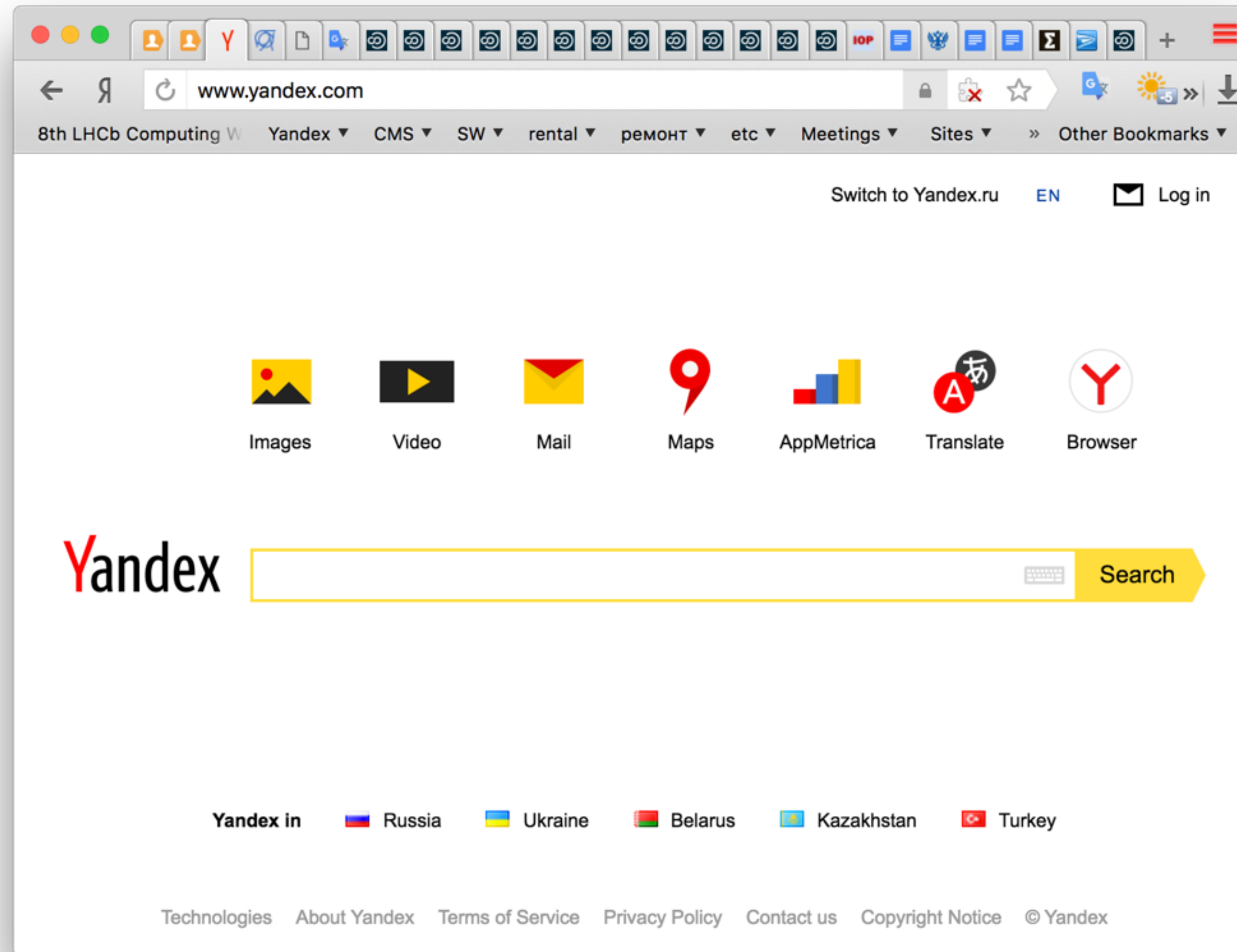


Yandex Activities in the CERN OpenLab Scope

CERN, Dec.9, 2016

Fedor Ratnikov on behalf of the group

Motivation



Motivation

- Yandex School for Data Analysis

 - › education and research

- YSDA is a member of the LHCb collaboration, so has a full access to LHCb data

- However YSDA tries to develop generic approaches which are applicable to other experiments and different branches of sciences

- CERN OpenLab is an excellent venue for such projects

Projects in This Presentation



| Data Popularity

| Grid Simulation

| Detector Monitoring

| Data Anomaly Detection

Status Update

Data Popularity



Data Popularity Problem

Datasets

- RAW, Primary RECO, Central Skims, Group Skims, User Skims, AOD, User Ntuples, ...

Storage

- tape: 0, 1, 2 replicas

 - › long term storage, expensive to extract (write once, read never)

