

# Comprehensive PID Processor

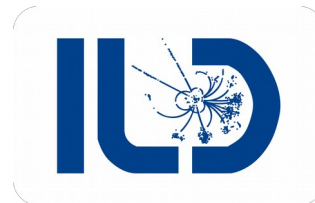
Beta version now online

Uli Einhaus

ECFA WG2 Reconstruction Workshop

12.07.2023

**HELMHOLTZ**  
RESEARCH FOR GRAND CHALLENGES

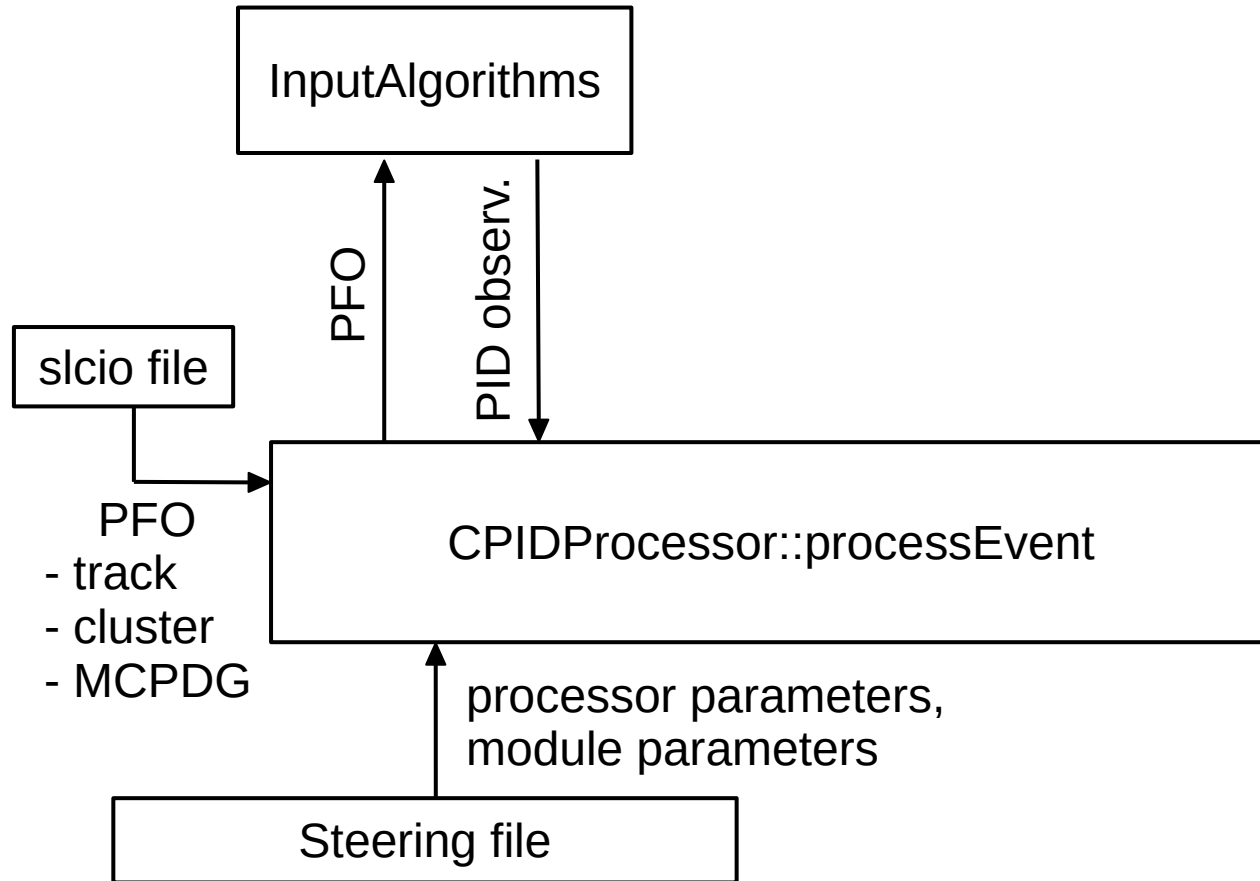


- Comprehensive Particle Identification (CPID) Processor
- Target: provide platform for future collider detectors to evaluate PID
- Approach: central book-keeping, modules for PID observables as well as training & inference
- Use Particle Flow Objects (PFOs),
- Currently Marlin processor using LCIO, usable in Gaudi via MarlinWrapper, goal is to have native implementation
- CPID (beta version) is [part](#) of the current iLCSoft release
  
