



Jet Flavour Tagging at e^+e^-

Franco Bedeschi, Loukas Gouskos, Michele Selvaggi
[full set of results: [arXiv:2202.03285](https://arxiv.org/abs/2202.03285)]

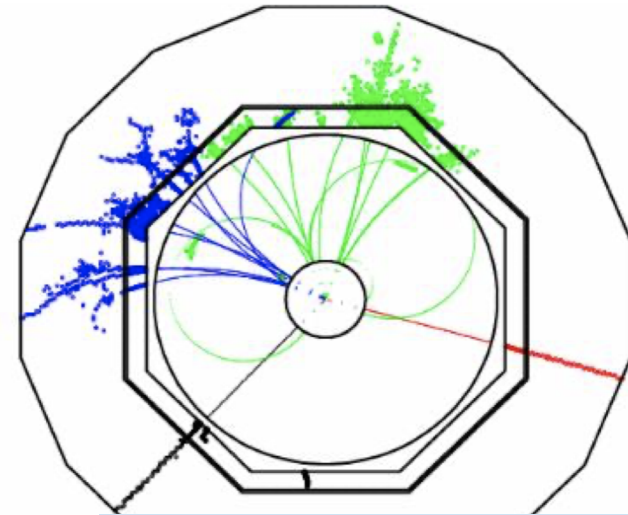
ECFA Higgs/Top/EW group Workshop
April 2022

Introduction

- Jet flavor identification (“tagging”): necessary tool to maximize physics outcome at colliders
- Today: Focus on e^+e^- colliders
 - provide a very clean environment
 - Much lower occupancy
 - no pileup compared to hadron colliders
- Scope of this work:

Build a general framework for developing flavor tagging algorithms for future colliders

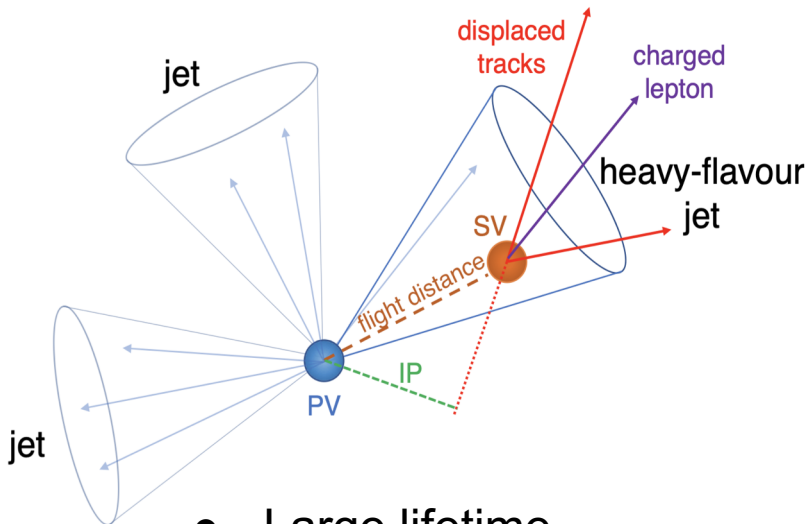
 - **Fast detector simulation**
 - Understand detector requirements/optimize design
 - e.g., Vertexing and PID capabilities of the FCCee detectors
 - **Develop a versatile jet flavor tagger for FCCee**
 - Identify with high purity gluon / light / strange / charm / bottom quarks



$e^+e^-: Z(\rightarrow\mu\mu)H(\rightarrow bb)$

Basics of flavour tagging

bottom/charm-tagging



- Large lifetime
- Displaced vertices/tracks
- Large track multiplicity
- Non-isolated e/ μ

Detector constraints:

Pixel/tracking detectors

- Little material, spatial resolution, precise track alignment

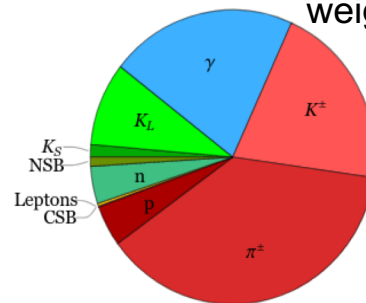
PID detectors:

- timing capabilities, energy loss (gas/silicon)

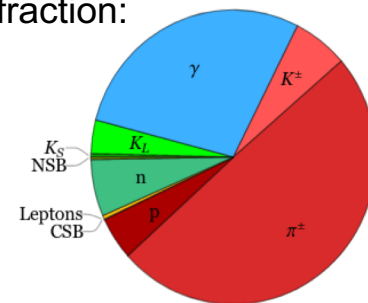
strange-tagging

[2003.09517]

Momentum weighted fraction:



Strange $p_T = 45$ GeV

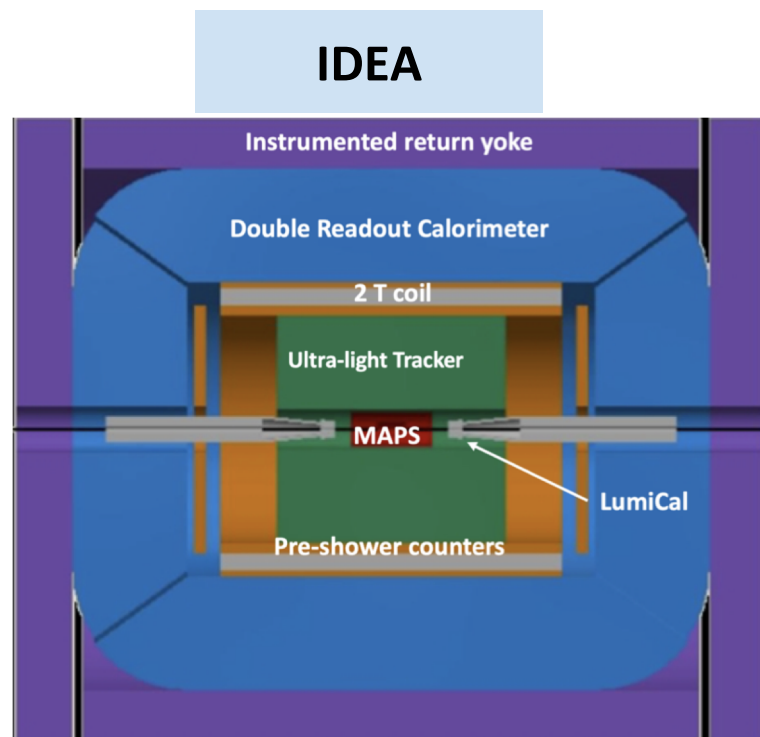


Down $p_T = 45$ GeV

- Large Kaon content
 - Charged Kaon as track:
 - K/pi separation
 - Neutral Kaons:
 - $K_S \rightarrow \pi\pi$, K_L

FCCEe detector

- Ideal for flavor identification [hence: measure Higgs couplings]
 - Impact parameter resolution
 - Low material budget tracker (minimise multiple scattering)
 - Small beam-pipe 1.5 cm -- investigating 1 cm
 - PID capabilities
 - dE/dx (Si tracker) -- Cluster counting dN/dx (Drift)
 - Time of flight -- timing layer



- Detector response based on Delphes:

- Including FastTrackCovariance
- Computes:
 - full track covariance matrix
 - Including multiple scattering
 - smeared track using the off-diagonal terms
 - path length and dN/dx for various gas mixes
- Allows fast turn-around when trying different detector options

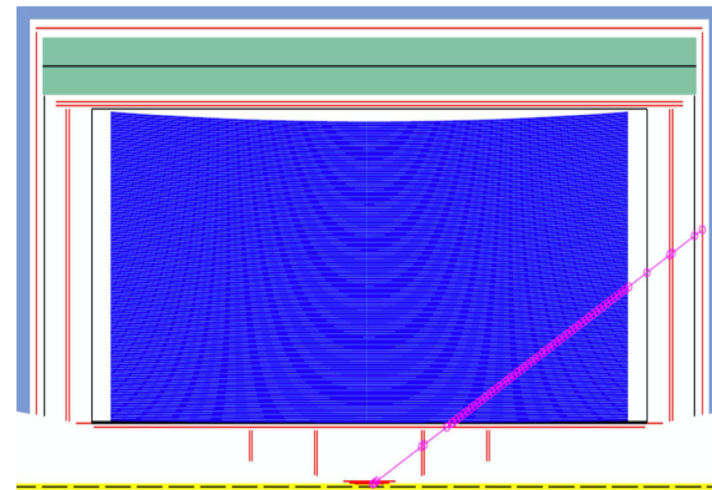
- MC Samples:

- MG5+Pythia8 used to generate:
 - $ee \rightarrow ZH \rightarrow \nu\nu XX$ events (X: g, ud, s, c, b)

- Jets clusters with the generalized-kT algorithm using $p=-1$

- Similar to the anti- k_T algorithm [IRC safe]

IDEA

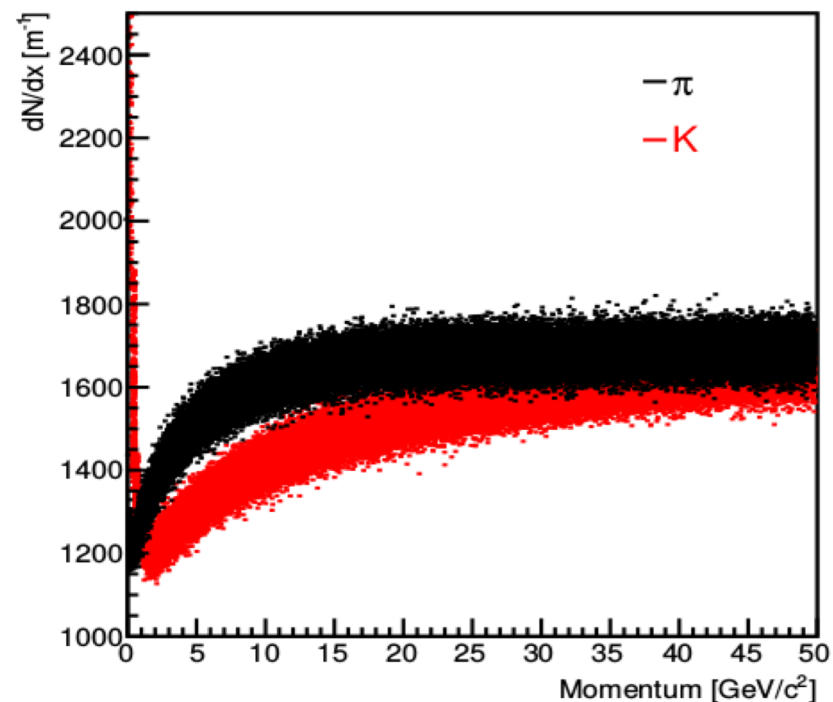


Cluster counting dN/dx

- Count number of **primary ionisation** clusters along track path
- Avoids large landau flukes (**poisson distributed**)
- Requires high granularity
- Module added in Delphes

