

Muon momentum estimation in scattering tomography

Fourth MODE Workshop on Differentiable Programming for Experiment Design
23-25 September 2024 Valencia

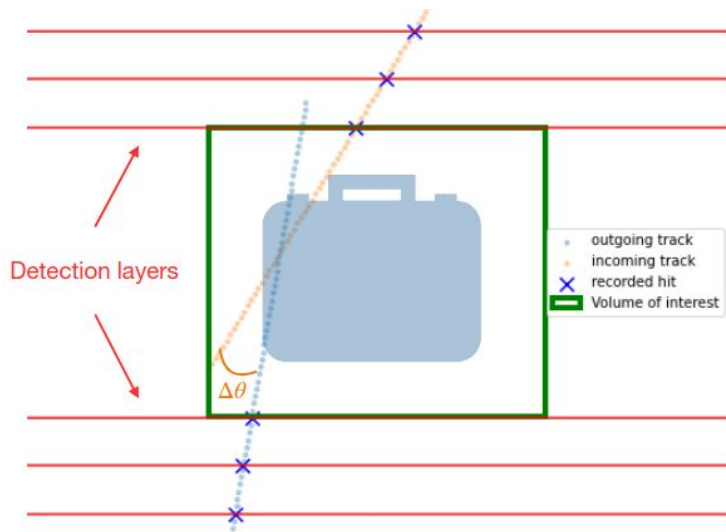
Florian Bury, Maxime Lagrange, for the TomOpt team





Thank you for being here at 9am !!!

Scattering tomography



For a given :

- Material with $X_0(A, \rho)$
- Momentum p

The deflection angle is a centred Gaussian with

$$\sigma_\theta = \frac{13.6 \text{ MeV}}{p} \sqrt{\frac{x}{X_0}}$$

Measuring track deflection gives access to X_0

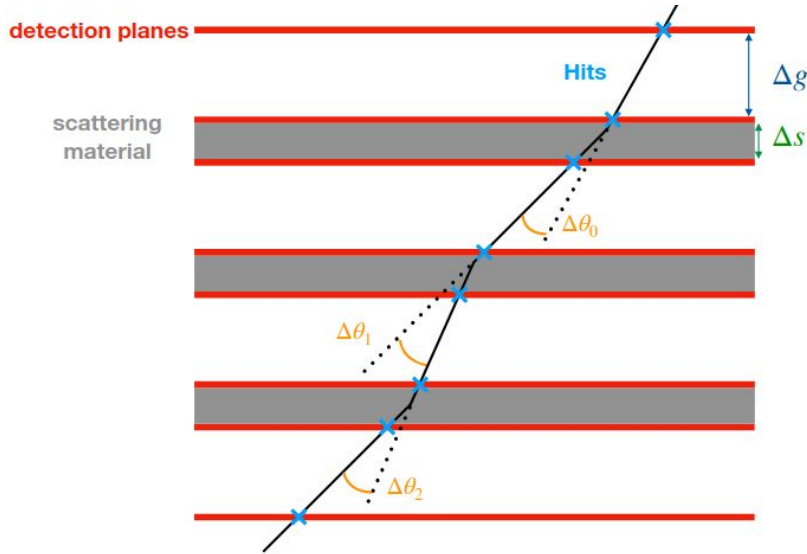
Measurement method : $\sigma_\theta \sim \theta_{RMS} = \frac{1}{N} \sum_i \Delta\theta_i^2$

Example : yesterday's great [talk](#) by Zahraa !

Scattering tomography

For a given :

- Material with $X_0(A, \rho)$
- Momentum p



The deflection angle is a centred Gaussian with

$$\sigma_\theta = \frac{13.6 \text{ MeV}}{p} \sqrt{\frac{x}{X_0}}$$

Measuring track deflection gives access to X_0

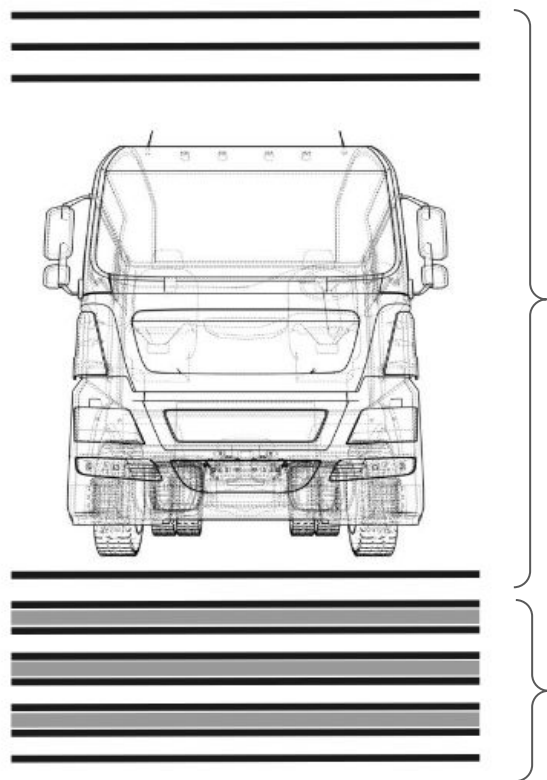
Measurement method : $\sigma_\theta \sim \theta_{RMS} = \frac{1}{N} \sum_i \Delta\theta_i^2$

... but this assumes knowledge of p

Possible solution : inverse the formula

$$p = \frac{13.6 \text{ MeV}}{\theta_{RMS}} \sqrt{\frac{x}{X_0}}$$

Full detector objective



Tomography section

$$\sigma_{\theta} \sim \theta_{RMS} = \frac{1}{N} \sum_i^N \Delta\theta_i^2$$

$$\sigma_{\theta} = \frac{13.6 \text{ MeV}}{p} \sqrt{\frac{x}{X_0}}$$

Measurement objective

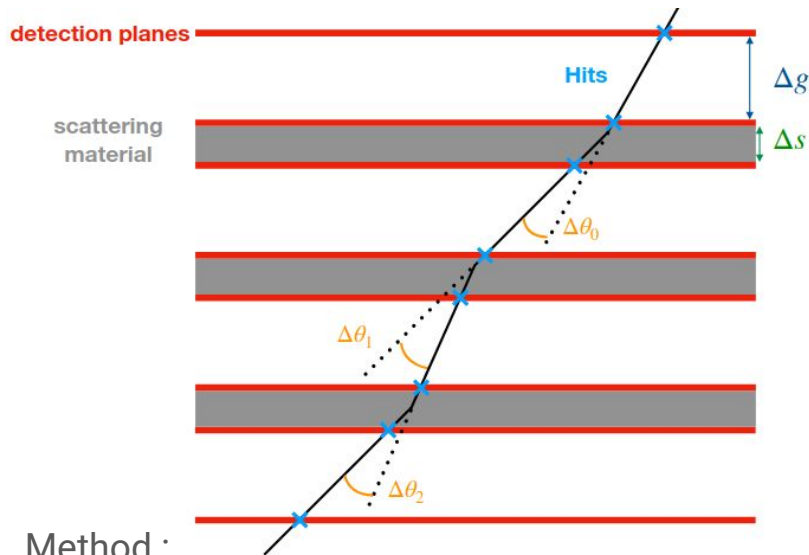
Momentum section

$$p = \underbrace{\frac{13.6 \text{ MeV}}{\theta_{RMS}}}_{\text{Measured}} \underbrace{\sqrt{\frac{x}{X_0}}}_{\text{Known by construction (lead, iron)}}$$

Measured

Known by construction
(lead, iron)

Momentum measurement

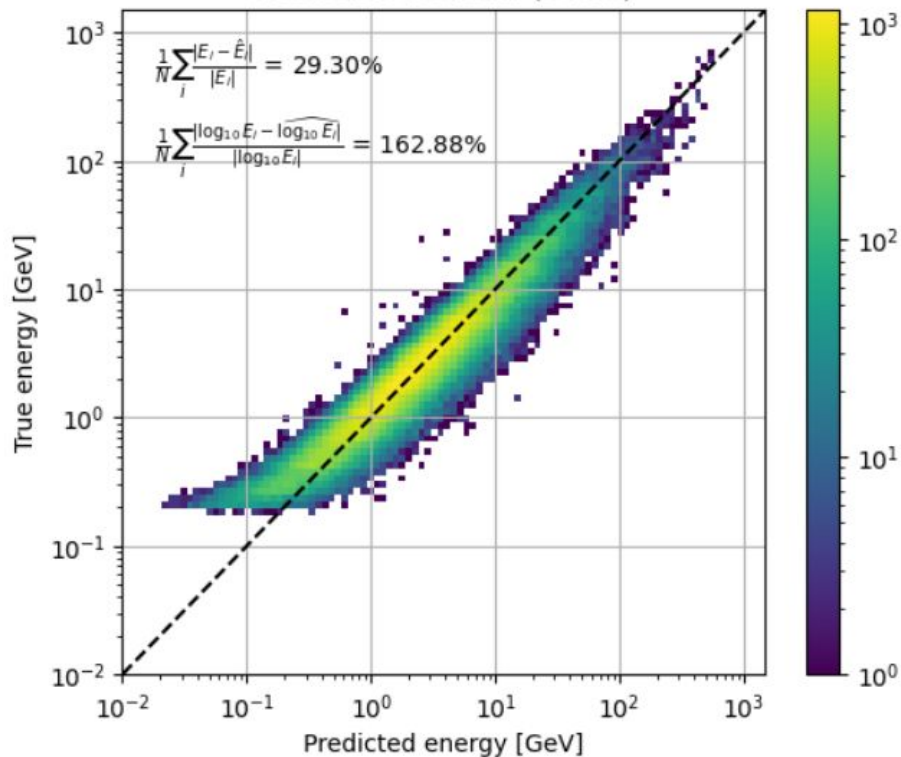


Method :

- From hits obtain segments
- From segments obtain scattering angles
- Take RMS and plug in

