

2019 Data Reconstruction Readiness

Norman Graf (SLAC)
HPS Collaboration Meeting
November 19, 2020

Issues

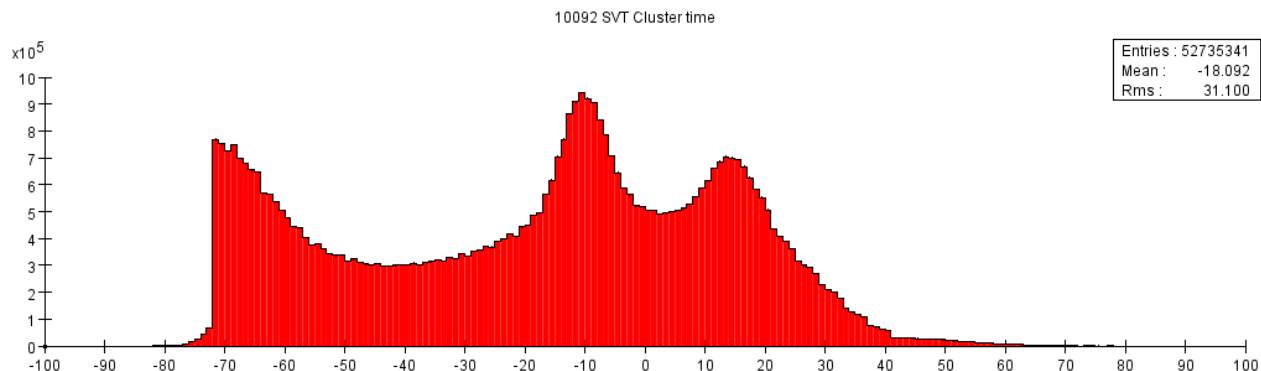
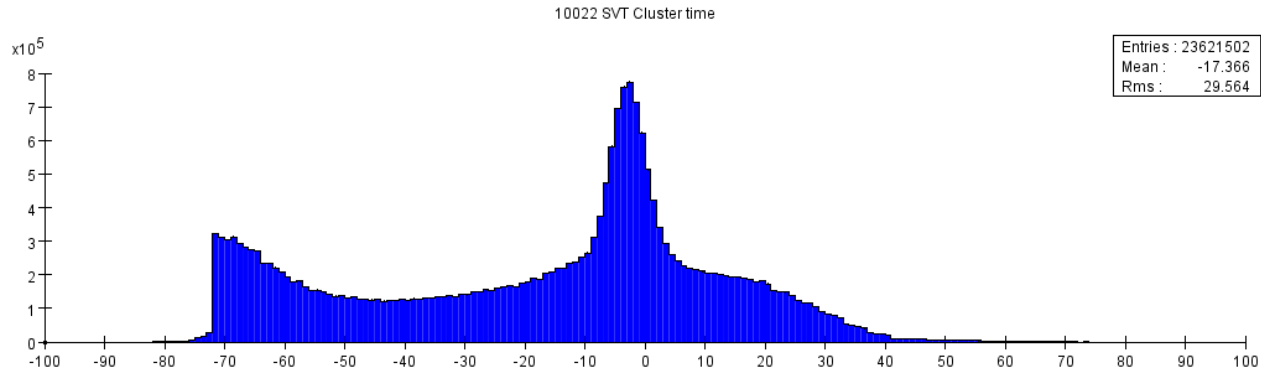
- We have promised our funding agencies that we will conduct a “timely” analysis of the 2019 data.
 - It’s been over a year since the run ended
- We have another data taking run coming up
 - Run is scheduled in less than a year
- What is needed and when is it needed to accomplish our goals?
 - Simulation (see Tongtong’s presentation)
 - Reconstruction
 - Analysis (see Cameron’s presentation)

