

# ML inference in Future Collider Software

IDT-WG3 topical meeting on running tools for ML inference

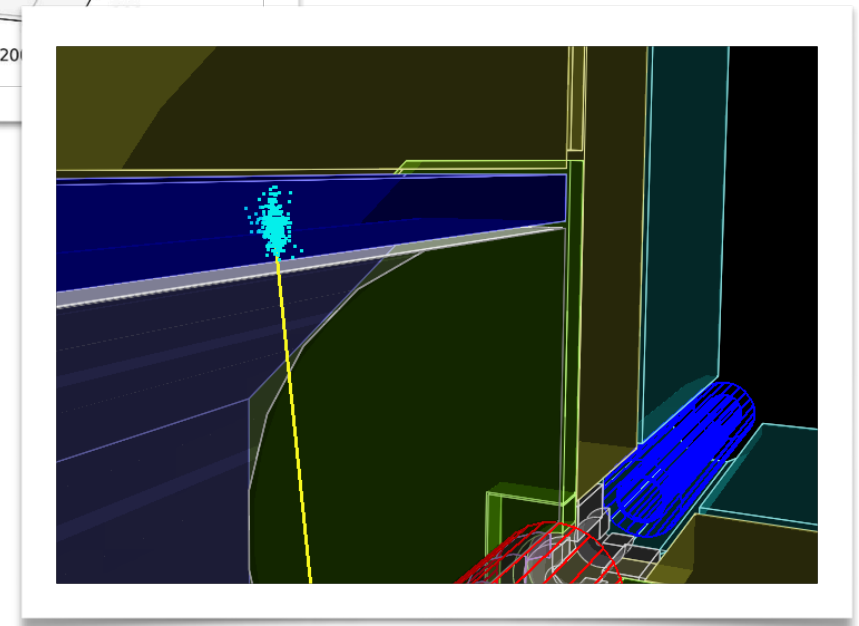
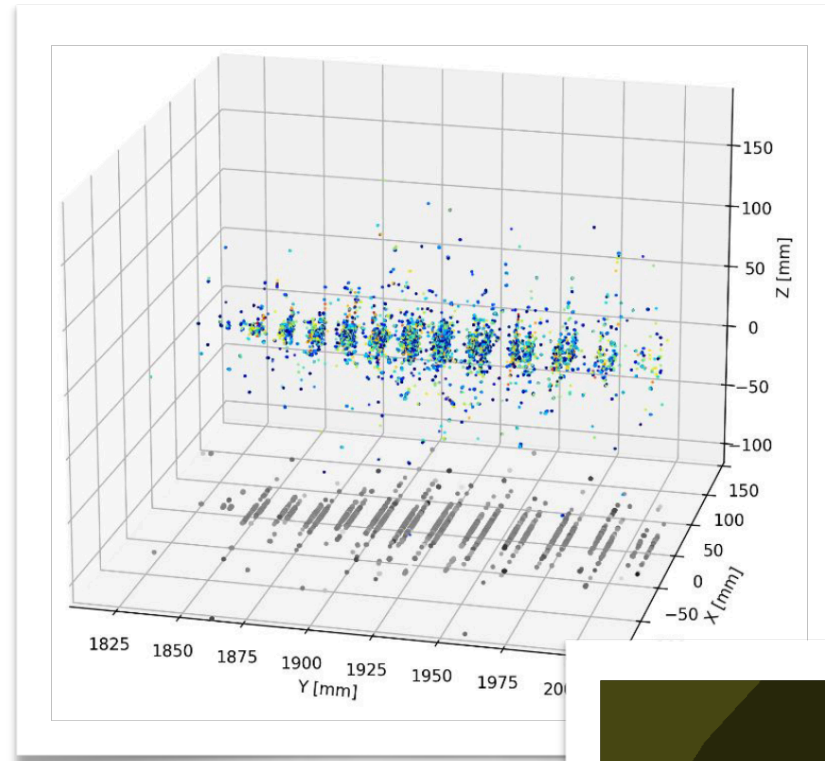
15.05.24