$$p = \frac{13.6 \text{ MeV}}{\theta_{RMS}} \sqrt{\frac{x}{X_0}}$$

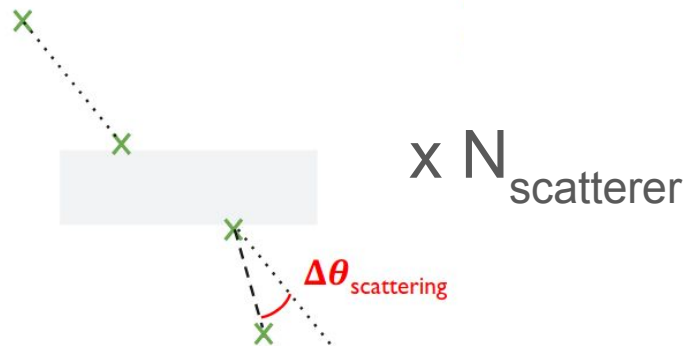
“Energy” = Kinetic energy



Maybe we can extract more than that

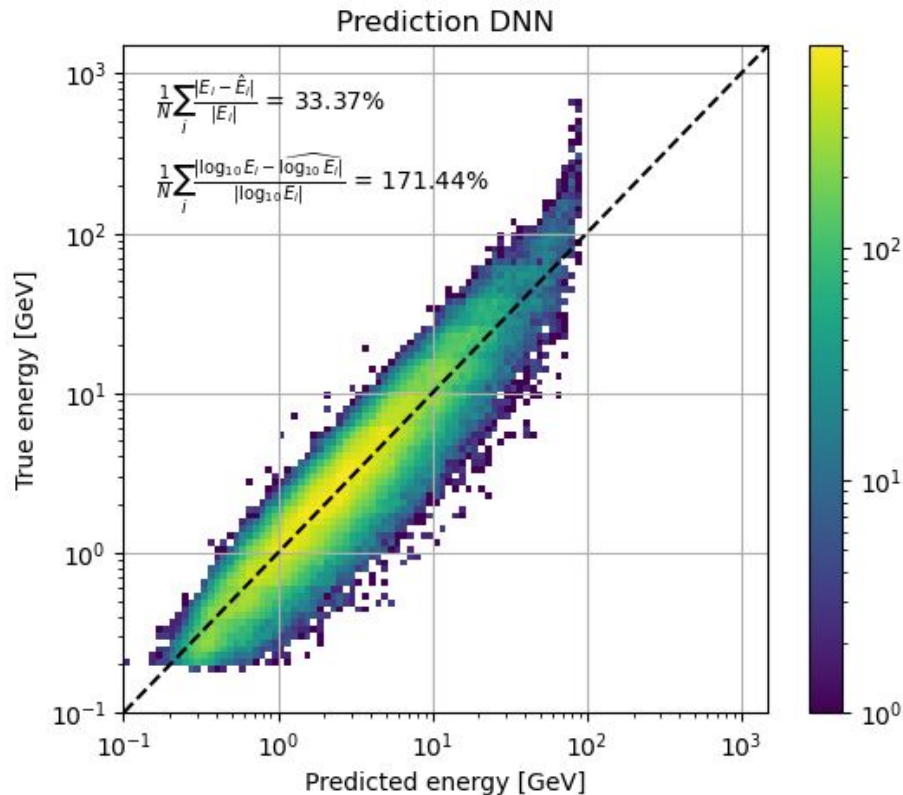
First ML attempt

Let us use a simple 3-layer DNN on the scattering angles

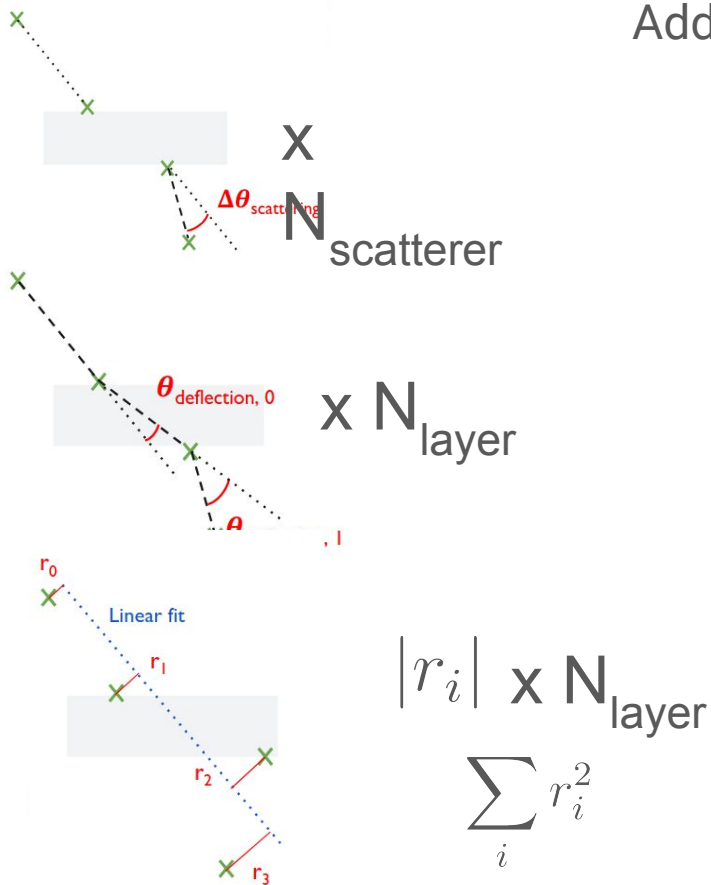


Improves at low energies
... issues at high energies

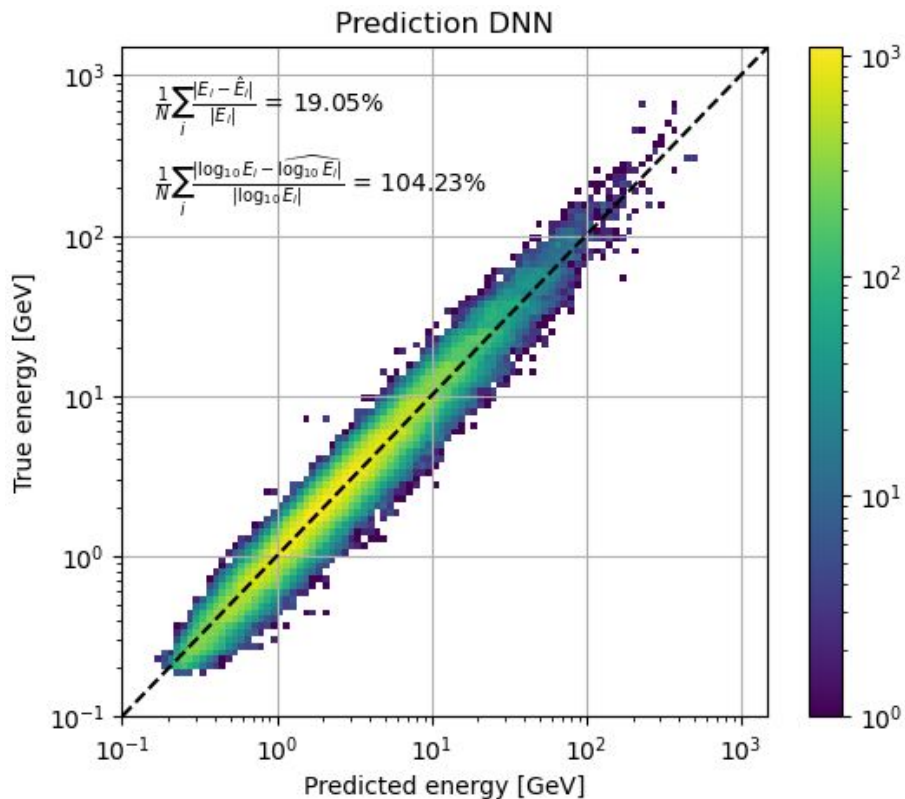
$$\Delta\theta_{\text{scattering}} \xrightarrow{E \rightarrow \infty} 0$$



Improved attempt

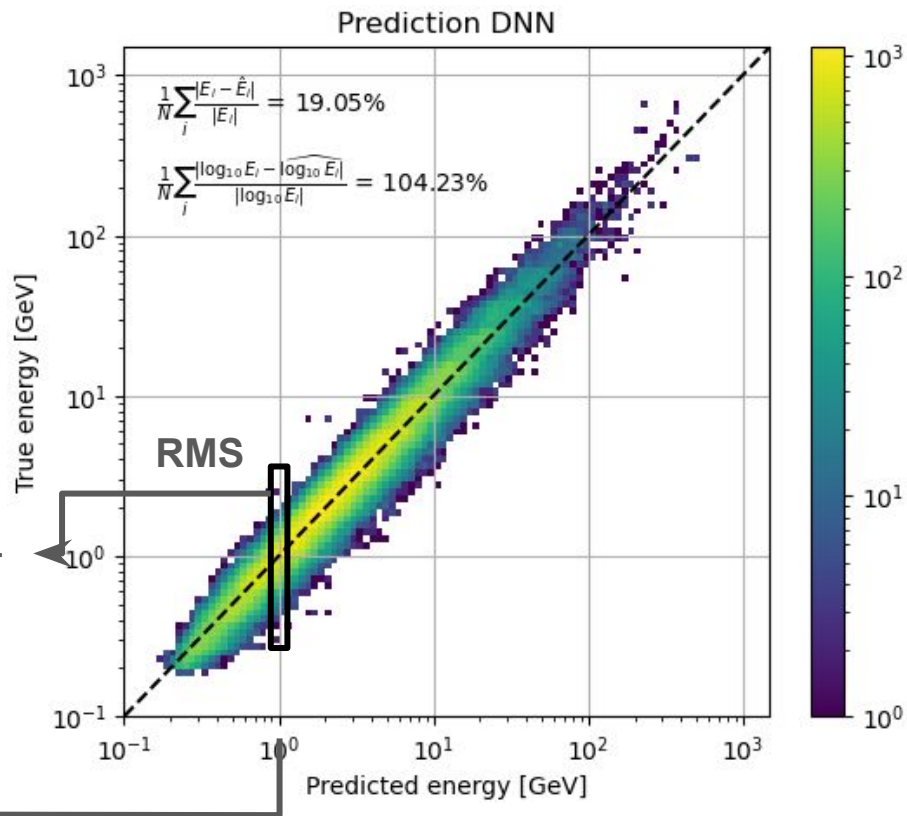
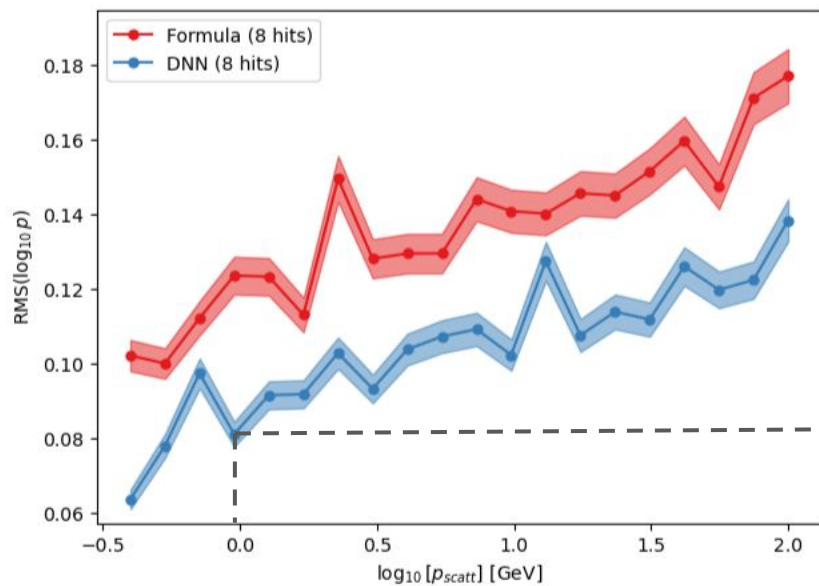


Adding more interesting variables

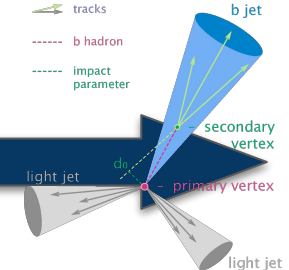


Improved attempt

A more quantitative comparison



Intermezzo : on the story of jet taggers



LHC run-1

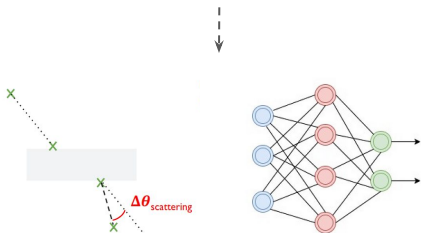
Jet probability taggers [1]

- Purely analytical
- Provides “likelihood” value

$$p = \frac{13.6 MeV}{\theta_{RMS}} \sqrt{\frac{x}{X_0}}$$

Combined Secondary Vertex (csv) [1]

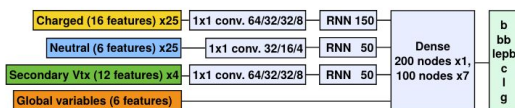
- Shallow neural network
- Provide a score based on several inputs
- Later on, went deeper (DeepCSV [2])



LHC run-2

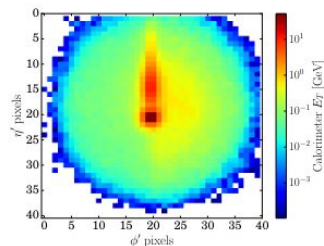
DeepJet [3]

- More complex architecture (CNN,RNN)
- Variable-length inputs (tracks)



Jet as images [4]

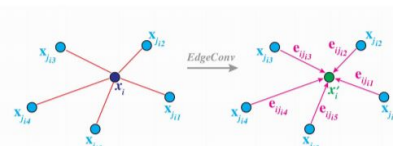
- Jet as image
- CNN architecture
- End-to-end



LHC run-3

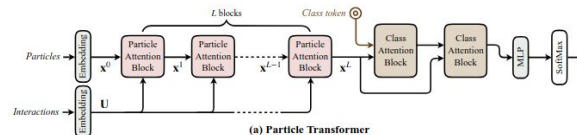
ParticleNet [5]

- Graph-based architecture
- Variable-length
- Sparse inputs



ParticleTransformer [6]

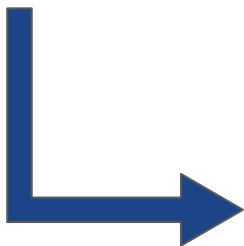
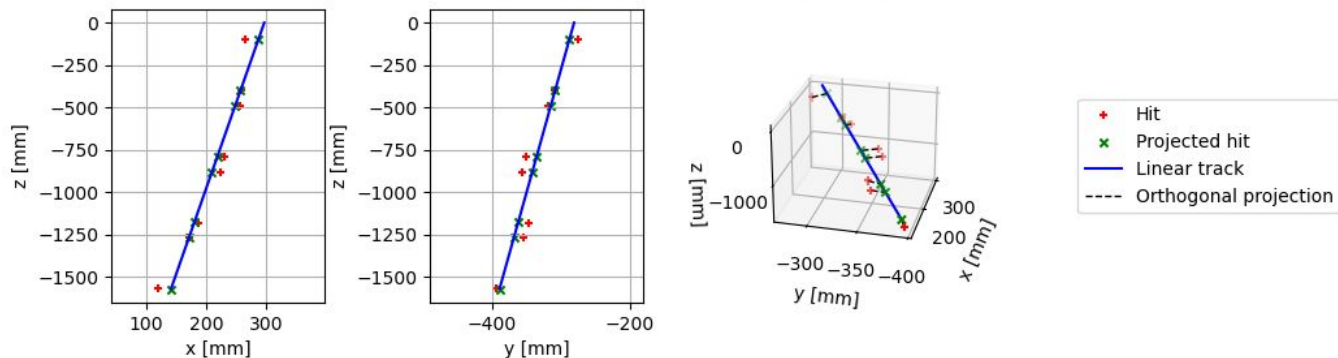
- Transformer (eg GPT) architecture
- Attention mechanism
- Variable-length



