## Data processing of a "typical" FAIR experiment (based on some personal experiences)

Johan Messchendorp, April 17, 2019







#### **Origin of Mass & Color Confinement**

- •Higgs addresses only a tiny (1%) fraction of the mass of all visible matter.
- •The rest is due to the *strong color force.*





## Physics Data crunching







## Basic analysis ingredients



- Clustering/pattern recognition: identify tracks from sensor data
- Kalman fitter/vertex finder: reconstruct momenta and decay vertex of each track
- Multi-variate analysis/Machine Learning: determine particle-type probabilities of each track, background reduction
- Kinematic constrain fitting: determine event topologies, improve resolutions
- Multi-parameter fitting (Partial Wave Analysis): determine properties of new discovered particle
- Monte Carlo simulations: study efficiencies, analysis techniques, systematic errors, etc.

## Basic analysis ingredients



- Clustering/pattern recognition: identify tracks from sensor data
- Kalman fitter/vertex finder: reconstruct momenta and decay vertex of each track
- Multi-variate analysis/Machine Learning: determine particle-type probabilities of each track, background reduction
- Kinematic constrain fitting: determine event topologies, improve resolutions
- Multi-parameter fitting (Partial Wave Analysis): determine properties of new discovered particle
- Monte Carlo simulations: study efficiencies, analysis techniques, systematic errors, etc.



 Computational complexity: machine learning, fitting, iterative procedures, Monte Carlo studies, etc..

- Computational complexity: machine learning, fitting, iterative procedures, Monte Carlo studies, etc..
- Data rates "high": up to 10 MHz.

- Computational complexity: machine learning, fitting, iterative procedures, Monte Carlo studies, etc..
- Data rates "high": up to 10 MHz.
- Data processing "slow": "second/event.

- Computational complexity: machine learning, fitting, iterative procedures, Monte Carlo studies, etc..
- Data rates "high": up to 10 MHz.
- Data processing "slow": "second/event.
- Data volume "large": few PBytes/year after filtering.

- Computational complexity: machine learning, fitting, iterative procedures, Monte Carlo studies, etc..
- Data rates "high": up to 10 MHz.
- Data processing "slow": ~second/event.
- Data volume "large": few PBytes/year after filtering.
- User community "large", O(103) and divers.

- Computational complexity: machine learning, fitting, iterative procedures, Monte Carlo studies, etc..
- Data rates "high": up to 10 MHz.
- Data processing "slow": "second/event.
- Data volume "large": few PBytes/year after filtering.
- User community "large", O(103) and divers.
- "Software" development has a "low reputation".

 <u>Centralized</u> "online" processing of detector data with dedicated hardware/software: event selection -> "raw data" (RAW).

- <u>Centralized</u> "online" processing of detector data with dedicated hardware/software: event selection -> "raw data" (RAW).
- <u>Centralized</u> calibration and reconstruction of raw data to "analysis objects" (ESD/AOD).

- <u>Centralized</u> "online" processing of detector data with dedicated hardware/software: event selection -> "raw data" (RAW).
- <u>Centralized</u> calibration and reconstruction of raw data to "analysis objects" (ESD/AOD).
- <u>Centralized</u> production and pre-processing of MC data.

- <u>Centralized</u> "online" processing of detector data with dedicated hardware/software: event selection -> "raw data" (RAW).
- <u>Centralized</u> calibration and reconstruction of raw data to "analysis objects" (ESD/AOD).
- <u>Centralized</u> production and pre-processing of MC data.
- Usage of <u>common software framework</u> (version control) and training of users to handle AOD to study their channel of interest.

- <u>Centralized</u> "online" processing of detector data with dedicated hardware/software: event selection -> "raw data" (RAW).
- <u>Centralized</u> calibration and reconstruction of raw data to "analysis objects" (ESD/AOD).
- <u>Centralized</u> production and pre-processing of MC data.
- Usage of <u>common software framework</u> (version control) and training of users to handle AOD to study their channel of interest.
- <u>Decentralized</u> and distributed data of users. Validation via internal reviewing process based on analysis memos, leading to publication.

#### How to deal with the challenges? Physics results (public) Analysis review Preprocessed Raw data Analysis Data Researchers data data data processing processing ESD AOD

Experiment and Monte Carlo

"Waterfall model"





Monte Carlo

"Waterfall model"







Limitation of computing model (and where ESAP could help) • Analysis methods and results by individual researchers not optimally shared. Reproducibility? Limitation of computing model (and where ESAP could help) • Analysis methods and results by individual researchers not optimally shared. Reproducibility?

Lack of a common user-data management.

- Analysis methods and results by individual researchers not optimally shared. Reproducibility?
- Lack of a common user-data management.
- Analysis workflow concentrated to one single experiment, no real export/import mechanism from other experiments foreseen.

- Analysis methods and results by individual researchers not optimally shared. Reproducibility?
- Lack of a common user-data management.
- Analysis workflow concentrated to one single experiment, no real export/import mechanism from other experiments foreseen.

Missing the concept of a modern information system.

"Target"

'Waterfall'

#### "Target"

backward chaining "target" architecture driven by user query process on-the-fly users pull data information system dynamic archive user data is sacred

#### "Waterfall"

forward chaining "tier" architecture driven by raw data process in pipeline operators push data results in releases static archive raw data is obsolete

Deploy ESCAPE-EOSC on analysis data!

Deploy ESCAPE-EOSC on analysis data!
"Realistic" Monte Carlo data available.

Deploy ESCAPE-EOSC on analysis data!
"Realistic" Monte Carlo data available.

 Exploit "data lake" prototype (WP2) between GSI and RuG.

Deploy ESCAPE-EOSC on analysis data!

"Realistic" Monte Carlo data available.

 Exploit "data lake" prototype (WP2) between GSI and RuG.

 Good news: we have time for proof-of-principle study!