

# Advancing ATLAS DCS Data Analysis with a Modern Data Platform

CHEP2024

Luca Canali (CERN), Andrea Formica (Université Paris-Saclay),

**Michelle Ann Solis (University of Arizona)**

on behalf of the ATLAS Computing Activity



# Analysis of DCS Data

- Detector Control System (DCS) data mostly used for online operations
  - We only read small chunks of data for monitoring, or when we find detector problems
  - Data archived into Oracle database using Siemens WinCC OA application



## Can we get insights by analyzing large chunks of DCS data?

- Working with detector experts for the analyses: ex ATLAS New Small Wheel (NSW)



## Toolset

- Let's not use the Oracle DB for this (avoid overloading of ATLAS online DB, cost - Oracle licensed per CPU to scale up, performance)
- Solution: import data into a separate and optimized platform for data analysis



DCS data in  
Oracle



# Offloading Data from Oracle for Analytics

- Moving data from databases to analytics
  - Export from DBs into files (Parquet format)
  - Query data with scalable engines (Apache Spark)
  - Used by ATLAS DCS and ATLAS NSW and other projects



# DCS Data Offloaded to Parquet

- Data is copied from **Oracle** to **Parquet** files
  - Files hosted on Hadoop in general-purpose cluster called Analytix at CERN
- DCS Data
  - Time series data where each row contains a timestamp and an element ID associated with specific measured data values
  - Apply daily partitioning to parquet files by timestamp values for efficient querying and data import
  - Data imported incrementally daily due to large size and frequent updates

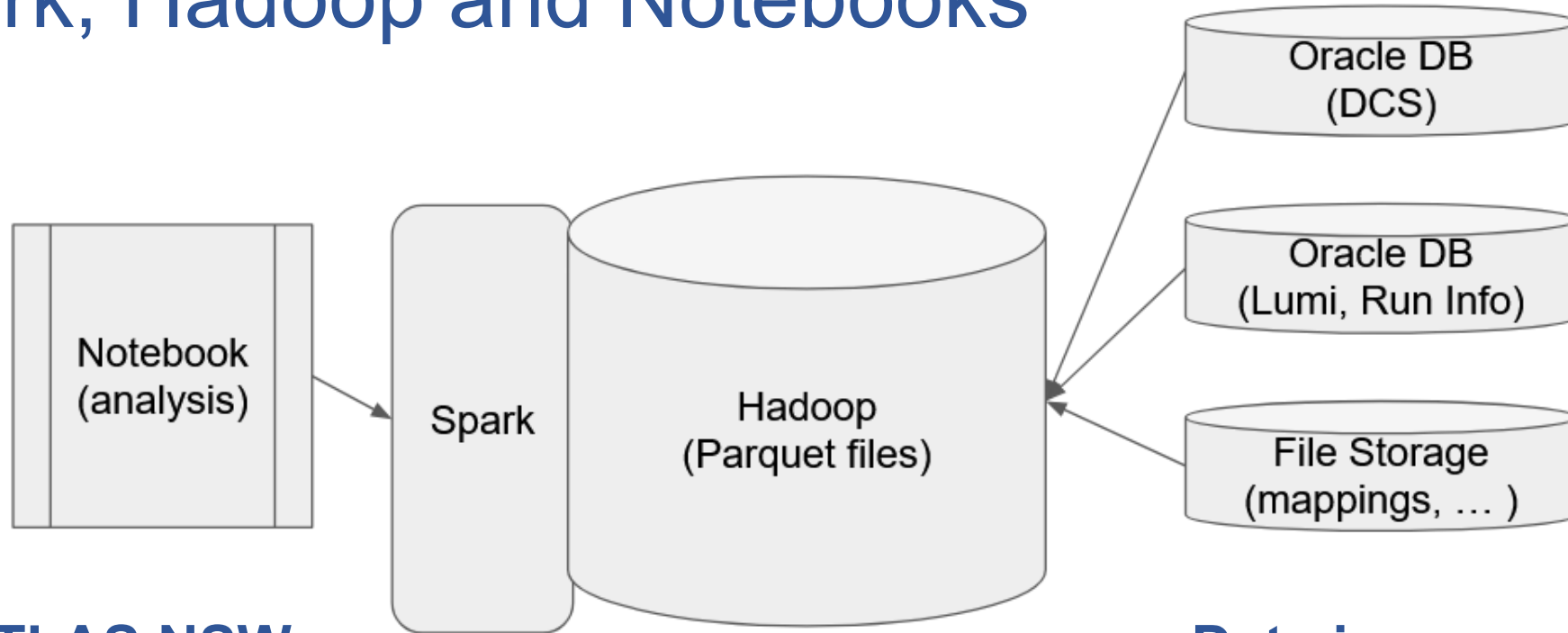
# Technology: why Apache Parquet?

