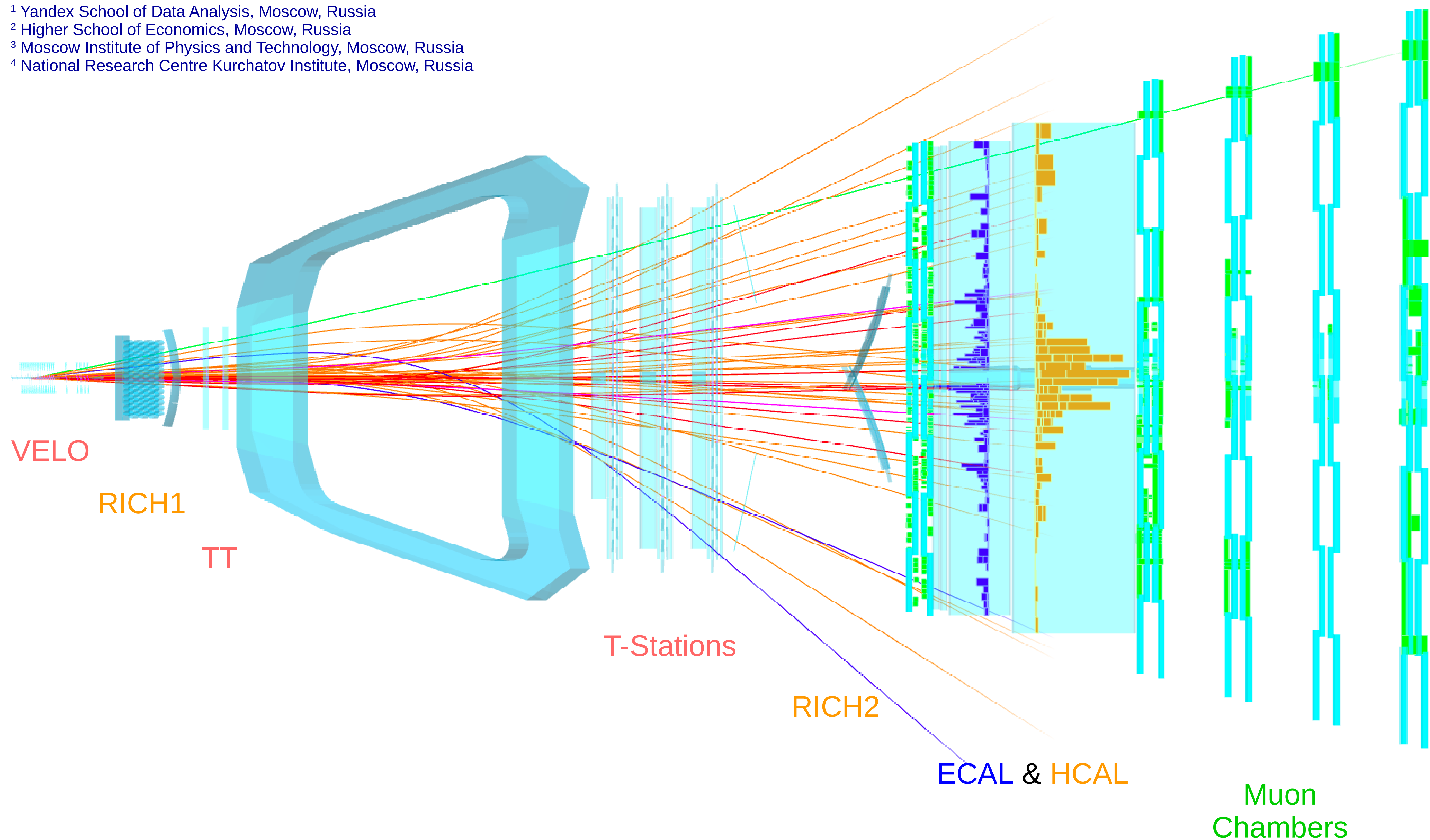


# Machine Learning Based Global Particle Identification at LHCb

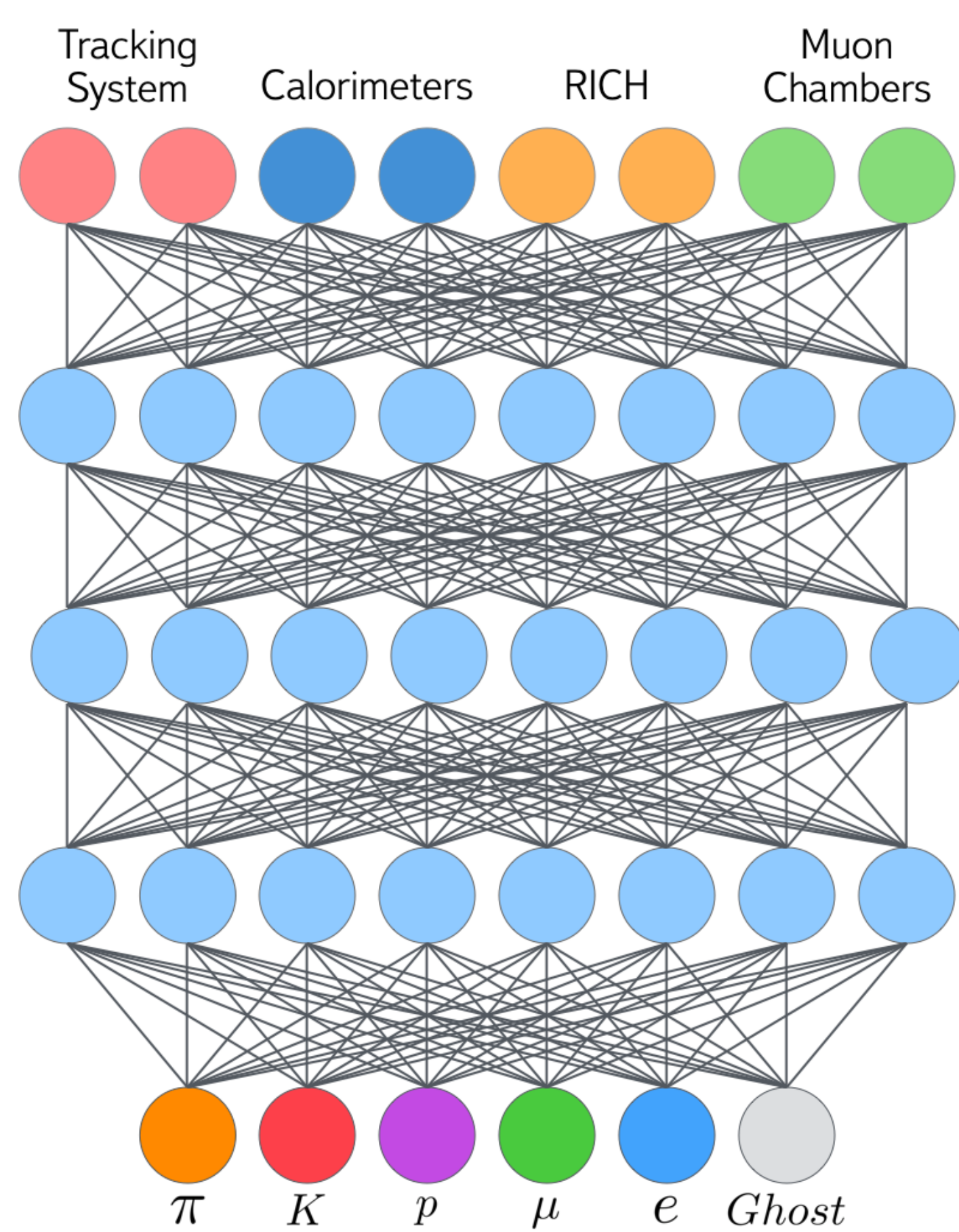
Denis Derkach<sup>1,2</sup>, Mikhail Hushchyn<sup>1,3</sup>, Tatiana Likhomanenko<sup>1,2,4</sup>, Alex Rogozhnikov<sup>1,2</sup>, Nikita Kazeev<sup>1,2</sup>, Victoria Chekalina<sup>1</sup>, Radoslav Neychev<sup>1,3</sup>, Stanislav Kirillov<sup>1</sup>, Fedor Ratnikov<sup>1,2</sup> on behalf of the LHCb collaboration

