

**Warsaw University
of Technology**

Domain adaptation techniques in PID

Michał Kurzynka

Presentation agenda

2

- Particle identification (PID) and data description
- Domain adaptation
- DA Models:
 - DANN
 - WDGRL
- Validation methods
- Final results

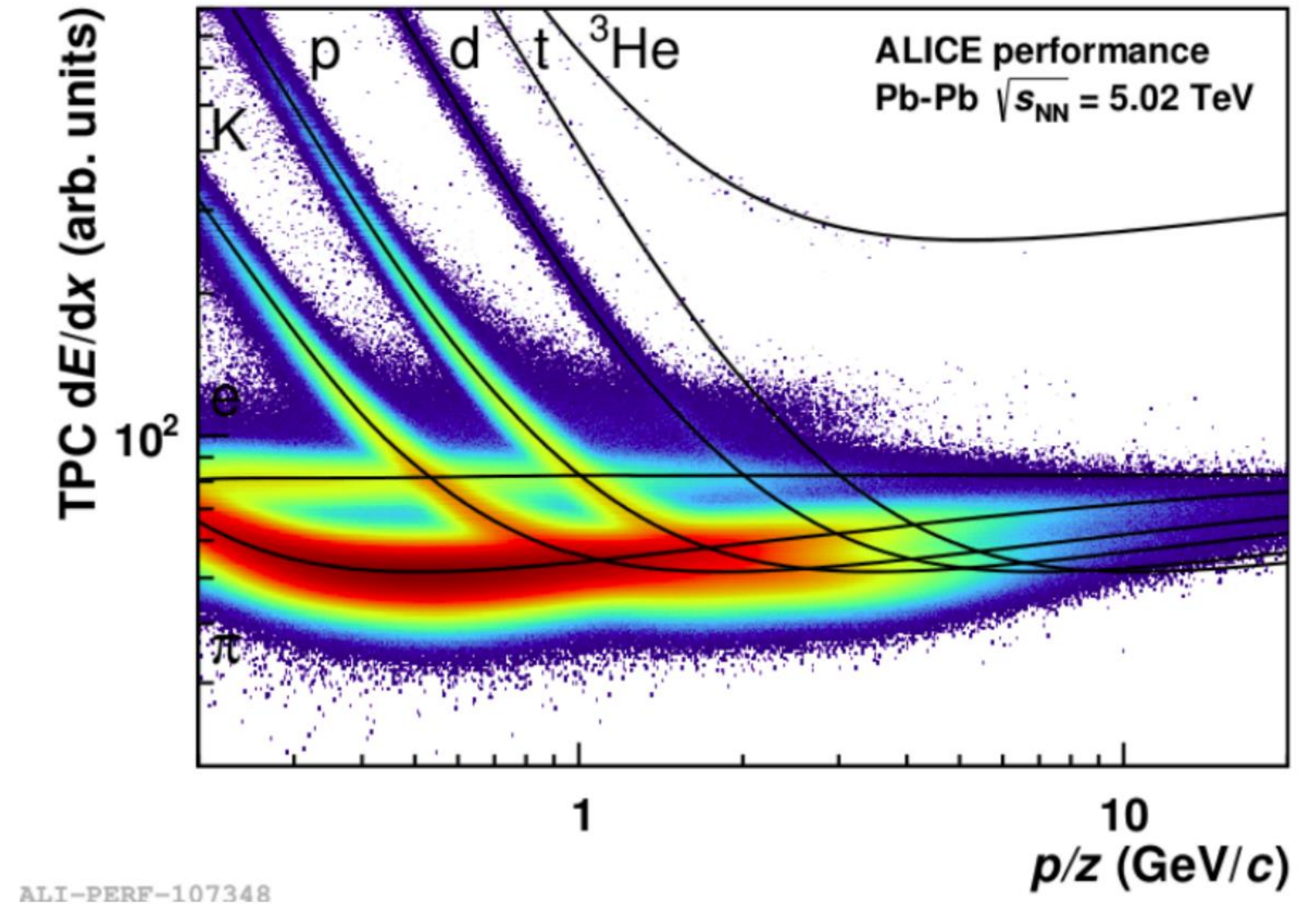
Particle identification (PID)

3

- Particle identification can be done using signals from ALICE detectors such as TPC, TOF and ITS.
- Currently two methods are used:
 - Linear classifiers (production data),
 - Non-linear classifiers (simulation data only).

Linear classifiers

- Build for multiple ranges of momentum.
- Precision drops with the momentum increase.
- Requires time to fine-tune.

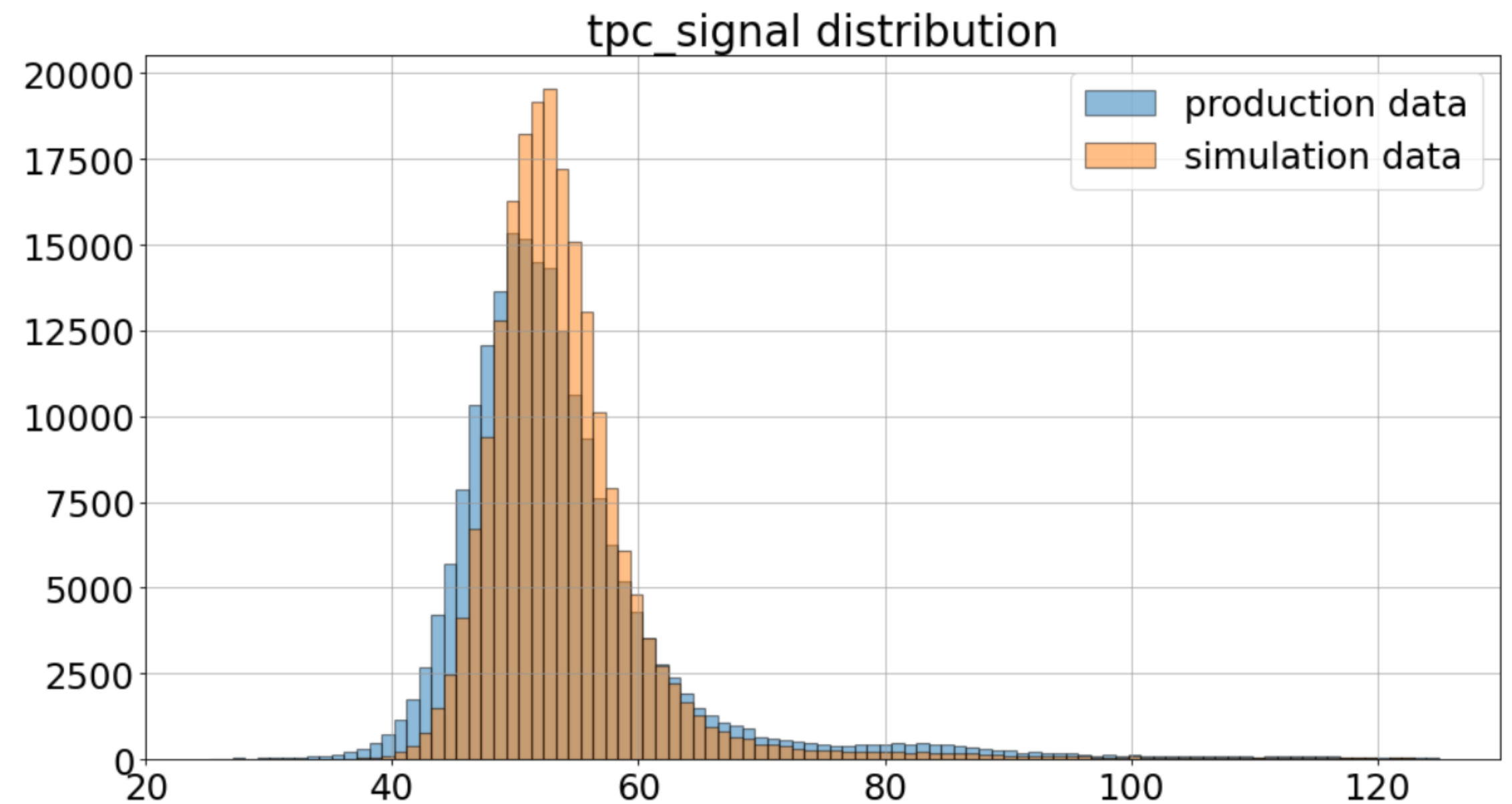
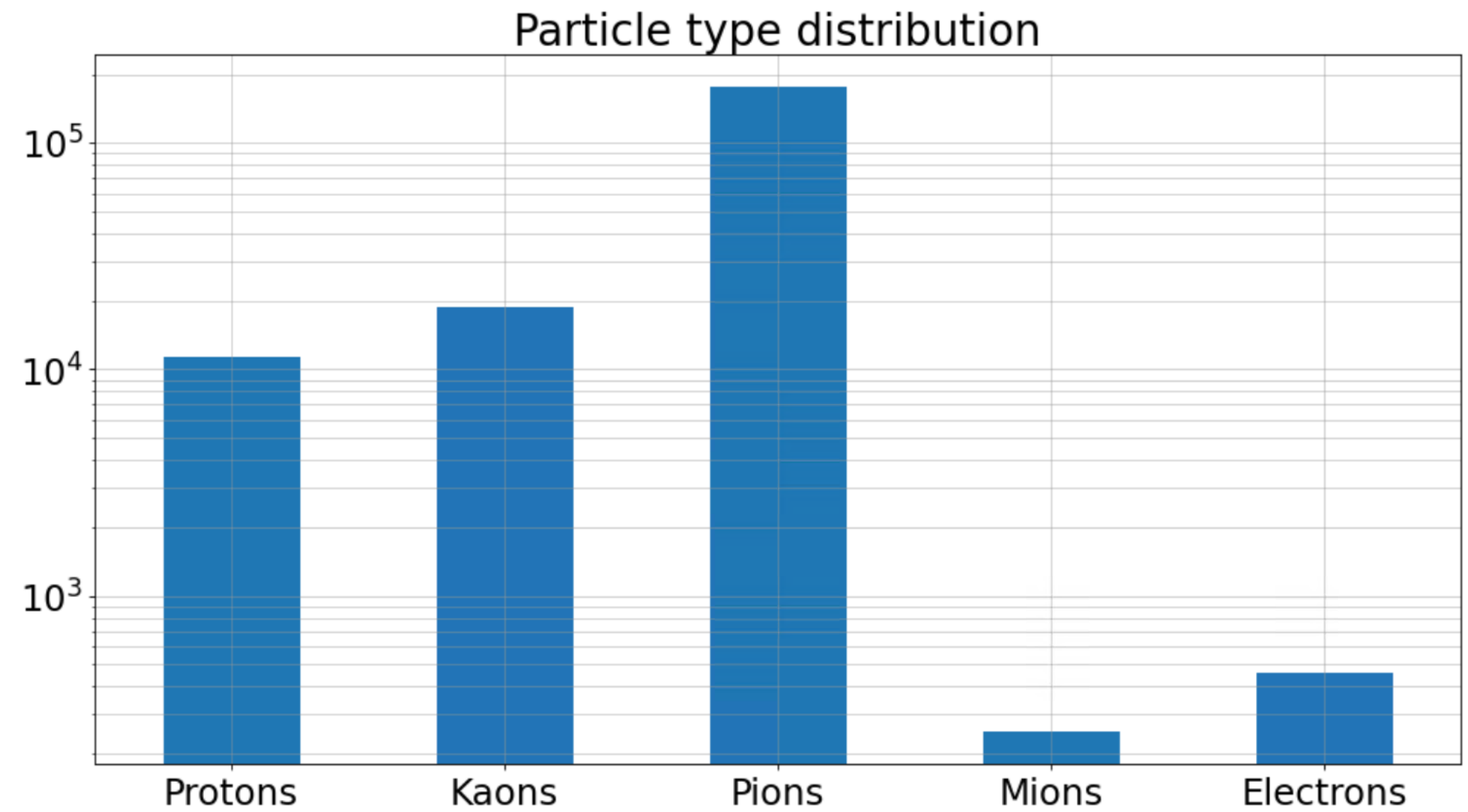


