# Models and Algorithms for (the future of) High-Energy Physics

Maurizio Pierini

CERN

European Research Council

erc

# In this talk

- *Big-data in real time*
  - *"data scouting" with trigger-level analysis*

- *Big-Data tools and High-Energy physics (HEP) workflows*

- *HPC centres & HEP computing workflows*
  - *opportunistic processing*
  - *distributed training for Machine Learning*

I am replacing M. Zanetti here. I am trying to follow his initial idea about this talk, but with some personal point of view.
Not sure this matches 100% what this talk was about and your expectations.
I hope this will be useful nevertheless.

**https://www.youtube.com/watch?v=jDC3-QSiLB4**

# The LHC Big Data problem

HLT farm

Offline computing

Data Analysis

**Data Flow**

- 40 MHz in / 100 KHz out
- ~ 500 KB / event
- Processing time: ~10 μs
- Based on coarse local reconstructions
- FPGAs / Hardware implemented

# The LHC Big Data problem

L1 trigger

HLT farm

Offline computing

Data Analysis

**Data Flow**

- 100 KHz in / 1 KHz out
- ~ 500 KB / event
- Processing time: ~30 ms
- Based on simplified global reconstructions
- Software implemented on CPUs

# The LHC Big Data problem

L1 trigger

HLT farm

Offline computing

Data Analysis

**Data Flow**

- 1 KHz in / 1.2 kHz out

- ~ 1 MB / 200 kB / 30 kB per event

- Processing time: ~20 s

- Based on accurate global reconstructions

- Software implemented on CPUs

# The LHC Big Data problem

L1 trigger

HLT farm

Offline computing

Data Analysis

**Data Flow**

- Up to ~ 500 Hz In / 100-1000 events out

- <30 KB per event
- Processing time irrelevant
- User-written code + centrally produced selection algorithms

# HEP, Big Data & Real Time Processing

# The Trigger Problem

◉ *Too many data, too large data → need to filter online*

◉ *Filters based on pheno bias:* **<u>we might be loosing good events</u>**



100 KHz

1 KHz
1 MB/evt

40 MHz

L1 trigger

High-Level Trigger

▸ *L1 trigger: local, hardware based, on FPGA, @experiment site*

▸ *HLT: local/global, software based, on CPU, @experiment site*

▸ *Offline: global, software based, on CPU, @CERN T0*

▸ *Analysis: user-specific applications running on the grid*

# Doing more with less



**_Real-time new physics search with large datasets_**

- ◉ _Run reconstruction in the trigger farm_

- ◉ _Avoid resource limitations: write less information (a few floats) for more events_

- ◉ _Probes **unexplored territory**, previously left behind_

**_Problem: practical (so far) only for specific topologies_**

# Why did we do this?

◉ *In Run I, dijet search was the first BSM analysis published by CMS*

   ◉ *Quick improved results from Tevatron in a wide range of mass spectra*

   ◉ *Quickly forced to reduce mass range under investigation, due to increasing trigger rates vs limited resources*

◉ *Scouting was introduced to recover the lost territory (500 to 1100 GeV)*

**3 pb⁻¹ @7 TeV in 2010**



**1 fb⁻¹ @7 TeV in 2011**

# The first attempt

# The first attempt

# What we accomplished

- *Recovered sensitivity to 500 GeV resonances*

- *Reached limitation of L1 seed-> need to improve our hardware trigger (more on this later)*

- *Now extending the method to more final states (collected x3 more data than the rest of CMS in 2017)*

# What we accomplished

- *Kept sensitivity to 500-1500 GeV resonances*

- *Current limitation is L1 efficiency*

- *Can probe lower couplings by collecting more data*

# An established approach

A search for narrow resonances decaying into dijet final states is performed on data from proton-proton collisions at a center-of-mass energy of 8 TeV, corresponding to an integrated luminosity of 18.8 fb$^{-1}$. The data were collected with the CMS detector using a novel technique called data scouting, in which the information associated with these selected events is much reduced, permitting collection of larger data samples. This technique enables CMS to record events containing jets at a rate of 1 kHz, by collecting the data from the high-level-trigger system. In this way, the sensitivity to low-mass resonances is increased significantly, allowing previously inaccessible couplings of new resonances to quarks and gluons to be probed. The resulting dijet mass distribution yields no evidence of narrow resonances. Upper limits are presented on the resonance cross sections as a function of mass, and compared with a variety of models predicting narrow resonances. The limits are translated into upper limits on the coupling of a leptophobic resonance $Z'_B$ to quarks, improving on the results obtained by previous experiments for the mass range from 500 to 800 GeV.

