

2<sup>nd</sup> EuCAIFCon, Cagliari, Italy

# **DINAMO: Dynamic and INterpretable Anomaly MOnitoring for Large-Scale Particle Physics Experiments**



# Arsenii Gavrikov<sup>1, \*</sup>, Julián García Pardiñas<sup>2</sup>, Alberto Garfagnini<sup>1</sup>

<sup>1</sup>University & INFN Padova, Italy; <sup>2</sup>Massachusetts Institute of Technology (MIT), USA

\*agavriko@cern.ch

### I. Data quality monitoring

- ✤ In particle physics experiments, DQM is a gatekeeper against corrupted, anomalous data
- Traditional DQM is manual: shifters compare data with references provided by experts:
  - High cost of person power
  - Limited accuracy
  - $\clubsuit$  Any detector or software updates  $\rightarrow$ challenges to adapt to changes in operational conditions



- ✤ An automated DQM to be:
  - ♦ Accurate → as high as possible
  - Specific  $\rightarrow$  where is the problem
  - ★ Interpretable → "why" the algorithm decided so
  - $\diamond$  **Dynamic**  $\rightarrow$  adaptability to changing conditions
  - ↔ **Fast**  $\rightarrow$  although analyzing vast amount of data
    - Automated DQM will help to reduce shifters burden and the amount of inconsistencies

# **III. Synthetic data generator**

Realistic synthetic data:

- **Solutions** Based on **1d Gaussian distributions** as histograms
- Main focus is the modeling of changes in conditions

Implemented features:

Slow drifts  $\rightarrow \mu$  evolve gradually (sinusoidal drift)



• Abrupt shifts  $\rightarrow$  sudden changes in  $\mu$  or/and  $\sigma$ 

### **II. DINAMO: Dynamic and Interpretable Anomaly Monitoring** [1]



Build the reference dynamically to account to the changing conditions based on Exponentially Weighted Moving Average (EWMA) method:

- 1. Include also **the weights of each run** according to their statistical noise:  $\omega_t = \frac{1}{2}$
- 2. Iterative update of sum of weights:  $W_{t+1} = \alpha W + (1 \alpha)\omega_t$
- 3. Iterative update of weighted sum:  $S_{t+1} = \alpha S_t + (1 \alpha)\omega_t X_t$
- 4. Iterative update of **sum of weighted squared residuals** to

#### **Substitutes** the EMWA-based part with the transformer encoder:



- Training via Online learning:
- Takes mini-batch of K good last runs
  - **Per each run** in the mini-batch we build **a context of** *M* preceding good runs

- Varying events statistics  $\rightarrow$  initial number of events is sampled from a uniform distribution
- $\Rightarrow$  Systematic uncertainty  $\rightarrow$ accidental increase or decrease of events in the right half of the histogram's window using binomial distribution



bad run

Anomalies:

**\* extra distortion** in μ or/and σ \* dead bins: random number of bins (up to 20) get missed

good run

content

### **IV. Main metrics**

 $\Rightarrow$  Balanced accuracy  $\rightarrow$  to balance the uneven class ratio

 $\Rightarrow$  Adaptation time  $\rightarrow$  an average amount of good runs to be misclassified before the algorithm adapts



model the uncertainty:  $S_{\sigma,t+1} = \alpha S_{\sigma,t} + (1-\alpha)\omega_t (X_t - \mu_t)^2$ 

5. Compute the reference mean and uncertainty:

$$u_{t+1} = \frac{S_{t+1}}{W_{t+1}} \qquad \sigma_{\mu_{t+1}} = \sqrt{\frac{S_{\sigma,t+1}}{W_{t+1}}}$$

Compute test statis 6. to compare the two:

stics  

$$\chi^2_{
u} = \frac{1}{N_b} \sum_{i=1}^{N_b} \frac{\left(\tilde{x}_{i,j} - \hat{\mu}_{i,j}\right)^2}{\sigma^2_{\tilde{x}_{i,j},p} + \hat{\sigma}^2_{\mu_{i,j}}}$$

Parameter α is a hyperparameter to control EWMA

- Learning to predict the bin-by-bin means and widths for future runs
- Anomaly detection as follows:
- **Context creation**: for a new run, we identify up to *M* most recent good runs as input to the transformer
- **Predict reference**: output  $\mu$  and  $\sigma$  by the model
- **Anomaly score**: compute the test statistic the same as with DINAMO-S

• Gaussian negative log-likelihood as loss:  $NLL = \frac{1}{2N_bK} \sum_{k=1}^{K} \sum_{j=1}^{N_b} \left[ \left( \frac{\tilde{x}_{k,j} - \hat{\mu}_{k,j}}{\hat{\sigma}_{k,j}} \right)^2 + \log(\hat{\sigma}_{k,j}^2) \right]$ 

 $\Rightarrow$  Uncertainty coverage  $\rightarrow$  how well the actual variability of the good runs is described using Jaccard distance between the true good runs and the predicted references in the z-score space

### ◆ More complex → slower performance but **more accurate** and quicker in **adapting**

# V. Results with synthetic data



- The comprehensive evaluation on synthetic data demonstrates that both **DINAMO algorithms** successfully address the core challenges of automated DQM
- DINAMO-ML outperformed DINAMO-S in all our metrics thanks to it more complex nature
- OINAMO-ML is particularly better in adaptation speed and in

#### balanced accuracy

**Aggregated results**: 1000 synthetic datasets with different seeds



Bal. acc.  $\uparrow$  Specif.  $\uparrow$  Sensit.  $\uparrow$  Jaccard D. for  $\sigma \downarrow$  Adapt. time  $\downarrow$ 

DINAMO-S	$0.947\substack{+0.020\\-0.033}$	$0.943\substack{+0.028\\-0.058}$	$0.956\substack{+0.029\\-0.075}$	$0.139_{-0.041}^{+0.069}$	$2.02^{+3.24}_{-1.13}$
DINAMO-ML	$0.966\substack{+0.012\\-0.018}$	$0.969\substack{+0.015\\-0.037}$	$0.966\substack{+0.024\\-0.044}$	$0.134\substack{+0.057\\-0.028}$	$1.61\substack{+0.87 \\ -0.61}$

#### VI. Conclusions

- ◆ We present **DINAMO** → a novel approach to **automate DQM** for large particle physics experiments:
  - EWMA-based, "standard" version: DINAMO-S
  - Transformer encoder-based, machine learning version: DINAMO-ML
- OINAMO-ML outperforms the standard version in all our metrics

#### Key advantages:

Adaptability to changes in operational conditions Interpretability through the references' dynamic creation (+ uncertainties)

**Relative simplicity** to enhance maintainability

DINAMO-S is already being commissioned at the LHCb experiment for offline DQM

[1] A. Gavrikov, J. García Pardiñas, A. Garfagnini, ''DINAMO: Dynamic and INterpretable Anomaly MOnitoring for Large-Scale Particle Physics Experiments", arXiv:2501.19237 (2025)

Arsenii Gavrikov is supported by the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie Grant Agreement No.101034319 and from the European Union – NextGenerationEU.

