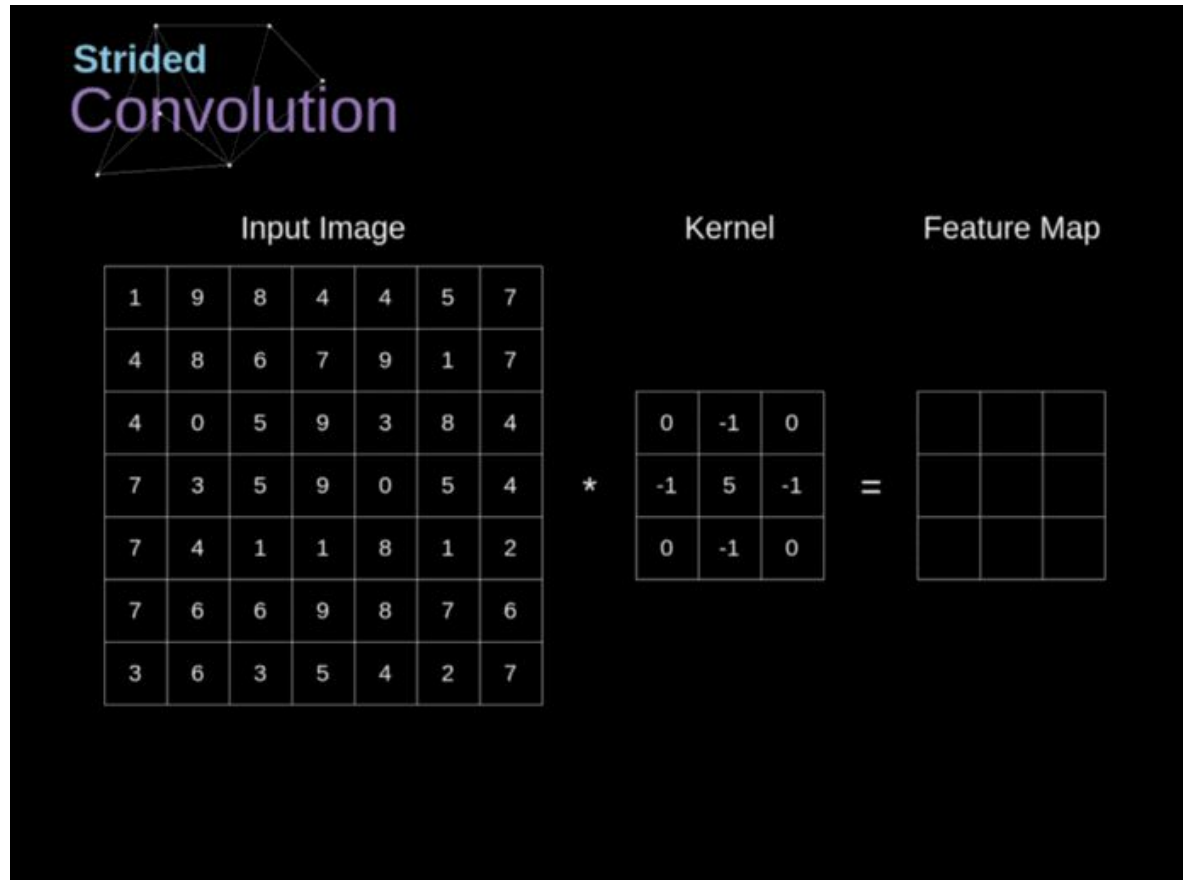
# CNN (kernels)

<https://towardsdatascience.com/>

- Different kernels produce different outputs (features).
  - so you are training a CNN layer to extract different features from the data



- How far to move your matrix window



# CNN (pooling)

- Many types: max, min, avg...

**Max Pooling Example**

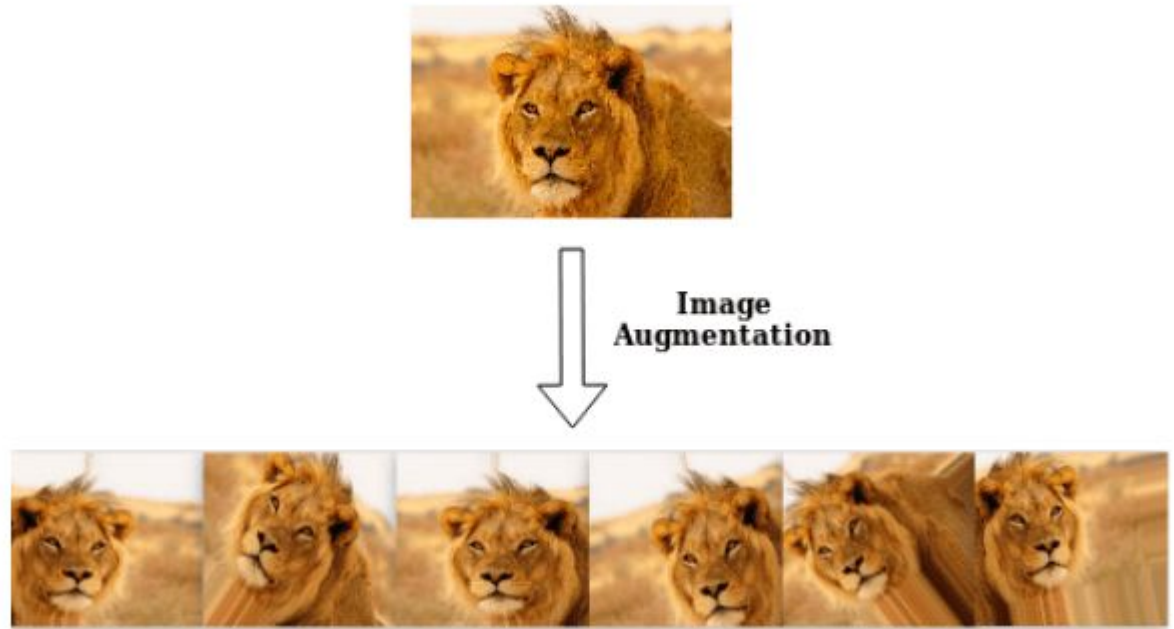
Input

3	0	1	5	1	3
5	7	3	4	4	6
7	7	1	8	3	5
6	1	7	0	0	5
0	4	5	5	7	2
3	2	0	2	0	2

Output


# Image Augmentation

- increases dataset and can make a model more accurate
  - can include color shifting etc.