Data Reconstruction Software Issues

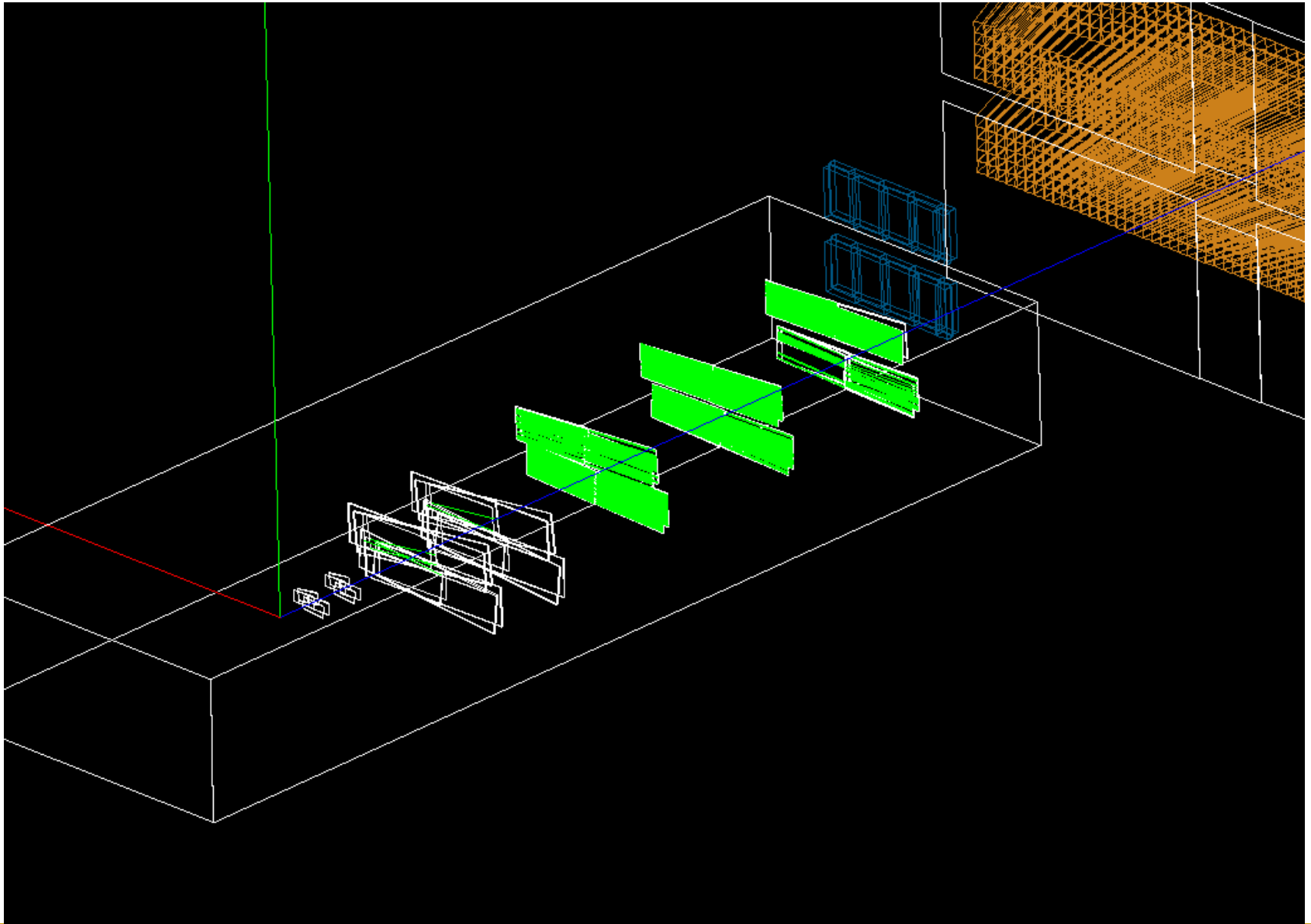
- ECal has finished up gains, sampling fractions, cluster position corrections and timing work ongoing.
- Hodoscope software OK.
- SVT trigger phase needs to be fixed for certain runs.
- Need a 2019 Event Flag Filter to remove obviously bad events
 - flag wrong SVT position, wrong SVT voltage, etc.
 - identify recoverable “monster” SVT events
 - identify, skip and drop truly “monster” SVT events
- SVT APV25 waveform fitting
 - Is the current fitting sufficient for our track timing?
 - Replacing simplex with migrad improves fitting, gives uncertainties, takes much more time.
 - Need to study this ASAP, as we plan to drop raw data from output.
 - Can it be improved and/or sped up?
- SVT actively working on alignment/calibration
 - See PF’s talk yesterday.
 - Will we need more than one alignment?
- Tracking group actively improving CPU performance
 - See PF’s talk today.
 - Need more work to optimize track-finding strategies
- Need characterization and performance evaluation of tracking software
 - Track-finding efficiencies
 - Momentum scale & resolution
- Output lcio files are bloated with extraneous data.
 - Remove extraneous Drivers
 - Need to prune our data tree and remove unnecessary collections from lcio output
- Memory footprint needs to be below 1GB to be efficient at JLab.

SVT Trigger Phase

- For certain runs, SVT locked on to the wrong trigger phase.
- Runs have been identified, just need to fix this.