Frank Gaede, DESY



# Outline

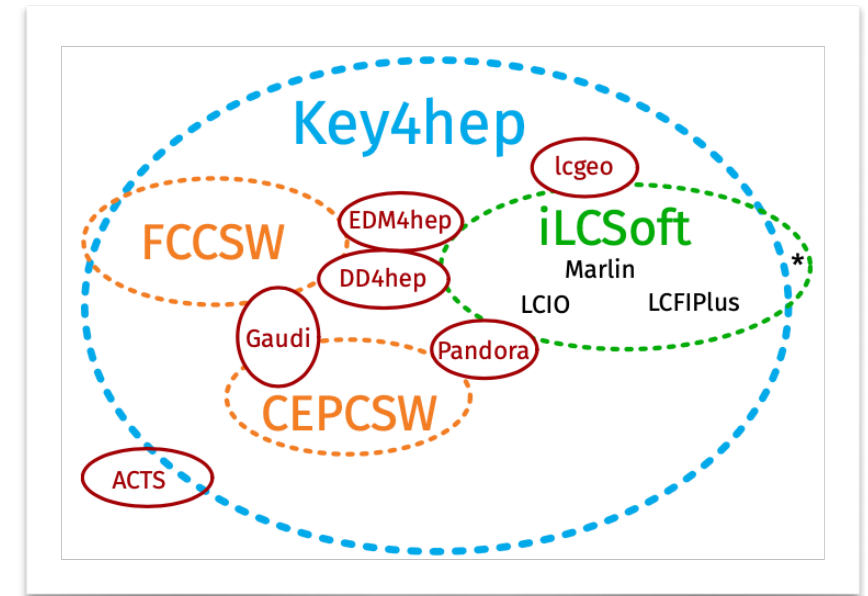
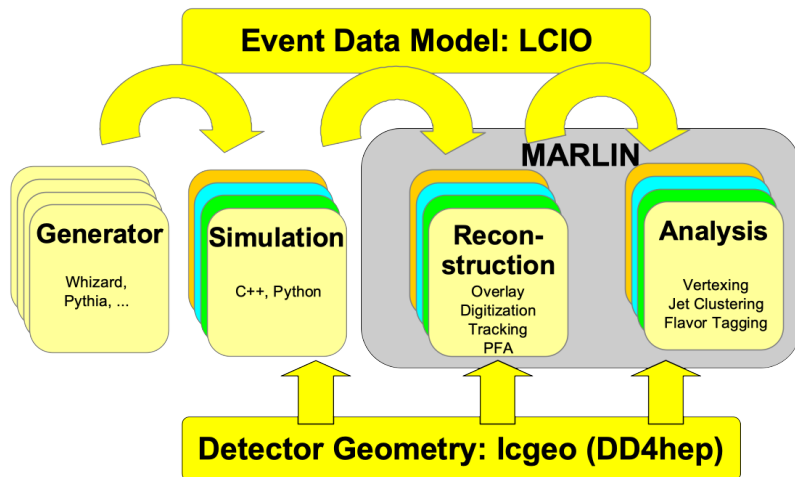
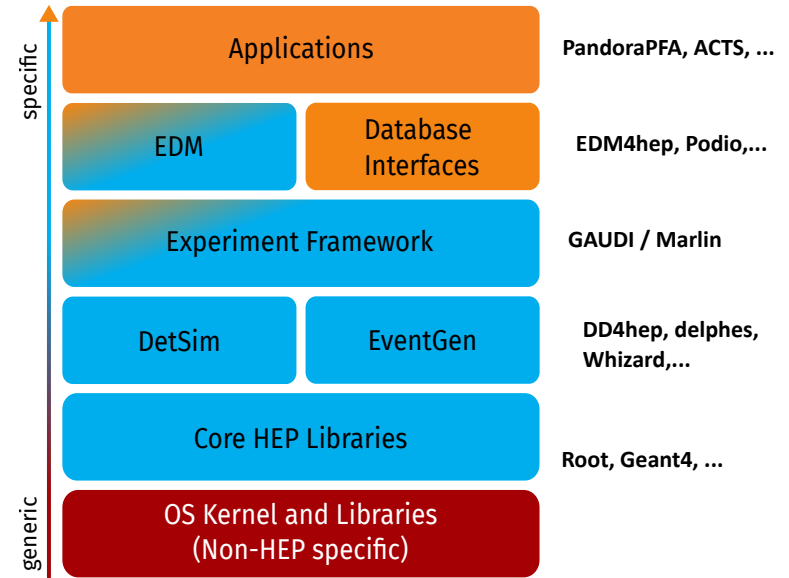
- Introduction
- Generative ML for Simulation
- ML inference for Reconstruction
- Future Developments
- Summary and Outlook



# Introduction

## turnkey software stack for all future colliders

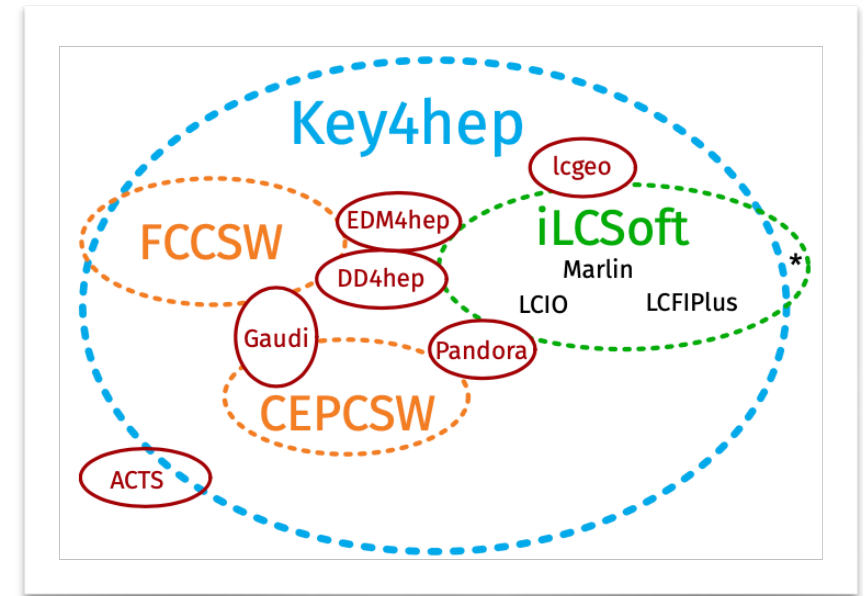
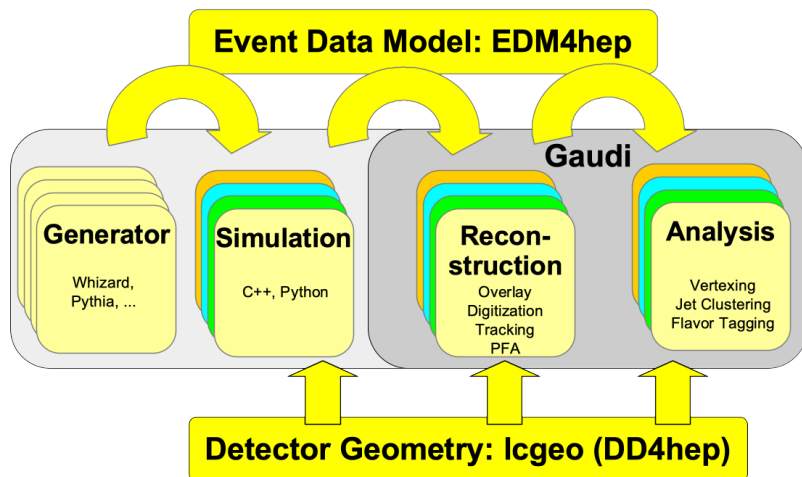
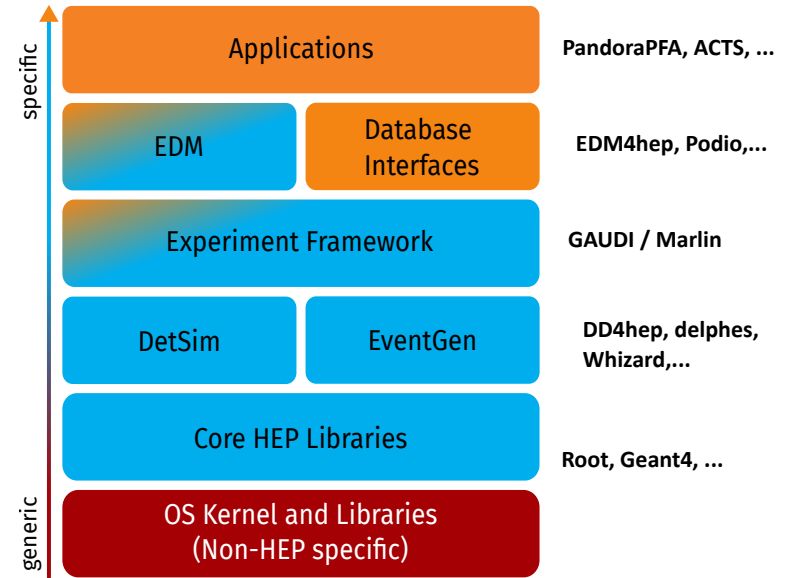
- HEP community decided to develop a **common turnkey software stack** – for future collider studies: **Key4hep**
- supported by **HSF** and **CERN EP-R&D** and *AIDAinnova*
- involved communities/contributors:
  - CEPC, CLIC, FCC, EIC, ILC, LUXE, Muon Collider ...
- **all future collider ML inference tools should be developed in Key4hep**
  - or iLCSoft as transition for linear collider still ongoing