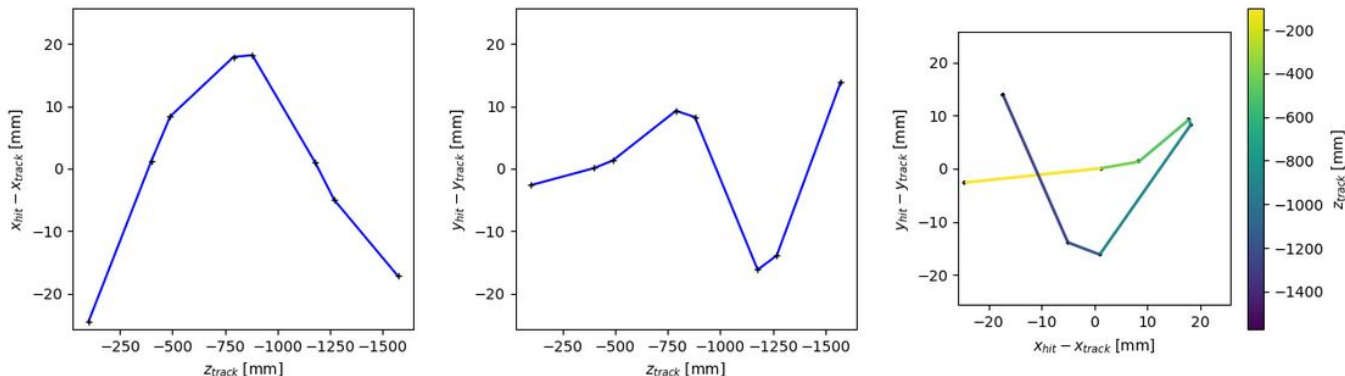
More complex architectures

Only deviations from a linear track carry information about the energy

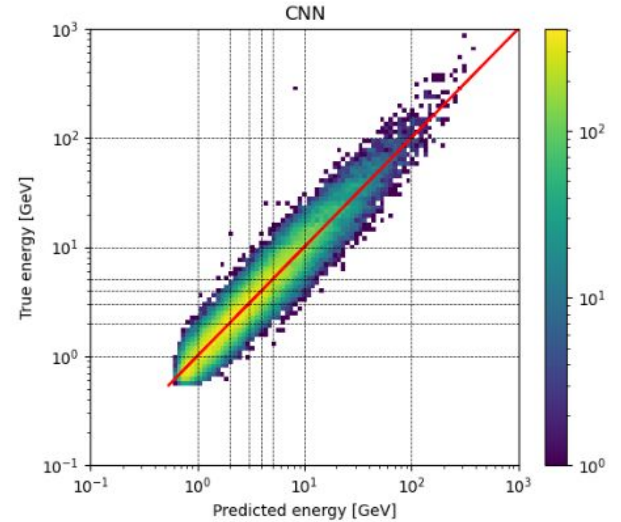
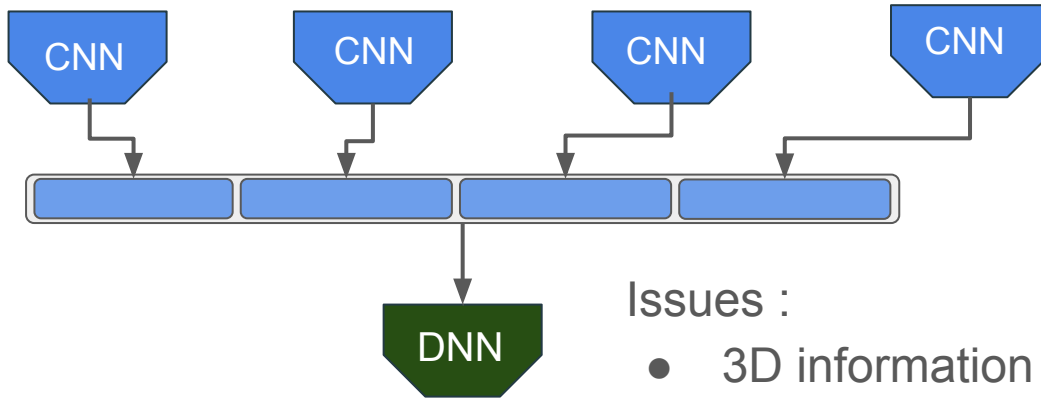
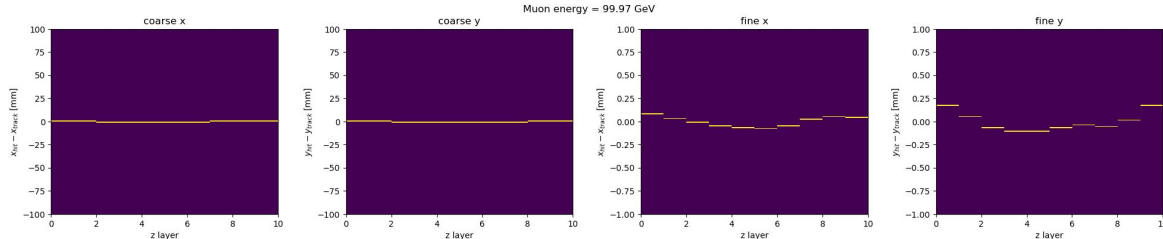
Hits for event 268364 [Energy = 0.19 GeV]



Deviations from linear track



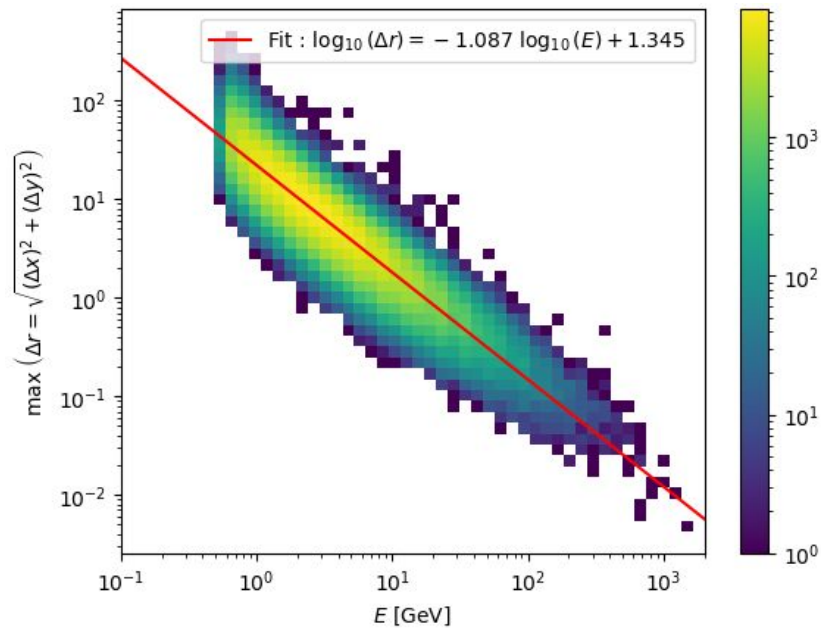
More complex architectures : deviations as an image



Issues :

- 3D information projected into 2D
- Pixelized image, need different scales
- Very sparse input (energy/time inefficient)
- 3D images would be too heavy

Dealing with scale

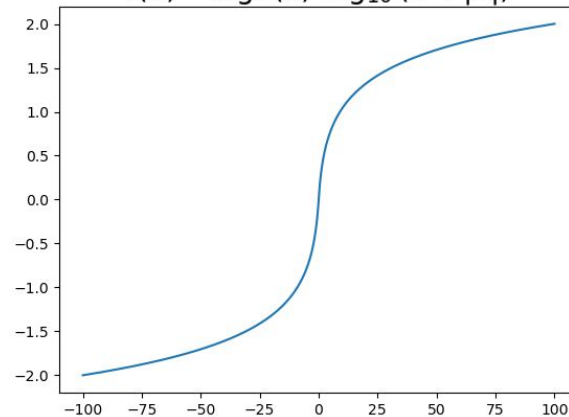


$$\sigma_\theta = \frac{13.6 \text{ MeV}}{p} \sqrt{\frac{x}{X_0}}$$

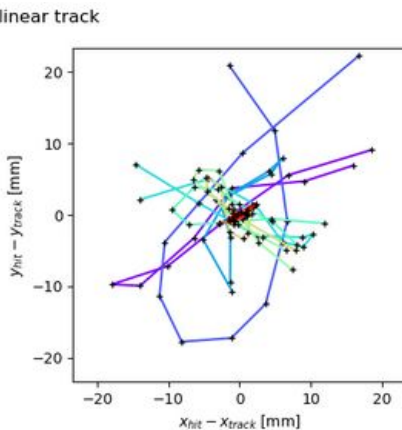
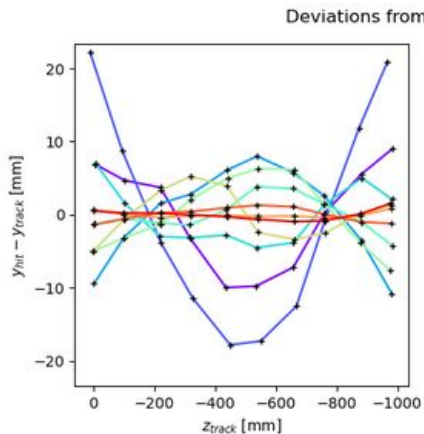
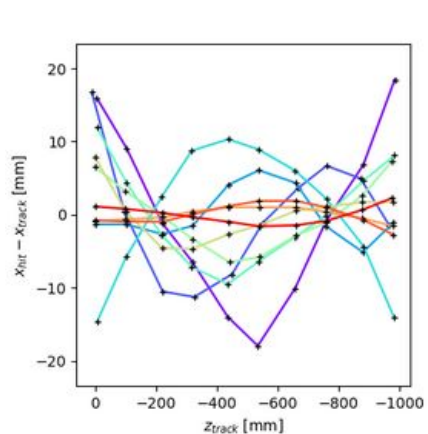


$$\max_{\text{layer}} \Delta r \propto \frac{1}{p}$$

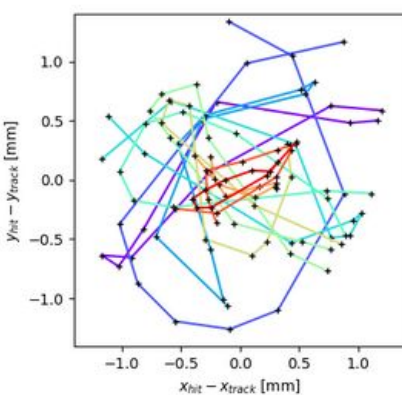
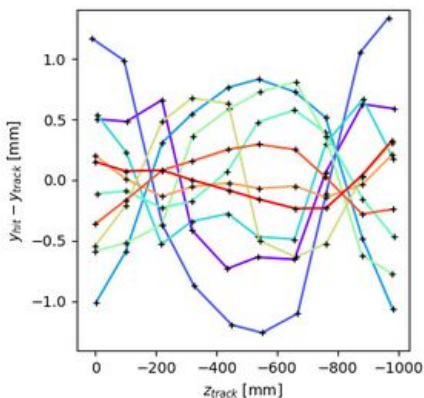
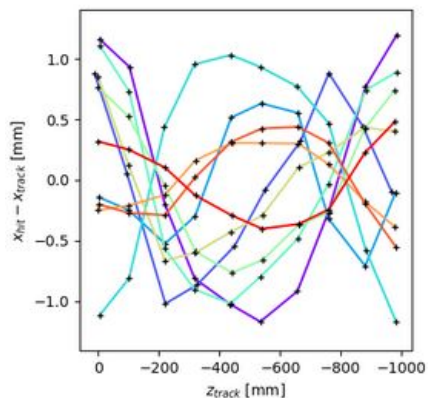
$$f(x) = \text{sign}(x) \log_{10}(1 + |x|)$$



Dealing with scale

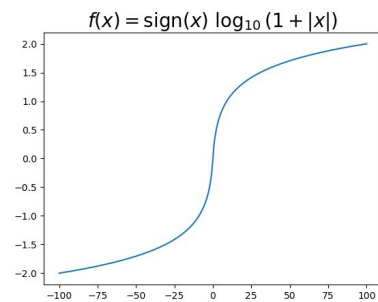


Deviations from linear track



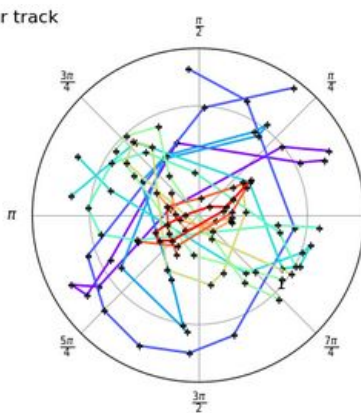
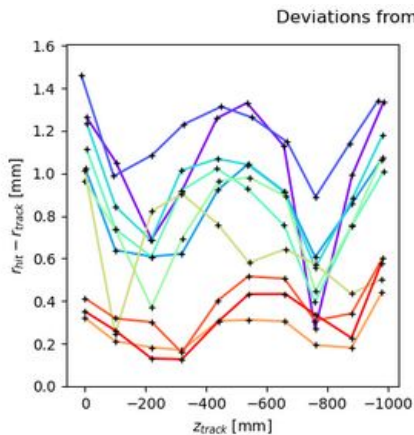
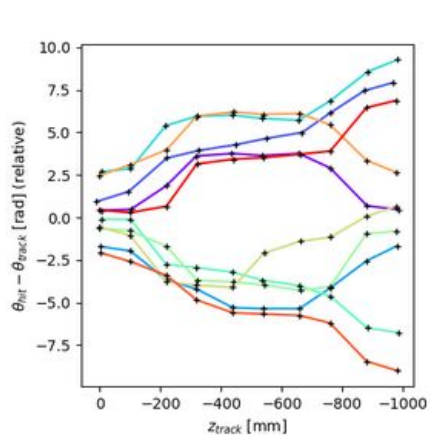
Deviations from linear track

- Track E = 0.80 GeV
- Track E = 1.52 GeV
- Track E = 1.61 GeV
- Track E = 1.72 GeV
- Track E = 2.00 GeV
- Track E = 2.43 GeV
- Track E = 2.55 GeV
- Track E = 9.79 GeV
- Track E = 11.19 GeV
- Track E = 18.40 GeV

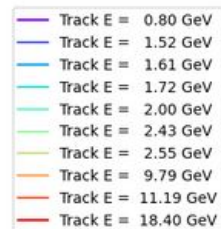
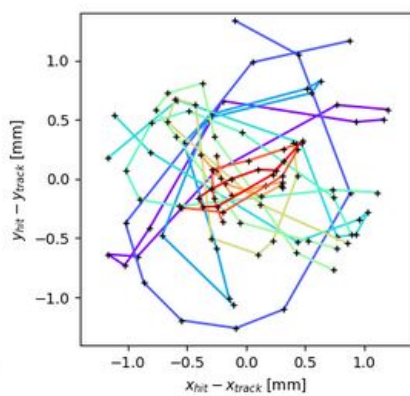
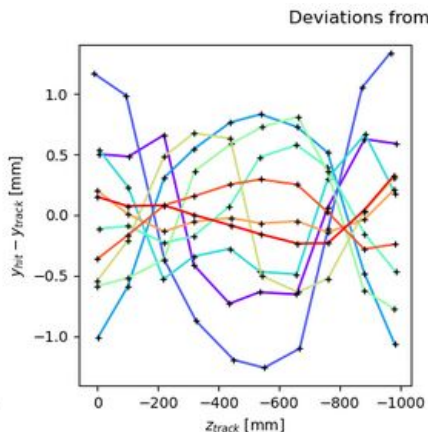
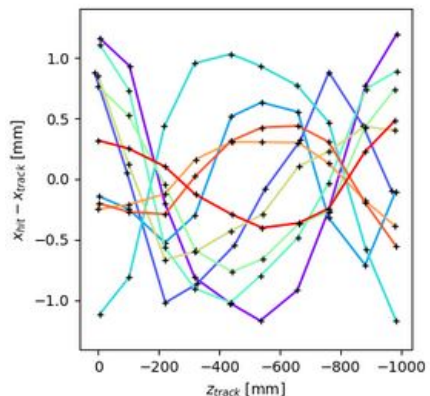


Apply f to Δr
(keeping same direction)