SVT “Monster” Event

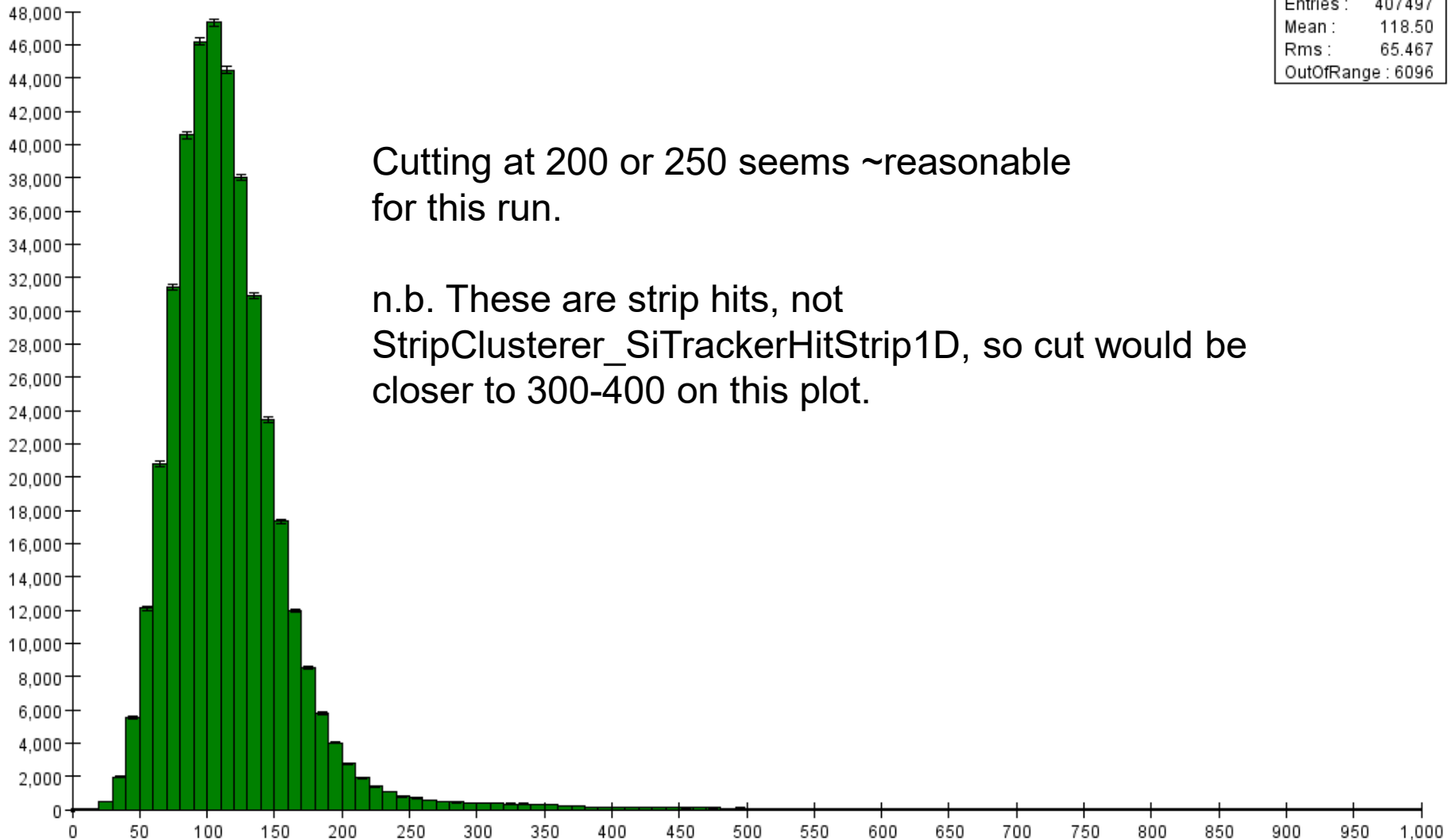


SVT “Monster” Events

- A number of *ad hoc* solutions have been implemented to deal with these events.
 - Matt cuts at less than 250 SVT clusters for SeedTracker
 - Robert cuts at fewer than 200 SVT clusters in Kalman Filter
 - PF cuts at less than 200 SVT clusters in DataTrackerHitDriver
- All of these were based on global number of hits
 - May work for some situations, but clearly fails for processing all our runs.
- Need to identify specific types of “monster” events and handle them individually.
- Need to do this early on (before fitting of waveforms)
- Truly “monster” events should be abandoned and not even processed further nor written out.

