



Machine Learning in Baikal-GVD:  
Noise Filtering and EAS Suppression via Neural Networks  
with Domain Adaptation

This work was supported by the Russian Science Foundation under grant no. 24-72-10056.

**Matseiko Albert**  
MIPT, INR RAS  
matseiko.av@phystech.edu

Quarks-2026, May

# Outline

0. *Baikal-GVD experiment: brief description*
1. *Baikal-GVD ML processing pipeline*
2. *Developed NN models: metrics and real data validation*
  - *Background EAS Events (pre)filtering*
  - *Signal Hits Filtering*
  - *$\nu$ -candidate extraction*

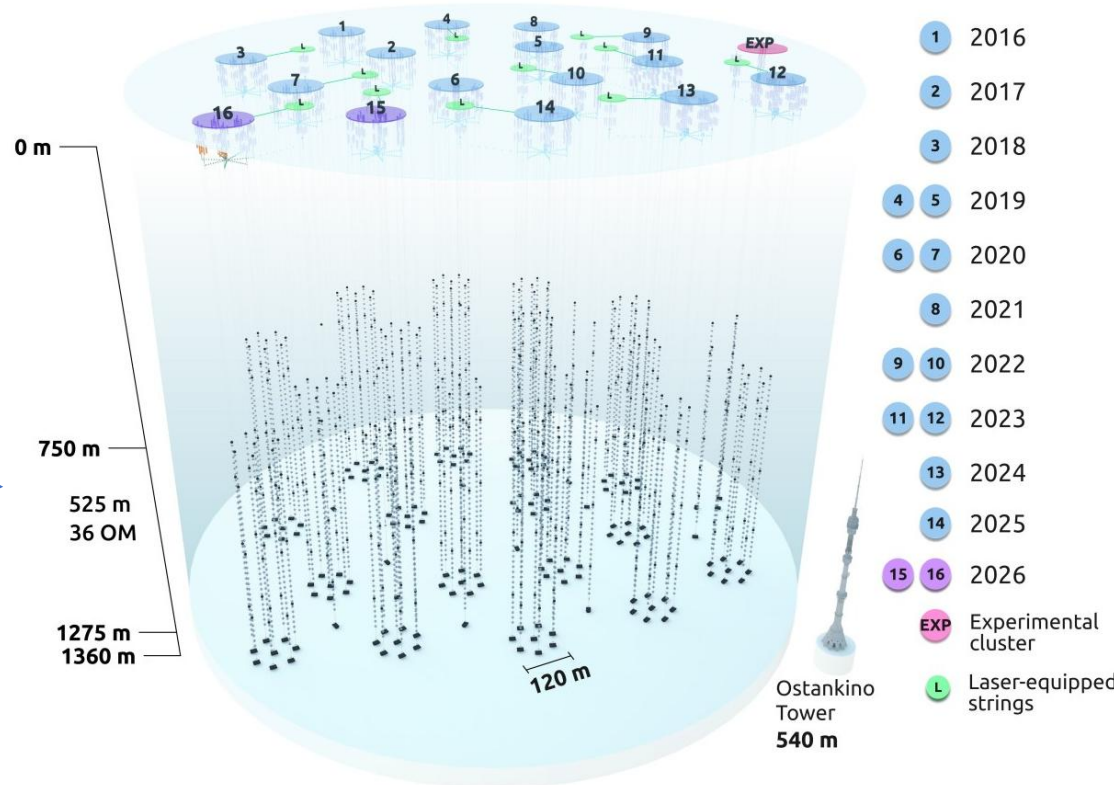
# Baikal-GVD

## Setup process

Optical modules

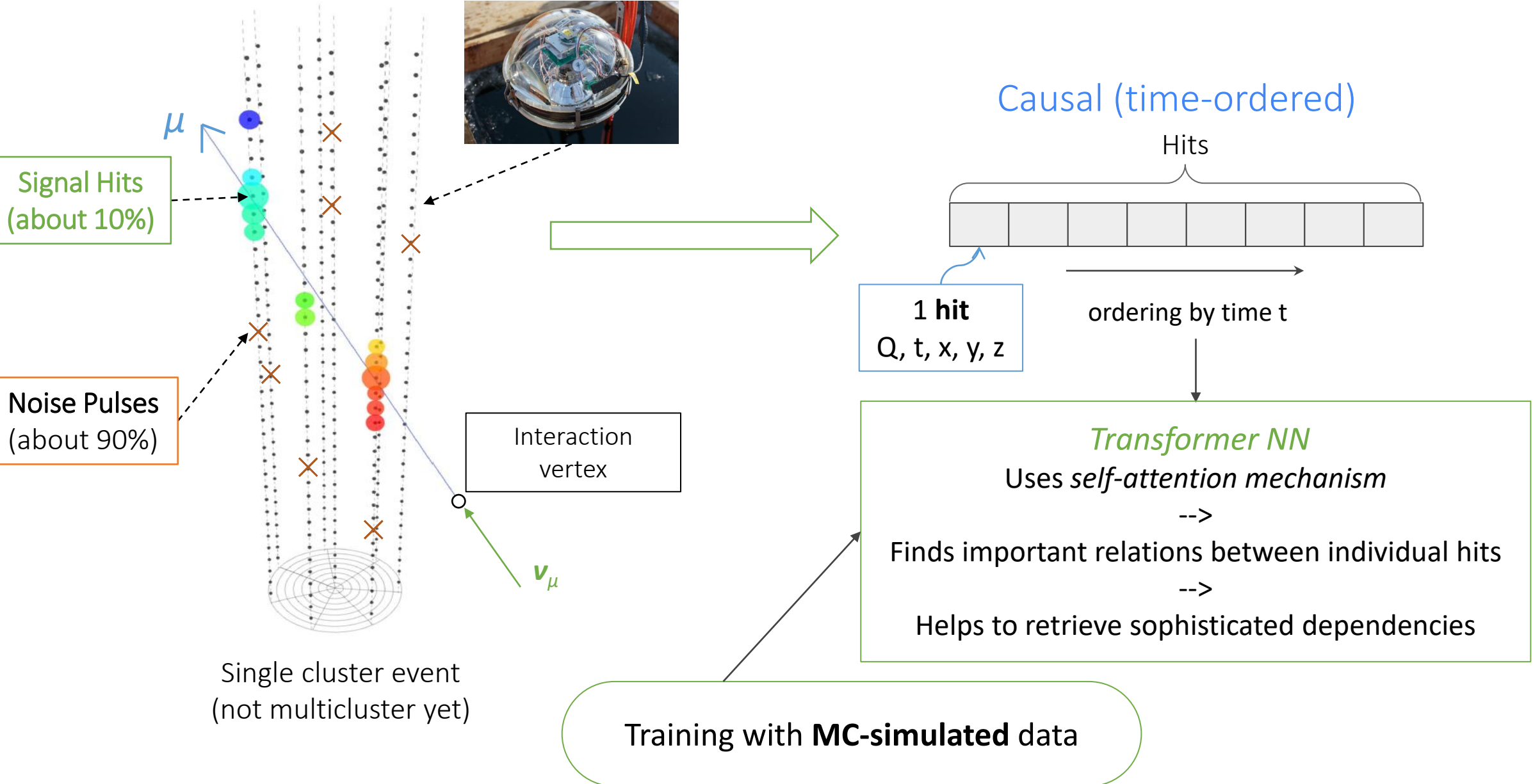


Form «strings»  
→  
Strings form «clusters»