Search for narrow resonances in dijet final states at $\sqrt{s} = 8\,\text{TeV}$ with the novel CMS technique of data scouting

The CMS Collaboration*

16

# Next-step:Trigger-less

- LHCb & ALICE soon to start a detector & online-infrastructure upgrade. Final goal is to

  - Read ALL collisions

  - Process them in real time

  - Align & calibrate detector at the same time

- The ultimate extrapolation of the scouting paradigm: Can take more data -> increase detector precision

**LHCb Upgrade Trigger Diagram**

**30 MHz inelastic event rate (full rate event building)**

**Software High Level Trigger**

Full event reconstruction, inclusive and exclusive kinematic/geometric selections

Buffer events to disk, perform online detector calibration and alignment

Add offline precision particle identification and track quality information to selections

Output full event information for inclusive triggers, trigger candidates and related primary vertices for exclusive triggers

**2-5 GB/s to storage**

bad        good

Misaligned        Aligned

LHCb OT Preliminary     $t_0$ Updated   $t_0$ Not updated

23/4/2016 - 9/9/2016

Run number [a.u.]

European Research Council

erc

# HEP, Cloud & HPCs

# The WLCG Grid

◉ *170 centres in 42 countries, for central processing and analysis-related user jobs*

◉ *1M cores*

◉ *1 EB storage*

◉ *>2M jobs & 3 PB moved /day*



~1M Cores



3PB/day

- alice
- atlas
- cms
- lhcb

European Research Council

# The challenge ahead

- *The evolving conditions of the machine are drifting the experiments to more prohibitive environments (luminosity comes with a cost)*

- *More (& bigger) events to handle*

- *More noise from pileup interactions*

- *Increase in resources will not scale with needs*

- *Flat (or decreasing?) budget*

- *(Non linearly) increasing demand ys to do*

# The challenge ahead



- *Event complexity, volume, and number will challenge the current paradigm*

- *Assuming flat budget, we simply cannot keep doing things as we do now*

# More CPU: Cloud

- *The growing complexity of LHC events is forcing us to look for more resources, particularly for the computation-heavy central reconstruction*

  - *CERN extended the T0 center by adding a site in Wigner (Hungary)*

  - *Similar approach used by T1s (e.g., CNAF T1 extended with CPUs in Bari)*

- *Paradigm extended opportunistically to Cloud services and HPC sites*

# More CPU: Cloud

◉ *The growing complexity of LHC events is forcing us to look for more resources, particularly for the computation-heavy central reconstruction*

◉ *CERN extended the T0 center by adding a site in Wigner (Hungary)*

◉ *Similar approach used by T1s (e.g., CNAF T1 extended with CPUs in Bari)*

◉ *Paradigm extended opportunistically to Cloud services and HPC sites*

# More CPU: HPCs

◉ *Similar tests done on HPC sites (NERSC Cori)*

◉ *x86 machines, in very different setup than T0/T1/T2/T3 sites*

◉ *Challenge stands in working out all details and finding workarounds to incompatible setups (e.g., not-supported components)*