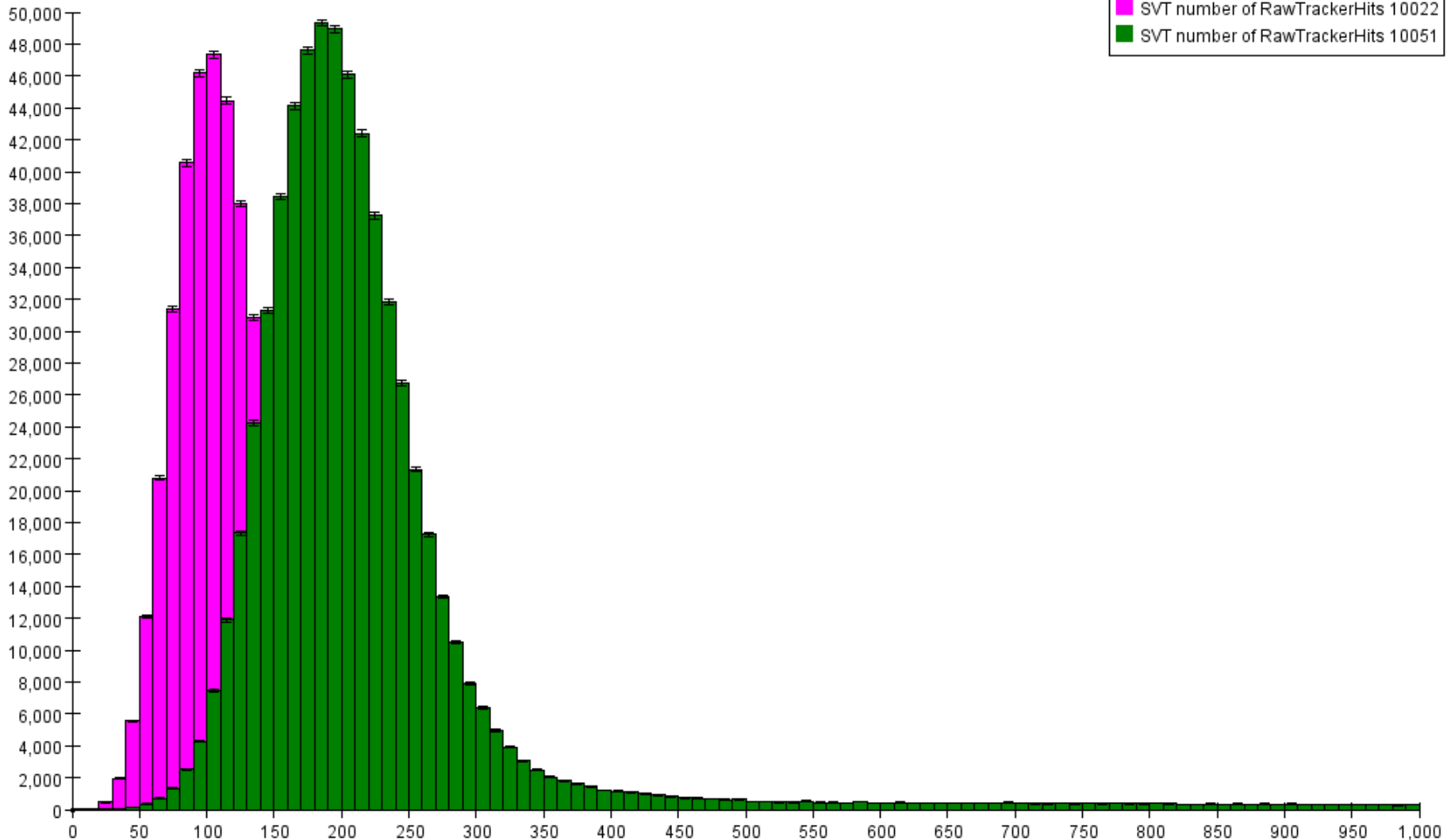
SVT Number of Raw Tracker Hits

SVT number of RawTrackerHits 10022



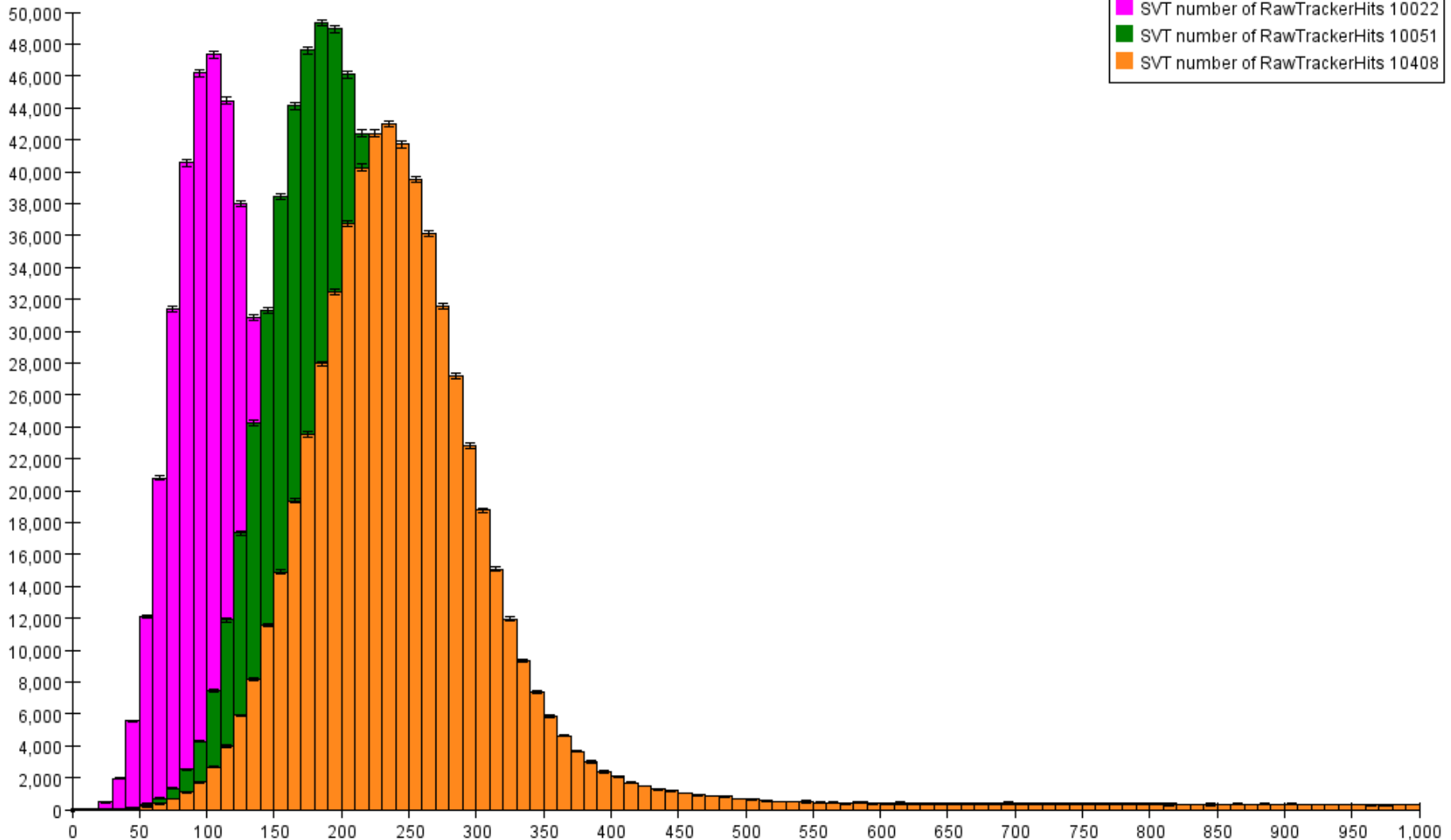
SVT Number of Raw Tracker Hits

eventPlots20201117.aida



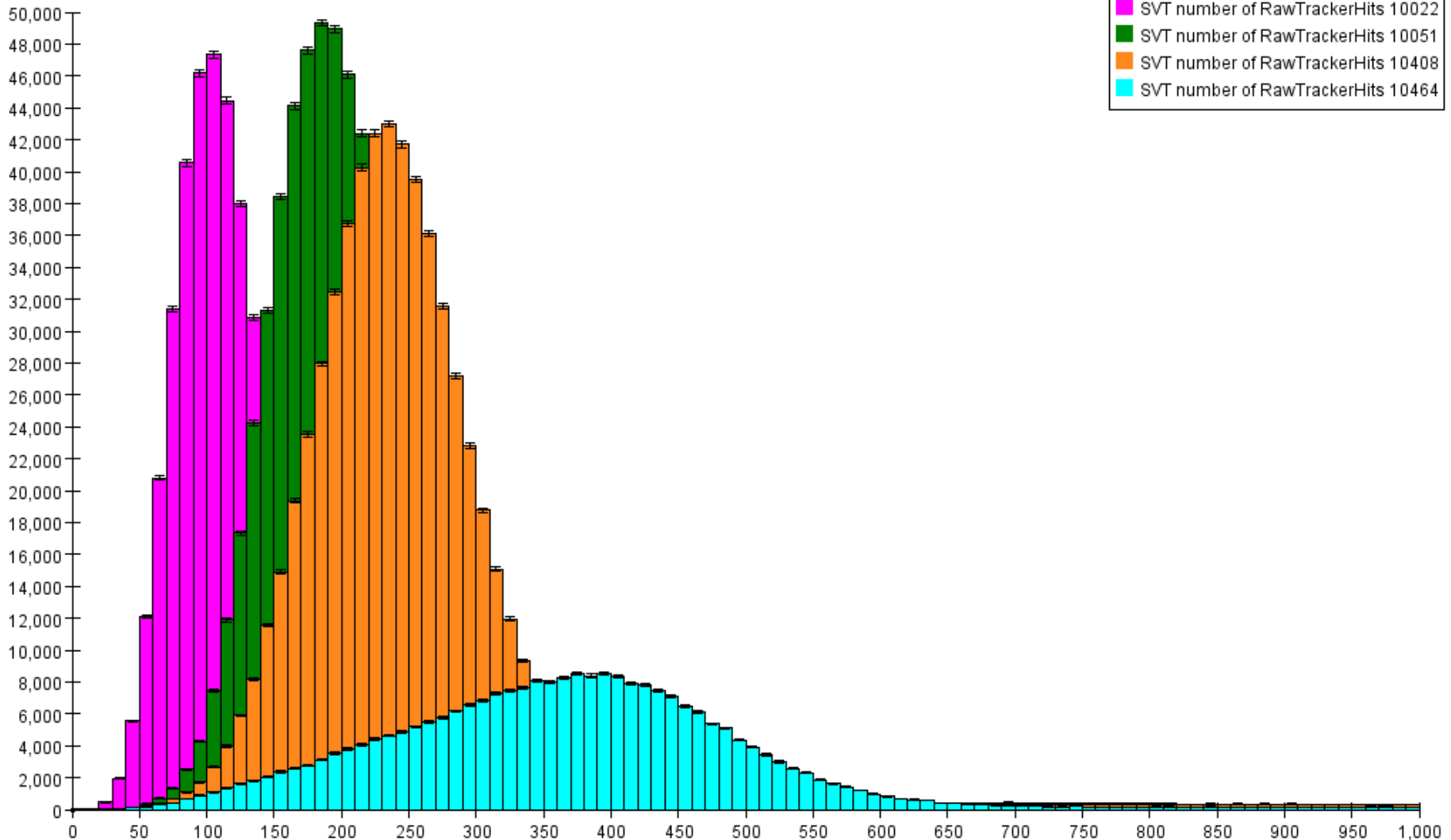
SVT Number of Raw Tracker Hits

eventPlots20201117.aida



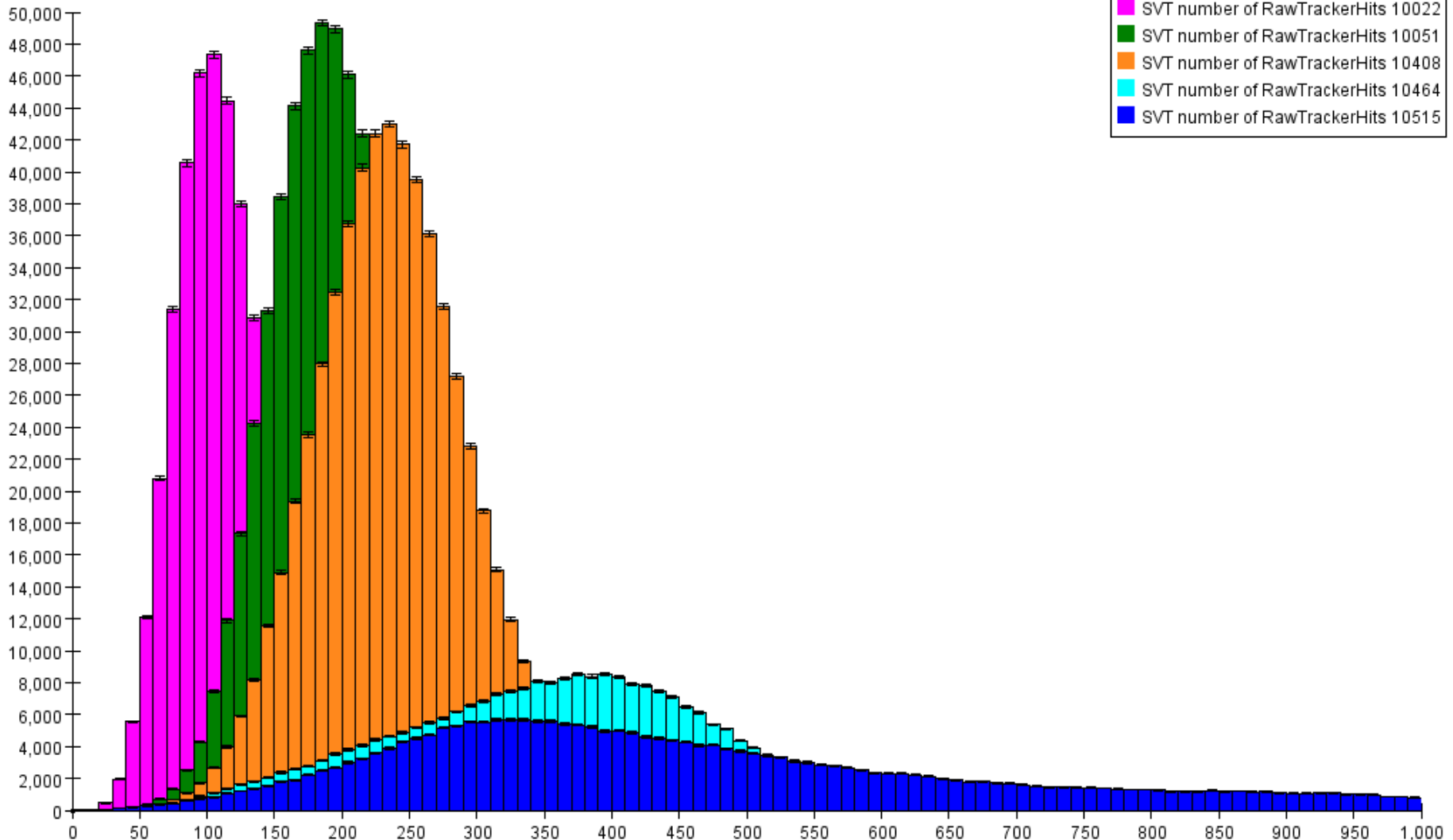
SVT Number of Raw Tracker Hits

eventPlots20201117.aida



SVT Number of Raw Tracker Hits

eventPlots20201117.aida

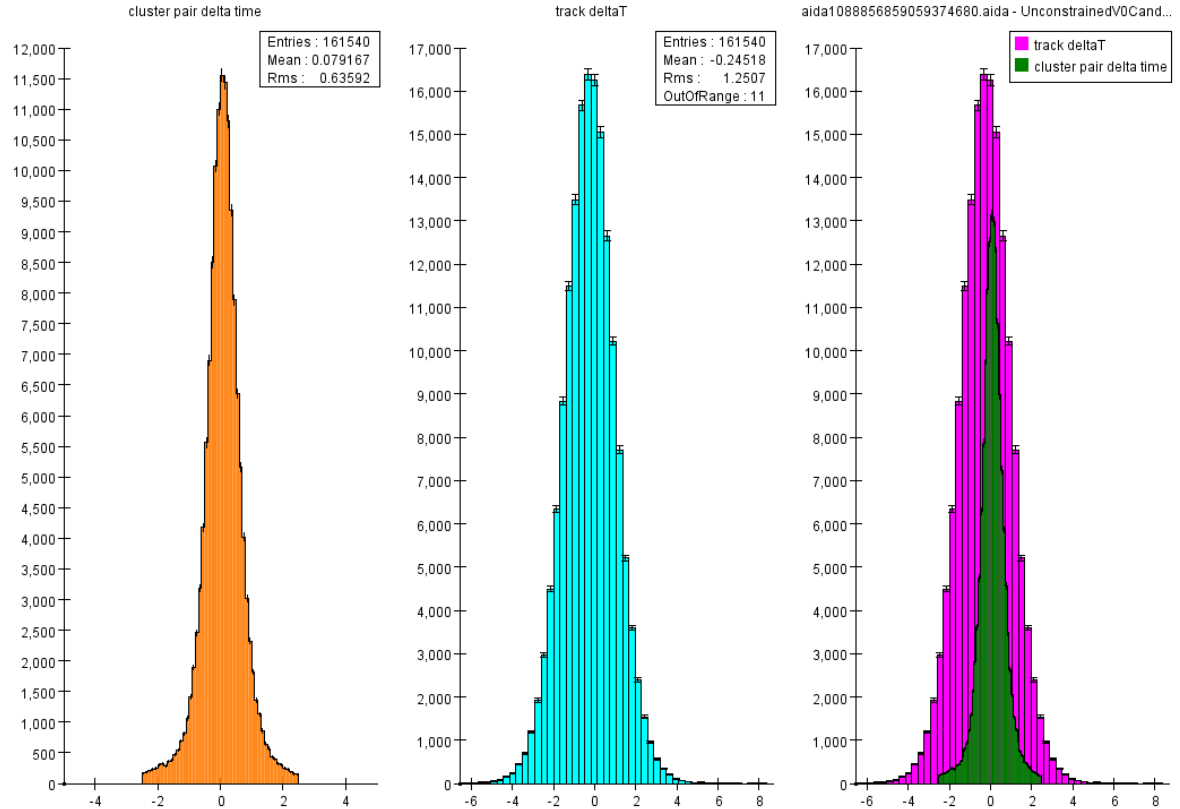


SVT “Monster” Events

- Plan is to identify and then **skip both processing and writing out the event.**
 - This is new behavior as in the past we simply flagged such events.
- A skim of events containing more than 250 SVTRawTrackerHits is available to characterize the issues, develop the algorithms and test the efficiency of the cuts.
- git issue [iss731](#) addresses this.

Track Timing vs Ecal Timing

- Select V0 candidates with Ecal Clusters associated with each track.
- Track timing resolution a factor of two worse.