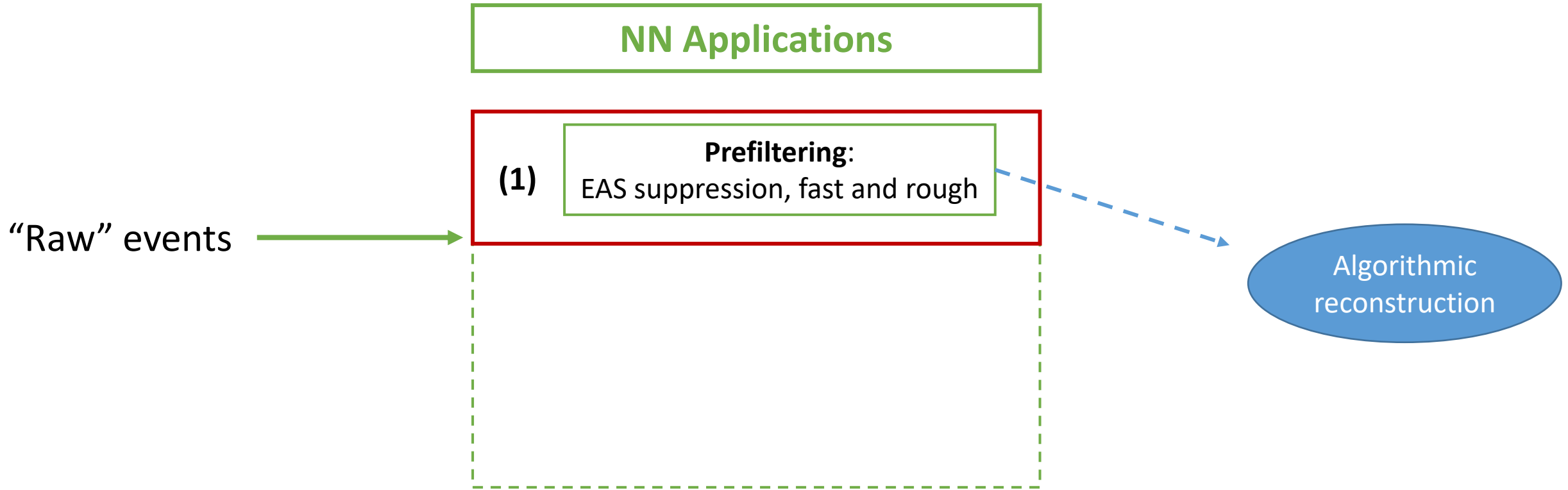


>16 clusters are setup by 2026

# Event representation for NN



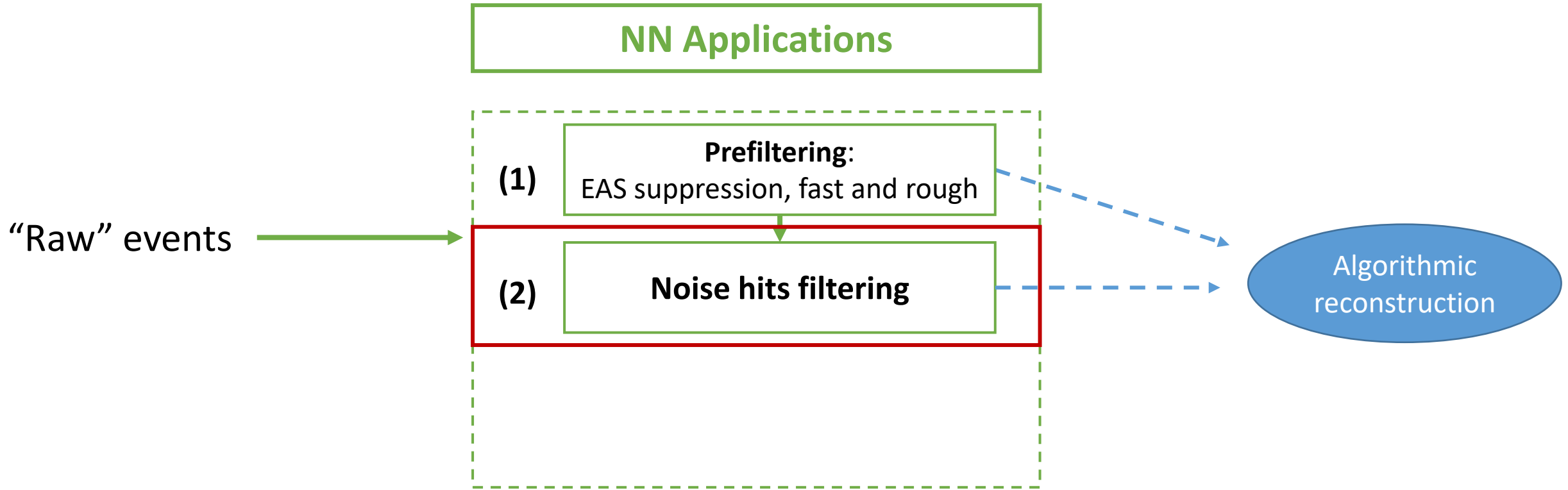
# ML processing pipeline



(1)  
*Problem: 1  $\nu$  per  $10^6$ - $10^7$  events.*  
*NN solution: Prefilter Model*  
suppresses background by factor 1000 while keeping 99%  $\nu$ -events

Also may speed up standard reconstruction!

# ML processing pipeline



(2)

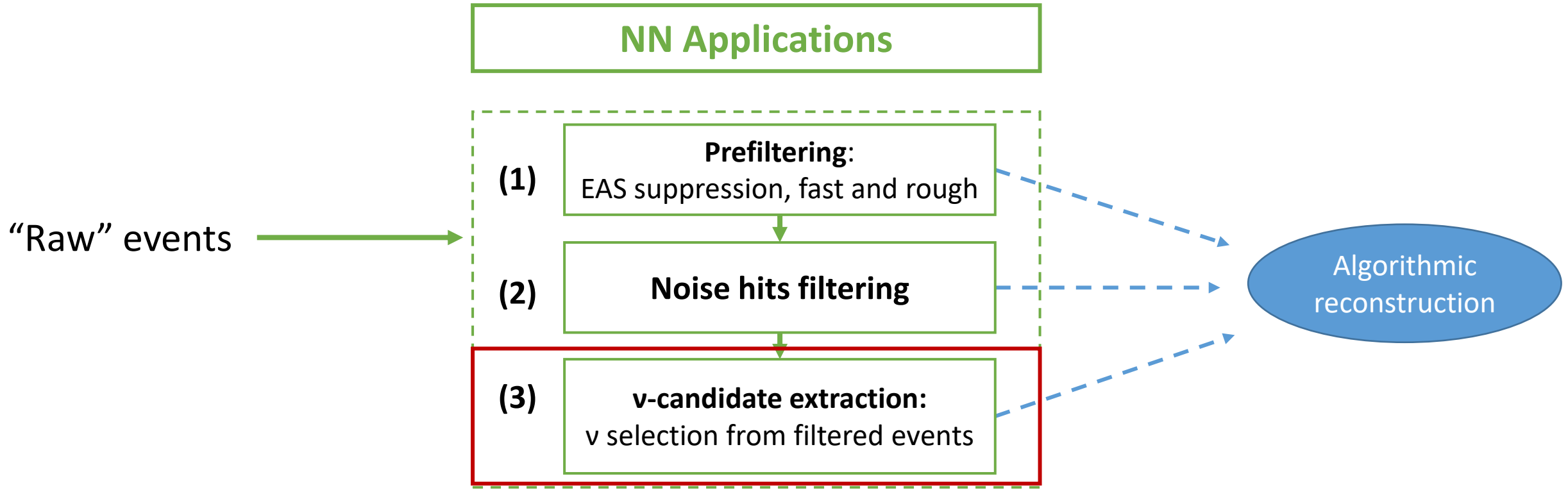
*Problem: 85%-90% of hits are due to water luminescence.*

*NN solution: Hit-level Filter*

suppresses noise hits with Recall=Precision=95%

Also may improve usual reconstruction!

# ML processing pipeline



(3)

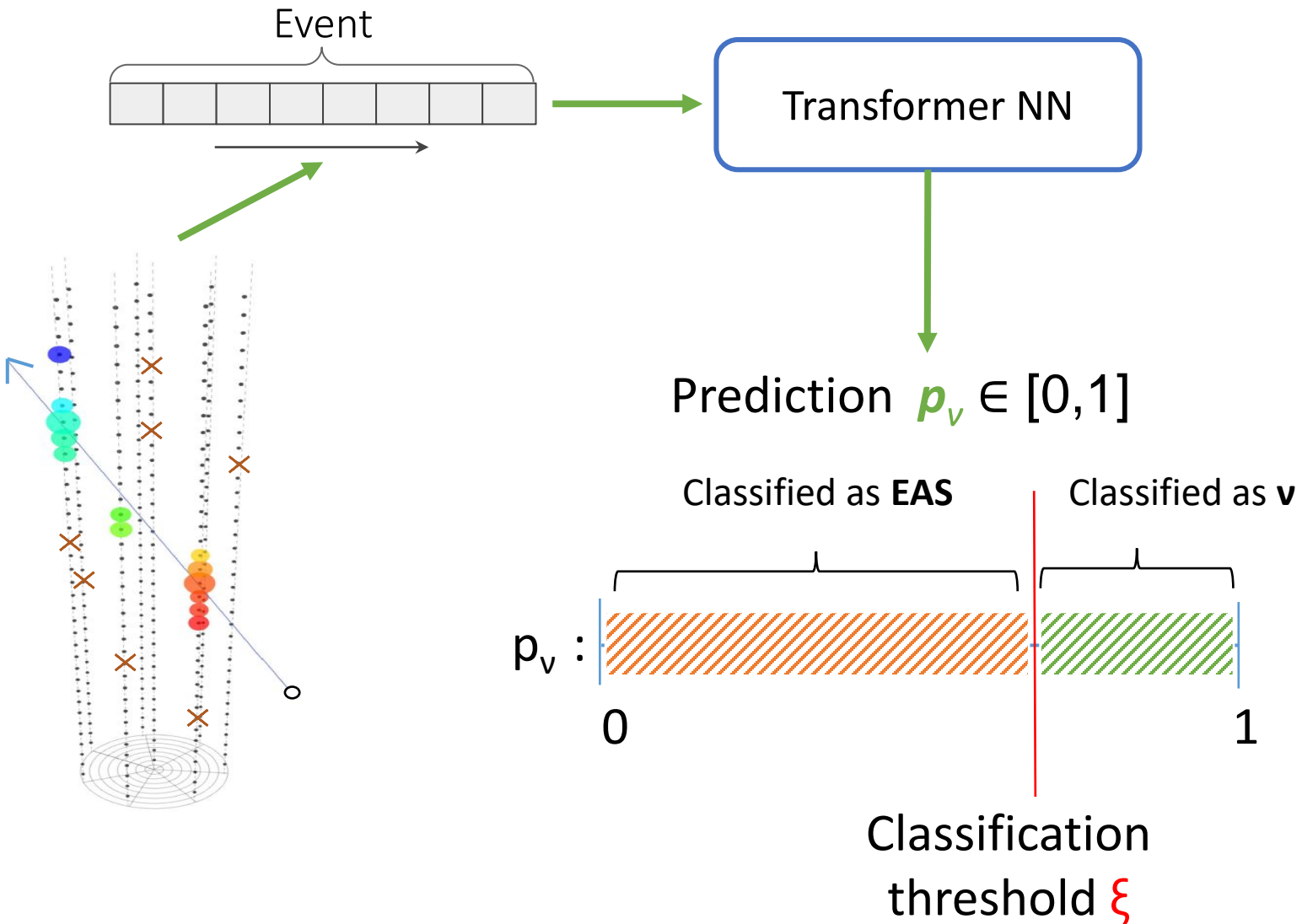
*Problem: still many background events in data.*