- With Apache Spark plus Parquet files we are building a DB for large scale analysis
  - High **adoption** in industry and open source for building data lakes and data warehouses
- Parquet
  - Is a **columnar** data format
  - Optimized for storing and querying data for large-scale analysis
  - Uses encoding and compression
  - Data is stored together with its schema
  - Works well with Apache Spark, Pandas, and many other tools
  - Provides a simple way to map DB tables into files

# IT Services: SWAN, Spark, and Hadoop

- Data analysis platform, interactive and at scale
  - Running on many CPUs like batch, but interactive (notebooks)
- Storage and compute (**Hadoop, Spark**) from CERN IT
  - Shared cluster capacity: 1500 (physical) cores, 20 PB
- How much do we use?
  - DCS data: used so far (Oct 2024) ~ 3 TB
  - Compute used: Sporadic CPU-intensive queries
- **SWAN**, a CERN service for Python notebooks
  - Integrates LCG releases and Spark
  - Easy to get started, build from examples

# End-to-End Data Analysis Platform with Spark, Hadoop and Notebooks



## Use case: ATLAS NSW

Done from user perspective via simple notebooks where the appropriate libraries are pre-installed

## Clusters for execution:

“See” via spark as relational tables, so able to easily make join-like requests on imported files

## Data is prepared:

This system is used as a dynamic database, importing needed information from external data sources into parquet files



# ATLAS New Small Wheel (NSW)

MicroMegas (**MMG**) +  
small Thin Gap Chambers  
(**sTGC**)

- Muon end-cap spectrometer
- Phase I upgrade





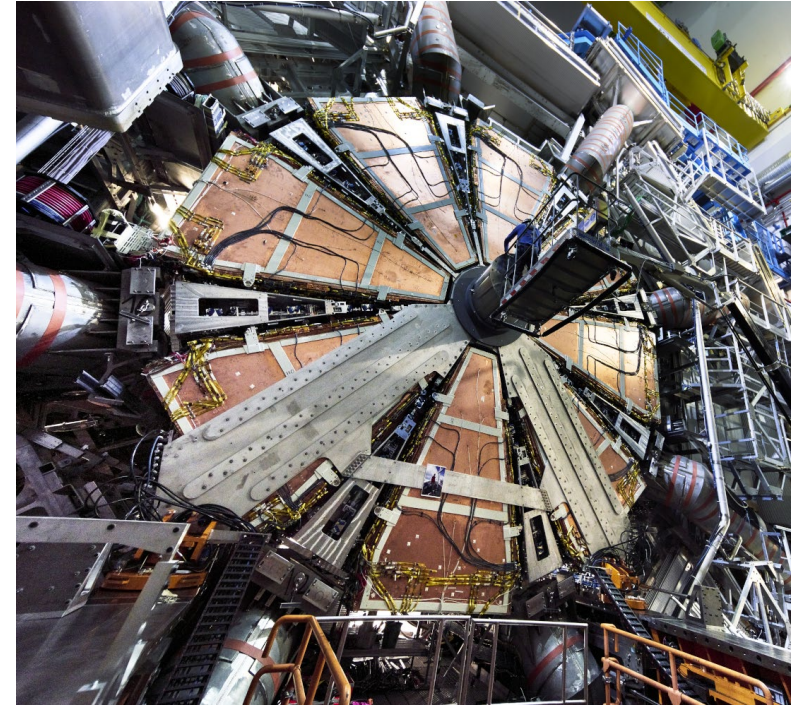
# ATLAS NSW Introduction

## DCS Monitoring

- HV & LV
- Gas
- Electronics
  - Temperatures
  - Voltages/Currents
  - online status and configuration