IDEA detector:

90% He / 10 % Isobutane



```
#####
# Cluster Counting
#####

module ClusterCounting ClusterCounting {

  add InputArray TrackSmearing/tracks
  set OutputArray tracks

  set Bz $B

  ## check that these are consistent with DCHCANI/DCHNANO parameters in TrackCovariance module
  set Rmin $DCHRMIN
  set Rmax $DCHRMAX
  set Zmin $DCHZMIN
  set Zmax $DCHZMAX

  # gas mix option:
  # 0: Helium 90% - Isobutane 10%
  # 1: Helium 100%
  # 2: Argon 50% - Ethane 50%
  # 3: Argon 100%

  set GasOption 0

}
```



Time-of-flight

- Allows for good K/pi separation at low momenta:

$$t_{\text{flight}} \equiv t_F - t_V = \frac{L}{\beta} = \frac{L\sqrt{p^2 + m^2}}{p}$$

- Need to make assumption on vertex time (crucial for highly displaced K_S): A.U.

```
#####
# Time Of Flight Measurement
#####

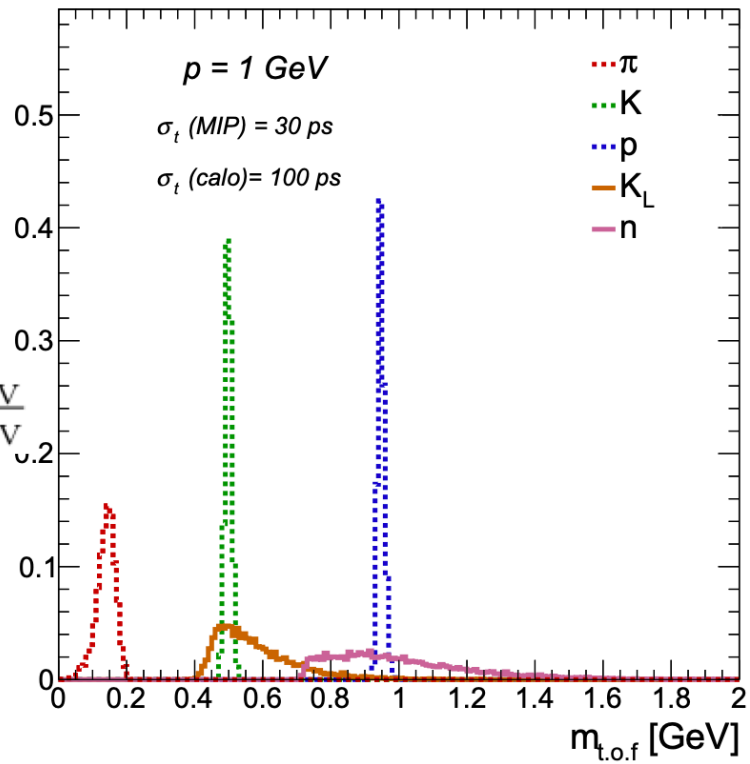
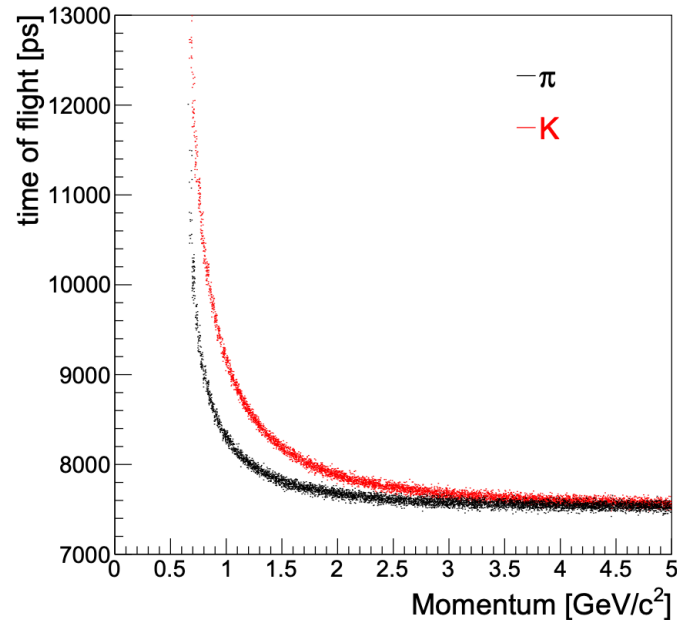
module TimeOfFlight TimeOfFlight {
  set TrackInputArray TimeSmearing/tracks
  set VertexInputArray TruthVertexFinder/vertices

  set OutputArray tracks

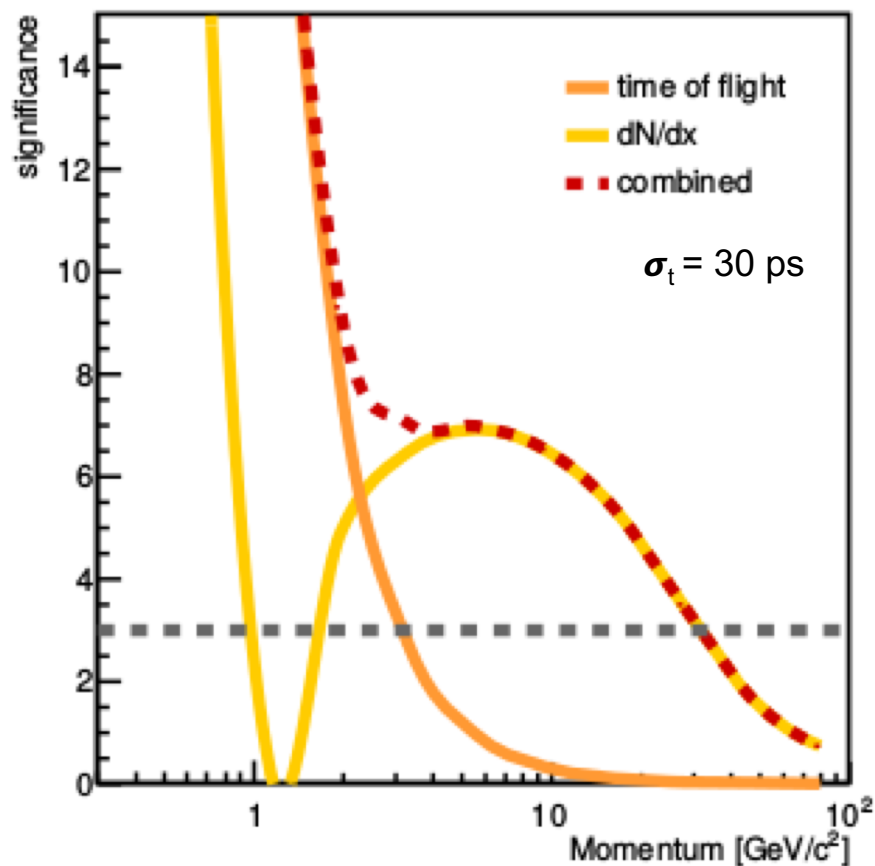
  # 0: assume vertex time tV from MC Truth (ideal case)
  # 1: assume vertex time tV = 0
  # 2: calculate vertex time as vertex TOF, assuming tPV=0

  set VertexTimeMode 2
}
```

$$t_V = \frac{r_V}{\beta_V}$$



Combined PID

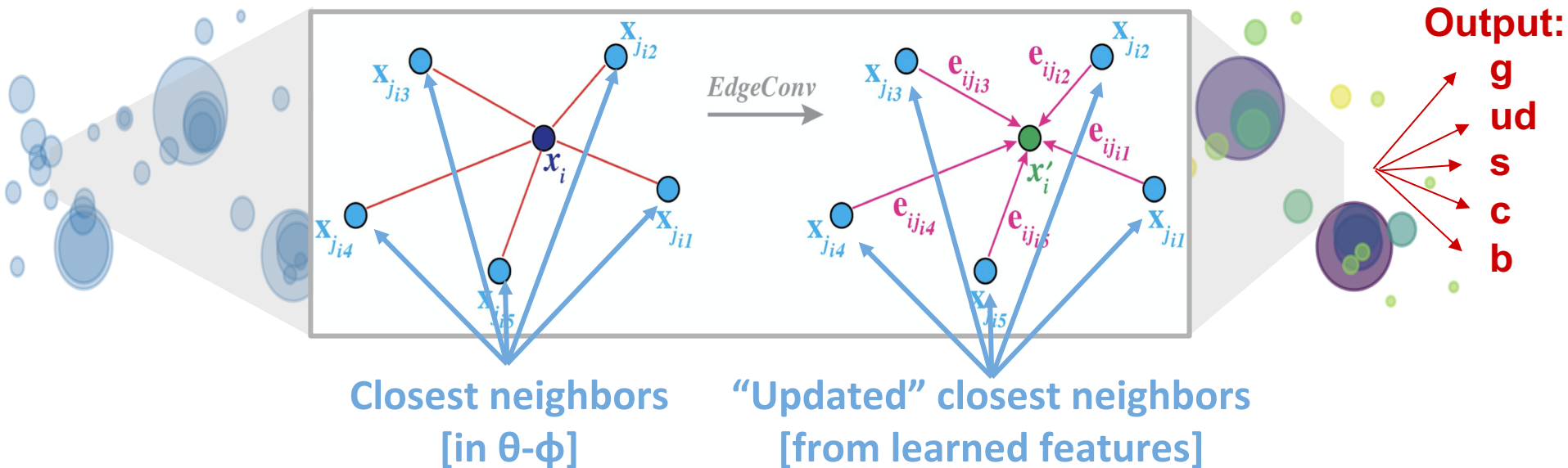


3 std deviation K/pi separation for tracks with $p < 30$ GeV

- Flavour tagging algorithm based on ParticleNet
 - Jet is represented as a “particle cloud”
- Follow a hierarchical learning approach:
 - **First:** Learn “local” structures; **Then:** move to more “global” features
 - Treat the particle cloud as a graph
 - **Particles** are the **vertices** of the graph
 - Relationships** between the particles are the **edges** of the graph