*NN solution: v-candidate Extractor*

uses signal hits to select v-events precisely

Selected v-candidates may be passed through usual reconstruction!

# Pre-filter model



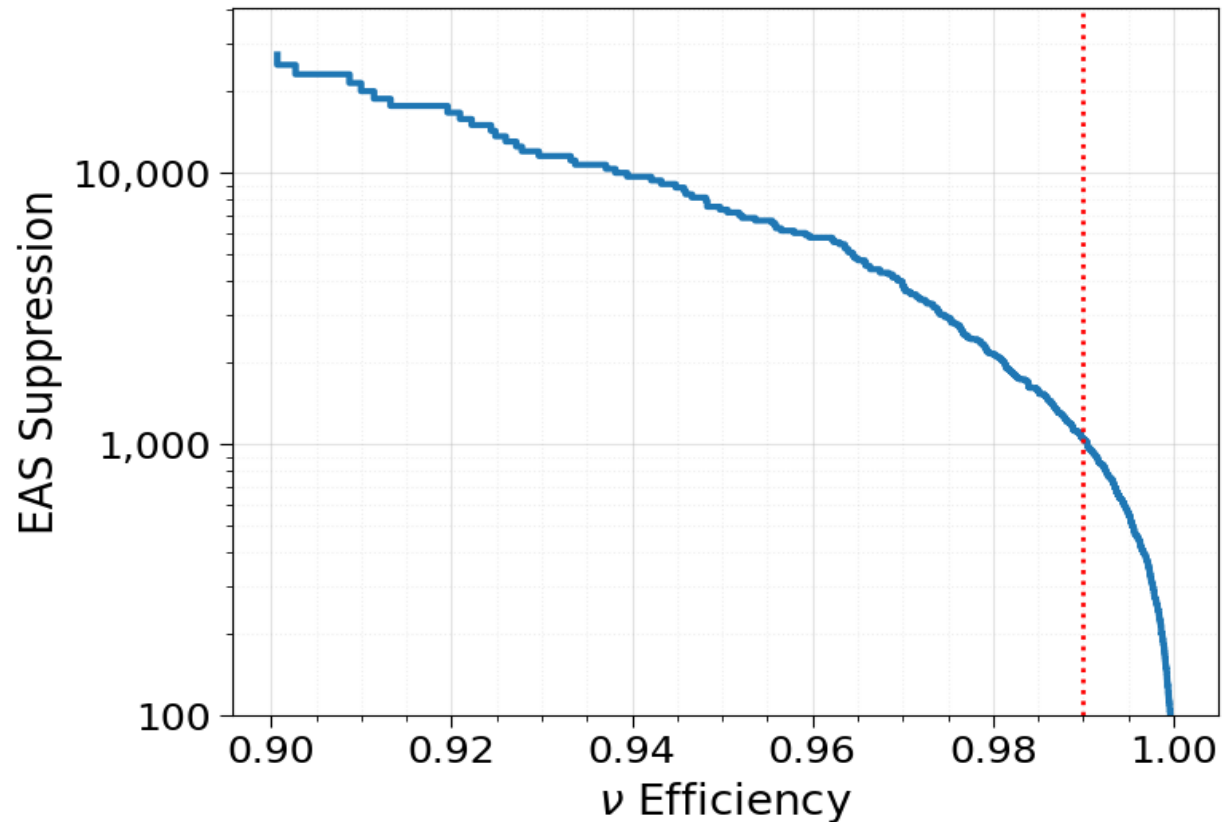
## MC datasets events #

	<b>EAS</b>	$\nu_{atm}$	$\nu_{cosmo}$
<b>Train</b>	1M	0.5M	0.5 M
<b>Val</b>	50k	25k	25k
<b>Test</b>	600k	100k	100k

Metrics having  $\xi$  fixed:

- **EAS suppression factor ( $FPR^{-1}$ )**
- **$\nu$  events efficiency (TPR)**

# Prefilter model: metrics

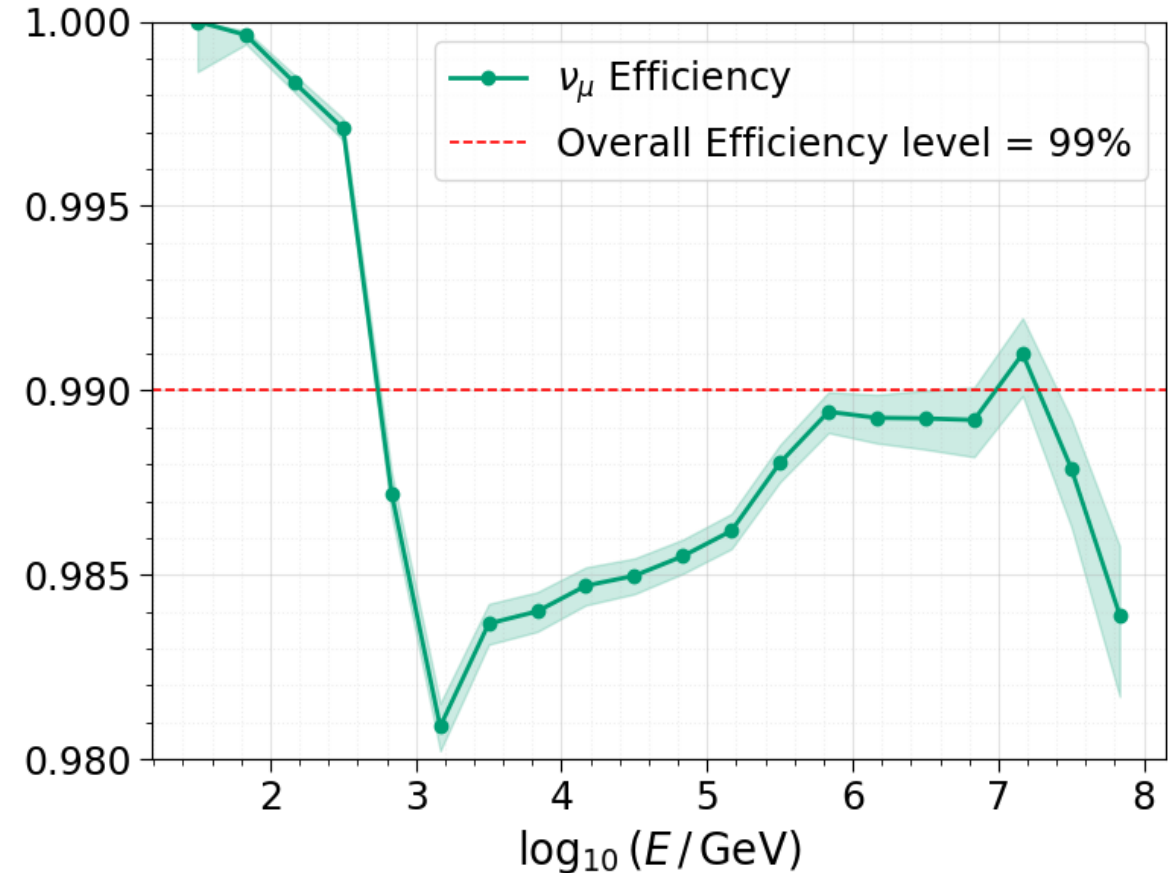


*EAS Suppression VS  $\nu$  Efficiency (scanning  $\xi$ )*

*At  $\xi = 0.761$ : suppresses EAS 1000 times while saving 99%  $\nu$*

# Results for MC ground truth quality cuts:

- # signal hits  $\geq 8$
- # signal strings  $\geq 2$



*$\nu$  Efficiency VS  $\nu$  Energy  
@  $\xi = 0.761$*

# Pre-filter model: Domain Adaptation (DA) from MC to real data

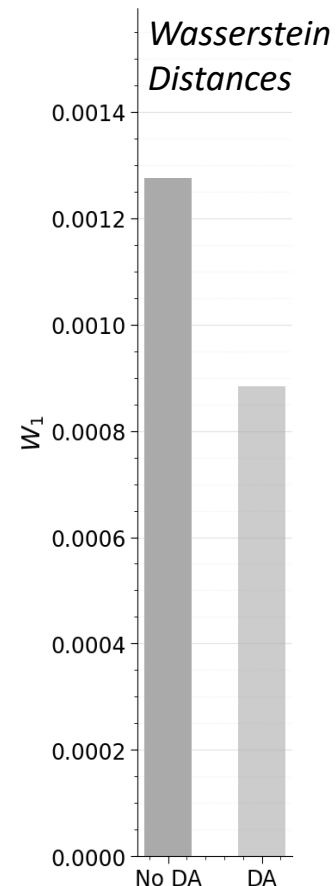
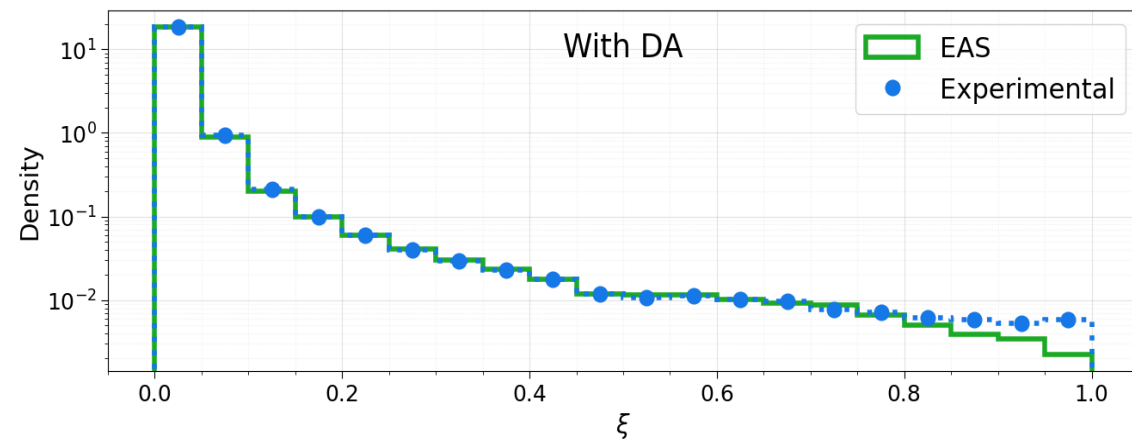
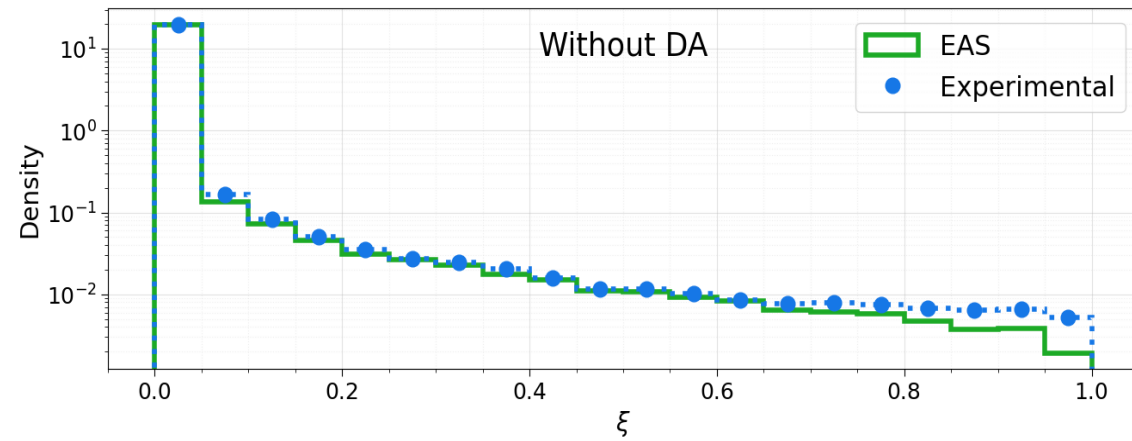
DA technique from [arXiv:1505.07818](https://arxiv.org/abs/1505.07818)

Uses (indirectly) real data while training on MC!

	Real Events
Train	255k
Val	45k

Applying standard cuts via **algorithmic** reconstruction (to select physically same EAS subsets):

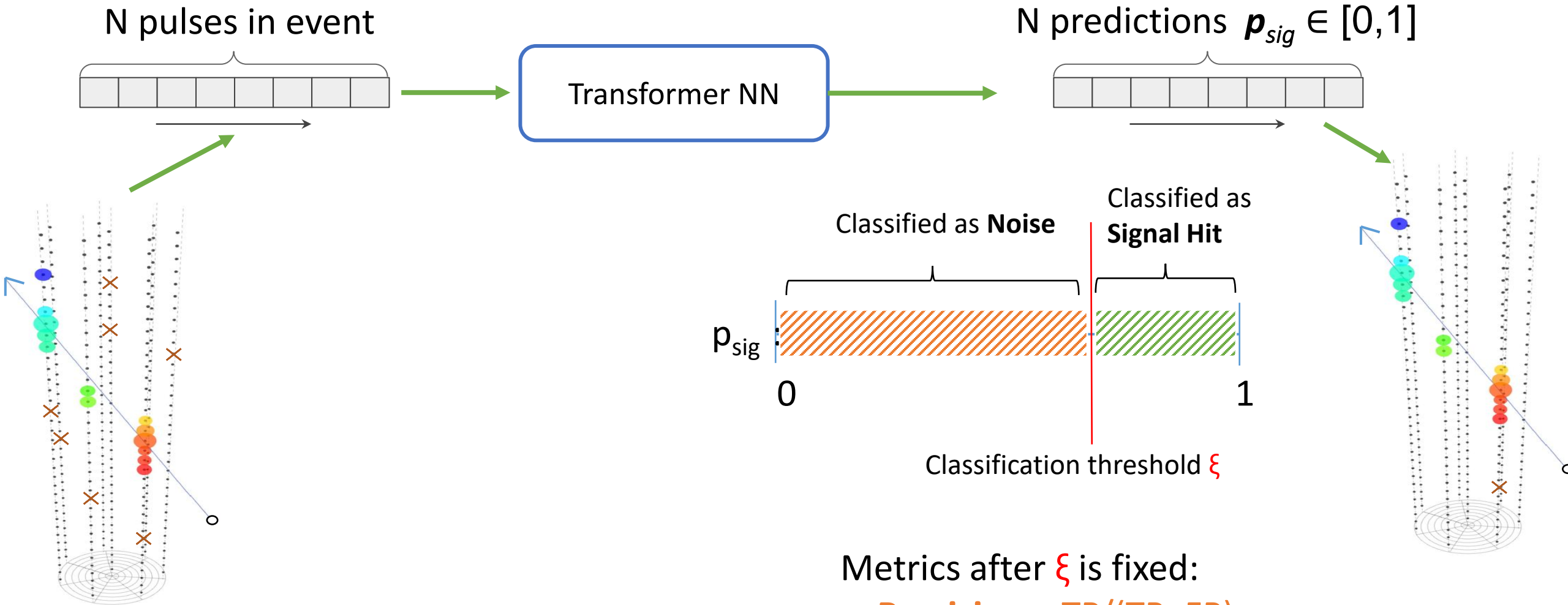
- # of signal hits  $\geq 8$  and # of signal strings  $\geq 3$ ;
- covariance matrix convergence status equal to 3 (fully successful fit);
- number of minimizer calls  $< 350$  (convergence quality);
- $\log_{10} \text{phit} > -9$ ;
- $\log_{10}(F/(n\text{hits} - 5)) < 2$ , where F is the minimizer objective value;
- $n\text{triplets}/n\text{hits} > 0.1$  (track-topology requirement);
- $z\text{dist} > 70$  m;
- $z\text{center} < 220$  m;
- $\theta_{\text{rec}} \in (80^\circ, 180^\circ)$ ;
- $\sigma\theta < 2.5$  rad.