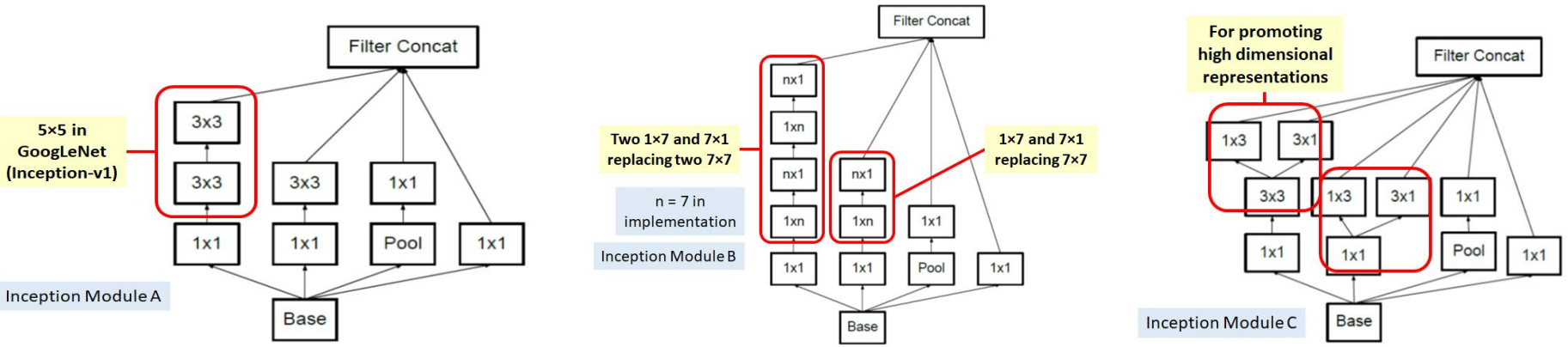


\*We don't actually want this for looking at our images as the images are all very very uniform

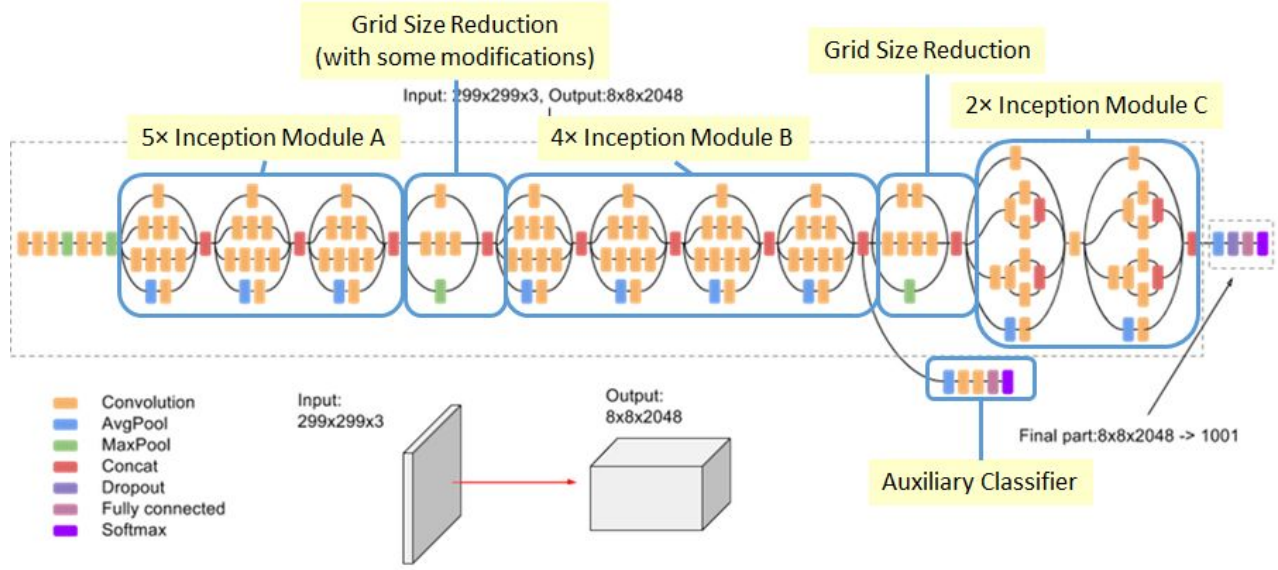
...well maybe color shifting to capture new scales.....