Non-linear classifiers

- Build only on simulation data.
- Require special treatment to classify particles on the production data (reweighting, data distribution fit).
- Can incorporate widely used ML methods such as neural networks.

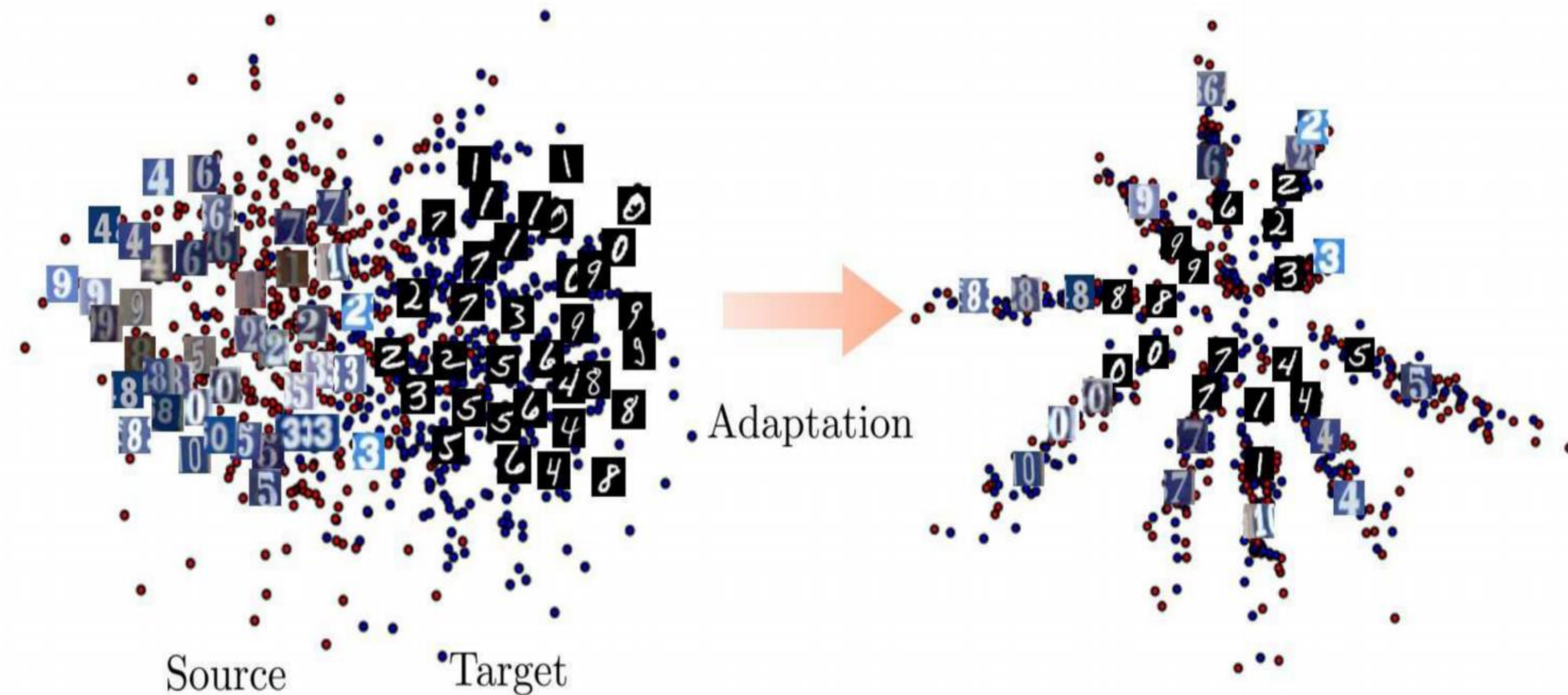
PID - Data

- We use ESD data for PID.
- Data classes are highly imbalanced.
- Two sources of data with different distributions (simulation and production).
- Production data set does not contain labels.