- Today: structure, how to use, module overview, PID performance



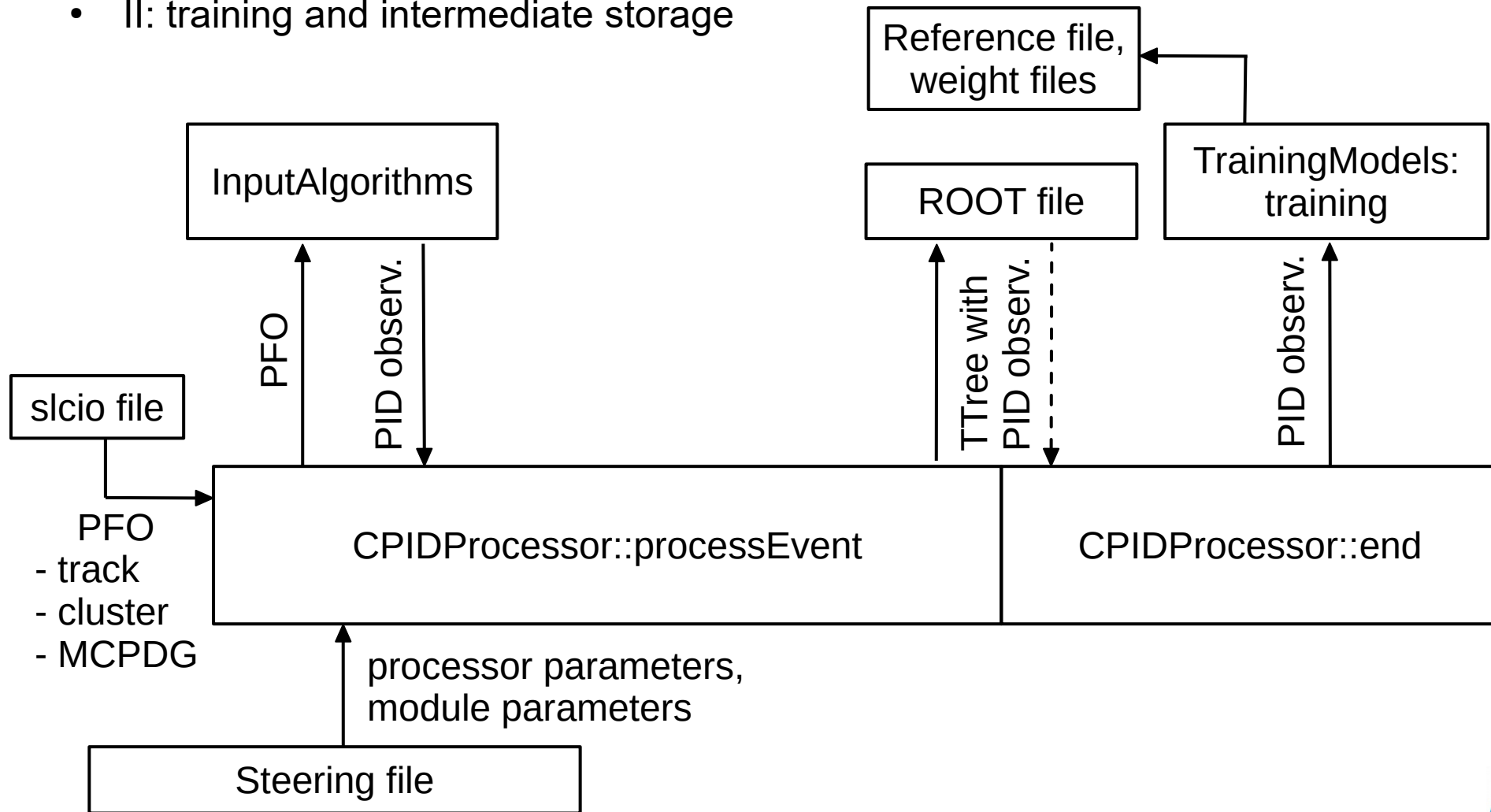
# Structure of the CPID workflow

- I: set up and observable extraction



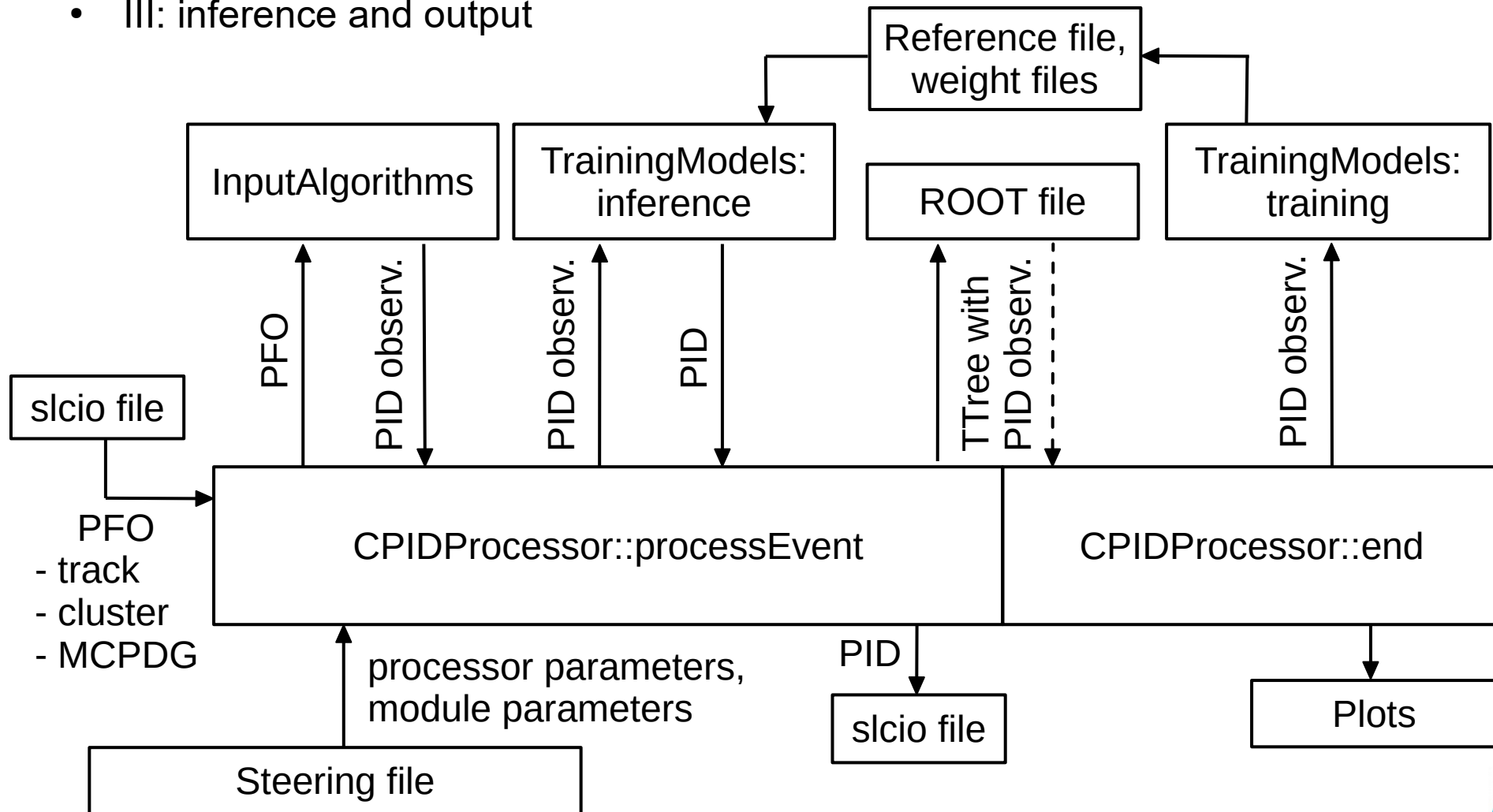
# Structure of the CPID workflow

- II: training and intermediate storage