◉ *Result: storage-less site added to the CMS grid as yet another Tx*



dashboard

**Running Job Cores**
168 Hours from 2016-10-04 to 2016-10-11 UTC

■ T3_US_NERSC

European Research Council
erc

# A convenient new Paradigm?

| Grid | Cloud | HPC |
|------|-------|-----|
| • Virtual Organizations (VOs) of users trusted by Grid sites<br><br>• VOs get allocations ➔ **Pledges**<br>  –Unused allocations: opportunistic resources<br><br>**"Things you borrow"** | ▪ Community Clouds - Similar trust federation to Grids<br><br>▪ Commercial Clouds - **Pay-As-You-Go** model<br>  •Strongly accounted<br>  •Near-infinite capacity ➔ **Elasticity**<br>  •Spot price market<br><br>**"Things you rent"** | ▪ Researchers granted access to HPC installations<br><br>▪ Peer review committees award **Allocations**<br>  •Awards model designed for individual PIs rather than large collaborations<br><br>**"Things you are given"** |
| **Trust Federation** | **Economic Model** | **Grant Allocation** |

CERN

🔶 Fermilab

erc — European Research Council

HEP, Big Data & "offline" Processing

# The foreseen analysis workflow

- **Central processing:** *Runs @T0. Start from RAW data and creates a collection of "Primary" Datasets then distributed to T1s*

- **Data skimming:** *Runs @T0 or T1s. From the Primary Datasets, produce "Secondary datasets" by removing events (so why did you take them to start with?) or reducing the information (data compression)*

- **Data analysis:** *runs on Secondary Datasets, applying analysis specific selection, reconstructing high-level objects on which signal-to-background discriminating quantities are computed. Runs on T3s, on the Grid, etc*

- **Result extraction:** *typically a ML fit, based on data distributions in signal region and control region + prediction from MC simulation (runs on laptops)*



[DATA HANDLING]
**TOO MUCH INFORMATION**

# It didn't really go like that

◉ *Disk issue is less (but still quite) serious than anticipated:*

  ◉ *We (all) introduced **AODs** (500-1000 kB/evt)  compressed version of RECO data format. We saved disk, so we just distributed Primary Datasets rather than using the (very bad) Secondary datasets*

  ◉ *With gain detector understanding, we (CMS) then moved forward to **miniAODs** (30 kB/evt) and **nanoAOD** (3 kB/evt), compressed data formats with top-bottom object definition, serving >80% of the analysis use cases*

◉ *Large **demand of CPU** faced breaking the paradigm rigidity:*

  ◉ *T1s and T2s interconnection was improved. Now one runs a job somewhere accessing data somewhere else*

◉ ***Still, we would use more disk & CPU if we had it …***

# Big Data tools & HEP

- *A lot is happening outside HEP*

  - *full data-scientist echo-system*

  - *big-data handling tools*

- *But we have specific tools (ROOT)*

  - *optimized on our use cases*

  - *very competitive on I/O point of view*

  - *long-term future guaranteed (we develop it for ourselves)*

- *A big effort to integrate ROOT & outside-world big-data tools is ongoing, with promising results*

# BigData tools integration

◉ *Effort to modernise approach to data analysis by integrating/creating data-analytics tools for physics analyses*

◉ *Goals:*

  ◉ *Reduce number of intermediate processing+storage steps*

  ◉ *Allow analyses to run on (mini)AODs way down to the publication-ready plots in a data-science framework*

# The Final Goal

- *Develop a CMS analysis workflow in Apache Spark:*

- *Full ROOT -> Spark analysis workflow with*

  - *Event selections*

  - *Data-Simulation comparison*

  - *Data reduction scheme in Spark*

- *Provided services*

  - *Machine-Learning toolkit*

  - *Data in memory for fast training*

- ***Benchmarking all that and compare performance/results with standard workflow***

# How It works

- ◉ *Dedicated libraries to implement the workflow:*

  - ◉ *XrootD connector to access files on CERN EOS filesystem: (So far) from public area. Authentication via certificate is under developement*

  - ◉ *Spark-root: Read ROOT object collections and automatically infer their class schema*

  - ◉ *Histogrammar (by DIANA-HEP): To fill histograms passing lambda functions and use them in the same way as transformations are used in Apache Spark*

- ◉ *100% data-science echosystem, compatible with ROOT I/O but w/o ROOT installation*

