# Status and Future of Analysis
# The ROOT Perspective
# (covering Python and C++)
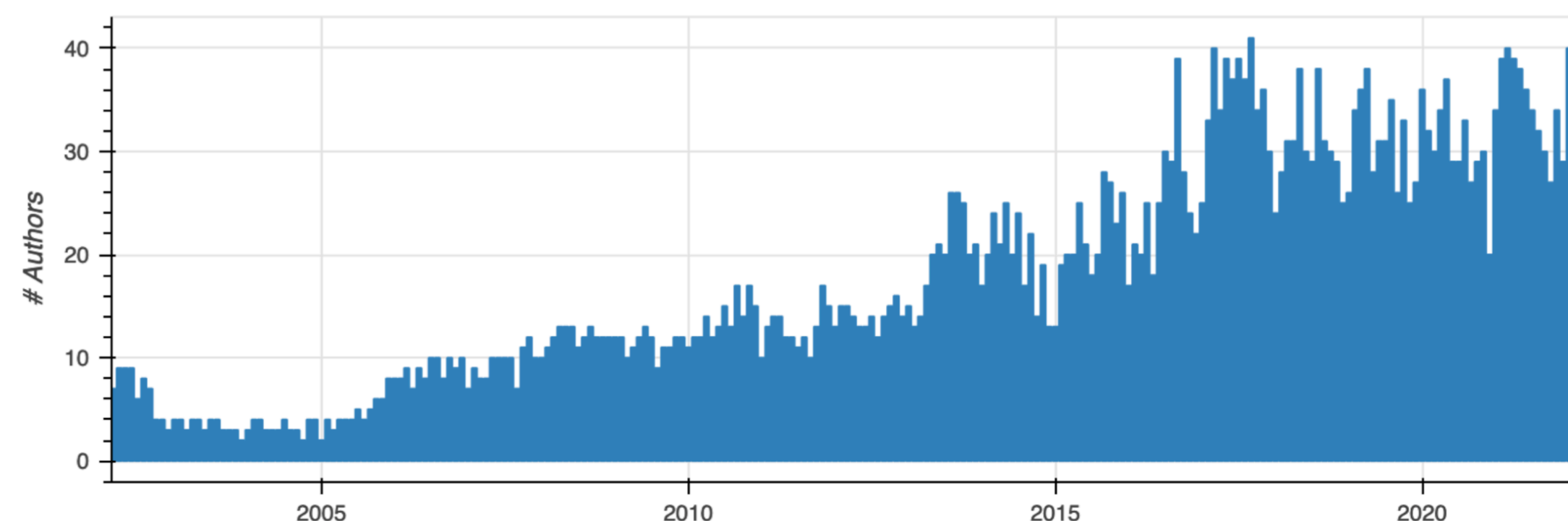
*Lorenzo Moneta*

*(CERN/EP-SFT )*

28/9/2022

# Introduction

- ROOT is an open source, community-driven project
  - C++ under the hood, but looks like a Python module
- ROOT is used by all HEP experiments
  - more than 1EB of data stored in ROOT I/O format.
- Serving both analysis and experiments software
  - excellent interactions with users (~1000 Root Forum post/month
- Despite existing for many years ROOT is still very active
  - high number of contributors/month
- Ongoing massive innovation
  - I/O, analysis interfaces, python, RooFit, ML, web-graphics, etc..

# ROOT Core Components

- I/O : reading and writing data efficiently
- Analysis interfaces: histograms and RDataFrame
- Interactivity: C++ interpreter and Python bindings
- Math libraries: PRNG, numerical algorithms (e.g. Minuit)
- Statistics and modeling: fitting
- Graphics: scientific visualization
- Machine learning: TMVA, model evaluations and interoperability

Major recent developments in some of these components and the future plans

# Outline

- Analysis landscape:
  - I/O developments
    - RNtuple
  - Python interfaces:
    - PyROOT
  - data exploration, declarative analysis
    - RDataFrame
  - machine learning, interoperability and analysis integration
    - TMVA/SOFIE
  - statistical analysis and modelling
    - RooFit
- What we have in ROOT now and the planned future developments

*With lots of material presented coming from R. Giraud's ICHEP'22 talk*

# Future Analysis Landscape

- Expect a higher number of events and more complex problems
  - larger data sets
  - higher dimensional problems requiring more CPU (or GPU)
- ROOT is well suited for this
  - ROOT columnar data (TTree and new RNTuple) designed for analysis on large data sets
  - New versatile analysis interfaces (RDataFrame)
  - Language and tool interoperability (e.g. with ML libraries)
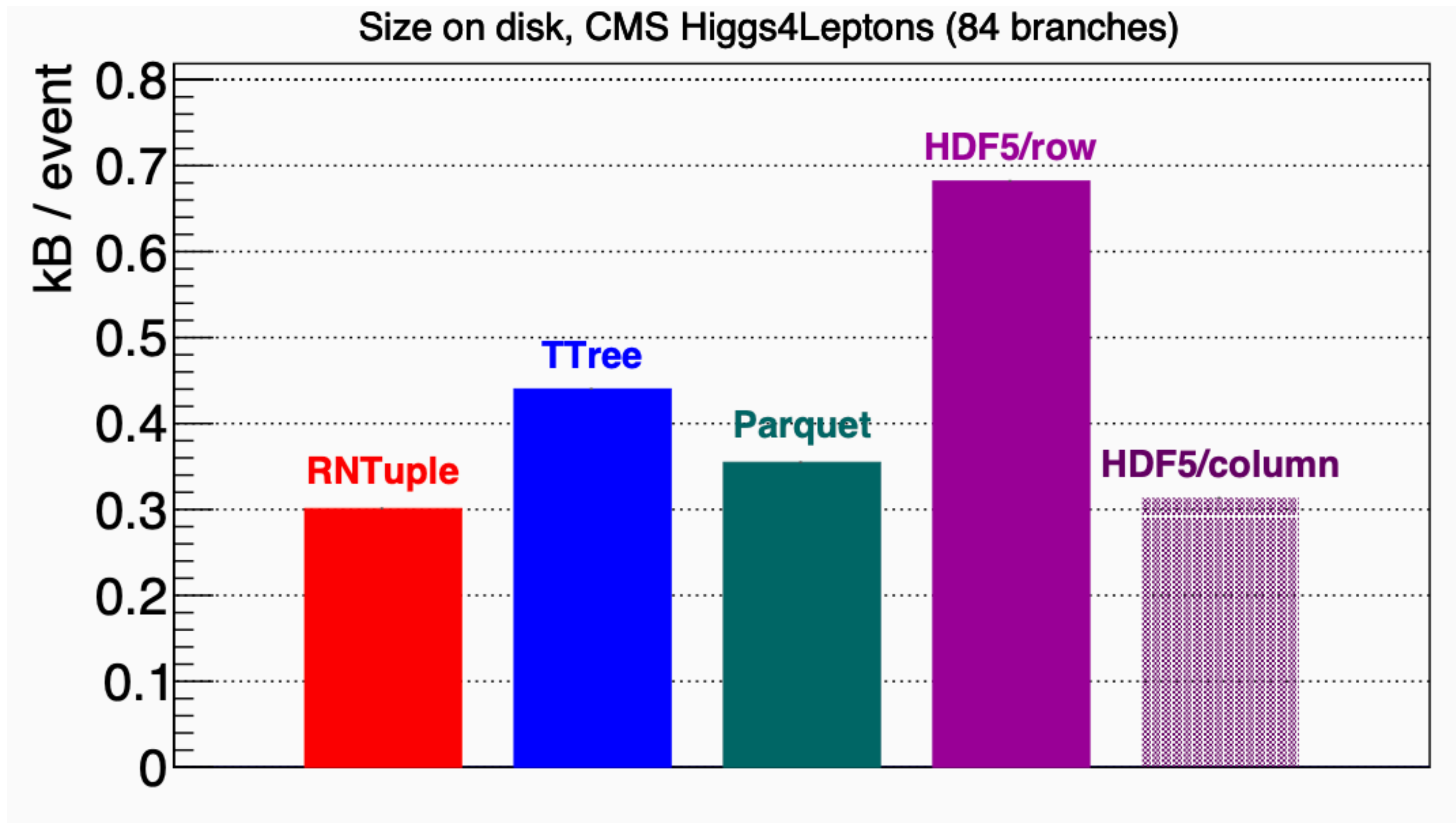  - Statistical analysis of complex problems (RooFit)