- disks: 0, 1, 2, ... replicas

 - › readily available for processing

Target

- Optimise replicas distribution over GRID sites

 - › make popular datasets readily available on necessary sites

 - › keep total used disk space within agreed site pledges

Machine Learning Approach

Based on LHCb data popularity

Use data logs

Short term forecast

> ~ 1 month

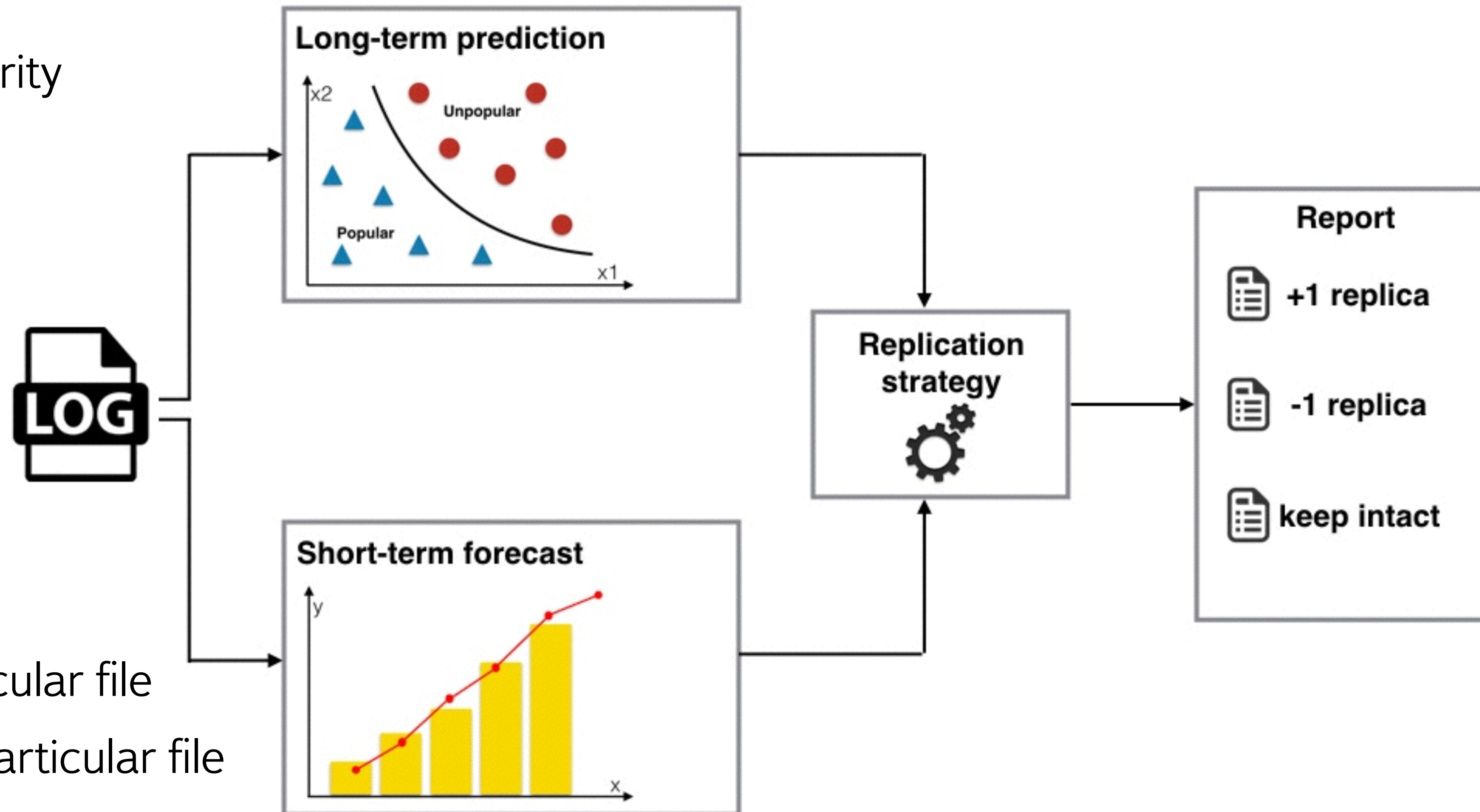
Long term prediction

> \approx 2 months

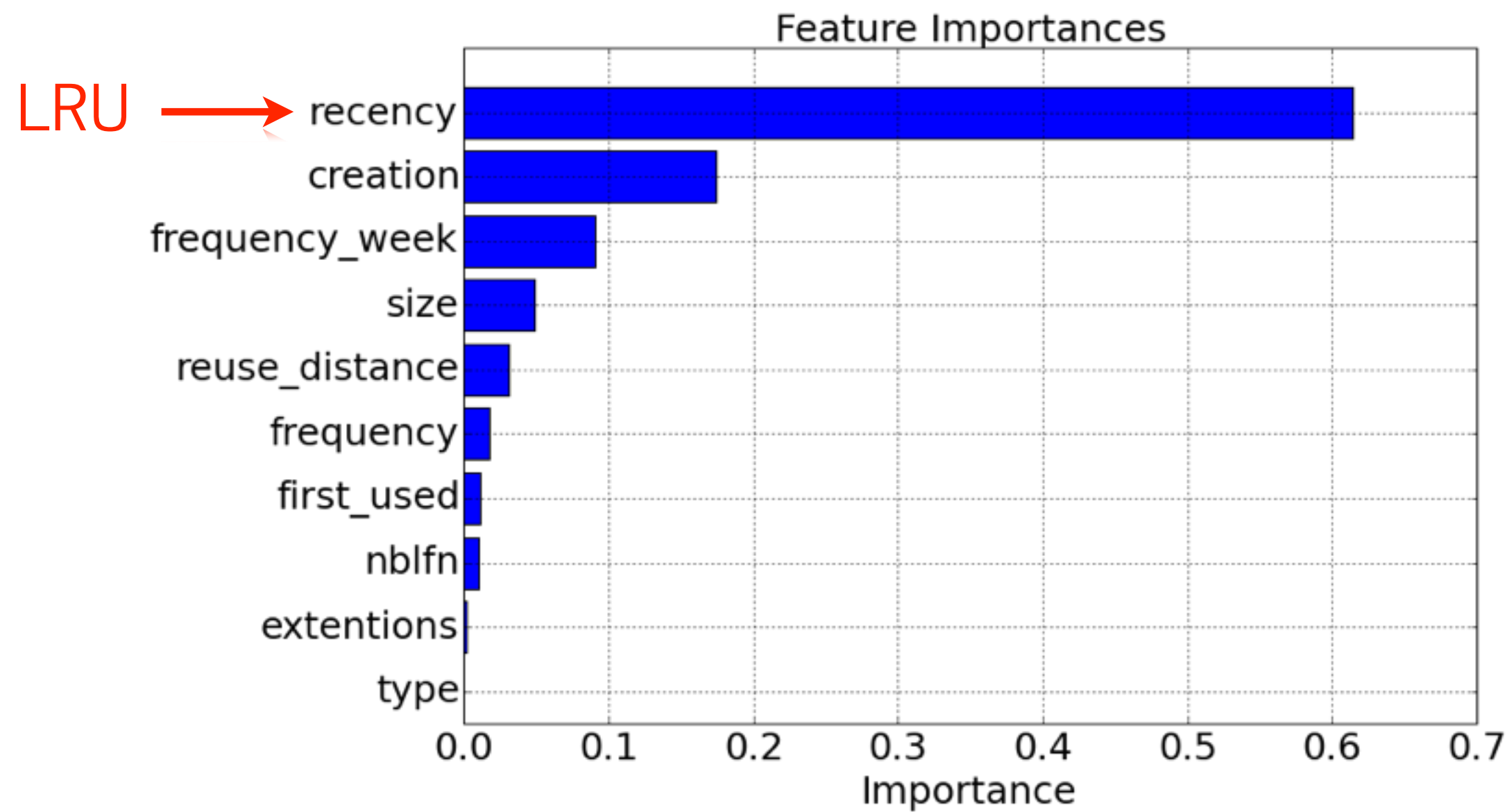
Decisions

> add one replica of particular file

> remove one replica of particular file



Long-term Prediction

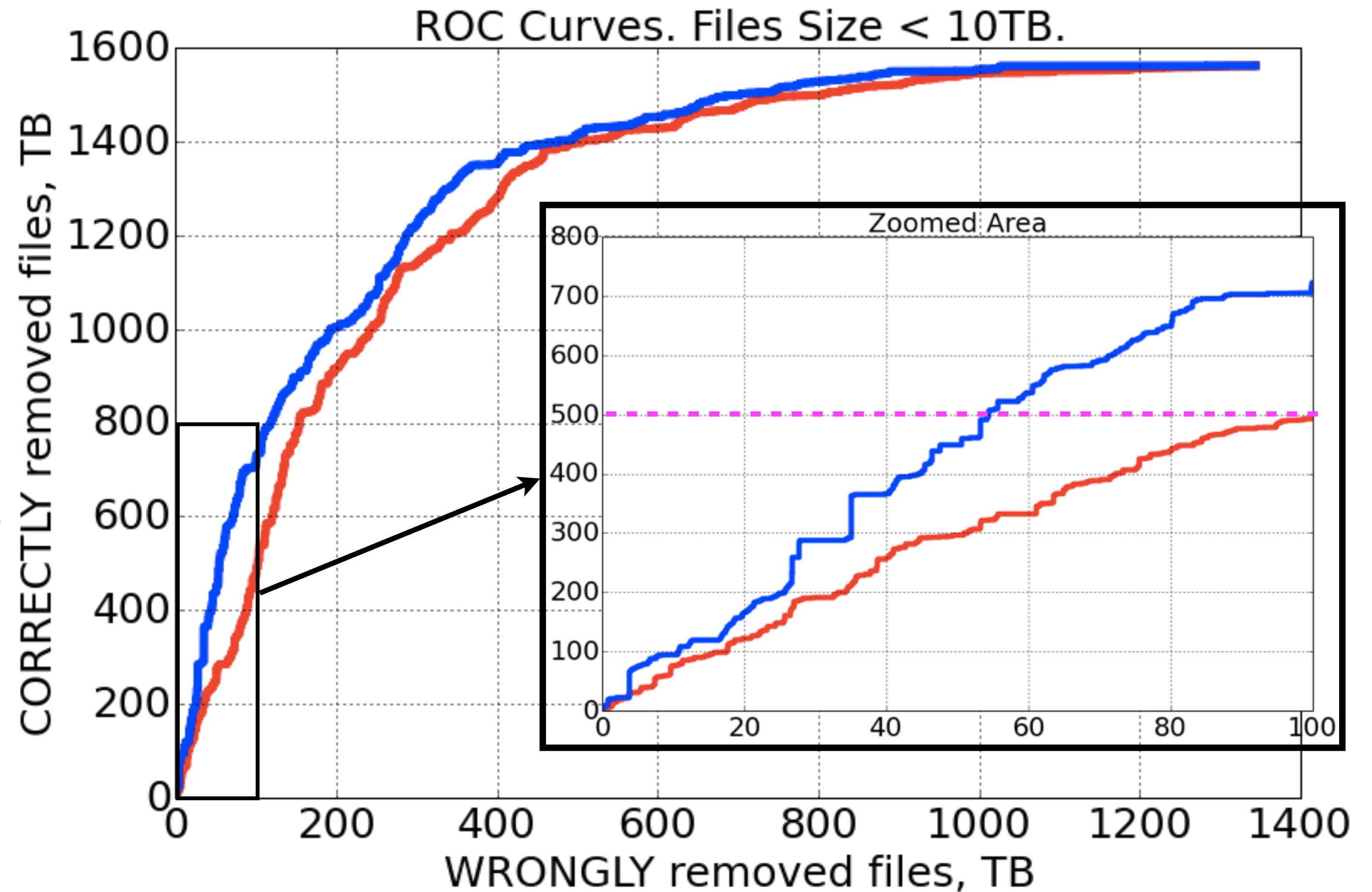


Use 2.5 years of records

Train classifier on (part of) historic data

Compare performance with Least Recently Used algorithm

Including more features allows to significantly decrease missed hit rate



Short-term forecast

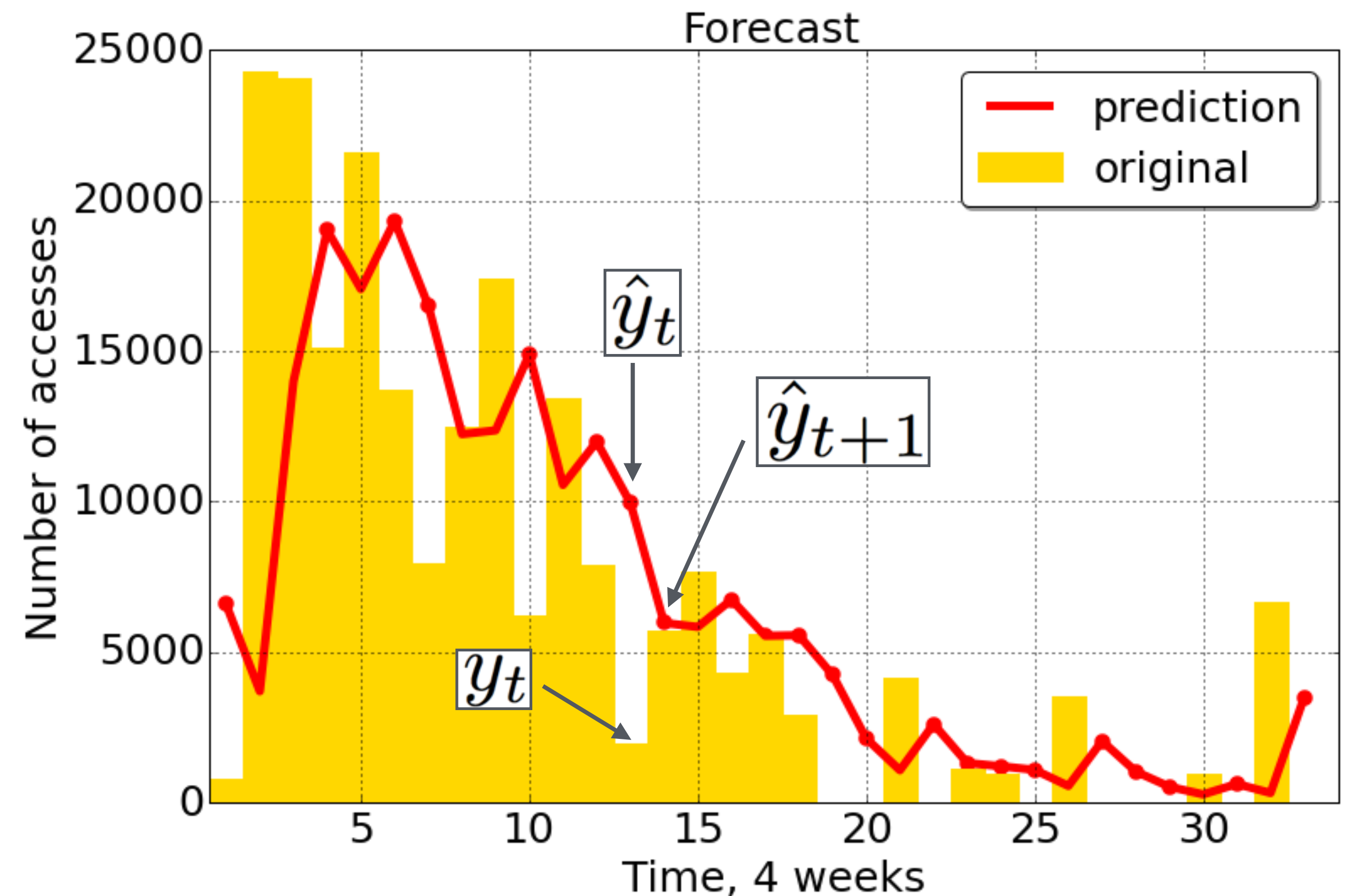
For a data file number of accesses prediction Brown's simple exponential smoothing model is used.