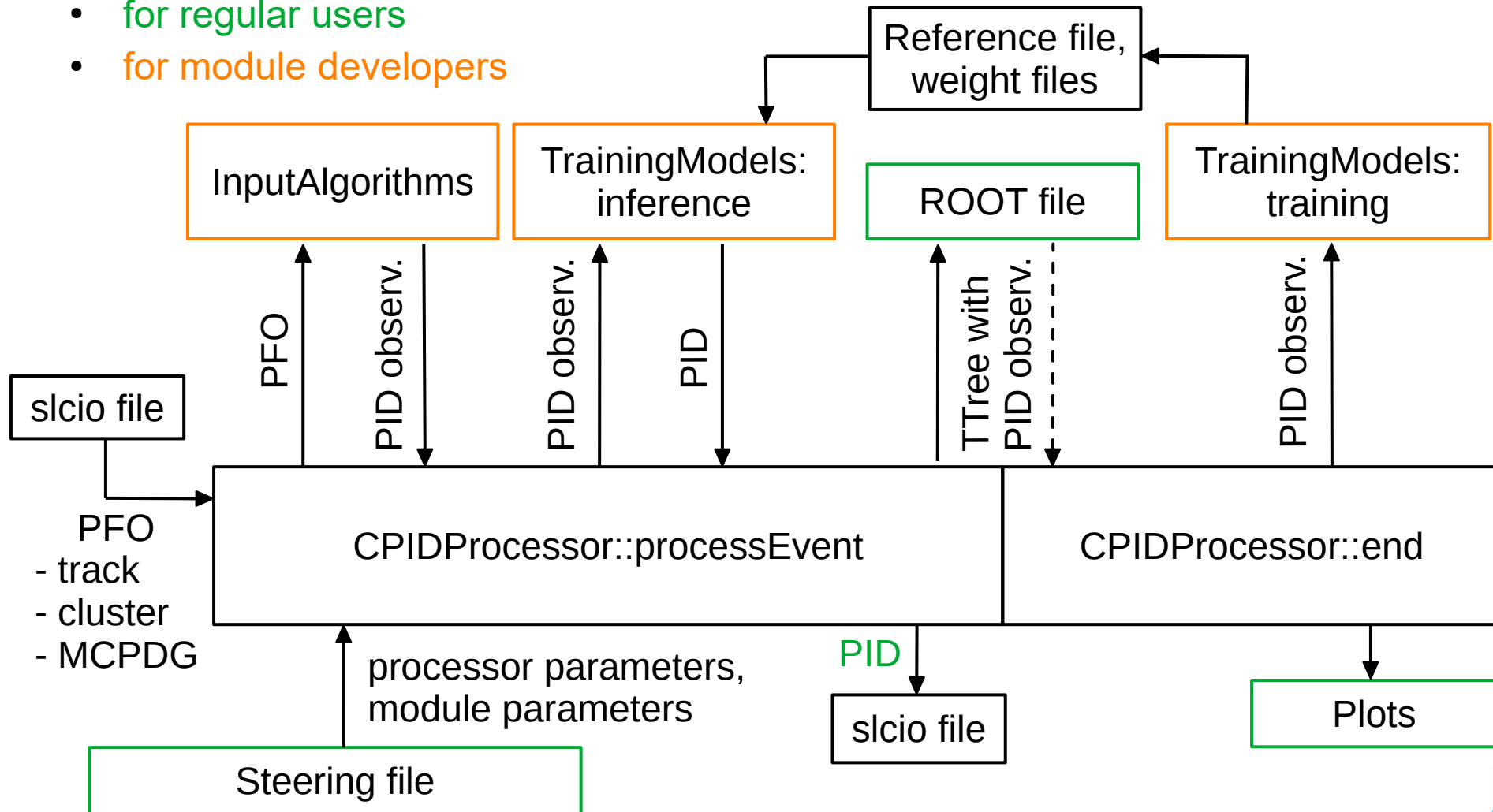
# Structure of the CPID workflow

- III: inference and output



# Structure of the CPID workflow

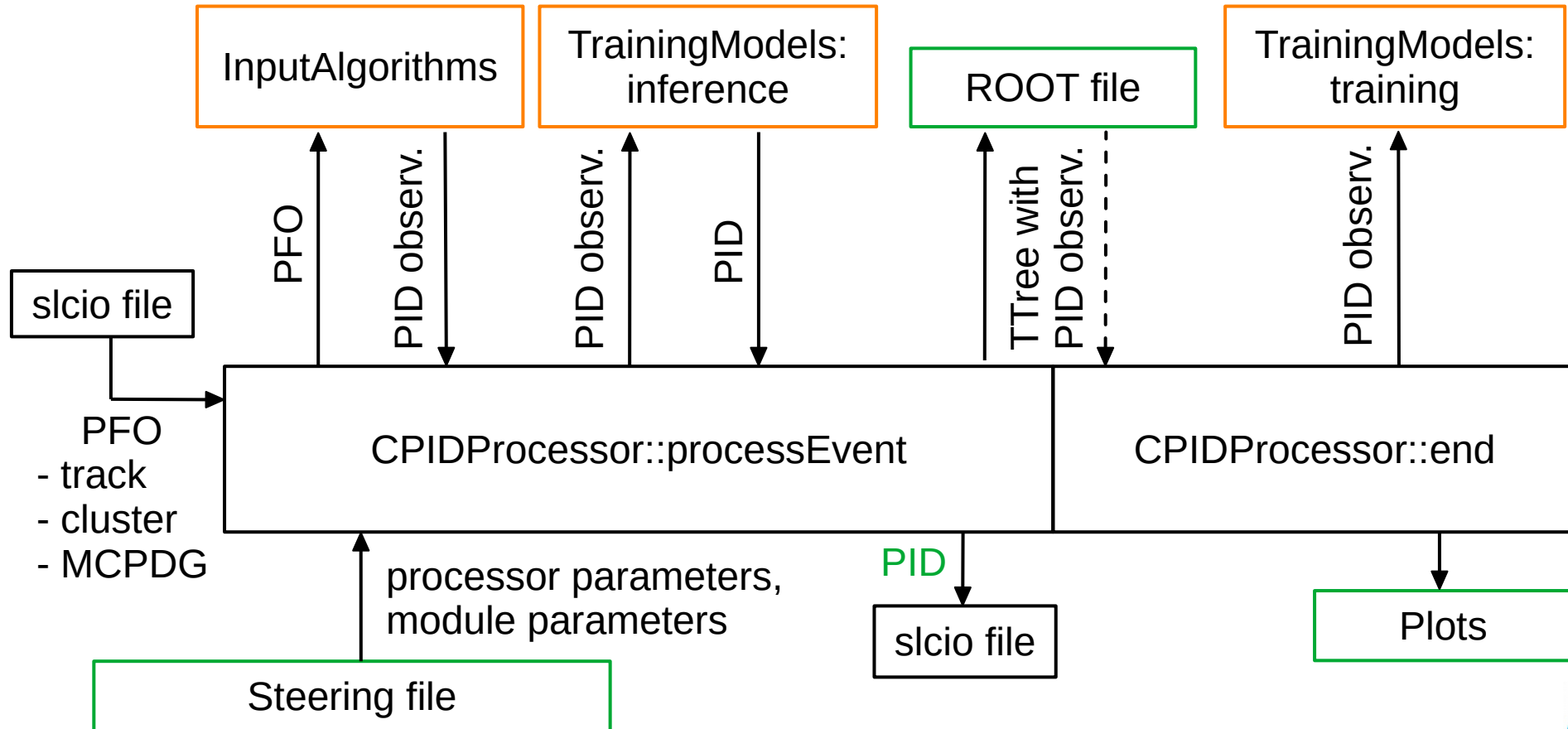
- for regular users
- for module developers



