

# Analyzing SiD PFAs: A path for improvement

Ron Cassell – SLAC  
ALCPG Physics and Detector weekly  
meeting  
10/18/07

# Disclaimer

- The opinions presented in this talk are those of the speaker, and not in any way intended to represent the opinions of an expert.
- Even those opinions attributed to another person or group of people may have been misinterpreted.

# Outline

- Brief PFA review
- What's the problem?
- How do we fix it?
- Ugly details
- Summary

# From Mat

What is the goal?

To produce lists of reconstructed final-state particles without cheating which are good enough to use in physics benchmarking & analysis.

This immediately throws up questions:

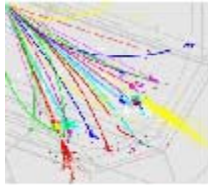
- Physics benchmarking: Which channels?
- Good enough: How good is that?
- Without cheating: What do we do in the meantime?
- Final-state particles: A whole other can of worms...
- Which detector? Digital or analog readout? etc etc etc

# What's the problem?

- Want a minimum of 4-5%  $dM/M$  rms90 for the delta Z mass in 500 GeV  $ZZ \rightarrow qq\nu\nu$  evts. Currently  $\sim 6.5\%$ .
- For 100 GeV light quark jets, see  $dE/E$  using rms90  $\sim 6\%$ .
- Why? Bad detector design? Bad algorithms? Bad idea to use PFA approach?

# Evidence it's NOT the PFA approach

- PFA approach: Mark Thompson reports Pandora PFA  $\sim 4.4\%$  dE/E using rms90 for 45 GeV jets, and  $\sim 3\%$  dE/E for 100 GeV jets.
- Although I haven't seen plots of delta Zmass from ZZ  $\rightarrow$  qqnnu events at 500 GeV, we can infer an rms90  $< \sim 4.5\%$ .
- (semi) independent confirmation of results from Marcel Stanitzki. (Very preliminary, results could change)



# Results

Configuration	n/sqrt(E)	Jet energy
LDC00Sc	30.5	45
LDC00Sc 5T	31.2	45
LDC00Sc 30 layer ECAL	32.4	45
LDC00Sc Sid-ish 4T	32.6	45
LDC00Sc Sid-ish 5T	32.0	45
LDC00Sc Sid-ish 6T	33.8	45
LDC00Sc	36.7	100
LDC00Sc Sid-ish 4T	42.7	100
LDC00Sc Sid-ish 5T	41.0	100
LDC00Sc Sid-ish 6T	39.8	100

Errors  $\pm 0.2-0.3$

**100 GeV Numbers very preliminary**



# Evidence it's NOT the detector

- See previous slide.
- But wait, big difference in Hcal. (SiD01 RPC's while LDC is analog scintillator)
- Scintillator and RPCs have been compared using perfect pattern recognition, with small differences. However, scintillator was treated digitally. I could easily do the comparison for straight analog, but would need help from NIU for the 2 bit analog option.
- Although perfect pattern recognition, difficult to believe pattern recognition any easier in scintillator.



