

# Update of Electron Identification Performance Based on BDTs

Hai-Jun Yang  
University of Michigan, Ann Arbor  
(with T. Dai, X. Li, A. Wilson, B. Zhou)

BNL Analysis Jamboree  
December 18, 2008

# Motivation

- Lepton ( $e$ ,  $\mu$ ,  $\tau$ ) Identification with high efficiency is crucial for new physics discoveries at the LHC
- Great efforts in ATLAS to develop the algorithms for electron identification:
  - Cut-based algorithm: IsEM
  - Multivariate algorithms: Likelihood and BDT
- Further improvement could be achieved with better treatment of the multivariate training using the Boosted Decision Trees technique

# Electron ID Studies with BDT

## Select electrons in two steps

- 1) Pre-selection: an EM cluster matching a track
- 2) Apply electron ID based on pre-selected samples with different e-ID algorithms (IsEM, Likelihood ratio, AdaBoost and **EBoost**).

## New BDT e-ID development at U. Michigan (Rel. v12)

- H. Yang's talk at US-ATLAS Jamboree on Sept. 10, 2008

<http://indico.cern.ch/conferenceDisplay.py?confId=38991>

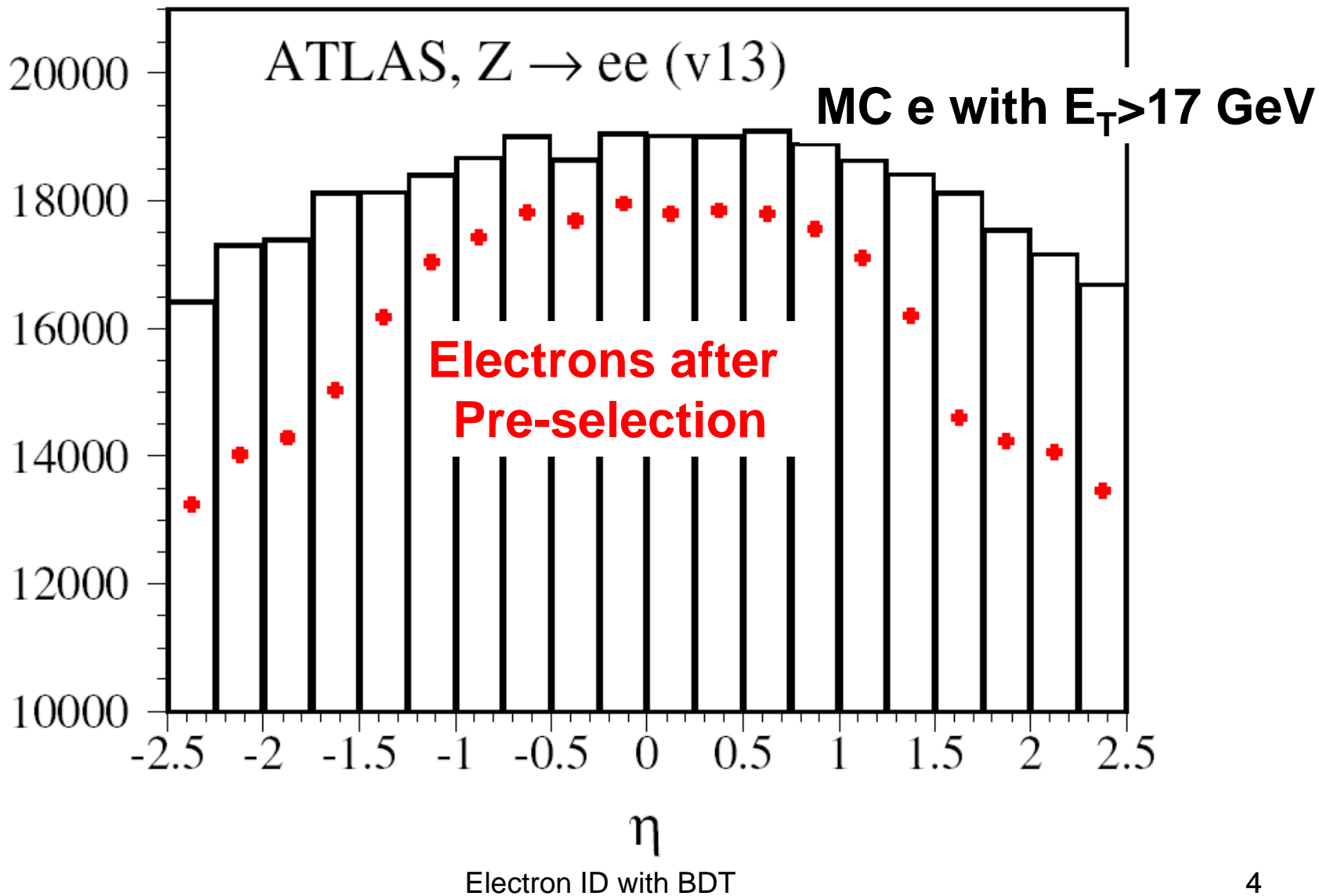
## New BDT e-ID (**EBoost**) based on Rel. v13

- H. Yang's talk at ATLAS performance and physics workshop at CERN on Oct. 2, 2008

<http://indico.cern.ch/conferenceDisplay.py?confId=39296>

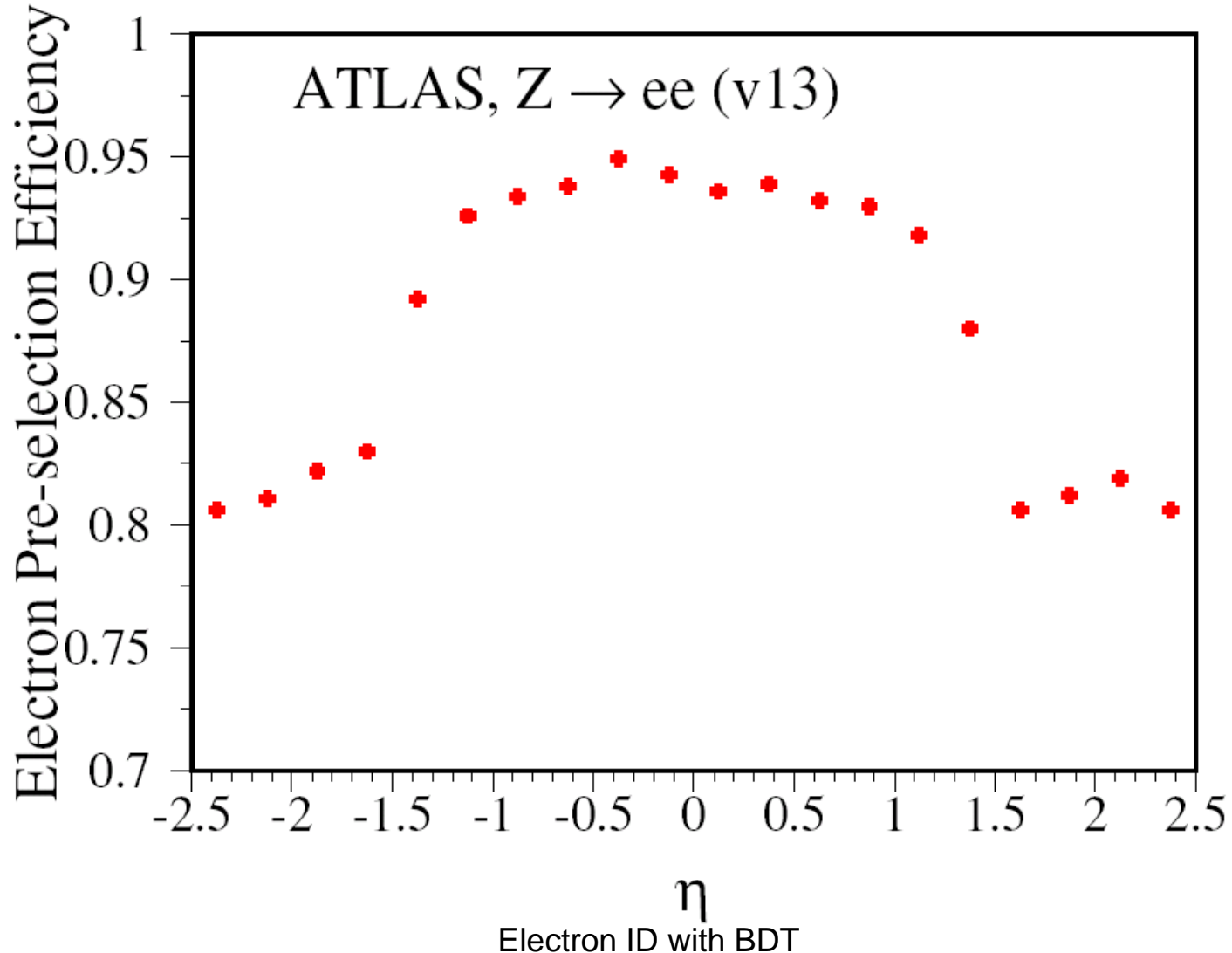
## Implementation of **EBoost** in EgammaRec (Rel. v14)

# Electrons



# Electron Pre-selection Efficiency

The inefficiency mainly due to track matching



# BDT e-ID (EBoost) Training

- BDT multivariate pattern recognition technique:
  - [ H. Yang et. al., NIM A555 (2005) 370-385 ]
- BDT e-ID training signal and backgrounds (jet faked e)
  - $W \rightarrow e\nu$  as electron signal (DS 5104, v13)
  - JF17 (DS 5802, v13)
- Using the same e-ID variables as IsEM for training (see variable list in next page)
- BDT e-ID training procedure
  - Apply additional cuts on the training samples to select hardly identified jet faked electron as background for BDT training to make the BDT training more effective.
  - Apply event weight to high  $P_T$  backgrounds to effectively reduce the jet fake rate at high  $P_T$  region. Event weight training technique reference, [ H. Yang et. al., JINST 3 P04004 (2008) ]

