

Pre-thesis internship

# Study of hadronic interactions in calorimeter with high granularity

The author :  
Yaroslav Nikolaiko

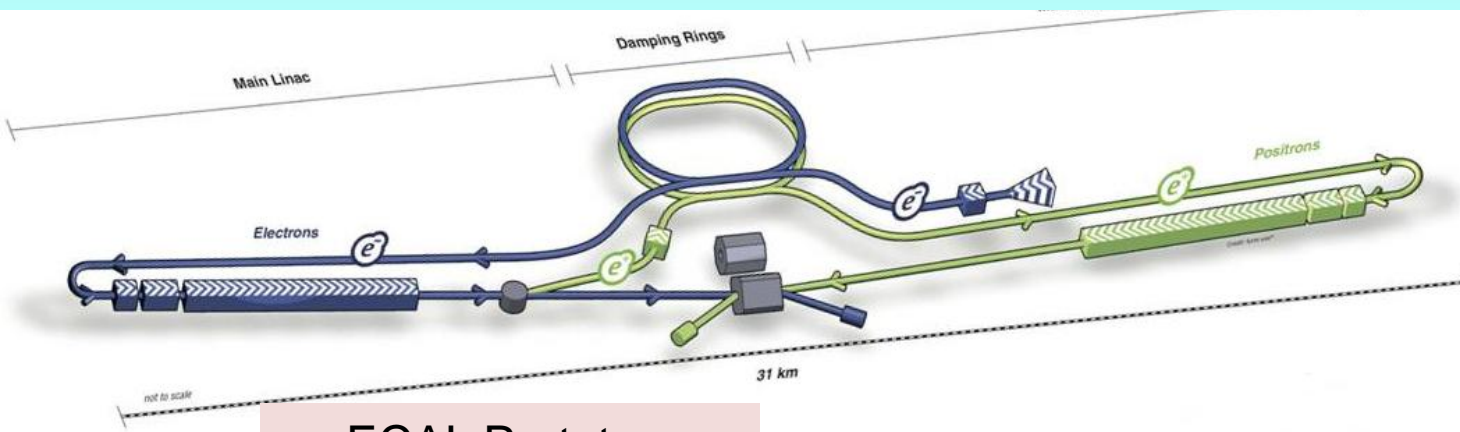
Supervision :  
**Roman Poeschl** and **Naomi van der Kolk**,  
LAL, ILC group-CNRS/IN2P3

# Content

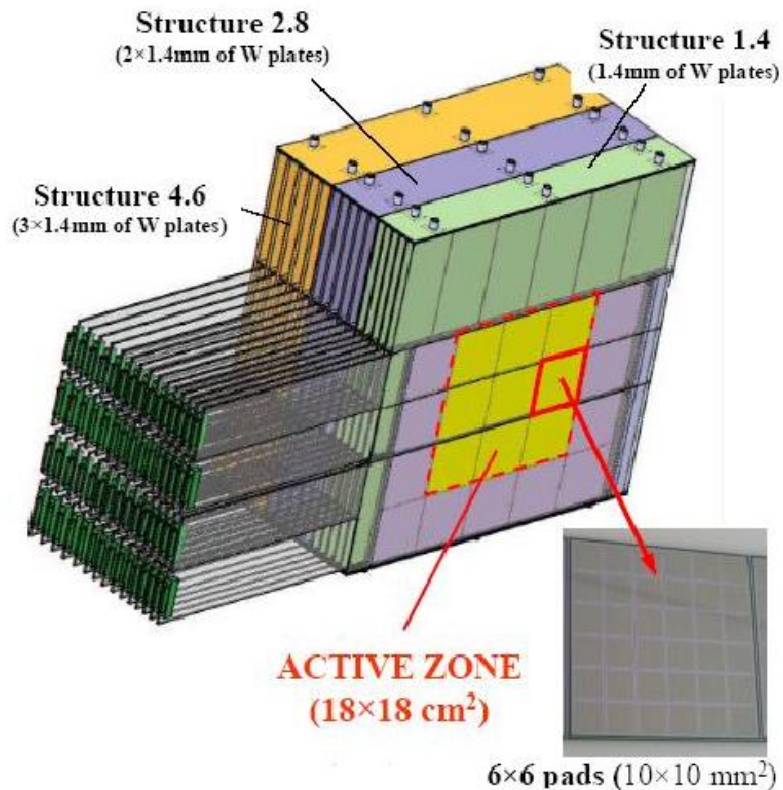
- ◆ Introduction : Future Linear Collider and Particle Flow
- ◆ Machine learning as tool for PF Algorithms improvement
- ◆ Hadronic Interaction in Geant4
- ◆ Trivial classification of interactions
- ◆ Summary

# Introduction

ILC  
0.5-1 TeV  $e^+e^-$  collider



## ECAL Prototype




- Precise study of SM parameters and beyond.
- Requires excellent energy resolution to separate multiple jet events.


- 30 layers sandwich structure (Si – active, W- absorber materials)
- Total depth –  $24 X_0$  or  $1 \lambda_i$
- Followed by HCAL with 48 longitudinal samples

# Particle Flow

# Particle Flow

$$\text{Energy resolution} = \text{Software} + \text{Hardware}$$

  
**Confusion term**  
 $\approx 2\%$

  
Calorimeter resolution  
 $\approx 2\%$

# Particle Flow

Energy resolution

= Software + Hardware

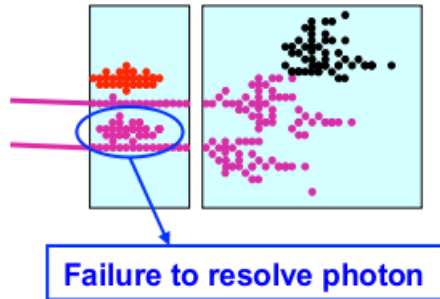
**Confusion term**

$\approx 2\%$

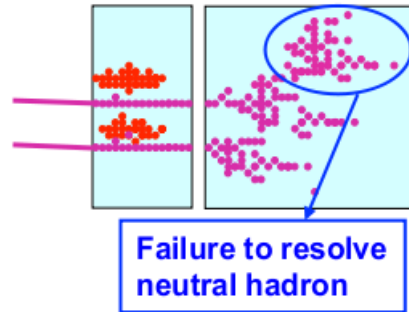
Calorimeter resolution

$\approx 2\%$

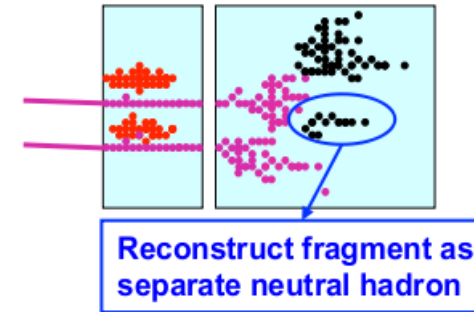
i) Photons



ii) Neutral Hadrons



iii) Fragments



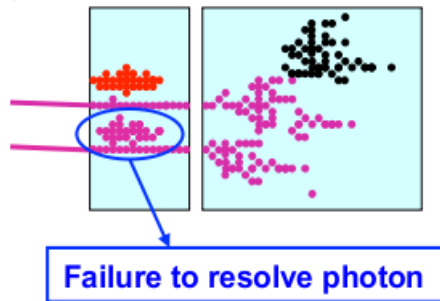
# Particle Flow

Energy resolution = Software + Hardware

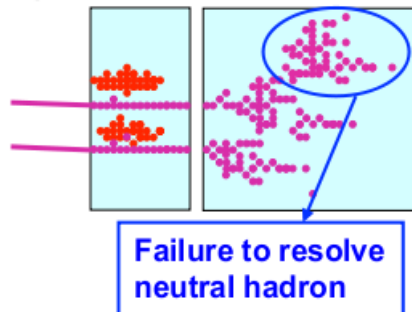
$\downarrow$   
**Confusion term**  
 $\approx 2\%$   
 $\downarrow$

$\downarrow$   
 Calorimeter resolution  
 $\approx 2\%$

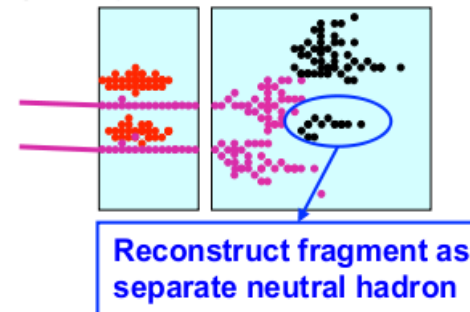
i) Photons



ii) Neutral Hadrons

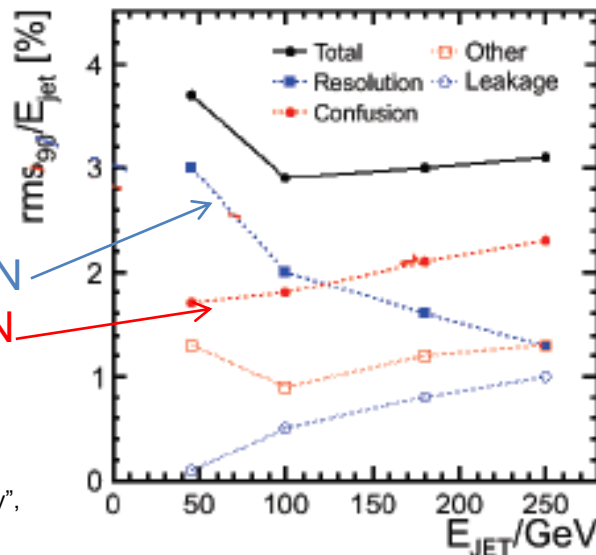


iii) Fragments



## ❖ Contribution to resolution

- Low energy jets: RESOLUTION
- High energy jets: CONFUSION



PandoraPFA performance

$$\sigma_E / E_j \approx 4\% (45\_GeV)$$

Our goal – remove confusion term

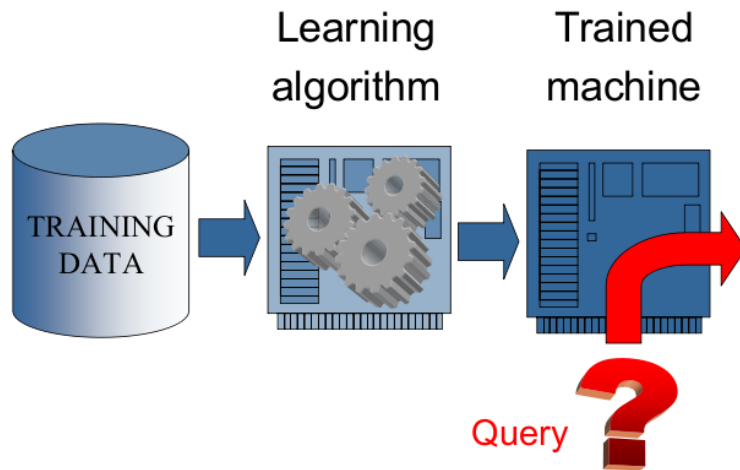


$$\sigma_E / E_j \approx 2\%$$

from M. Thomson, "Particle Flow Calorimetry", Mainz, February 2013

# Supervised and Semi-Supervised Machine Learning

## Machine learning concept



The goal of supervised learning is to infer a function  $g : \mathcal{X} \rightarrow \mathcal{Y}$  from a data set