# RNTuple

- Successor of TTree
  -  columnar storage of event data optimised for selective reads
- Integration with LHC experiments is ongoing, to be production ready for HL-LHC
- Based on 25+ years of TTree experience,
- Redesigned I/O subsystem providing:
  - Less disk and CPU usage for the same data content
  - 15% smaller files, × 2–5 better single-core performance
  - 10 GB/s per box and 1 GB/s per core sustained end-to-end throughput
- And also:
  - Systematic use of exceptions to prevent silent I/O errors
  - Efficient support of modern hardware (built for multi-threading and async I/O)
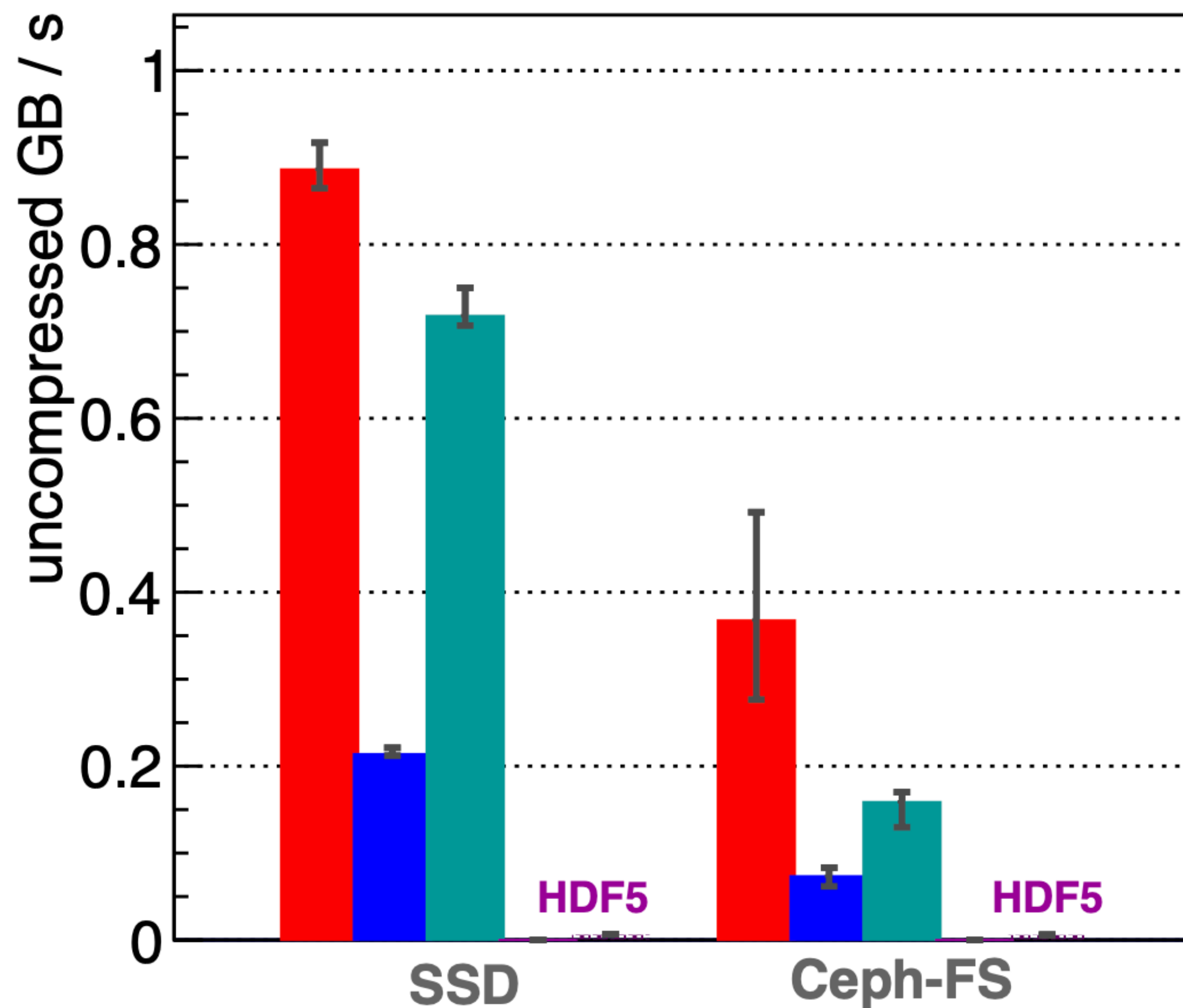  - Native support for object stores
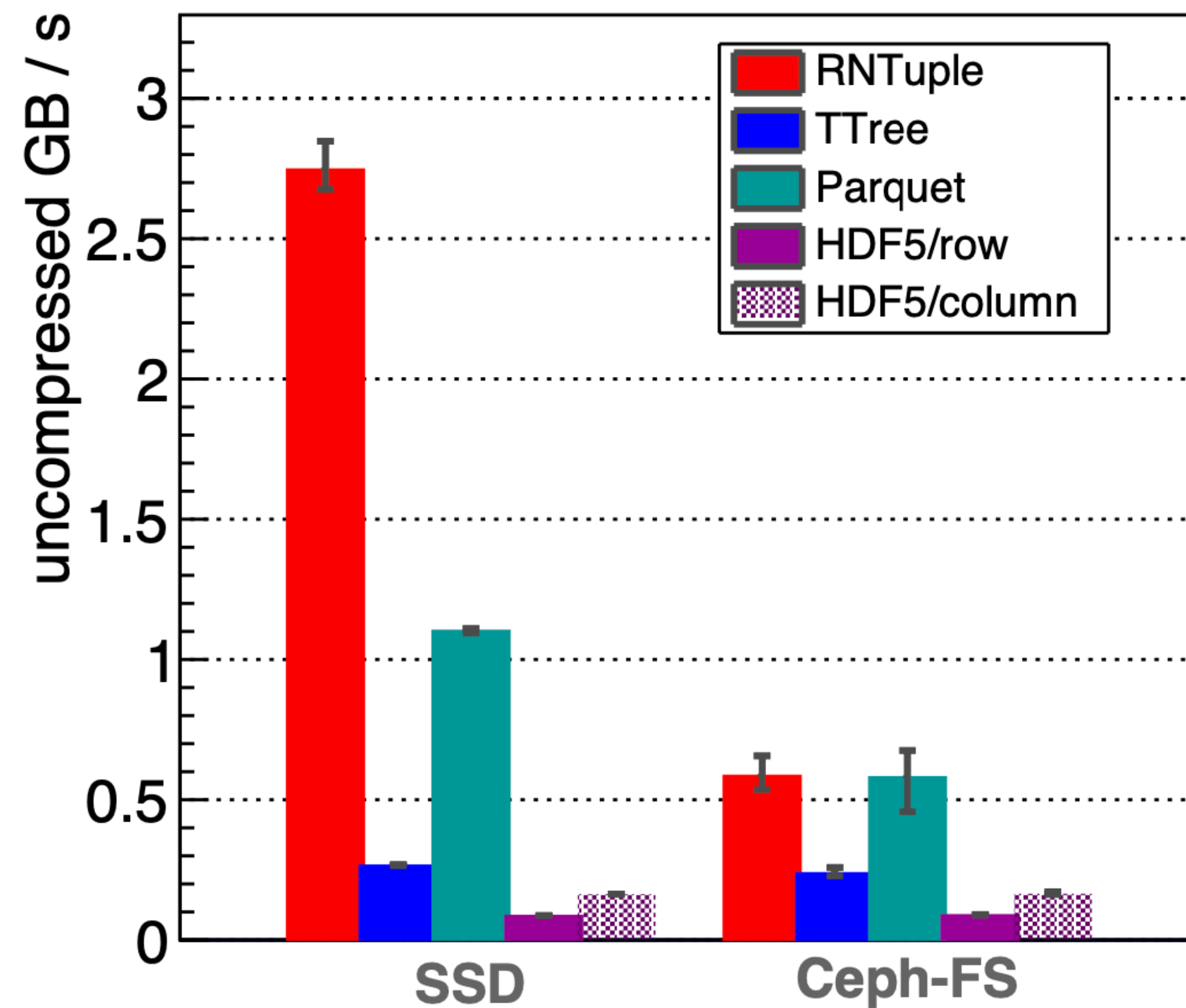
# RNTuple Performances: I/O Size