Dealing with scale

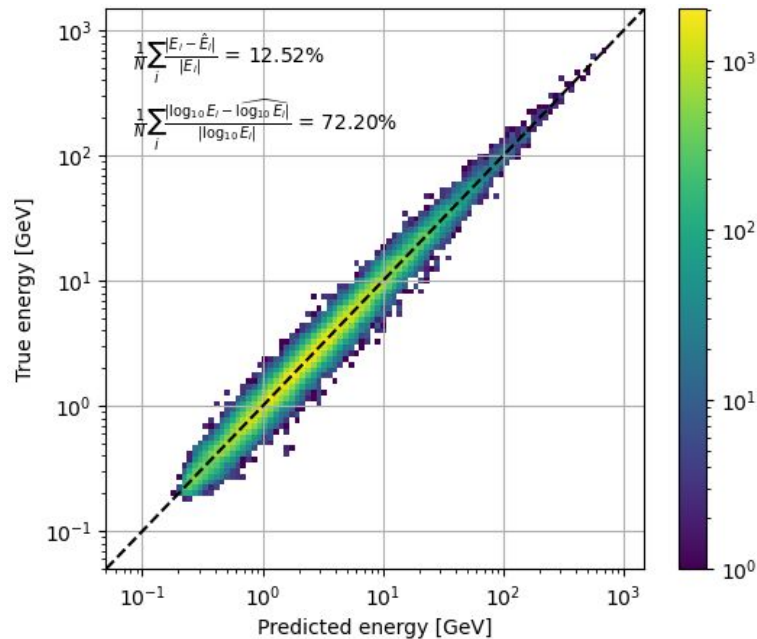
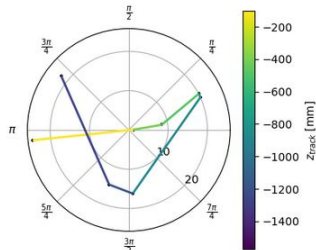
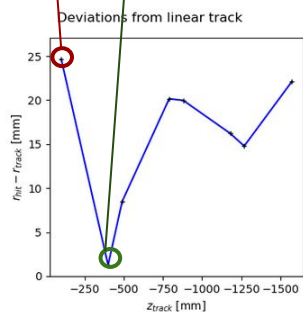
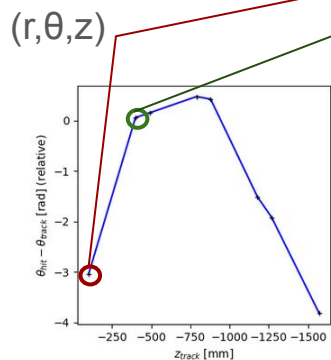
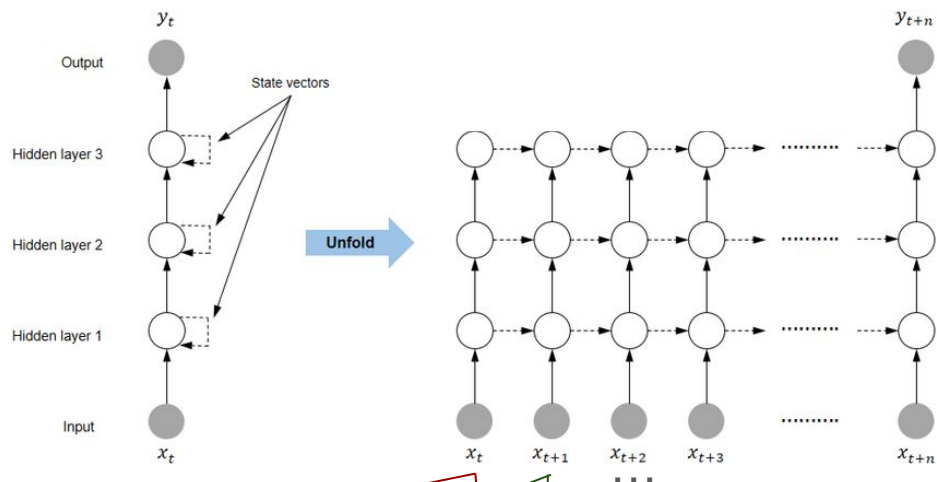


Moving to cylindrical coordinates
(r, θ, z)



More complex architectures : deviations as a sequence

Recurrent neural network (RNN) “Many-to-one” scheme

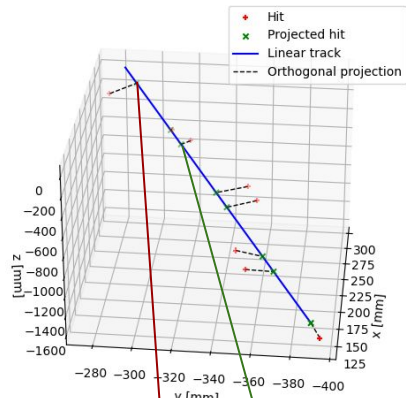


Issues :

- Sequential processing (not parallel, slow)
- No geometrical info

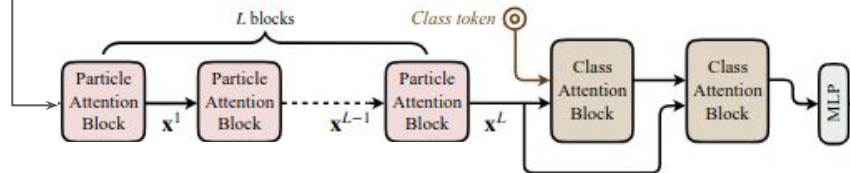
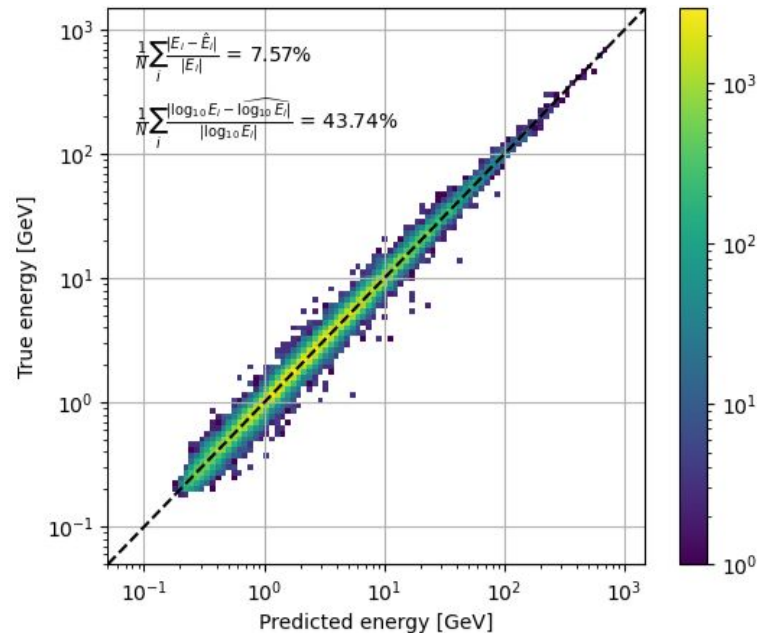
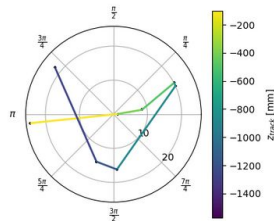
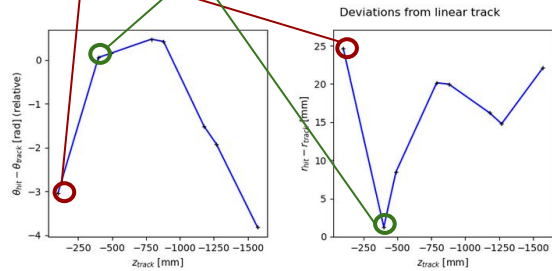
More complex architectures : deviations as a sequence

Transformer encoder

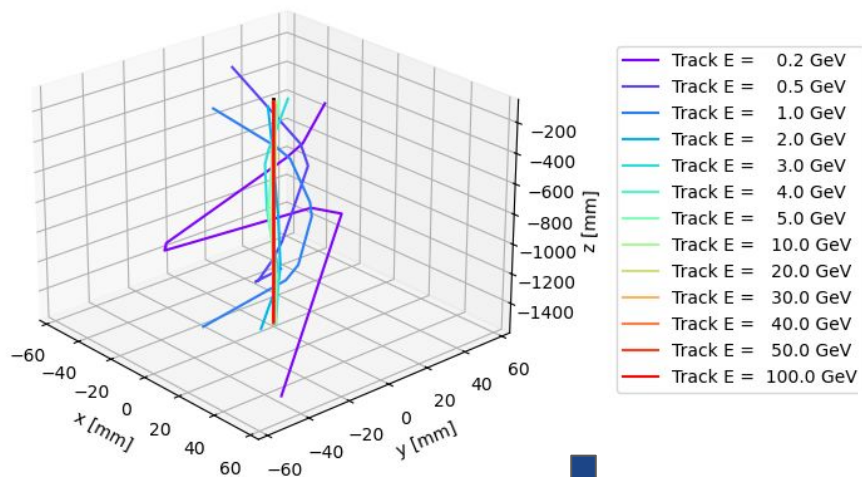


Compared to RNN :

- Parallel processing
- Attention mechanism
- Long+short distance interactions



More complex architectures : deviations as a graph

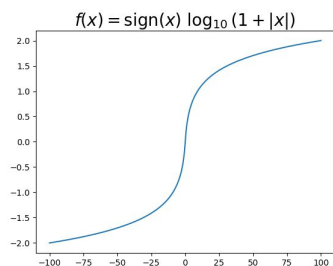


Graph is set of $\{\Delta x_i, \Delta y_i, z_i\}$

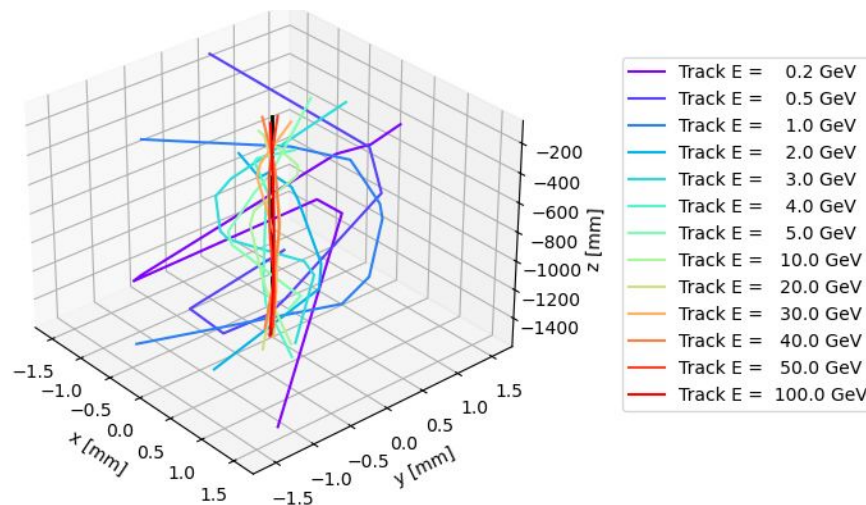
($\Delta x_i, \Delta y_i$ are the deviations at height z_i)

Same problem as before :

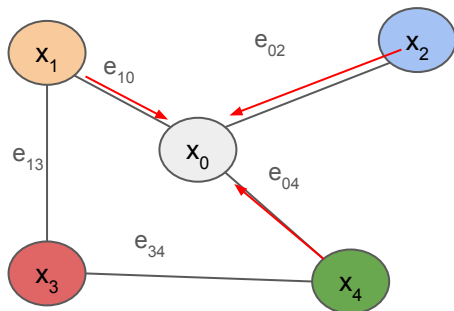
$$\max \Delta x_i \propto \frac{1}{p}$$



(Applied on the radius, keeping same angle)



More complex architectures : deviations as a graph



$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left(\mathbf{x}_i^{(k-1)}, \underbrace{\bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)}(\mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{ij})}_{\text{“Combine” operation}} \right)$$

Optional embedding of the concatenated previous step with aggregate+combine

“Aggregate” operation

- Permutation invariant
- Typically min,max, mean, std, ...
- Can be more complex (ML)

“Combine” operation

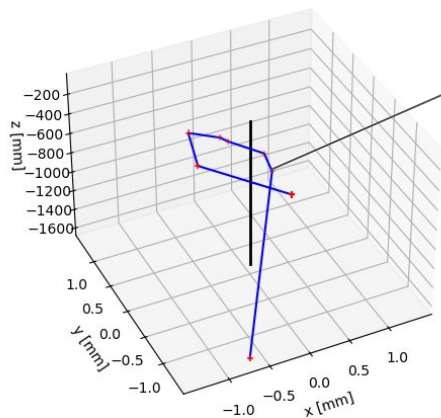
- Applied on the node features and the ones of his immediate neighbours
- Can include edge features
- Usually some ML operation

After N steps : $h_G = \bigoplus_{i \in G} x_i^{(N)}$ then fed to a DNN for regression

(\oplus can again be simple (eg, mean) or complex (ML))

More complex architectures : deviations as a graph

Track deviations for event 9 [Energy = 0.5 GeV]



$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left(\mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \underbrace{\phi^{(k)}(\mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{ij})}_{\text{Combine}} \right)$$

Optional embedding of the concatenated previous step with aggregate+combine

“Aggregate” operation

- Permutation invariant
- Typically min,max, mean, std, ...
- Can be more complex (ML)

“Combine” operation

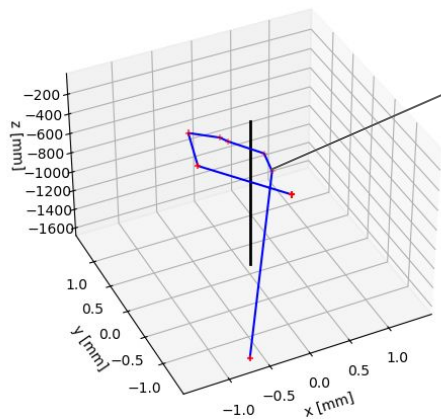
- Applied on the node features and the ones of his immediate neighbours
- Can include edge features
- Usually some ML operation

After N steps : $h_G = \bigoplus_{i \in G} x_i^{(N)}$ then fed to a DNN for regression

(\oplus can again be simple (eg, mean) or complex (ML))

More complex architectures : deviations as a graph

Track deviations for event 9 [Energy = 0.5 GeV]



$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left(\mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \underbrace{\phi^{(k)}(\mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{ij})}_{\text{Combine}} \right)$$

Optional embedding of the concatenated previous step with aggregate+combine

“Aggregate” operation

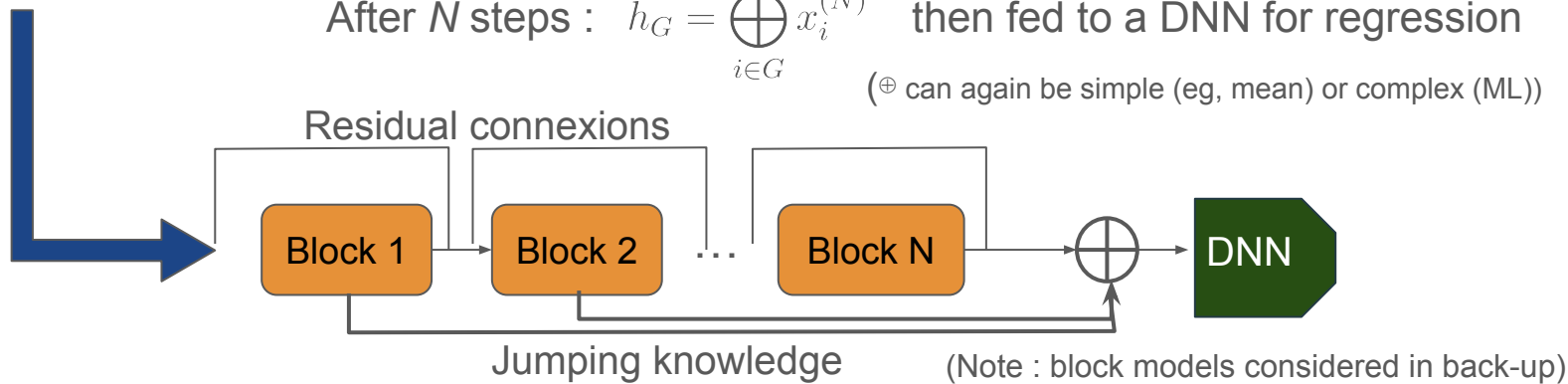
- Permutation invariant
- Typically min,max, mean, std, ...
- Can be more complex (ML)