Prediction distributions for MC (EAS) and real (Experimental) dataset

Better correspondence with DA (?)

# Noise hits filter



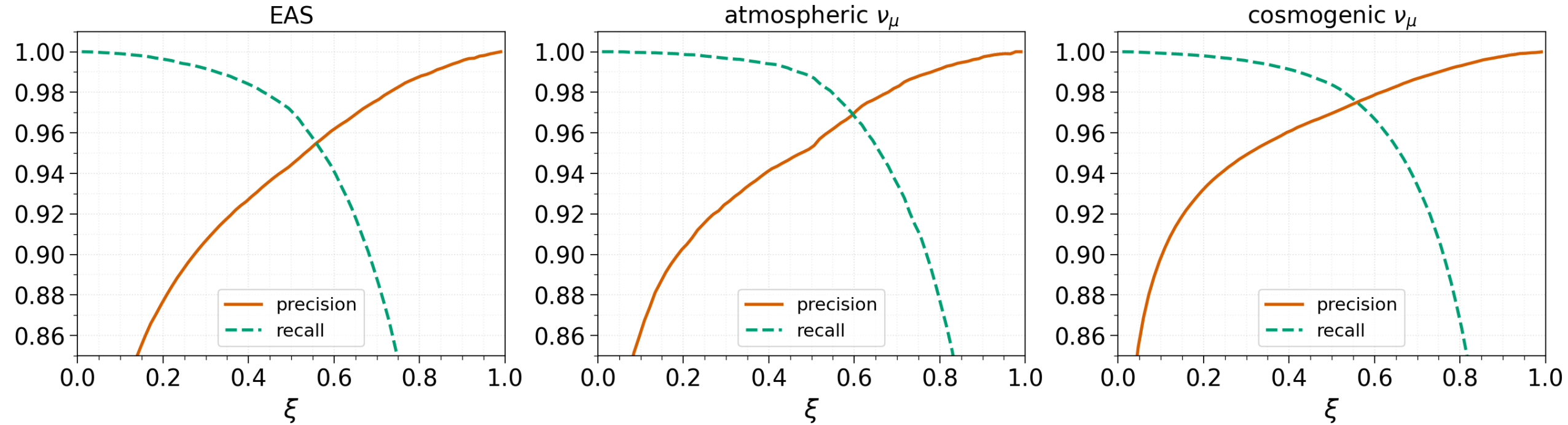
Metrics after  $\xi$  is fixed:

- **Precision** =  $TP / (TP + FP)$
- **Recall** =  $TP / (TP + FN)$

# Noise hits filter: metrics

Results for quality cuts **via the model**:

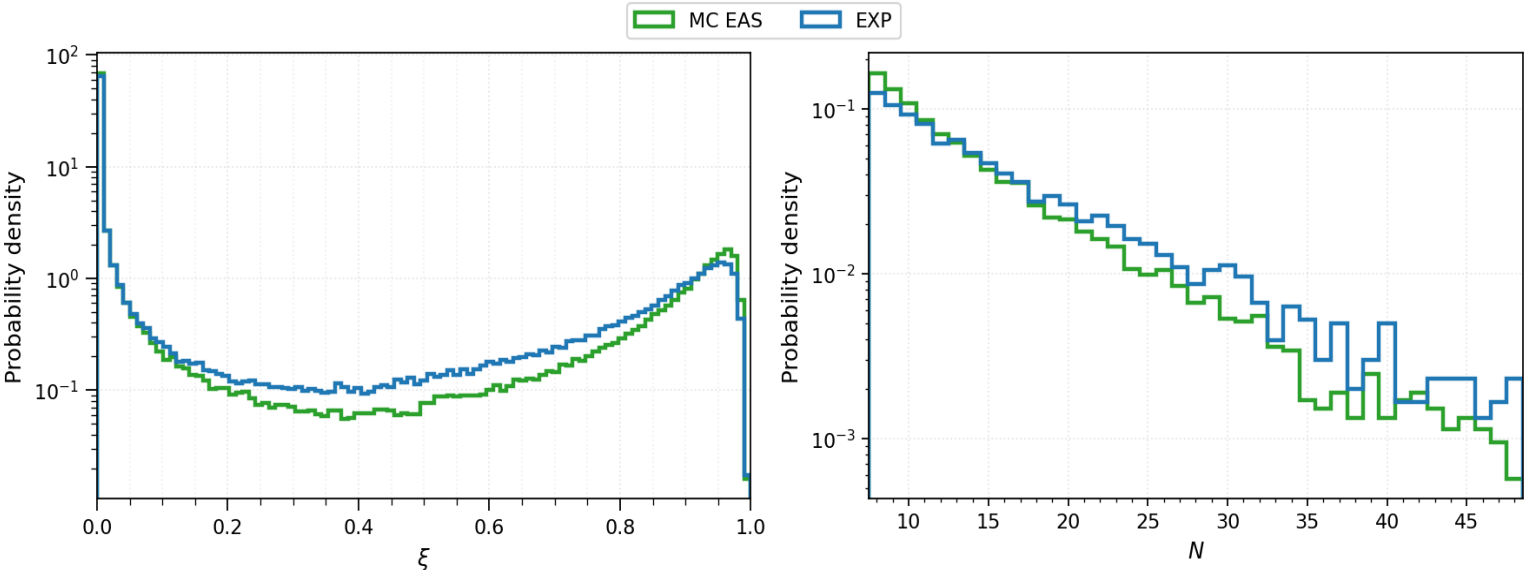
- # signal hits  $\geq 8$
- # signal strings  $\geq 2$



Both precision and recall exceed 95% at  $\xi = 0.55$  for all event types.

Surpasses standard “scan-fit” algorithm [arXiv:2108.00208](https://arxiv.org/abs/2108.00208)

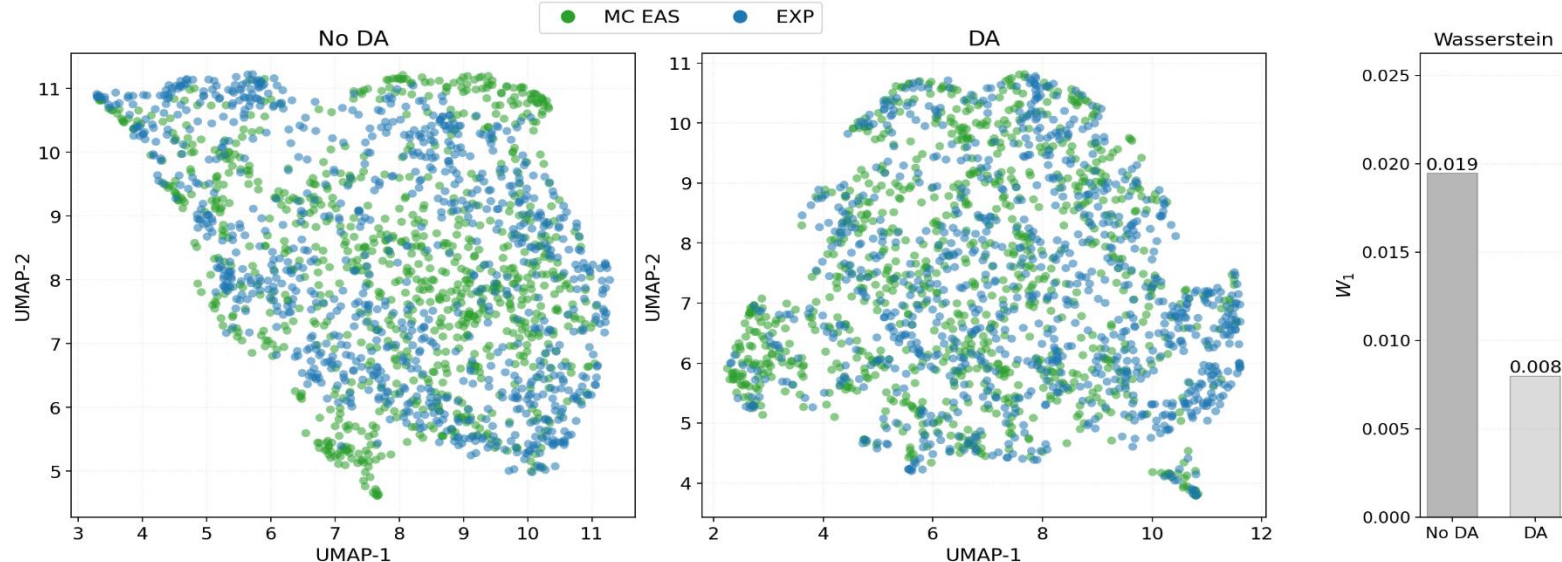
# Noise hits filter: Domain Adaptation (DA) from MC to real data