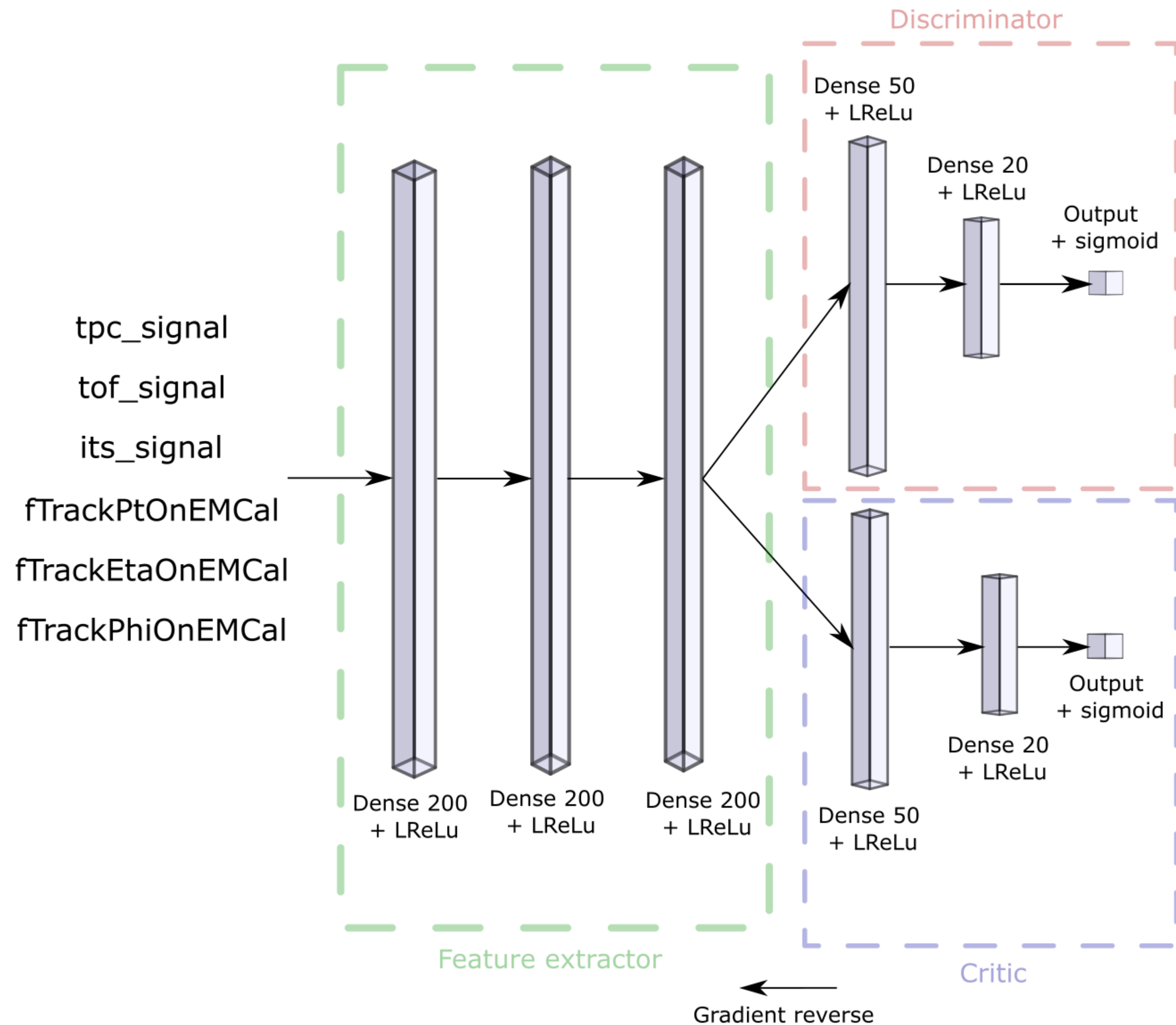


Unsupervised domain adaptation

- The main idea of the unsupervised domain adaptation is to use two sources of data (source domain and target domain).
- There are several variants of unsupervised domain adaptation.
- We will focus on the unsupervised domain adaptation method based on the feature mapping to the common latent space.