# Variables Used for BDT e-ID (EBoost)

The same variables for IsEM are used

## ▶ `egammaPID::ClusterHadronicLeakage`

fraction of transverse energy in TileCal 1<sup>st</sup> sampling

## ▶ `egammaPID::ClusterMiddleSampling`

Ratio of energies in 3\*7 & 7\*7 window

Ratio of energies in 3\*3 & 7\*7 window

Shower width in LAr 2<sup>nd</sup> sampling

Energy in LAr 2<sup>nd</sup> sampling

## ▶ `egammaPID::ClusterFirstSampling`

Fraction of energy deposited in 1<sup>st</sup> sampling

Delta E<sub>max2</sub> in LAr 1<sup>st</sup> sampling

E<sub>max2</sub>-E<sub>min</sub> in LAr 1<sup>st</sup> sampling

Total shower width in LAr 1<sup>st</sup> sampling

Shower width in LAr 1<sup>st</sup> sampling

F<sub>side</sub> in LAr 1<sup>st</sup> sampling

## ▶ `egammaPID::TrackHitsA0`

B-layer hits, Pixel-layer hits, Precision hits

Transverse impact parameter

## ▶ `egammaPID::TrackTRT`

Ratio of high threshold and all TRT hits

## ▶ `egammaPID::TrackMatchAndEoP`

Delta eta between Track and egamma

Delta phi between Track and egamma

E/P – egamma energy and Track momentum ratio

## ▶ `Track Eta and EM Eta`

## ▶ `Electron isolation variables:`

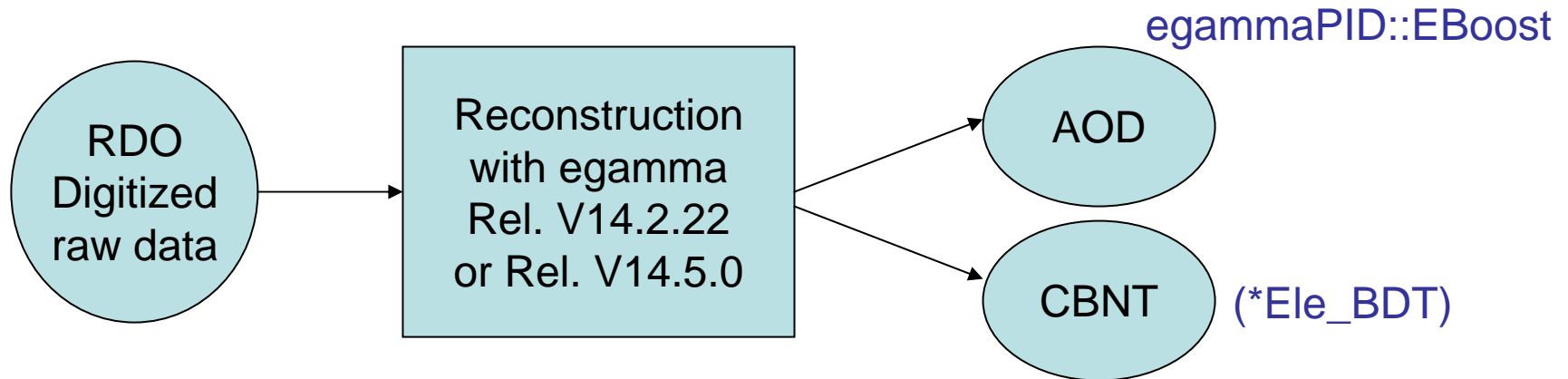
*Number of tracks ( $\Delta R=0.3$ )*

*Sum of track momentum ( $\Delta R=0.3$ )*

*Ratio of energy in EtCone45 / E<sub>T</sub>*

# Implementation of BDT Trees in Egamma Package and Test

- E-ID based on BDT has been implemented into egamma reconstruction package (private).
- We successfully run through the reconstruction package based on v14.2.22 and v14.5.0 to test the BDT e-ID.

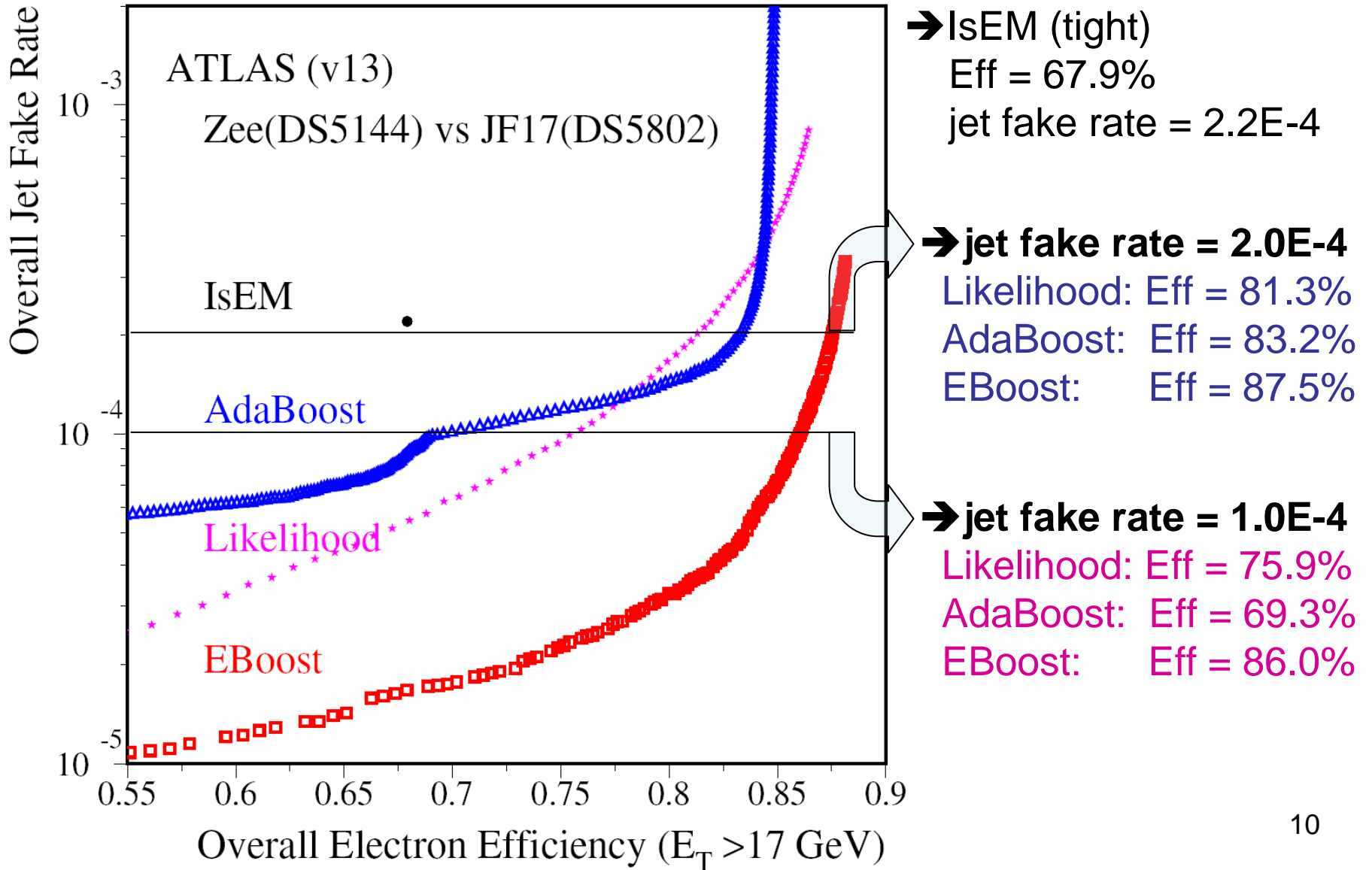




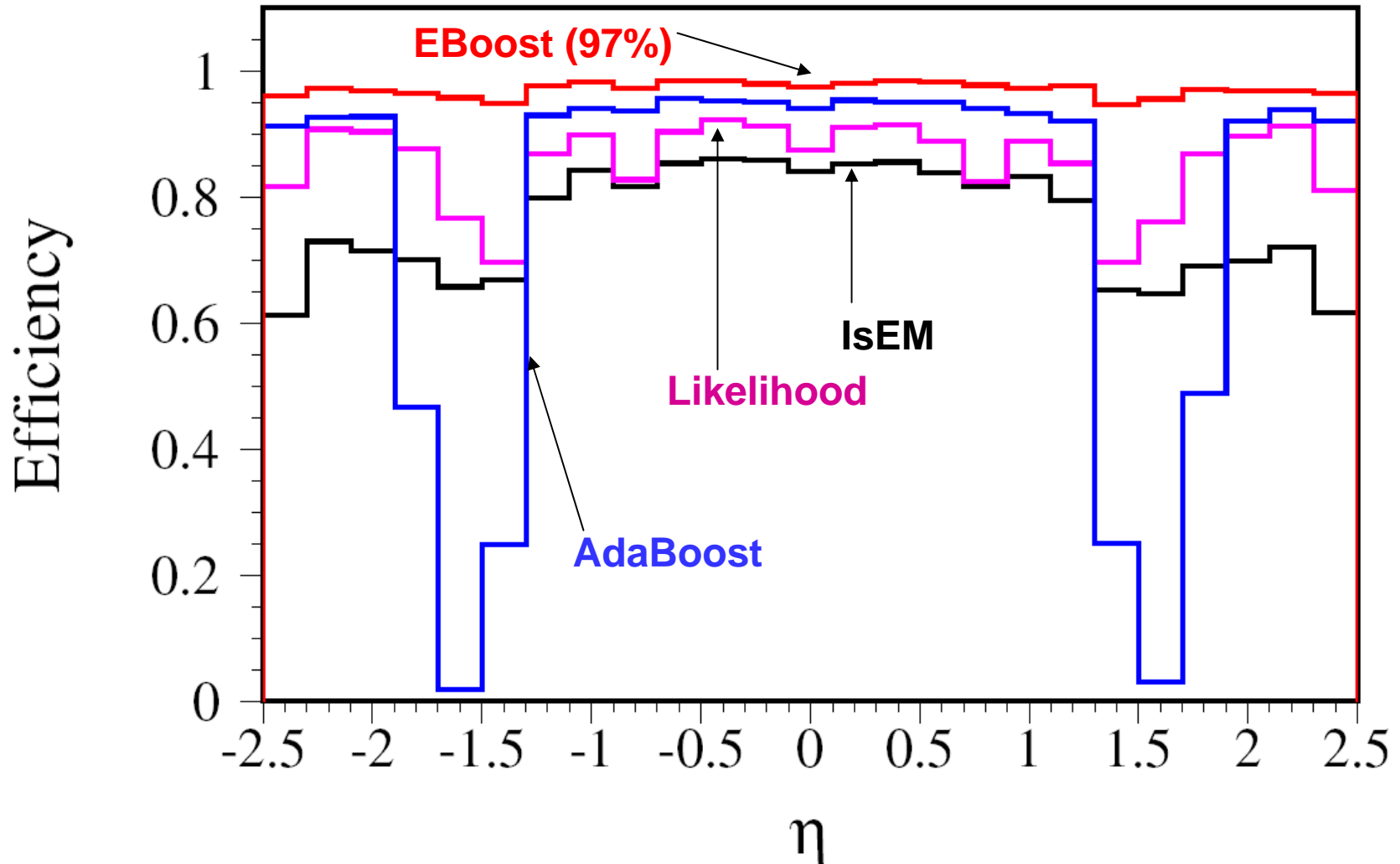
# E-ID Testing Samples Produced at $\sqrt{s} = 14 \text{ TeV}$ (v13)