“Combine” operation

- Applied on the node features and the ones of his immediate neighbours
- Can include edge features
- Usually some ML operation

After N steps : $h_G = \bigoplus_{i \in G} x_i^{(N)}$ then fed to a DNN for regression

(\oplus can again be simple (eg, mean) or complex (ML))

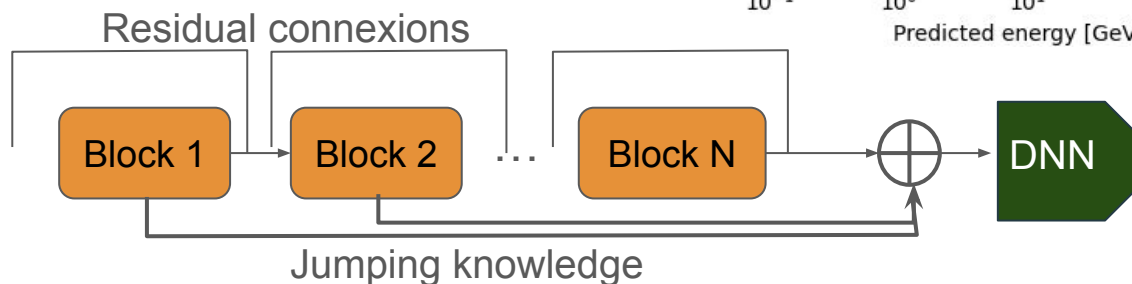
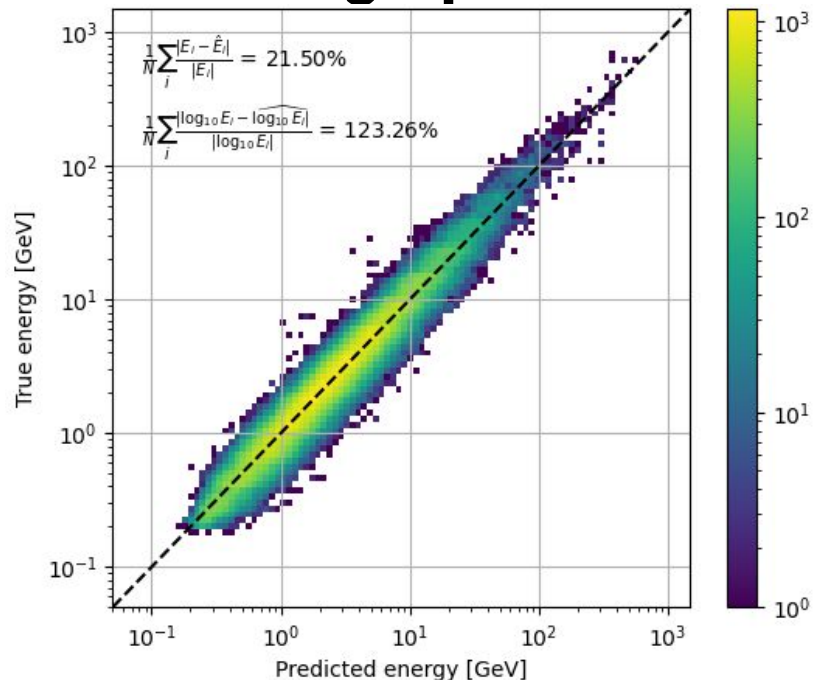


More complex architectures : deviations as a graph

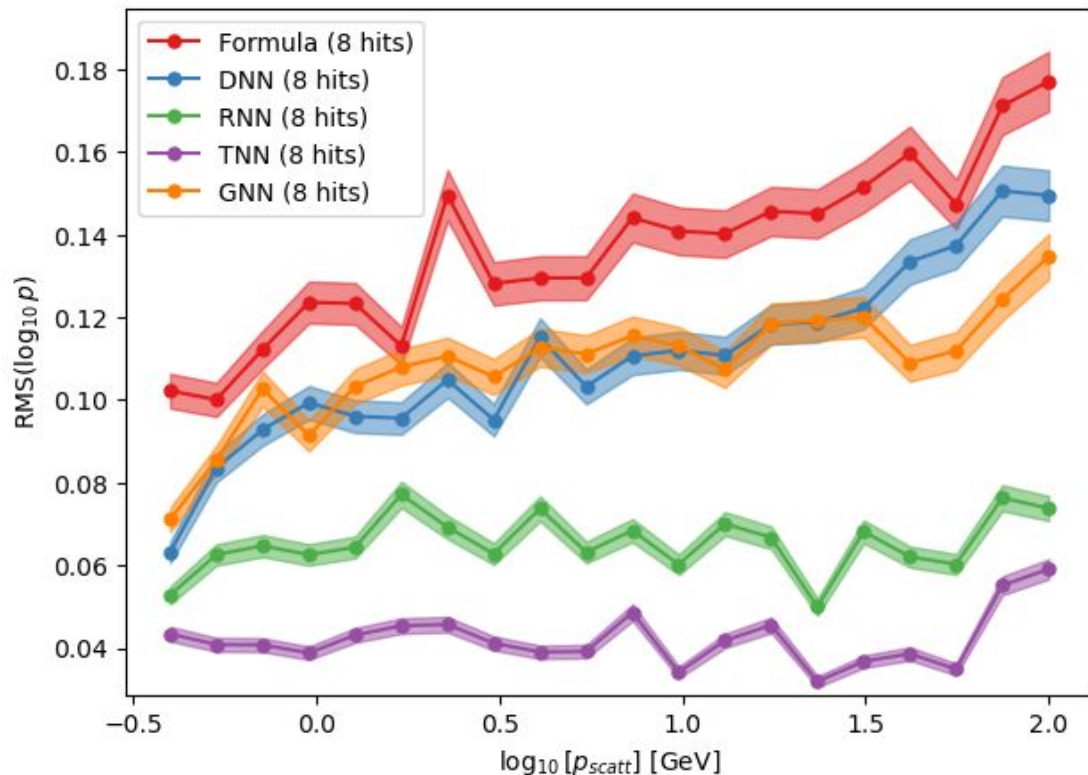
$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left(\mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)}(\mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{ij}) \right)$$

Here PointGCN
(other block models in back-up)

Not satisfactory, still have to play with
aggregation and models



More complex architectures : comparing performance



Model summary

- Model size
 - DNN : 18K params / 0.07 MB
 - RNN : 165K params / 0.62 MB
 - TNN : 430K params / 1.64 MB
 - GNN : 100K params / 0.40 MB
 - Computation time
(batch size = 1024, GPU GeForce RTX 2080)
 - DNN : 21 μs / muon
 - RNN : 60 μs / muon
 - TNN : 24 μs / muon
 - GNN : 80 μs / muon
- (note : CPU batching not optimised)

Real-life requirements

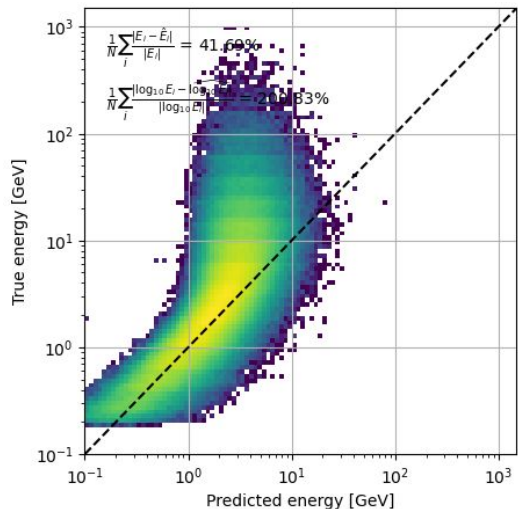
Effect of spatial resolution

Using a Gaussian smearing on x and y hit position

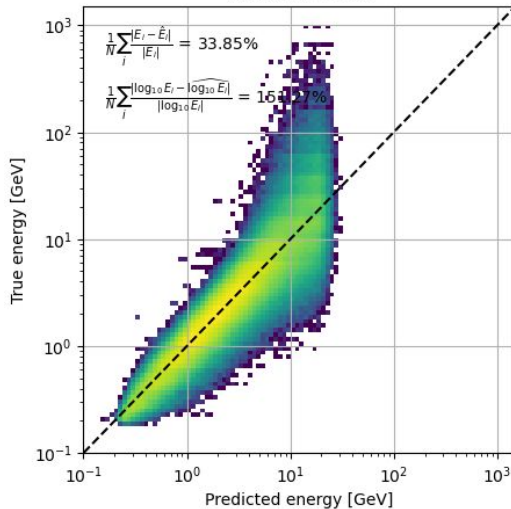
$$x, y \sim \mathcal{N}(\mu = 0, \sigma = 1 \text{ mm})$$

(typical resolution of RPC detectors)

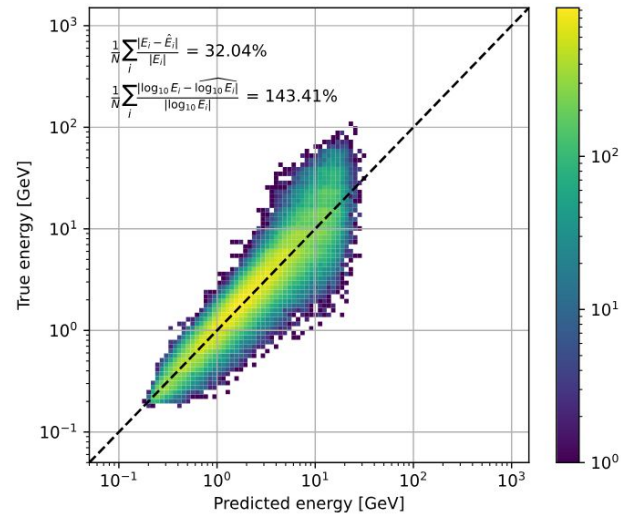
Scattering formula



DNN



TNN



ML methods can limit the issue of the limited precision ... up to a point

Real-life requirements

Effect of hit efficiency

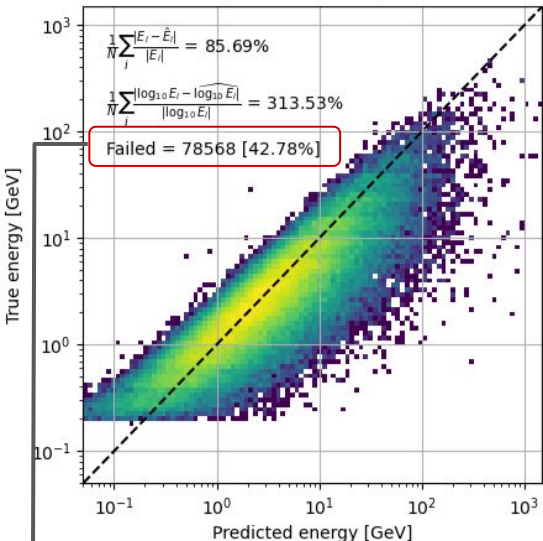
$$p = \frac{13.6 \text{ MeV}}{\theta_{RMS}} \sqrt{\frac{x}{X_0}}$$

For 8 detector planes (3 scatterer volumes) :

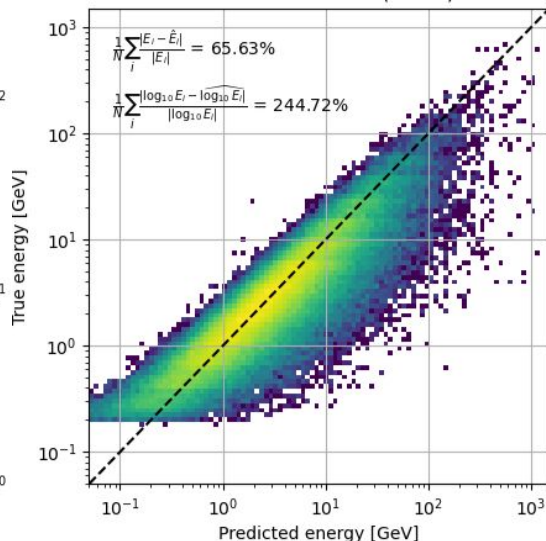
- $\epsilon = 90\%$: $N_{8 \text{ hits}} = 43\%$, $N_{7 \text{ hits}} = 38\%$, $N_{6 \text{ hits}} = 15\%$
- $\epsilon = 95\%$: $N_{8 \text{ hits}} = 66\%$, $N_{7 \text{ hits}} = 28\%$, $N_{6 \text{ hits}} = 5\%$
- $\epsilon = 99\%$: $N_{8 \text{ hits}} = 92\%$, $N_{7 \text{ hits}} = 7\%$, $N_{6 \text{ hits}} = <1\%$

3 angles \rightarrow can be reduced when hit is missed
(but worse performance)

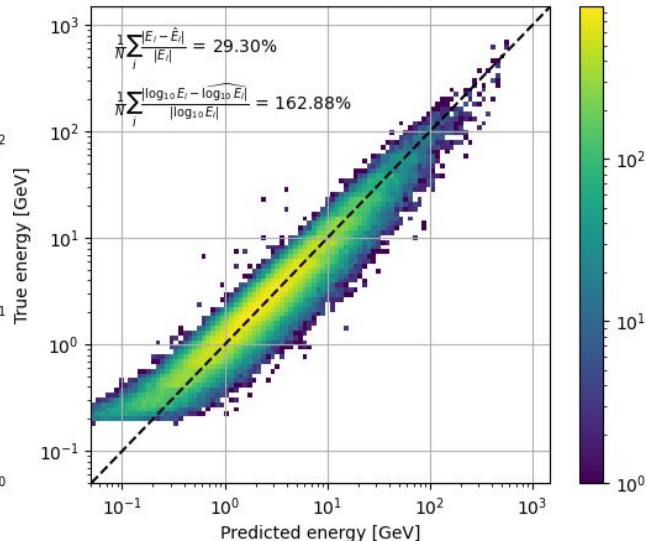
2 missed hits (total = 6)