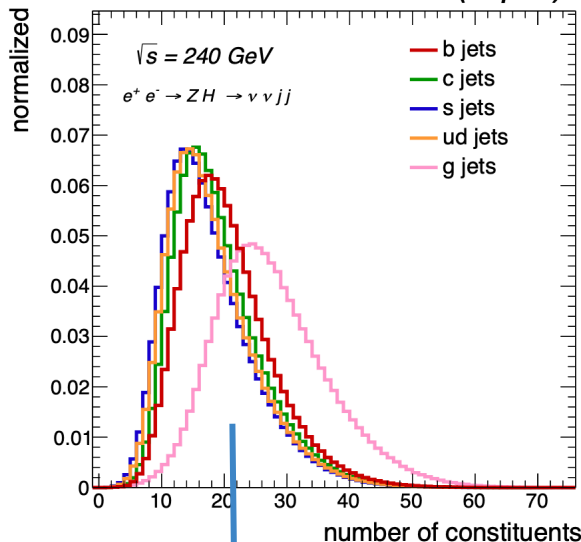
Jet:
As particle cloud

Identify “neighboring” particles



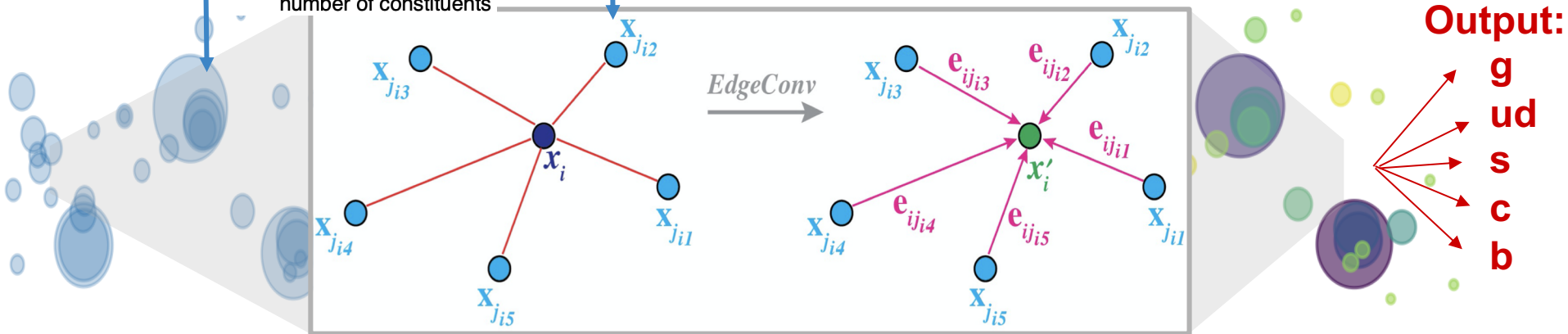
Flavour tagging using ParticleNet (II)

FCC-ee simulation (Delphes)



Particle features:

Particle kinematics, particle charge, Impact parameter (d_0 , d_z) and $O(20)$ /particle significance, particle type (el, mu, γ ,...)



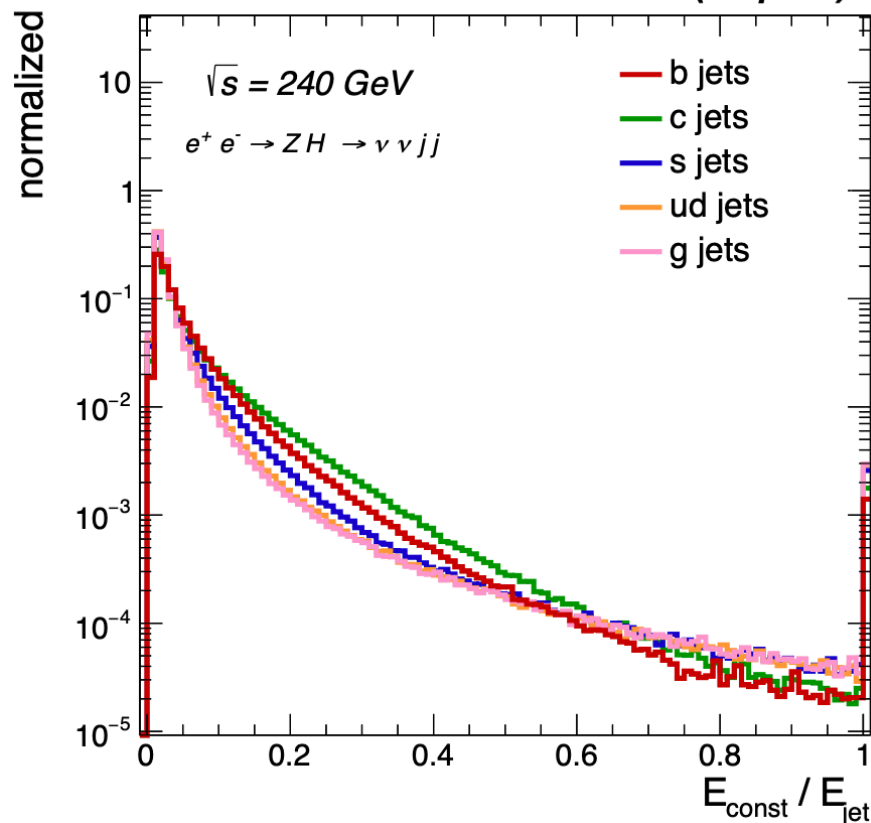
Inputs:
75 particles/jet

Input variables

- Comparison of input distributions for different jet flavors

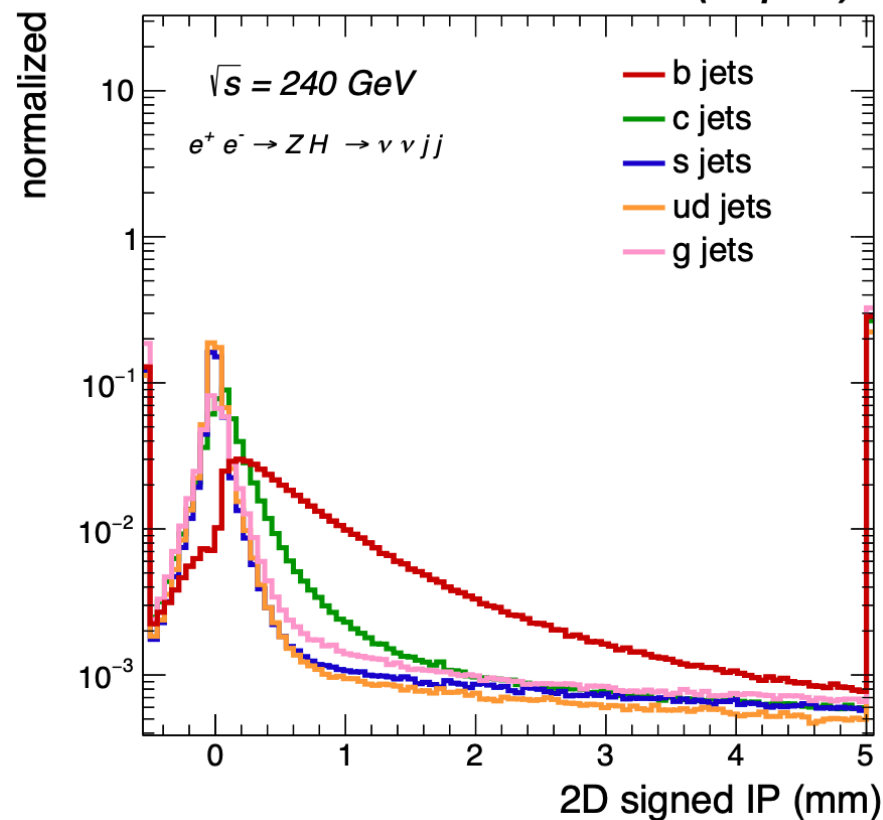
Constituent relative energy

FCC-ee simulation (Delphes)



Impact parameter (d_0)

FCC-ee simulation (Delphes)

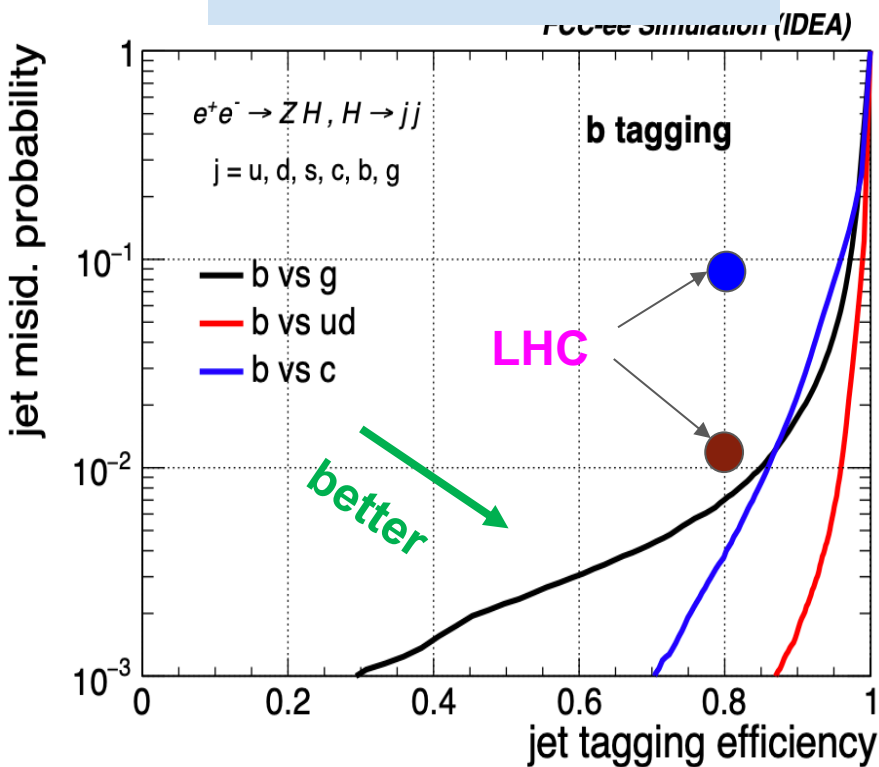


- More comparisons:

<https://selvaggi.web.cern.ch/selvaggi/FCC/FCCEe/FlavourTagging/>

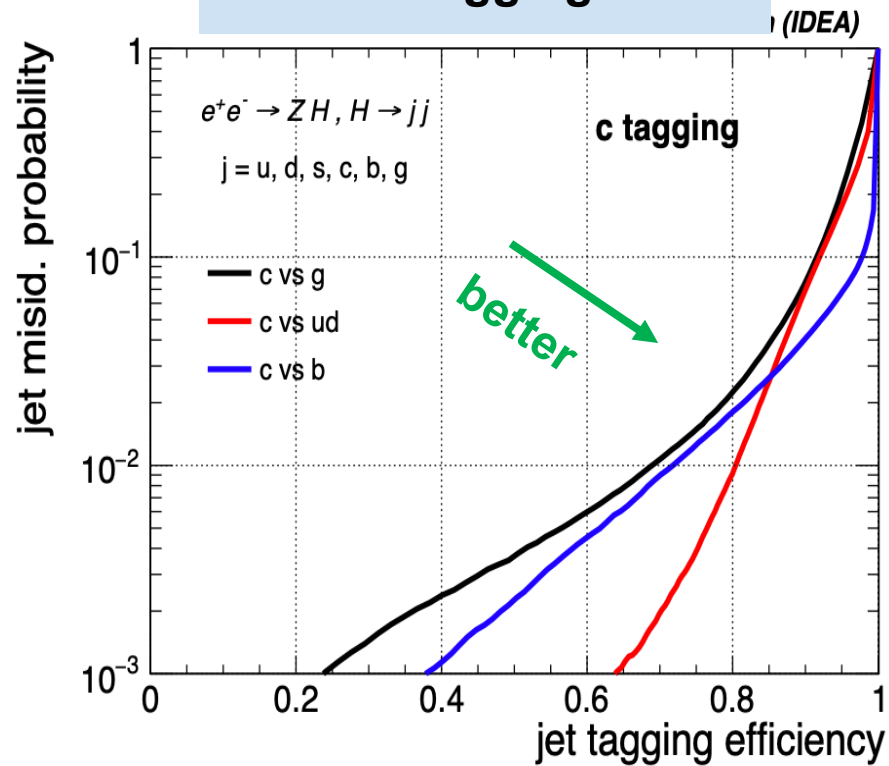
Performance (b/c)

b-tagging



WP	Eff (b)	Mistag (g)	Mistag (ud)	Mistag (c)
Loose	90%	2%	0.1%	2%
Medium	80%	0.7%	<0.1%	0.3%

c-tagging

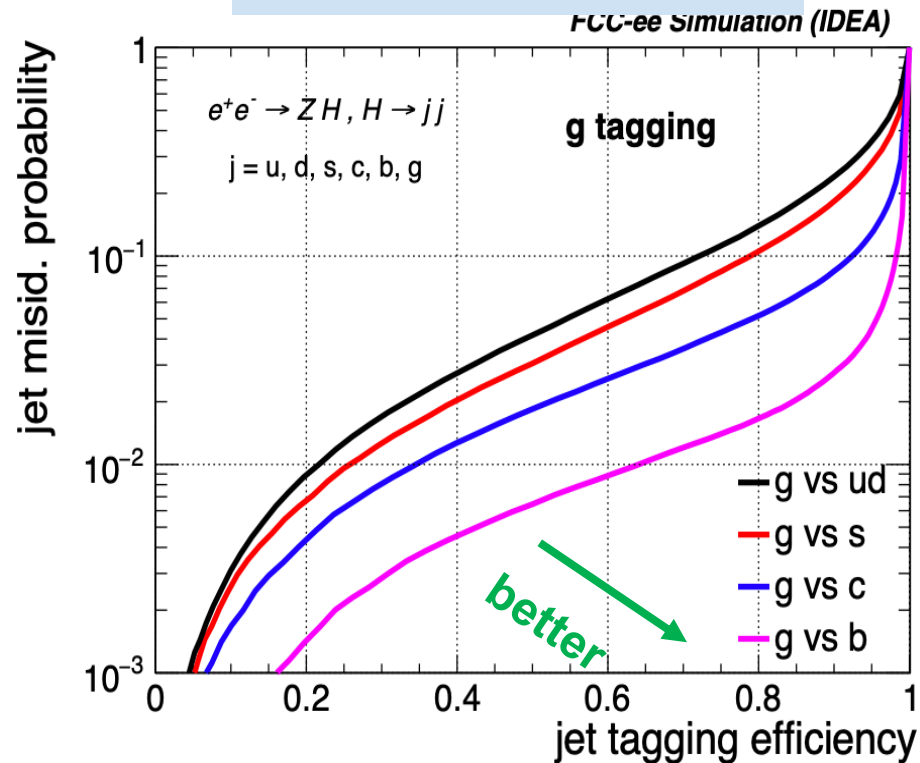
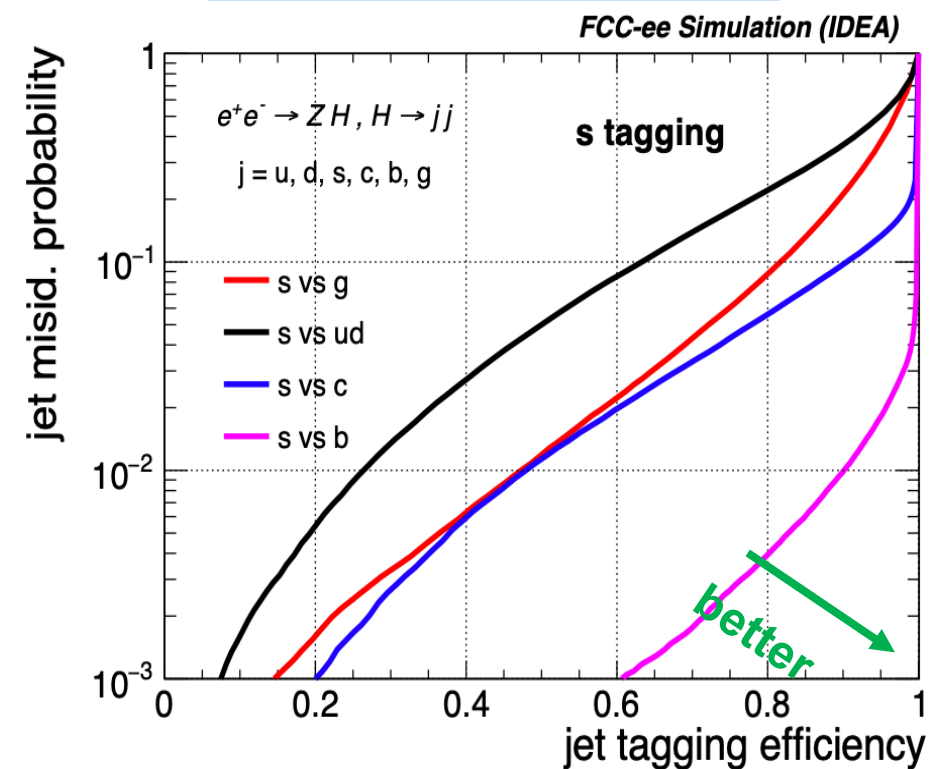


WP	Eff (c)	Mistag (g)	Mistag (ud)	Mistag (b)
Loose	90%	7%	7%	4%
Medium	80%	2%	0.8%	2%

Performance (strange/gluon)

strange-tagging

gluon -tagging

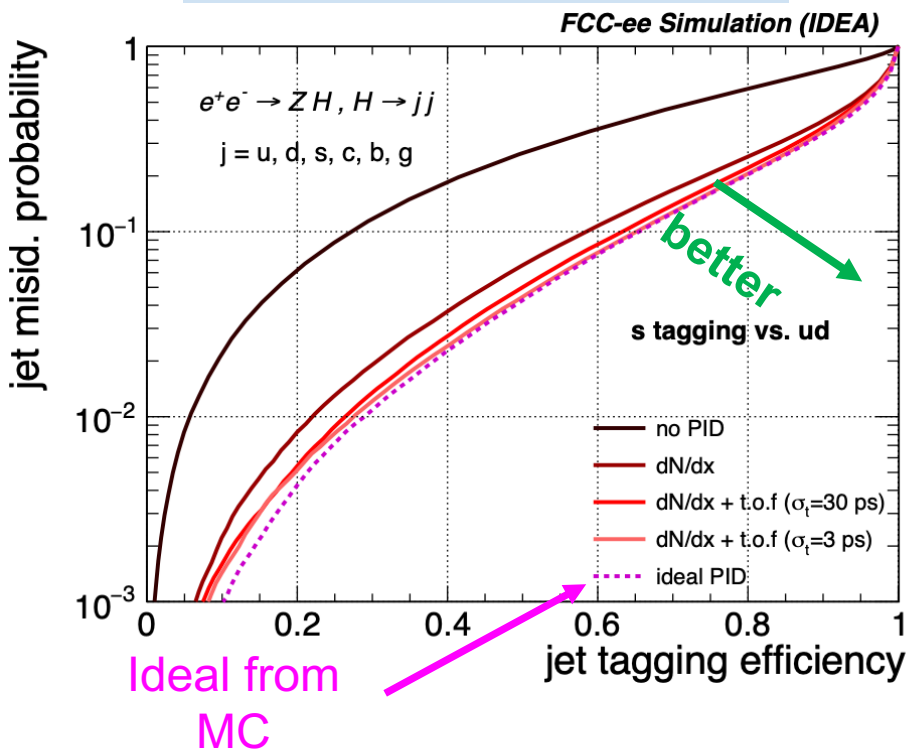


WP	Eff (s)	Mistag (g)	Mistag (ud)	Mistag (c)	Mistag (b)
Loose	90%	20%	40%	10%	1%
Medium	80%	9%	20%	6%	0.4%

WP	Eff (g)	Mistag (ud)	Mistag (c)	Mistag (b)
Loose	90%	25%	7%	2.5%
Medium	80%	15%	5%	2%

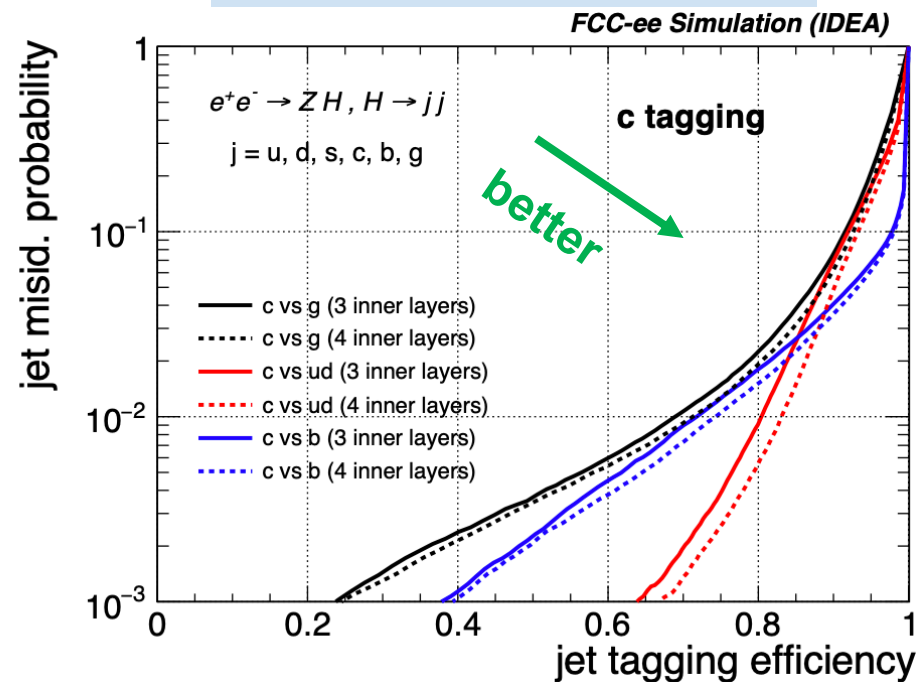
Impact of detector configurations

Strange tagging [PID]