- Good enough?



Track Finding

- Currently running both SeedTracker and Kalman Filter for track finding.
- Need to characterize the performance and optimize track-finding efficiency and CPU time.
- May need condition-dependent strategies
 - e.g. no need to include a “dead” layer in either the seed or confirm stage of pattern recognition.
- Characterizing track-finding efficiencies
 - ECal FEEs → clean sample of high-energy electrons
 - ECal WABs → clean sample of lower-energy electrons
 - Positron trigger (coincidence of hodoscope hit with calorimeter cluster) provides clean sample of positrons

Logistics

- We need good estimates of our CPU needs to process the full 2019 “good” data sample
 - ~50 Billion events
 - Goal is better than 10Hz with a memory footprint of less than 1 GB
- We need good estimates of the amount of computing power we can rely on.
 - will be competing with other experiments for processing resources
- We need good estimates of our storage needs
 - ~600TB of evio data
 - How large is the recon output?
- Will inform the overall HPS data processing plan

Sample Partitions

- Essential to test on a faithful subset of the full 2019 data run.
- Partitions sampled ('.*04[12]\$') from the list of 282 “good” runs:
 - Roughly 150M events (142937602)
 - Roughly 3‰ of the 2019 “good” data (857 / 276339)
- Using:
 - EvioToLcio
 - pass1-dev_fix + iss732-refactor
 - PhysicsRun2019_pass0_recon.lcsim
 - HPS_TY_iter4

Output File Size Status

- Status of the reconstruction is still in flux.
- Little (no?) effort has been devoted to limiting content or file size.
- Effort concentrated on understanding efficiency, resolution, etc. i.e. "physics" performance.
- Production Reconstruction will differ substantially.
 - Latest "pass0" steering file includes both SeedTracker/GBL & Kalman Filter to enable comparison of tracking.
- Nevertheless...

Recon Output

LCSim Event
Run: 10486 Event: 5973679

Event

- BeamspotConstrainedV0Candidates
- BeamspotConstrainedV0Candidates_KF
- BeamspotConstrainedV0Vertices
- BeamspotConstrainedV0Vertices_KF
- EcalCalHits
- EcalClusters