- Wenu: DS5104 (Eff\_precuts = 88.6%)
  - 42020 electrons with  $E_t > 17 \text{ GeV}$ ,  $|\eta| < 2.5$
  - 37230 electrons after pre-selection cuts
- Zee: DS5144 (Eff\_precuts = 88.6%)
  - 181281 electron with  $E_t > 17 \text{ GeV}$ ,  $|\eta| < 2.5$
  - 160615 electrons after pre-selection cuts
- JF17: DS5802 (Eff\_precuts = 2.4%)
  - 1946968 events, 7280046 reconstructed jets
  - 176727 jets after pre-selection

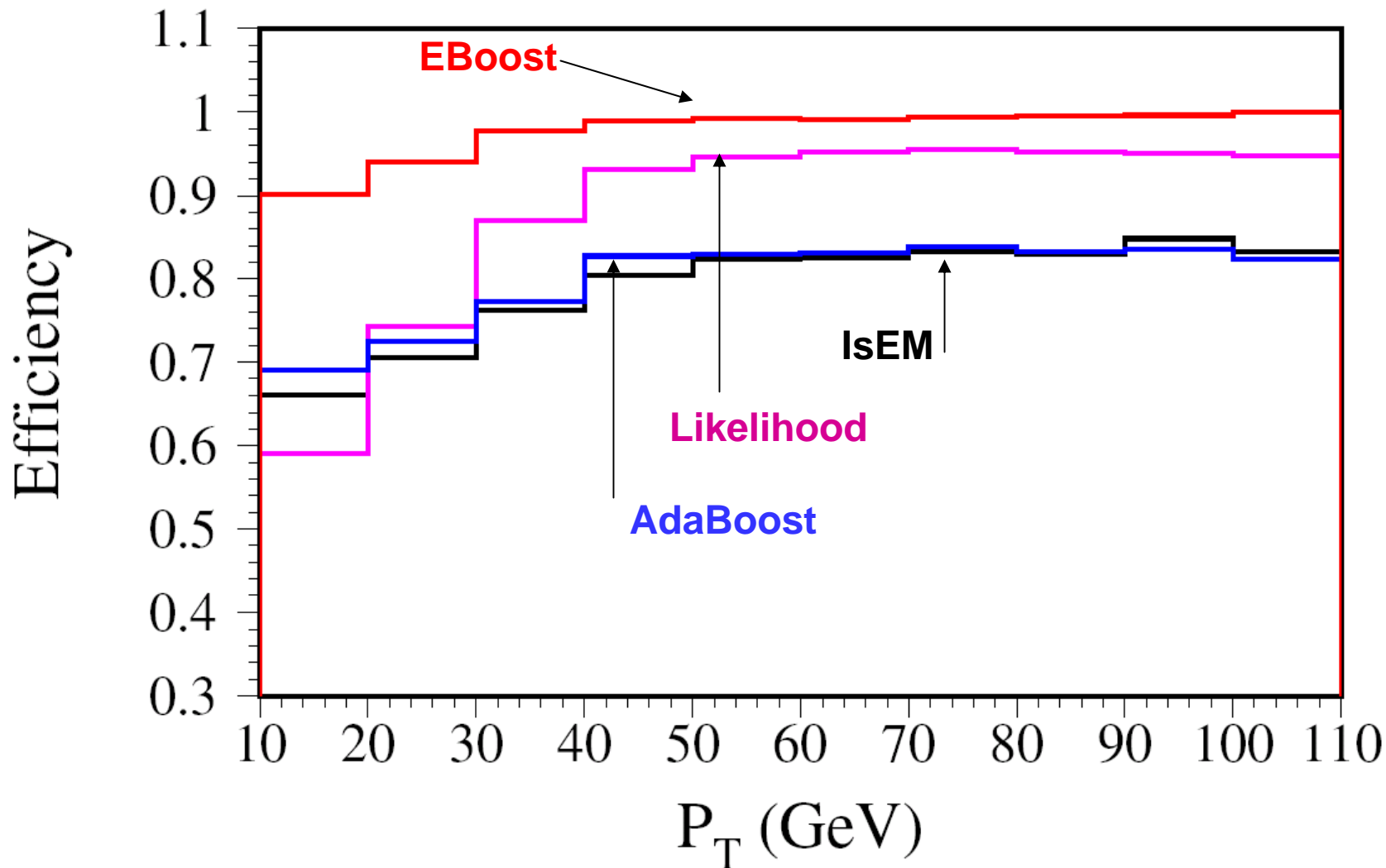
# Comparison of e-ID Algorithms (v13)



# E-ID Efficiency after pre-selection vs $\eta$ (v13, jet fake rate=1.0E-4)



# E-ID Efficiency after pre-selection vs $P_T$ (v13, jet fake rate=1.0E-4)



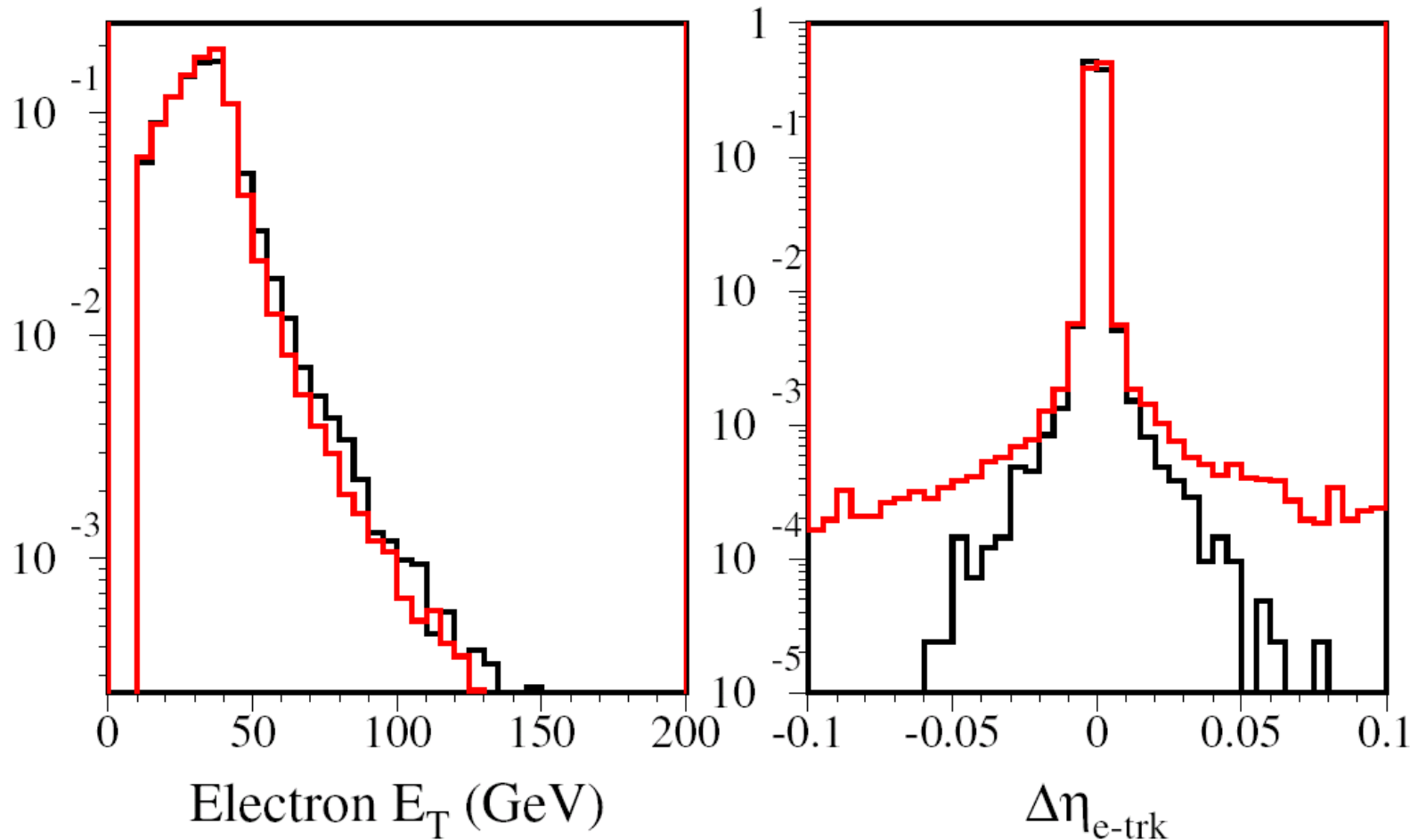
# E-ID Testing Samples Produced at $\sqrt{s} = 10 \text{ TeV}$ (v14)

- **Wenu: DS106020 (Eff\_precuts = 86.7%)**
  - 58954 electrons with  $E_t > 17 \text{ GeV}$ ,  $|\eta| < 2.5$
  - 51100 electrons after pre-selection cuts
- **Zee: DS106050 (Eff\_precuts = 86.7%)**
  - 108550 electrons with  $E_t > 17 \text{ GeV}$ ,  $|\eta| < 2.5$
  - 94153 electrons after pre-selection cuts
- **JF17: DS105802 (Eff\_precuts = 2.34%)**
  - 237950 events, 896818 reconstructed jets
  - 20994 jets after pre-selection cuts