# What Can It Do?

I/O from structured ROOT files (e.g., experiment-specific data files)

```
-- patMuons_slimmedMuons__RECO_: struct (nullable = true)
    |-- present: boolean (nullable = true)
    |-- patMuons_slimmedMuons__RECO_obj: array (nullable = true)
        |-- element: struct (containsNull = true)
            |-- ... struct (nullable = true)
                |-- vertex_: struct (nullable = true)
                    |-- fCoordinates: struct (nullable =
true)
                        |-- fX: float (nullable = true)
                        |-- fY: float (nullable = true)
                        |-- fZ: float (nullable = true)
                |-- p4Polar_: struct (nullable = true)
                    |-- fCoordinates: struct (nullable = true)
true)
                        |-- fPt: float (nullable = true)
                        |-- fEta: float (nullable = true)
                        |-- fPhi: float (nullable = true)
                        |-- fM: float (nullable = true)
                |-- qx3_: integer (nullable = true)
                |-- pdgId_: integer (nullable = true)
                |-- status_: integer (nullable = true)
```

```python
# read in the data
df = sqlContext.read\
    .format("org.dianahep.sparkroot.experimental")\
    .load("hdfs://path/to/files/*.root")

# count the number of rows:
df.count()

# select only muons
muons =
df.select("patMuons_slimmedMuons__RECO_.patMuons_slim
medMuons__RECO_obj.m_state").toDF("muons")

# map each event to an invariant mass
inv_masses = muons.rdd.map(toInvMass)

# Use histogrammar to perform aggregations
empty = histogrammar.Bin(200, 0, 200, lambda row: row.mass)
h_inv_masses = inv_masses.aggregate(empty,
    histogrammar.increment,
    histogrammar.combine)
```

On-the-fly Feature Engineering

# HPC, HEP & Deep Learning

http://www.asimovinstitute.org/neural-network-zoo

- ◉ *With Deep Learning gaining territory in HEP$^{(*)}$, NN training will become soon a new workflow for large HEP experiments*

  - ◉ *Experiments will want to maximise performances*

  - ◉ *Fast turn-around for new trainings, as long as new data are collected*

  - ◉ *Need dedicated hardware to be effective (TPU, GPU, etc)*

- ◉ *Ideal use case to integrate HEP workflows into network of HPC sites*



J(w)

Initial weight

Gradient

Global cost minimum
$J_{min}(w)$

*(*) more on this in afternoon seminar*

# Hyper-parameter optimization

◉ *Not only the best set of parameters, but also the best network overall:*

  ◉ *how many layers?*

  ◉ *how many nodes/layers?*

  ◉ *which activation function?*

◉ *Answers to be find by Optimization algorithm*

  ◉ **Bayesian Optimization**

  ◉ **Evolutionary Algorithms**

  ◉ *...*

◉ **One extra reason to train production-ready algorithms @HPC sites**



https://tinyurl.com/yc2phuaj

- Objective function is approximated as a multivariate gaussian
- Measurements provided one by one to improve knowledge of the objective function
- Next best parameter to test is determined from the acquisition function

- Using the python implementation from https://scikit-optimize.github.io



- Chromosome crossover:
  - Let Parent A be more fit than Parent B
  - For each parameter $p$, generate a random number $r$ in $(0, 1)$ to find $p_{child}$

$$p_{child} = (r)(p_{Parent\ A} - p_{Parent\ B}) + p_{Parent\ A}$$

- Non-uniform mutation (Michalewicz):
  - In generation $g$ out of a total $G$ generations, for each parameter $p$ in a child, generate random numbers $r_1, r_2 \in (0, 1)$ to define a mutation $m$:

$$m = \left(1 - r_1^{\left(1-\frac{g}{G}\right)^3}\right) * \begin{cases} (p_{MAX} - p_{child}) & IF \quad r_2 > 0.5 \\ (p_{LOW} - p_{child}) & IF \quad r_2 \leq 0.5 \end{cases}$$
$$p_{child} = p_{child} + m$$

# K-folding cross validation

- *Assigning uncertainties to training-goodness figures of merit to establish ACTUAL improvements*