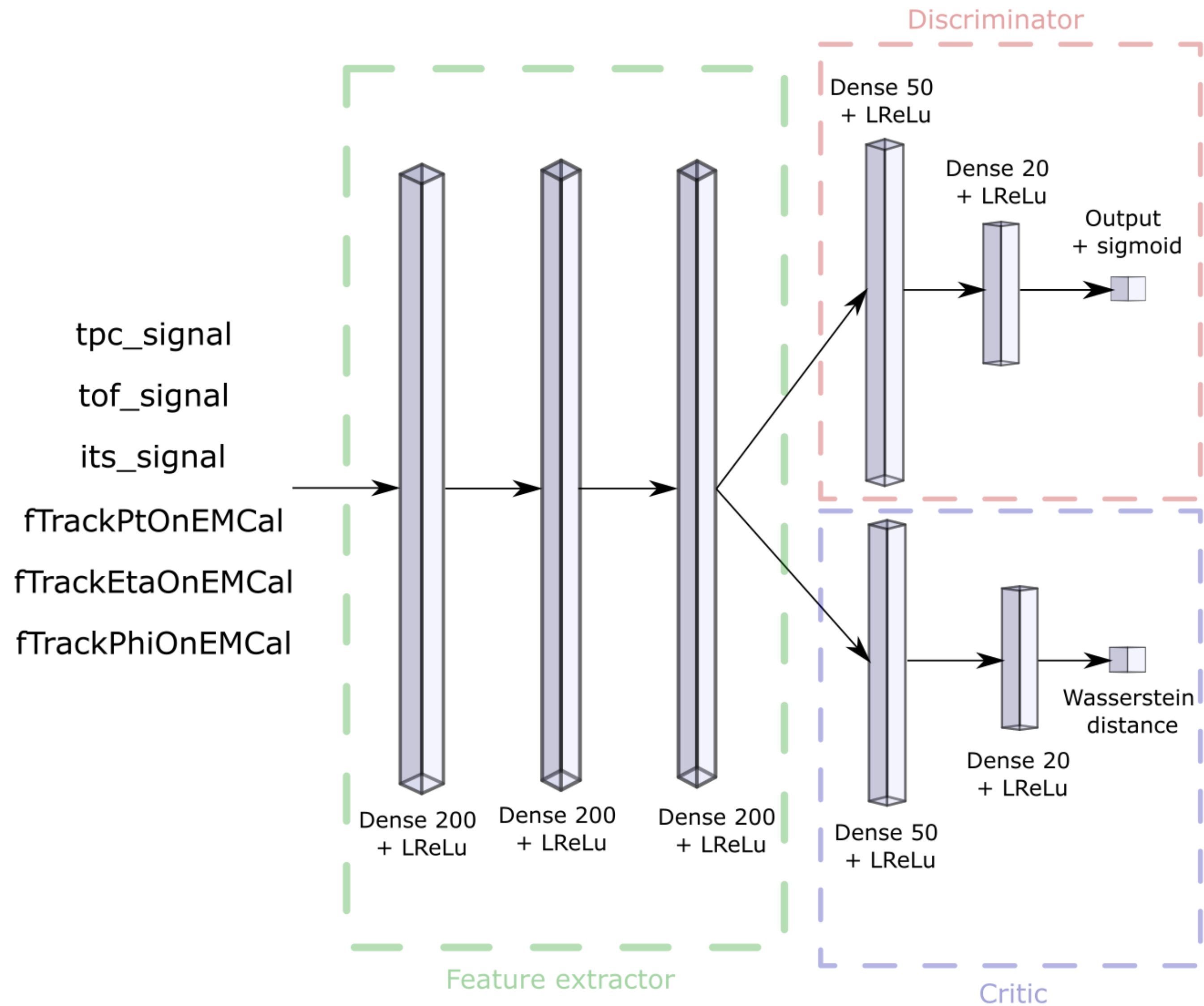


Domain adversarial adaptation: DANN



$$\min_{E,D,C} \{L_D + \lambda L_C\}$$

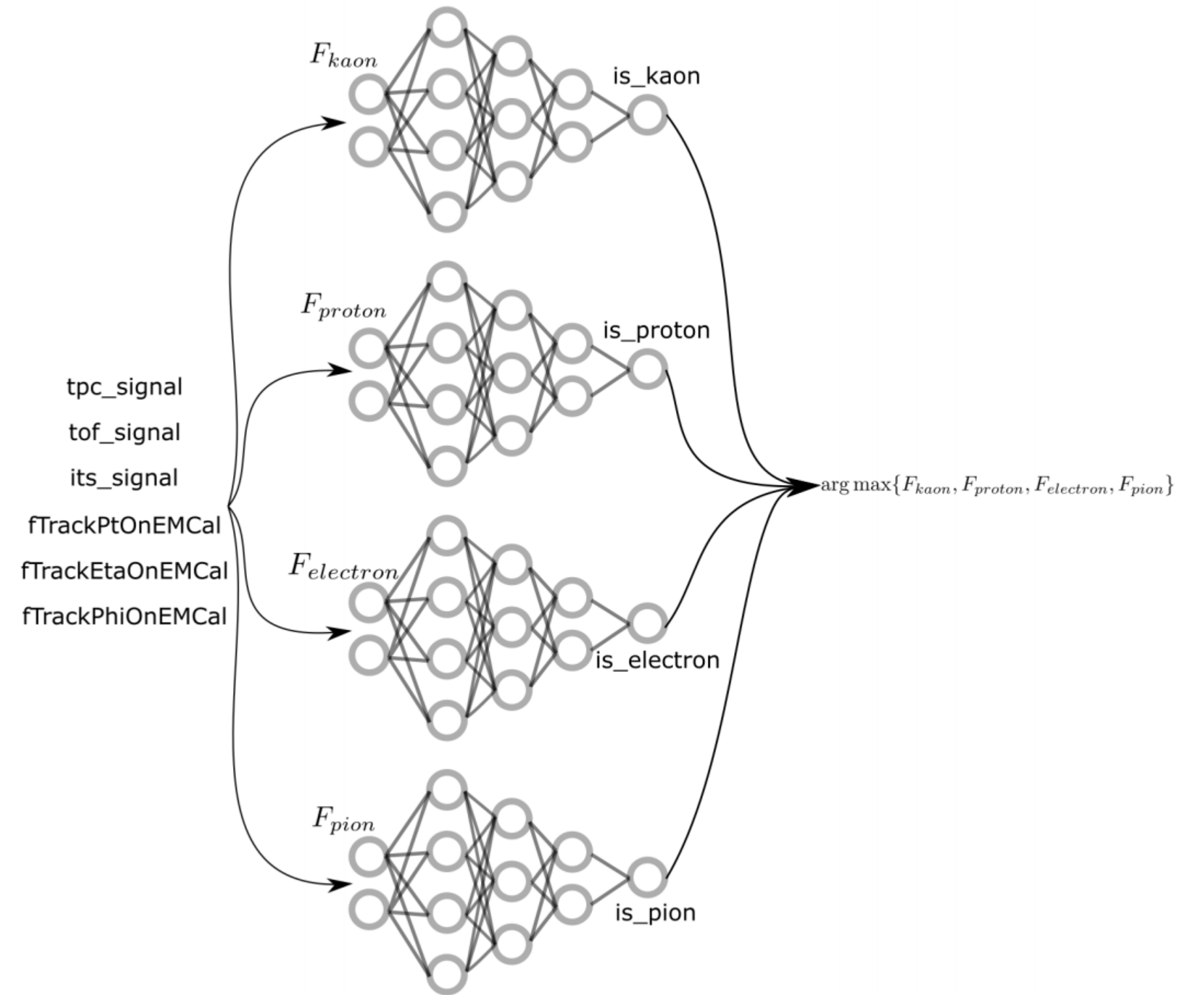
Domain adversarial adaptation: WDGRL



$$\min_{E,D} \{L_D + \lambda \max_C \{L_w - \gamma L_{grad}\}\}$$

Training of the final model

- Create OVA classifier for each particle type.
- Pre-train classifier on the simulation data set taking into account class imbalance.
- Adapt domain.
- Aggregated all OVA models into one classifier.



- Due to lack of labels in the production data we split validation procedure into two tasks:
 - Invariant mass validation on the production data for low momentum particles.
 - Validation on the perturbed data set.
- As a validation metrics due to severe class imbalance we propose:
 - Precision – ratio of true positive to all positively classified:

$$PPV = \frac{TP}{TP + FP}$$

- F1 score – harmonic mean of precision and recall (TP and all positive):