# Variable distribution Comparison

## 14 TeV vs 10 TeV

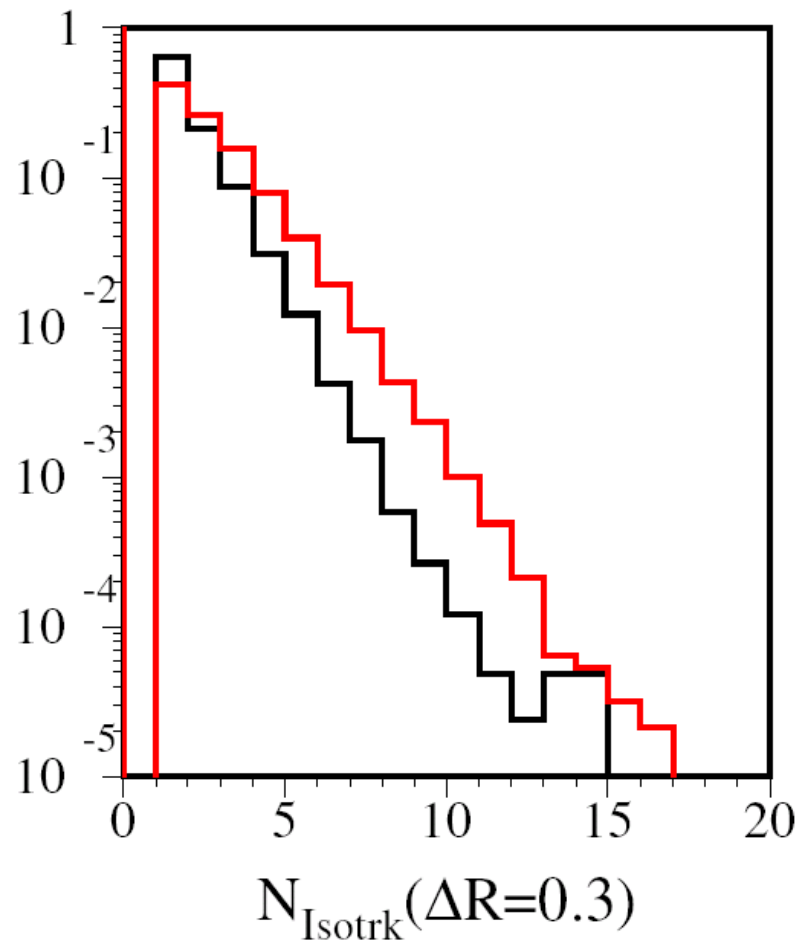
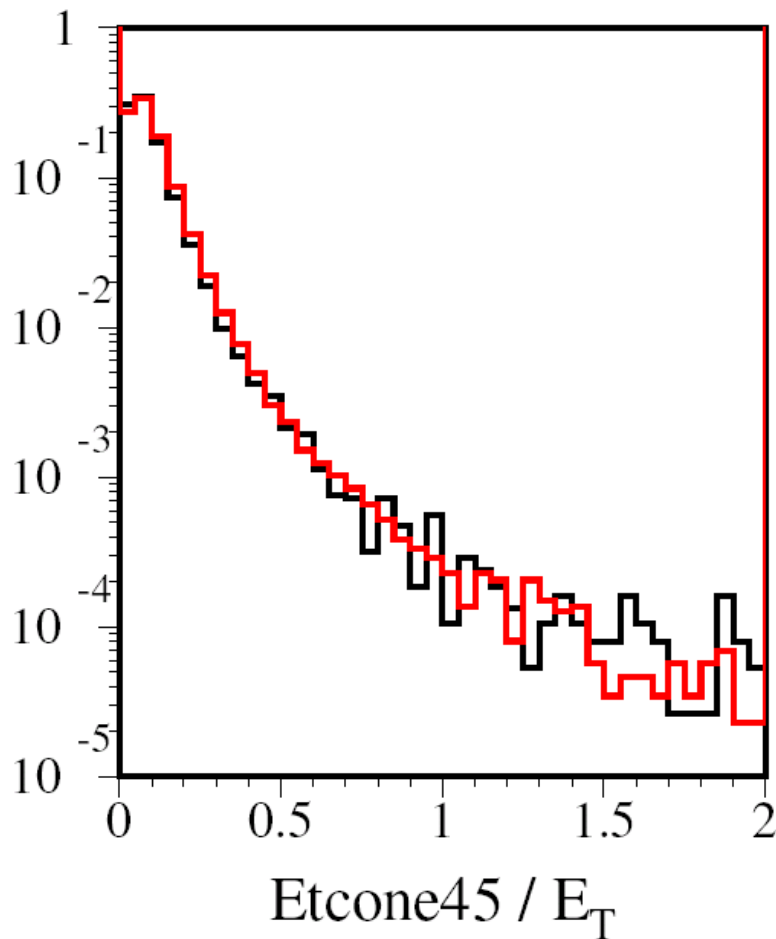
$W \rightarrow e\nu$ , DS5104(14TeV,black) vs DS106020(10TeV,red)



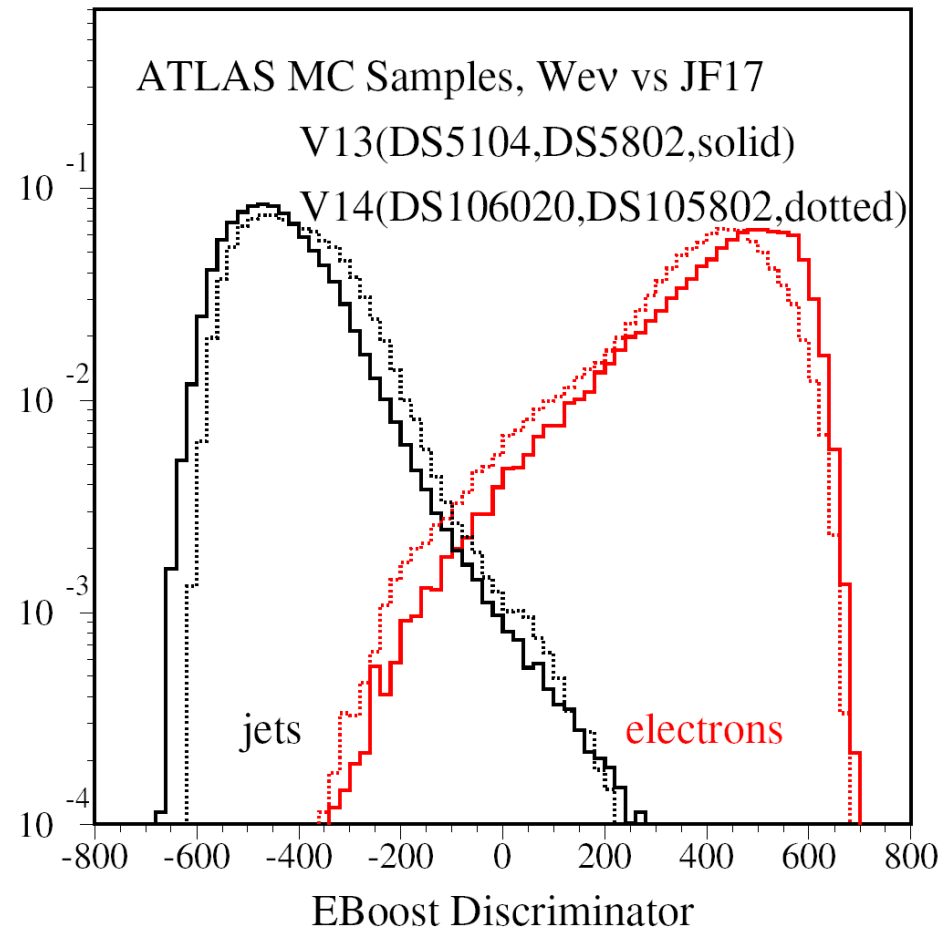
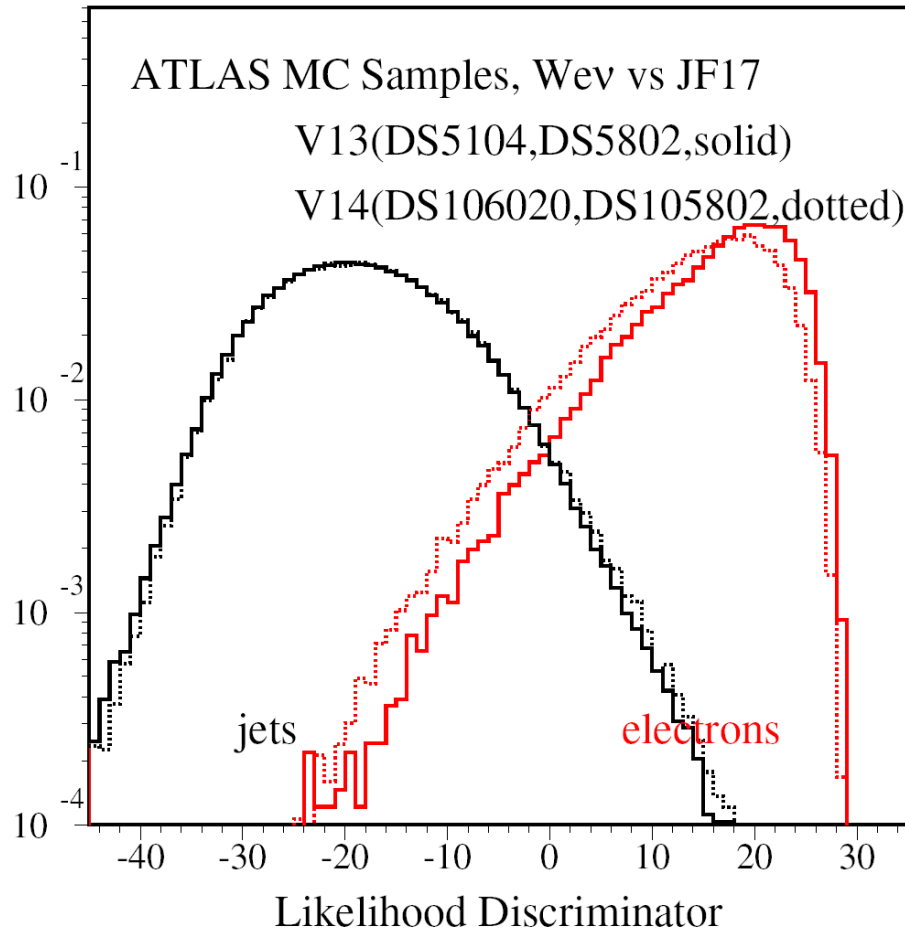
# Variable distribution Comparison

## 14 TeV vs 10 TeV

$W \rightarrow e\nu$ , DS5104(14TeV,black) vs DS106020(10TeV,red)