18th International Workshop on Advanced Computing and Analysis Techniques in Physics Research, Seattle, August 21-25, 2017

<sup>1</sup> Yandex School of Data Analysis, Moscow, Russia  
<sup>2</sup> Higher School of Economics, Moscow, Russia  
<sup>3</sup> Moscow Institute of Physics and Technology, Moscow, Russia  
<sup>4</sup> National Research Centre Kurchatov Institute, Moscow, Russia



## Best-Efficiency Models



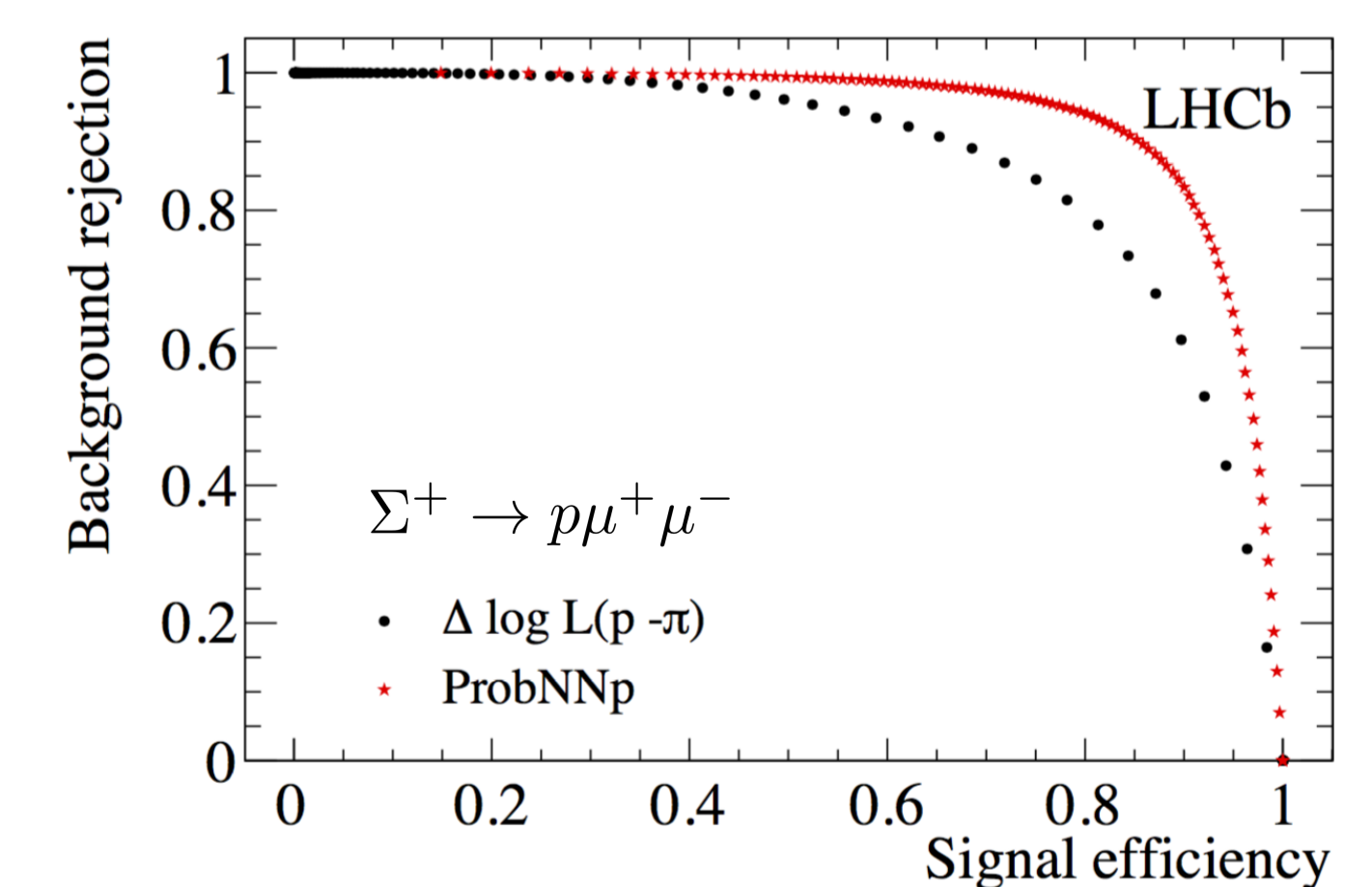
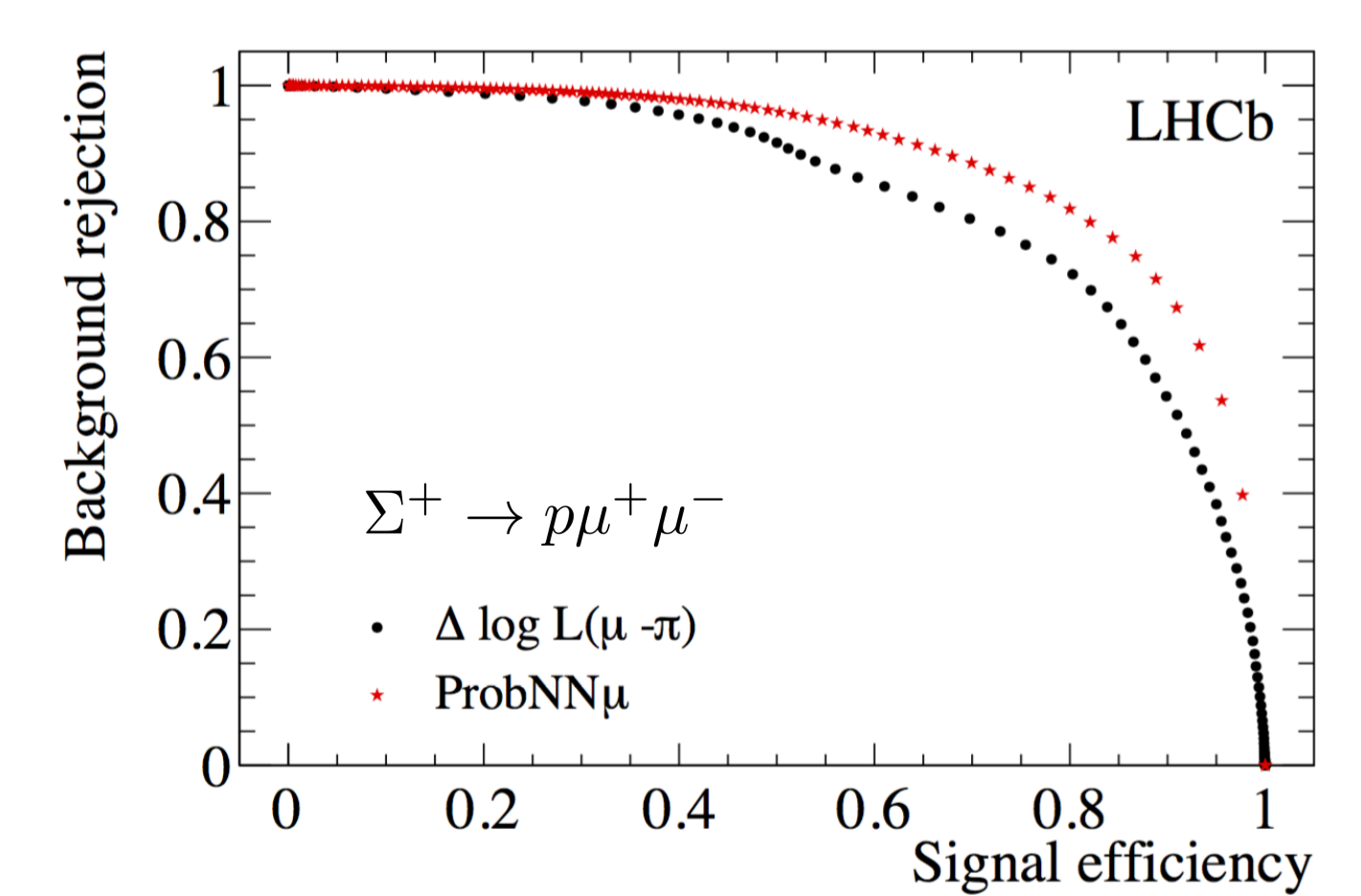
The problem is to identify the charged particle type which a given track is associated. There are five particle types: **electron**, **muon**, **pion**, **kaon**, **proton** plus the "not classifiable" **ghost** for a total of 6 hypotheses. This problem can be considered as a multiclass classification problem. For PID information from Ring-Imaging Cherenkov Detector (RICH), Electromagnetic Calorimeter, Hadronic Calorimeter and Muon Chambers sub-detectors and track reconstruction is used.