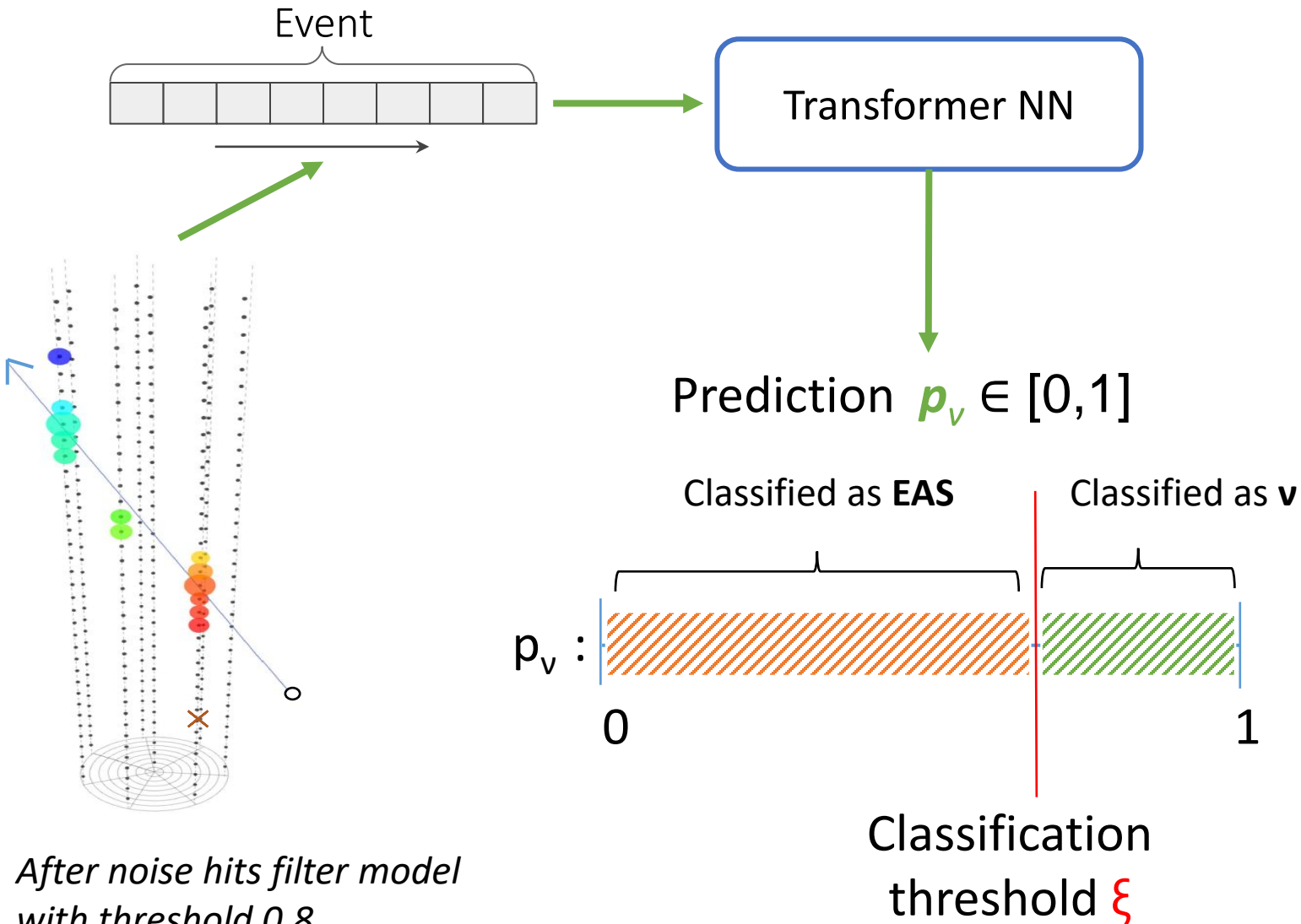
$p_{sig}$  distributions and signal hits multiplicity for MC (EAS) and real data after cuts h8s2 **by the model itself** (trained **with DA**)

**UMAP** of the last encoder layer output to compare cases with and without DA

Wasserstein distance is computed between  $p_{sig}$  distributions



# $\nu$ -candidate Extractor



Soft cut applied: signal hits  $\geq 5$

MC datasets events #

	<b>EAS</b>	$\nu_{atm}$	$\nu_{cosmo}$
<b>Train</b>	2.4M	1.2M	1.2 M
<b>Val</b>	200k	100k	100k
<b>Test</b>	23M	100k	100k

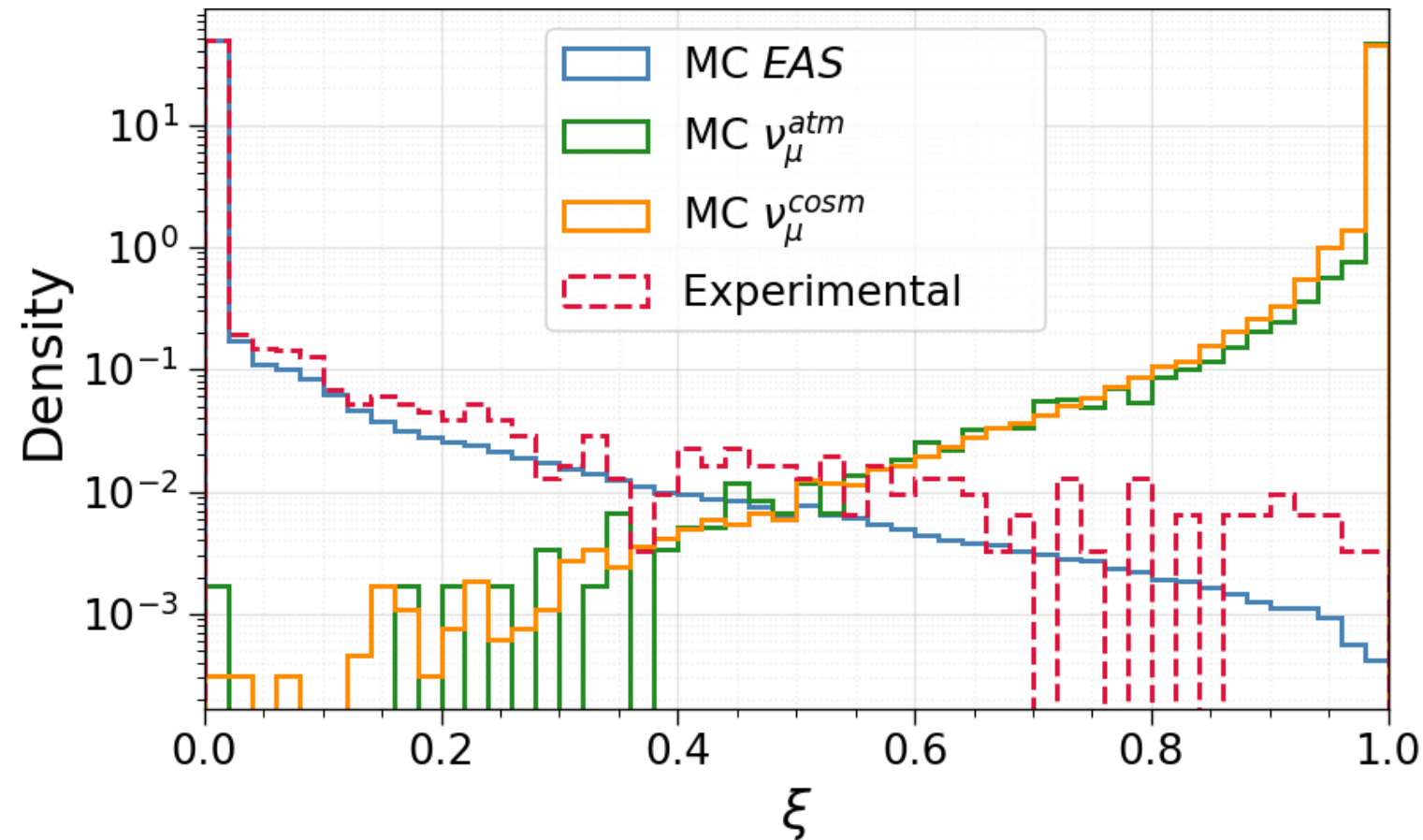
Metrics having  $\xi$  fixed:

- **EAS** suppression factor ( $FPR^{-1}$ )
- $\nu$  events efficiency (TPR)

# $\nu$ -candidate Extractor: real data (preliminary)

Same DA technique from [arXiv:1505.07818](https://arxiv.org/abs/1505.07818) was applied

Quality cuts: # signal hits  $\geq 8$ , # signal strings  $\geq 3$ . MC EAS suppression of 1M times working point.