# Introduction

turnkey software stack for all future colliders

- HEP community decided to develop a **common turnkey software stack** – for future collider studies: **Key4hep**
- supported by **HSF** and **CERN EP-R&D** and *AIDAinnova*
- involved communities/contributors:
  - CEPC, CLIC, FCC, EIC, ILC, LUXE, Muon Collider ...
- **all future collider ML inference tools should be developed in Key4hep**
  - or iLCSoft as transition for linear collider still ongoing



# Introduction

## modern (generative) machine learning

- modern machine learning development is almost exclusively done in the **Python** world using the **PyTorch** toolkit and to a much lesser extend TensorFlow/KERAS
  - in HEP traditionally also ROOT-TMVA used ( for more classic approaches)
- mostly using **HDF5** file format for training, test and validation data
- to run such new models in standard HEP frameworks such as Key4hep/iLCSoft some **suitable glue code or middleware** is needed -> this talk
- in particular we need C++ inference tools to call the ML model during sim/rec/ana in framework
  - here we focus on
    - **ONNX** - Open Neural Network Exchange (<https://onnx.ai/>)
    - **libtorch** (<https://pytorch.org/cppdocs/>)

# Existing ML inference tools

in Key4hep - iLCSoft and (briefly) covered today

generative ML in simulation:

- **DDFastShowerML (DDML)**
  - generative ML architectures for fast shower simulation

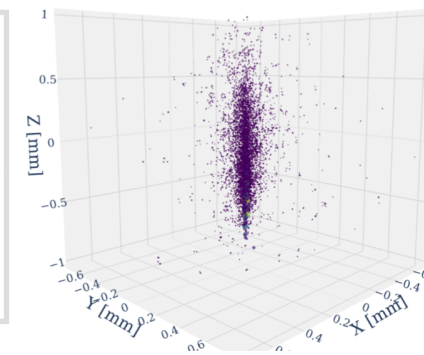
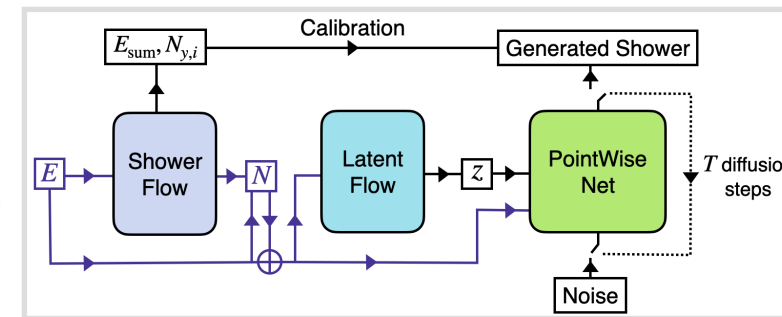
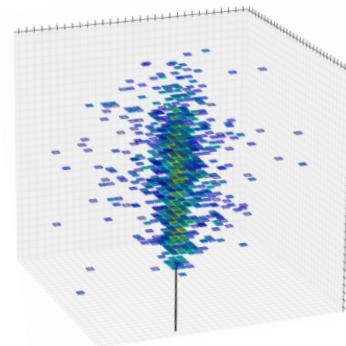
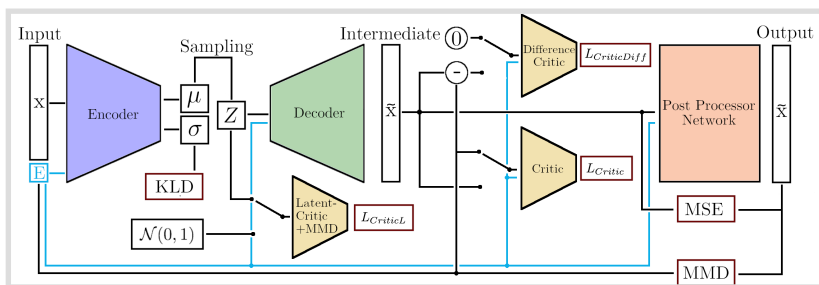
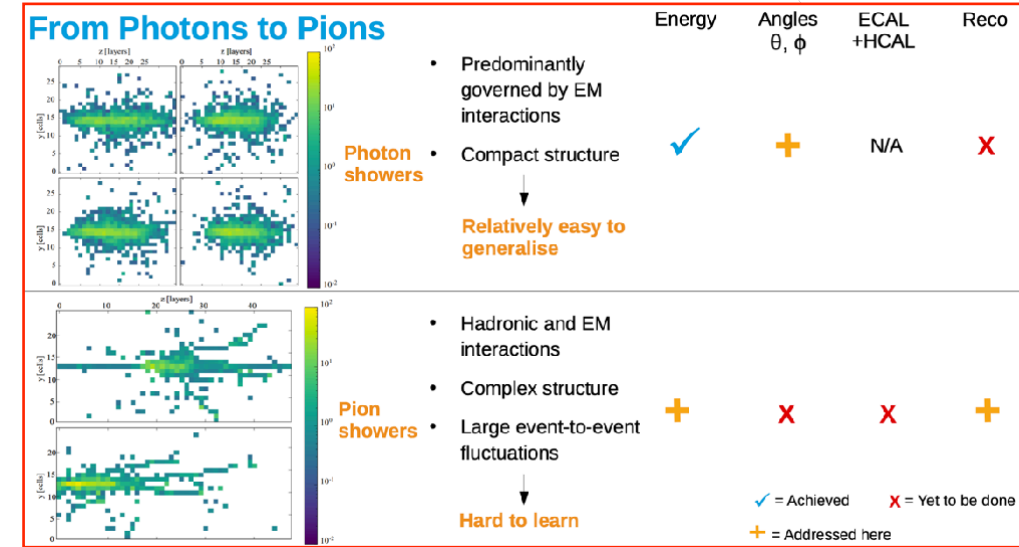
running inference for (high level) reconstruction:

- **CPID**
  - comprehensive particle identification for future colliders
- **MarlinML (FlavourTaggingML)**
  - machine learning flavour tagging for future higgs factories

# Fast shower simulation

## with generative ML

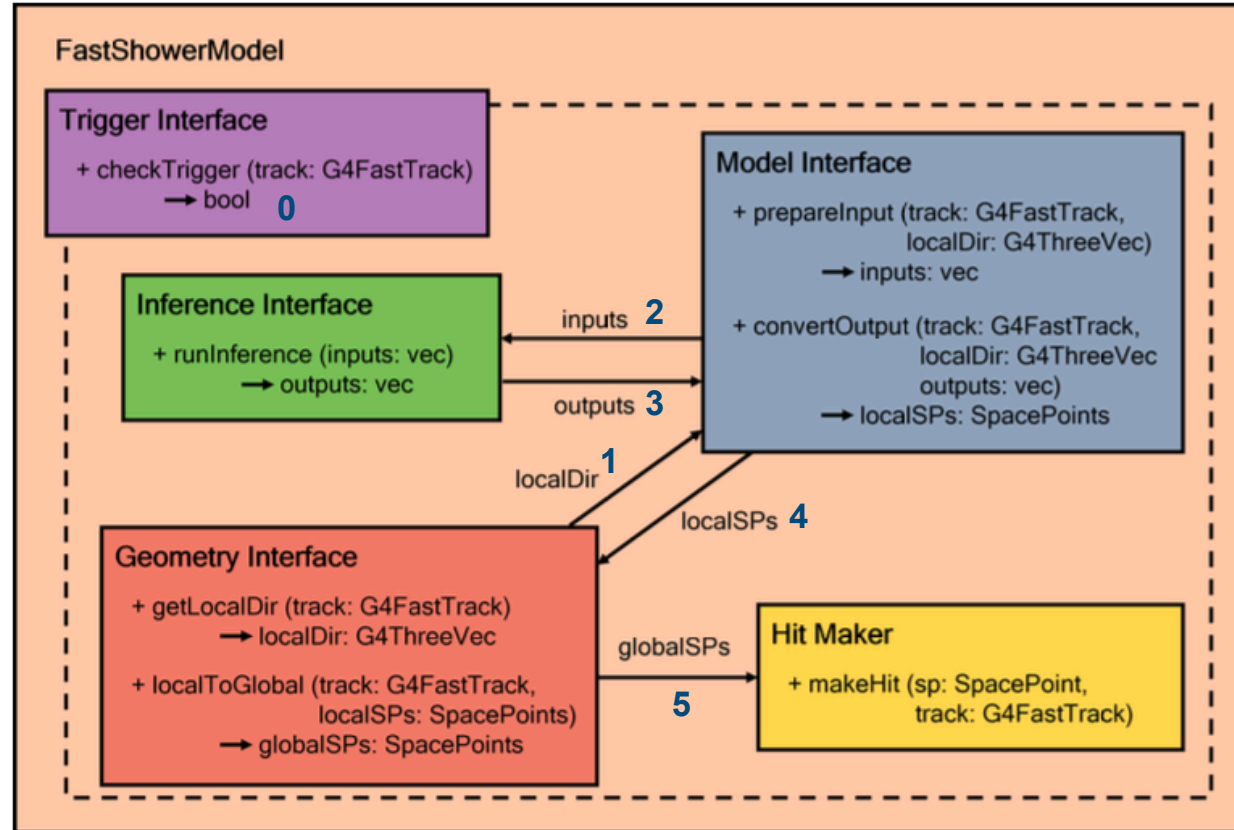
- shower generation with ML offers great potential  $O(10e3)$  for faster and more sustainable simulation in HEP
- studying many generative models: (W)GANs, VAE, BIB-AEs, Normalising Flows, Point Cloud Diffusion Models,...
- working on fixed, regular 3d grids or with point clouds of energy depositions
- context of highly granular calorimeters for future colliders ideal context to study this
- and develop a **dedicated common tool** for this...