# Structure of the CPID workflow

- for regular users
- for module developers

Dynamic loading of modules means module developers don't need to touch the actual processor (analogous to Marlin processors and actual Marlin)



# How To Use

- Example steering file

```
<processor name="MyComprehensivePIDProcessor" type="ComprehensivePIDProcessor">
```

```
<parameter name="PFOCollection" type="string" value="PandoraPFOs"/>  
<parameter name="RecoMCTruthLink" type="string" value="RecoMCTruthLink"/>
```

collections

```
<parameter name="modeExtract" type="bool" value="true" />  
<parameter name="modeTrain" type="bool" value="true"/>  
<parameter name="modeInfer" type="bool" value="false"/>
```

mode selection

```
<parameter name="TTreeFileName" type="string" value="TTreeFile.root"/>  
<parameter name="reffile" type="string" value="Ref.12bins.txt"/>  
<parameter name="signalPDGs" type="FloatVec" value="11 13 211 321 2212"/>  
<parameter name="backgroundPDGs" type="FloatVec" value=""/>  
<parameter name="plotFolder" type="string" value="."/>  
<parameter name="fileFormat" type="string" value=".png"/>
```

miscellaneous

```
<parameter name="momMin" type="float" value="1"/>  
<parameter name="momMax" type="float" value="100"/>  
<parameter name="momLog" type="bool" value="true"/>  
<parameter name="momNBins" type="float" value="12"/>
```

momentum bins → separate  
model run per bin

```
<parameter name="cutD0" type="float" value="0"/>  
<parameter name="cutZ0" type="float" value="0"/>  
<parameter name="cutLamMin" type="float" value="0"/>  
<parameter name="cutLamMax" type="float" value="0"/>  
<parameter name="cutNTracksMin" type="int" value="1"/>  
<parameter name="cutNTracksMax" type="int" value="-1"/>
```

PFO cuts, not used  
for training





# How To Use

- Example steering file

```
<parameter name="inputAlgoSpecs" type="StringVec">
  dEdx_RCD:dEdx_RCD
  TOF:TOF50
  Pandora: Pandora
</parameter>
```

```
<parameter name="TOF0.S" type="StringVec" value="TOFEstimators0ps" />
<parameter name="TOF10.S" type="StringVec" value="TOFEstimators10ps"/>
<parameter name="TOF50.S" type="StringVec" value="TOFEstimators50ps"/>
```

```
<parameter name="dEdx_RCD.F" type="FloatVec">
-1.28883368e-02  2.72959919e+01  1.10560871e+01 -1.74534200e+00  -9.84887586e-07
 6.49143971e-02  1.55775592e+03  9.31848047e+08  2.32201725e-01  2.50492066e-04
 6.54955215e-02  8.26239081e+04  1.92933904e+07  2.52743206e-01  2.26657525e-04
 7.52235689e-02  1.59710415e+04  1.79625604e+06  3.15315795e-01  2.30414997e-04
 7.92251260e-02  6.38129720e+04  3.82995071e+04  2.80793601e-01  7.14371743e-04
 1
</parameter>
```

```
<parameter name="trainModelSpecs" type="StringVec">
  TMVA_BDT_MC:TMVA_BDT_MC_12bins
</parameter>
<parameter name="trainingObservables" type="StringVec"> </parameter>
```

```
<parameter name="TMVA_BDT_MC_12bins.S" type="StringVec">
!V:!Silent:Color:DrawProgressBar:Transformations=I;D;P;G,D:AnalysisType=multiclass
SplitMode=Random:NormMode=NumEvents:!V
!H:!V:NTrees=100:BoostType=Grad:Shrinkage=0.10:UseBaggedBoost:BaggedSampleFraction=0.50
dEdx_RCD_piDis>-900&&dEdx_RCD_kaDis>-900
</parameter>
```

specify InputAlgorithms  
[type]:[name]

give module parameters  
[name].S, [name].F

-  
these depend on the  
individual modules

specify TrainingModel(s)  
[type]:[name]

give module parameters  
[name].S, [name].F

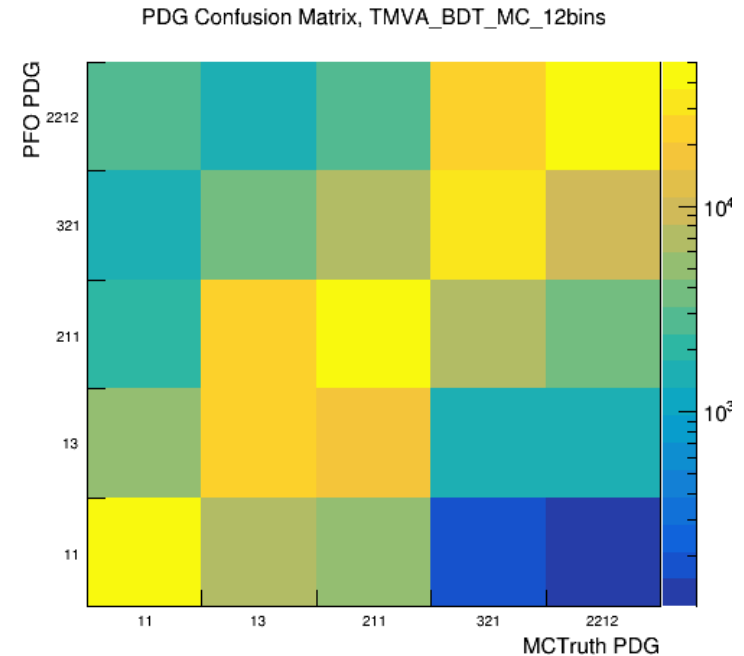


- Direct connection with observables TTree, run 'standard' ROOT TMVA BDT
- Take 4 string inputs, corresponding to TMVA loader and factory options
- TMVA\_BDT uses simple sig/bkg BDT, trains `_signalPDGs` vs. `_backgroundPDGs`
- TMVA\_BDT\_MC uses multiclass BDT, trains all `_signalPDGs` against each other  
→ used for performance plots on following slides