<https://towardsdatascience.com/>

# The Inception v3 Network

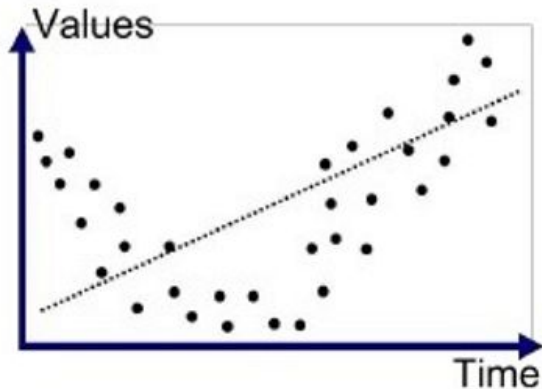


Already in Keras and Tensorflow

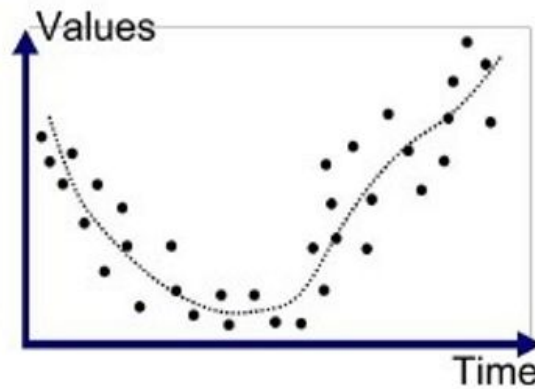


# Overfitting

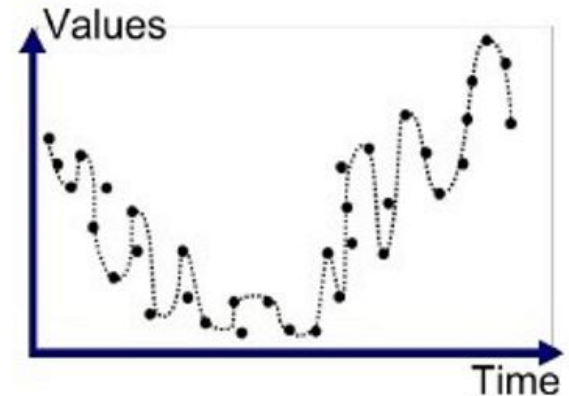
- Overfitting a big question



Underfitted



Good Fit/Robust



Overfitted

- We have the ability to generate “fake” data
  - May not be feasible to train on
- If the points above were the totality of possibilities...
  - the right most model would be a great, albeit unsatisfying, model



# Introducing Hydra

- Hydra aims to be an extensible framework for training and managing A.I. for near real time monitoring
  - If you need it to tell a dog from cat I can have hydra do that, without system modification, now
- Most importantly, Hydra allows me to embrace my inner sloth:

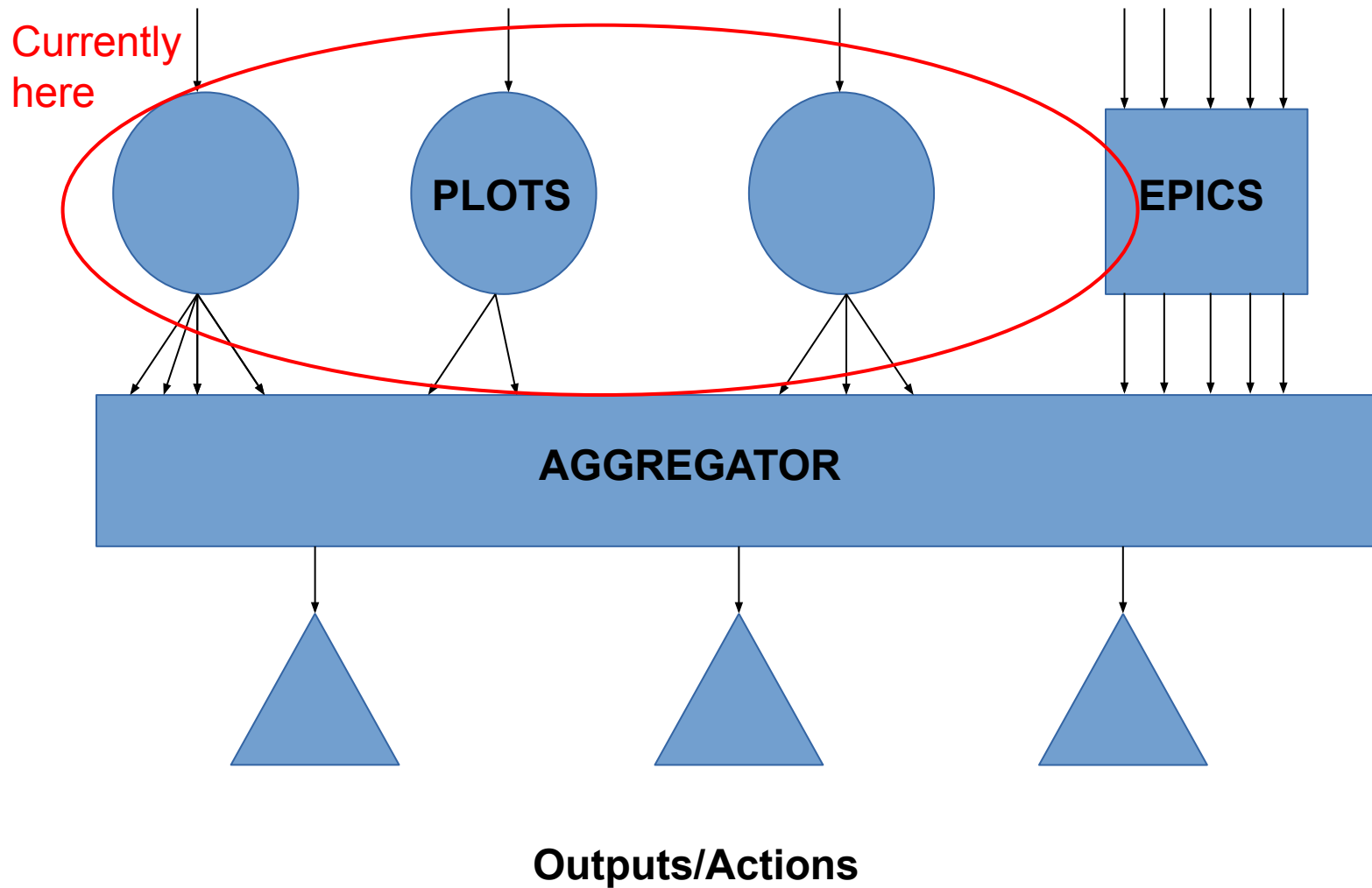


Koboldpress.com

NOPE  
NOT TODAY



# Future Topology



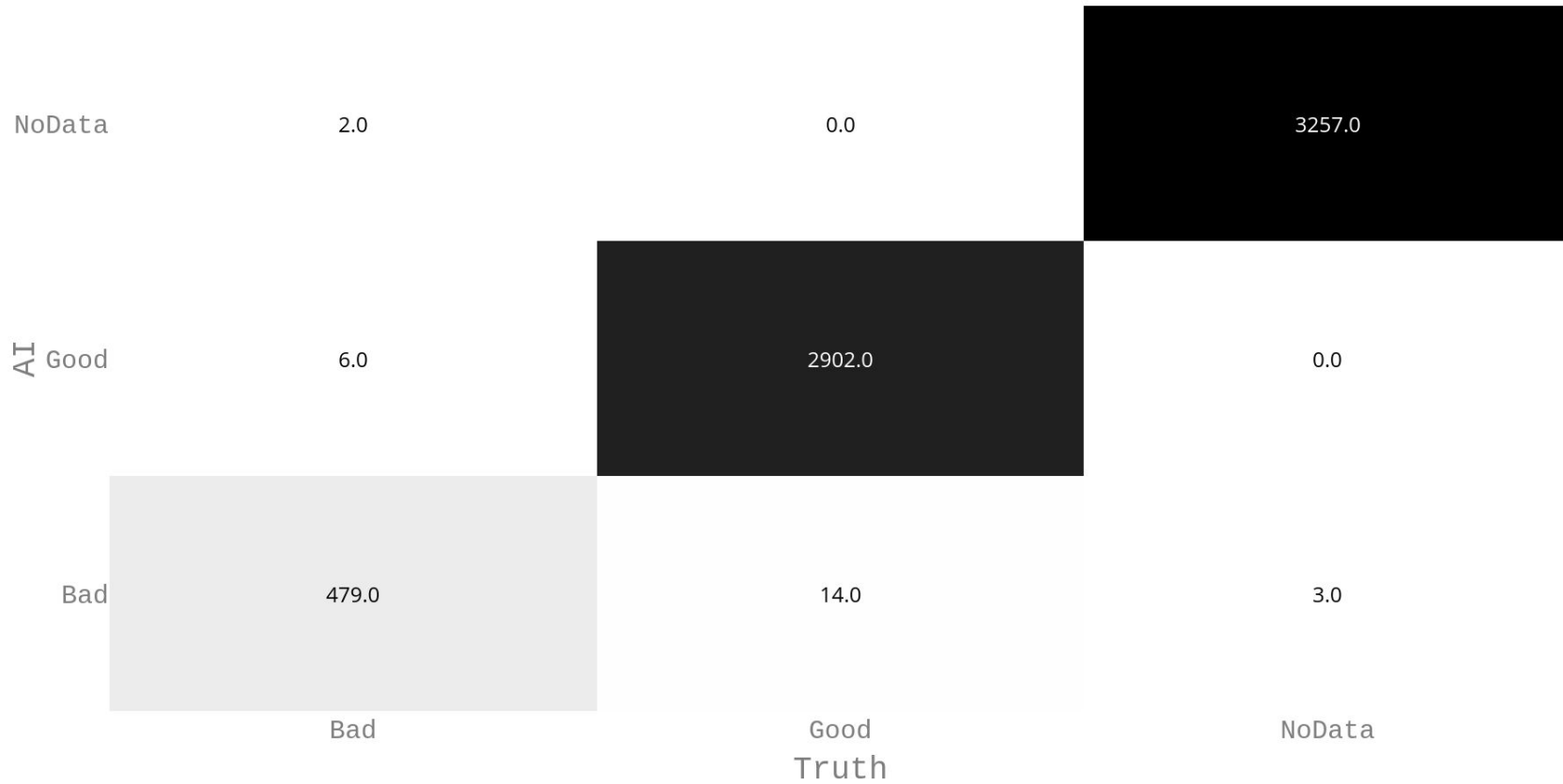
# Predictions

- Prediction performed by a script call
  - **3 modes:**
    - **Datum** (full path to image)
    - **Directory** (non-blocking attempts to analyze all data it finds)
    - **Crawl** (over all saved data in the DB)
- Writes and publishes (0MQ) **JSON report** including each label the model knows about and it's confidence in the label for each piece of data analyzed
  - Can be used in downstream processes
- Automatically captures examples for future training/analysis
  - "Bad" examples
  - Uncertain plots
  - Every Nth example