# DDFastShowerML

## fast simulation with ML inference in DD4hep

- based on FastSim hooks implemented in Geant4/DDG4
- split different necessary steps into dedicated components that can be - partly - reused:
  - **Trigger**: decide for w/ particles to use ML
  - **Model**: prepare input and output
    - typical for model specific
  - **Inference**: generic inference w/ **ONNX** and **libtorch**
    - common for all models
  - **Geometry**: local to global conversions of cell positions
    - barrel and endcap w/ some reuse between models (uses `ddrec::LayeredCalorimeter`)



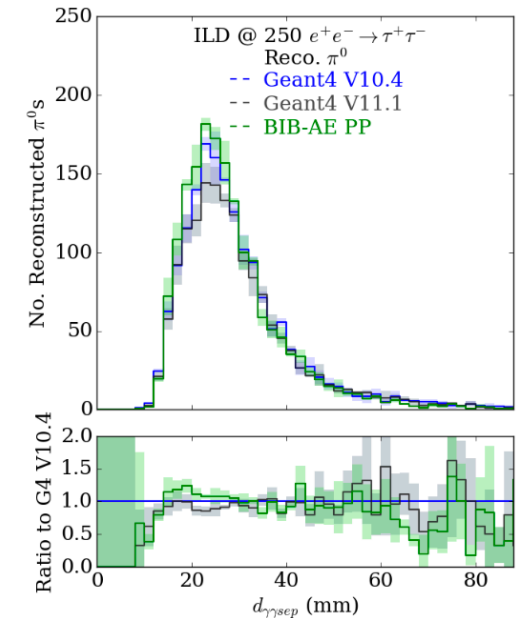
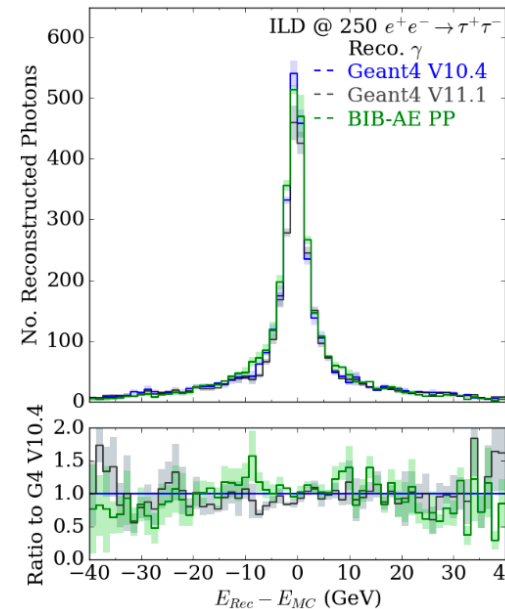
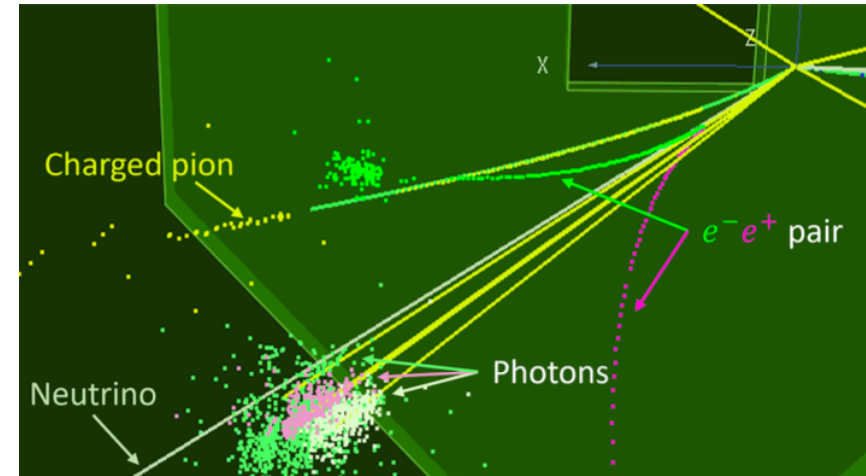
<https://gitlab.desy.de/ilcsoft/ddfastshowerml>



# DDFastShowerML

## example usage: tau-pair events in ILD

- with DDFastShowerML can include generative ML in **standard full simulation** to replace part of the showers with generative fast sim
- prerequisites:
  - need to be able to create ONNX or torch binary files:
    - typically have the whole ML network architecture running in one large **PyTorch** script
    - have a **DD4hep detector model** (with DDRRec extensions)
  - example:  $e^+e^- \rightarrow \tau^+ \tau^- \rightarrow \text{XXX} + \pi^0 \rightarrow \text{gamma gamma}$  at 250 GeV in ILD using a BIB-AE model conditioned on two angles