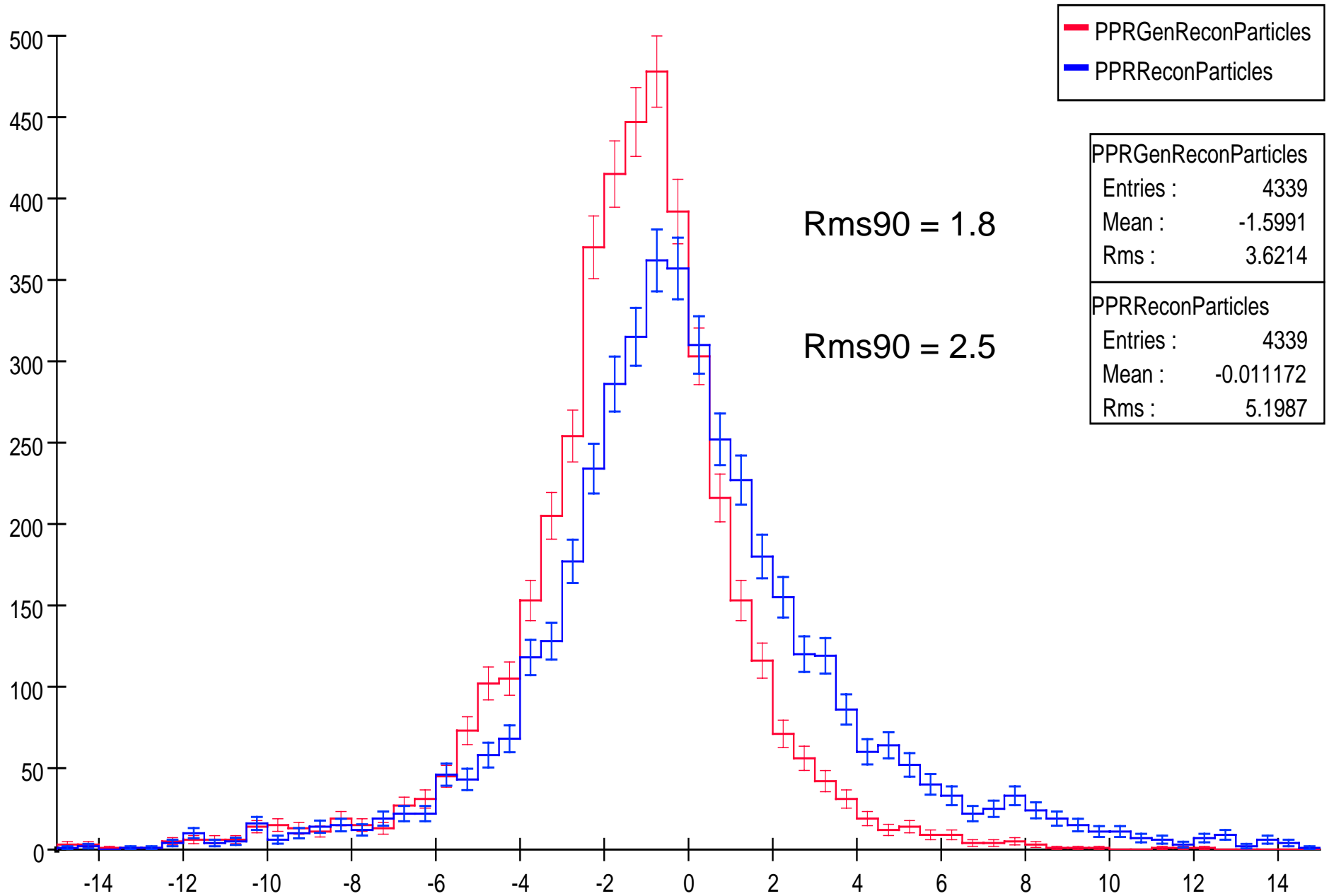
# Algorithms

- Must be the algorithms. Where do they go wrong? How do we approach fixing them?
- (admitting you have a problem is the first step toward recovery???)
- Modular approach (Template): Excellent idea in principle, somewhat more difficult in practice to change.
- Each implementation has a different approach, and I would like to see all of them in CVS to study what works and what doesn't.

# The road to improvement

- This is not new and shocking information, so what have we been doing the last few months? (Since I can't speak for everyone, I'll describe the approach I've taken)
- Define the final state particles. (What are we trying to reconstruct?) This is not as trivial as it sounds. Interactions and decays before the calorimeter make significant differences in the Generator final state particles and what gets reconstructed in the detector. To illustrate the point, the following plot shows  $\Delta Z_{\text{mass}}$  with perfect hit assignments back to the generator FS particles, compared to our current definition of final state particles.