# False Positives/Negatives

CDC\_occupancy Confusion Matrix

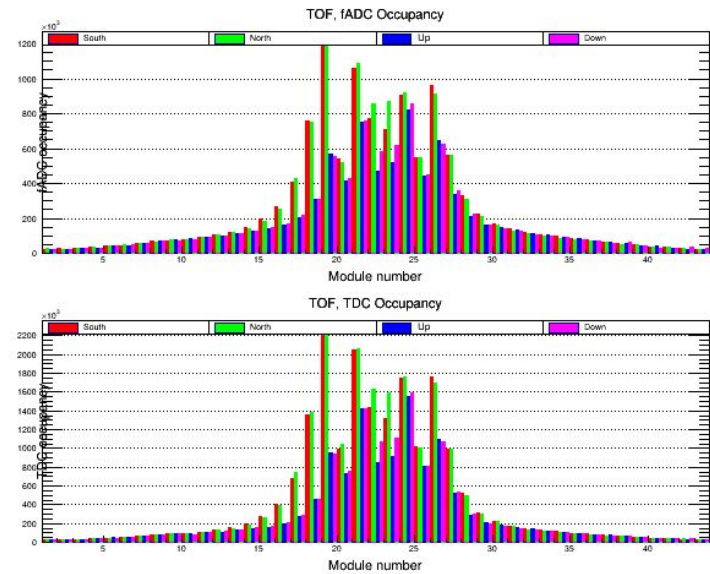
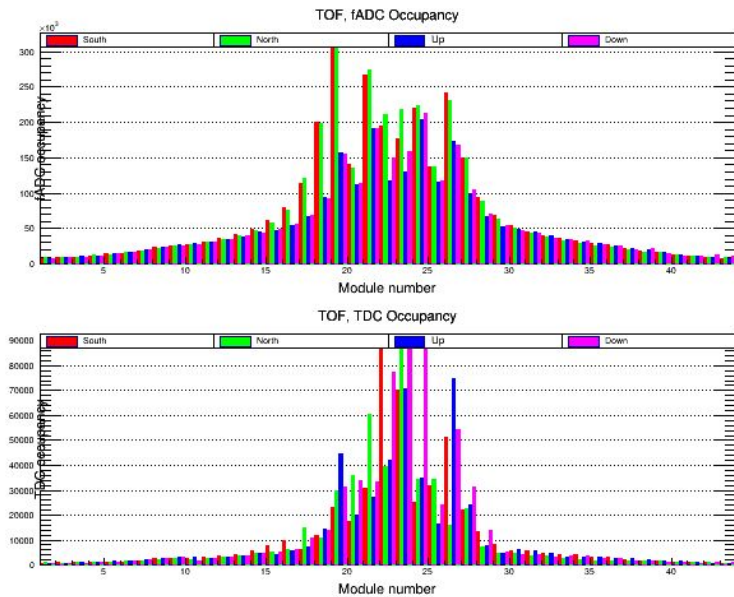