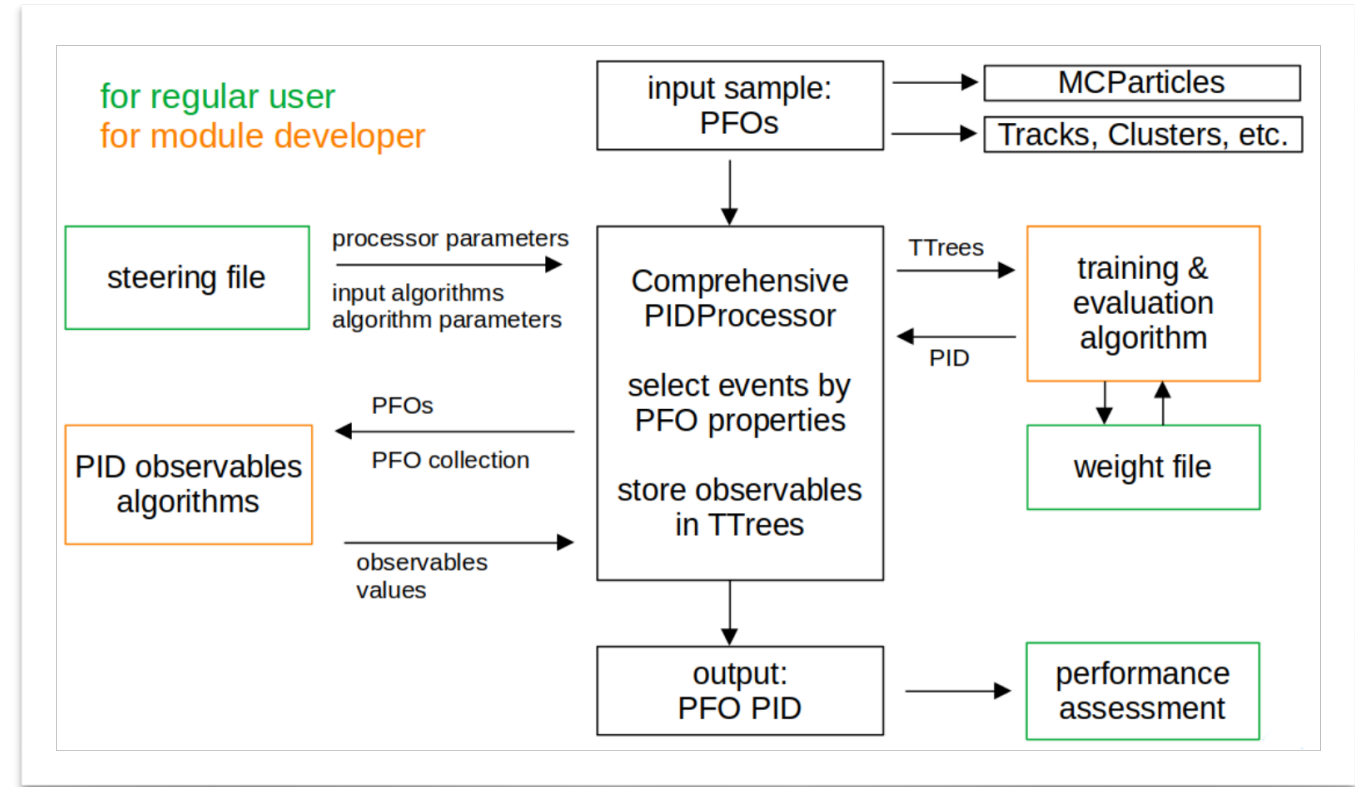


taus decaying in ILD detector w/ ML photons (PhD theses Peter McKeown)

# Particle ID with ML

## CPID: Comprehensive Particle identification Framework for Future Colliders

- implemented in LCIO and Marlin
  - see: [MarlinReco/Analysis/PIDTools](#)
- using **ROOT::TTree** as data exchange format
  - provides direct interface to **ROOT-TMVA**
  - easy conversion to **HDF5** (if needed)
- current implementation uses **TMVA** only
- could maybe extended to use **libtorch/ONNX** ...



<https://arxiv.org/pdf/2307.15635>

# Flavour Tagging

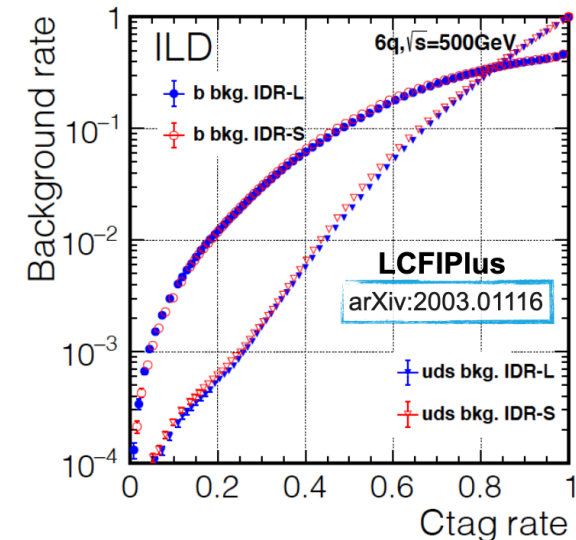
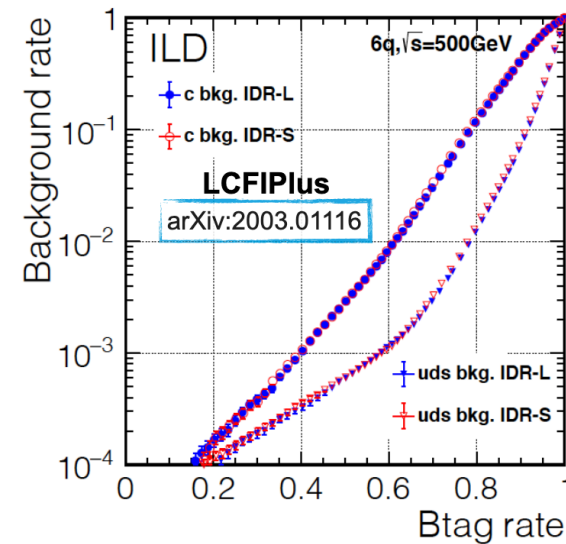
## With Deep Neural Networks

M.Meyer

- LCFIPlus flavour tagging workhorse for ILC studies for ~20 years
- recent progress in deep ML applied to LHC suggests significant performance increases possible
- study application of CMS DeepJet and ParticleNet to ILD ( in iLCSoft)

- **current standard** for heavy flavour tagging at ILD: **LCFIPlus**
- based on TMVA (BDTs)

arXiv:1506.08371,  
<https://github.com/lcfiplus/LCFIPlus>



➔ Can the **heavy flavour tagging** be improved by replacing the BDTs used in LCFIPlus with (deep) NNs?