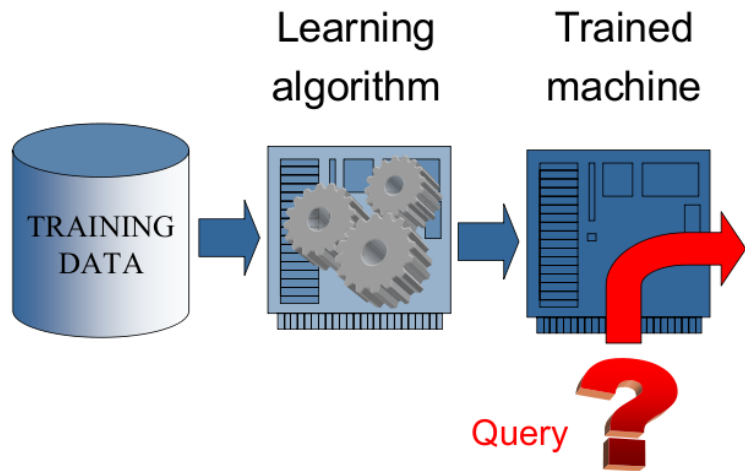
$$\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$$

The quality of  $g$  on an arbitrary pair  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is measured by an *error* or *loss* function  $L(g, (x, y))$



# Supervised and Semi-Supervised Machine Learning

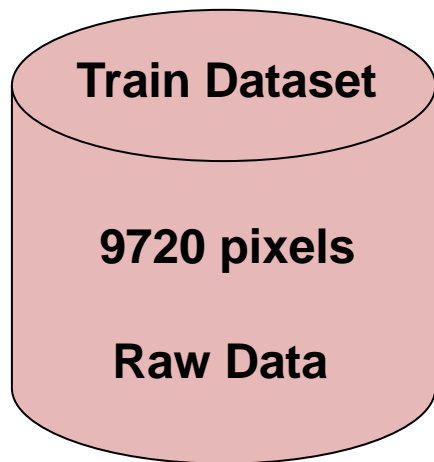
## Machine learning concept



The goal of supervised learning is to infer a function  $g : \mathcal{X} \rightarrow \mathcal{Y}$  from a data set

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$$

The quality of  $g$  on an arbitrary pair  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is measured by an *error* or *loss* function  $L(g, (x, y))$

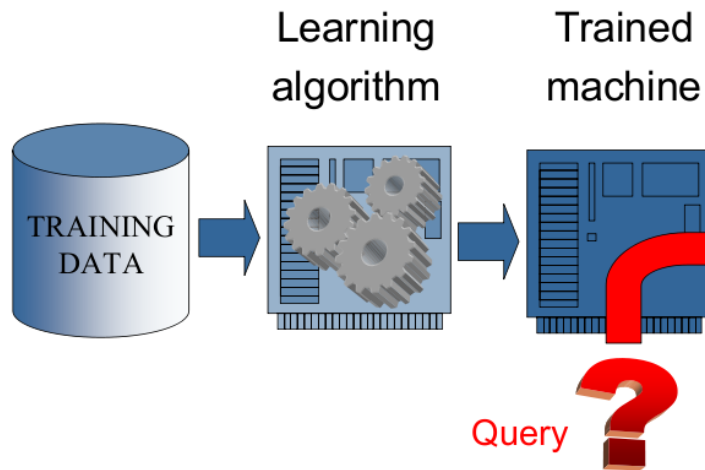


# Supervised and Semi-Supervised Machine Learning

## Machine learning concept

The goal of supervised learning is to infer a function  $g : \mathcal{X} \rightarrow \mathcal{Y}$  from a data set

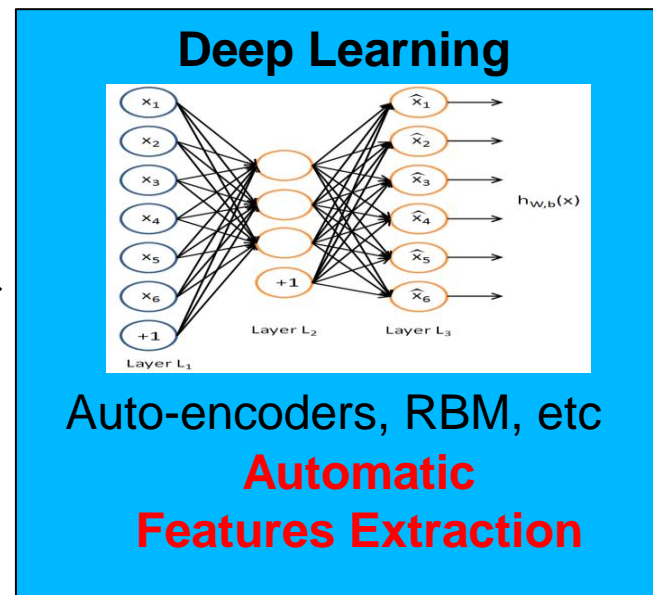
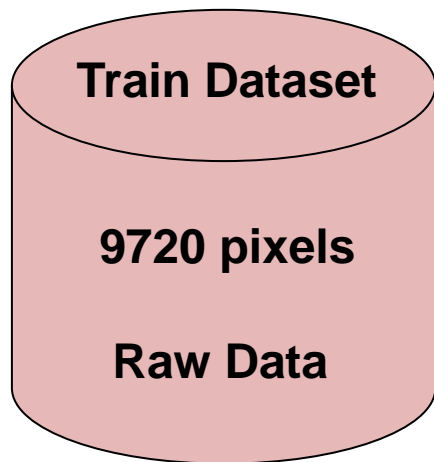
$$\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$$



Answer

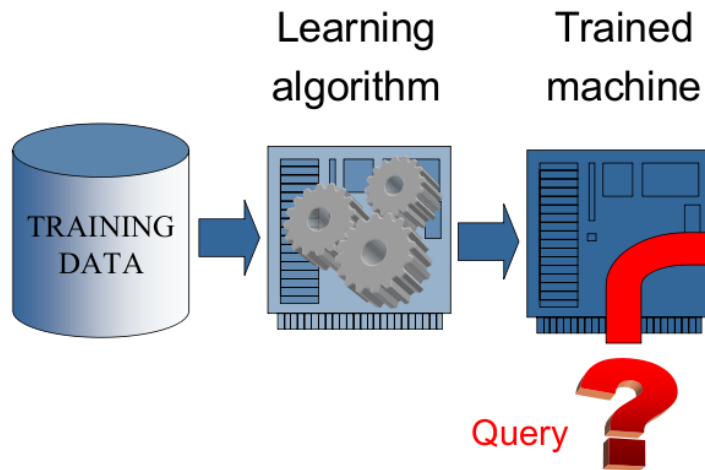
The quality of  $g$  on an

arbitrary pair  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is measured by an *error* or *loss* function  $L(g, (x, \bar{y}))$



# Supervised and Semi-Supervised Machine Learning

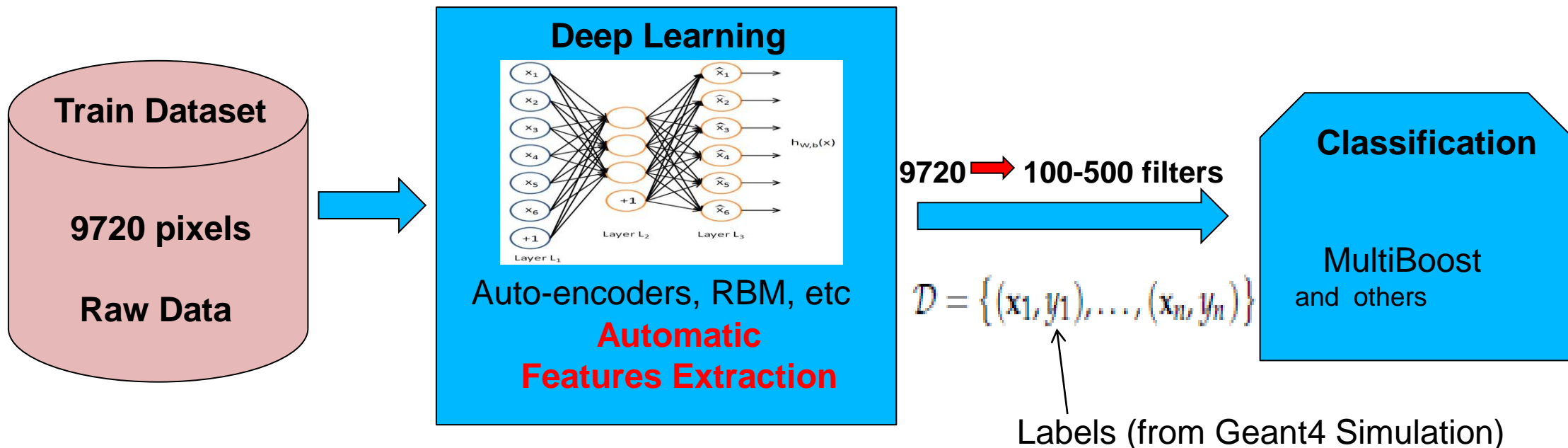
## Machine learning concept



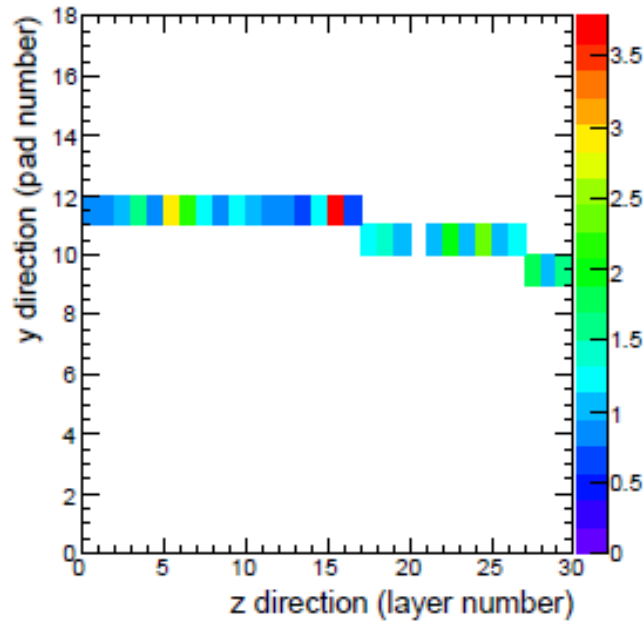
The goal of supervised learning is to infer a function  $g : \mathcal{X} \rightarrow \mathcal{Y}$  from a data set

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$$

The quality of  $g$  on an arbitrary pair  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  is measured by an *error* or *loss* function  $L(g, (x, \bar{y}))$



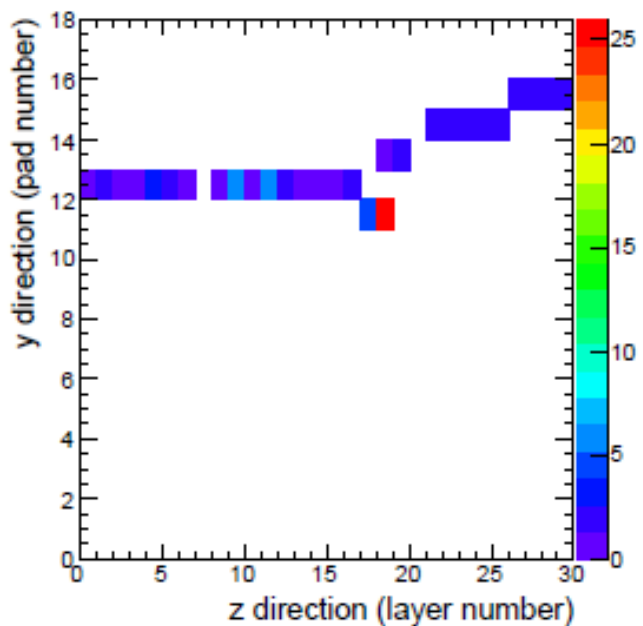
# Test-Beam Events Interaction Classification