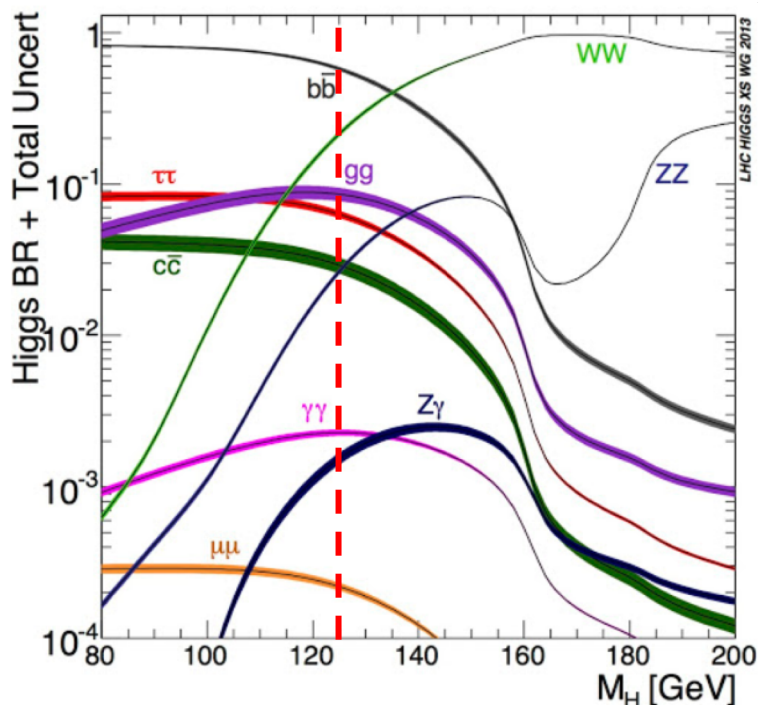
- dN/dX brings most of the gain; additional gain w/ TOF (30ps)
 - TOF (3ps) brings marginal improvement
 - dN/dX+TOF(30ps): very close to perfect PID

c-tagging [PIX layers]



- Improvement up to 2x for charm tagging
 - marginal/no improvement in b-tagging

Higgs couplings: $H \rightarrow cc$



\sqrt{s} (GeV)	240		365	
Luminosity (ab^{-1})	5		1.5	
$\delta(\sigma\text{BR})/\sigma\text{BR}$ (%)	HZ	$\nu\bar{\nu}$ H	HZ	$\nu\bar{\nu}$ H
$H \rightarrow \text{any}$	± 0.5		± 0.9	
$H \rightarrow b\bar{b}$	± 0.3	± 3.1	± 0.5	± 0.9
$H \rightarrow c\bar{c}$	± 2.2		± 6.5	± 10
$H \rightarrow gg$	± 1.9		± 3.5	± 4.5
$H \rightarrow W^+W^-$	± 1.2		± 2.6	± 3.0
$H \rightarrow ZZ$	± 4.4		± 12	± 10
$H \rightarrow \tau\tau$	± 0.9		± 1.8	± 8
$H \rightarrow \gamma\gamma$	± 9.0		± 18	± 22
$H \rightarrow \mu^+\mu^-$	± 19		± 40	
$H \rightarrow \text{invis.}$	< 0.3		< 0.6	

Ref: P. Janot talk at the CDR Symposium; March 2019

FCce: $\sigma_{ZH} \sim 200\text{fb}$, $L \sim 5 \text{ ab}^{-1}$ (2 IP): $\sim 1\text{M ZH}$
 [600k $H \rightarrow b\bar{b}$, 100k $H \rightarrow gg$, 30k $H \rightarrow c\bar{c}$]

Use Loose WP:

[c-tag: 90%, b-mistag: 4%, g-mistag: 7%]

- **Scenario 1:** $Z(\rightarrow \nu\nu)H$

$\delta(\sigma_{\text{xBR}})/\sigma_{\text{xBR}}$ (%) ~ 1.5 [no systematics]

- **Scenario 2:** $Z(\rightarrow \text{all})H$

$\delta(\sigma_{\text{xBR}})/\sigma_{\text{xBR}}$ (%) ~ 0.7 [no systematics]

- **Stat limit [i.e. no BKG]:**

$\delta(\sigma_{\text{xBR}})/\sigma_{\text{xBR}}$ (%) $\sim 0.6\%$

- **No BKG rejection:**

$\delta(\sigma_{\text{xBR}})/\sigma_{\text{xBR}}$ (%) $\sim 2.9\%$

Results look promising



Higgs couplings: $H \rightarrow ss$

$$\text{BR}(H \rightarrow ss) = \text{BR}(H \rightarrow cc) (m_s/m_c)^2 \sim 2.3 \cdot 10^{-4}$$

FCCEe: $\sigma_{ZH} \sim 200 \text{ fb}$, $L \sim 5 \text{ ab}^{-1}$ (2 IP): **$\sim 1 \text{ M ZH}$**

[600k $H \rightarrow bb$, 100k $H \rightarrow gg$, 30k $H \rightarrow cc$, **200 $H \rightarrow ss$**]

Use Tight WP:

[s-tag: 60%, **g-mistag**, **c-mistag** and **b-mist**: negligible]

- The most challenging BKG is ZZ

with one Z off-shell $\sim 125 \text{ GeV}$ [$\sim 10\%$ of the Higgs signal]

- **Optimistic assumption:**

- 100% of the Higgs events (i.e. the 1M events above) are reconstructed
- 100k ZZ events; (BR for $Z \rightarrow s\bar{s}$) $\sim 15\%$
- 15k ZZ events. After applying the Tight WP of the tagger:
- 5.4k events $\rightarrow 88/\sqrt{5400} = 1.2\sigma$

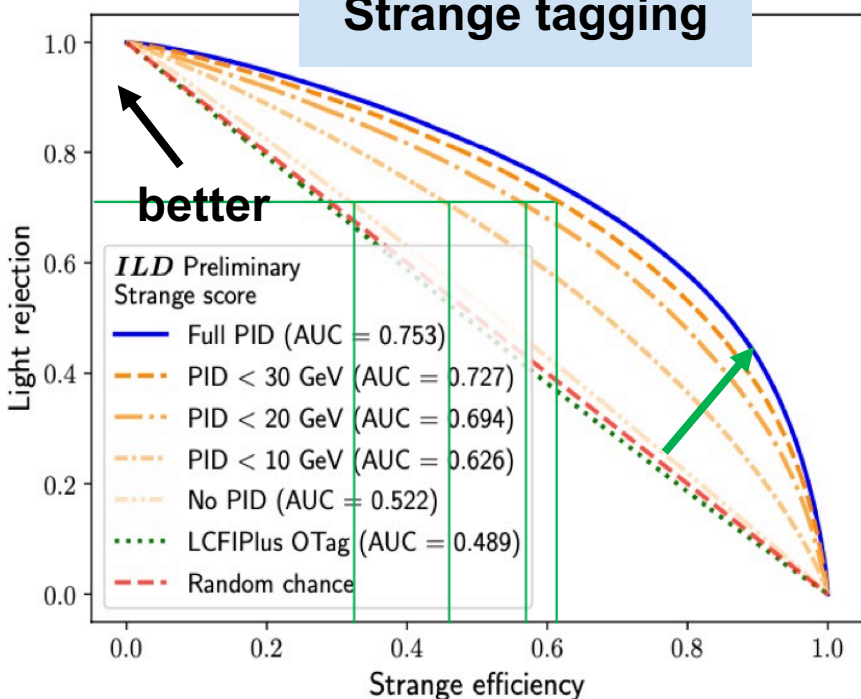
Rough numbers
from FCC
Workshop 2022
[\[slides\]](#)

Back-of-the envelope estimates

THOROUGH STUDIES NEEDED

- Jet tagging using Recurrent Neural Net (RNN)
 - **Inputs:** jet-level info + particle-level info [10 highest- p_T particles]
 - **Multiclass output:** b, c, s, ud, gluon
 - Designed for **ILD**; uses **FullSim**

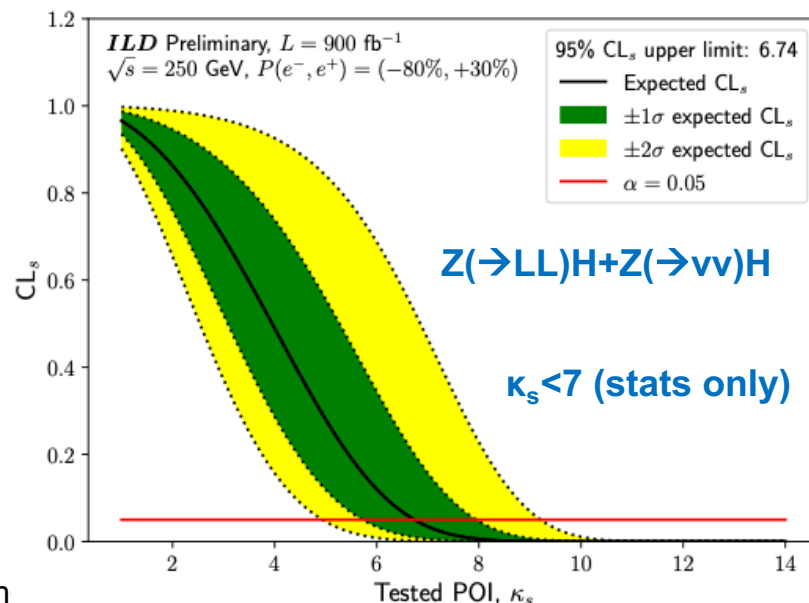
Strange tagging



- Significant improvement compared to previous ILD tagging algorithm (LCFIPlus)
 - PID capabilities up to $p_T \sim 30$ GeV important