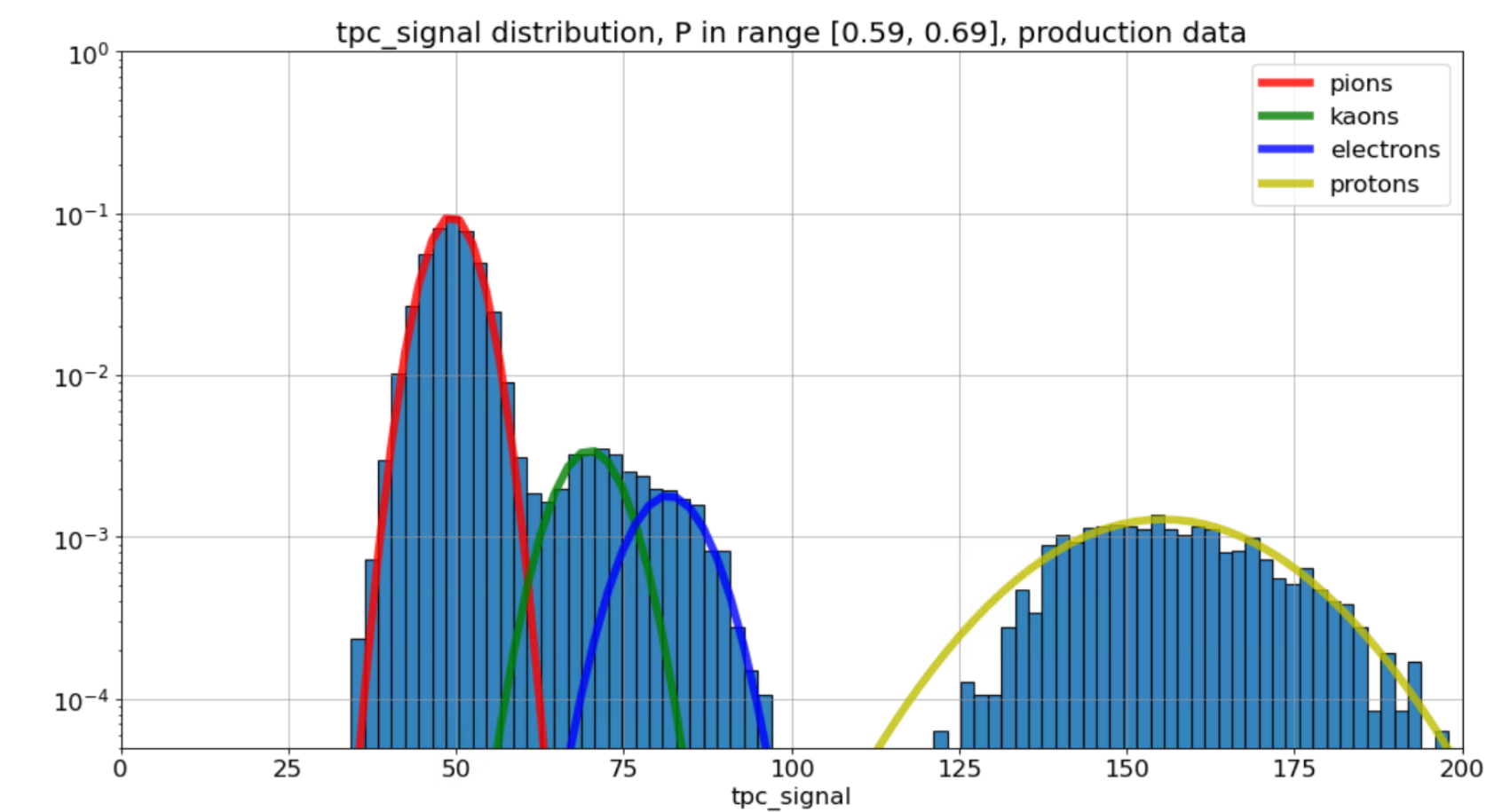
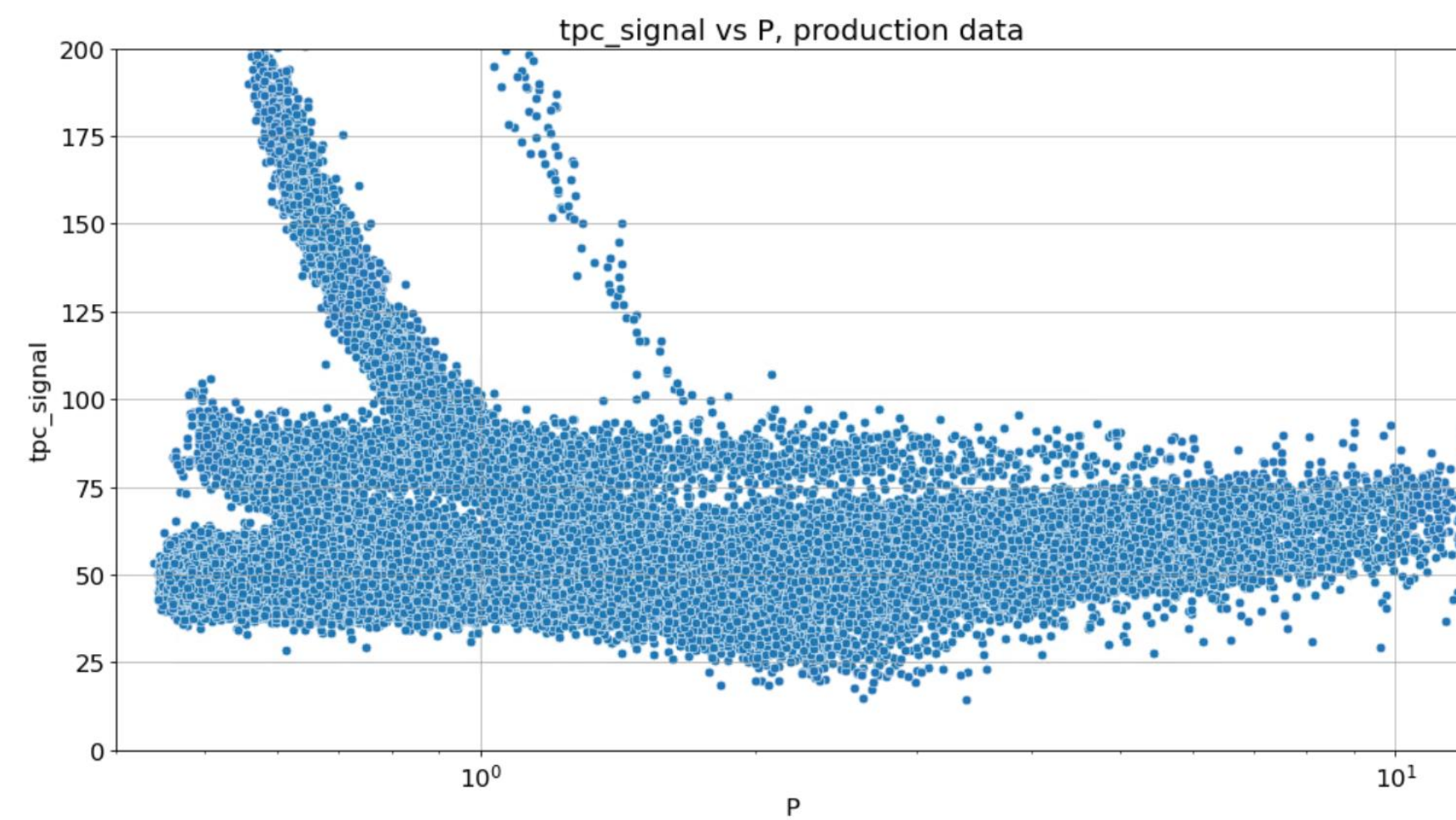
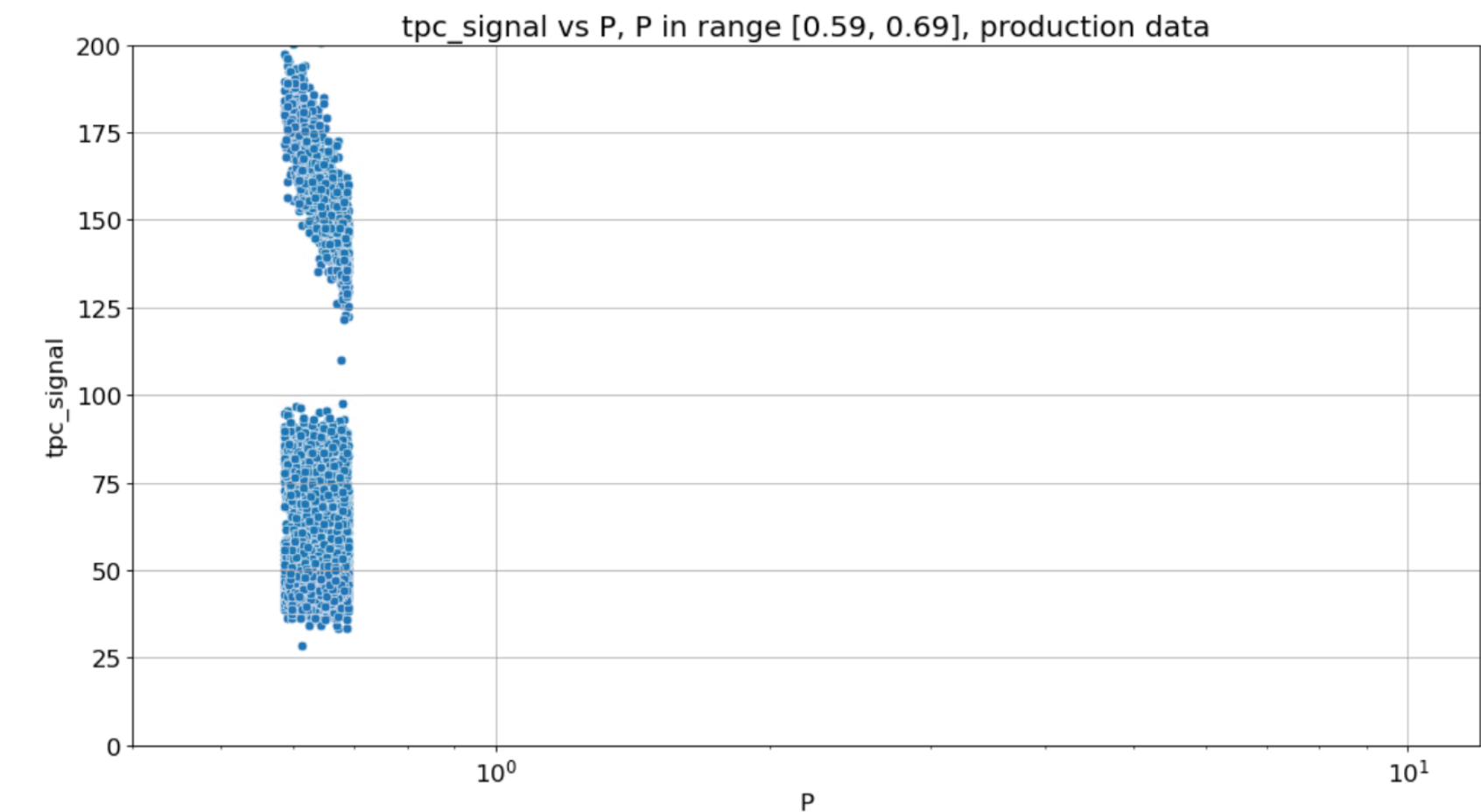
$$F_1 = 2 \frac{PPV \cdot TPR}{PPV + TPR}$$

- Efficiency:

$$E = \frac{TP_{prod}}{\int_{-3\sigma}^{3\sigma} f_{\mu,\sigma}(x)}$$

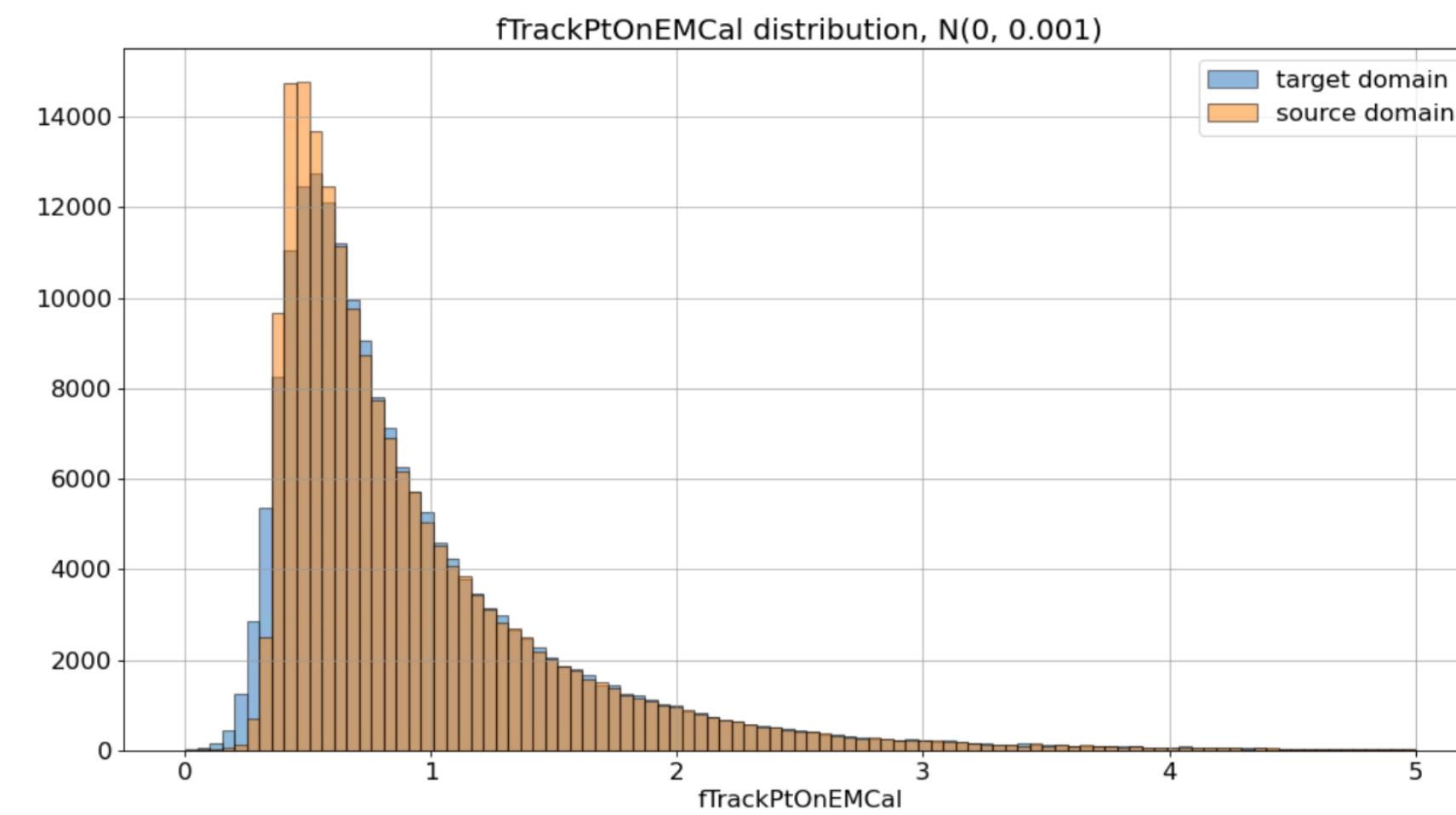
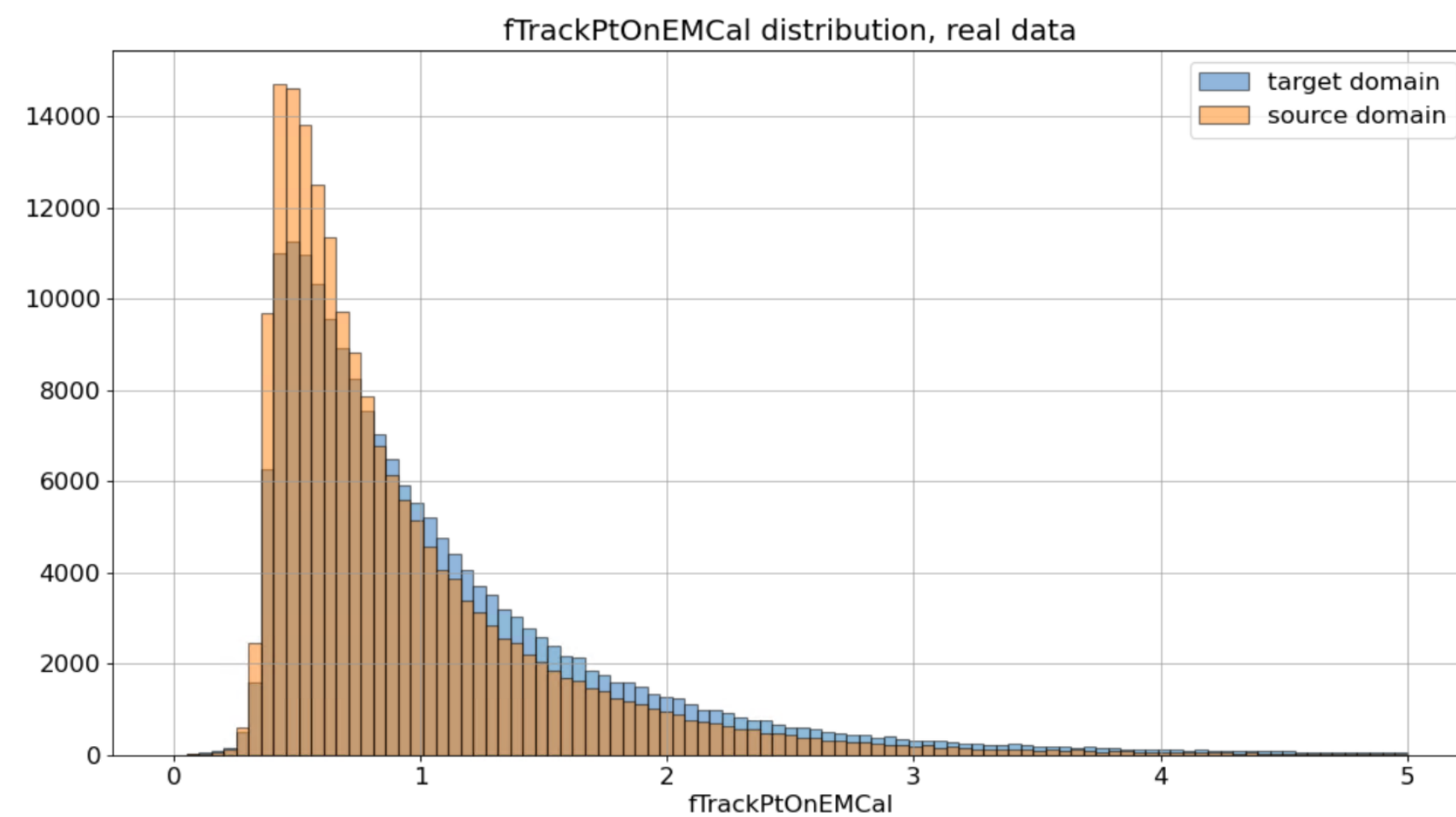
Mass invariant (Linear classifiers)

- Only for low momentum particles ($p < 0.9$ GeV).
- Split particles into several momentum ranges,
- For each range fit Gaussian curves to the production data,
- Check if positively classified particles lie beneath corresponding curve.



Perturbed data set

- Create data set with similar attributes distribution to the production data set.
- We can use simulation data with noise from the natural distribution that mimics target domain distortions.
- Perturbed data set contain true labels because it is created using simulation data set.



Final results

- We have compared DANN and WDGRL performance on the perturbed data to determine which model should be used on the production data.