# The Worst Case

CDC\_occupancy Confusion Matrix

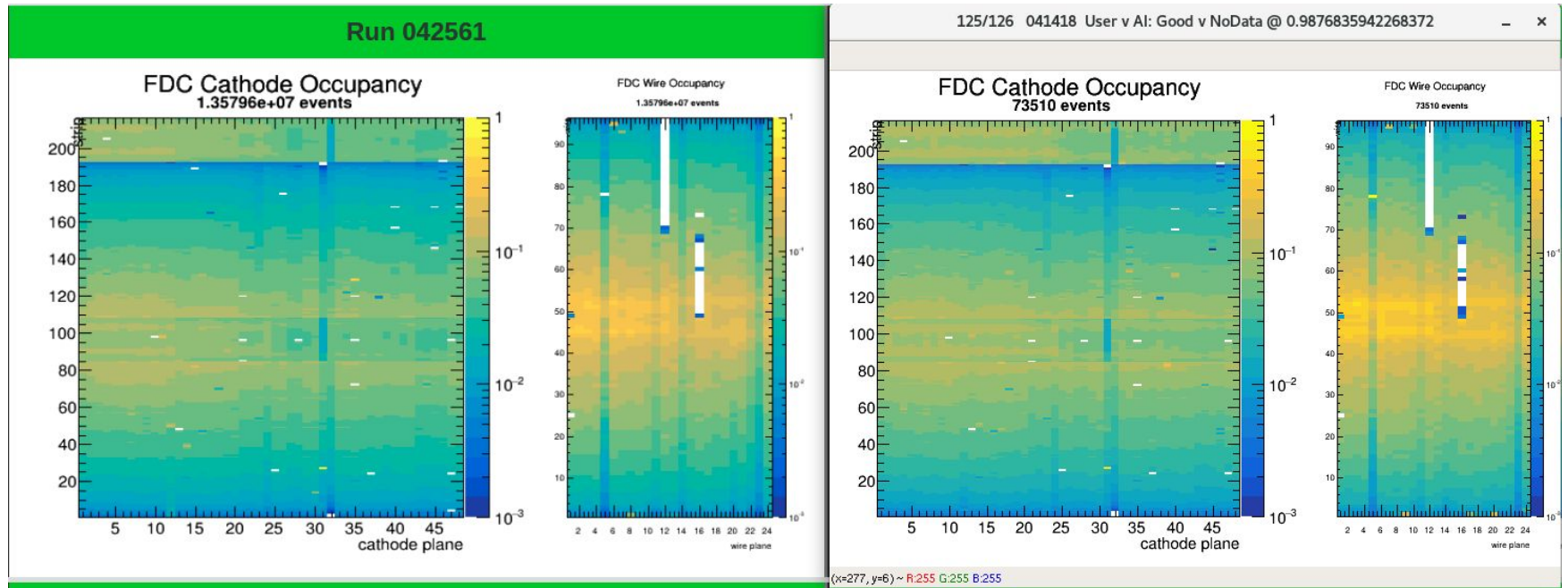


# An Anecdote



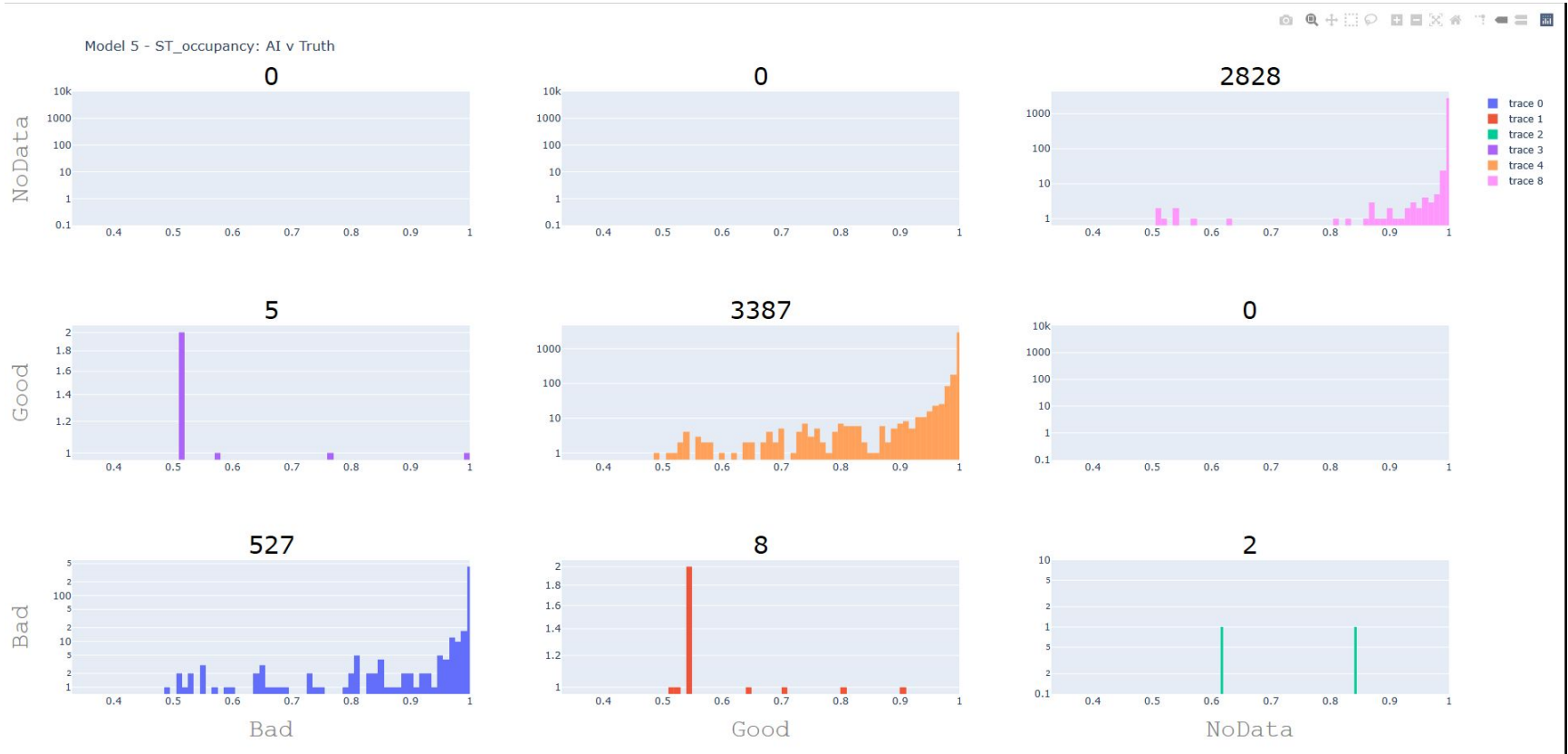
- Both of these look “good” at first glance (both initially labeled good)
  - The one on the left is actually bad (the A.I. caught it)
- A.I. seems to be able to look at subtle differences in shape and **maybe even “read” the y-axis**

# Another Anecdote



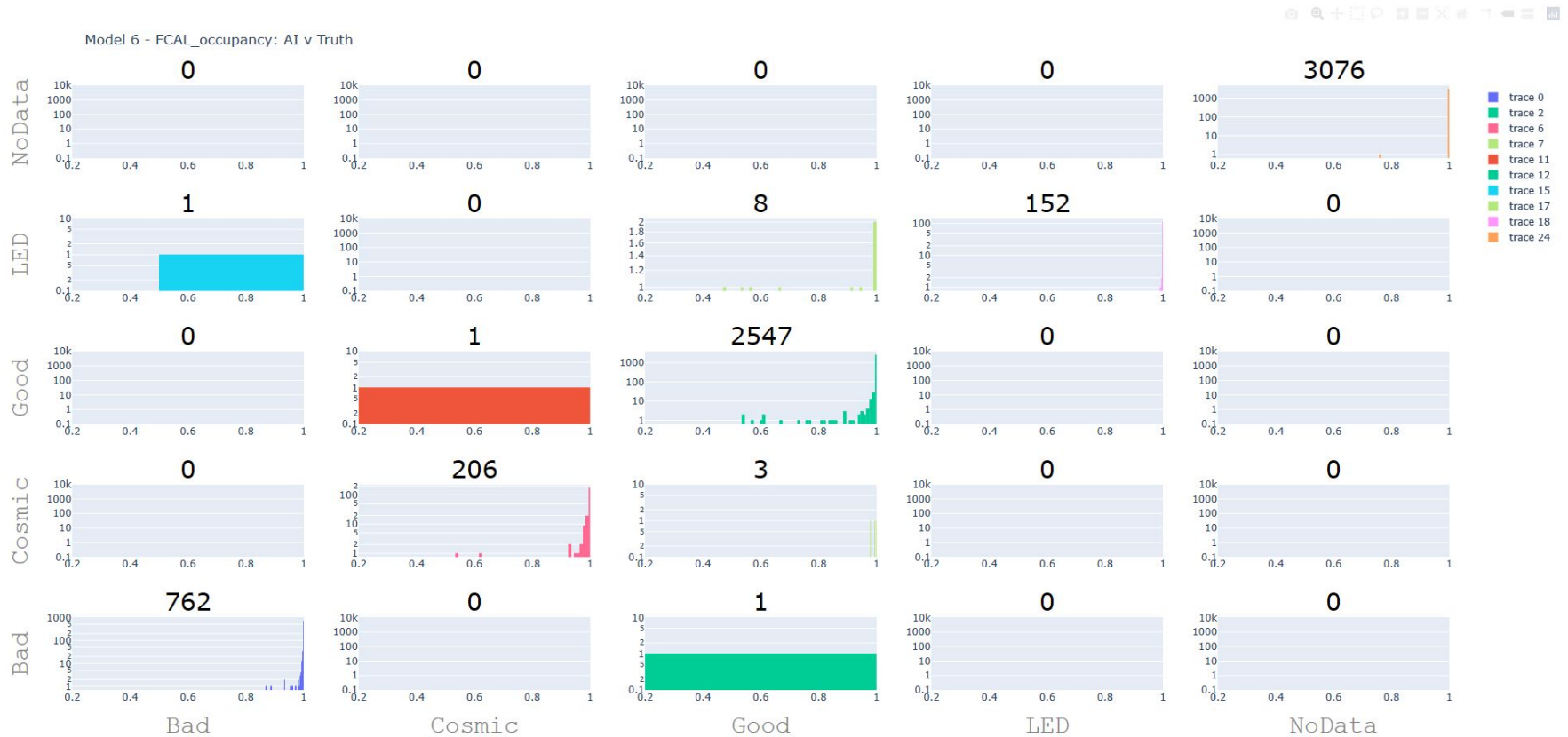
- The labeler was instructed by the detector expert to label any plot containing **fewer than 100k events as “NoData”**. This is one example of several in which the labeler labeled as “Good” and the **A.I. predicted “NoData”**...the true label given the number of events