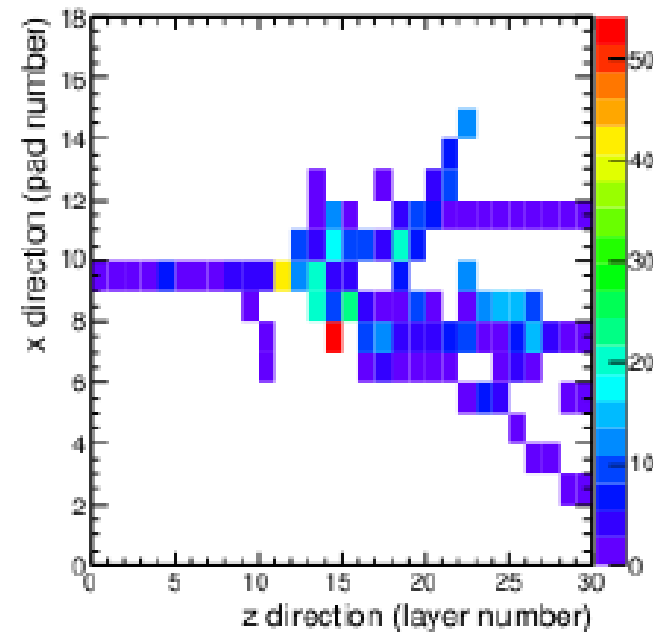
“Elastic”

Performance of ECAL physical prototype

- Pions with energy 2-10 GeV



“Point-like”



“FireBall”

# Geant4(simulation toolkit) Set-Up

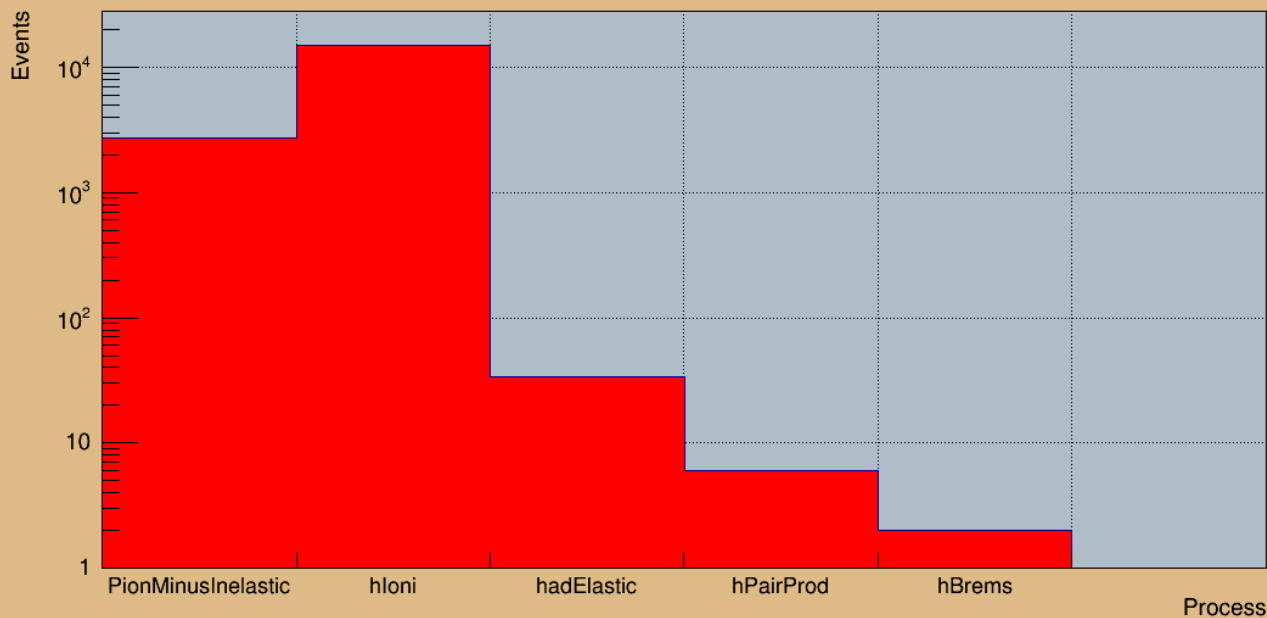
Ionization  
(MIP)

Elastic

Inelastic

Pair  
Production

Hadronic processes for pions in Geant4



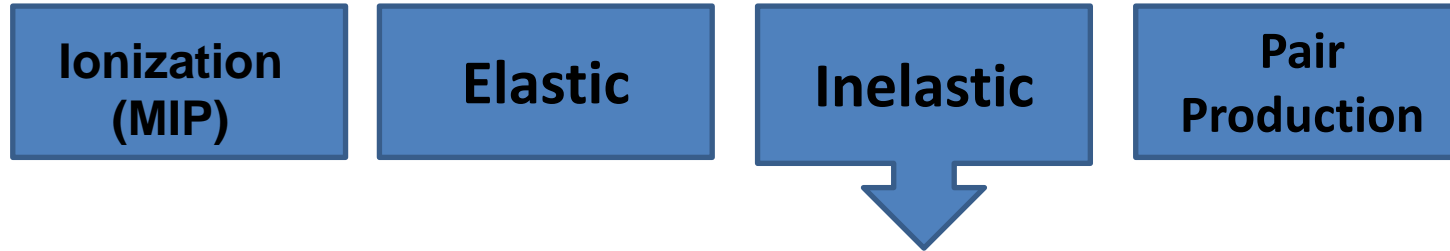
➤  $\pi^-$  10 GeV

➤ QGSP\_BERT physics list

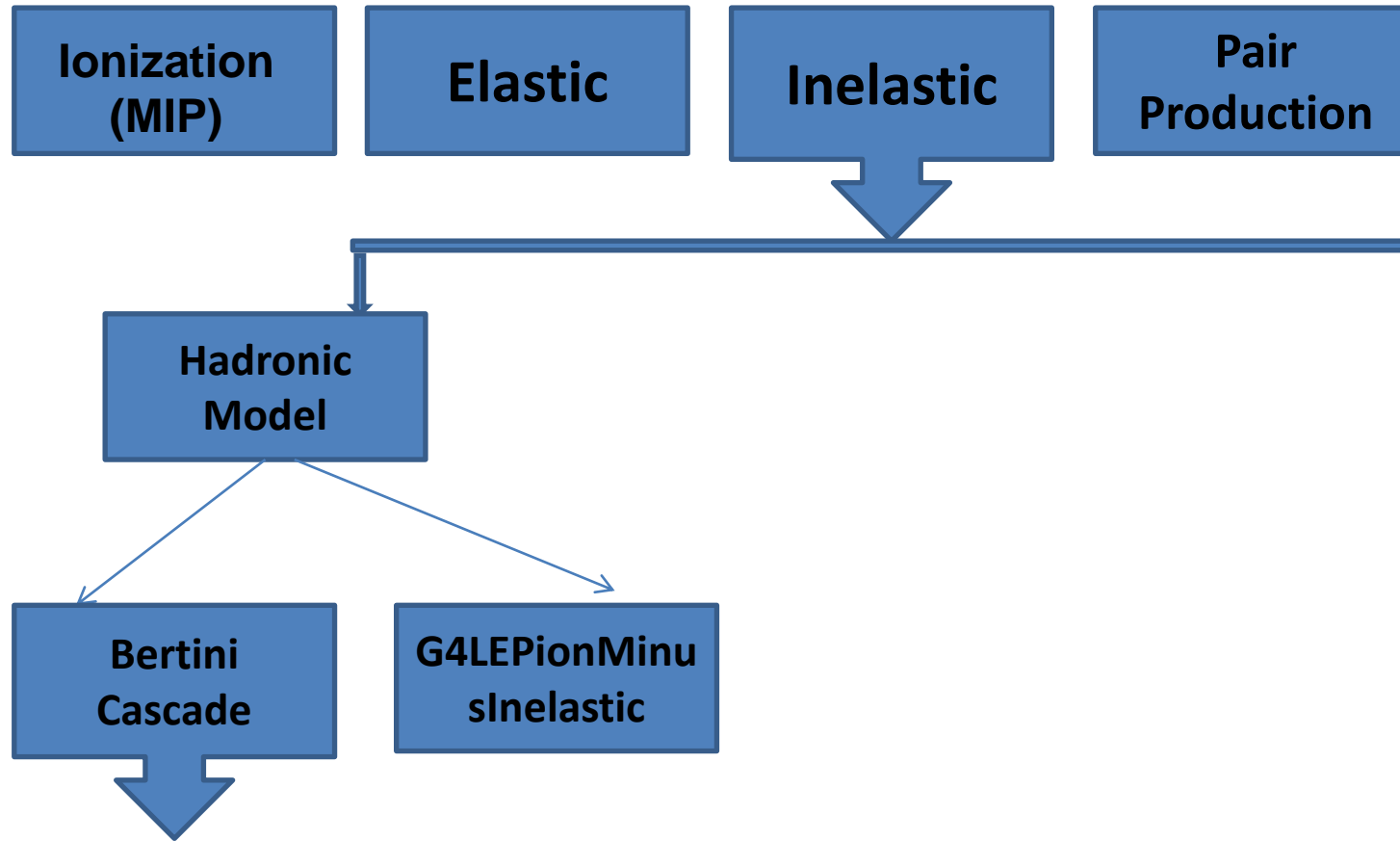
- ✓ quark-gluon string model of protons, neutrons, kaons and nuclei
- ✓ Bertini cascade + parameterized model for low energy

➤ Full Detector geometry simulation (Mokka software)