1 missed hit (total = 7)



0 missed hit (total = 8)

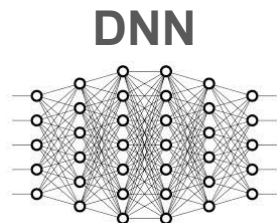


Two missed hits at the wrong places, and no segment can be built around scatterers

Real-life requirements

Effect of hit efficiency

Missed variables are zero-padded

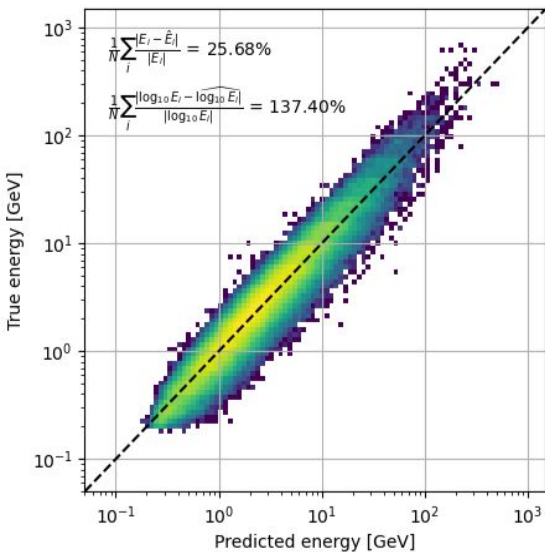


For 8 detector planes (3 scatterer volumes) :

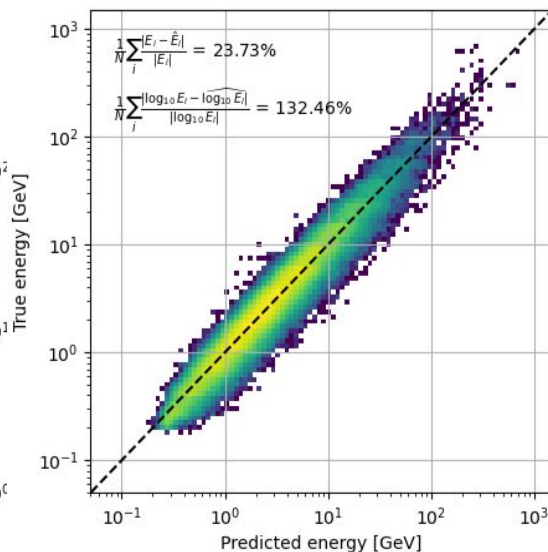
- $\epsilon = 90\%$: $N_{8 \text{ hits}} = 43\%$, $N_{7 \text{ hits}} = 38\%$, $N_{6 \text{ hits}} = 15\%$
- $\epsilon = 95\%$: $N_{8 \text{ hits}} = 66\%$, $N_{7 \text{ hits}} = 28\%$, $N_{6 \text{ hits}} = 5\%$
- $\epsilon = 99\%$: $N_{8 \text{ hits}} = 92\%$, $N_{7 \text{ hits}} = 7\%$, $N_{6 \text{ hits}} = <1\%$

Model is trained on a mix of datasets
(fit and deviations done only on recorded hits)

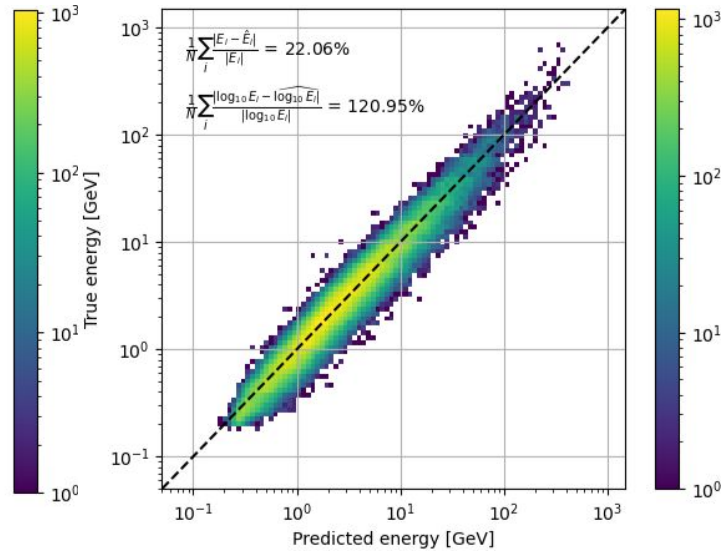
2 missed hits (total = 6)



1 missed hit (total = 7)



0 missed hit (total = 8)

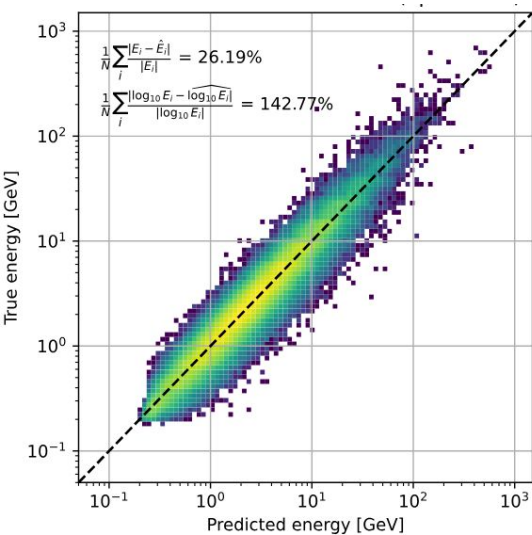


Real-life requirements

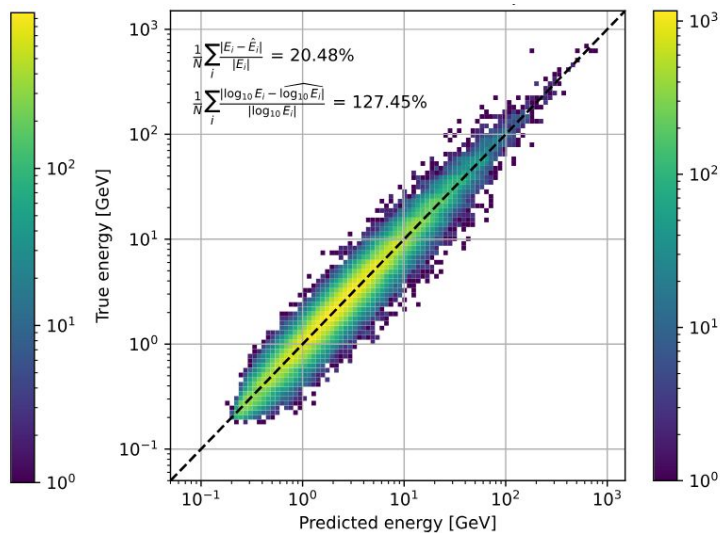
Effect of hit efficiency

Using the
variable-length feature

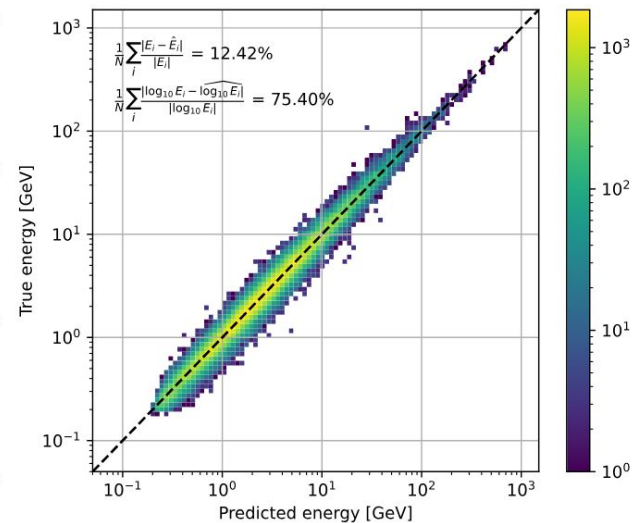
2 missed hits (total = 6)



1 missed hit (total = 7)



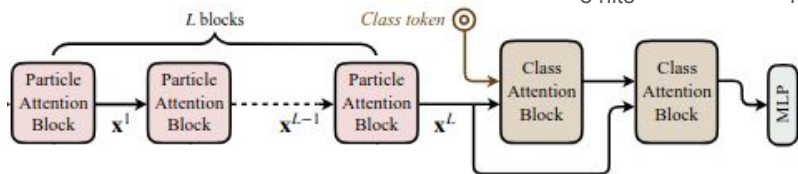
0 missed hit (total = 8)



For 8 detector planes (3 scatterer volumes) :

- $\epsilon = 90\%$: $N_{8 \text{ hits}} = 43\%$, $N_{7 \text{ hits}} = 38\%$, $N_{6 \text{ hits}} = 15\%$
- $\epsilon = 95\%$: $N_{8 \text{ hits}} = 66\%$, $N_{7 \text{ hits}} = 28\%$, $N_{6 \text{ hits}} = 5\%$
- $\epsilon = 99\%$: $N_{8 \text{ hits}} = 92\%$, $N_{7 \text{ hits}} = 7\%$, $N_{6 \text{ hits}} = <1\%$

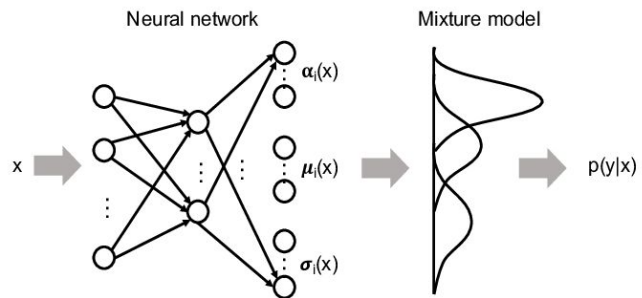
TNN



Real-life requirements

Resolution estimation

Mixture density network (MDN)

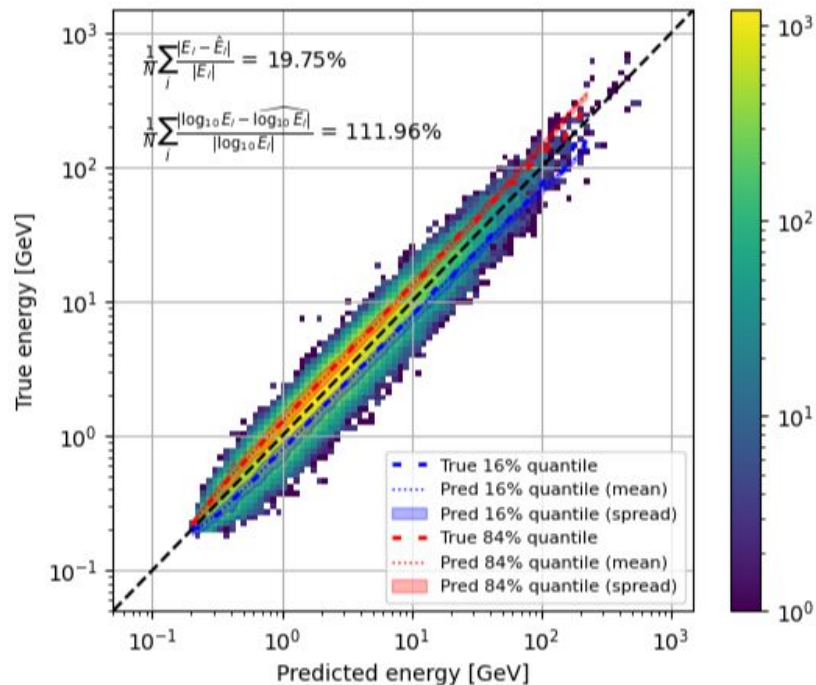


$$p(E|x, \theta) = \sum_i \pi_i \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x-\mu_i)^2}{2\sigma_i^2}}$$

$$\mu_i \equiv \mu_i(x|\theta), \sigma_i \equiv \sigma_i(x|\theta), \pi_i \equiv \pi_i(x|\theta)$$

$$L(t, \{\mu_i, \sigma_i, \pi_i\}) = -\log \left(\sum_i \exp(\log(\pi_i) - \log \sigma_i - \frac{1}{2} \left(\frac{t - \mu_i}{\sigma_i} \right)^2) \right)$$

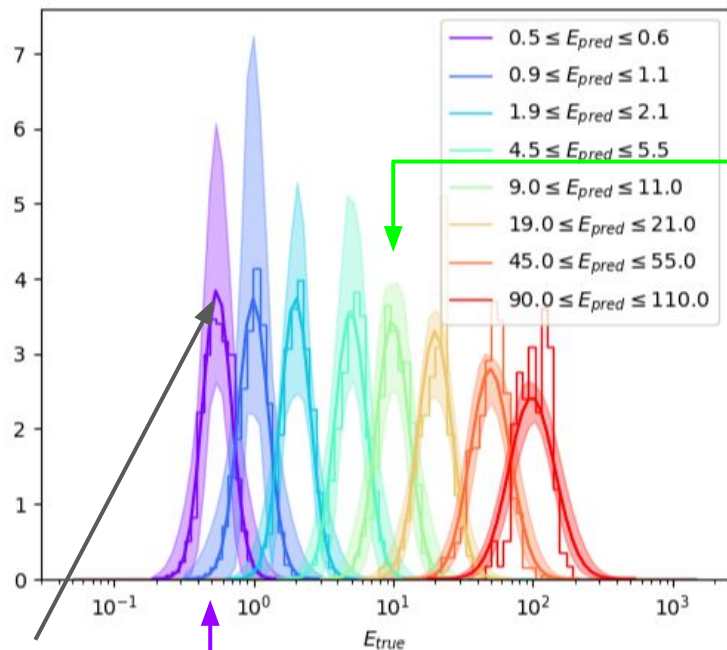
Correctly predicts the 1σ range of E_{true}