Current particle ID approaches:

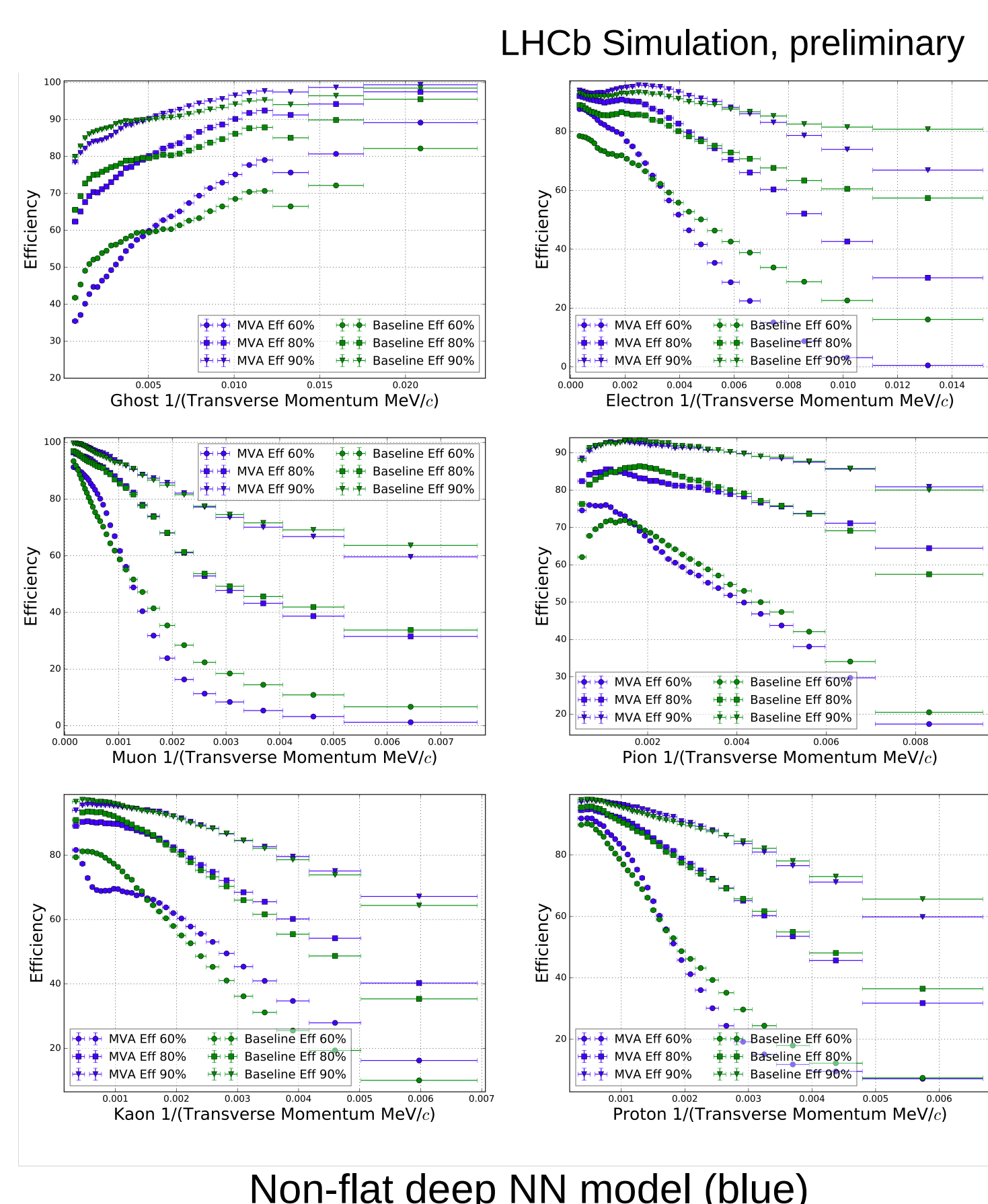
- $\Delta \log L$ : Estimate likelihood of a particle type based on a subdetector response. Likelihood of the subdetectors are combined into global likelihood of the particle type.
- **ProbNN** (baseline): The subdetector responses are combined using one hidden layer neural network (TMVA MLP) in one-particle-versus-rest mode.