Size on disk, CMS Higgs4Leptons (84 branches)

CMS Higgs4Leptons (10/84 branches)

LHCb B2HHH (10/26 branches)

# RNTuple Plans

- Smaller files and significantly faster reads compared to TTree

- Modern and robust API

- Capable of making efficient use of modern devices and storage systems (such as SSD, object stores, and many cores)

- More results and developments will be presented at the coming ACAT22

- RNTuple is work in progress in *ROOT::Experimental*

  - the on-disk format is still subject to small changes!

  - happy to get your feedback!

| | LS 2 | | LHC Run 3 | | | | LS 3 | | | Run 4 (HL-LHC) |
|---|---|---|---|---|---|---|---|---|---|---|
| 2020 | 2021 | 2022 | 2023 | 2024 | 2025 | 2026 | 2027 | 2028 | 2029 | |

RNTuple work in progress in ROOT::Experimental          RNTuple goes production, adoption phase

# PyROOT

- **ROOT is not only C++**
- It can be used as a Python module thanks to Python-C++ bindings

```
> import ROOT
> h1 = ROOT.TH1D("h1","my hist",100,0.,10.)
```

- All C++ functionality is available in Python
  - can also access user-defined C++ code from Python
- ROOT-specific pythonizations added on top of C++ API
  - make ROOT classes more "pythonic" to use
- PyROOT is now built on top of Cppyy *(from W. Lavrijsen)*
  - based on Cling and ROOT type system
  - support modern C++

# ROOT Pythonizations

- Make C++ classes easier to use in Python
- Available for some classes:
  - TFile, TTree
  - RDataFrame, RVec
  - Main RooFit classes
  - TMVA

**PyROOT**

In the same way as for **TDirectory**, it is possible to inspect the content of a **TFile** object from Python as if the directories and objects it contains were its attributes. For more information, please refer to the **TDirectory** documentation.

In addition, **TFile** instances can be inspected via the Get method, a feature that is inherited from **TDirectoryFile** (please see the documentation of **TDirectoryFile** for examples on how to use it).

In order to write objects into a **TFile**, the WriteObject Python method can be used (more information in the documentation of **TDirectoryFile**).

**PyROOT** modifies the **TFile** constructor and the **TFile::Open** method to make them behave in a more pythonic way. In particular, they both throw an OSError if there was a problem accessing the file (e.g. non-existent or corrupted file).

This class can also be used as a context manager, with the goal of opening a file and doing some quick manipulations of the objects inside it. The **TFile::Close** method will be automatically called at the end of the context. For example:

```python
from ROOT import TFile
with TFile("file1.root", "recreate") as outfile:
    hout = ROOT.TH1F(...)
    outfile.WriteObject(hout, "myhisto")
```

- Documented in reference guide (in class documentation)
- Working on pythonizations of Histograms and Graph classes
  - better inter-operability with Numpy
  - implementing <u>UHI indexing</u>:  `hist[bin1:bin2]`

# Examples of Pythonizations

- ## TTree and TFile

```python
import ROOT

f = ROOT.TFile('myfile.root')

t = f.mytree
# vs f.GetObject('mytree')
```

TFile is a (dynamic) Python proxy of a C++ class

f is a (dynamic) Python proxy of a C++ object
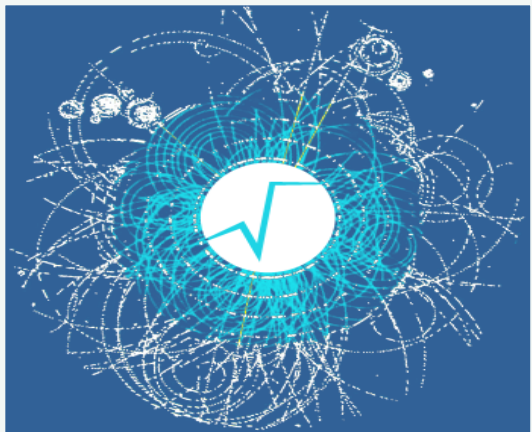
Pythonization: access tree as an attribute

- ## Process with RDataFrame and convert to Numpy arrays

```python
from ROOT import RDataFrame
df = RDataFrame('myTree', 'file.root')

# Apply cuts, define new columns
df = df.Filter('x > 0').Define('z', 'x*y')

cols = df.AsNumpy() # cols is a dict of NumPy arrays
```

**Tutorial here**

# Analysis Landscape

## Analysis life cycle

skimming,
ntuple production

quick exploration,
first implementation

systematics,
scale out

statistical
analysis

## Platforms

laptop or PC

many-core machine

computing cluster
+ job submission

## Analysis languages

↓ ~50%   C++
↑ ~50%   Python

## Storage

local disk
fast-access network storage
EOS or other not-so-fast backend

*from E. Guiraud: RDF@ICHEP 2022*

# RDataFrame

## A Swiss army knife for data analysis

**Analysis life cycle**

skimming,
ntuple production

quick exploration,
first implementation

systematics,
scale out

**Platforms**

laptop or PC

many-core machine

computing cluster
+ job submission

**Analysis languages**

⬇ ~50%  C++
⬆ ~50%  Python

**Storage**

local disk
fast-access network storage
EOS or other not-so-fast backend

**ROOT.RDataFrame** is a modern analysis interface that addresses all these use cases with **one high-level programming** model that performs well, scales well and enables **HEP-specific ergonomics**, in C++ and Python.

*from E. Guiraud: RDF@ICHEP 2022*

# RDataFrame Code

**Example of RDF code in Python**

```python
df = ROOT.RDataFrame(dataset)  ················  on this (ROOT, CSV, …) dataset

df = df.Filter("x > 0")  ·······················  only accept events for which x > 0      event selection

       .Define("r2", "x*x + y*y")  ···············  define r2 = x² + y²      derived quantities,
                                                                             object selections

rHist = df.Histo1D("r2")  ·····················  plot r2 for events that pass the cut

df.Snapshot("newtree", "out.root")  ···········  write the skimmed data and r2      data aggregations
                                                 to a new ROOT file
```

Users can inject **arbitrary code** at all steps,
which makes this relatively simple API extremely versatile.