Note : can also be done with TNN or GNN as backbone

Real-life requirements

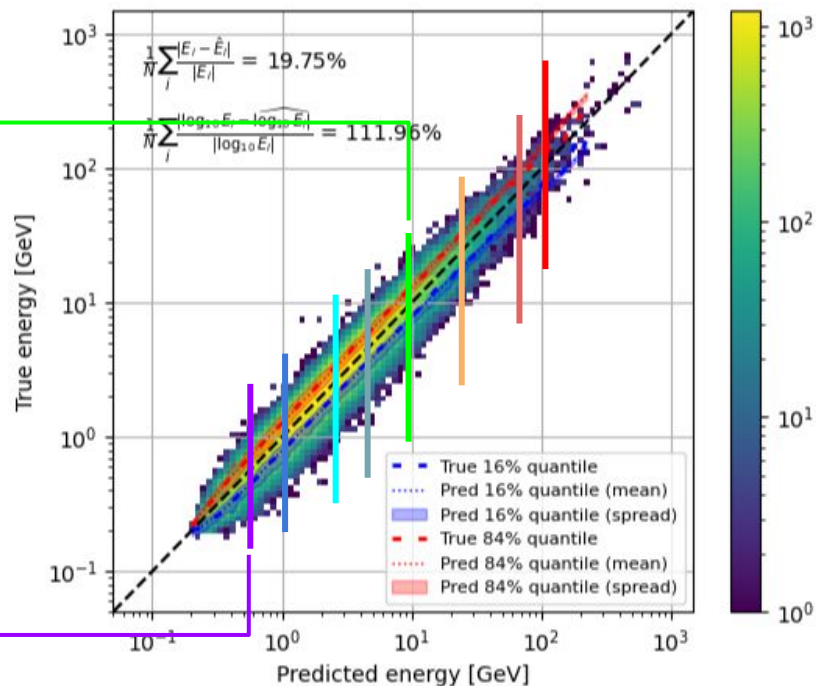
Resolution estimation



MDN mean and spread

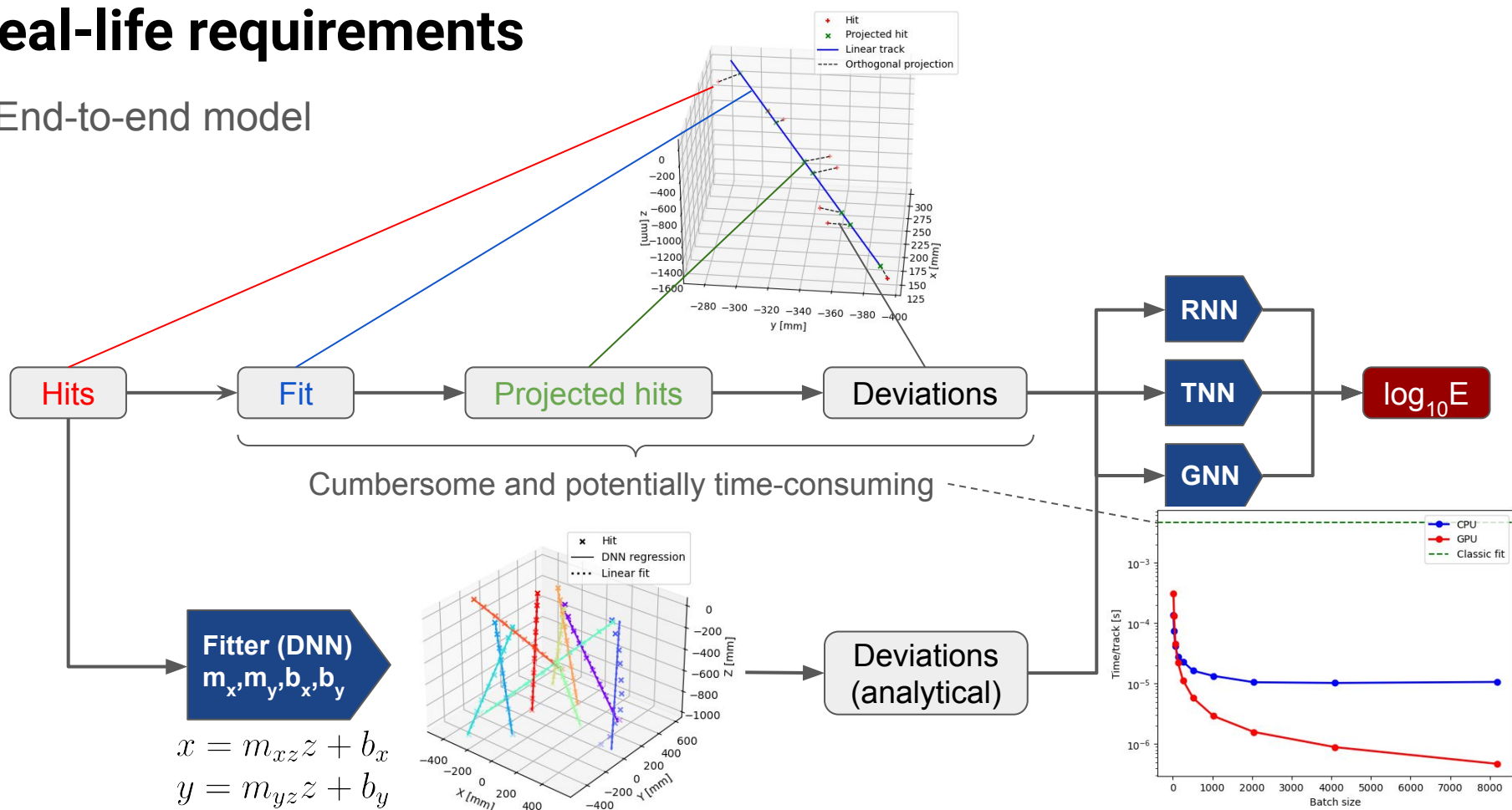
Good reconstruction of $P(E_{\text{true}}|E_{\text{pred}})$

Correctly predicts the 1σ range of E_{true}



Real-life requirements

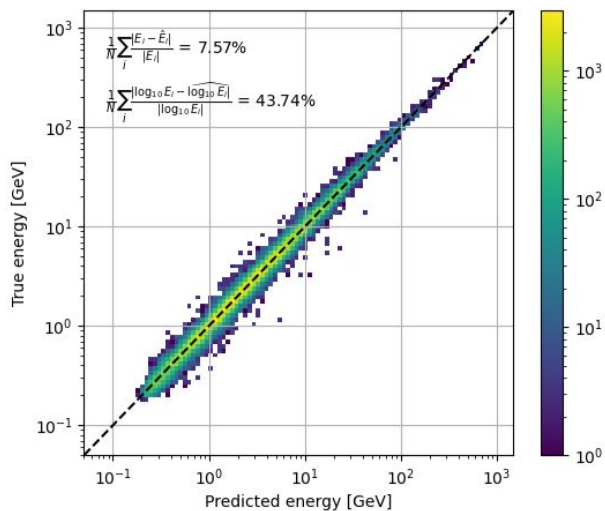
End-to-end model



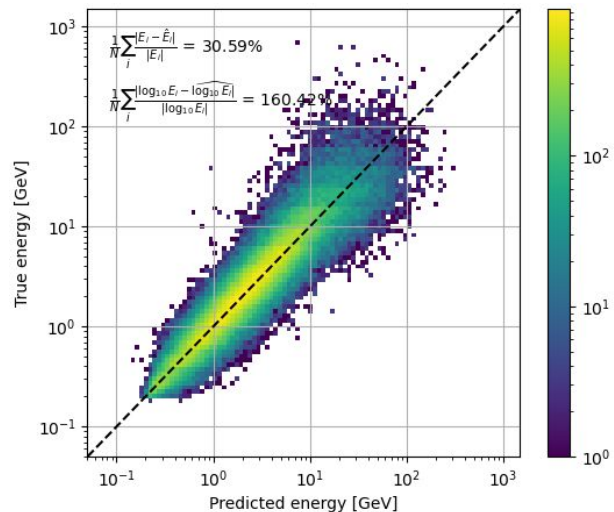
Real-life requirements

End-to-end model

Classic fit
+
TNN (trained on classic fit data)

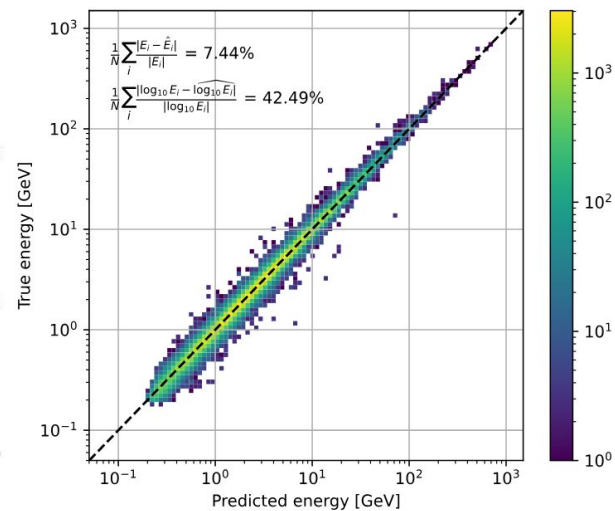


DNN fit
+
TNN (trained on classic fit data)



DNN fitter not perfect
(especially at high energies)

DNN fit
+
TNN (trained on DNN fit data)



Can be recovered by specific training

Conclusion

$$p = \frac{13.6 \text{ MeV}}{\theta_{RMS}} \sqrt{\frac{x}{X_0}}$$

Advanced ML architectures allow extracting more information from the scattering

Detector-specific predictions

Done :

- Hit efficiency
- Hit resolution
- Posterior $P(E_{\text{true}} | E_{\text{pred}})$
- End-to-end model

To-dos :

- Modify E distributions
- Different detector setups
- Optimize architectures

Detector agnostic predictions

- Exploit symmetries more explicitly (contrastive learning)
- Generative pre-training
- Test within TomOpt

Conclusion

Advanced ML architectures allow extracting more information from the scattering

Experimental requirements are met :

- Hit efficiency
- Hit resolution
- Posterior $P(E_{\text{true}}|E_{\text{pred}})$
- End-to-end model

Future checks :

- Modify E spectrum in training
- Hyperparameter scans
- Test different setups (#planes, X_0)

Future developments :

- Include symmetry (contrastive learning)
- Generative pre-training

$$p = \frac{13.6 \text{ MeV}}{\theta_{RMS}} \sqrt{\frac{x}{X_0}}$$

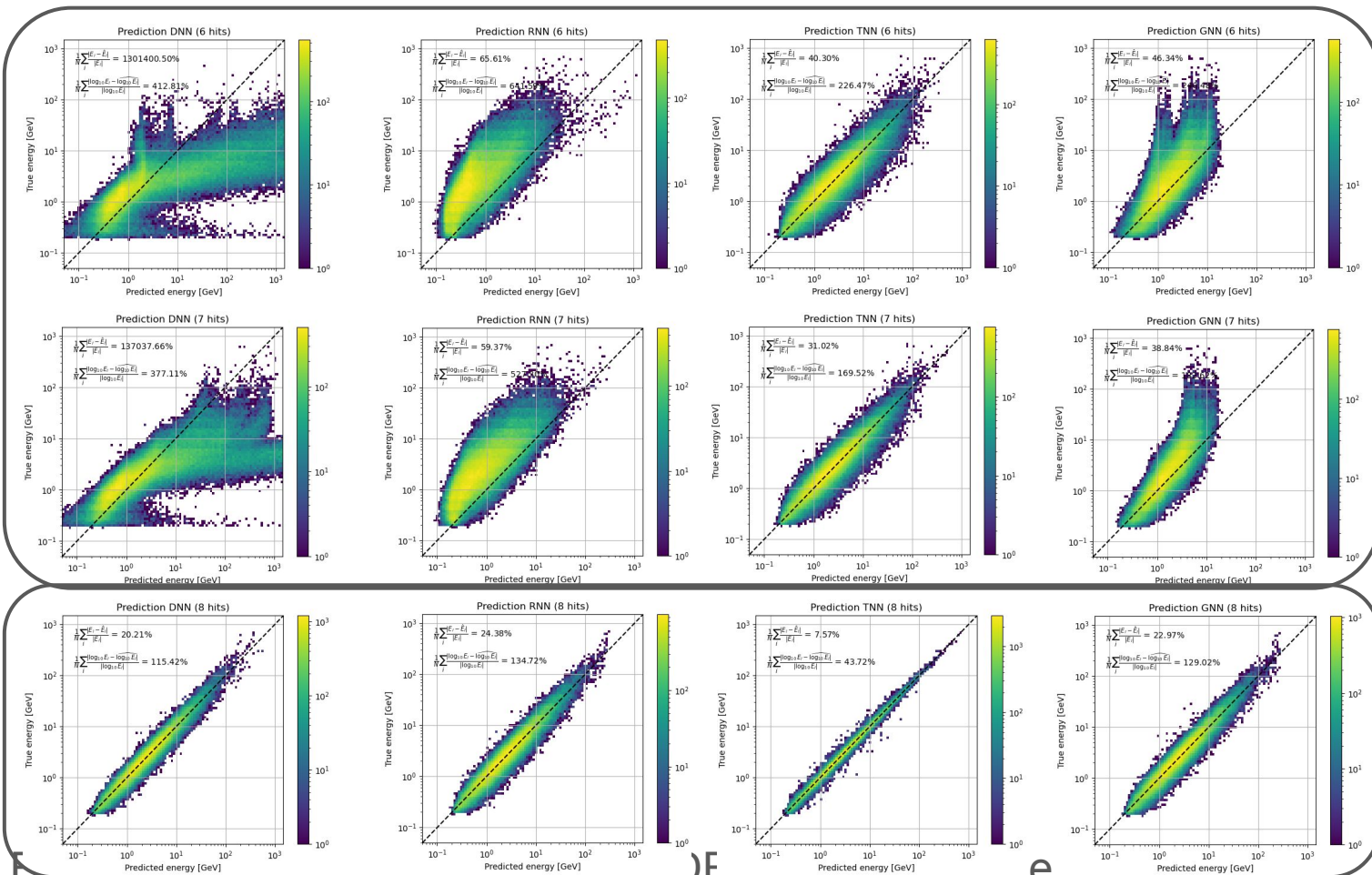
Specific detector

Detector-agnostic predictions
(TomOpt)