# ReconCheaterTest.aida - Delta Zmass



# Improvement (cont)

- The PFA reconstructions take time to run, so we need to be able to persist the results for analysis. Several problems arose, and as far as we know all have been solved.
- With the template approach, we need to be able to replace an algorithm and hopefully show improvement. Photon finding seemed the obvious choice, since it should be “easy”, and indications were we weren’t doing a very good job. So I wrote a photon finder. Although it has  $> 95\%$  eff on single photons, when applied to the ZZ dataset, both eff and pur were in the mid 80’s. Further study found this to be due to overlaps, even with a NN111 clustering algorithm in the Ecal. Since the current algorithms were giving  $\sim 60\%$  eff, this was still usable as a test example.
- We also need to be able to look in detail at the calorimeter hit assignments. So I developed an analysis package for such studies.

# Improvement

- So now I had an example of a possible improvement, and I wanted to find the problems in actually implementing the change in a template algorithm.
- Not to pick on Mat, but as of a few weeks ago his algorithm was the only one in CVS producing the results he was reporting. So this is my example.
- I was able to run the reconstruction, analyze it, change the photon finder, and repeat.

# Analyzing NonTrivialPFA (Mat)

- Rather than show plots, I'll explain in words what was done and Mat can object when I get it wrong.
- Run PFA, find  $\sim 6.5\%$   $\Delta M/M$ . For cal hit assignments found:
  - Photons: eff = 63%, pur = 83%
  - Nhad: eff = 82%, pur = 27%
  - Chhad: eff = 58%, pur = 92%
- Replacing the photon finder, the photon eff and purity were both  $\sim 85\%$ . Not surprisingly, the mass width didn't change, since 40% of the tracks were being measured with the calorimeter!
- So it looked easy! The charged track association was very poor, so fixing that should be a big help

# oops

- So I replaced the track association algorithm with a cheater (and so did Mat independently) and while all the hit assignment numbers improved drastically, the mass width didn't change.
- Since I have no rational explanation for this, a detailed study is needed. I plan to start on this with Norman and Mat immediately. The hope is to do the same in parallel with all the PFA implementations. (Steve, Lei, NIU) I only need a public version of the code that produces a list of ReconstructedParticles with the CalorimeterHits attached (clusters).

# Some possibilities

- The clustering isn't good enough. This is the most popular theory, and combined with breaking apart large clusters could be the problem. With Mat's algorithm and cheating on the track associations, > 90% eff and pur were obtained for the calorimeter energy assignments for charged tracks. Perhaps this is not near good enough.
- Making neutral hadrons from the remaining hits: another good possibility. Not a trivial procedure, and is being checked with single particles. If the reconstructed neutral hadron energy has an extra 50% +/- 100%, narrowing the mass width would be futile.
- We don't know how to calculate mass: seems unlikely, but if there are errors going from track parameters to momentum to reconstructed particle 4-vectors, could cause the problem.
- ...



# Ugly details

- Implementation: Many details vary, and even the ones in common may not be ideal. The 2 I'm most familiar with are Mat's and Steve's. Some of the fundamental differences are the in the track association and the building of the clusters. Mat's dropping of tracks that are not associated with clusters may not be ideal. Both remove hits from consideration that are identified as photons, but how pure do these need to be for this to work? Neither attempt to break up large clusters, and this is the largest strategic difference in Steve's and PandoraPFA.
- These are just a few examples, but studies are needed to try to quantify the effects, both on mistakes in energy association and energy and mass resolution.

# More disclaimers...

- The SiD PFA's aren't giving "good enough" results to work with. But in no way am I trying to downplay the difficulty or the amount of effort put in to develop and IMPLEMENT the algorithms. Starting from scratch, even with a great new idea, is not a viable option. Analyzing the problem areas of existing implementations and improving them is.

# LOI

- Can we improve existing algorithms to perform at the current PandoraPFA level? Of course.
- Can we do it in time to help produce the LOI? With a rational explanation for the difference in Mat's implementation and the perfect pattern recognition, I would have said of course. But with that lack of understanding, I lean toward probably.

- There is a vast arsenal of tools available to allow detailed studies and analysis to be done quickly. The ease of working within the org.lcsim framework with SLIC simulations is a real tribute to the developers. There is a large amount of full simulation data available for an enormous number of detectors, and requested new datasets are generally available within 24 hrs.

# Summary

- When you're knee deep in alligators, it's difficult to remember the objective is to clear the swamp.