**Physicists describe analysis ingredients, ROOT handles technicalities**

# RDataFrame Code: MT

## Switch to multi-thread execution (Python):

```
ROOT.EnableImplicitMT()
```
Run a multi-thread event loop

```
df = ROOT.RDataFrame(dataset)
```
on this (ROOT, CSV, ...) dataset

```
df = df.Filter("x > 0")
```
only accept events for which $x > 0$

```
    .Define("r2", "x*x + y*y")
```
define $r2 = x^2 + y^2$

```
rHist = df.Histo1D("r2")
```
plot r2 for events that pass the cut

```
df.Snapshot("newtree", "out.root")
```
write the skimmed data and r2 to a new ROOT file

# RDataFrame Code: DistRDF

## Switch to distributed execution (Python):

```python
cluster = dask_jobqueue.HTCondorCluster(n_workers=64)
df = RDataFrame(dataset, daskclient=Client(cluster))     ········· connect to
                                                                   HTCondor via Dask
df = df.Filter("x > 0")

        .Define("r2", "x*x + y*y")     ···········································  other code
                                                                   stays the same
rHist = df.Histo1D("r2")

df.Snapshot("newtree", "out.root")
```

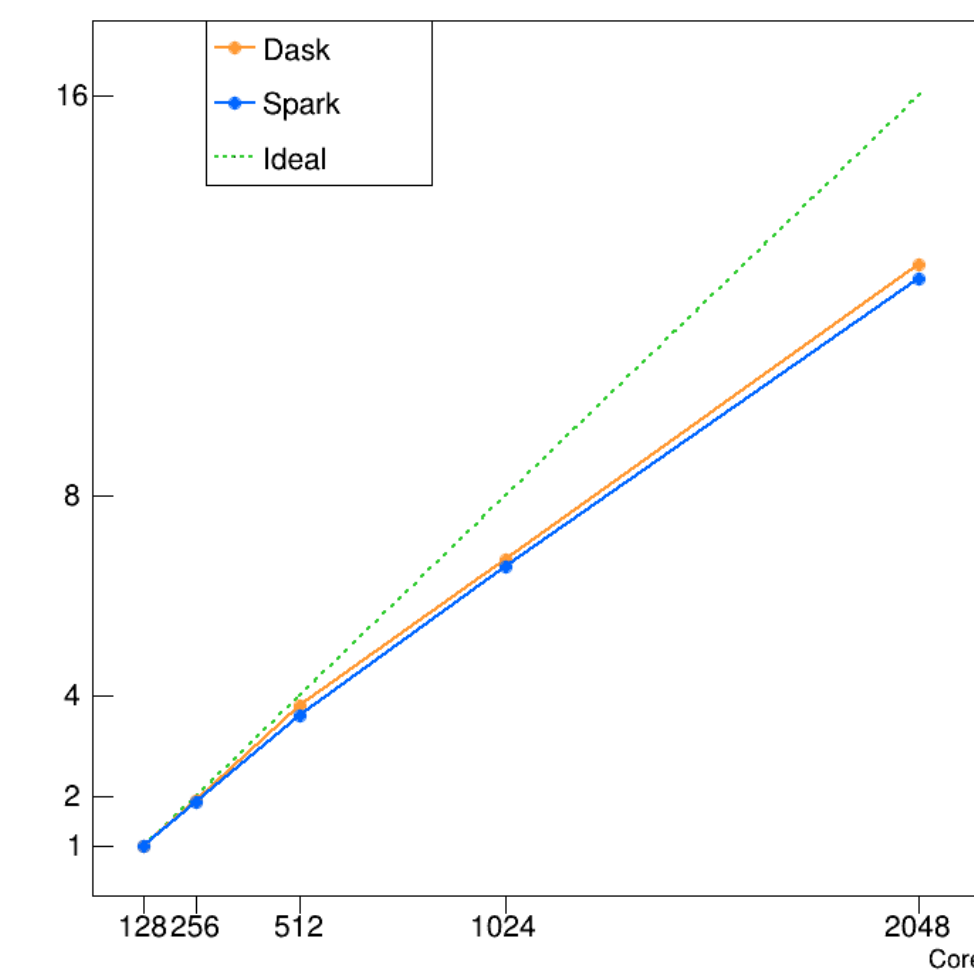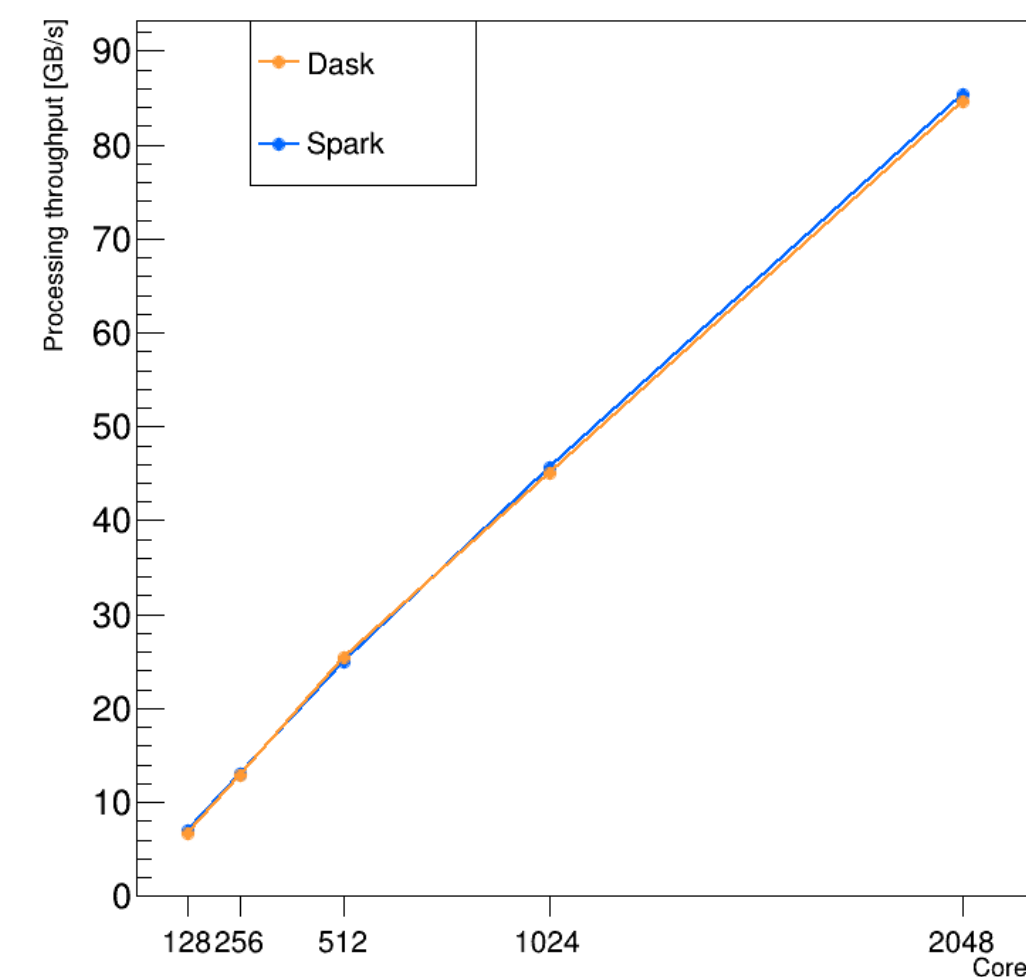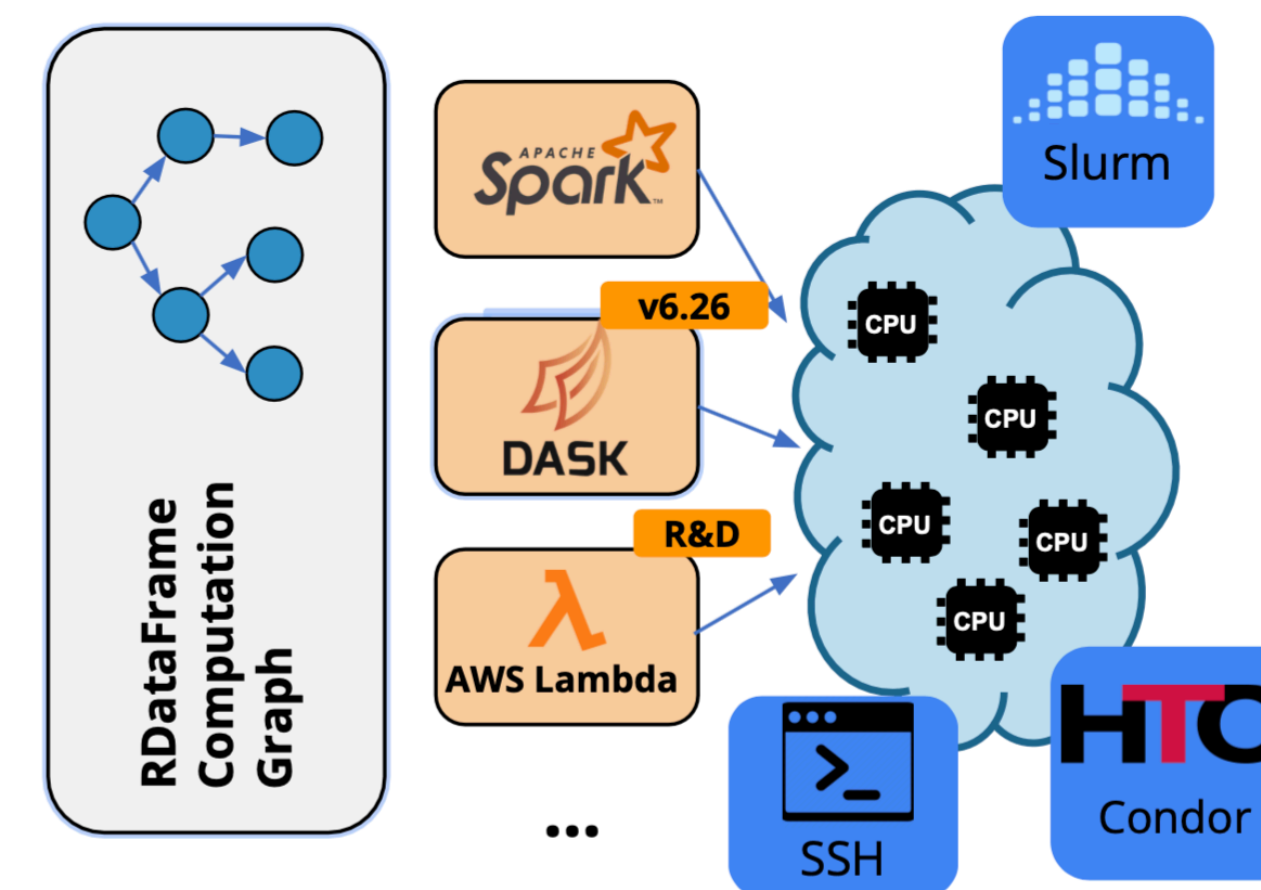**Available since v6.26 (experimental)**     Also see this tutorial, the docs, the recent ATTF talk