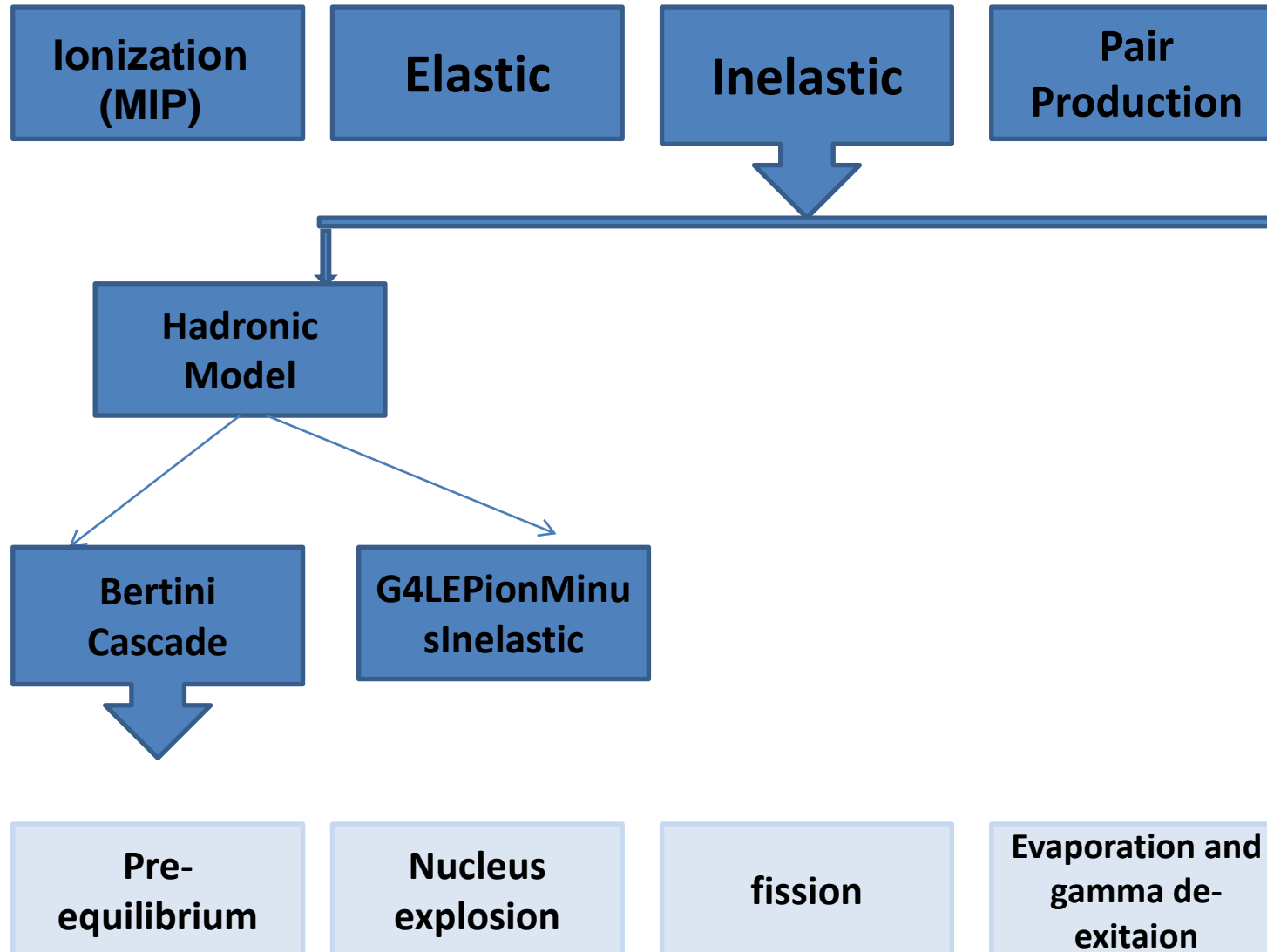
# Geant4 Set-Up



# Geant4 Set-Up

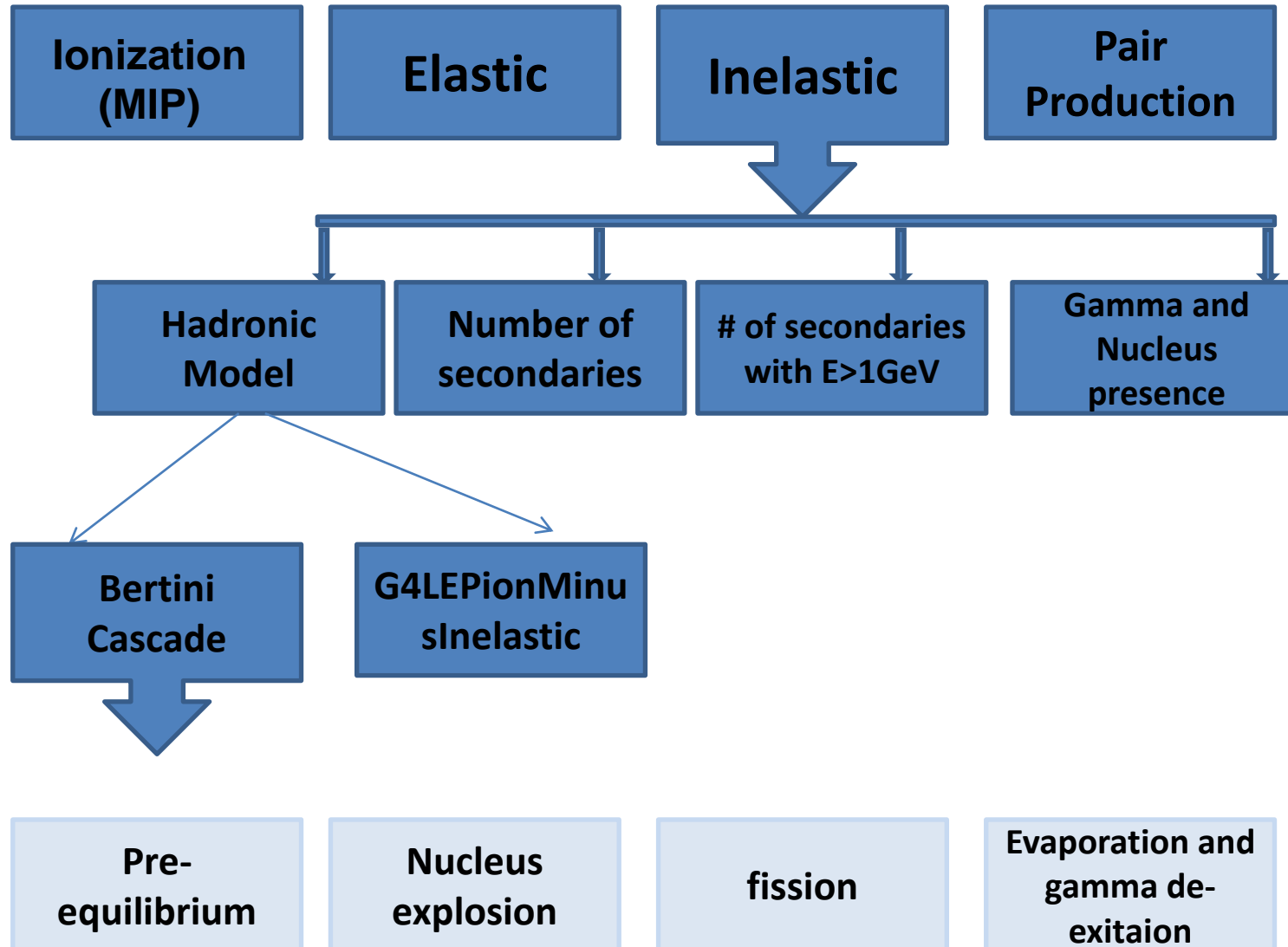


# Geant4 Set-Up

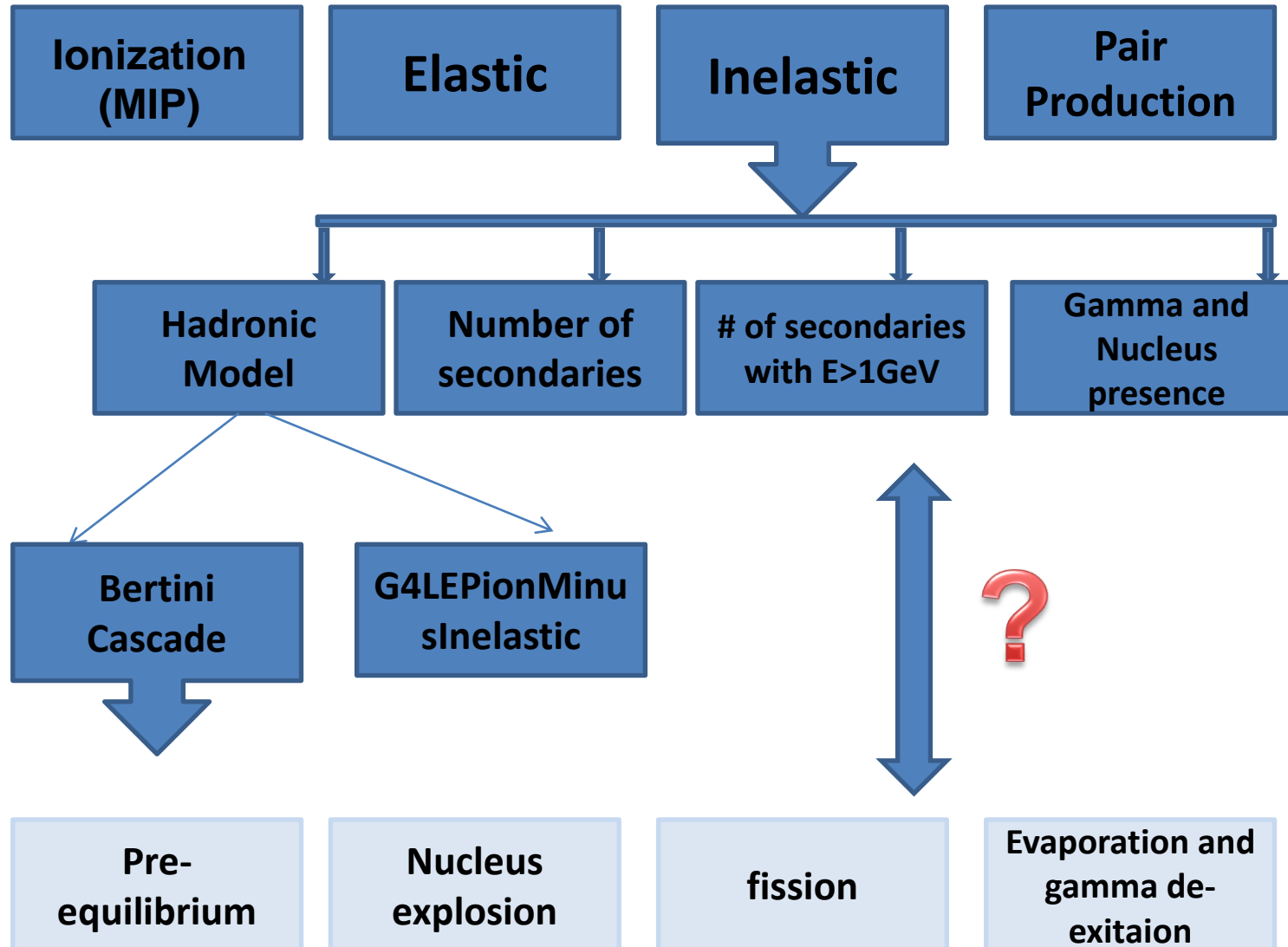




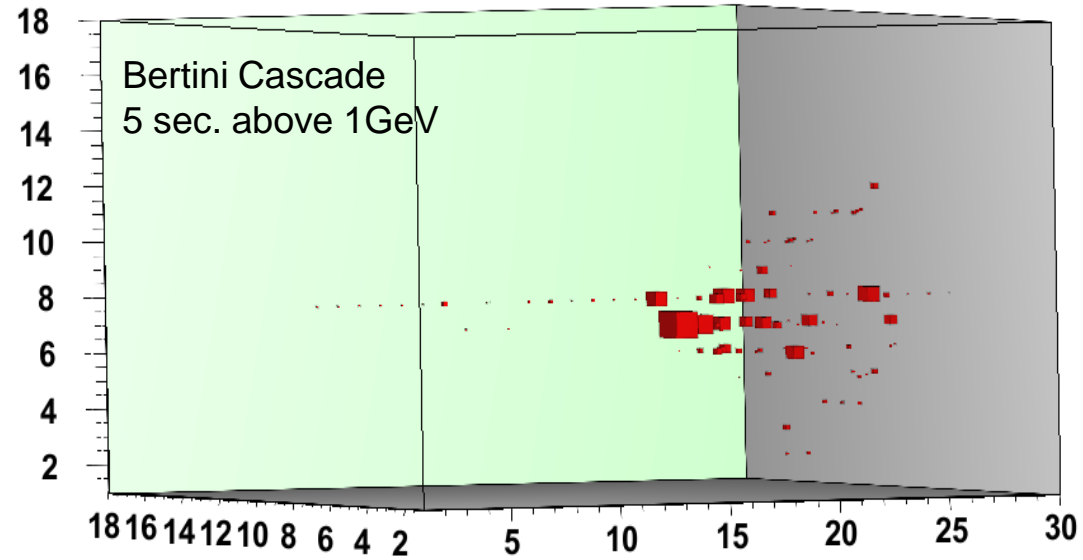