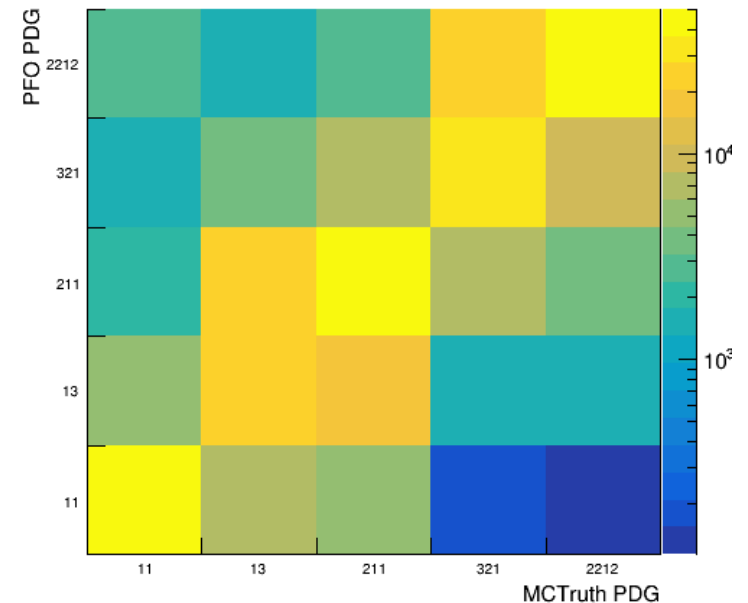
# Default Performance Assessment

- Standard plot: 5 charged particles (e,  $\mu$ ,  $\pi$ , K, p) confusion matrix
- Single particles, flat in  $\log(p)$  and  $\cos(\theta)$ , integrated over  $1 \text{ GeV} < p < 100 \text{ GeV}$
- Trained and evaluated with multiclass BDT with 12 bins in  $\log(p)$
- Variables: observable(s) of the corresponding InputAlgorithm +  $p + \lambda$



- RCD = reference curve distance, i.e. distance to Bethe-Bloch curves, removes momentum dependence compared to using dE/dx value directly
- Takes (fully reconstructed) dE/dx value from Compute\_dEdxProcessor, and reference curves from dEdxAnalyser as float parameter input
- Optional: adjustment of dE/dx value by scaling of distance to true curve, emulates better/worse dE/dx resolution
- Gives 5 observables: eDis, muDis, piDis, kaDis, prDis

PDG Confusion Matrix, TMVA\_BDT\_MC\_12bins

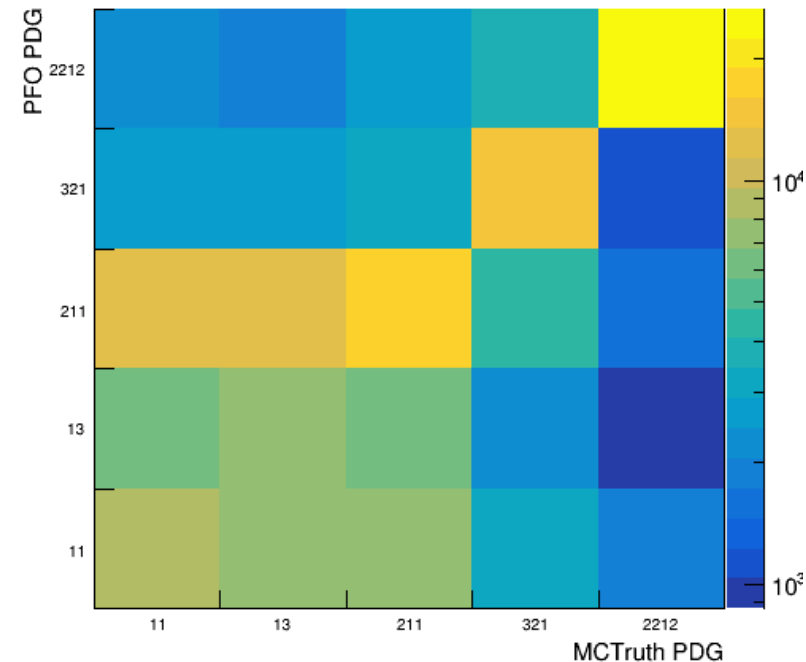


# Input Algorithms: TOF and TOF223

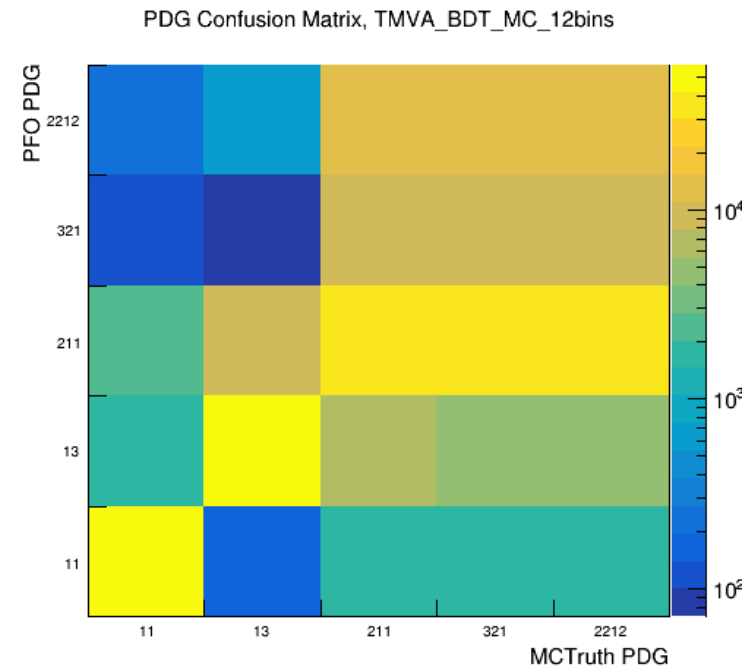
- TOF:
  - time of flight, using initial implementation
  - various issues leading to limited purity
  - available in last large MC production
  - returns track beta ( $v/c$ )
- TOF223:
  - time of flight (in iLCSoft v02-02-03), using B. Dudar's new TOF [implementation](#) including track length estimation
  - only available in newest simulation
  - runs on REC files, i.e. can't be done retroactively
  - returns reconstructed PFO mass
- Both use output of TOFEstimator, only need corresponding estimator name to find in PFO PIDHandler