The model is defined as:

$$\hat{y}_{t+1} = \hat{y}_t + \alpha(y_t - \hat{y}_t)$$

$$\alpha = \operatorname{argmin}_t \sum (\hat{y}_t - y_t)^2$$

$$\hat{y}_0 = \frac{1}{n} \sum_{i=0}^n y_i, \quad \alpha \in (0, 1)$$



Replication strategy

The results of the long-term prediction and short-term forecast are used to calculate one of the following metrics:

$$M = \frac{\hat{y}_{t_{curr}+1}}{n_{replicas}} (\alpha + classifier_output)$$

short-term effect

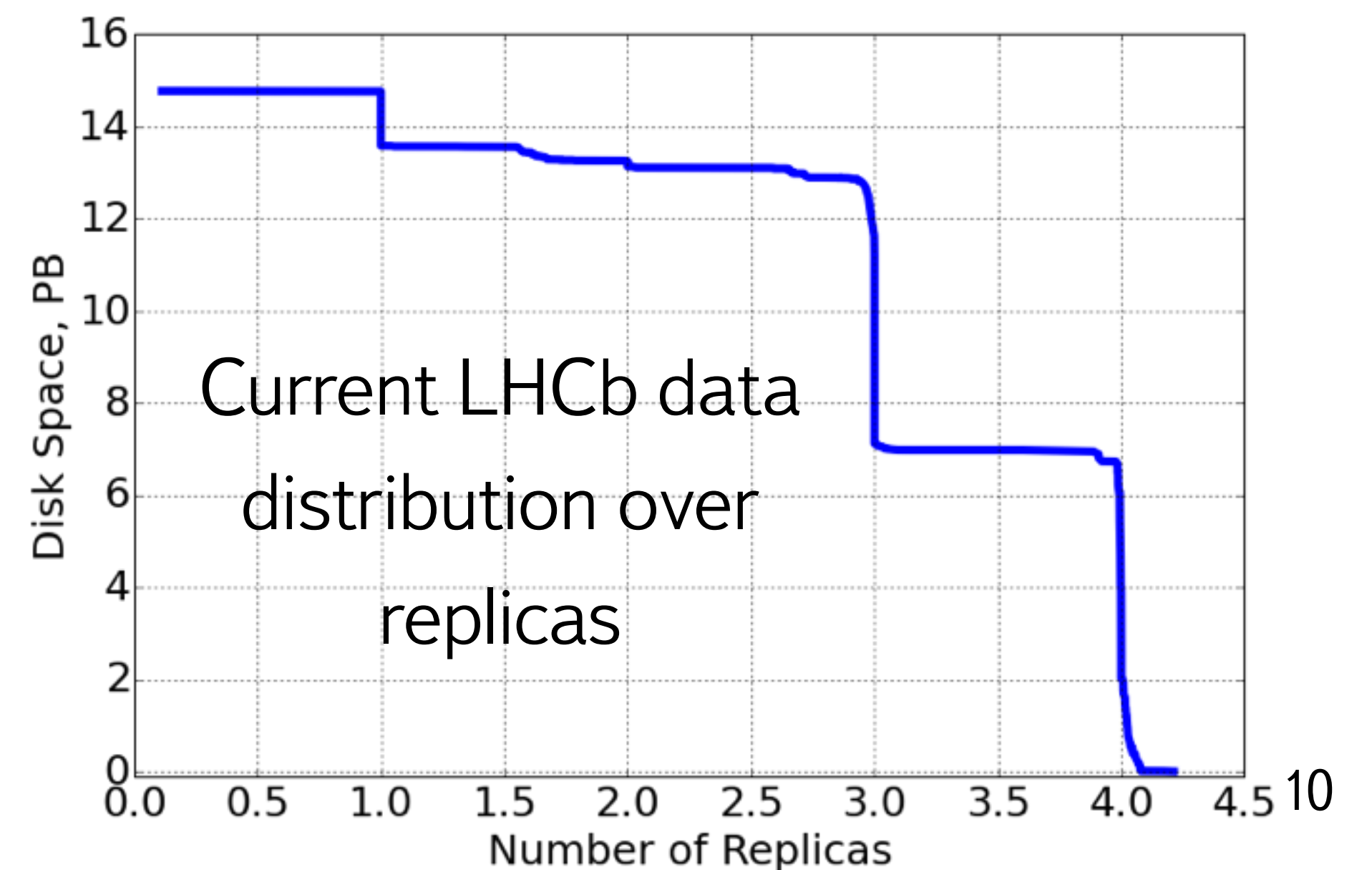
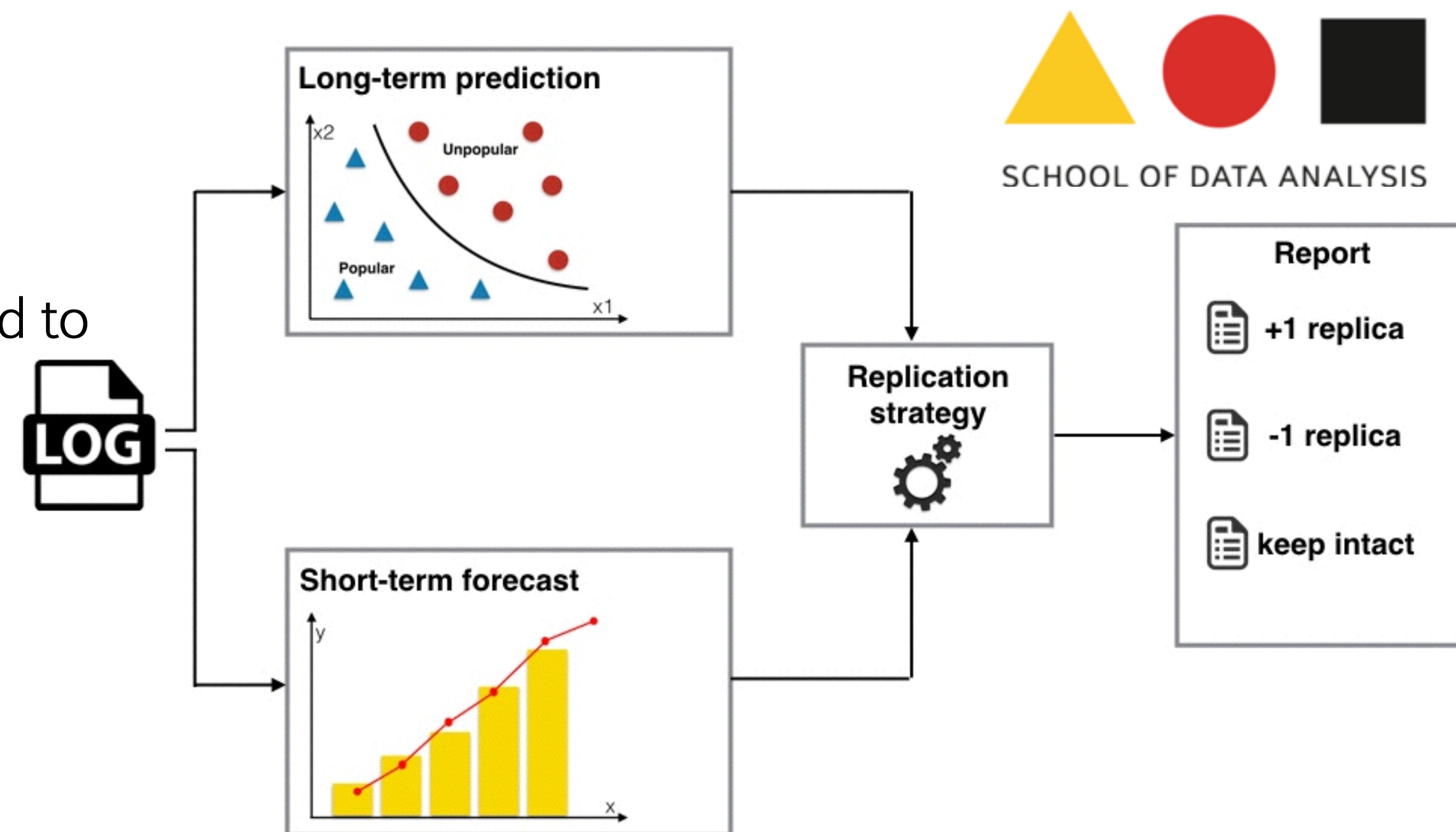
long-term effect

To save disk space:

- › remove replica of the dataset with minimal M
this will remove replicas for datasets which will be presumably less popular in the future.

To fill disk space:

- › add replica of the dataset with maximal M
this will add replicas for the datasets which will be presumably more popular in the future.



Replication Practice

When considering replicas to remove/add:

- which particular site should be used?
- how much missing hits cost?

Need to optimise data storage and data processing together

- need a simulation to estimate effect of variations of different data placement and data submission algorithms and metrics

Status Update

Grid Simulation



Workflow

The simulator includes the LHCb Tier0 at CERN and 7 Tier-1 with corresponding CPUs, disk and tape resources

Typical job and data varieties are simulated.

Use jobs and data statistics from the LHCb Dirac Web Portal as a reference.