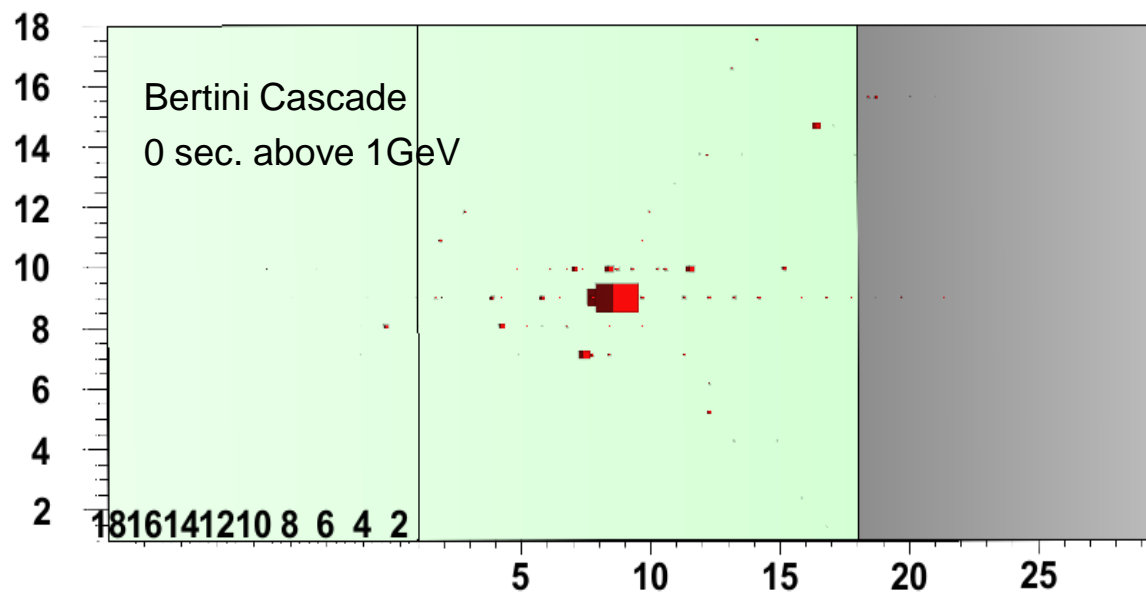
# Geant4 Set-Up



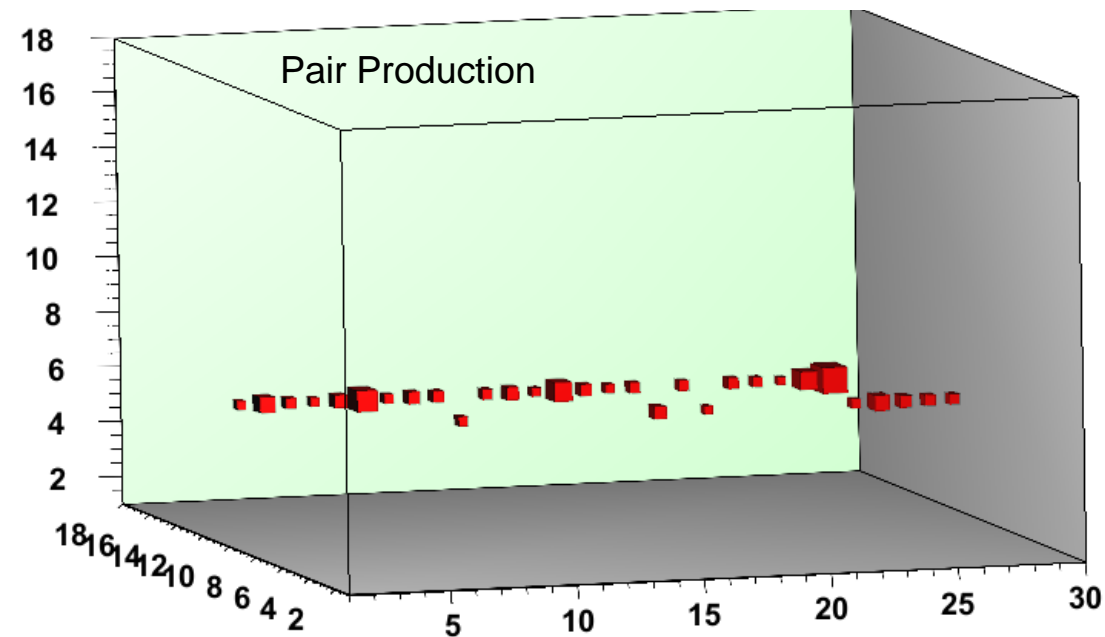
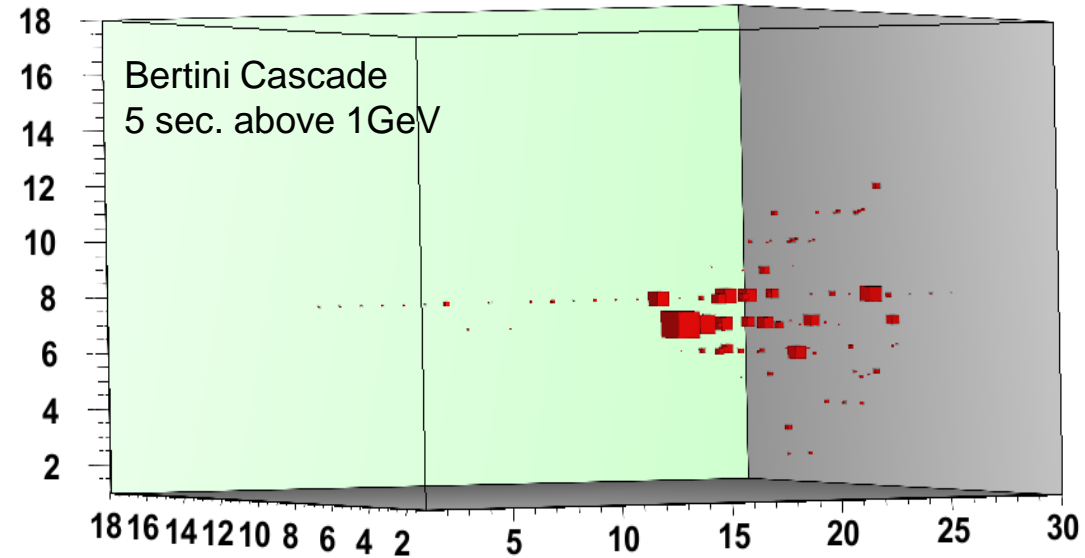
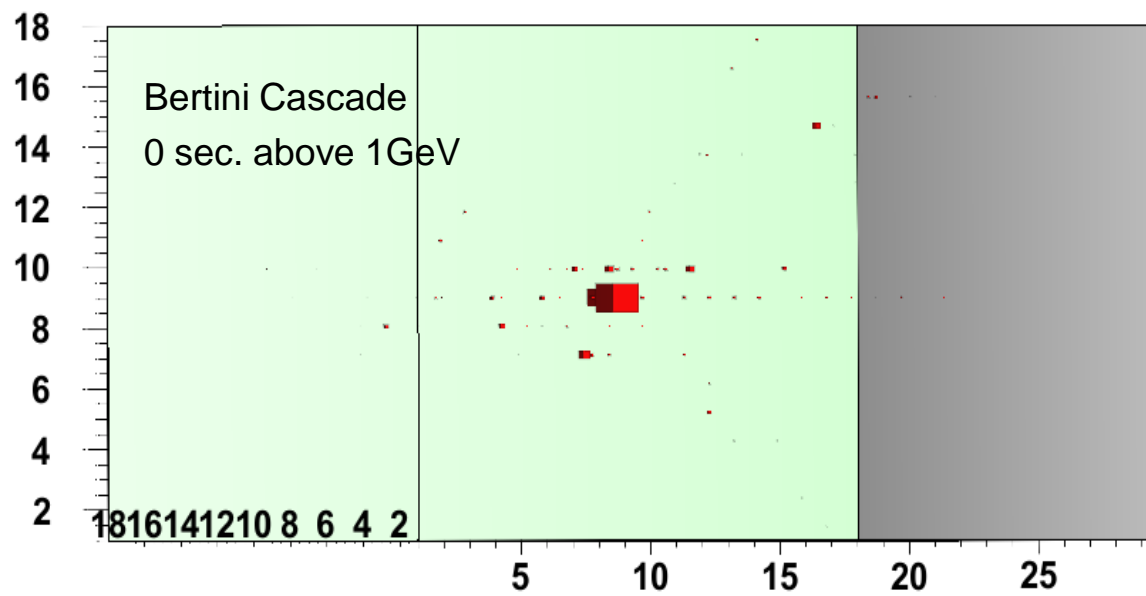
# Geant4 Set-Up