# Distributed RDF



- Enables **interactive large-scale distributed** data analysis

- Uses Python RDF API
  (C++ event loop)

- High throughput with native ROOT I/O

- **Spark/Dask/HTCondor/Slurm/ SSH**…. take care of scheduling and resource management

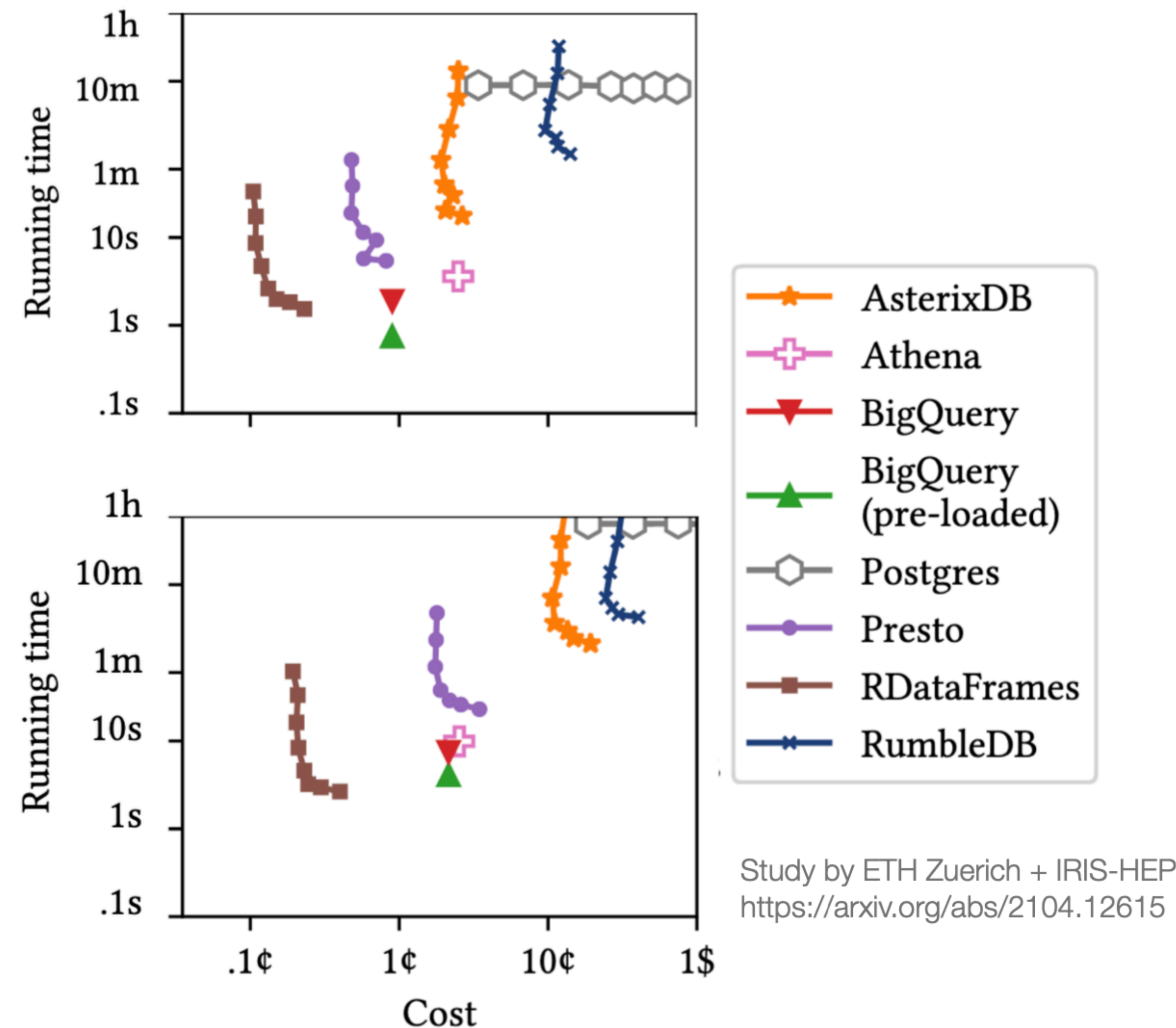- Transparently merges results coming from different computing nodes



**Distributed RDF benchmark (dimuon, 4000x data)**
**See this talk**

# RDF Performances

- **RDF enables fast turnaround** for complex analysis use cases
- **scales well** to many cores, many nodes, many histograms
- Independent study shows RDF analysis to be significantly **faster**
- Performance is always ongoing work:
  - we are constantly looking for feedback/use cases