- EcalClustersCorr
- EcalReadoutHits
- EcalUncalHits
- FADCGenericHits
- FinalStateParticles
- FinalStateParticles_KF
- GBLKinkData
- GBLKinkDataRelations
- GBLTracks
- HelicalTrackHitRelations
- HelicalTrackHits
- HodoCalHits
- HodoGenericClusters
- HodoReadoutHits
- KFGBLStripClusterData
- KFGBLStripClusterDataRelations
- KFTrackData
- KFTrackDataRelations
- KalmanFullTracks
- MatchedToGBLTrackRelations
- MatchedTracks
- OtherElectrons
- RFHits
- RotatedHelicalTrackHitRelations
- RotatedHelicalTrackHits
- SVTFittedRawTrackerHits
- SVTRawTrackerHits
- SVTShapeFitParameters
- StripClusterer_SiTrackerHitStrip1D
- TSBank
- TargetConstrainedV0Candidates
- TargetConstrainedV0Candidates_KF
- TargetConstrainedV0Vertices
- TargetConstrainedV0Vertices_KF
- TrackData
- TrackDataRelations
- TriggerBank
- UnconstrainedV0Candidates
- UnconstrainedV0Candidates_KF
- UnconstrainedV0Vertices
- UnconstrainedV0Vertices_KF
- UnconstrainedVcCandidates
- UnconstrainedVcCandidates_KF
- UnconstrainedVcVertices
- UnconstrainedVcVertices_KF
- VTPBank