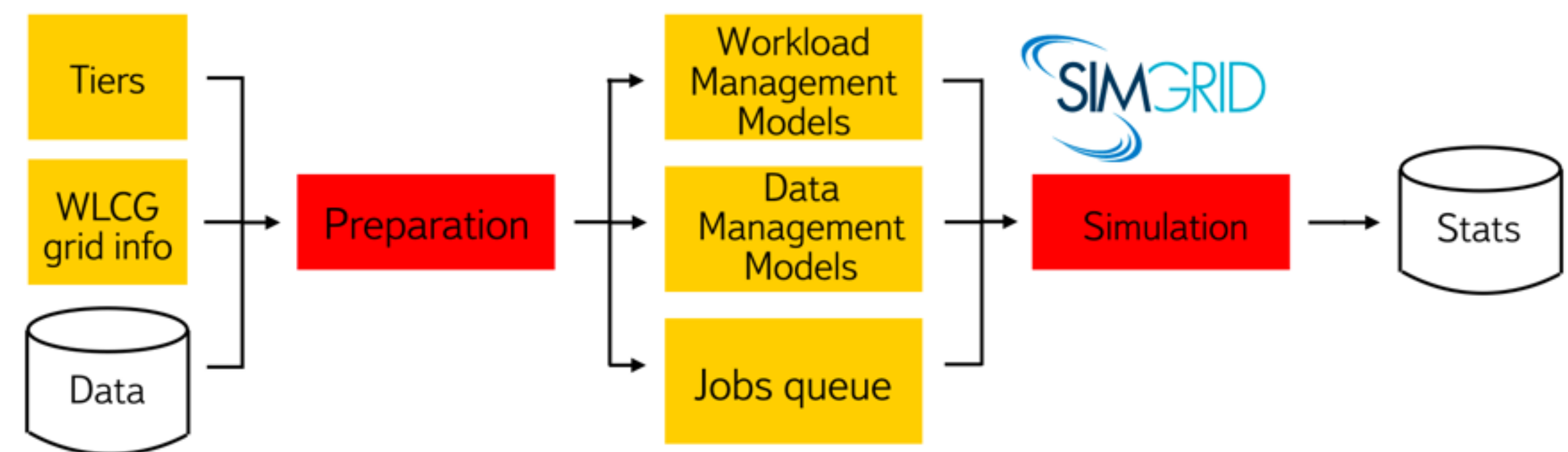
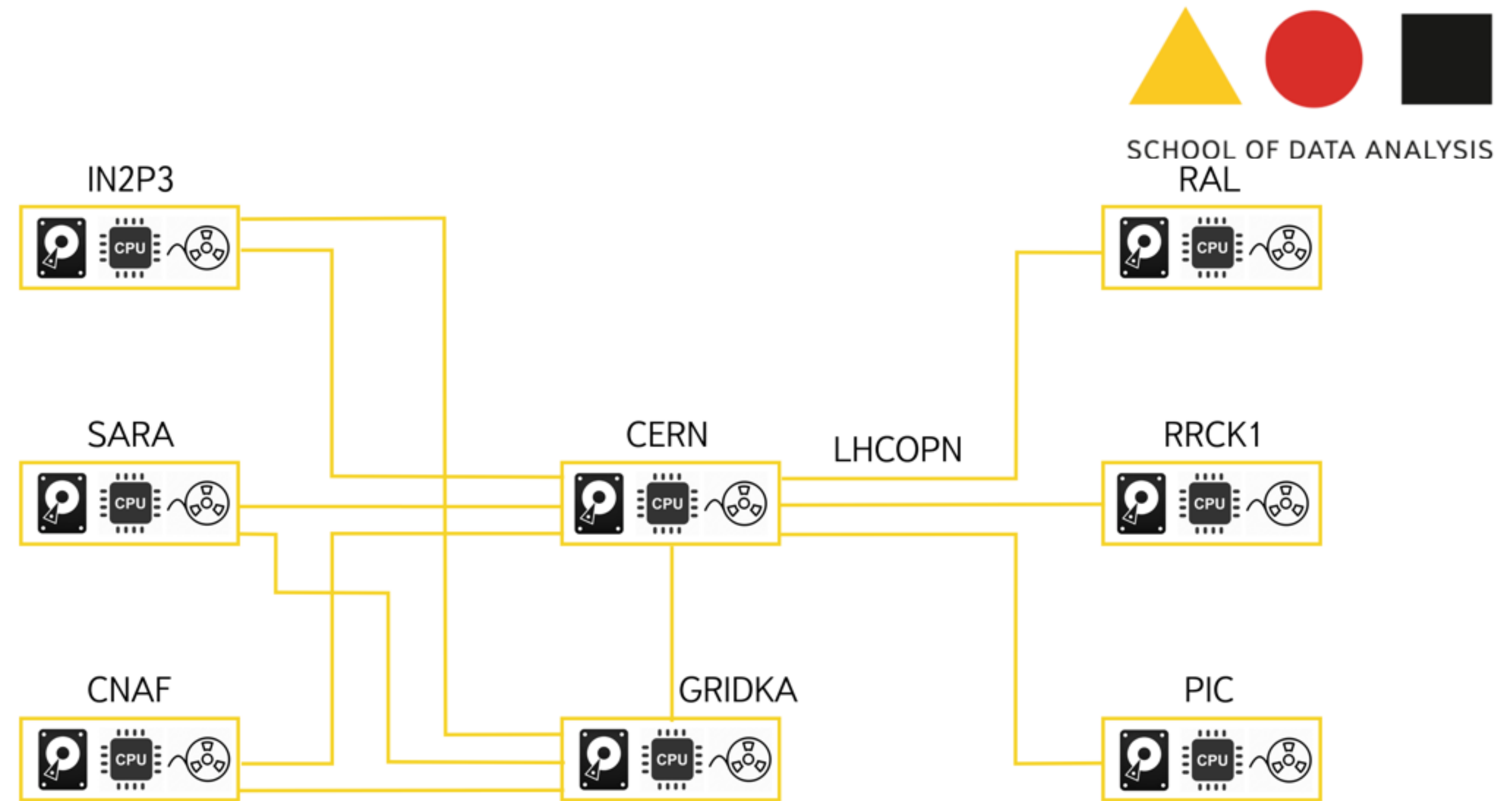
Two jobs scheduling models are considered:

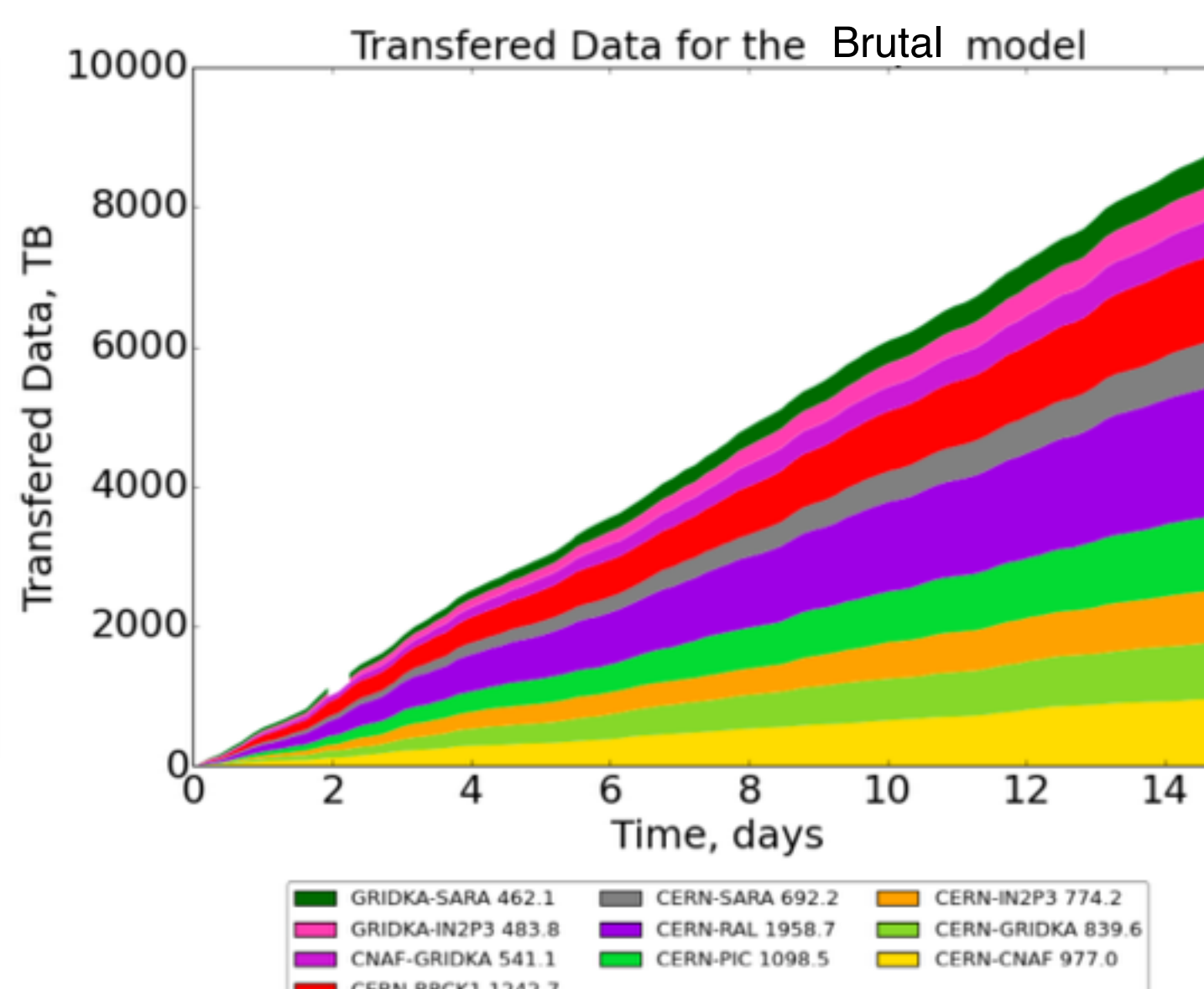
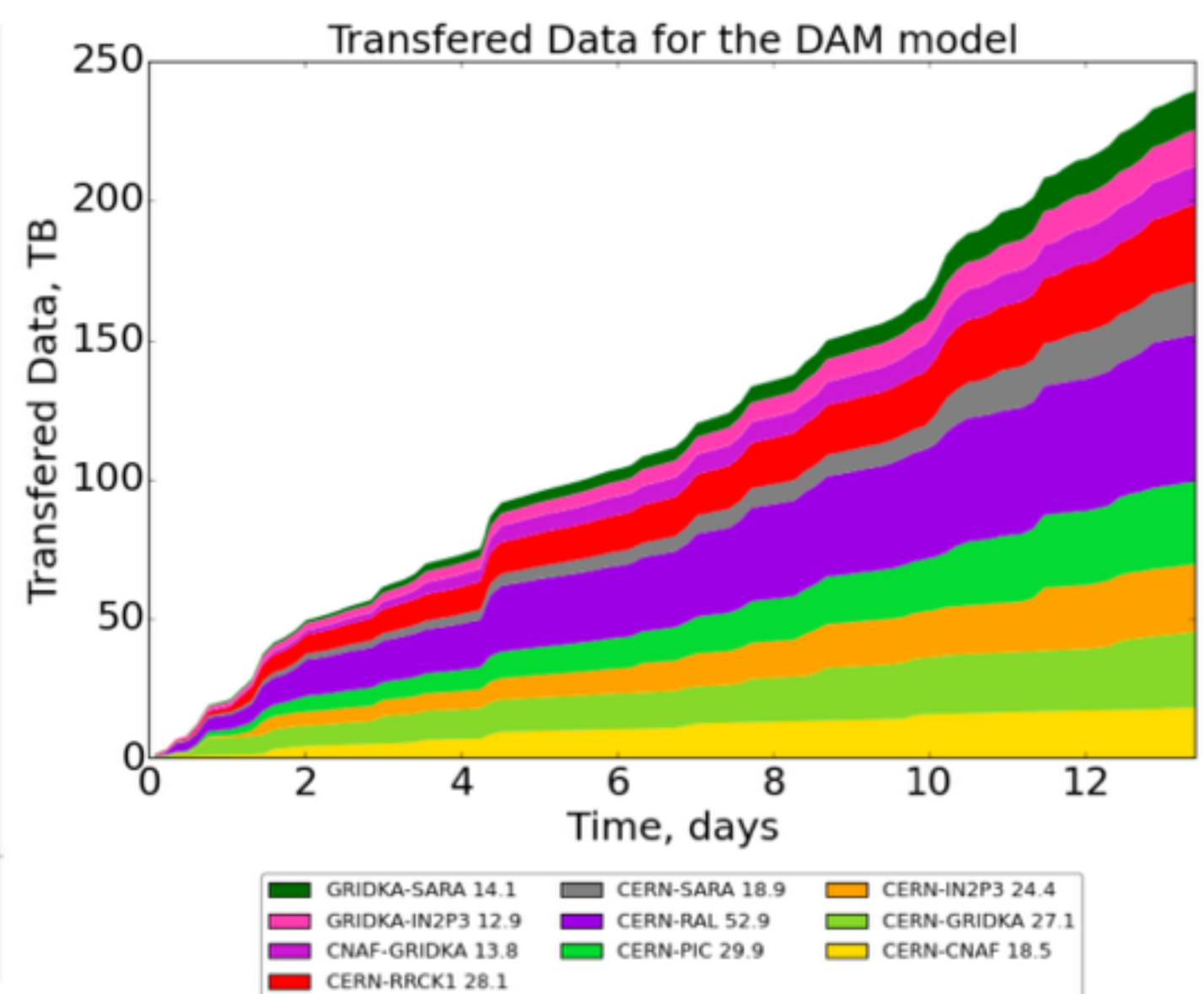
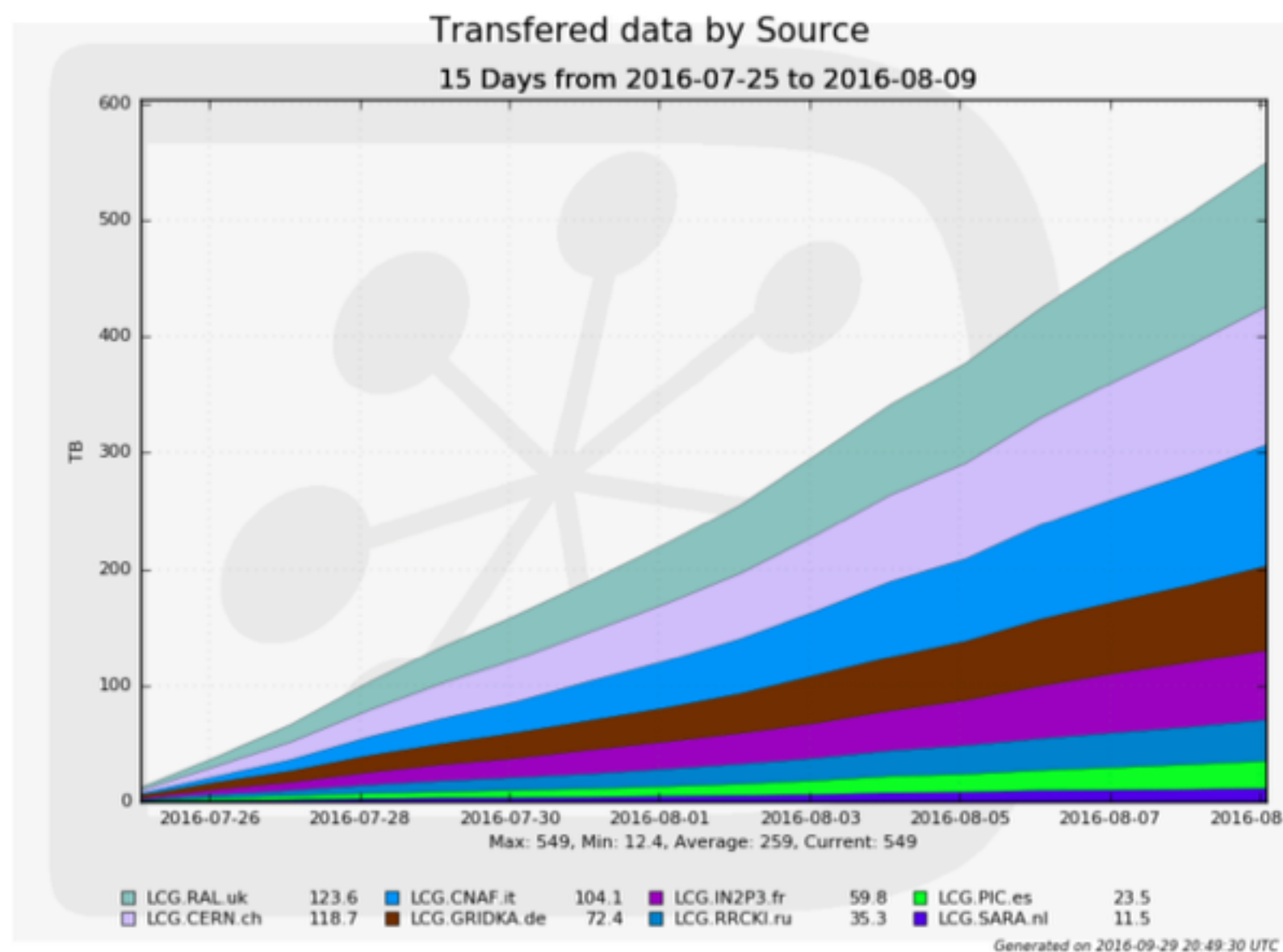
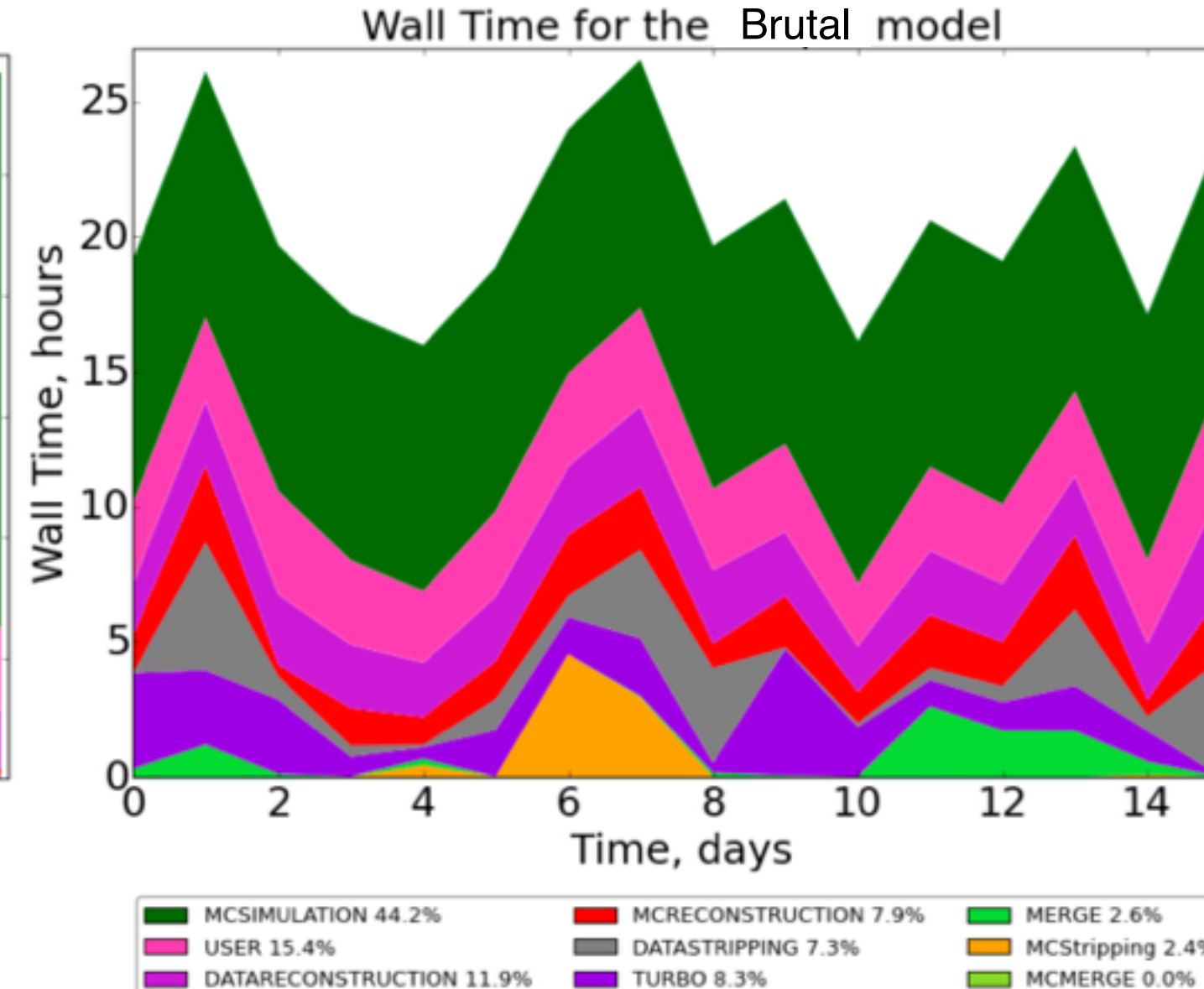
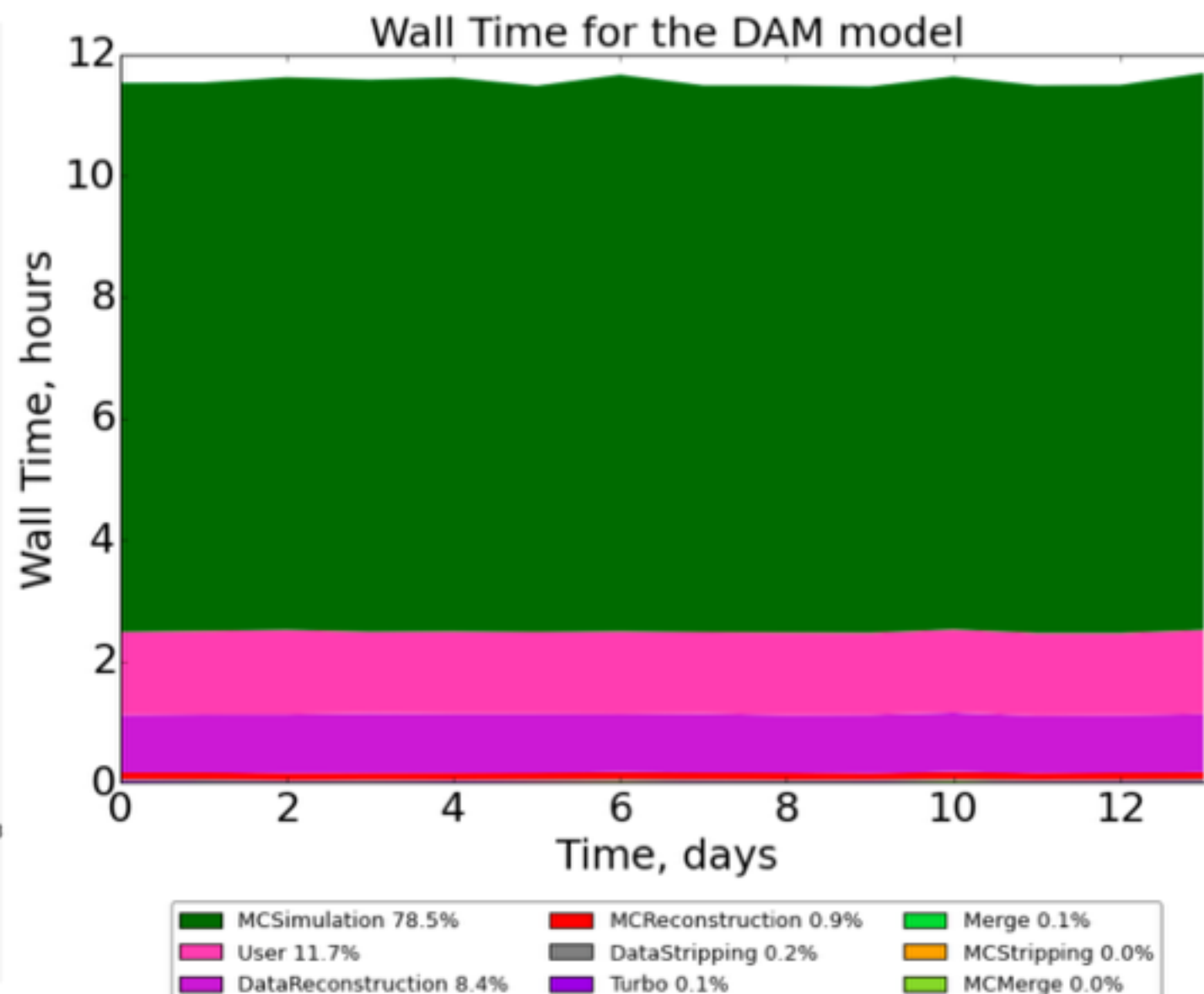
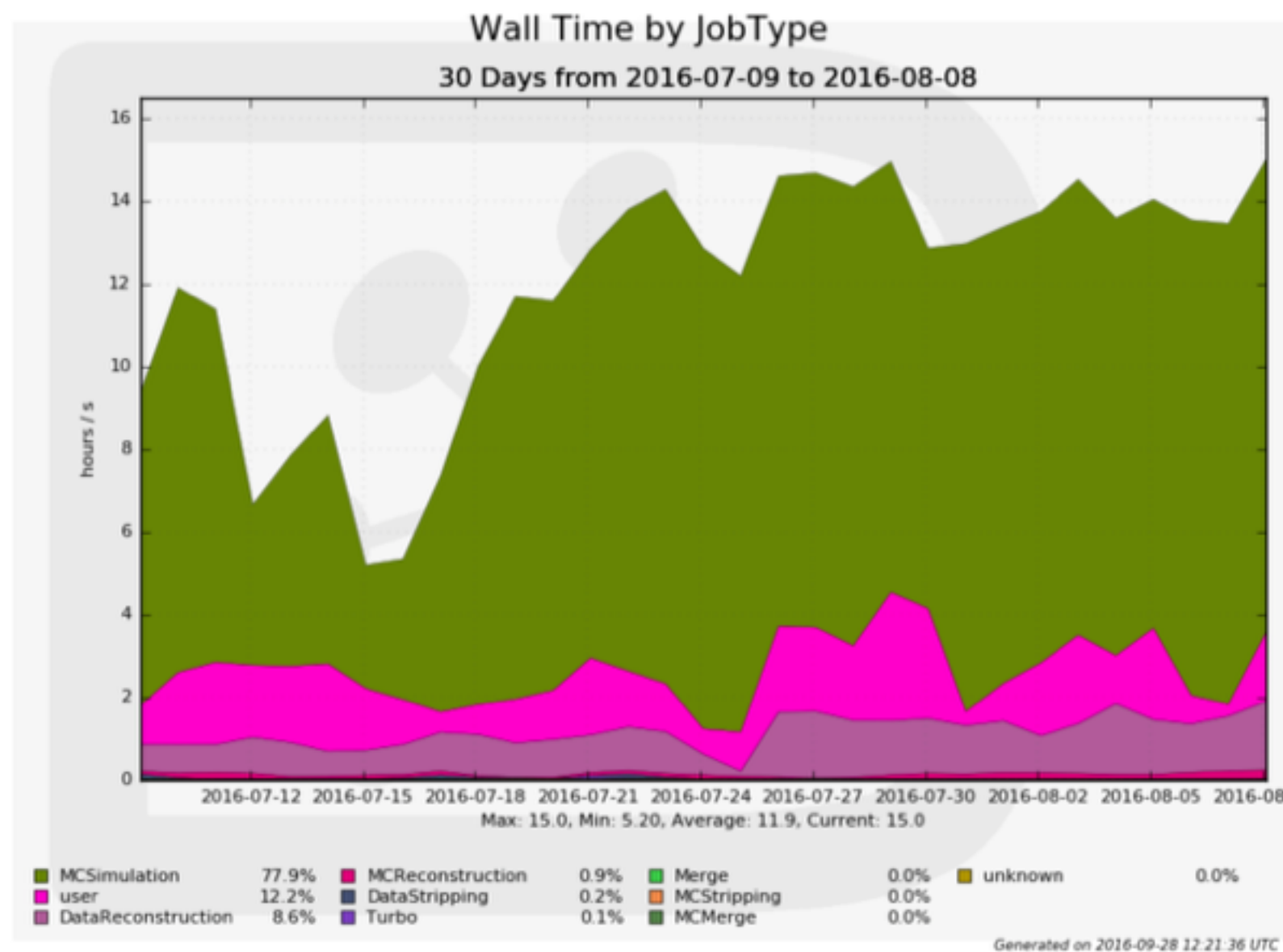
Brutal:

- › Job is sent to the node regardless the data availability

Data Availability Matching (DAM):

- › Job is sent to the node hosting required data





Intermediate Conclusion

1. Work in progress
2. Data replication optimisation using data popularity prediction allows to reduce miss rate for deletions by factor of 2 on actual data
3. Simulation of dedicated features of GRID system is necessary to estimate realistic effect of improvements
4. Simple simulation model reasonably reproduces historical GRID behaviours

Status Update

LHCb Data Quality Monitoring



LHCb's "presenter"

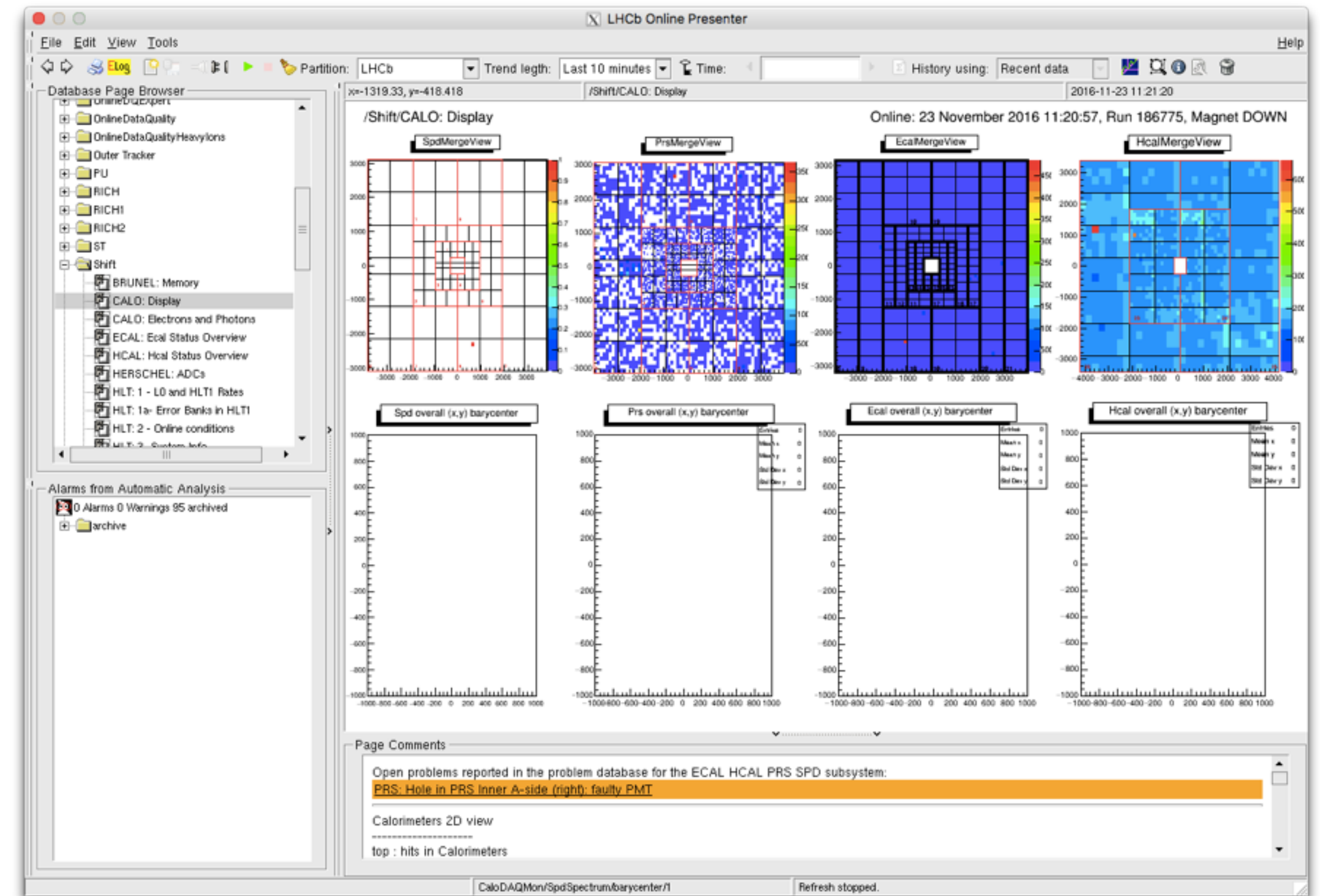
- Previous LHCb application for data quality monitoring (DQM)