Note: here  $1 \text{ GeV} < p < 10 \text{ GeV}$

PDG Confusion Matrix, TMVA\_BDT\_MC\_12bins

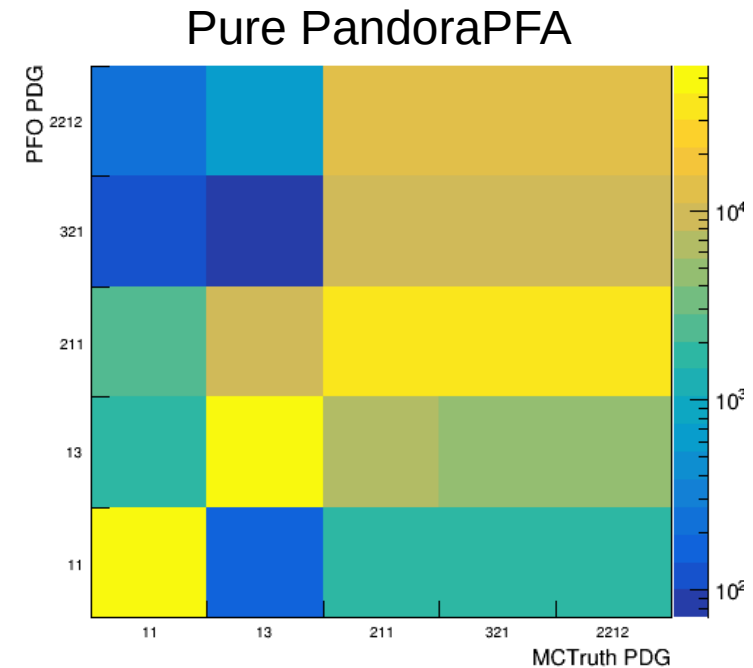
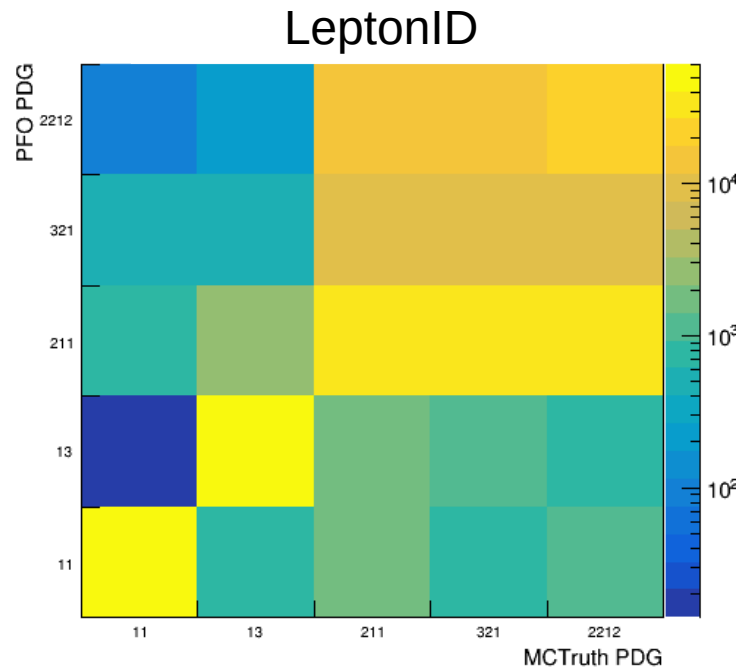


- Takes the PID assigned by PandoraPFA (particle flow, PID based on cluster shapes) to the PFO
- Either  $e$ ,  $\mu$ ,  $\pi$ ,  $\gamma$  or  $n$ , returns 1 observable: PDG (11, 13, 211, 22, 2112)

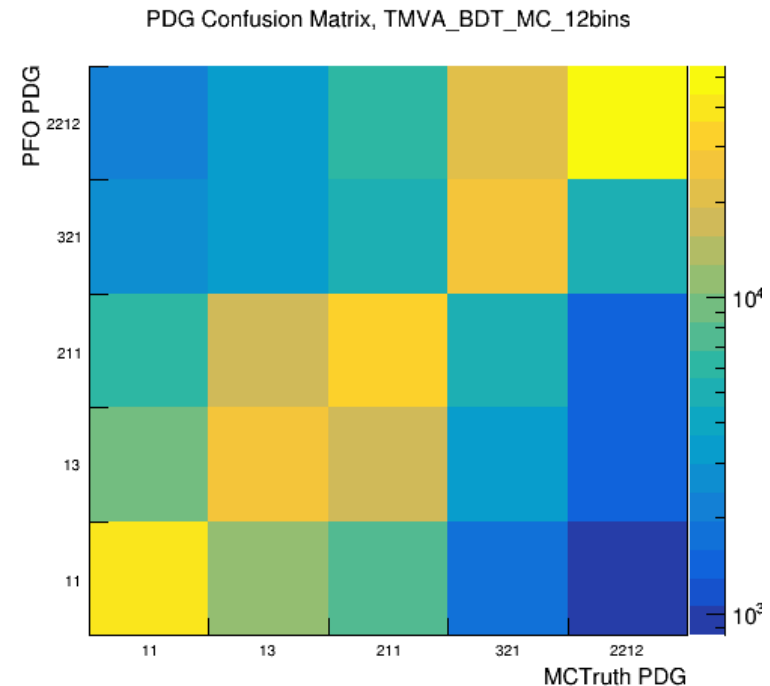


# InputAlgorithm: LeptonID

- L. Reichenbach's new [LeptonID](#) (e vs.  $\mu$  vs. hadrons, target: semileptonic b/c decays)
- Improved  $\mu/\pi$  separation compared to pure PandoraPFA
- Use LeptonID's BDT scores (based on  $\sim 20$  inputs) as input for CPID