# Flavour Tagging

## With Deep Neural Networks

M.Meyer



### DeepJet: input features

#### global variables

$p_{\text{jet}}, p_{\text{Tjet}},$   
 $N_{\text{charged jet const.}}, N_{\text{neutral jet const.}}, N_{\text{SV}}$   
additional global variables from LCFIPlus

21 input features

#### neutral jet constituents

$p_{\text{neutral const.}}, p_{\text{neutral const.}/p_{\text{jet}}}$   
 $\Delta R(\text{jet, neutral const.})$   
is photon?  
 $E_{\text{HCAL}}/E_{\text{HCAL+ECAL}}$

5 input features

#### charged jet constituents

$p_{\text{track}/p_{\text{jet}}}, p_{\text{Ttrack}}(\text{rel. jet}), \vec{p}_{\text{track}} \cdot \vec{p}_{\text{jet}/p_{\text{jet}}}$   
 $\Delta R(\text{track, jet})$   
impact parameter & significances  
track reconstructed in PV?  
lepton related variables  
pid variables  
 $\chi^2/\text{ndf}$

19 input features

#### secondary vertices

$m_{\text{SV}}$   
 $N_{\text{tracks in SV}}$   
 $\Delta R(\text{SV, jet})$   
 $E_{\text{SV}}/E_{\text{jet}}, E_{\text{SV}}$   
 $\cos(\text{flight direction}_{\text{SV}}, \vec{p}_{\text{SV}})$   
3D IP and significance  
 $\chi^2, \text{ndf}$

10 input features

### ParticleNet: input features

#### jet constituents: coordinates

$\Delta\eta, \Delta\Phi$

#### jet constituents: features

$\Delta\eta, \Delta\Phi$   
 $\log(p_{\text{T}}), \log(E), \log(p_{\text{T}}/p_{\text{Tjet}}), \log(E/E_{\text{jet}}),$   
 $\vec{p}_{\text{track}} \cdot \vec{p}_{\text{jet}/p_{\text{jet}}}$   
 $\Delta R$   
 $q$   
isElectron, isMuon, isChargedHadron,  
isNeutralHadron, isPhoton  
impact parameter & significances  
track used in PV?  
lepton related variables  
pid variables  
 $E_{\text{HCAL}}/E_{\text{HCAL+ECAL}}$   
 $\chi^2/\text{ndf}$

28 input features

#### secondary vertices: coordinates

$\Delta\eta, \Delta\Phi$

#### secondary vertices: features

$\Delta\eta, \Delta\Phi$   
 $\log(p_{\text{T}}), E_{\text{SV}}/E_{\text{jet}}, E_{\text{SV}}$   
 $\eta$   
 $m_{\text{SV}}$   
 $N_{\text{tracks in SV}}$   
 $\chi^2/\text{ndf}$   
impact parameters & significances  
 $\cos(\text{flight direction}_{\text{SV}}, \vec{p}_{\text{SV}})$

14 input features

2 SVs & all jet constituents considered, no ordering of inputs

- dedicated Marlin processor to prepare input variables from event (jet constituents):
  - VariablesForDeepMLFlavorTagger
- writes ROOT::Three file for training and assigns parameters (as PIDParameters) to PFOs for inference