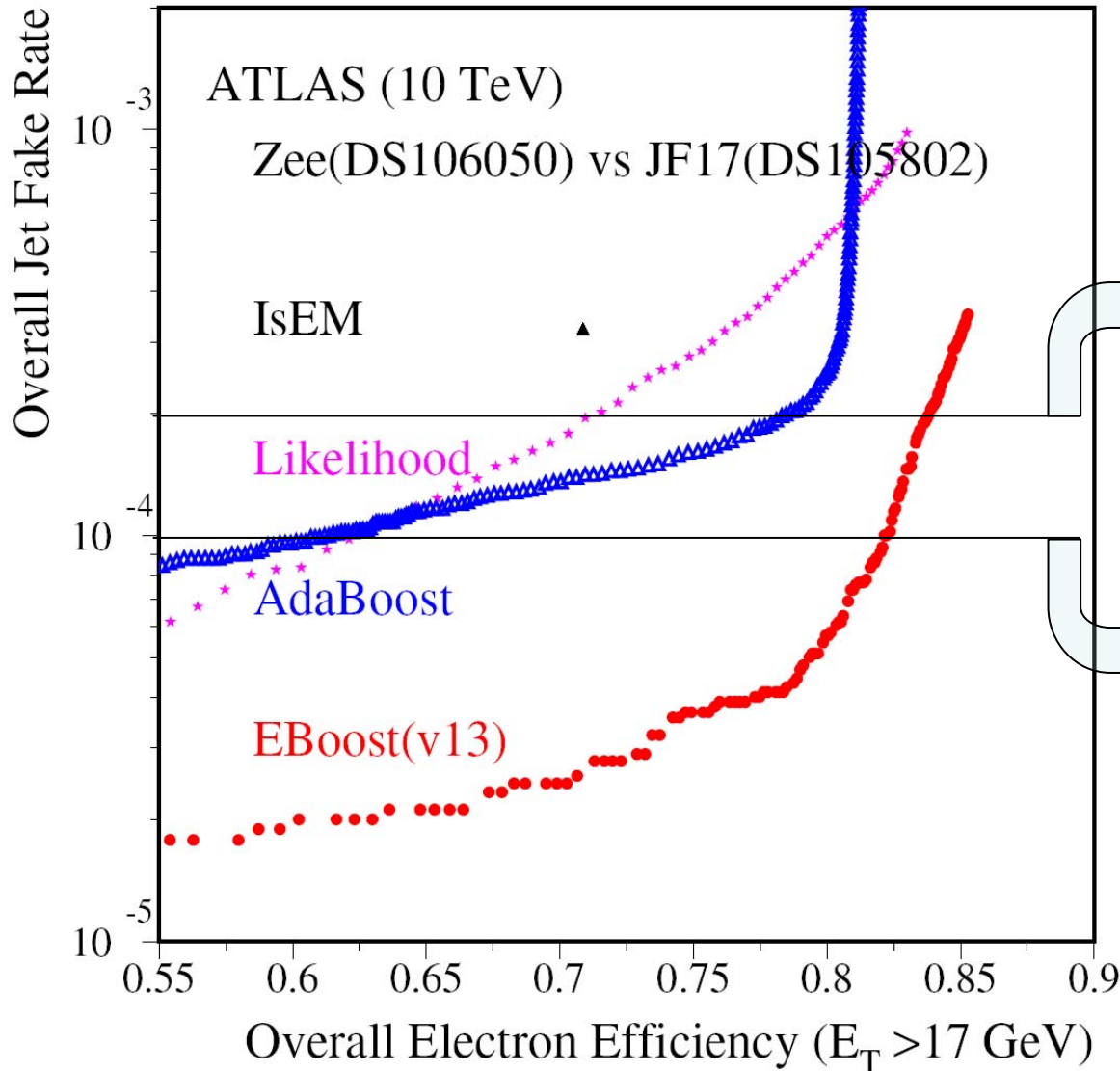


# E-ID Discriminators with no retraining for 10 TeV MC Samples





# Comparison of e-ID Algorithms (v14)



→ IsEM (tight)  
Eff = 70.9%  
jet fake rate =  $3.2E-4$

→ jet fake rate =  $2.0E-4$   
Likelihood: Eff = 71.6%  
AdaBoost: Eff = 78.5%  
EBoost: Eff = 83.9%

→ jet fake rate =  $1.0E-4$   
Likelihood: Eff = 62.9%  
AdaBoost: Eff = 61.2%  
EBoost: Eff = 82.2%

# Robustness of Multivariate e-ID

( $\sqrt{s} = 14$  TeV vs. 10 TeV without retraining)

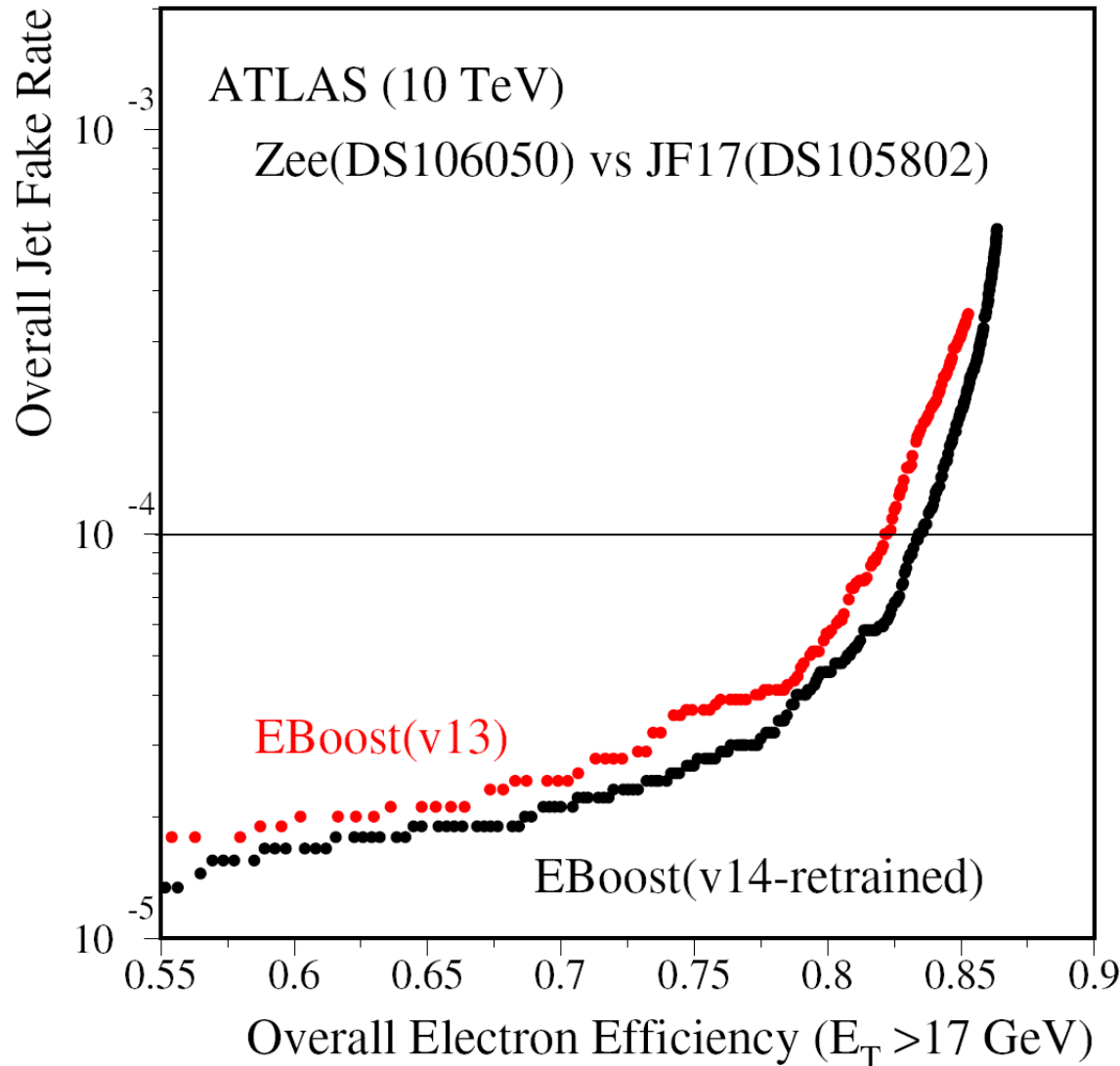
Test MC	Precuts	Likelihood	AdaBoost	EBoost
Z $\rightarrow$ ee (v13) $\sqrt{s} = 14$ TeV	88.6%	75.9%	69.3%	86.0%
Z $\rightarrow$ ee (v14) $\sqrt{s} = 10$ TeV	86.7%	62.9%	61.2%	82.2%
Eff. Change after pre-sel	-1.9%	-13.0%	-8.1%	-3.8%
JF17 (v13) $\sqrt{s} = 14$ TeV	2.4E-2	1.0E-4	1.0E-4	1.0E-4
JF17 (v14) $\sqrt{s} = 10$ TeV	2.3E-2	1.0E-4	1.0E-4	1.0E-4

# Robustness of Multivariate e-ID

( $\sqrt{s} = 14$  vs 10 TeV without retraining)

Test MC	Precuts	Likelihood	AdaBoost	EBoost
Z $\rightarrow$ ee (v13) $\sqrt{s} = 14$ TeV	88.6%	81.3%	83.2%	87.5%
Z $\rightarrow$ ee (v14) $\sqrt{s} = 10$ TeV	86.7%	71.6%	78.5%	83.9%
Eff. Change after pre-sel	-1.9%	-9.7%	-4.7%	-3.6%
		-7.8%	-2.8%	-1.7%
JF17 (v13) $\sqrt{s} = 14$ TeV	2.4E-2	2.0E-4	2.0E-4	2.0E-4
JF17 (v14) $\sqrt{s} = 10$ TeV	2.3E-2	2.0E-4	2.0E-4	2.0E-4

# Improvement with EBoost Re-training using $\sqrt{s} = 10$ TeV MC Samples



→ EBoost (retrained)  
Eff = 82.2 → 83.4%  
jet fake rate =  $1.0E-4$

→ For each major release multivariate e-ID should be retrained to obtain optimal performance

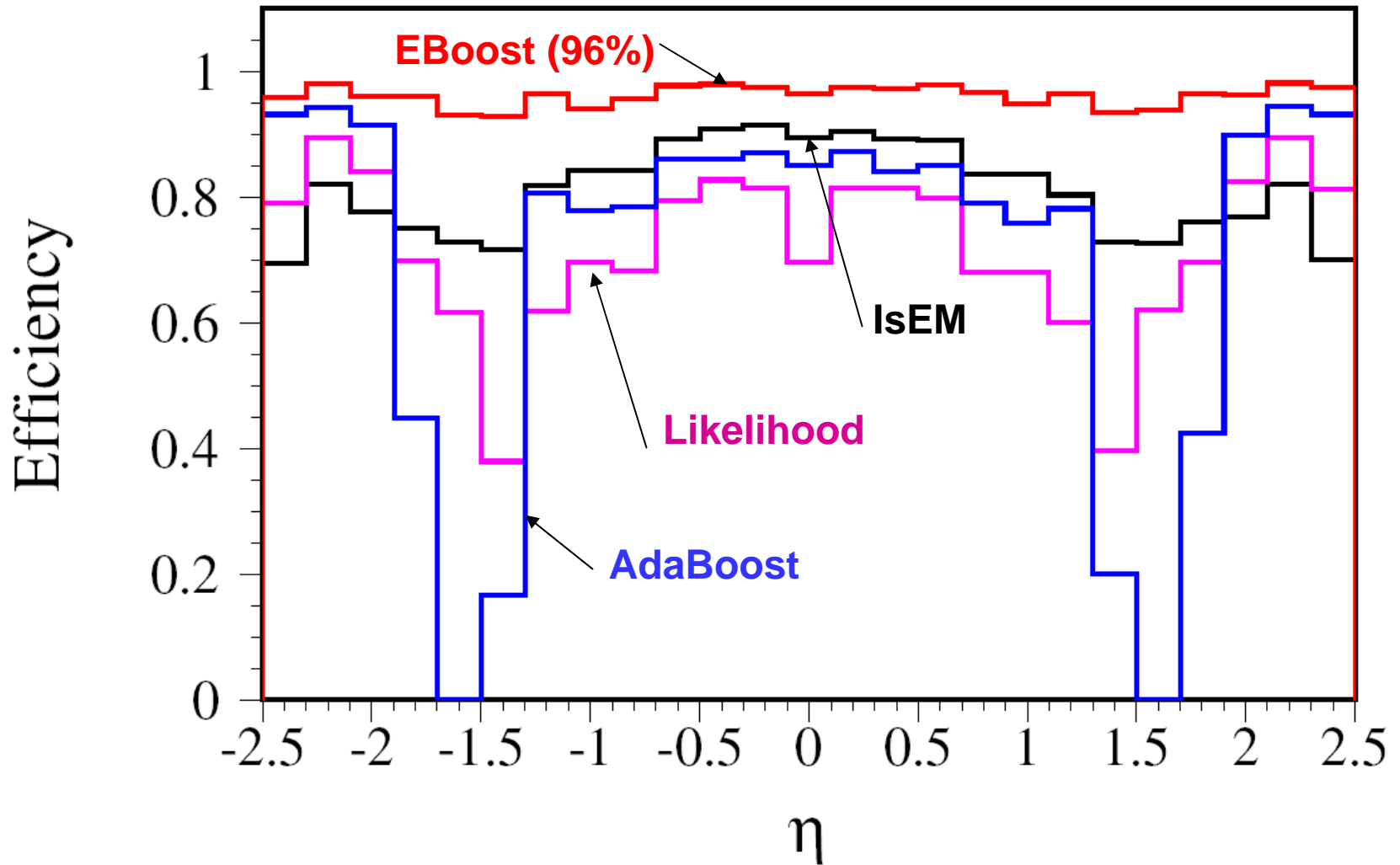
→ All multivariate e-ID should be retrained using real data