- Implemented in C++

- Uses X-Window for GUI

- Problematic to support and extend

- Only usable from LHCb pit



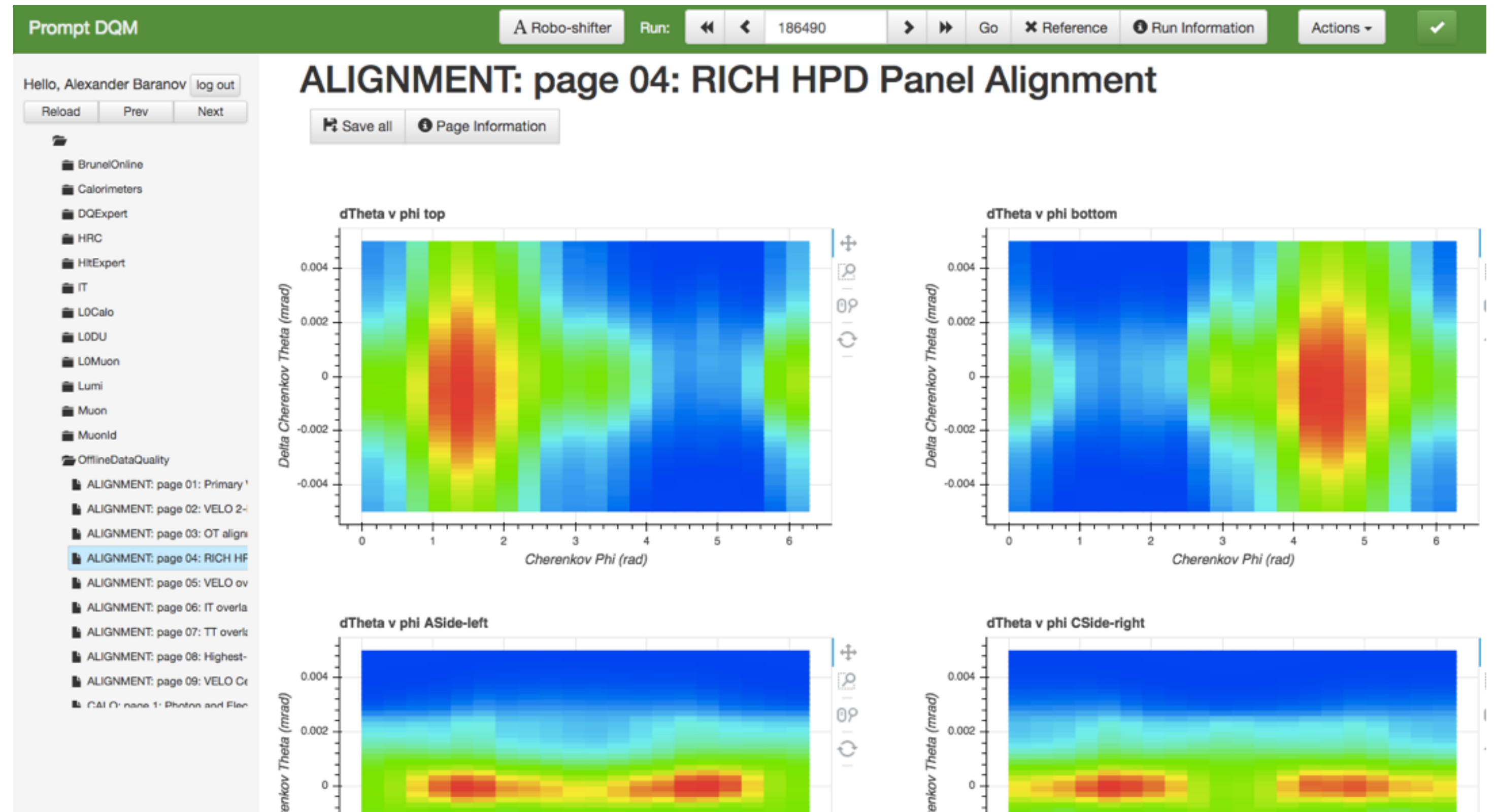
New Presenter Ecosystem: Monet

New app for LHCb DQM

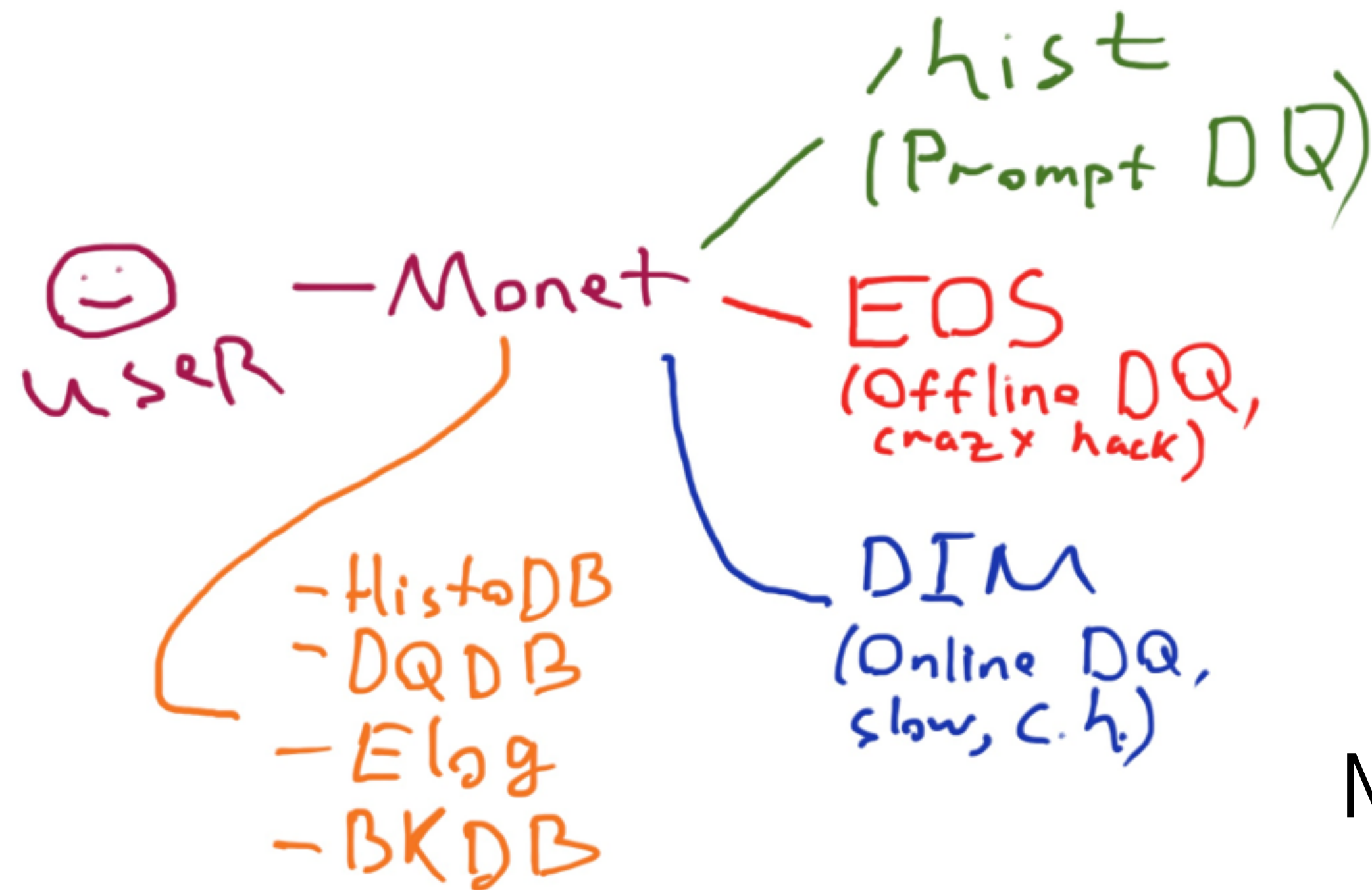
Web application

Implemented in Python

Already extended for Prompt, Offline and Simulation DQ



Connected data sources



Currently Monet is able to obtain the data from:

- > DIM (needed for Online Monitoring)
- > EOS stored root files
- > /hist/ or other mounted directories
- > histoDB

More data sources are foreseen:

- > MongoDB/ElasticSearchDB (data storage)
- > ZeroMQ (data transfer)