## Lots of data to analyze!

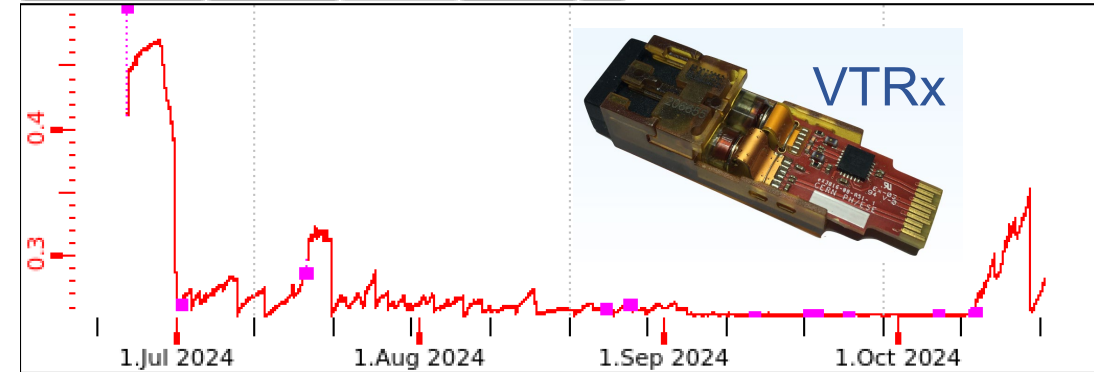
- **23M** rows per day, **6.6 billion** rows in 2024 for MMG so far
- Too large to query Oracle directly



- 64 sectors, 8 layers per sector
- 5k front-end electronics boards
- over 2M readout channels

# NSW DCS Monitoring

## Motivations for querying DCS

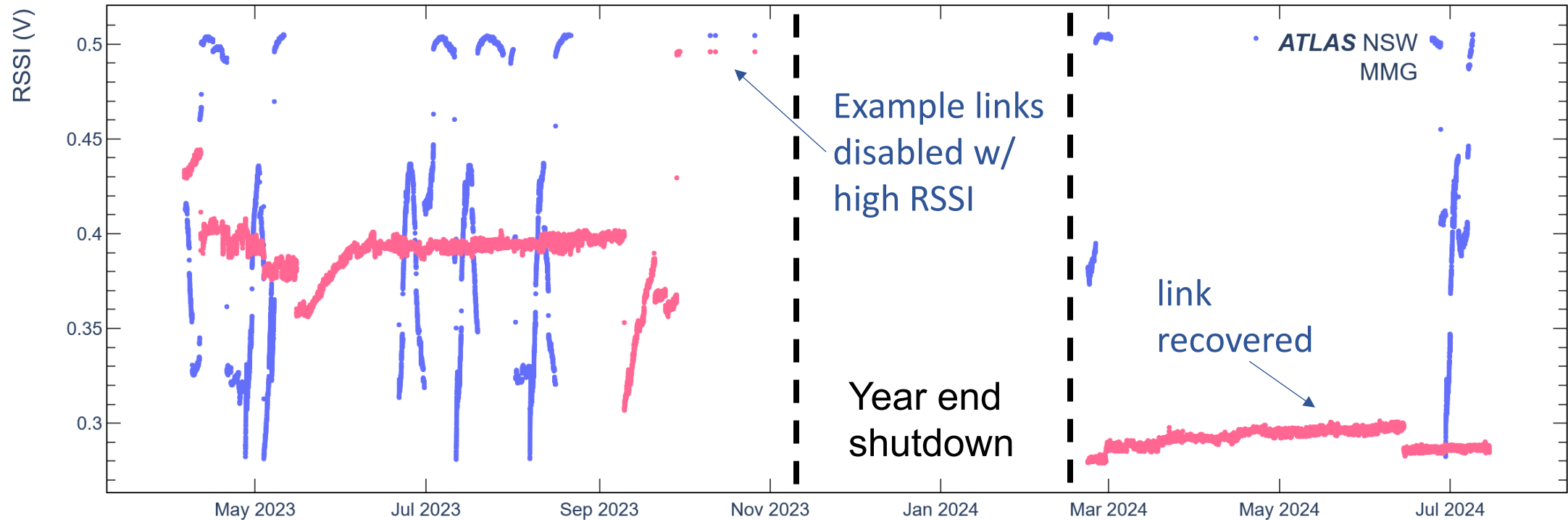


High/unstable Received Signal Strength Indicator (**RSSI**) sign of failing links