# Future Plan

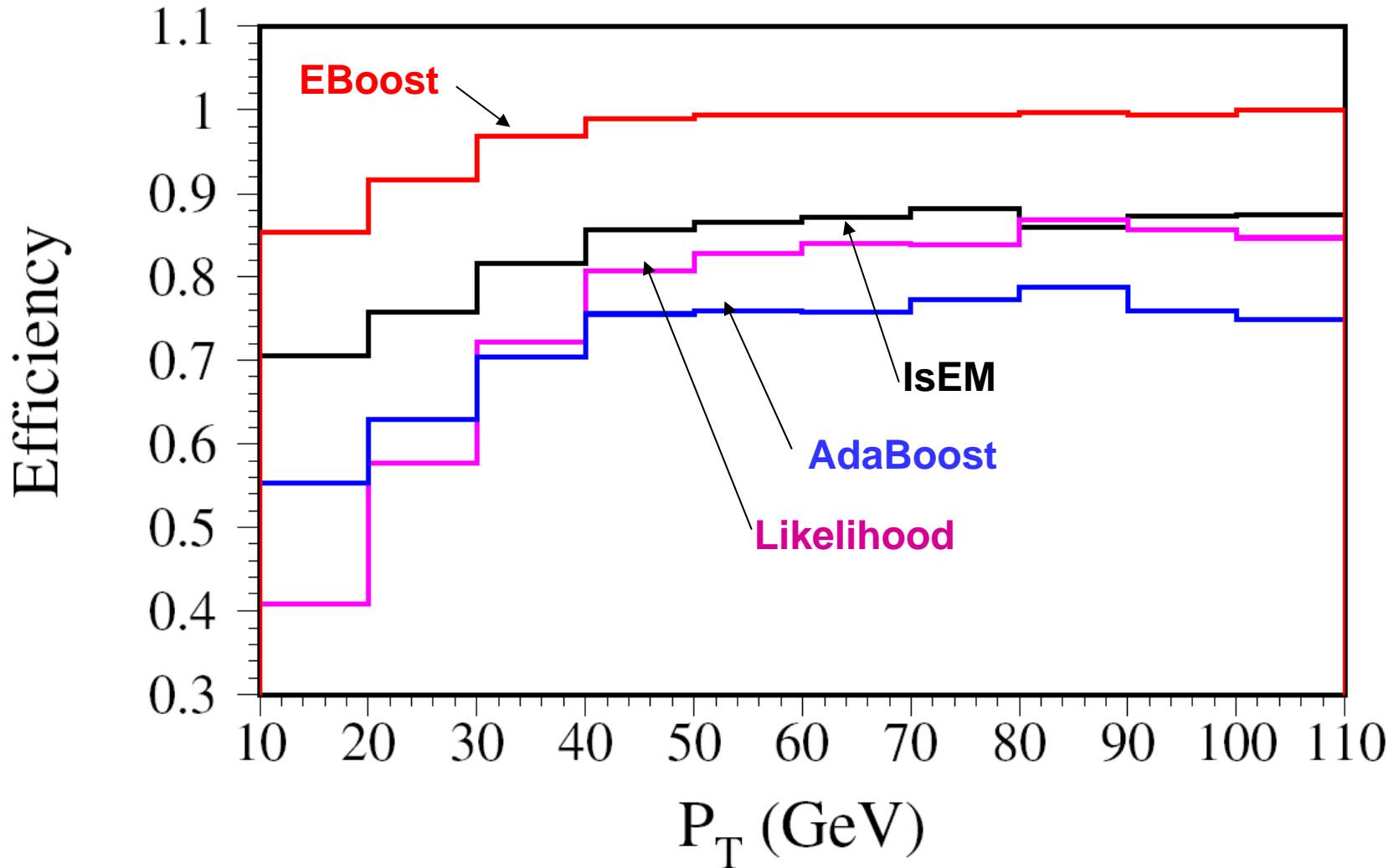
- We have requested to add the EBoost in ATLAS official egamma package and make EBoost discriminator variable available for more test and for physics analysis.
- We have proposed to provide EBoost trees to ATLAS egammaRec for each major software release
- We will explore new variables to further improve e-ID by suppressing  $\gamma$  converted electron etc.