Study by ETH Zuerich + IRIS-HEP
https://arxiv.org/abs/2104.12615

*from E. Guiraud: RDF@ICHEP 2022*

# Systematic Variations



```python
                                    Python
nominal_hx =
    df.Vary("pt", "RVecD{pt*0.9, pt*1.1}", ["down", "up"])
      .Filter("pt > k")
      .Define("x", someFunc, ["pt"])
      .Histo1D("x")


hx = ROOT.RDF.VariationsFor(nominal_hx)
hx["nominal"].Draw()
hx["pt:down"].Draw("SAME")
```

*attach an up/down variation to "pt"*

*proceed as usual, as if working with nominal values only*

*obtain all variations*

**Variations automatically propagate** to selections, derived quantities and results.
**Multi-thread** and **distributed** execution **just works.**
Only needed quantities are re-computed, all in **one event loop.**

# RDF: Status and Plans

- **and many more features, e.g.**
  - transparent support for RNTuple
  - machine learning inference:
    - SOFIE, RBDT
  - interoperability with Numpy
- **and coming soon:**
  - bulk (batch) processing of data
    - improved performance
    - batch generator for ML tools
- **RDF keeps evolving**
  - cooperation with user community
- **Critical to focus on the right features**
  - **need your feedback and help !**

| Lazy action | Description |
|---|---|
| Aggregate() | Execute a user-defined accumulation operation on the processed column values. |
| Book() | Book execution of a custom action using a user-defined helper object. |
| Cache() | Cache column values in memory. Custom columns can be cached as well, filtered entries are not cached. Users can specify which columns to save (default is all). |
| Count() | Return the number of events processed. Useful e.g. to get a quick count of the number of events passing a Filter. |
| Display() | Provides a printable representation of the dataset contents. The method returns a ROOT::RDF::RDisplay() instance which can print a tabular representation of the data or return it as a string. |
| Fill() | Fill a user-defined object with the values of the specified columns, as if by calling Obj.Fill(col1, col2, ...). |
| Graph() | Fills a TGraph with the two columns provided. If multi-threading is enabled, the order of the points may not be the one expected, it is therefore suggested to sort if before drawing. |
| GraphAsymmErrors() | Fills a TGraphAsymmErrors. If multi-threading is enabled, the order of the points may not be the one expected, it is therefore suggested to sort if before drawing. |
| Histo1D(), Histo2D(), Histo3D() | Fill a one-, two-, three-dimensional histogram with the processed column values. |
| HistoND() | Fill an N-dimensional histogram with the processed column values. |
| Max() | Return the maximum of processed column values. If the type of the column is inferred, the return type is double, the type of the column otherwise. |
| Mean() | Return the mean of processed column values. |
| Min() | Return the minimum of processed column values. If the type of the column is inferred, the return type is double, the type of the column otherwise. |
| Profile1D(), Profile2D() | Fill a one- or two-dimensional profile with the column values that passed all filters. |
| Reduce() | Reduce (e.g. sum, merge) entries using the function (lambda, functor...) passed as argument. The function must have signature T(T,T) where T is the type of the column. Return the final result of the reduction operation. An optional parameter allows initialization of the result object to non-default values. |
| Report() | Obtain statistics on how many entries have been accepted and rejected by the filters. See the section on named filters for a more detailed explanation. The method returns a ROOT::RDF::RCutFlowReport instance which can be queried programmatically to get information about the effects of the individual cuts. |
| Stats() | Return a TStatistic object filled with the input columns. |
| StdDev() | Return the unbiased standard deviation of the processed column values. |
| Sum() | Return the sum of the values in the column. If the type of the column is inferred, the return type is double, the type of the column otherwise. |
| Take() | Extract a column from the dataset as a collection of values, e.g. a std::vector<float> for a column of type float. |

| Instant action | Description |
|---|---|
| Foreach() | Execute a user-defined function on each entry. Users are responsible for the thread-safety of this callable when executing with implicit multi-threading enabled. |
| ForeachSlot() | Same as Foreach(), but the user-defined function must take an extra unsigned int slot as its first parameter. slot will take a different value, 0 to nThreads - 1, for each thread of execution. This is meant as a helper in writing thread-safe Foreach() actions when using RDataFrame after ROOT::EnableImplicitMT(). ForeachSlot() works just as well with single-thread execution: in that case slot will always be 0. |
| Snapshot() | Write the processed dataset to disk, in a new TTree and TFile. Custom columns can be saved as well, filtered entries are not saved. Users can specify which columns to save (default is all). Snapshot, by default, overwrites the output file if it already exists. Snapshot() can be made lazy setting the appropriate flag in the snapshot options. |

**Queries**

These operations do not modify the dataframe or book computations but simply return information on the RDataFrame object.

| Operation | Description |
|---|---|
| Describe() | Get useful information describing the dataframe, e.g. columns and their types. |
| GetColumnNames() | Get the names of all the available columns of the dataset. |
| GetColumnType() | Return the type of a given column as a string. |
| GetColumnTypeNamesList() | Return the list of type names of columns in the dataset. |
| GetDefinedColumnNames() | Get the names of all the defined columns. |
| GetFilterNames() | Return the names of all filters in the computation graph. |
| GetNRuns() | Return the number of event loops run by this RDataFrame instance so far. |
| GetNSlots() | Return the number of processing slots that RDataFrame will use during the event loop (i.e. the concurrency level). |
| SaveGraph() | Store the computation graph of an RDataFrame in DOT format (graphviz) for easy inspection. See the relevant section for details. |

*from E. Guiraud: RDF@ICHEP 2022*

# Machine Learning in ROOT

- TMVA: Toolkit for Multi-Variate Analysis

  - Provides a collection of several ML algorithms including Decision Trees, but also Deep Learning tools with GPU support

- Shifted Focus now on :

  - interoperability with Python ML libraries

    - developing glue between ROOT I/O and Python ML ecosystem
    - designing also new interfaces in TMVA for better interoperability

  - optimal model inference

    - enable physicists to deploy easily their ML models on large datasets
    - integration of model inference with analysis tools (RDF)

# SOFIE: Fast ML Inference in ROOT

- New ROOT tool for inference code generation for DL models
  - input, a trained ML model:
    - ONNX (new standard for ML)
    - TF/Keras
    - PyTorch
  - Output:
    - C++ code
      - minimal dependency (BLAS)
      - can be compiled on the fly (e.g. with Cling)

**ROOT 6.26 (experimental)**

**Stored model**                    **C++ source**



```
using namespace TMVA::Experimental::SOFIE;
RModelParser_ONNX parser;
RModel model = parser.Parse("Model.onnx");
```

```
// generate text code internally (with some options)
model.Generate();
// write output header file and data weight file
model.OutputGenerated();
```

# ML inference with RDataFrame

- ML model evaluation can be integrated in RDF event loop

  - SofieFunctor:  adapter for using SOFIE within RDF

```cpp
auto h1 = df.DefineSlot("DNN_Value",
        SofieFunctor<7,TMVA_SOFIE_higgs_model_dense::Session>(nslots),
        {"m_jj", "m_jjj", "m_lv", "m_jlv","m_bb","m_wbb","m_wwbb"}).
        Histo1D("DNN_Value");
```