Precision

Particle type	DANN	WDGRL
Pions	97.6%	97.8%
Kaons	68.7%	74.8%
Electrons	95.6%	96.7%
Protons	84.9%	86.9%

F1 score

Particle type	DANN	WDGRL
Pions	97.3%	97.5%
Kaons	73.4%	75%
Electrons	93.8%	97%
Protons	85%	89.5%

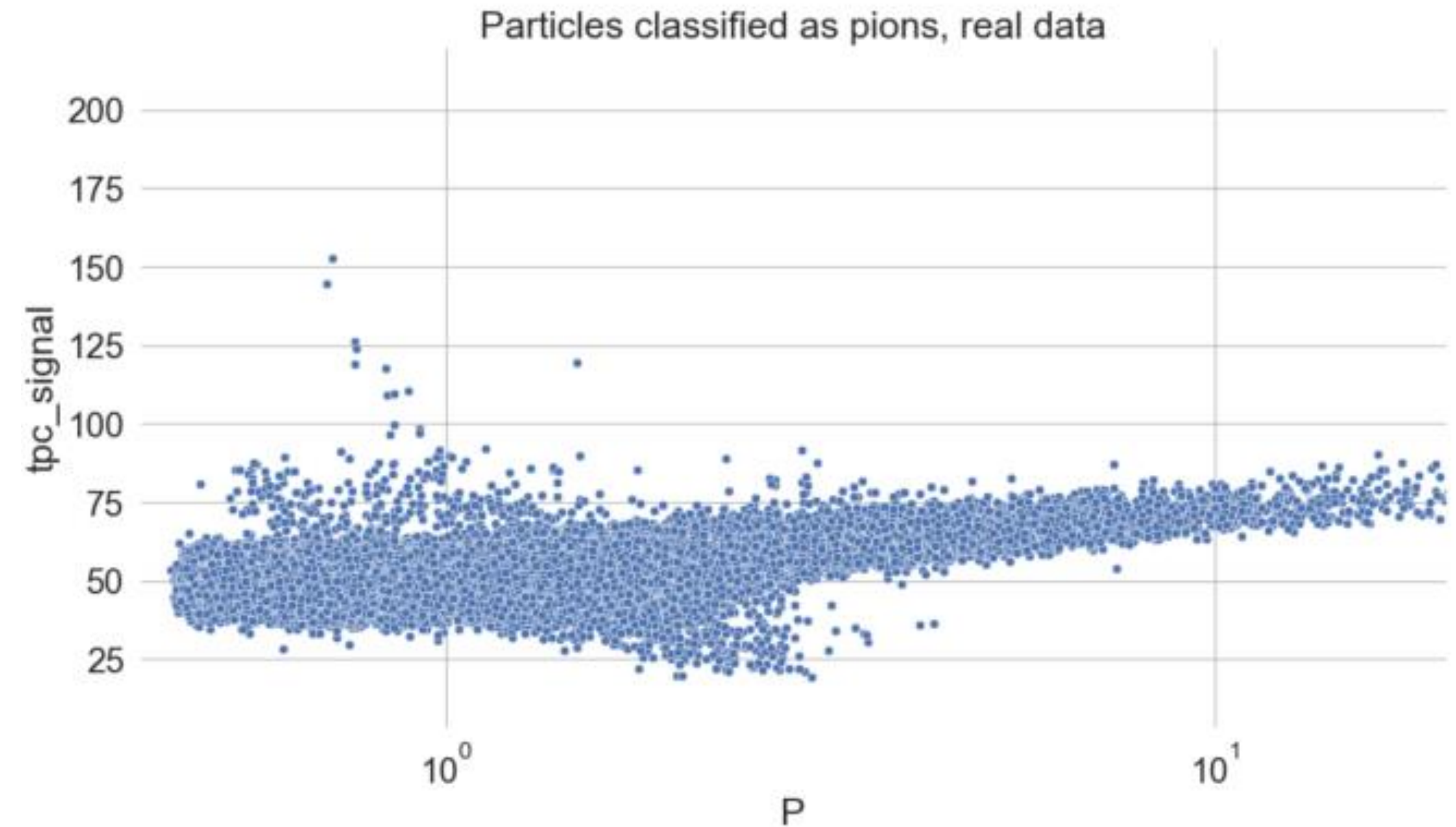
Final results: Pions

- Classification quality on perturbed data set:

Precision	97.8%
F1 score	97.5%

- Classification quality on production data set for particles with $p < 0.9$ GeV:

Precision	99.56%
Efficiency	0.91



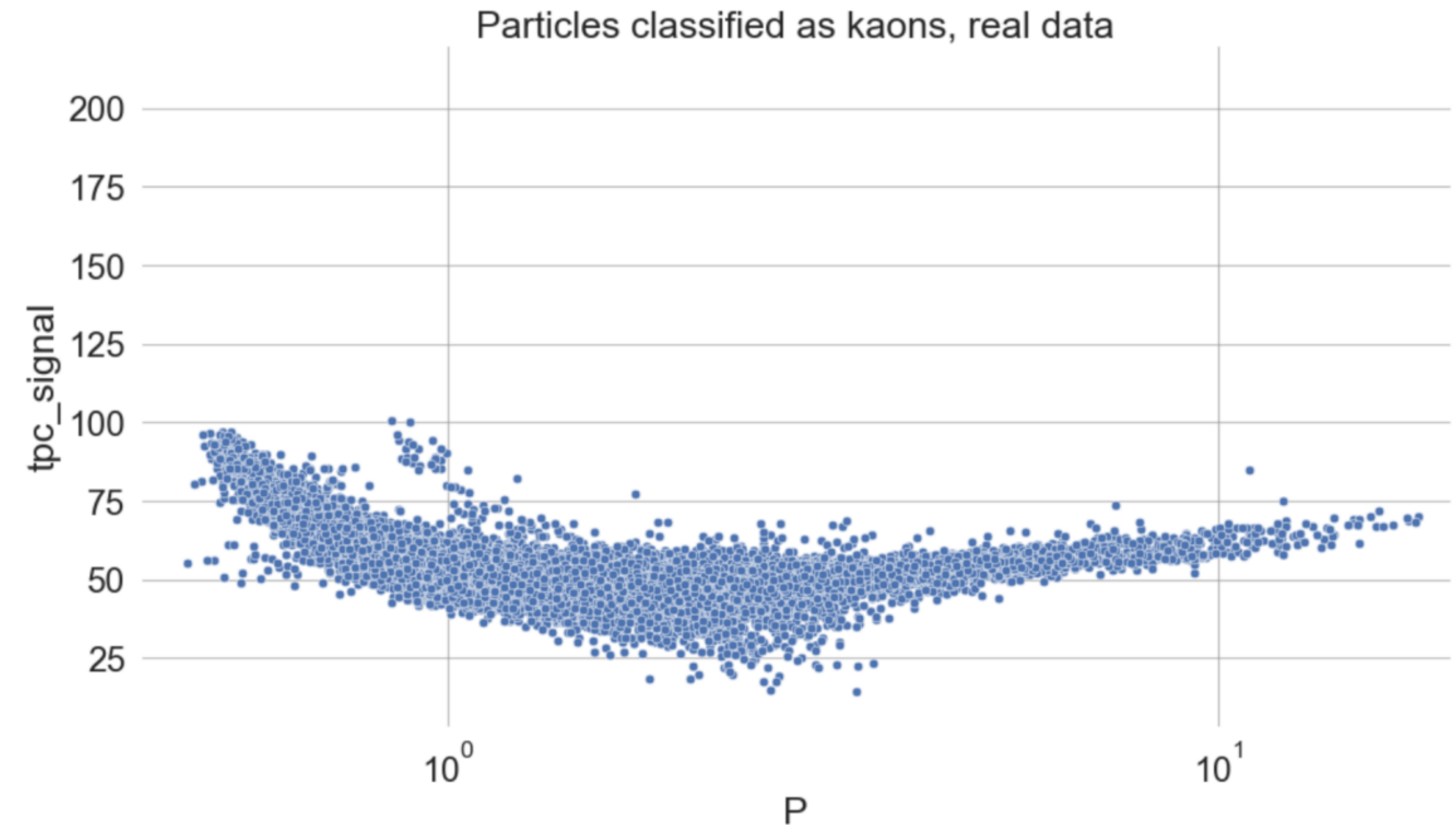
Final results: Kaons

- Classification quality on perturbed data set:

Precision	74.8%
F1 score	75%

- Classification quality on production data set for particles with $p < 0.9$ GeV:

Precision	98.68%
Efficiency	0.98



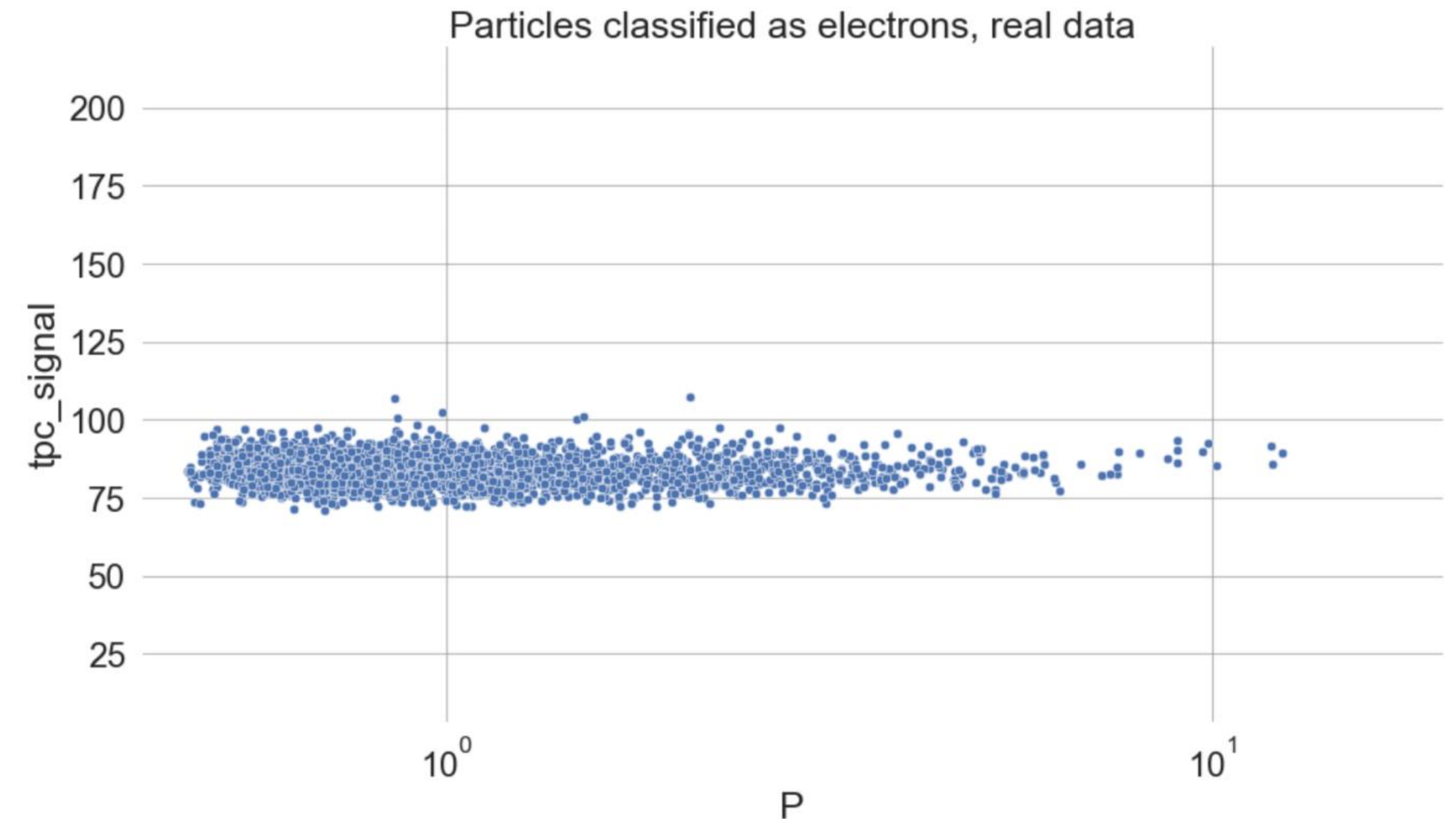
Final results: Electrons

- Classification quality on perturbed data set:

Precision	96.7%
F1 score	97%

- Classification quality on production data set for particles with $p < 0.9$ GeV:

Precision	98.69%
Efficiency	0.98



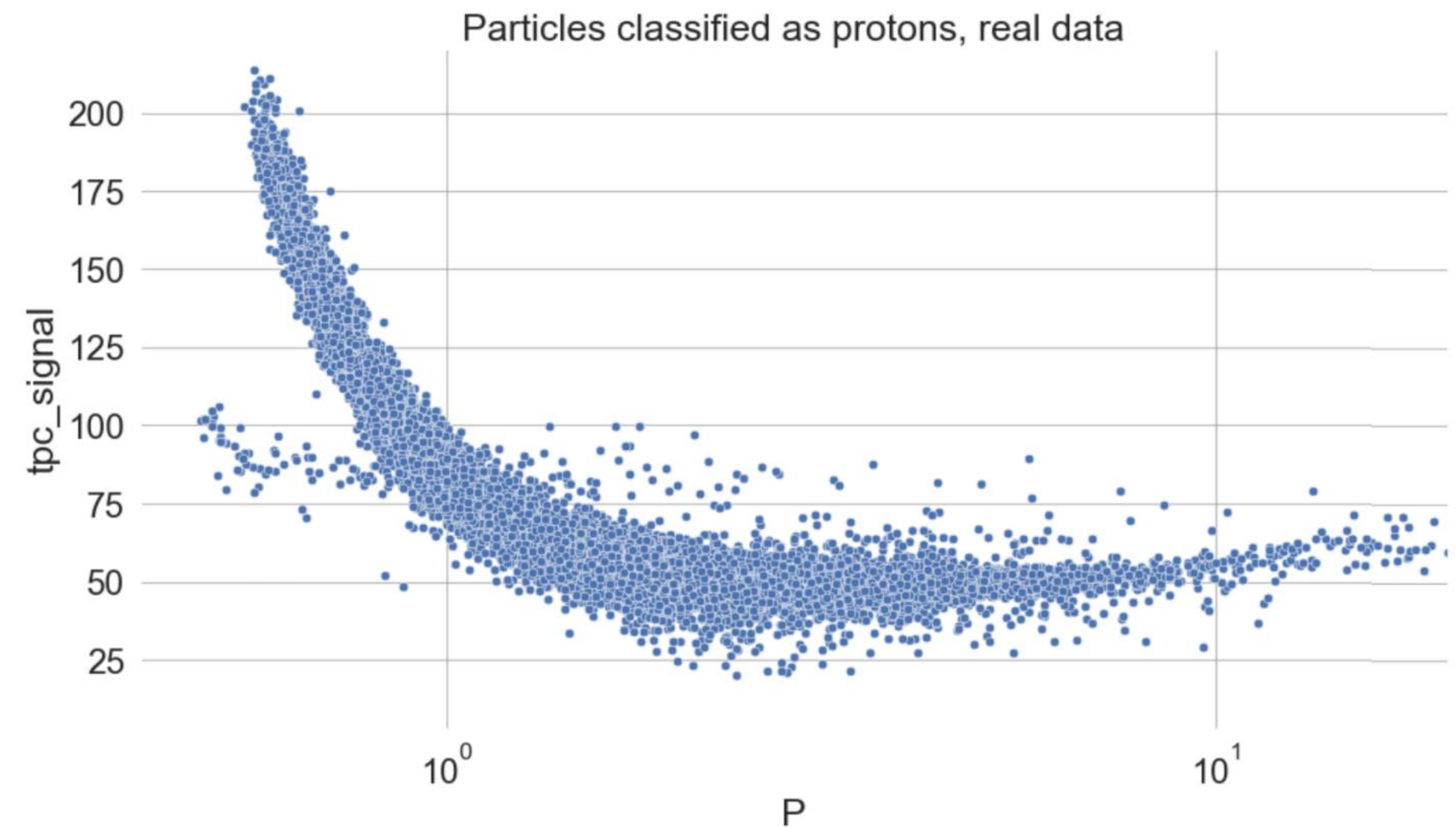
Final results: Protons

- Classification quality on perturbed data set:

Precision	87%
F1 score	90%

- Classification quality on production data set for particles with $p < 0.9$ GeV:

Precision	98.35%
Efficiency	0.98



- New method for PID using both sources of ALICE data.
- WDGRL method yields in this case much better results and achieve more stable training process.
- Splitting the training process into two stages helps with severe class imbalance.
- Further steps:
 - Further evaluation of model on the production data.
 - Evaluation of the model on the samples with high purity.

Collaborators

20

- Tomasz Piotr Trzciński, Warsaw University of Technology
- Georgy Kornakov, Warsaw University of Technology
- Kamil Rafał Deja, Warsaw University of Technology

Backup slides

Wasserstein distance and 1-Lipshitz constrain

- Loss of the critic consists of Wasserstein distance and 1-Lipshitz constrain (limiting the gradient of function).

$$\|f_w(x_1) - f_w(x_2)\| \leq \|x_1 - x_2\|$$

$$L_{grad}(h) = (\|\nabla_{\theta_w} f_w(\hat{h})\|_2 - 1)^2$$

- If the constrain is satisfied we can approximate Wasserstein distance as:

$$L_w = \frac{1}{n^s} \sum_{x^s \in X^s} (f_w(f_e(x^s))) - \frac{1}{n^t} \sum_{x^t \in X^t} (f_w(f_e(x^t)))$$