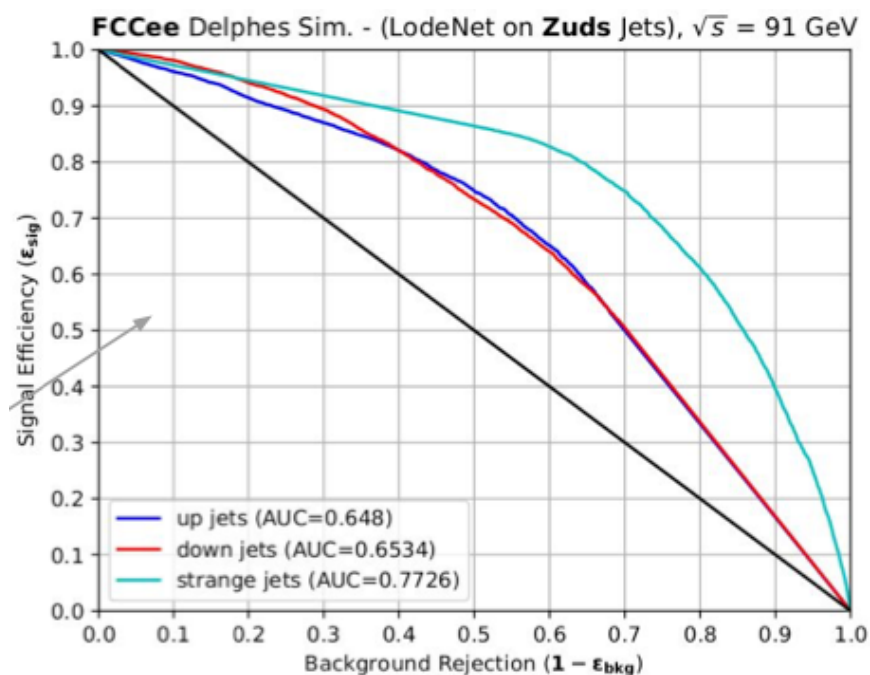
Application on $H \rightarrow ss$

- ILC @ 250 GeV; 900 fb^{-1}
- Signal: $Z(\rightarrow \nu\nu)H$ and $Z(\rightarrow LL)H$
- Analysis design: selection on evt-level vars
- Signal extraction: fit strange-tagging discriminant



- Jet tagging using CNN-2D [focusing on strange-tagging]
 - **Inputs:** jet images; several channels: K , π , e , μ , γ ...
 - **Multiclass output:** s , u , d
 - **IDEA** detector; use **FastSim**
 - **Bonus:** Improved Jet flavor assignment

Strange tagging



Impact of PID

Signal Efficiency	10% fake rate	5% fake rate	1% fake rate
Generator	47.2%	27.7%	7.5%
PF only	17.7%	9.7%	2.0%
PF + Ks	21.9%	12.9%	4.4%
PF + Ks + K [±]	39.5%	24.8%	7.0%

- Up to 2x improved $\epsilon(\text{SIG})$ with K_S / K^{\pm} separation

Summary & outlook

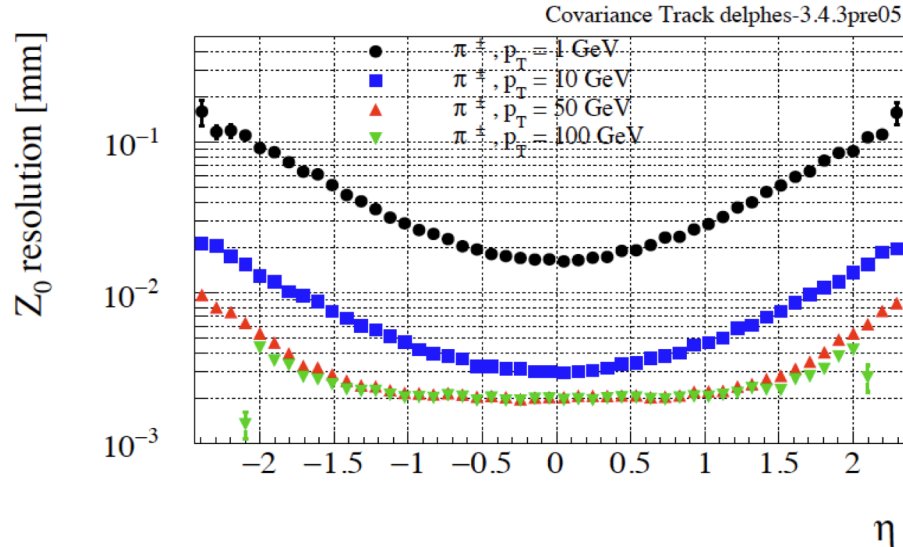
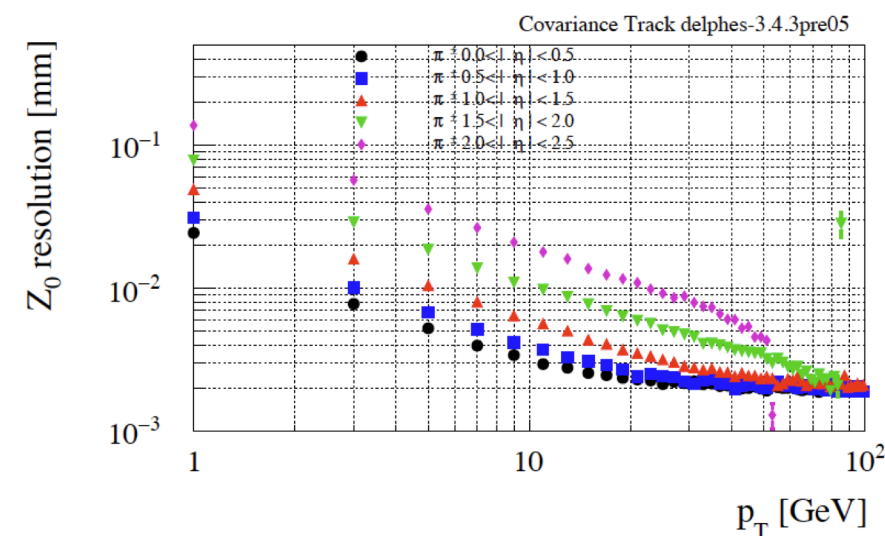
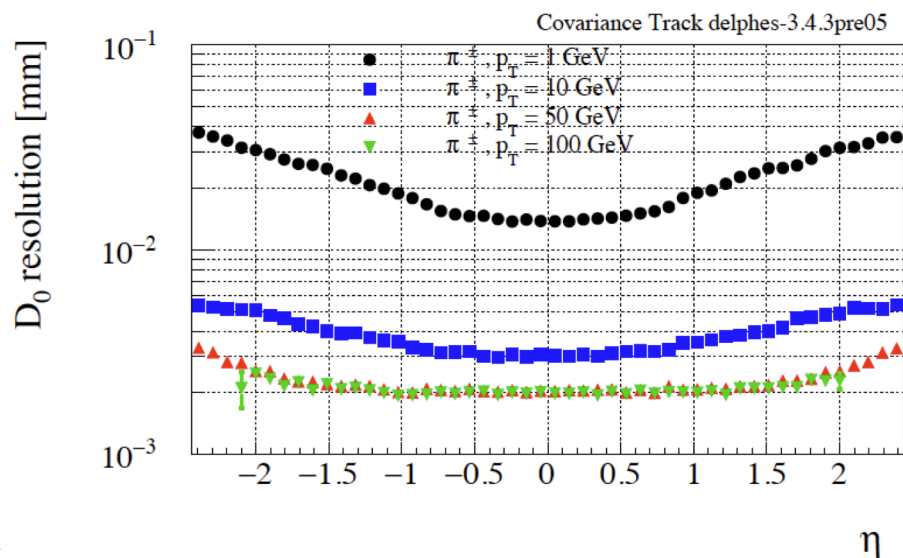
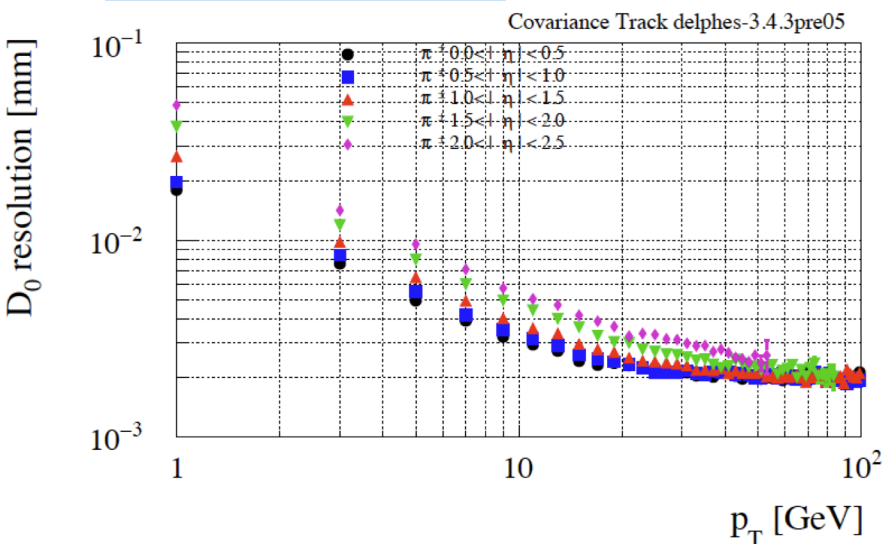
- Powerful jet flavour identification essential for the success of the e^+e^- physics program
- A first version of a jet identification algorithm based on **PF candidates** and **PID** and **advanced ML** in place
 - Multi-class classifier $b/c/s/ud/g$; Results extremely promising
 - More details: [arXiv:2202.03285](https://arxiv.org/abs/2202.03285)
- Conclusions:
 - adding an additional vertex layer does not tremendously improve b-tagging performance (resolution of $\sim 2\mu\text{m}$ already outstanding)
 - but improves charm tagging
 - Some room for improved in strange tagging with more powerful PID
- Next steps/work in progress
 - Implementation in FCCSW
 - Test performance using FullSim
 - Application on physics analyses [e.g., $H \rightarrow cc$, $H \rightarrow ss$]

Backup

Impact parameter performance

Credits to Sylvie Braibant

IDEA detector:



2 μ m IP resolution at high- p_T

Designing a jet flavour tagging algorithm

A point cloud



Source: <https://news.voyage.auto/an-introduction-to-lidar-the-key-self-driving-car-sensor-a7e405590cff>

- Point cloud (Wikipedia):
 - A set of data **points** in space
 - Produced by 3D scanners, which measure a large number of points on the external surfaces of objects around them

From point clouds to particle clouds

A point cloud

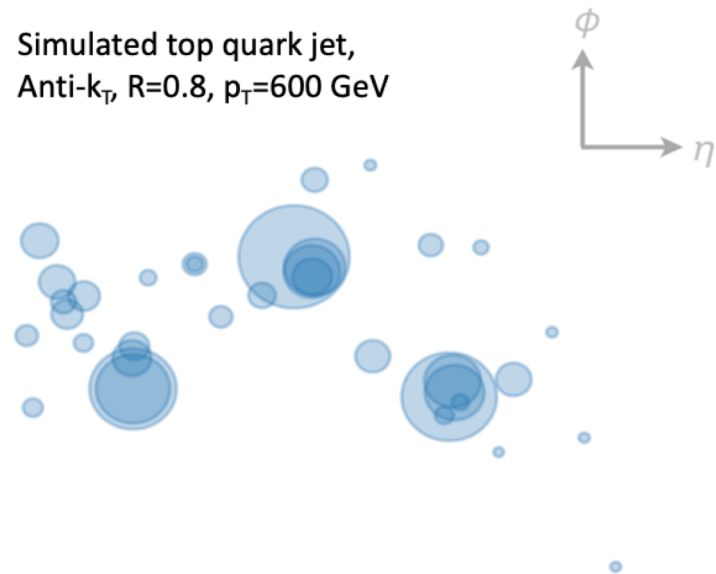


Source: <https://news.voyage.auto/an-introduction-to-lidar-the-key-self-driving-car-sensor-a7e405590cff>