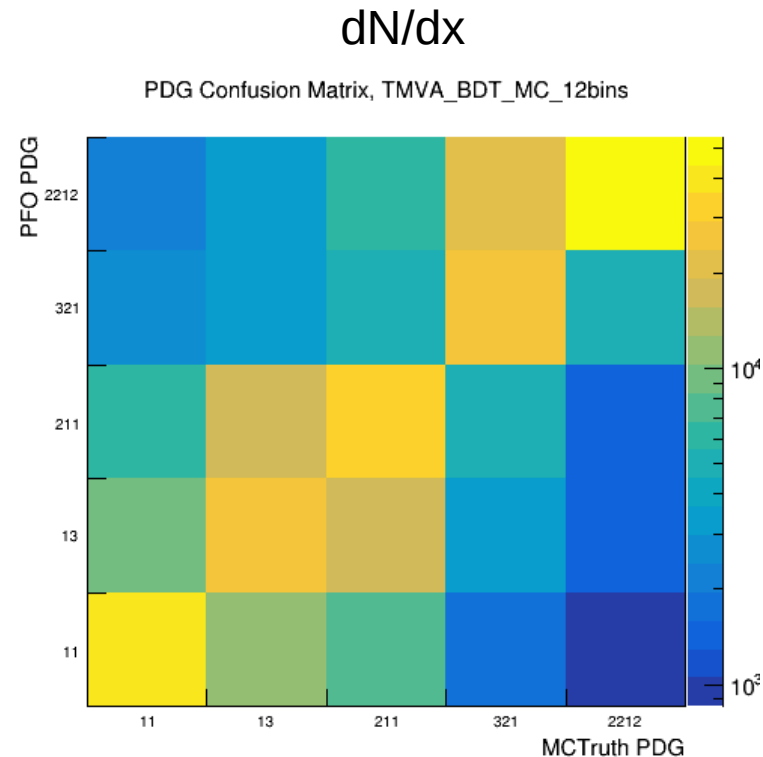
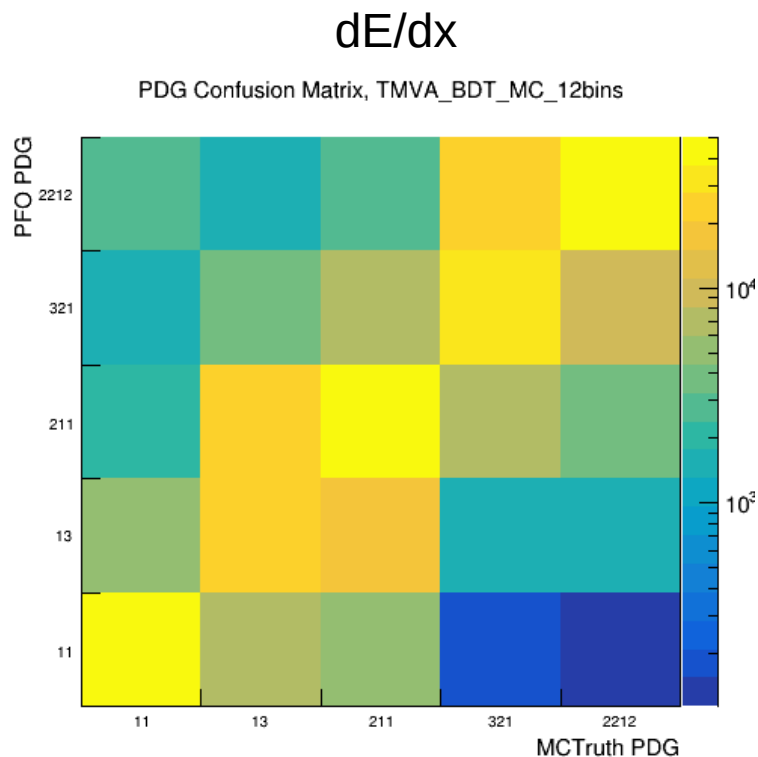


- Request to include [cluster counting](#) by IDEA DC community
- Using parametrisation of cluster counting  $dN/dx(\beta\gamma)$  from Delphes
- Take first track in PFO, calculate track length in TPC or parametrised cylinder, get Poisson(average #clusters), get distance to reference curves
- Also here added parameter to scale distance to true curve
- Gives 5 observables: eDis, muDis, piDis, kaDis, prDis



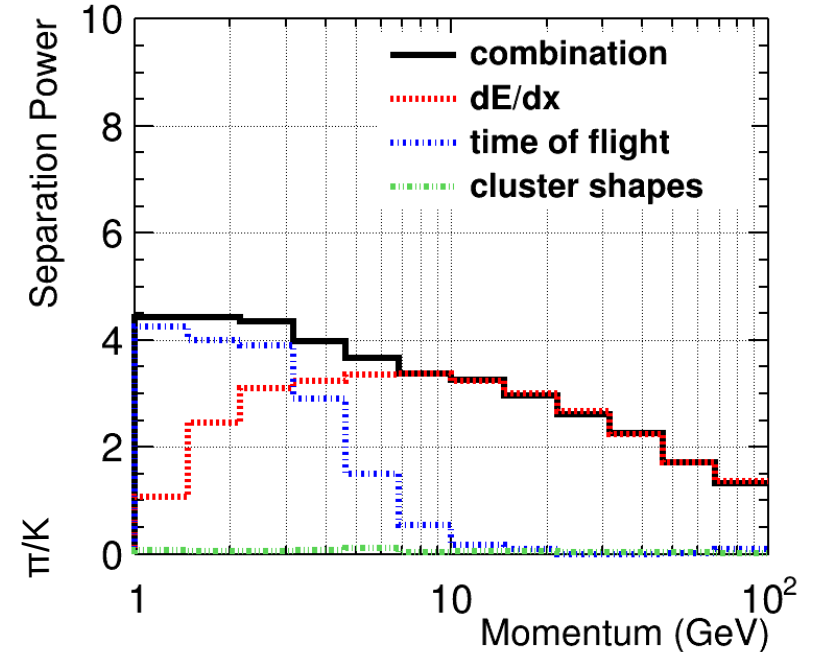


- Direct comparison: dN/dx with simulated effective resolution of  $\sim 2\%$ , conventional dE/dx with 4.5 %