Prediction distribution for real and MC data

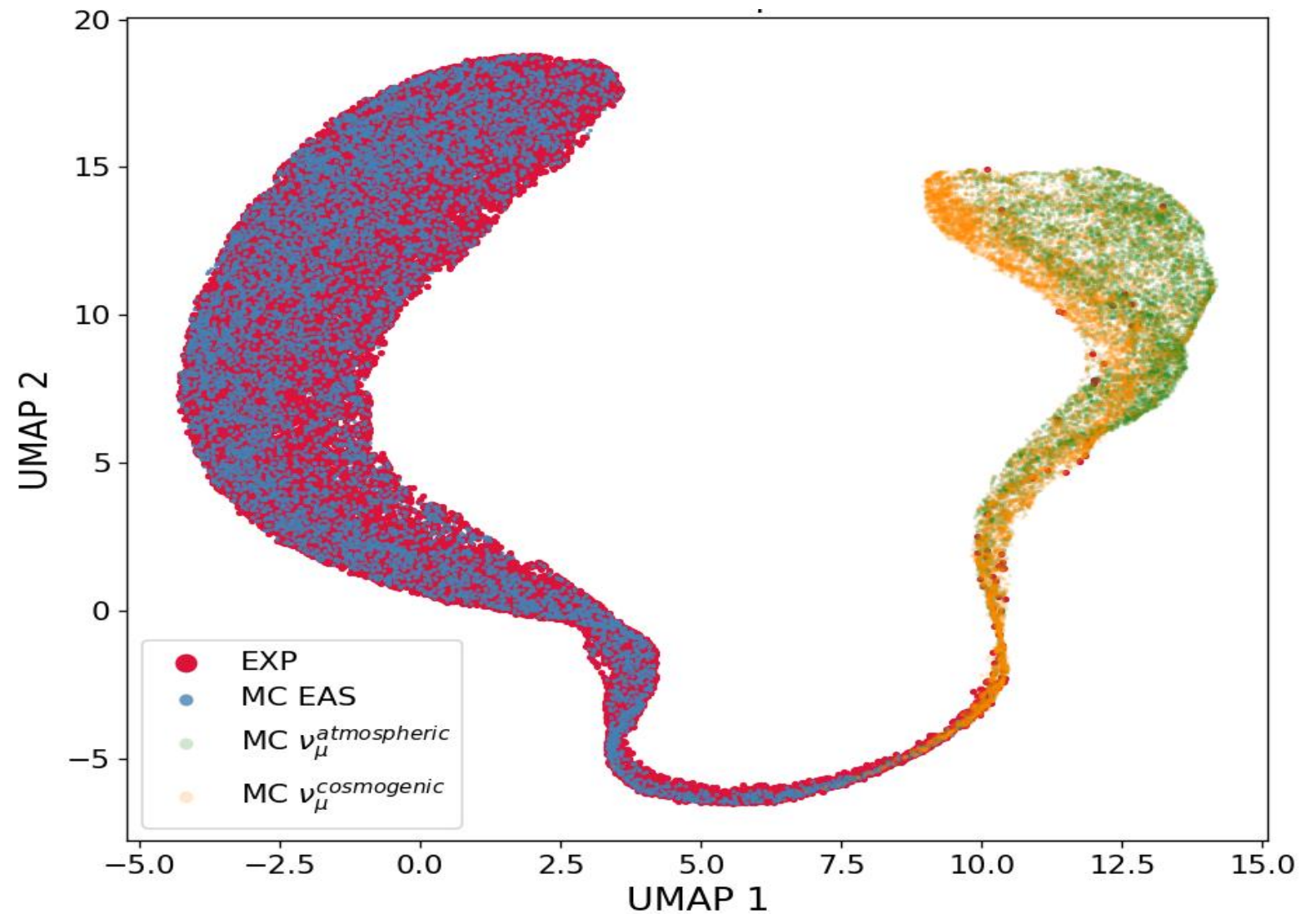
- At  $\xi=0.995$  we suppress MC EAS by factor  $10^6$ , saving 80% of MC  $\nu$  events
- Close MC EAS and real data distributions (pipeline works!)
- But not exactly: some false  $\nu$  candidates with  $\xi > 0.99$  in  $\approx 15,000$  experimental data events sample.

# $\nu$ -candidate Extractor: real data (preliminary)

Same DA technique from [arXiv:1505.07818](https://arxiv.org/abs/1505.07818) was applied

Quality cuts: # signal hits  $\geq 8$ , # signal strings  $\geq 3$ . MC EAS suppression of 1M times working point.

At the UMAP representation red “experimental” dots are found deep in the MC  $\nu$  events area  $\rightarrow$  strange sign to work with.



# Summary

## 1. Progress in 3 ML tasks was achieved (MC metrics)s

### ***I. Events Prefilter***

Preserves 99% h8s2 **v events**  
Suppresses **background** in 1000 times.

### ***II. Noise hits filter***

**95% Recall** and **Precision**  
for all hits

### ***III. v-candidate extractor***

Suppresses **background** by factor  $10^6$   
Preserves 80% **v events**

## 2. Made steps towards real data processing

- Domain Adaptation makes models more robust while switching from MC to real data
- However some domain mismatch was found for **v-candidate extraction** task

## 3. What's next?

- Figure out and solve the domain mismatch (MoE, soft labeling)
- To solve downstream reconstruction tasks:
  - Energy
  - Direction of Arrival



Thanks!  
Please ask your questions!

Contacts:

[matseiko.av@phystech.su](mailto:matseiko.av@phystech.su)

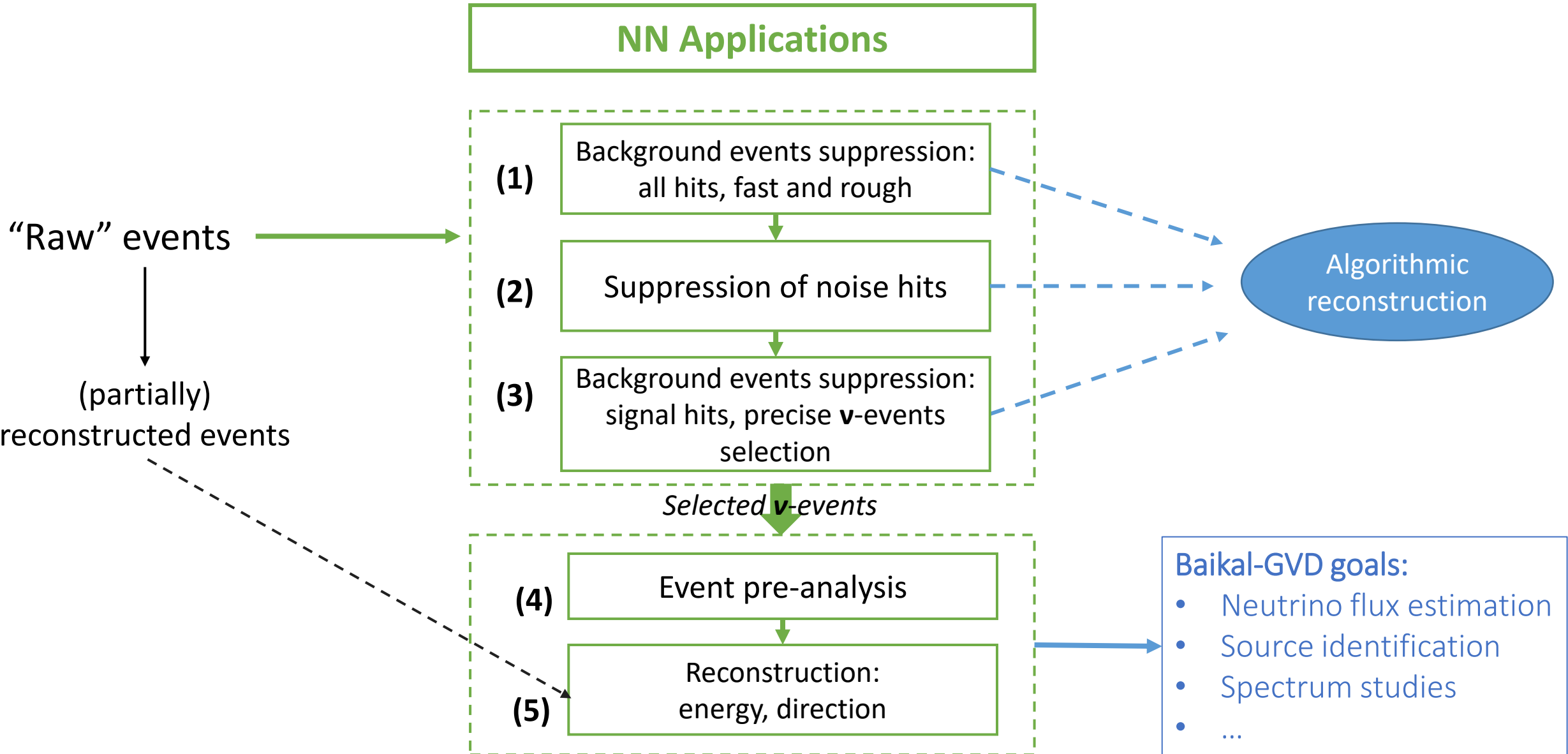
[t.me/AlbertMac280](https://t.me/AlbertMac280)