# Backup Slides

# E-ID Efficiency after pre-selection vs $\eta$ (v14, jet fake rate=1.0E-4)

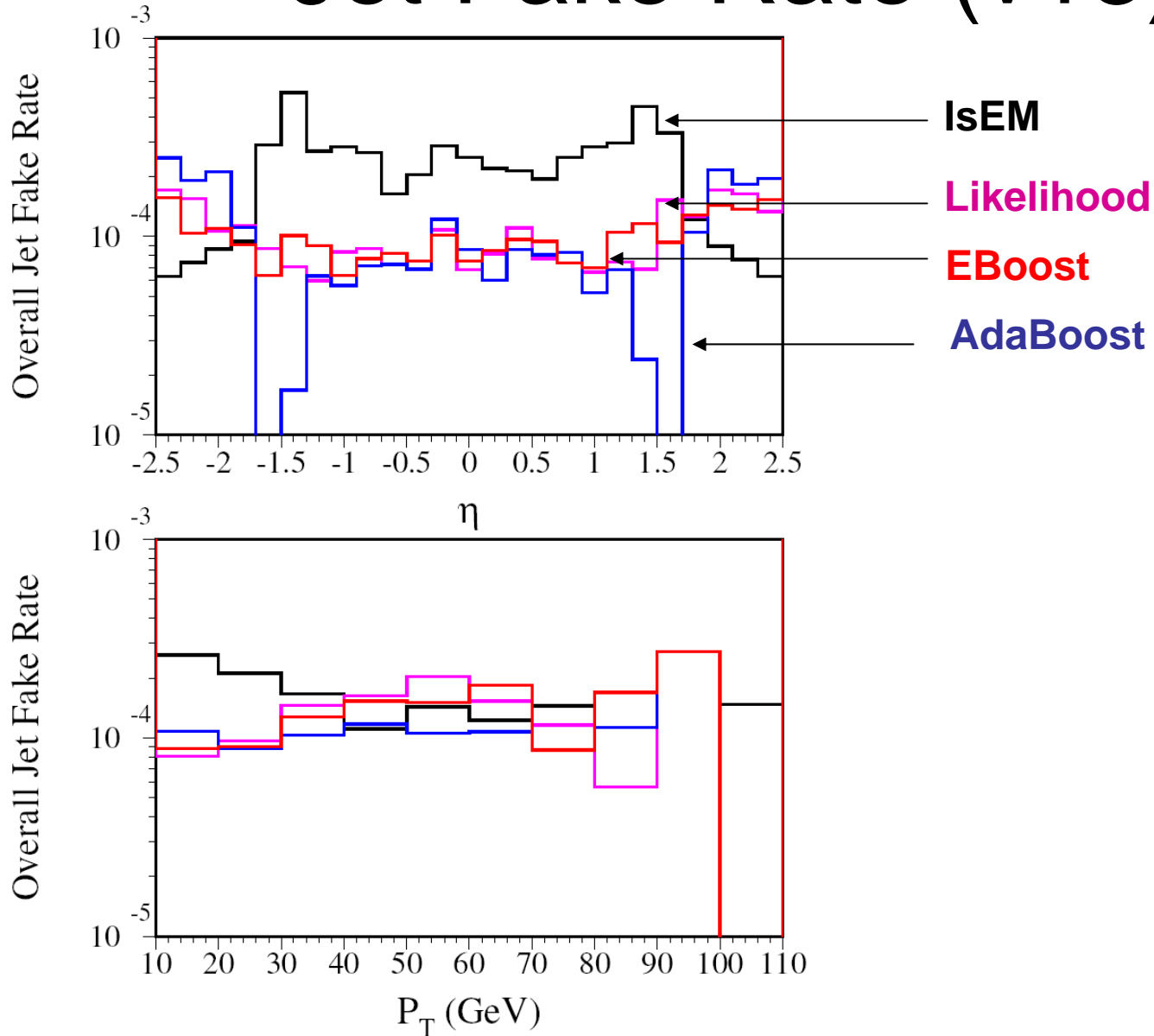


# E-ID Efficiency after pre-selection vs Pt (v14, jet fake rate=1.0E-4)

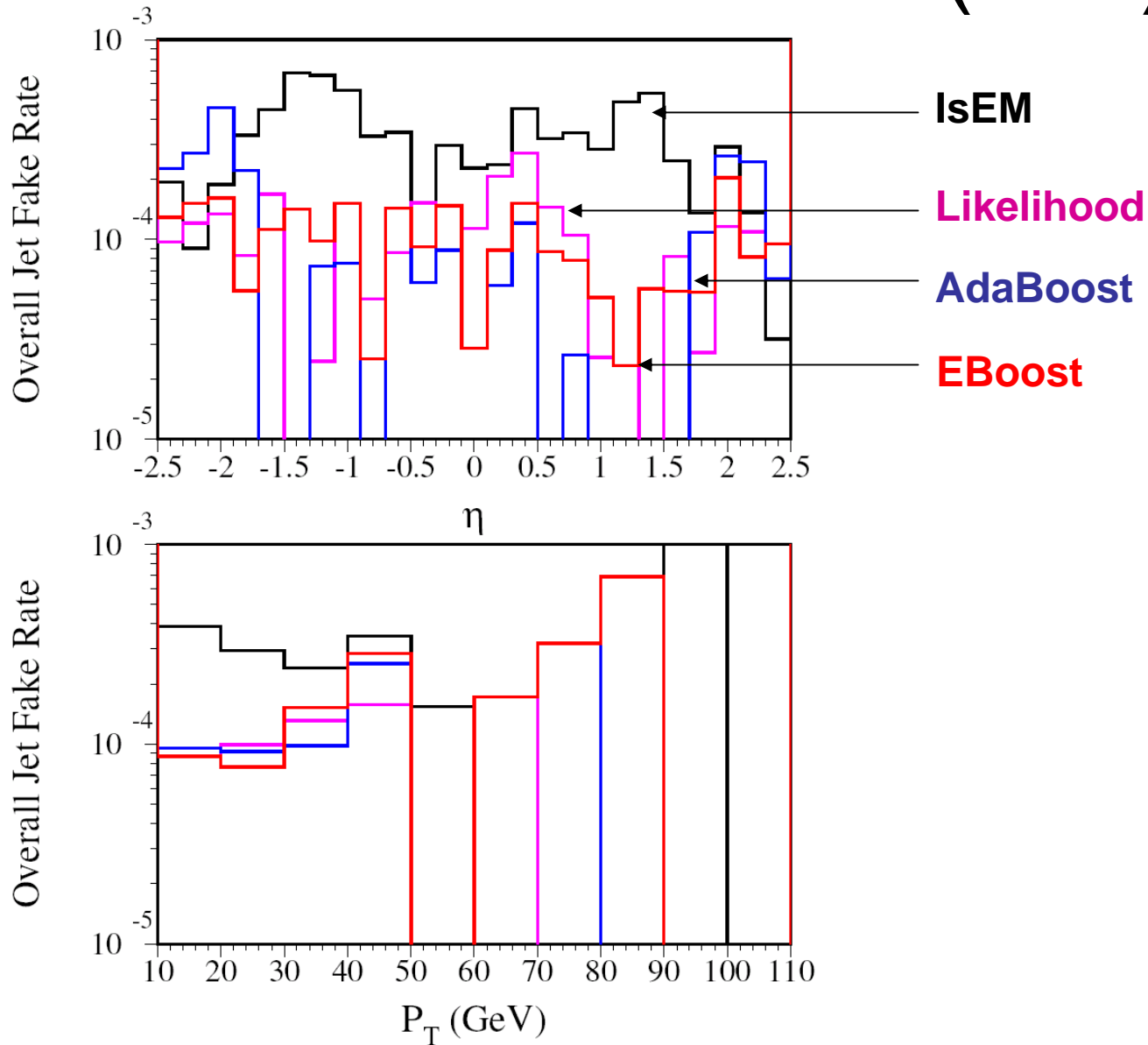




# Jet Fake Rate (v13)



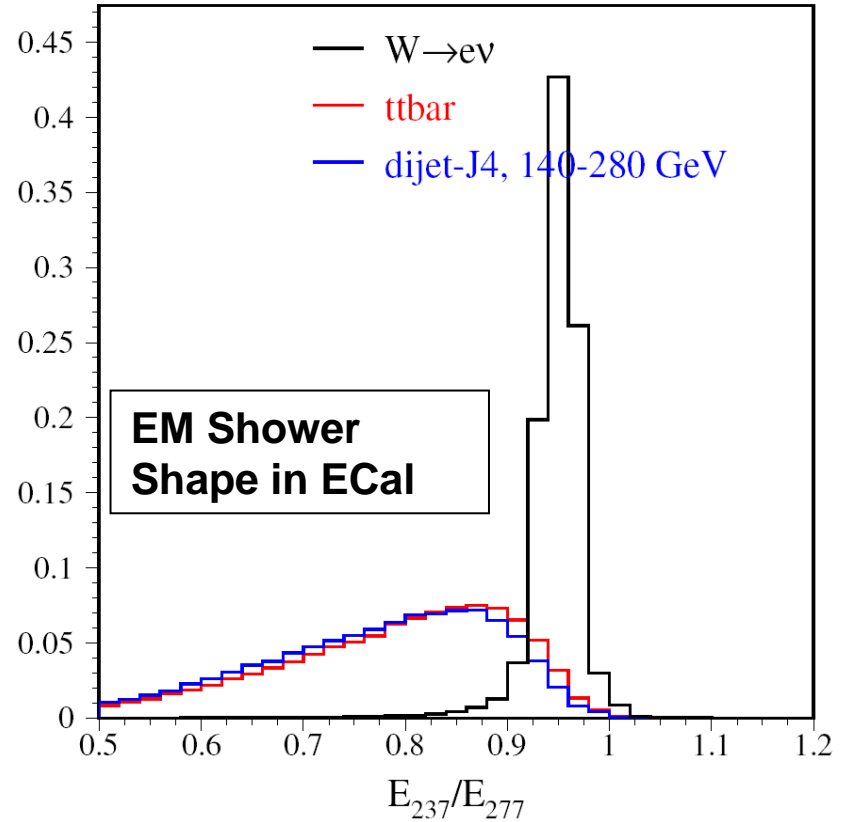
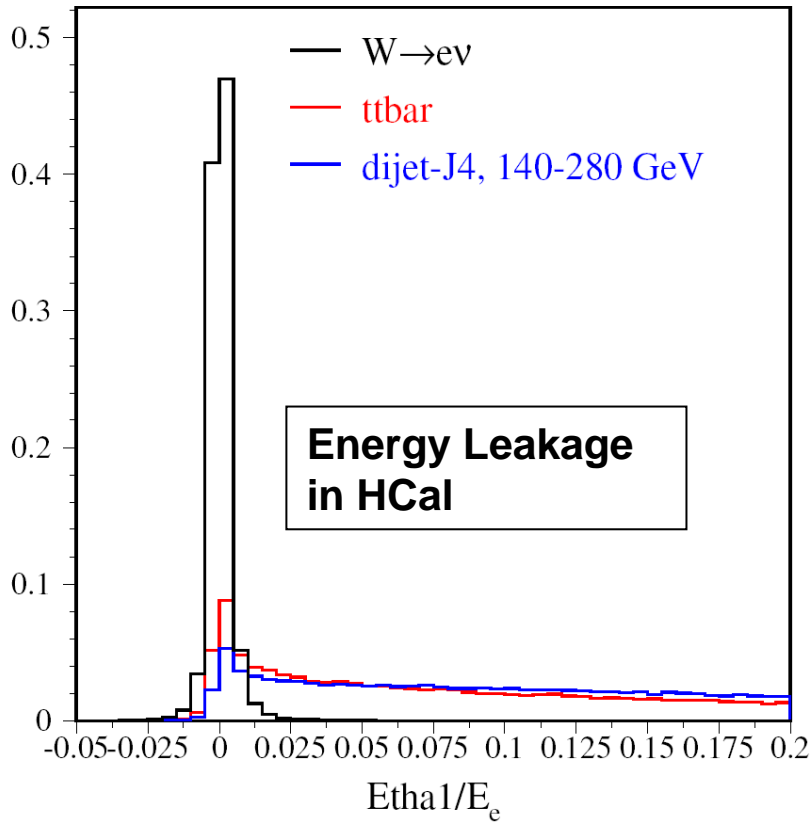
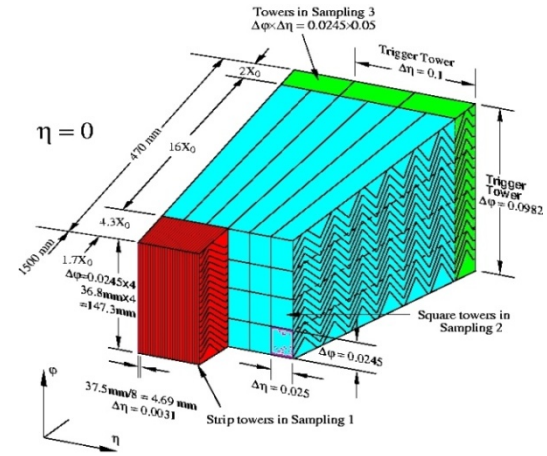
# Jet Fake Rate (v14)



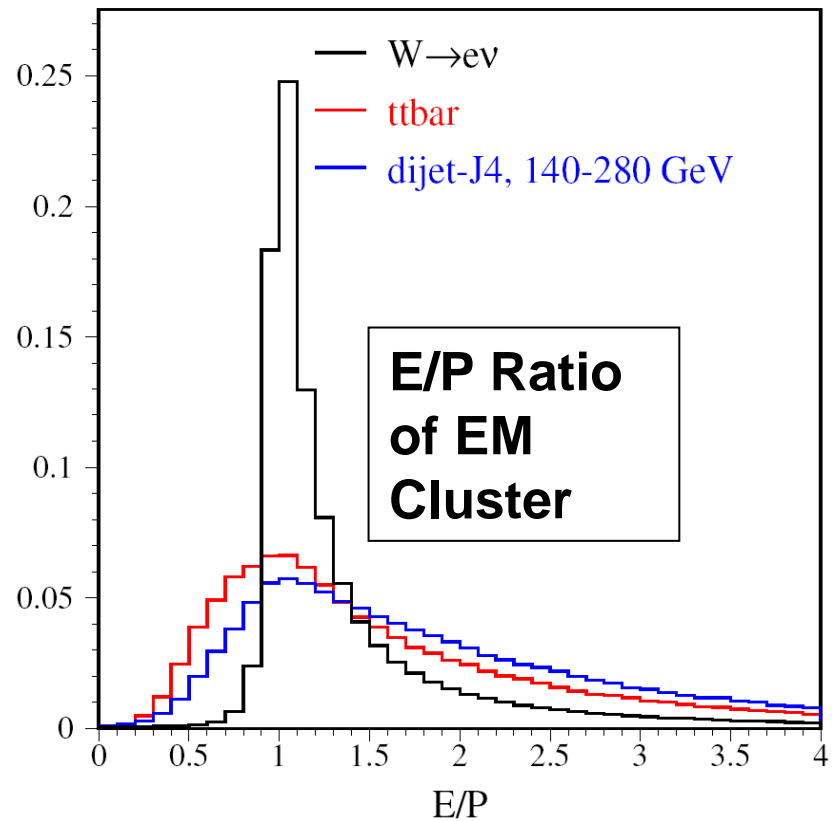
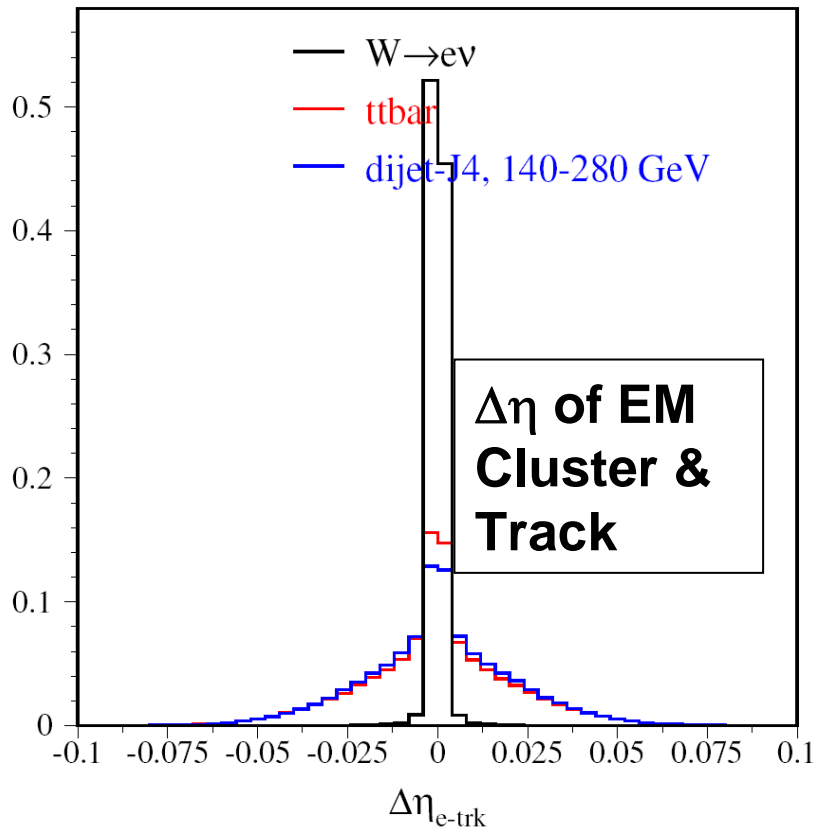
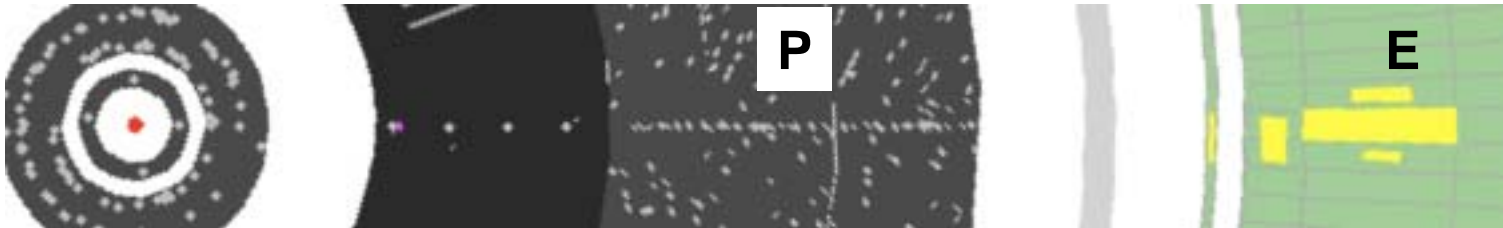
# List of Variables for BDT