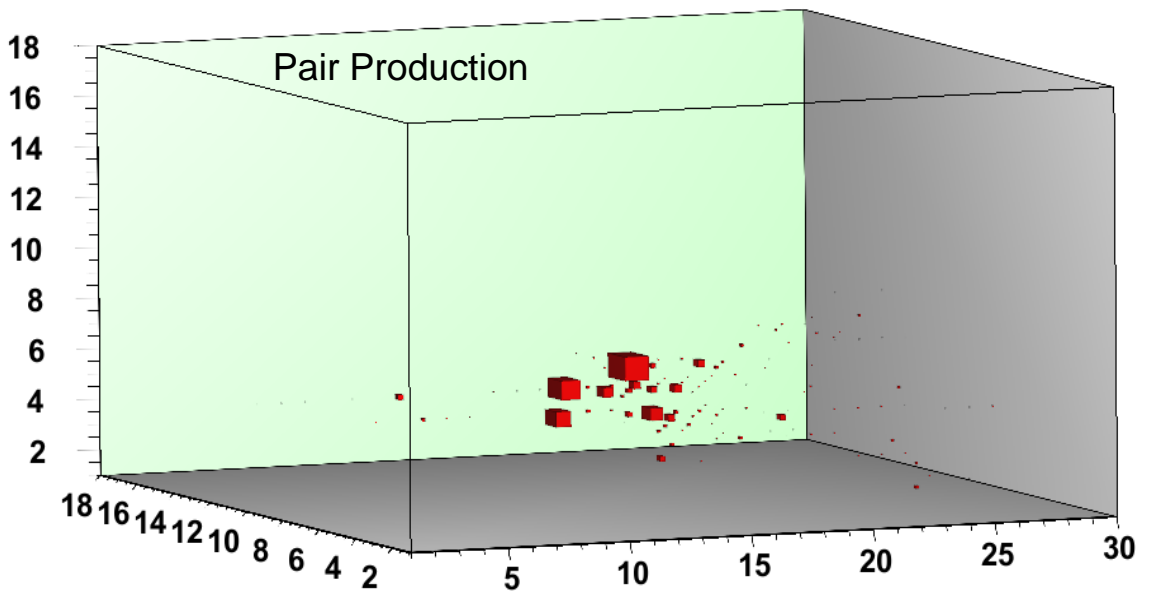
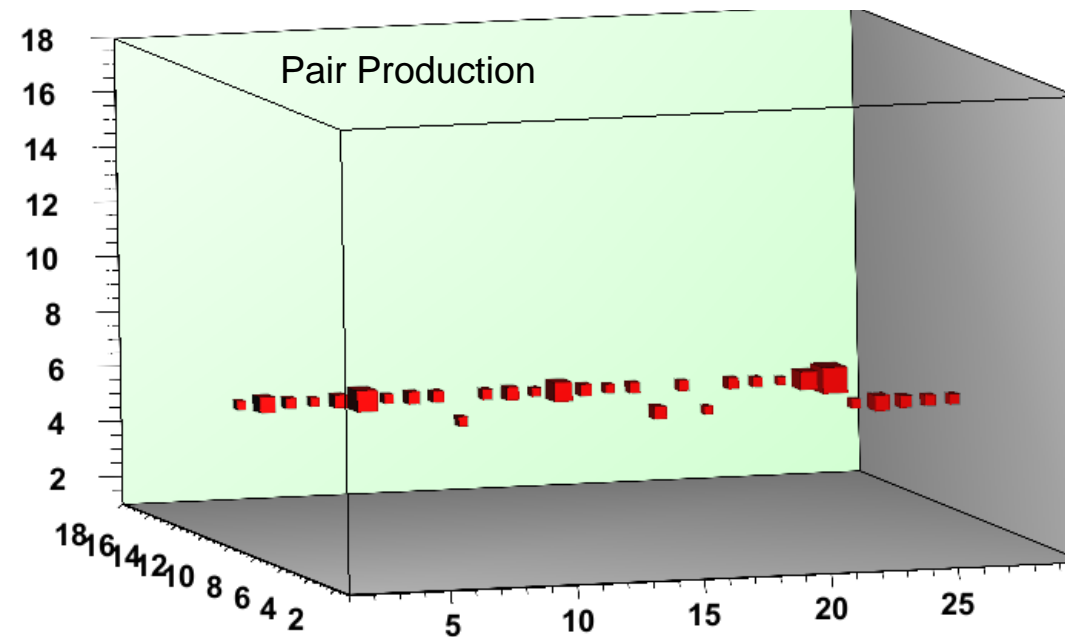
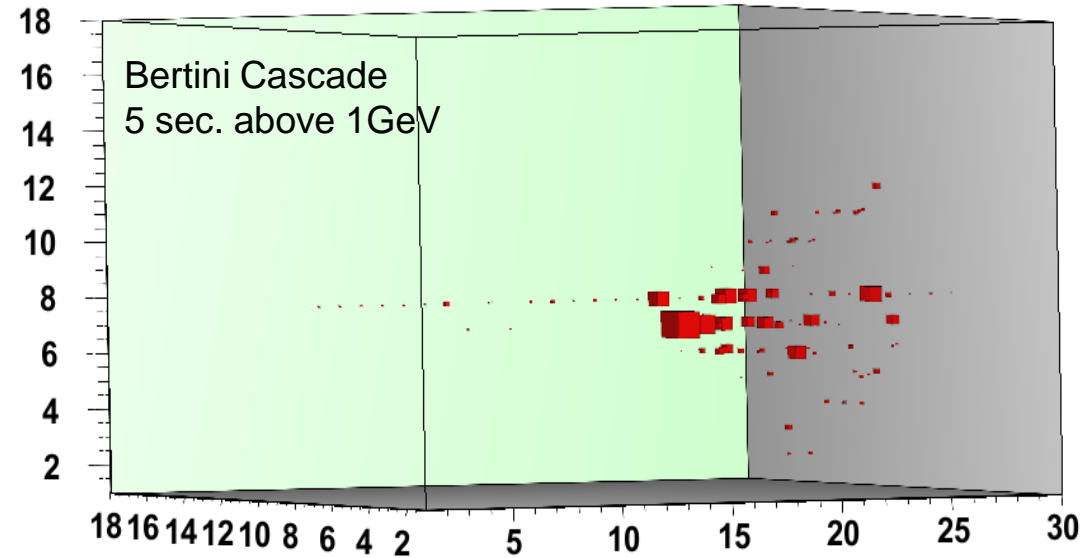
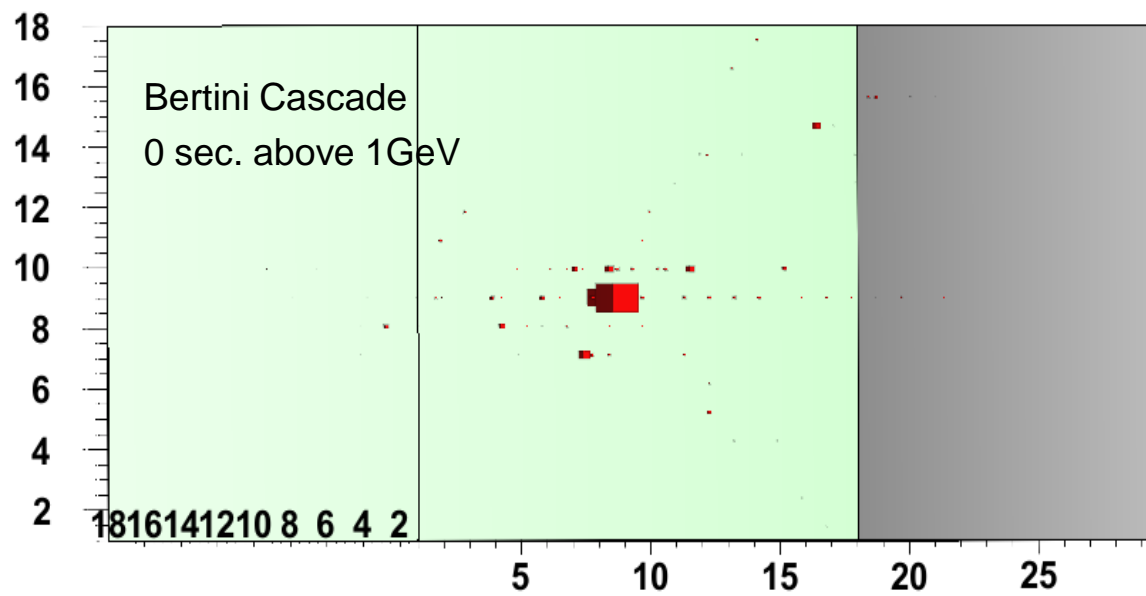
# Events examples