The PID was improved using gradient boosting algorithms (**XGBoost**, **CatBoost**) and deep neural network (**deep NN**) in multiclassification mode.

	(1-AUC)/(1-AUC <sub>baseline</sub> )					
	Ghost	Electron	Muon	Pion	Kaon	Proton
baseline						
deep NN	-29 %	-41 %	-52 %	-37 %	-20 %	-17 %
XGBoost	-24 %	-37 %	-50 %	-34 %	-18 %	-15 %
CatBoost	-30 %	-43 %	-54 %	-37 %	-20 %	-18 %



## Best-Flatness Models



Non-flat deep NN model (blue)

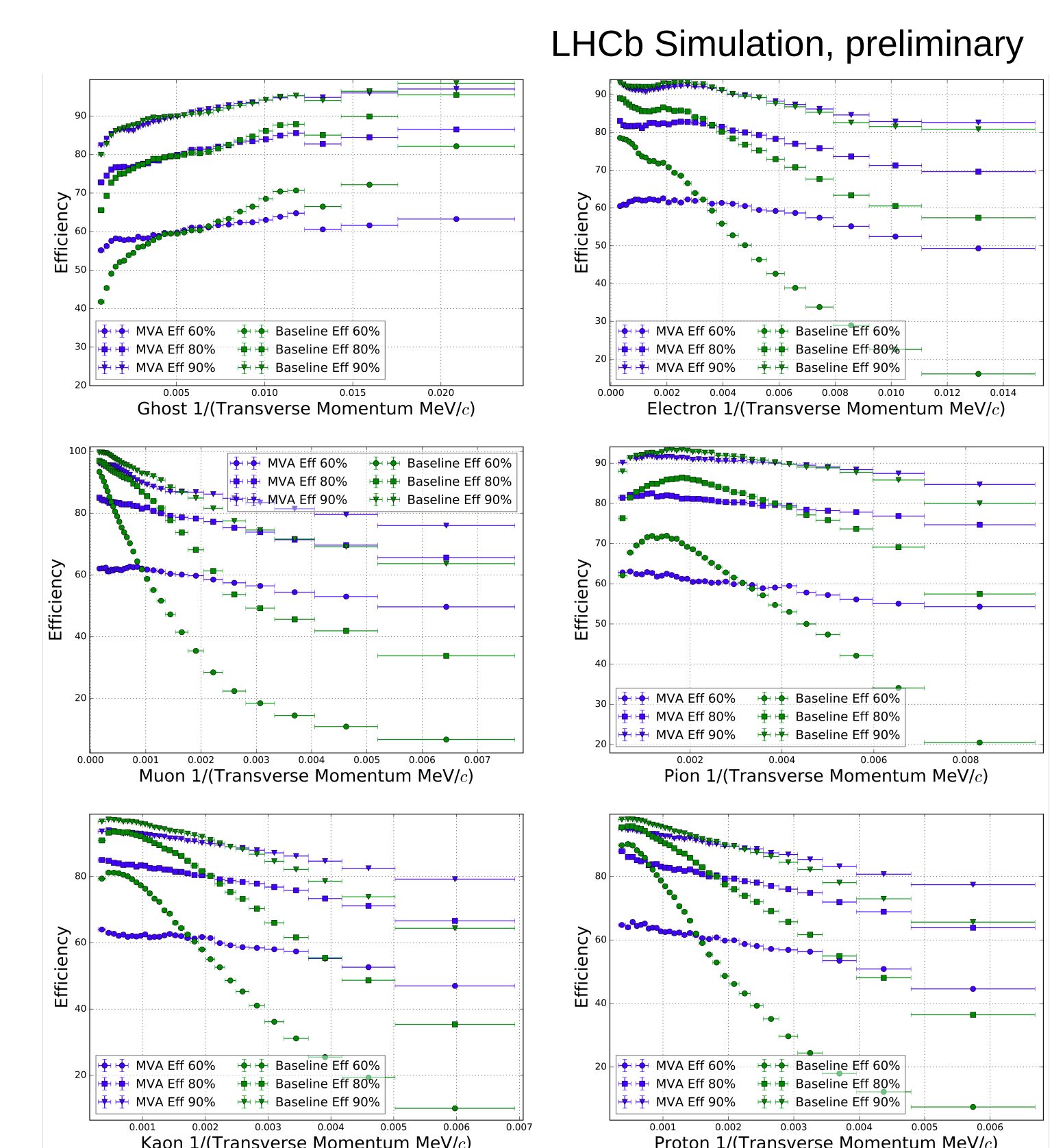
The PID information strongly depends on the kinematic variables. This relationship leads to strong dependency between PID efficiency and kinematic variables as it is shown in the left figure. In practice it could be helpful to have flat PID efficiencies along chosen control variables to reduce systematic effects. To provide uniformity along some observables models were trained using the modified loss function:

$$\mathcal{L} = \mathcal{L}_{ExpLoss} + \alpha \mathcal{L}_{FL}$$

where the first term is the classification loss function and the second term is the uniformity loss written in differential form of *Cramer-von Mises* measure. The right figure demonstrates the flatness improvement.

	(1-AUC)/(1-AUC <sub>baseline</sub> )					
	Ghost	Electron	Muon	Pion	Kaon	Proton
baseline						
p + p <sub>T</sub> flatness	-23 %	-20 %	-27 %	-26 %	+2 %	+5 %
2d(p, p <sub>T</sub> ) flatness	-21 %	-9 %	-13 %	-23 %	+12 %	+23 %
p + p <sub>T</sub> + η + nTracks flatness	-22 %	-13 %	-26 %	-24 %	+2 %	+6 %
4d(p, p <sub>T</sub> , η, nTracks) flatness	-21 %	-4 %	-13 %	-20 %	+10 %	+25 %

There is a trade off between the PID models efficiency and flatness. It is possible to archive any efficiency flatness of a model along any variable. However, the efficiency decreases with increasing its flatness.



Flat 4d model (blue)