Back-up

Effect of hit efficiency

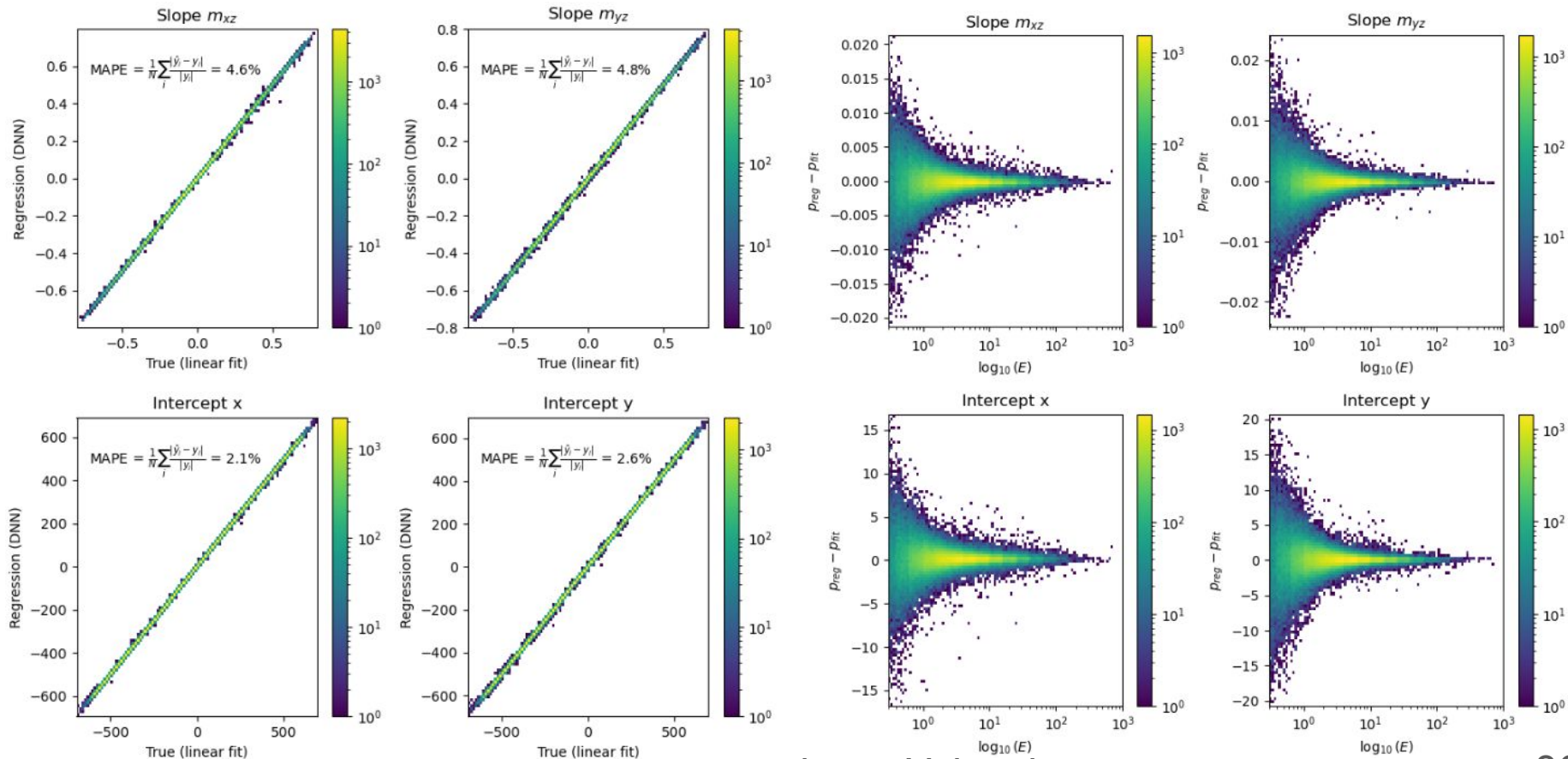
Model only trained on 8 hits, evaluated on 6, 7 and 8 hits



Not seen during training

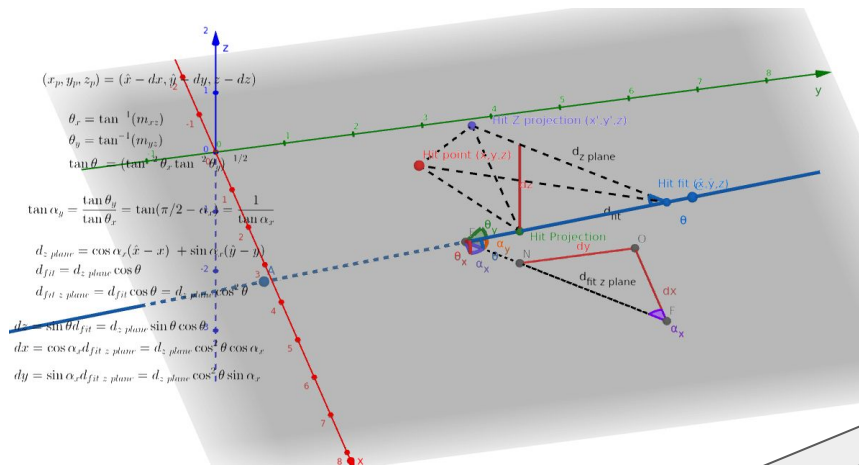
Seen during training

DNN fitter

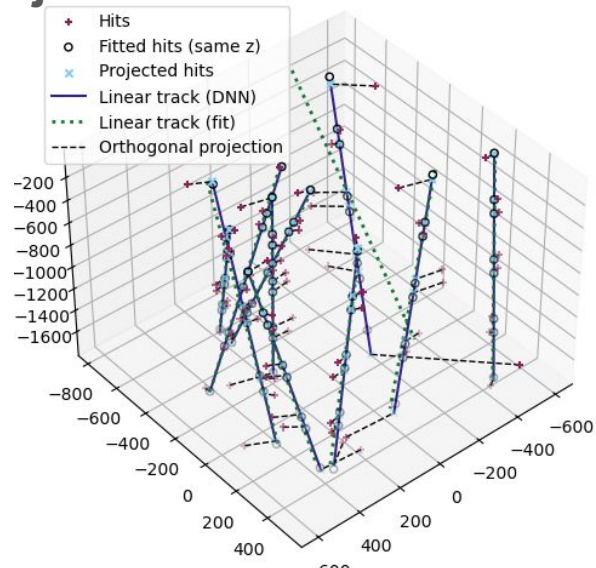


DNN fitter

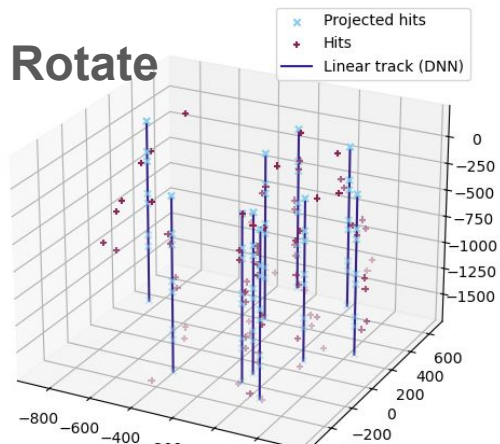
Calculate



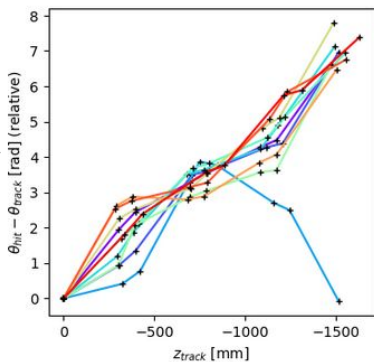
Project



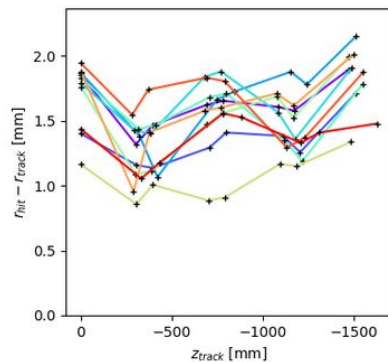
Rotate



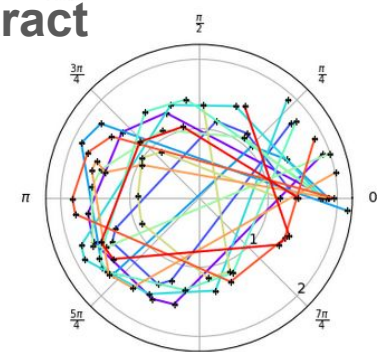
4th order



Deviations from linear track

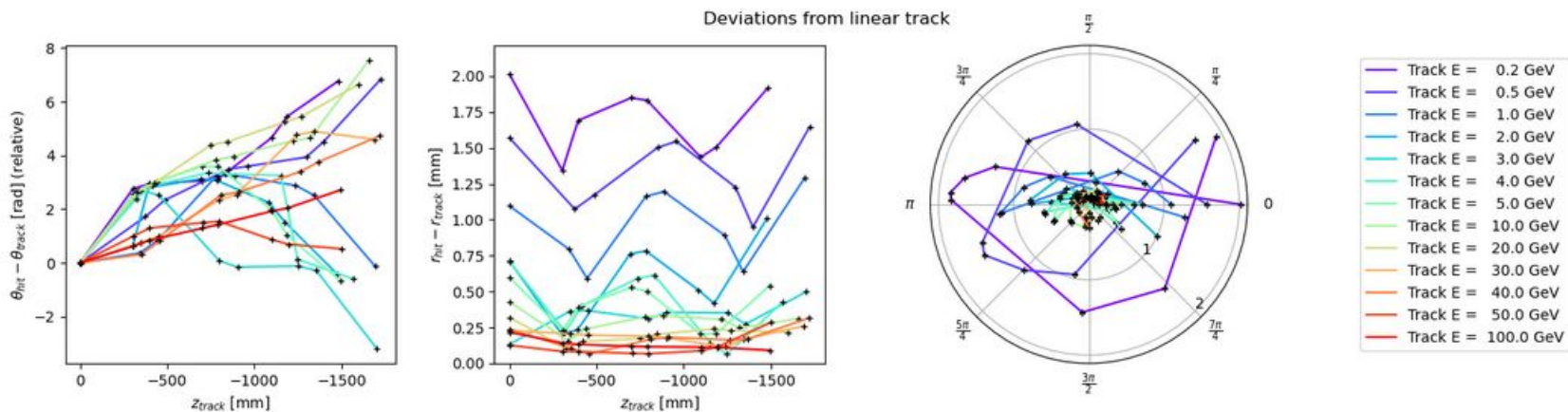


Extract

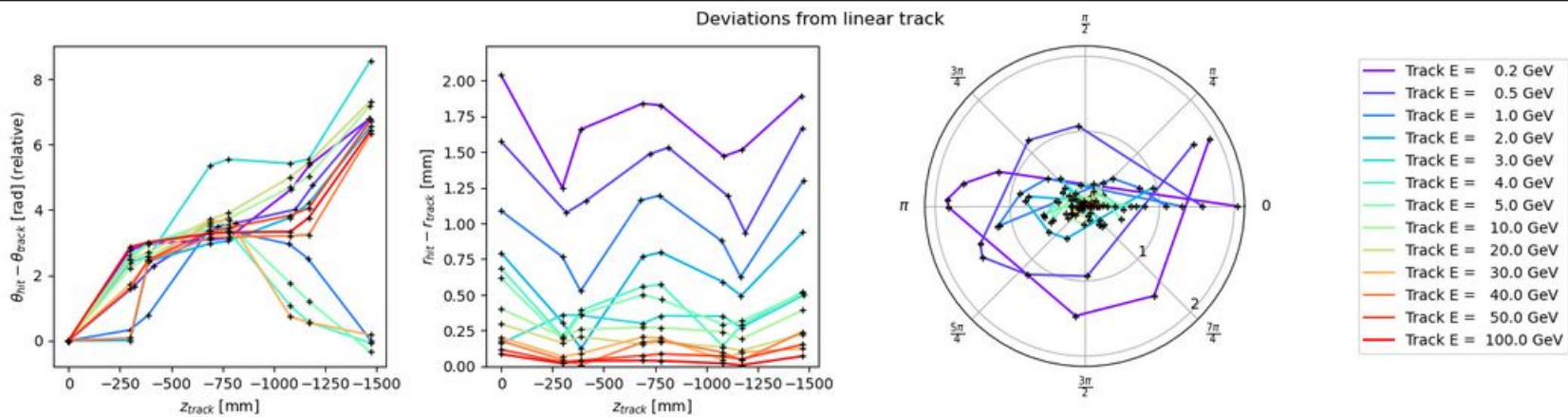


DNN fitter

DNN fitter



Classic fitter



More complex architectures : deviations as a graph

$$\mathbf{x}_i^{(k)} = \gamma^{(k)} \left(\mathbf{x}_i^{(k-1)}, \bigoplus_{j \in \mathcal{N}(i)} \phi^{(k)}(\mathbf{x}_i^{(k-1)}, \mathbf{x}_j^{(k-1)}, \mathbf{e}_{ij}) \right)$$

Many models on the market (see [pytorch-geometric](#)), the following are considered :

- Convolutional-based : in feature space \mathbf{x} (in my case $\mathbf{x}_i = \text{deviations } \{\Delta x_i, \Delta y_i, z_i\}$)
 - GCN [1] (2016) : $\oplus = \text{mean}$ and $\phi = \text{linear layer}$
 - GAT(v2) [2][3] (2017, 2021) : $\oplus = \text{mean (concat per head)}$ and $\phi = \text{attention mechanism}$
 - GraphSAGE [4] (2017) : $\oplus = \text{max}$, $\phi = \text{linear layer} + \text{ReLU}$ and $\gamma = \text{linear layer on concat}$
 - EdgeConv [5] (2018) : $\oplus = \text{mean/sum/max/...}$ and $\phi = \text{NN applied on concat}(\mathbf{x}_i, \mathbf{x}_i - \mathbf{x}_j)$
 - GIN [6] (2018) : $\oplus = \text{sum}$ and $\gamma = \text{NN on (variable) weighted sum of } \mathbf{x}_i$ and the aggregation
 - PNA [7] (2020) : $\oplus = \text{several aggregations combined and scaled}$, $\gamma, \phi = \text{NN}$
- Point-based : depends on positions \mathbf{p}_i (with potential addition of node features \mathbf{x}_i)
 - PointNet [8] (2016) : $\oplus = \text{max}$, $\gamma = \text{NN}$ and $\phi = \text{NN applied on concat}(\mathbf{x}_i, \mathbf{p}_i - \mathbf{p}_j)$
 - PointTransformer [9] (2020) : more complex attention mechanism based on \mathbf{x}_i and $\mathbf{p}_i - \mathbf{p}_j$
 - PointGCN [10] (2020) : more complex with 3 NNs based on \mathbf{x}_i and $\mathbf{p}_i - \mathbf{p}_j$