- **Monitoring of known hardware issues**
  - **VTRx** Versatile Transceiver optical link failures affecting ~5%
  - > 1400 VTRx on NSW
  - How many links affected? Recovered? Stable over time?
  - Determine any possible intervention needed during long shutdown (LS3)
- **DAQ link stability investigation**
  - ~10% of optical fibers showed issues in DAQ
  - Looked at several parameters monitored in the DCS to find correlations, ruling out hardware issues
- **Monitoring of HV status for performance efficiency**
  - Disabled, resistive, below nominal channels

# NSW VTRx Analysis

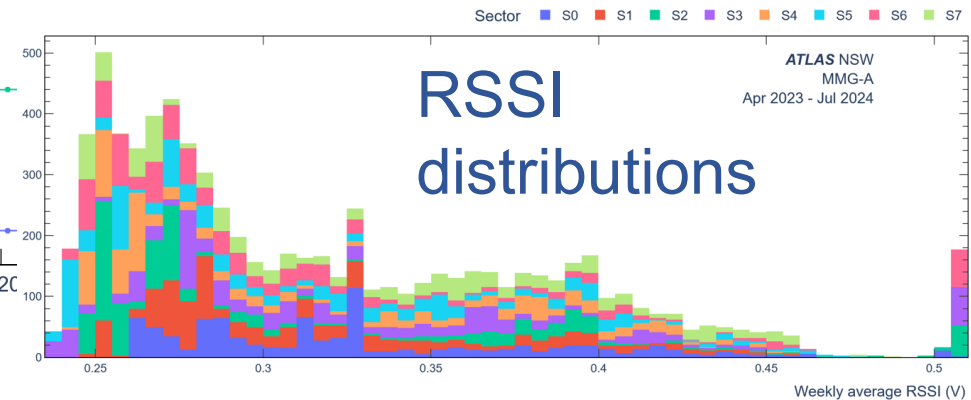
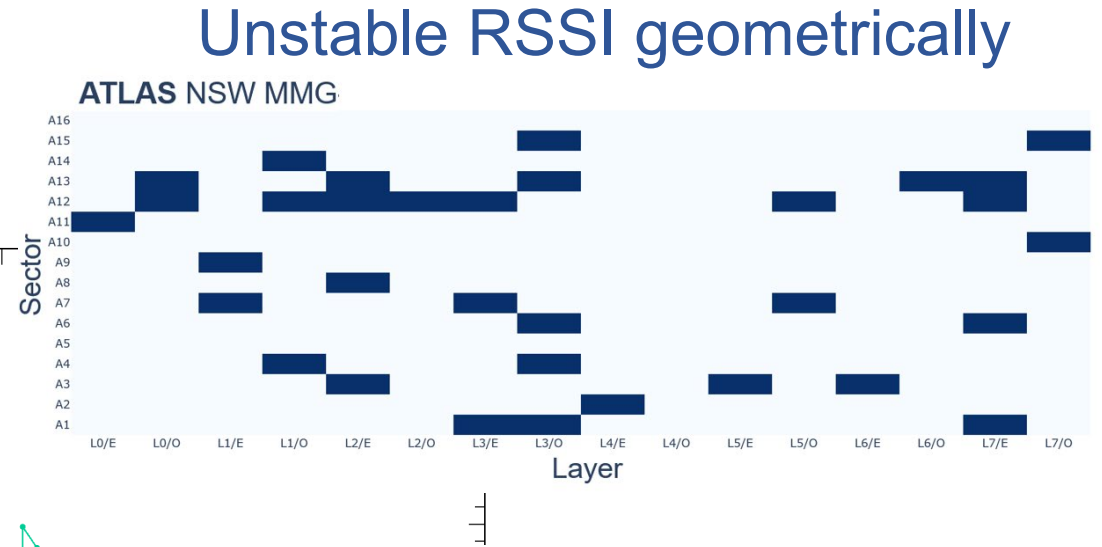
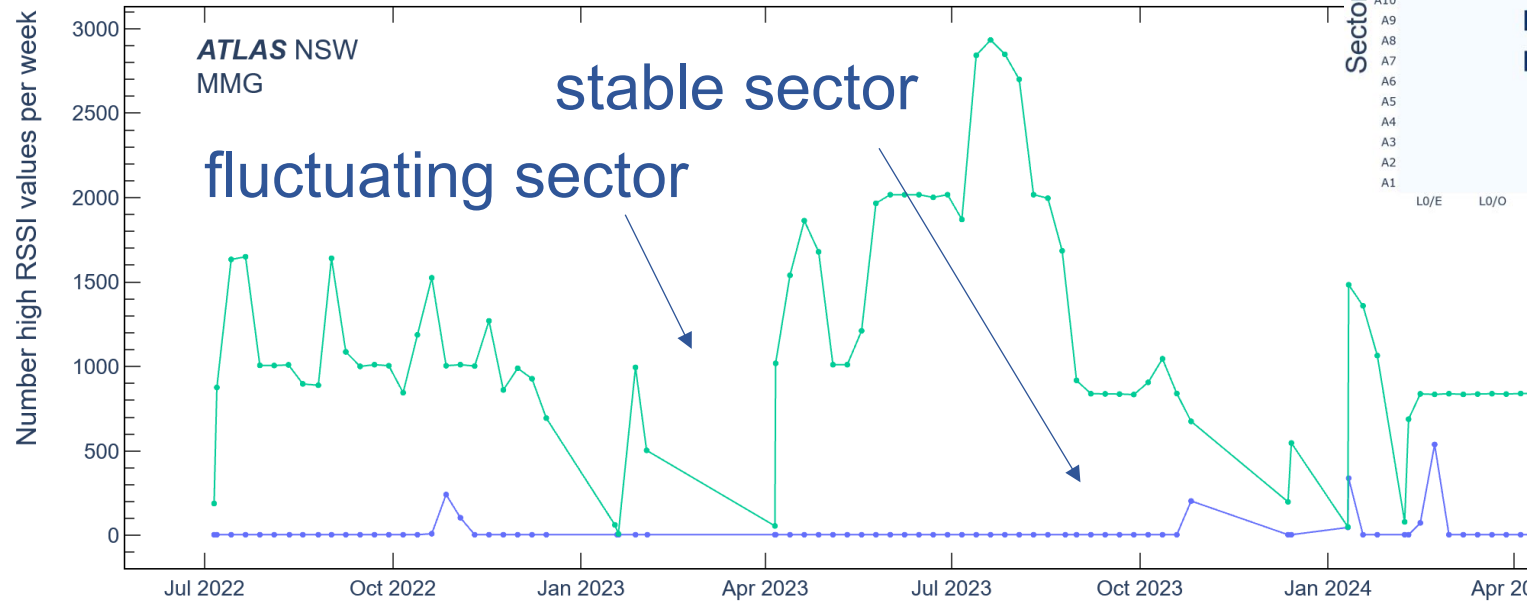
Quickly pick out problematic VTRx over time



- Queried parquet files for **all MMG RSSI data April 2023 to July 2024** into DataFrame → < 30 sec
- **Aggregated** time in 10 sec intervals for all links → 2 min
- **Selected** on geometry (side, sector, layer, etc) and identified all bad VTRx → < 1 min
- Quickly convert Spark DataFrames to Python's Pandas DataFrames to **plot**