  - support for multi-threaded evaluation
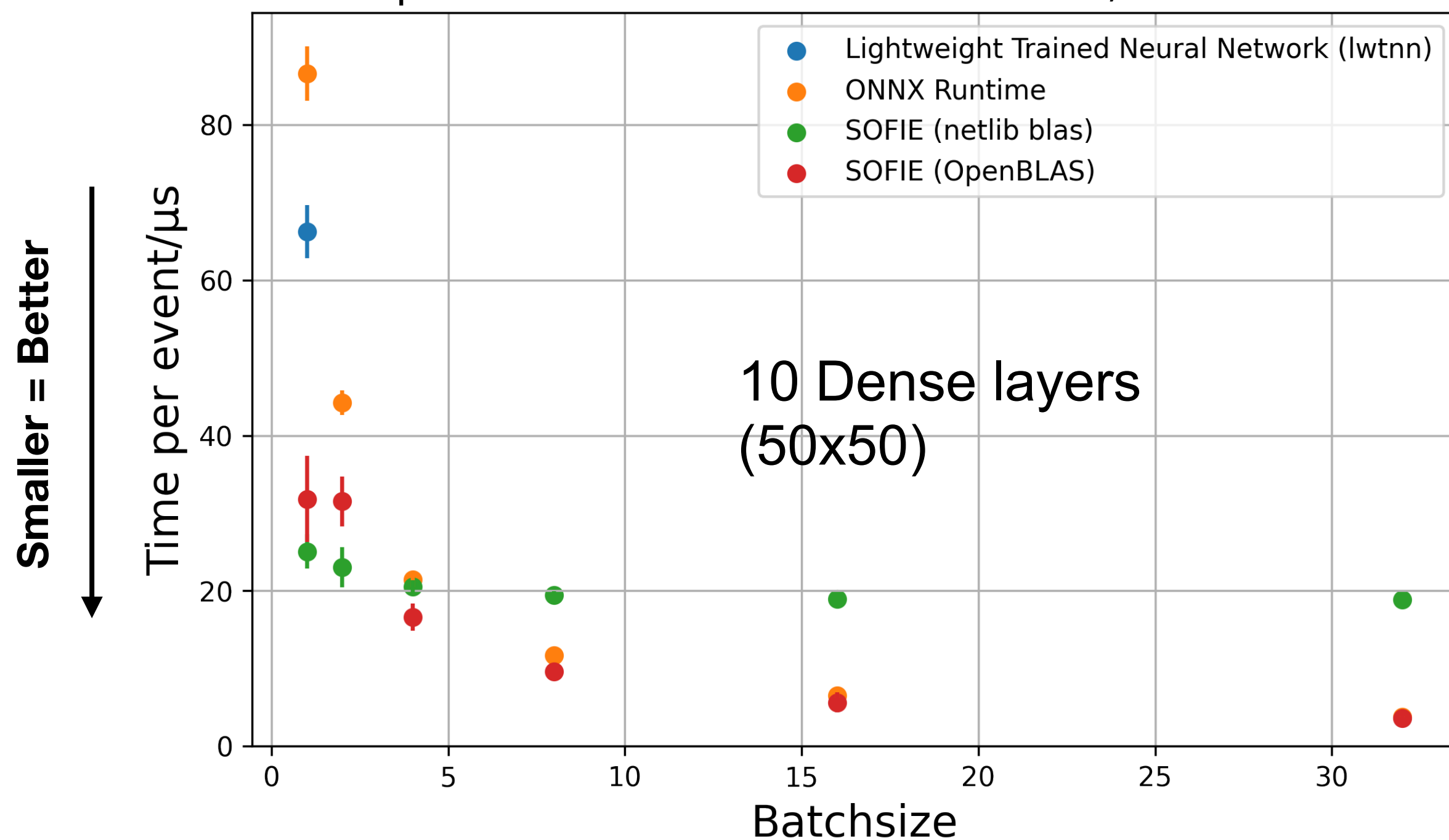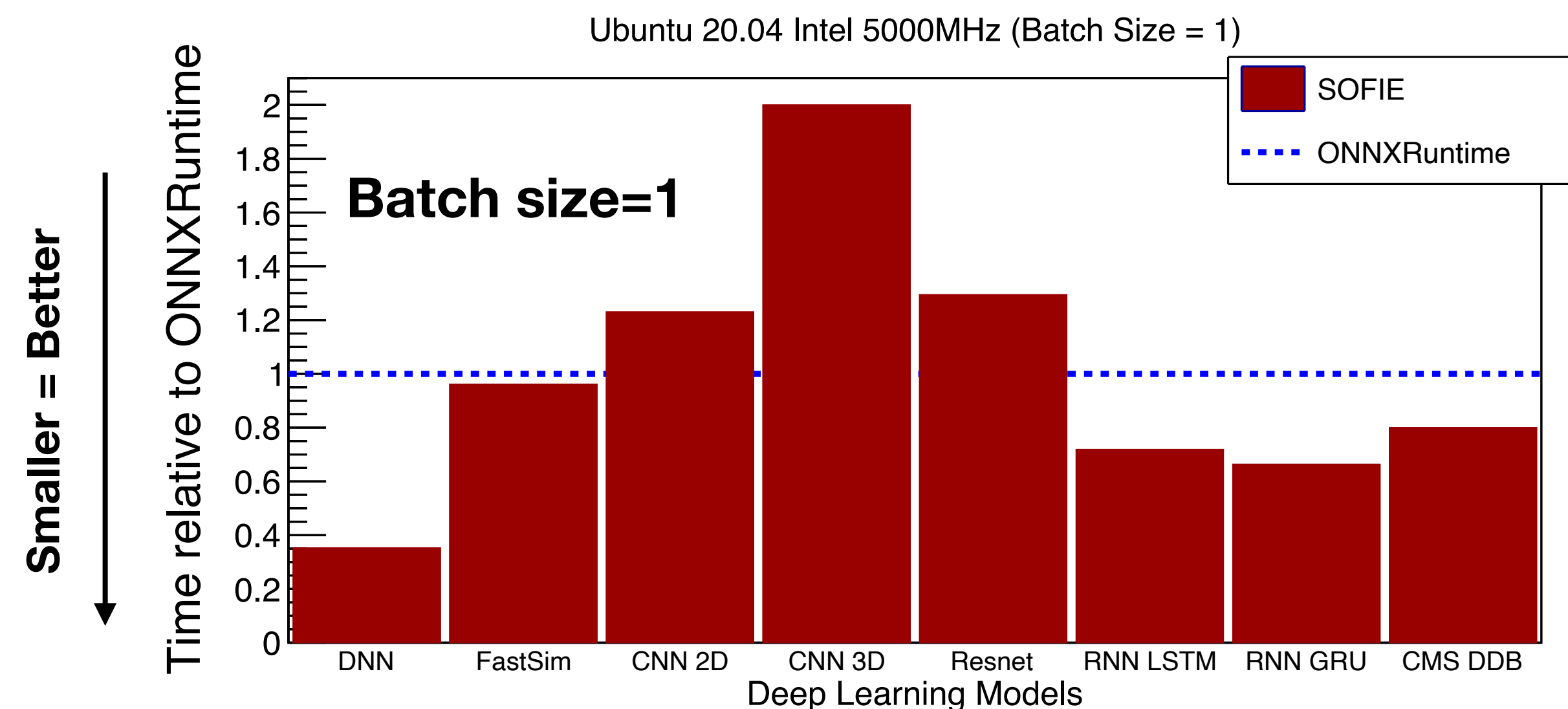
  - see full Example tutorial code in C++ or Python
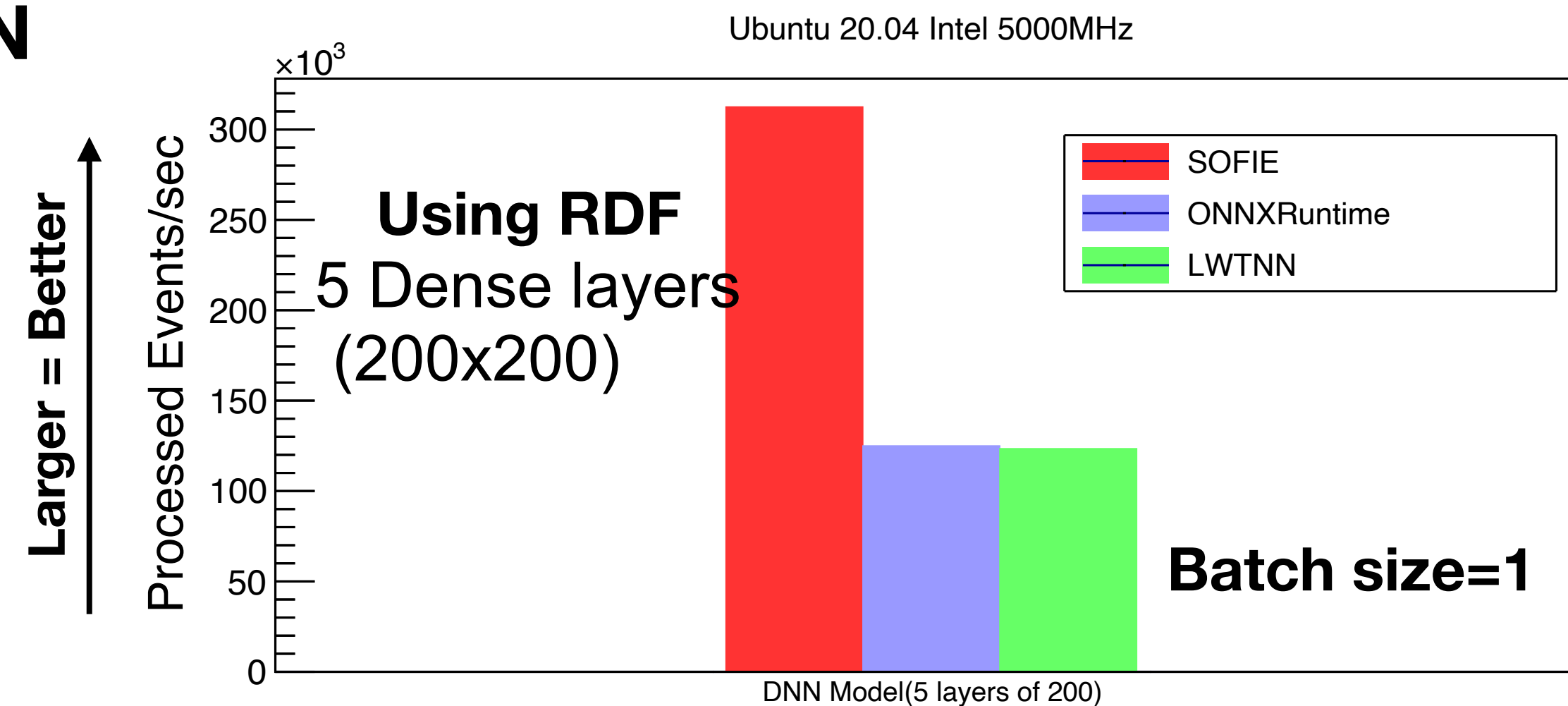
# SOFIE Benchmarks