LCIO Event Header

Run	10486
Event	5973679
Time Stamp	1003753987372
Detector Name	HPS_TY_iter4
Event Weight	1.0
IDRUP	0
SLIC Version	
Geant4 Version	

Collections

Name	Type	Size
SVTFittedRawTrackerHits	org.lcsim.event.LCRelation	264
SVTShapeFitParameters	org.lcsim.event.GenericObject	264
SVTRawTrackerHits	org.lcsim.event.RawTrackerHit	219
StripClusterer_SiTrackerHitStrip1D	org.lcsim.event.TrackerHit	115
HelicalTrackHitRelations	org.lcsim.event.LCRelation	92
RotatedHelicalTrackHits	org.lcsim.event.TrackerHit	46
HelicalTrackHits	org.lcsim.event.TrackerHit	46
RotatedHelicalTrackHitRelations	org.lcsim.event.LCRelation	46
FADCGenericHits	org.lcsim.event.GenericObject	37
HodoReadoutHits	org.lcsim.event.RawTrackerHit	32
HodoCalHits	org.lcsim.event.CalorimeterHit	30
EcalCalHits	org.lcsim.event.CalorimeterHit	17
EcalUncalHits	org.lcsim.event.CalorimeterHit	17
KFGBLStripClusterData	org.lcsim.event.GenericObject	16
KFGBLStripClusterDataRelations	org.lcsim.event.LCRelation	16
EcalReadoutHits	org.lcsim.event.RawTrackerHit	15
HodoGenericClusters	org.lcsim.event.GenericObject	6
KFTrackData	org.lcsim.event.GenericObject	2
KFTrackDataRelations	org.lcsim.event.LCRelation	2
TriggerBank	org.lcsim.event.GenericObject	2
VTPBank	org.lcsim.event.GenericObject	2
FinalStateParticles_KF	org.lcsim.event.ReconstructedParticle	2
KalmanFullTracks	org.lcsim.event.Track	2
RFHits	org.lcsim.event.GenericObject	1
EcalClusters	org.lcsim.event.Cluster	1
GBLKinkDataRelations	org.lcsim.event.LCRelation	1
TrackDataRelations	org.lcsim.event.LCRelation	1
EcalClustersCorr	org.lcsim.event.Cluster	1
FinalStateParticles	org.lcsim.event.ReconstructedParticle	1
GBLTracks	org.lcsim.event.Track	1
MatchedToGBLTrackRelations	org.lcsim.event.LCRelation	1
GBLKinkData	org.lcsim.event.GenericObject	1
TrackData	org.lcsim.event.GenericObject	1
TSBank	org.lcsim.event.GenericObject	1
MatchedTracks	org.lcsim.event.Track	1
UnconstrainedV0Candidates_KF	org.lcsim.event.ReconstructedParticle	0
UnconstrainedVcVertices	org.lcsim.event.Vertex	0
BeamspotConstrainedV0Vertices	org.lcsim.event.Vertex	0
UnconstrainedVcCandidates	org.lcsim.event.ReconstructedParticle	0
TargetConstrainedV0Candidates_KF	org.lcsim.event.ReconstructedParticle	0
TargetConstrainedV0Vertices	org.lcsim.event.Vertex	0
UnconstrainedV0Vertices_KF	org.lcsim.event.Vertex	0
BeamspotConstrainedV0Candidates_KF	org.lcsim.event.ReconstructedParticle	0
UnconstrainedV0Candidates	org.lcsim.event.ReconstructedParticle	0

Recon Output

- So, it's clear that there is a LOT of extra data included in this file.
 - For instance, we won't have both SeedTracker/GBL and Kalman Filter tracks and ReconstructedParticles.
- Won't try to analyze every collection here, but it's clear that we need to survey what's going into the output and justify what's there.

What is the role of the recon file?

- Historically we have kept all of the data, including the raw data, to enable re-reconstruction from the Icio files.
- At this point, we should be able to drop the raw waveforms and only save the fitted t0 and pulse area.
- Obviously we need to include all the information for subsequent “physics” analyses.
- But, do we need to save individual SVT readout channels, or can we live with just StripClusterer_SiTrackerHitStrip1D?

Output Data Size Reduction

- A number of strategies can gain us a substantial reduction in the size of our recon output files.
- Dropping the “raw” waveforms is easiest.
 - Are we satisfied with our current pulse fitting?
- Not running the SeedTracker is straightforward
 - Need to validate Kalman Filter.
 - More, better tracks faster is a requirement.
- Pruning un-needed collections is next.
- Can consider DST set of collections which is optimized for “physics” analysis.
 - Just ReconstructedParticles?

Software Issues

- Need a major effort to merge all of the extant git branches back into master
 - both hps-java and lcsim
 - Must retain backwards-compatibility for 2016/2019
- Need to tag and make a library release before production processing.

Summary

- This list of topics is long, but not exhaustive.
- There is still a LOT of work to be done before we can start the full production processing of the 2019 data.
- Current focus is still on characterizing the detector and “physics” performance
 - alignment, calibration, track-finding efficiency, energy & momentum scale and resolution, etc.
- Much work has gone into speeding up the code along the way, more will most likely need to be done.
- Memory requirements may be OK
 - may be due to dropping high-multiplicity events.
- ~No effort expended on reducing output file size
- Need more involvement and feedback from other members of the collaboration!

This means YOU!