# Flavour Tagging

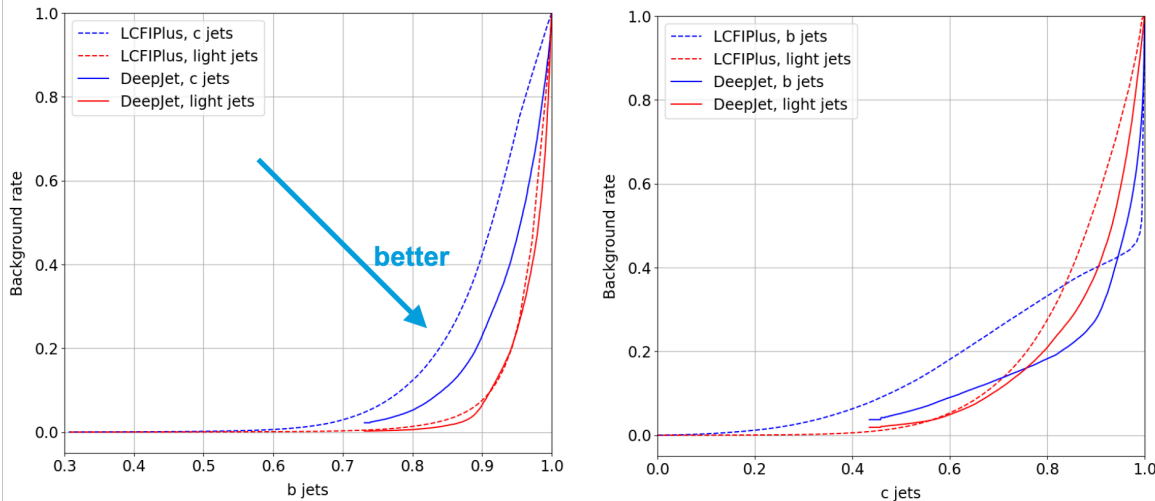
## With Deep Neural Networks

M.Meyer



### DeepJet: ROC curves - comparison to LCFIPlus

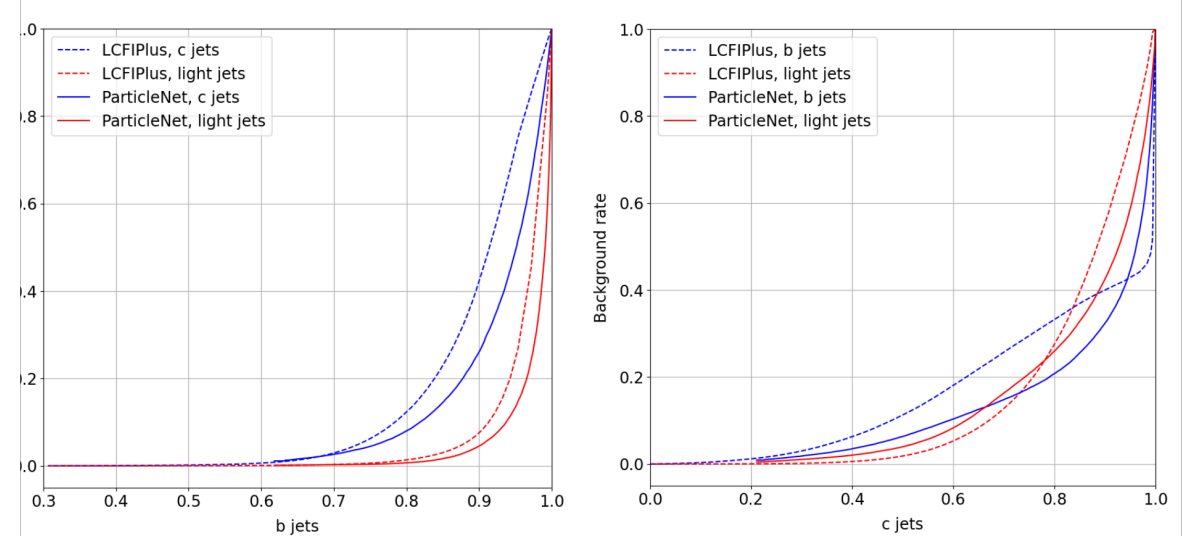
validation data



**better performance of DeepJet training over large parts of the b & c tagging efficiencies w.r.t default LCFIPlus used in ILD**

### ParticleNet: ROC curves - comparison to LCFIPlus

validation data



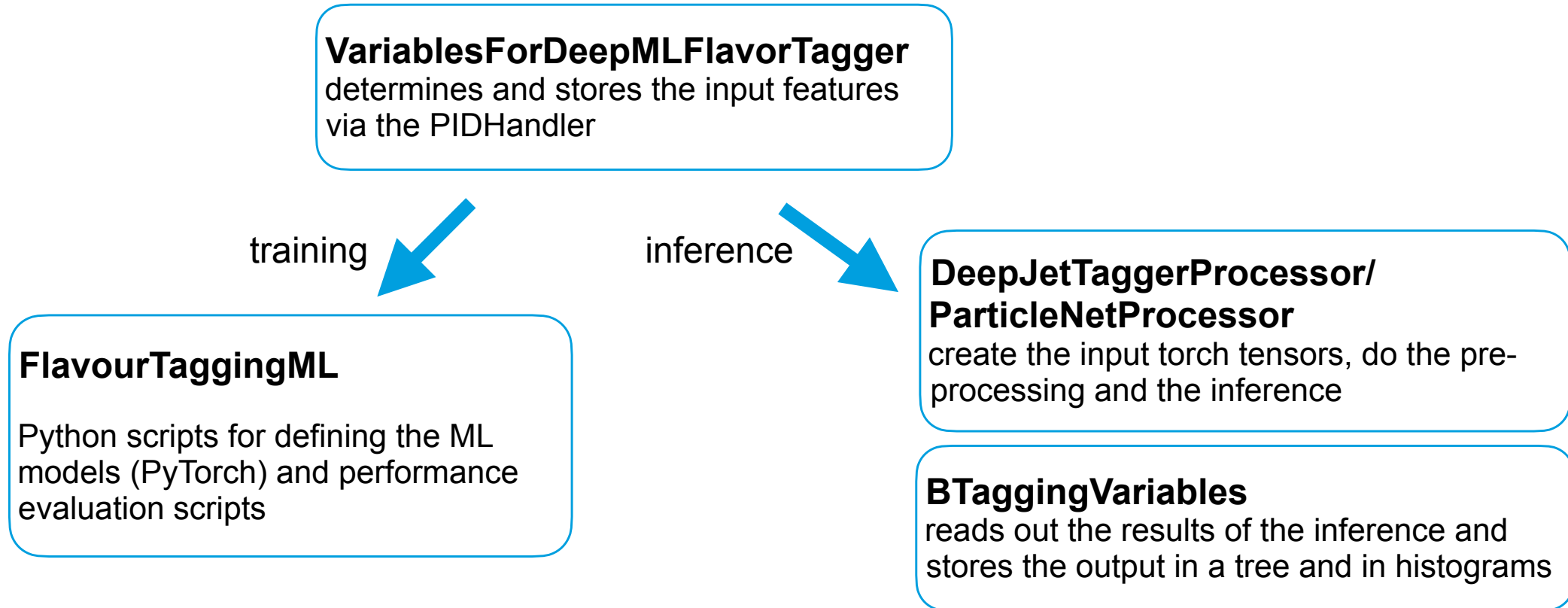
**better performance than LCFIPlus over large parts of the b and c tagging efficiencies**  
**one of the first trainings with this architecture, a lot of possibilities for optimization**  
(architecture, hyperparameters, features, over-training in c-jet category...)

- both deep deep learning methods achieve better performance than LCFIPlus

# Flavour Tagging

## With Deep Neural Networks - Packages and Modules

M.Meyer



- all packages available at <https://gitlab.desy.de/ilcsoft>
- serve as prototype example implementation and could be used in ILD standard reconstruction
- currently no manpower to continue working on this ...

# Future Developments

## and related open questions

- MarlinML with the JetTagger examples is still somewhat work in progress
  - however, serves as demonstrator and prototype for how ML inference can be done in Marlin for reconstruction and analysis
- need to also include other ongoing ML activities in Marlin, e.g.:
  - Particle Transformer flavour tagging and PFA w/ DNN ( see [talk T. Suehara 2nd ECFA H.EW ws](#))
- Open Questions:
  - can some code be reused as in DDML ?
    - generic ML inference (libtorch, ONNX, ...) -> probably straight forward
    - preparation of input variables (training/test/validation/inference data) -> probably quite some work needed to make this generic - maybe better to quickly write dedicated code ?
  - how can this be (re)used in Key4hep with Gaudi and EDM4hep ?
    - currently via MarlinWrapper
    - eventually need to create a **k4MLInference** package -> need EDM4hep prod. version

# Summary and Outlook

- being able to run ML inference in Key4hep/iLCSoft absolutely mandatory
- some prototypes and example implementations exist:
  - DDFastShowerML, MarlinML (JetTagger), CPID
- work in progress towards:
  - more generic tools - how many different applications can be combined in one package ?
  - proper Key4hep examples using EDM4hep and Gaudi
  - manpower limited ...
- your contribution and thoughts are welcome !