1. Ratio of  $E_t(\Delta R=0.2-0.45) / E_t(\Delta R=0.2)$
2. Number of tracks in  $\Delta R=0.3$  cone
3. Energy leakage to hadronic calorimeter
4. EM shower shape  $E_{237} / E_{277}$
5.  $\Delta\eta$  between inner track and EM cluster
6. Ratio of high threshold and all TRT hits
7. Number of pixel hits and SCT hits
8.  $\Delta\phi$  between track and EM cluster
9.  $E_{\max 2} - E_{\min}$  in LAr 1<sup>st</sup> sampling
10. Number of B layer hits
11. Number of TRT hits
12.  $E_{\max 2}$  in LAr 1<sup>st</sup> sampling
13.  $E_{\text{overP}}$  – ratio of EM energy and track momentum
14. Number of pixel hits
15. Fraction of energy deposited in LAr 1<sup>st</sup> sampling
16.  $E_t$  in LAr 2<sup>nd</sup> sampling
17.  $\eta$  of EM cluster
18.  $D_0$  – transverse impact parameter
19. EM shower shape  $E_{233} / E_{277}$
20. Shower width in LAr 2<sup>nd</sup> sampling
21.  $\text{Frac}_{s1}$  – ratio of  $(E_{7\text{strips}} - E_{3\text{strips}}) / E_{7\text{strips}}$  in LAr 1<sup>st</sup> sampling
22. Sum of track  $P_t$  in  $\Delta R=0.3$  cone
23. Total shower width in LAr 1<sup>st</sup> sampling
24. Shower width in LAr 1<sup>st</sup> sampling

# EM Shower shape distributions of discriminating Variables (signal vs. background)

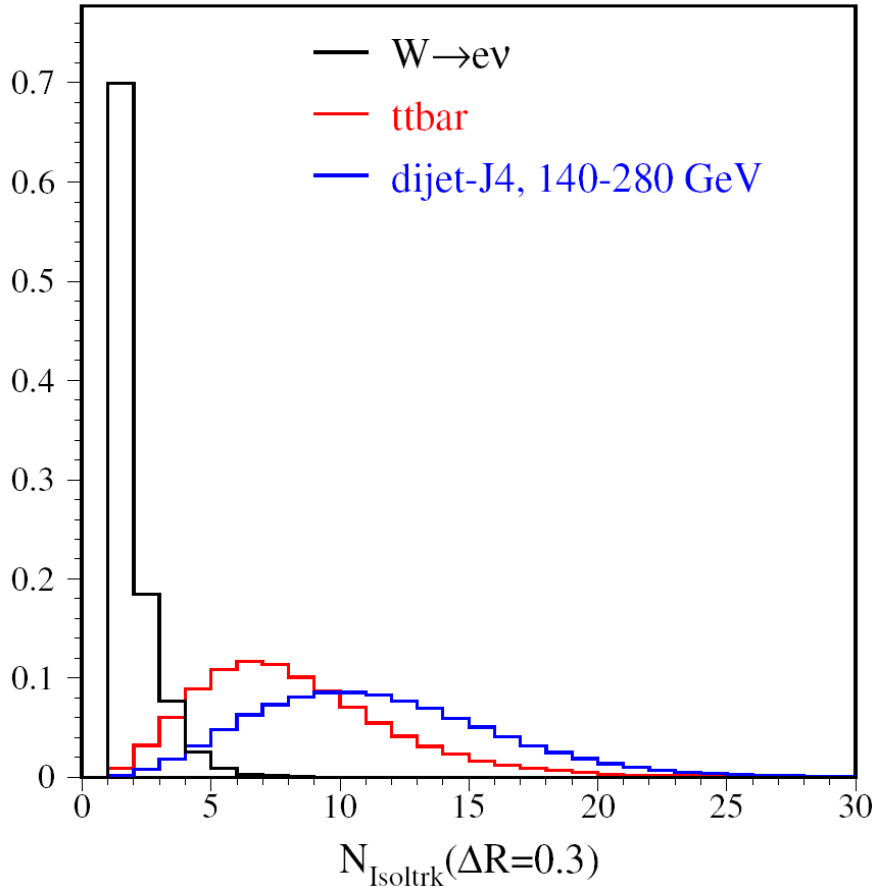


# ECal and Inner Track Match

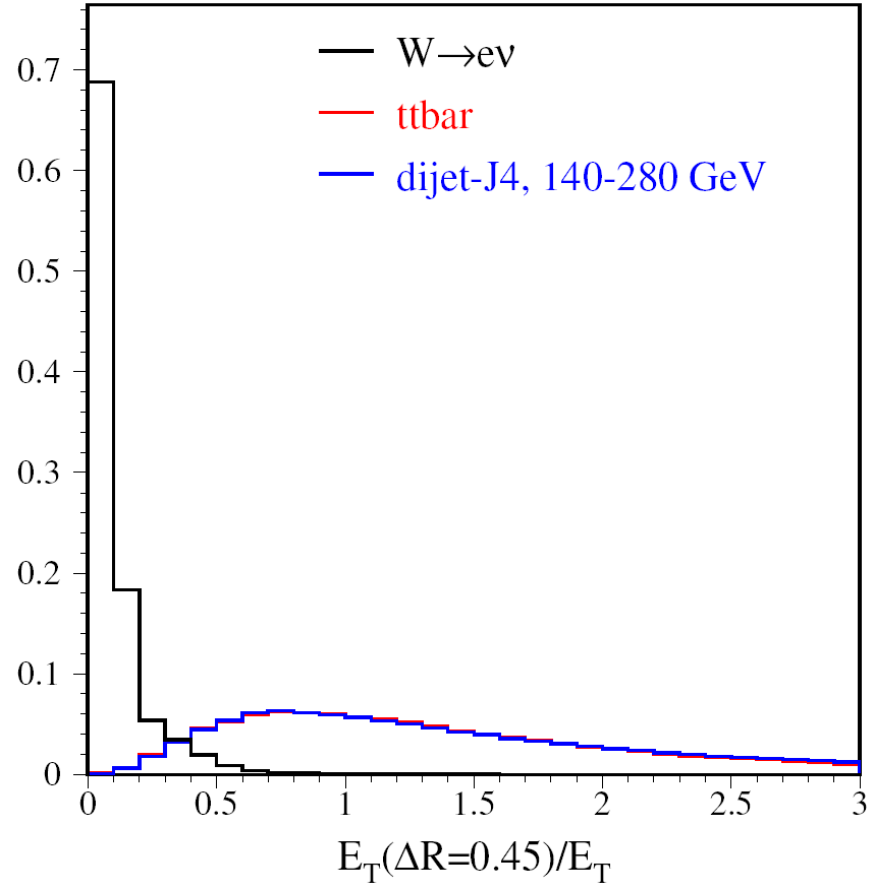


# Electron Isolation Variables

$N_{\text{trk}}$  around Electron Track



$E_T(\Delta R=0.2-0.45)/E_T$  of EM



# Signal Pre-selection: MC electrons

- MC True electron from  $W \rightarrow e\nu$  by requiring
  - $|\eta_e| < 2.5$  and  $E_T^{\text{true}} > 17 \text{ GeV}$  ( $N_e$ )
- Match MC e/ $\gamma$  to EM cluster:
  - $\Delta R < 0.2$  and  $0.5 < E_T^{\text{rec}} / E_T^{\text{true}} < 1.5$  ( $N_{\text{EM}}$ )
- Match EM cluster with an inner track:
  - $\text{eg\_trkmatchnt} > -1$  ( $N_{\text{EM/track}}$ )
- **Pre-selection Efficiency =  $N_{\text{EM/Track}} / N_e$**

# Pre-selection of Jet Faked Electrons

- Count number of reconstructed jets with
  - $|\eta_{\text{jet}}| < 2.5$  ( $N_{\text{jet}}$ )
- Loop over all EM clusters; each cluster matches with a jet
  - $E_{\text{T}}^{\text{EM}} > 17 \text{ GeV}$  ( $N_{\text{EM}}$ )
- Match EM cluster with an inner track:
  - $\text{eg\_trkmatchnt} > -1$  ( $N_{\text{EM/track}}$ )
- **Pre-selection Acceptance =  $N_{\text{EM/Track}} / N_{\text{jet}}$**