# NSW VTRx Analysis

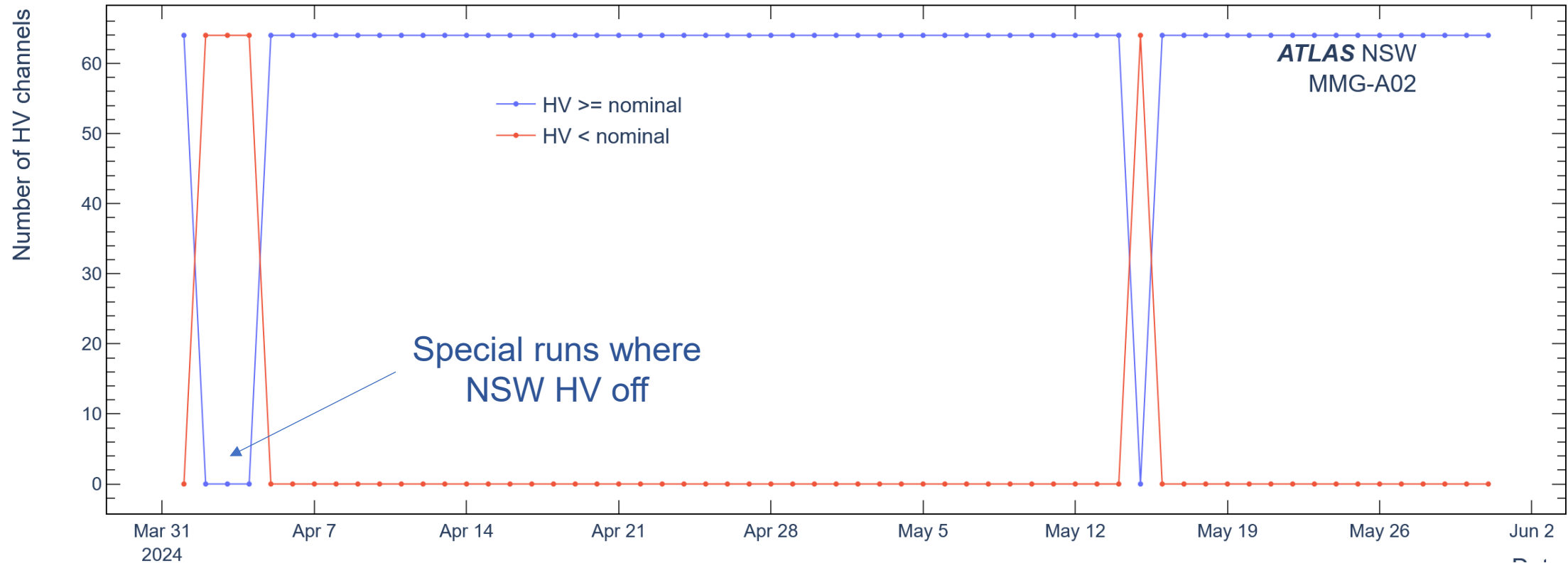
Notebooks and Spark DataFrame API provide easy, powerful way to analyze lots of data



- Over 500 optical links on MMG
- Perform **large scale data operations** in a few lines of code: extract data fields, filter, compute averages, derivatives, etc.

# NSW High Voltage Status

One sector during ATLAS combined runs



- Over 2000 MMG HV channels
- **Join** different sources with **DCS** data in the same platform, for ex:  
**T0 processed runs** to select times just during physics runs



# Wrap Up

- **DCS** data from Oracle to a platform for **analysis**
  - **Pipelines** are up and running, maintained by ATLAS-DBA using CERN IT infrastructure
- Analysis on CERN Jupyter Notebook service (**SWAN**) using Python
  - Expressive APIs, executed at scale on clusters
- **Example**
  - Monitoring of known hardware issues and high voltage
  - Understand readout and stability issues
  - Showed how NSW profited from this platform to analyze DCS monitoring data to better understand detector issues in an easy and accessible way

# References

- Gitlab project:
  - <https://gitlab.cern.ch/atlas-dba/dcs-offload>
- SWAN service:
  - <https://swan.web.cern.ch/swan/>
- Hadoop service:
  - <https://hadoop-user-guide.web.cern.ch/>
- Training material on Spark and SWAN
  - <https://sparktraining.web.cern.ch/>
- Acknowledgements:
  - ATLAS DBA and ADAM (ATLAS Database and Metadata) group, NSW experts, Service support for Oracle, SWAN and Spark, Hadoop