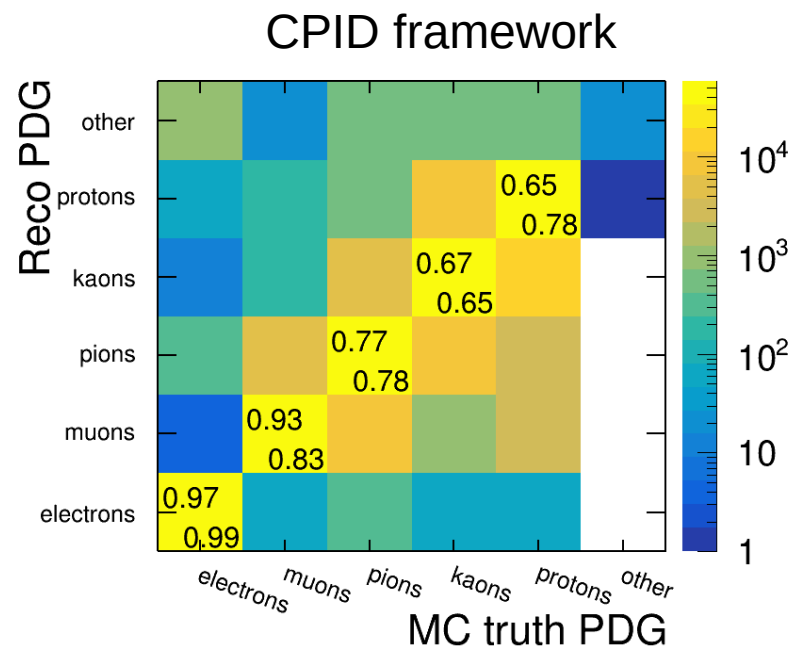
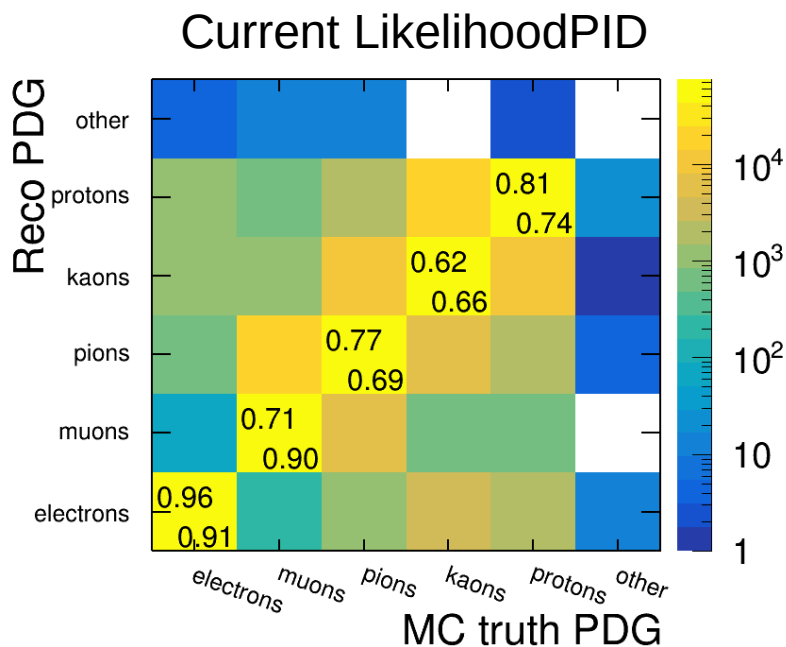
# $\pi/K$ Separation with Combined Observables

- $dE/dx$ , TOF10 and Pandora
- TMVA\_BDT with sig = K, bkg =  $\pi$   
train & eval per 12 mom bins and per used observable(s)
- TOF works at low momenta,  $dE/dx$  at moderate ones, Pandora not at all (as expected)
- Combination corresponds to independent observables



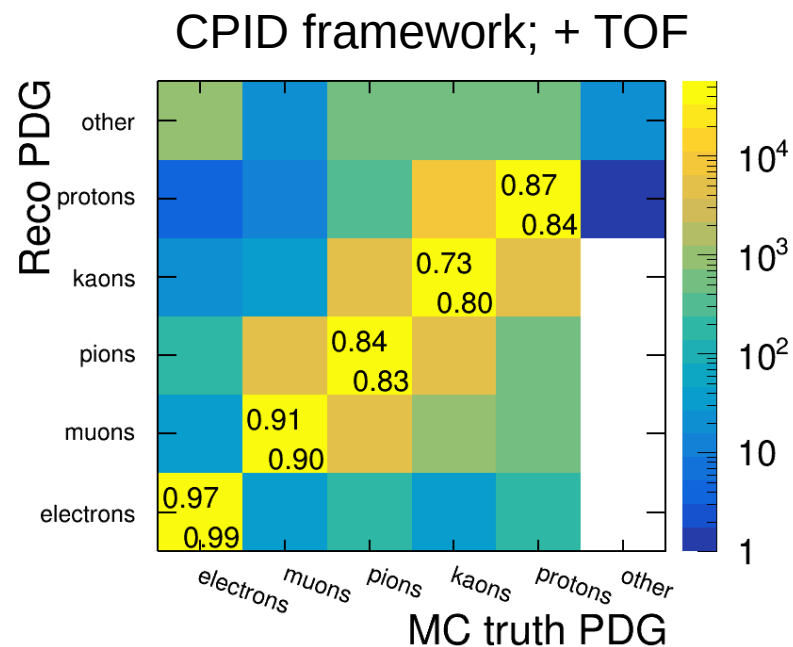
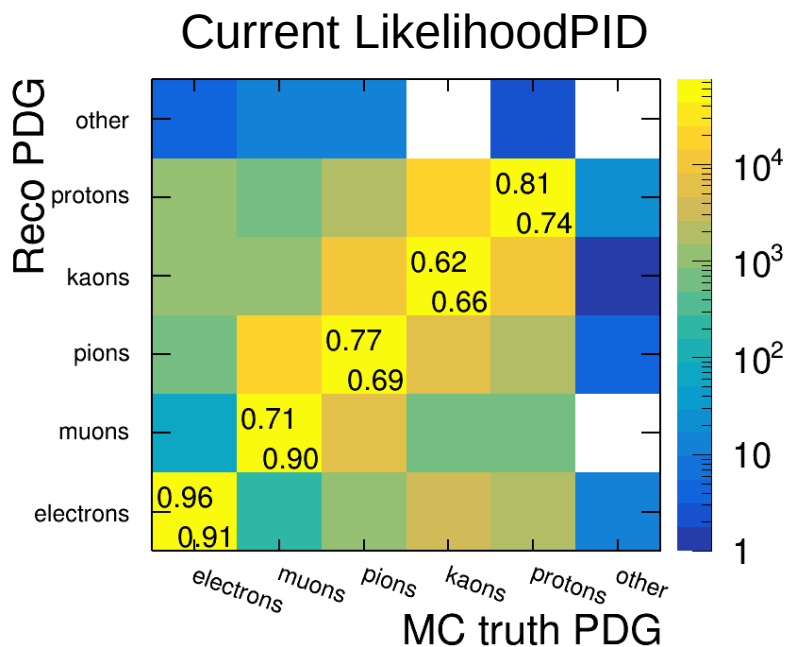
# Performance Comparison

- Here: multiclass BDT; confusion matrix with  $\text{eff}/\text{pur}$  on diagonal
- Simple BDT already generates similar performance to current LikelihoodPID
- Using dE/dx and Pandora (based on cluster shapes)



# Performance Comparison

- Here: multiclass BDT; confusion matrix with  $\text{eff} / \text{pur}$  on diagonal
- Simple BDT already generates similar performance to current LikelihoodPID
- Addition of TOF gives immediately better result – previously hard, easy in CPID



- First CPID version online – beta version, no warranty!
- Feel free to test and send feedback & feature requests!
- More features
  - (abstract) RICH algorithm?
  - neural network training, possible interface with pyTorch
  - more performance assessment output
  - more documentation...
- More plots → understand combinations, scalings, etc.
- Intention: add CPID to ILD standard high level reco and add PID estimator to the PFO, possibly a conservative and an ambitious one
- Make this available in native Key4HEP / Gaudi