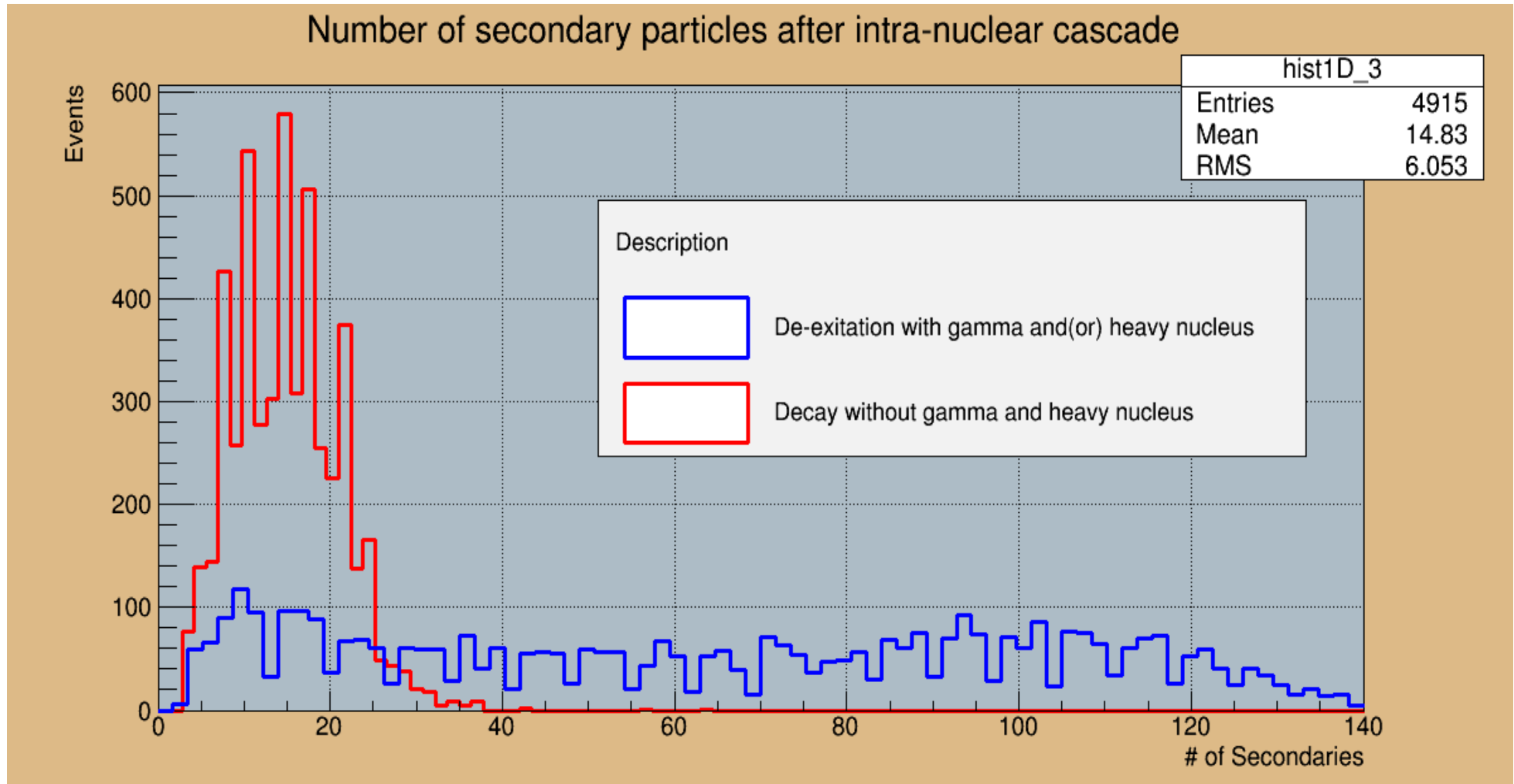
# Events examples



# Events examples

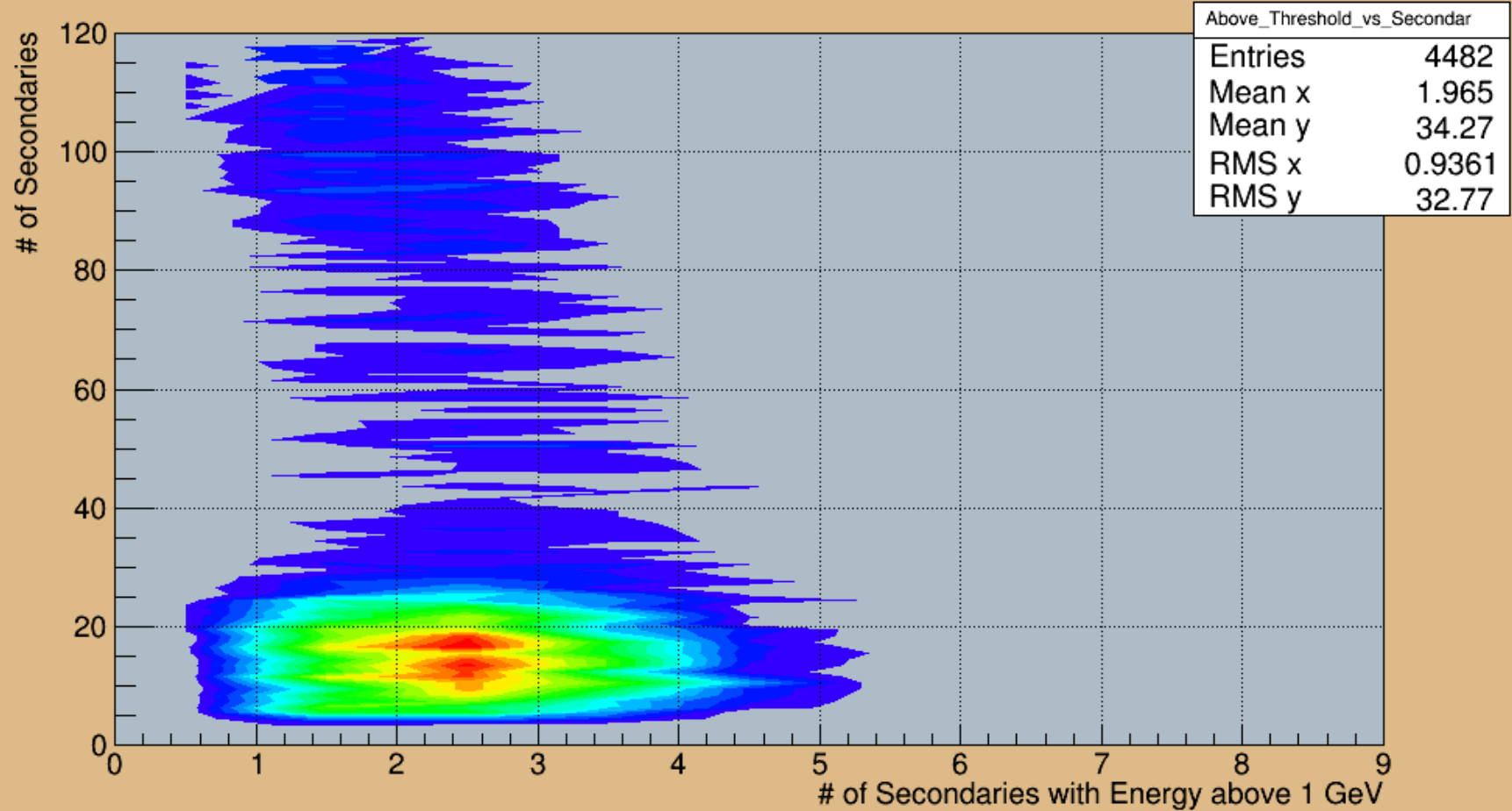


# Relations between labels



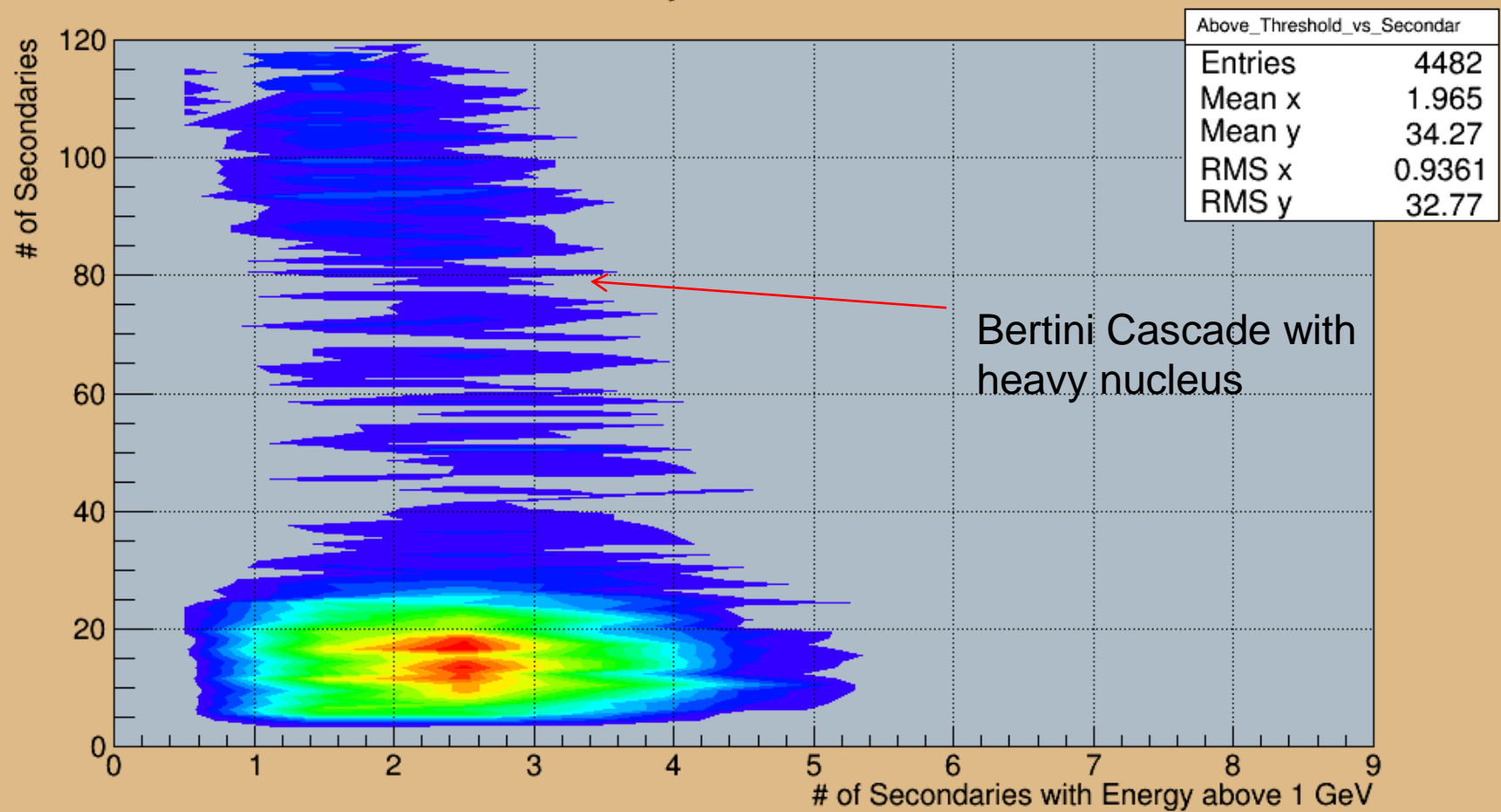
# Relations between labels

Number of secondary with  $E > 1$  GeV distribution



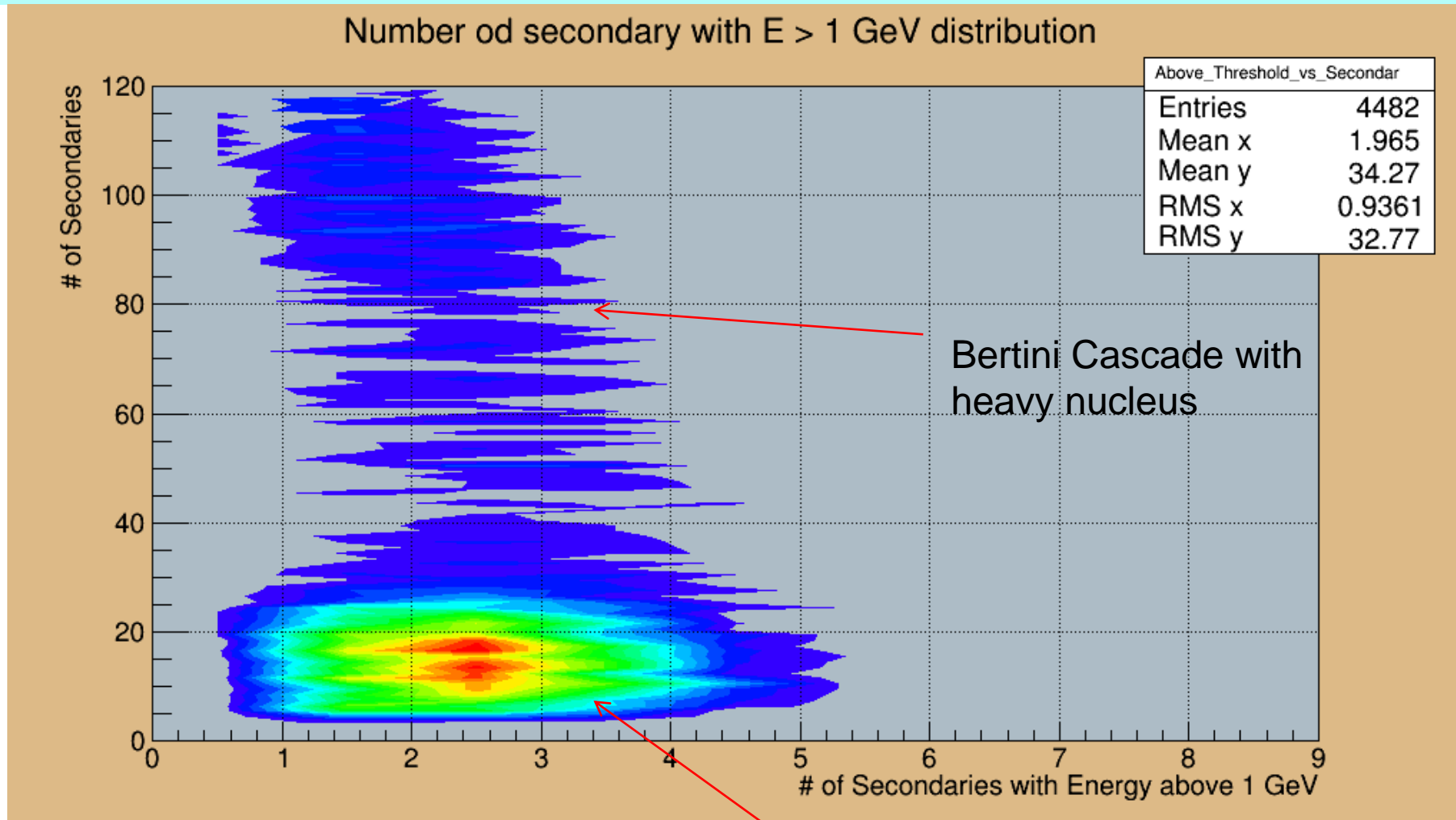
# Relations between labels

Number of secondary with  $E > 1$  GeV distribution





# Relations between labels



Parameterized model and Bertini Cascade without heavy nucleus

# MultiBoost Performance

Lerner type – Tree Learner

BaseLerner – Single Stump Learner

Number of Leaves – 2

Number of iterations 100-300

# MultiBoost Performance

Lerner type – Tree Learner

BaseLerner – Single Stump Learner

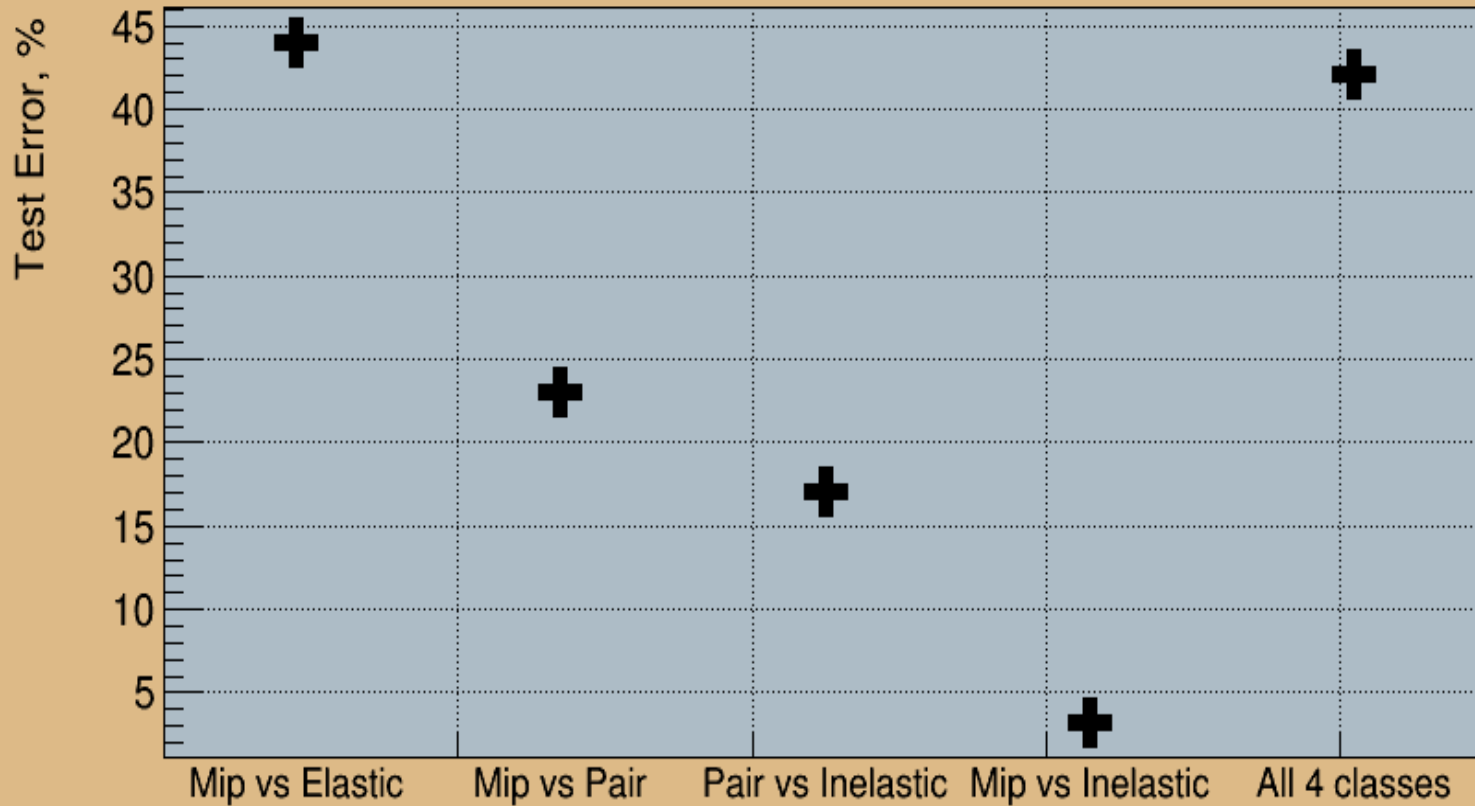
Number of Leaves – 2

Number of iterations 100-300

Manual Features – different combinations of longitudinal profile, energy per layer, radius per layer, Hough transform, etc

# MultiBoost Performance

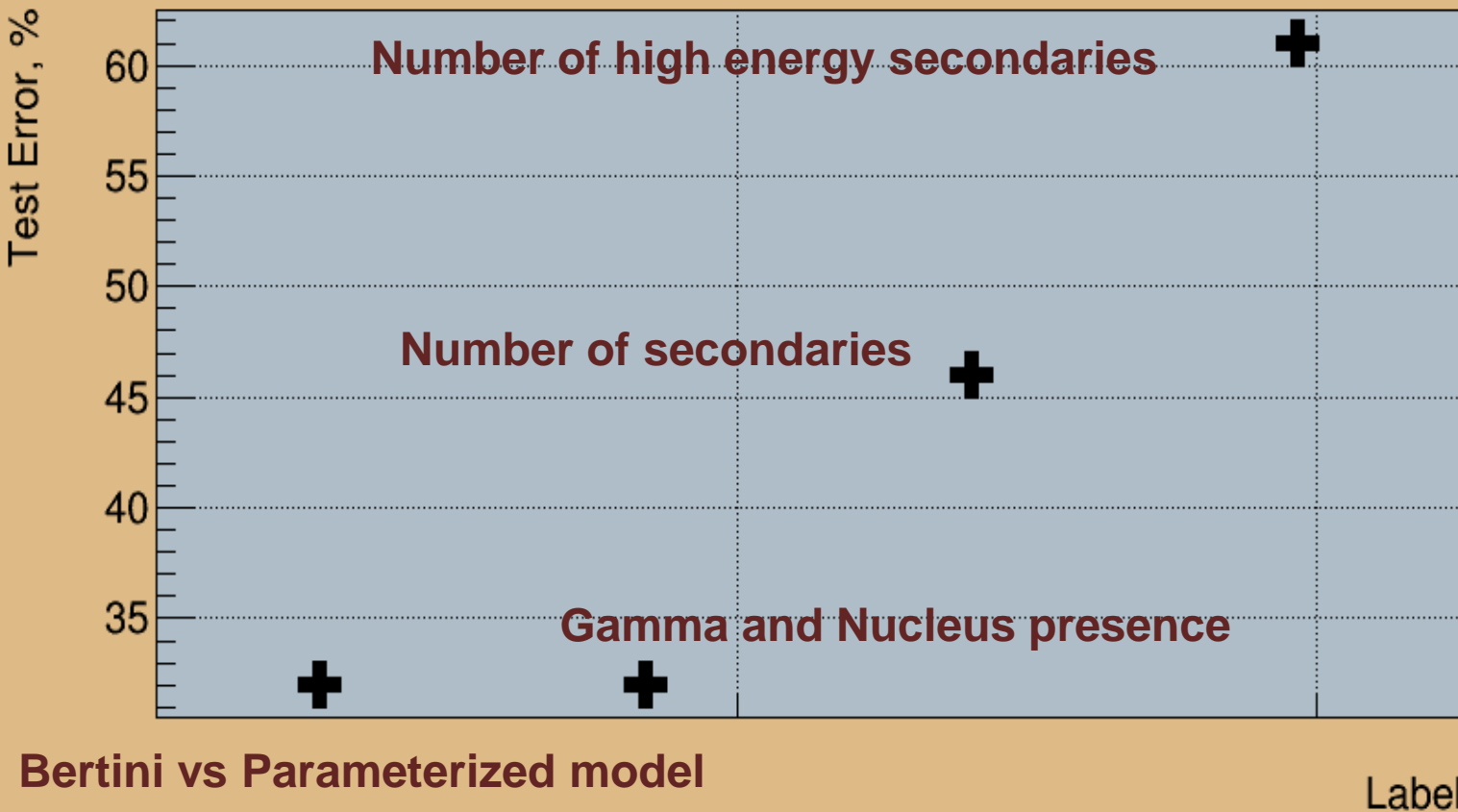
Test Error for basic processes classification



	Train set # of events	Test set # of events
Mip	3739	935
Pair	1760	440
Inelastic	3465	867
Elastic	2788	698

# MultiBoost Performance

Test Error for advance inelastic Labels



	Train set	Test set
Bertini	6565	1633
G4LEP	7429	1866
Gamma & Nucleus	6416	1602
No Gamma & Nucleus	7578	1897
# sec. 0-20	7398	1888
# sec. 20-40	2326	521
# sec. 40-70	1182	316
# sec. 70-100	1527	384
# sec. >100	1561	390
#>1GeV. 0	576	153
#>1GeV. 1	4430	1169
#>1GeV. 2	5446	1316
#>1GeV. 3	2815	675
#>1GeV. >3	727	186

Bertini vs Parameterized model

# Conclusions & Outlook

- Architecture of Geant4 Kernel was explored .
- The basic labels which corresponds to hadronic physical processes for 10 GeV pions was produces with QGSP\_BERT geant4 physics list
- More advanced investigation of pions Inelastic interactions with detector media was studied.
- Supervised learning set-up was tested for hadrons interaction classification and produced benchmarks for further study with deep learning technique. The results is satisfactory with basic classes(physical processes), however in advance inelastic interaction classification results is not significant.
- Further study will imply usage of deep deep learning for existing labels and creation of new ones, more significant.
- After interaction classification problem and precise tune of pipeline, semi-supervised learning will be applied for the study of multiple events and optimization of existing PFA with the aim to remove confusion term and improve jet energy resolution.