# Preliminary Results (Start Counter)





# Preliminary Results (FCAL)



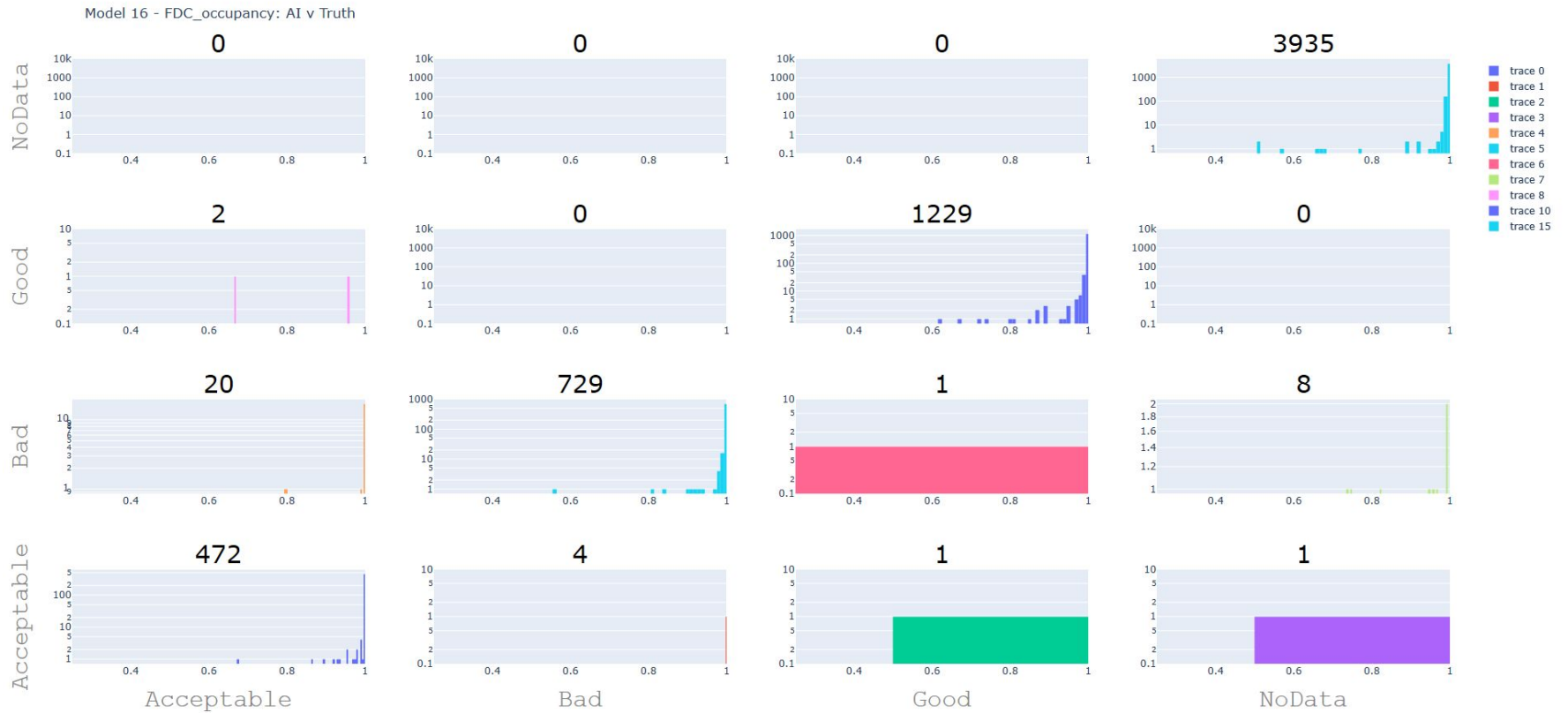
# Preliminary Results (DIRC)