Code and models:

[github.com/ml-inr/Baikal-ML](https://github.com/ml-inr/Baikal-ML)

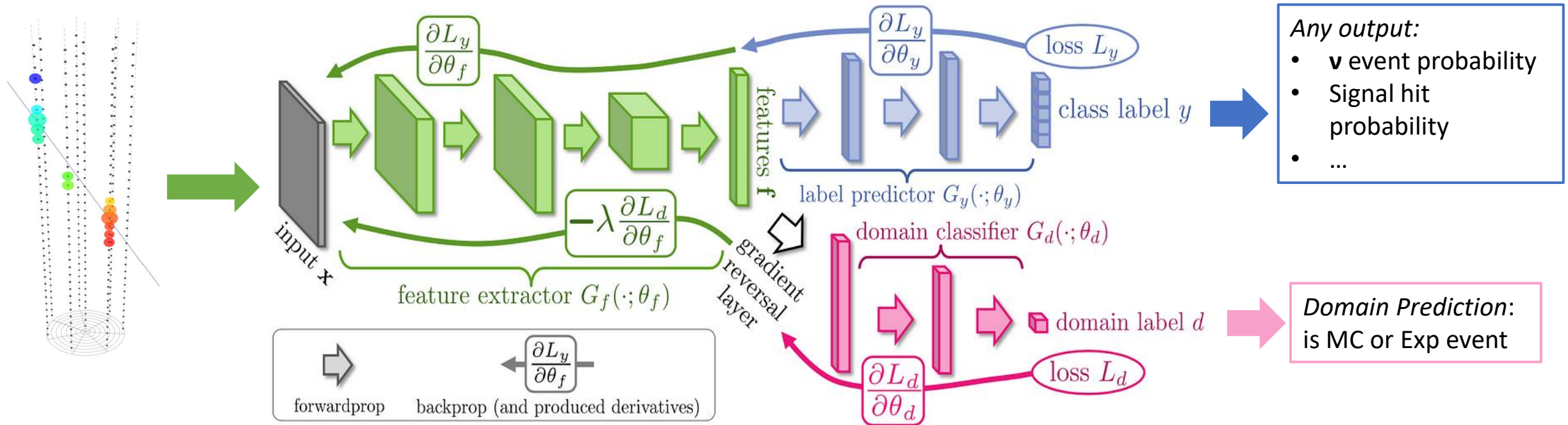
# Backup

# Full ML processing pipeline



# Domain Adaptation Technique

DA scheme:  
domain-adversarial training with a gradient reversal layer



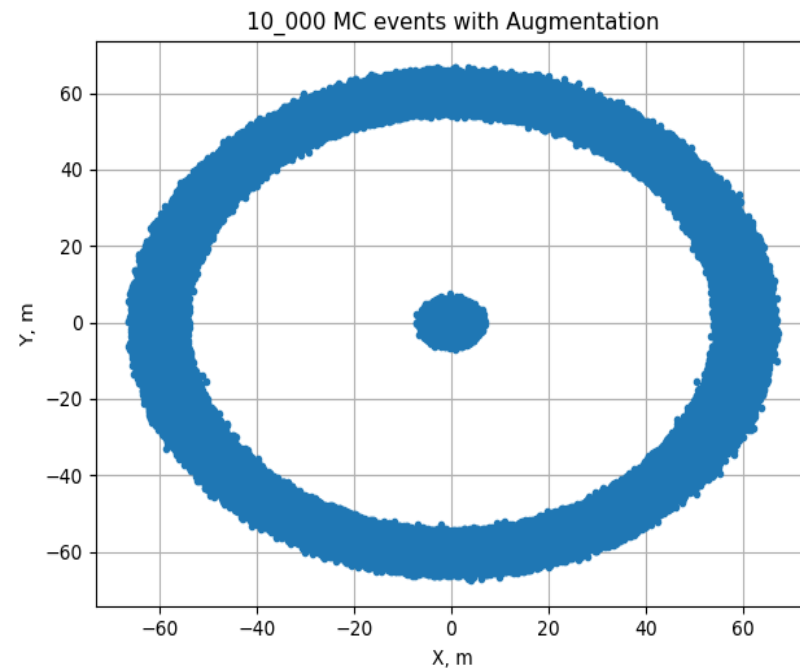
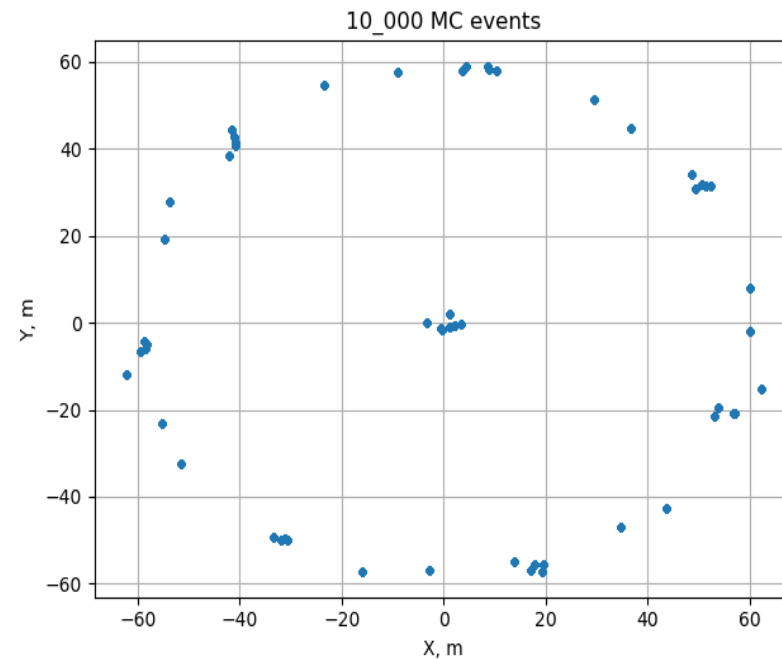
# Fast pre-filter: dataset

## Data augmentations:

- Random gaussian noise to all inputs.

STDs:

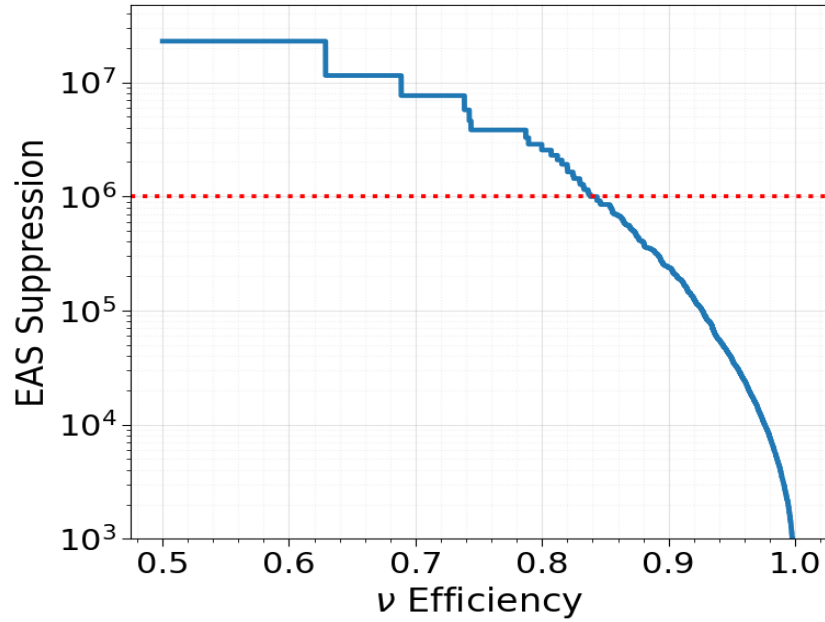
- 1 m for  $x$  and  $y$
- 2m for  $z$
- 5 ns for  $t$
- 0.05 pe for  $Q$
- Random rotations in x-y plane



# $\nu$ -candidate Extractor: MC metrics

Results for data after quality cuts:

- # signal hits  $\geq 8$
- # signal strings  $\geq 3$



*EAS Suppression VS  $\nu$  Efficiency (scanning  $\xi$ )*

*Capable of suppressing EAS by  $10^6$  times while saving  $>80\%$   $\nu$*

*$\nu$  Efficiency VS azimuth angle  $\Theta$  and Energy  
@ suppression factor of  $10^6$*