# Gratitude

- The author is thankful to ILC and AppStat team for their help, advises and friendly company. And especial gratitude to **Naomi van der Kolk, Balazs Kegl, Thibault Frisson, Philip Bambade, Franck Dubard.**

# Backup Slides

**e01** refers to the *one-error*

$$\widehat{R}_1 \left( \mathbf{f}^{(t)}, \mathbf{X}, \mathbf{Y} \right) = \frac{1}{n} \sum_{i=1}^n \mathbb{I} \left\{ y_{i, \ell_{\mathbf{f}^{(t)}}(\mathbf{x}_i)} < 0 \right\} = \frac{1}{n} \sum_{i=1}^n \mathbb{I} \left\{ \max_{\ell: y_{i, \ell} < 0} f_{\ell}^{(t)}(\mathbf{x}_i) \geq \max_{\ell: y_{i, \ell} > 0} f_{\ell}^{(t)}(\mathbf{x}_i) \right\}.$$

where

$$\ell_{\mathbf{f}^{(t)}}(\mathbf{x}_i) = \arg \max_{\ell'} f_{\ell'}^{(t)}(\mathbf{x}_i) \quad (4)$$

is the single label predicted by  $\mathbf{f}^{(t)}$ .<sup>3</sup> If only one label is positive (single-label multi-class), the error condition is equivalent to  $\ell_{\mathbf{f}^{(t)}}(\mathbf{x}_i) \neq \ell(\mathbf{x}_i)$  where  $\ell(\mathbf{x}_i)$  is the correct single label of  $\mathbf{x}_i$ . If more than one labels are positive (multi-label), it is sufficient for a good prediction if the predicted label (4) is *one* of the positive labels.



# Backup Slides

ADABOOST( $\mathcal{D}$ , BASE( $\cdot, \cdot$ ),  $T$ )

1  $\mathbf{w}^{(1)} \leftarrow (1/n, \dots, 1/n)$   $\triangleright$  initial weights

2 for  $t \leftarrow 1$  to  $T$

3  $h^{(t)} \leftarrow \text{BASE}(\mathcal{D}, \mathbf{w}^{(t)})$   $\triangleright$  base classifier

4  $\epsilon^{(t)} \leftarrow \sum_{i=1}^n w_i^{(t)} \mathbb{I} \{h^{(t)}(\mathbf{x}_i) \neq y_i\}$   $\triangleright$  weighted error of the base classifier

5  $\alpha^{(t)} \leftarrow \frac{1}{2} \ln \left( \frac{1 - \epsilon^{(t)}}{\epsilon^{(t)}} \right)$   $\triangleright$  coefficient of the base classifier

6 for  $i \leftarrow 1$  to  $n$   $\triangleright$  re-weighting the training points

7 if  $h^{(t)}(\mathbf{x}_i) \neq y_i$  then  $\triangleright$  error

8  $w_i^{(t+1)} \leftarrow \frac{w_i^{(t)}}{2\epsilon^{(t)}}$   $\triangleright$  weight increases

9 else  $\triangleright$  correct classification

10  $w_i^{(t+1)} \leftarrow \frac{w_i^{(t)}}{2(1-\epsilon^{(t)})}$   $\triangleright$  weight decreases

11 return  $f^{(T)}(\cdot) = \sum_{t=1}^T \alpha^{(t)} h^{(t)}(\cdot)$   $\triangleright$  weighted "vote" of base classifiers

# Backup Slides

## Mip vs inelastic



# Backup Slides

## Mip vs Elastic