- Point cloud (Wikipedia):
 - A set of data **points** in space
 - Produced by 3D scanners, which measure a large number of points on the external surfaces of objects around them

A particle cloud

Simulated top quark jet,
Anti- k_T , $R=0.8$, $p_T=600$ GeV



- Particle cloud :
 - A set of **particles** in space
 - Produced by clustering a large number of particles measured by the detectors



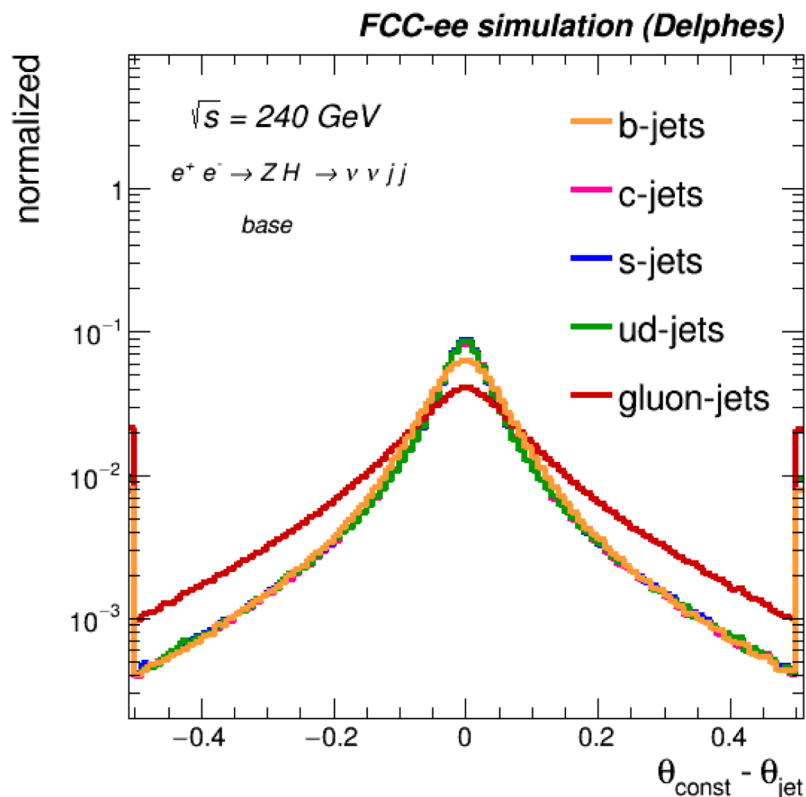
Full list of input variables

Variable	Description
Kinematics	
$E_{\text{const}}/E_{\text{jet}}$	energy of the jet constituent divided by the jet energy
θ_{rel}	polar angle of the constituent with respect to the jet momentum
ϕ_{rel}	azimuthal angle of the constituent with respect to the jet momentum
Displacement	
d_{xy}	transverse impact parameter of the track
d_z	longitudinal impact parameter of the track
$\text{SIP}_{2\text{D}}$	signed 2D impact parameter of the track
$\text{SIP}_{2\text{D}}/\sigma_{2\text{D}}$	signed 2D impact parameter significance of the track
$\text{SIP}_{3\text{D}}$	signed 3D impact parameter of the track
$\text{SIP}_{3\text{D}}/\sigma_{3\text{D}}$	signed 3D impact parameter significance of the track
$d_{3\text{D}}$	jet track distance at their point of closest approach
$d_{3\text{D}}/\sigma_{d_{3\text{D}}}$	jet track distance significance at their point of closest approach
C_{ij}	covariance matrix of the track parameters
Identification	
q	electric charge of the particle
$m_{\text{t.o.f.}}$	mass calculated from time-of-flight
dN/dx	number of primary ionisation clusters along track
isMuon	if the particle is identified as a muon
isElectron	if the particle is identified as an electron
isPhoton	if the particle is identified as a photon
isChargedHadron	if the particle is identified as a charged hadron
isNeutralHadron	if the particle is identified as a neutral hadron

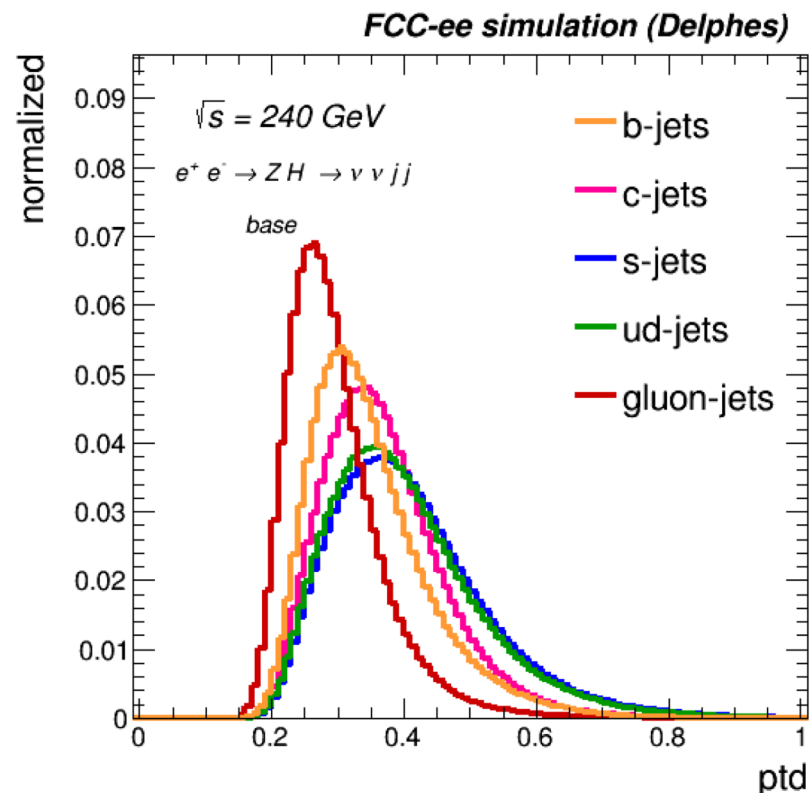
Input variables

- Comparison of input distributions for different jet flavors

Projection || to jet axis



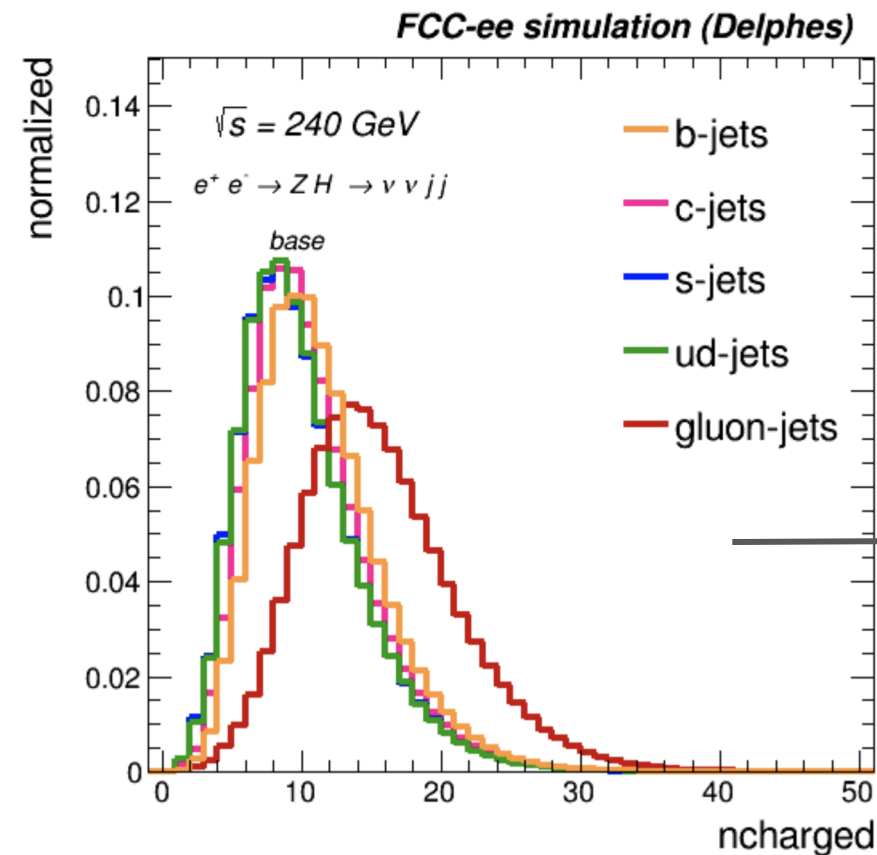
$p_{\text{T}}D$



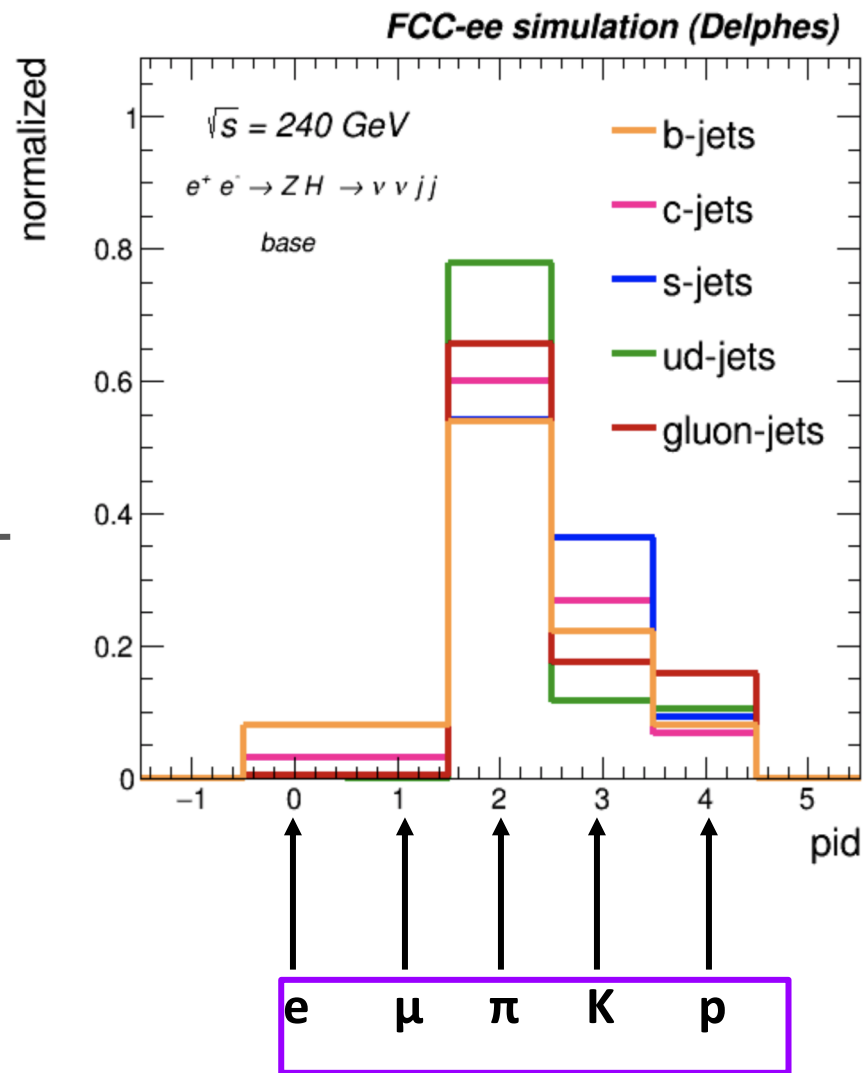
- More comparisons:

<https://selvaggi.web.cern.ch/selvaggi/FCC/FCCee/FlavourTagging/>

Performance w/ PID



no PID: only charge
realistic: $e, \mu, m_{\text{tof}}, dN/dx$
perfect PID: e, μ, π, K, p
 from MC truth

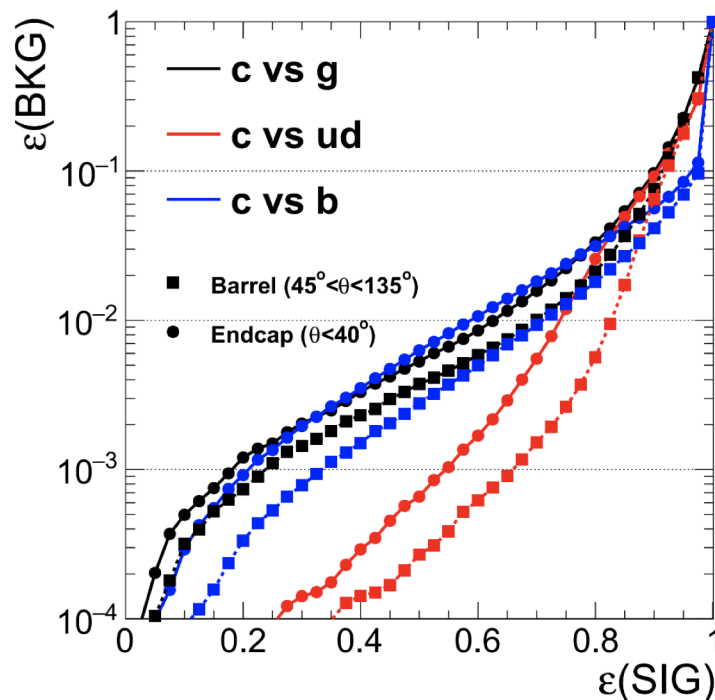
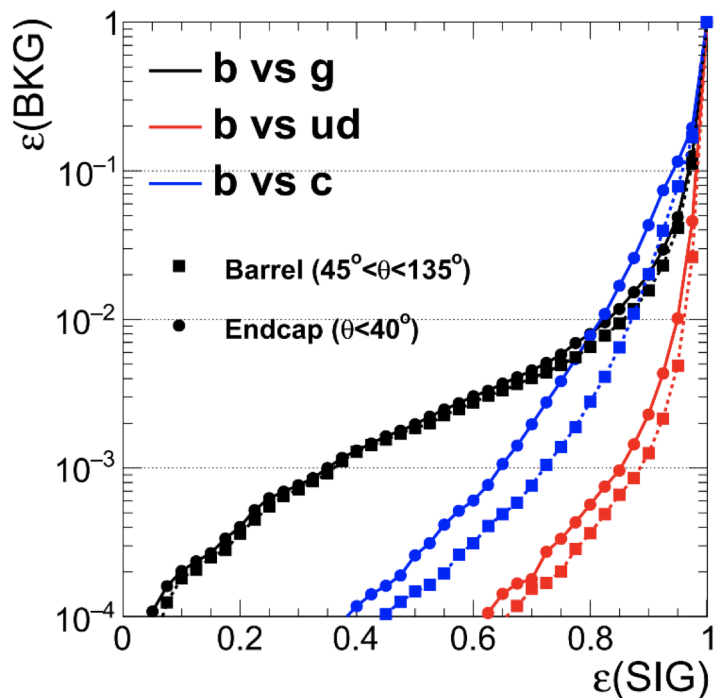


Performance vs theta (b/c)

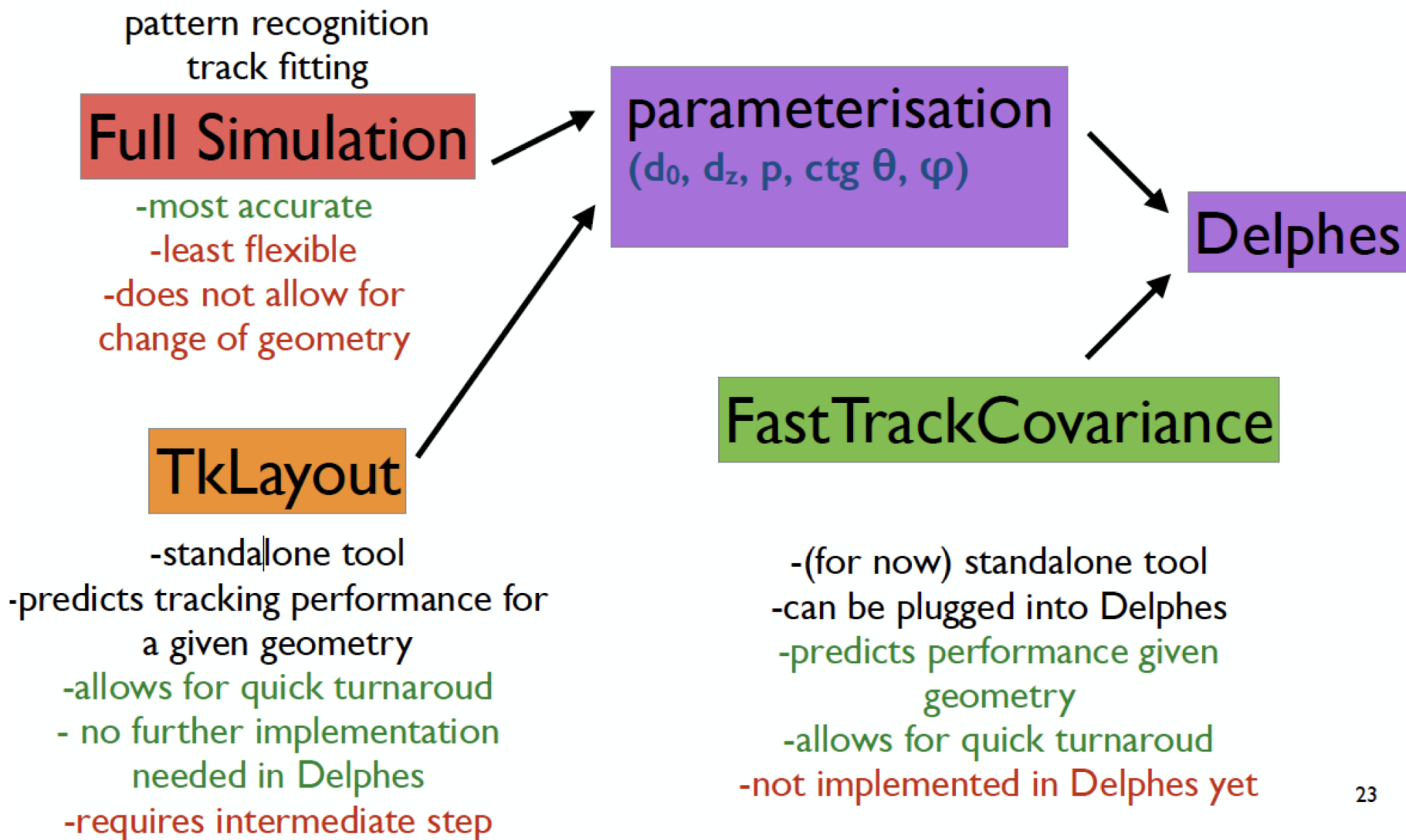
b-tagging

c-tagging

PRELIMINARY !! (LOW STATS TRAINING)



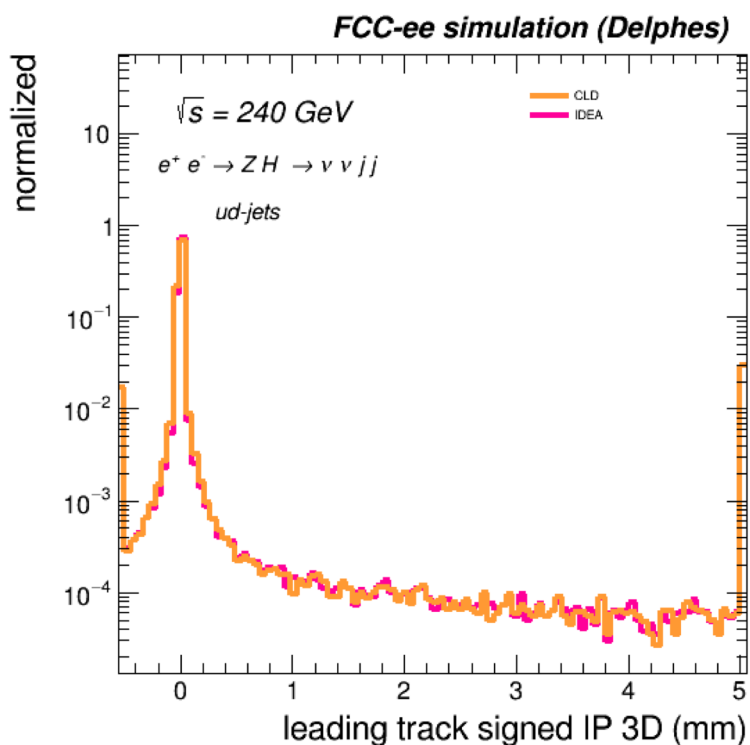
Tracking in Delphes



Comparison: IDEA vs. CLD

- No big differences between in input variables between IDEA & CLD
 - small difference in material budget observed on light jets since $dxy \sim 0$
 - expect slightly better performance for IDEA detector for discrimination vs light

ud-jets



c-jets