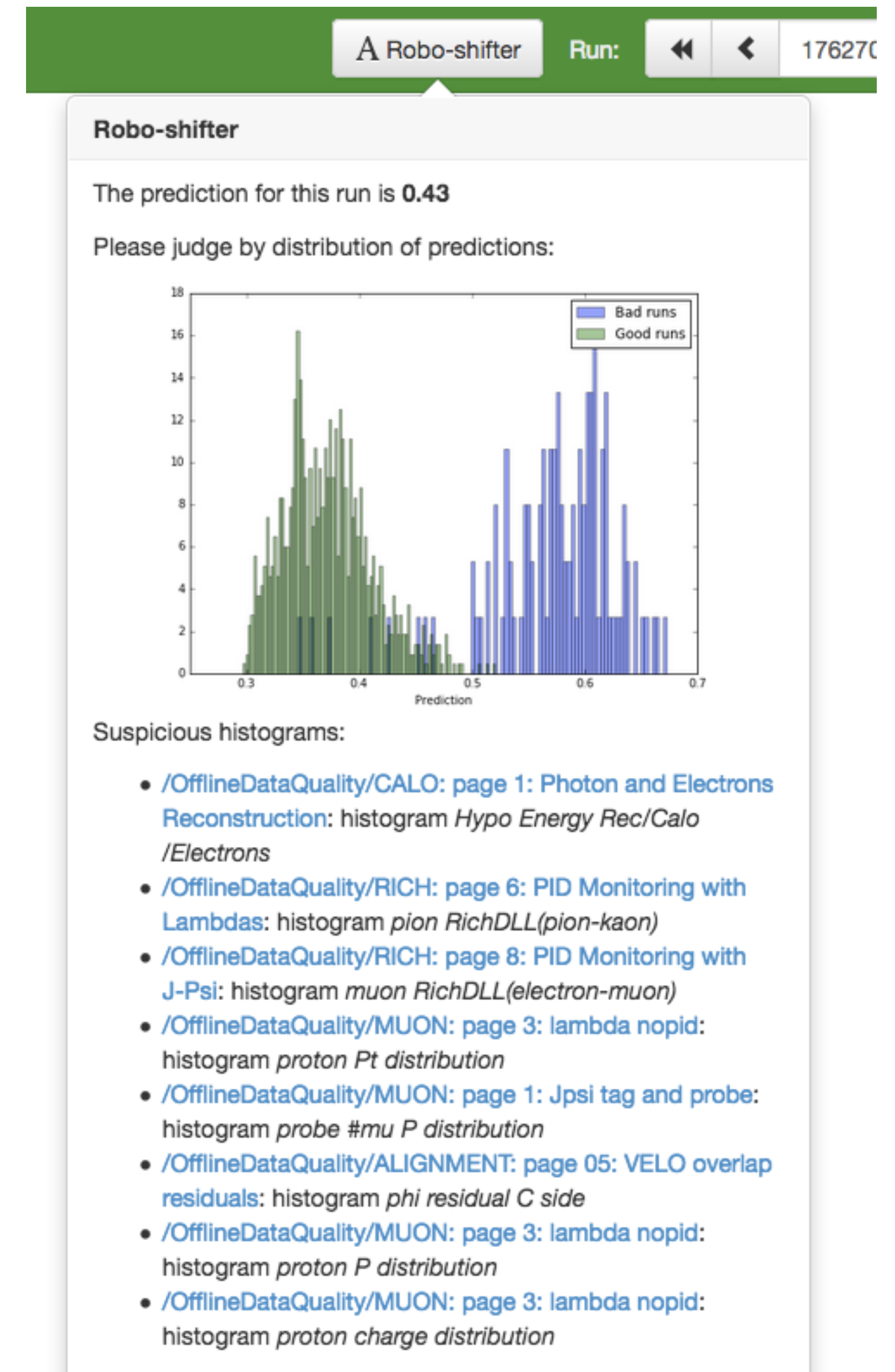
Status Update

Roboshifter



Roboshifter: Anomaly Detection

- ML-based assistant for DQ shifter
- Can predict probability of run being good or bad
- Based on high-level run data
- Provides potential problem sources extracted from BDT
- Integrated in Monet UI



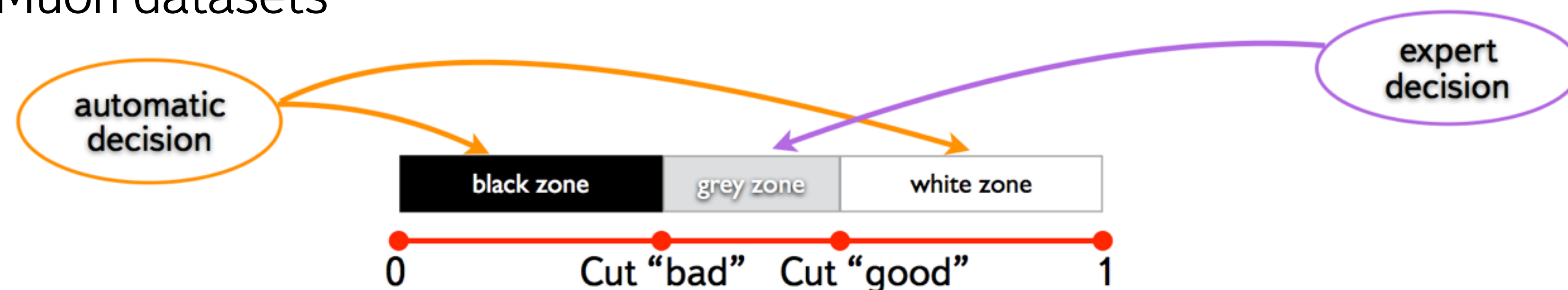
Status Update

Anomaly Detection in Detector Data



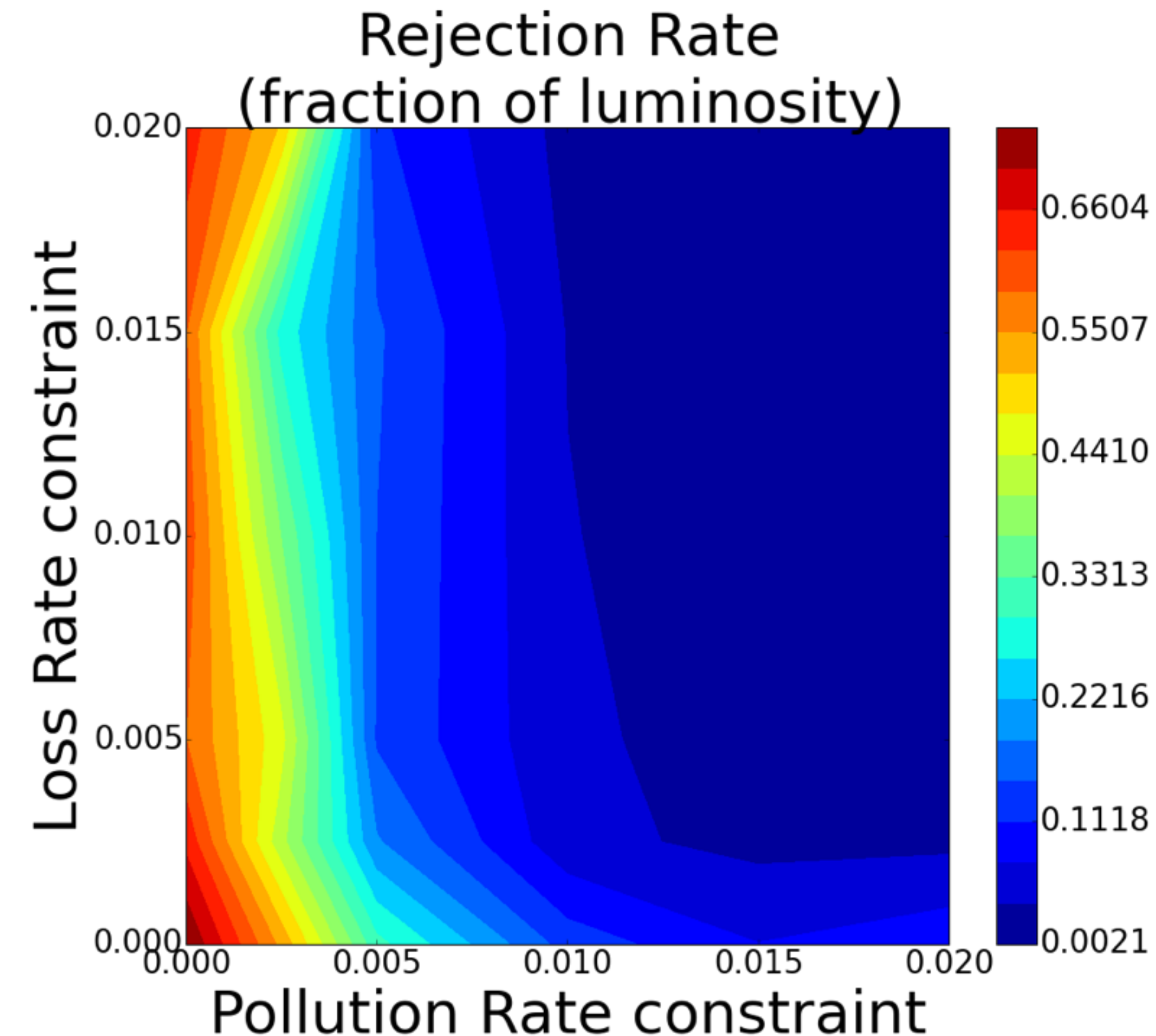
Supervised Learning

- › Use historical data processed by expert
- › ML algorithm learns the pattern that lead to the experts' decision
- › Combined effort with CMS
- › Until recently: use CMS 2010B run open data
- › Automated classification of LumiSections as “good” or “bad” using expert opinions on previous runs
- › Features: Particle Flow jets, Calorimeter Jets, Photons, Muons in Minbias, Photon, Muon datasets



Supervised Learning

- › The aim is to minimise the Manual work with:
 - low Loss Rate (“good” classified as “bad”)
 - and low Pollution Rate (“bad” classified as “good”)
- › ~80% saving on manual work is feasible for Pollution Rate at 5‰ and zero Loss Rate



From Proof of Concept to Routine Operation...

- | Promising result based on historical data

- | Close cooperation with CMS team is necessary to upgrade workflow to current data, incorporate into experiment Data Quality Monitoring framework

- | Launching it as an OpenLab project

- › with the ultimate goal to convert this proof of concept into routine operation

Title Text

Summary



Summary

- Several projects to deal with problems common for HEP experiments
- Emphasise use of Machine Learning in different applications
- Probe OpenLab as a venue for practical cooperation between LHC experiment and external team of experts