  - *Are ROC AUCs 98.754 and 97.998 actually different?*

- *Done by different training vs validation dataset splits*

- *Average performances and performance dispersion allow to "measure" mean and variance*

- *Multiply workflow computing needs by K*



K-folding Layout



- One master running the optimization. Receiving the average figure of merit over $N_F$ folds of the data
  - $N_G$ groups of nodes training on a parameter-set on simultaneously
    - $N_F$ groups of nodes running one fold each

07/12/8

# Parallelisms

● ___Data Parallelism:___ *master nodes handle parameter setting, receives gradients from workers and distribute new parameter values. Good for datasets with many events*



1) Compute gradient, send to Master

2) Update network weights
$$\vec{w} \rightarrow \vec{w} - \eta \nabla Q(\vec{w})$$

$\nabla Q(\vec{w})$

3) Send new weights to Worker

https://arxiv.org/abs/1712.05878



Training master group 0, subrank 0

Parameter-set group 0

TM1 12   TMN$_M$

TW0 0   TW0

TWN$_W$ N$_W$   TWN$_W$



https://github.com/duanders/mpi_learn

38

# Deployed @HPC centres

- Speed up in training generative adversarial networks on Piz Daint CSCS and Titan ORNL supercomputers
  - Using easgd algorithm with rmsprop
  - Speed up is not fully efficient. Bottlenecks to be identified



Low performance degradation



NVIDA P100 on Piz Daint, CSCS



NVIDA K20 at Titan, ORNL

# Parallelisms

◉ *Gradient Parallelism: parallelise gradient computation of a single batch on multiple workers. Good for datasets with large-size examples*



Some performance degradation
Mostly in the low energy regions for large batchsize





Sofia V. @ https://sites.google.com/nvidia.com/ai-hpc

**Communication through Horovod with fast GPU-GPU communication (nVidia NCCL)**

07/12/18

Deep Learning
Training & Optimization, J-R Vlimant, CHEP18

40

# Parallelisms

● *Model Parallelism: compute gradients for different parts of the networks on different workers. Good for large models*



**Requires good devide-to-device communication**
**Used TensorFlow native multi-device manager**
**Aiming to test this on machines with multi-gpu nodes (Summit)**

# Conclusions

- *Large-scale computing is a consolidated tradition in HEP*

- *Things didn't go as planned, since new developments made us more competitive at fixed/decreasing resource budget (e.g., scouting & real-time processing to do more with less)*

- *New challenges ahead call for qualitative and quantitative (more CPU + GPU + …) and qualitative (Big data tools integration) improvements*

- *Exploiting existing sites (HPC) rather than building dedicated facilities*

  - *Opportunistic cloud computing*

  - *Large-scale training as a service on GPU clusters*

- ***Challenges ahead, time for brave people to come out with ideas***

# Backup

# Distributed training

**Use keras 2.13 /Tensorflow 1.9 (Intel optimised)**

- AVX512 –FMA-XLA support
- Intel® MKL-DNN (with 3D convolution support)

**Optimised multicore utilisation**

- inter_op_paralellism_threads/intra_op_paralellism threads

**Horovod 0.13.4**

- Synchronous SGD approach
- MPI_AllReduce

**Run on TACC Stampede2 cluster:**

- Dual socket Intel Xeon 8160
- 2x 24 cores per node, 192 GB RAM
- Intel® Omni-Path Architecture

**Test several MPI scheduling configurations**

- 2,4, 8 processes per nodes.
- Best machine efficiency with 4 processes/node

14

Some performance degradation
Mostly in the low energy regions for large batchsize



Ratio of Ecal and Ep

Data
BatchSize=1000
BatchSize=4000
BatchSize=10000



**High Energy Physics: 3D GANS Training Scaling Performance**
Intel 2S Xeon(R) on Stampede2/TACC, OPA Fabric
TensorFlow 1.9+MKL-DNN+horovod, IMPI, Core Aff. BKMs, 4 Workers/Node

Sofia V. @ https://sites.google.com/nvidia.com/ai-hpc