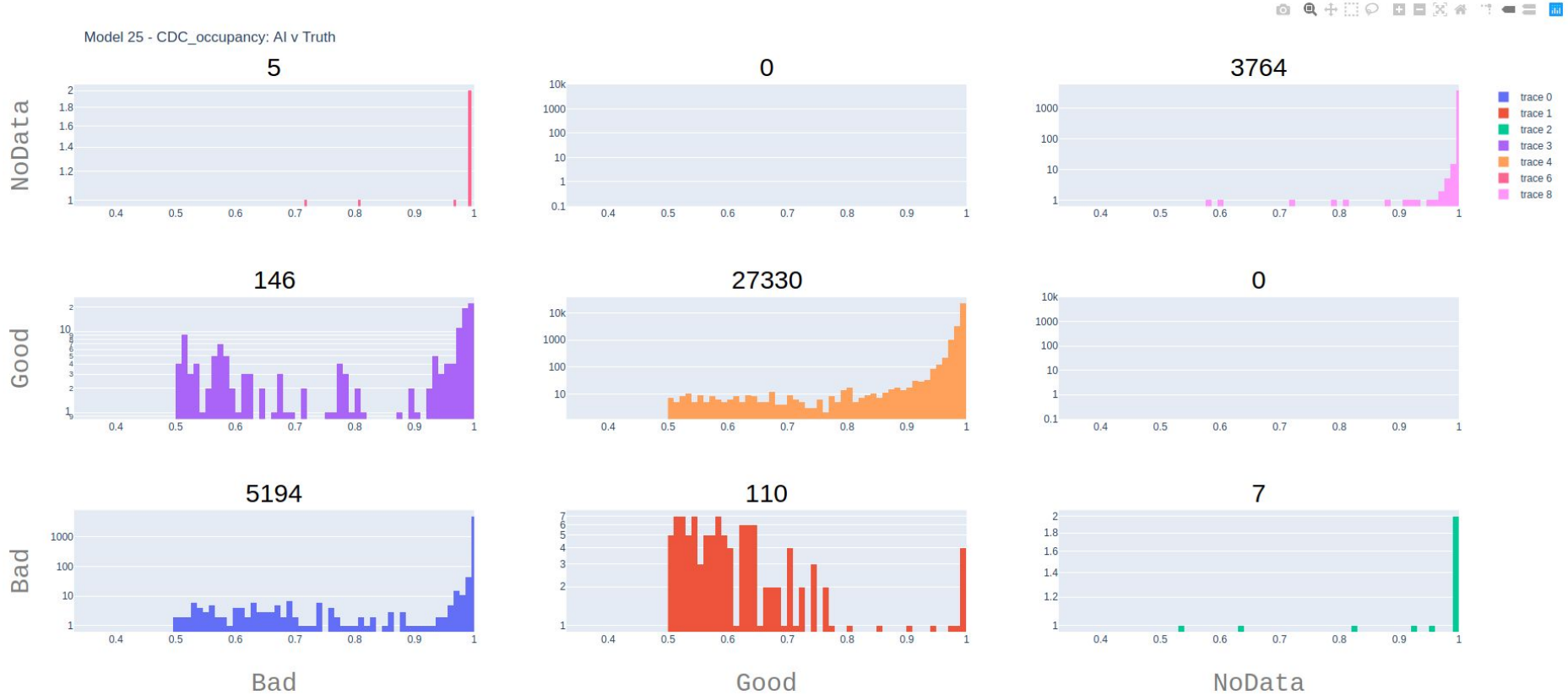
# Preliminary Results (BCAL)



# Preliminary Results (FDC)



# Preliminary Results (CDC)



# Preliminary Results (TOF)



# Some Things I Noticed on the Way

- Tensorflow 1.14 and their partial bundling of keras is weird
  - I ended up having to use some of the tf.keras functions and some functions from keras standalone to get it to work
    - I've since started trying to move to TF 2.2 which is much better in this regard
- It looks like TF does some library loading when the model is loaded and more the first time predict is called for a model. Because I have many models that get used "on-the-fly" the predict times were terribly long (first predict in **4 seconds**)
  - Had to run a **prediction over an arbitrary image** on model load. Then all first predicts became fast.

```
def PreloadModels():
    then=int(time.time()*1000.0)
    print("Preloading Models")
    data_to_analyze_q="SELECT * FROM Plot_Types where Active_Model_ID IS NOT
    dbcursor.execute(data_to_analyze_q)
    data_to_analyze= dbcursor.fetchall()
    for d in data_to_analyze:
        print("Loading head for "+str(d["Name"]))
        chunked=-1
        if(d["IsChunked"] == 1):
            chunked=1
        LoadModel(d["Name"],d["FileType"],chunked,d["Active_Model_ID"])
        headkey=d["Name"]
        #print("BREAKFAST: "+str(chunked))
        if chunked==1:
            headkey=headkey+"_1"
        #print("BREAKFAST HEADKEY: "+headkey)
        Breakfast(headkey) ←
    now=int(time.time()*1000.0)
    print("Finished preloading all active models in: "+str(now-then)+" ms")
```

# Hydra's Breakfast





# Conclusion

- It is possible to employ computer vision techniques to handle many of the monitoring procedures
- A.I. is trivial to setup and get running
  - The devils are all over the details and can cause a lot of problems
    - Maintaining provenance
    - Interpretability
    - Validation
    - Optimizations
- Hydra has been deployed and already has taken some burden off shift crews. A big benefit during single person shift taking during COVID
- Lots of avenues for further research and improvements
  - Improved generalization via curation?
  - HPO
  - Increased analysis (new plots/data types)
  - Etc etc

# Hot Off the Presses

- The prevailing wisdom is that, especially with deep learning, you need more data
  - This may be generically true but in certain situations less is more
- Label imbalance occurs when one or more labels are much greater in quantity than others
  - This can lead to poorly trained models as the optimum is heavily skewed
  - can input the smaller populations many times to “fake” a larger sample
    - this leads to longer training times
- Because the variance in “Bad” data is much greater than other labels we can undersample the larger samples to match the smaller sample
  - MUCH faster trainings for models that work just as well (or better)