## Comparison with ONNXRuntime and LWTNN

- 2-3 faster than ONNXRuntime for DNN with batch size=1
- 20% faster for RNN operators
- slower for CNN ( 20% for 2D, x2 for 3D)

Time per event for different batch size, cache flushed



10 Dense layers (50x50)

**varying batch size…**



Using RDF 5 Dense layers (200x200)

Batch size=1



Batch size=1

**using different models (DNN, CNN, RNN)**

# Supported Operators in SOFIE

| | |
|---|---|
| **Perceptron:** Gemm | Implemented and integrated (ROOT 6.26) |
| **Activations:** Relu, Seul, Sigmoid, Softmax, LeakyRelu | Implemented and  integrated |
| **Convolution** (1D, 2D and 3D and ConvTranspose | Implemented and integrated |
| **Recurrent:** RNN, GRU, LSTM | Implemented and integrated |
| BatchNormalization | Implemented and integrated |
| **Pooling:** MaxPool, AveragePool, GlobalAverage | Implemented and integrated |
| **Layer operations:** Add, Sum, Mul, Div, Reshape, Flatten,  Transpose, Squeeze, Unsqueeze, Slice, Concat, Identity | Implemented and integrated |
| User-defined operator | Implemented and integrated |
| **Gather (for embedding) and more depending on user needs** | Planned for next release |
| Direct GNN support (NEW) | Planned for next release |

# Fast Decision Tree Inference

- Inference engine taking model parameters from externally trained models
- Features:
  - Simple to use from Python and C++
  - Thread-safe
  - Zero-copy
  - Fast for single event and batch inference

- External training and model conversion

```
xgb = xgboost.BDTClassifier(options)
xgb.fit(x, y)

ROOT.TMVA.SaveXGBoost(xgb, "myBDT", "model.root")
```
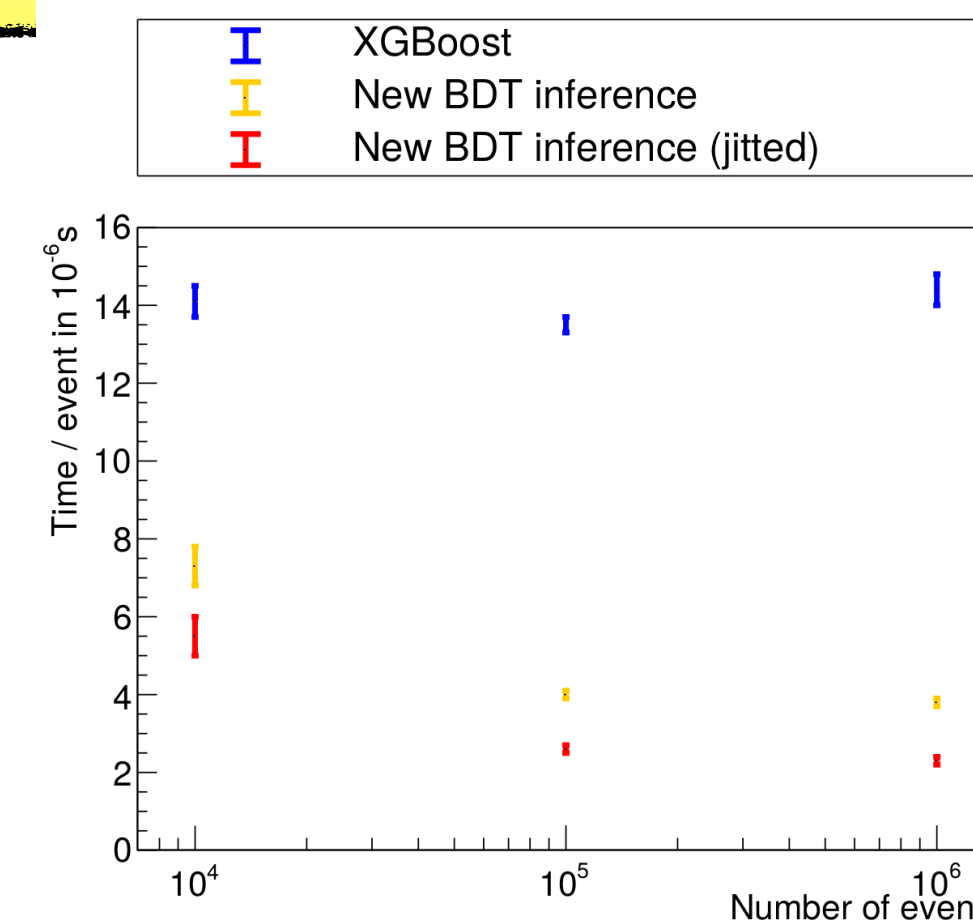
- Python application

```
xgb = xgboost.BDTClassifier(options)
xgb.fit(x, y)

ROOT.TMVA.SaveXGBoost(xgb, "myBDT", "model.root")
```

- C++ application

```
TMVA::RBDT bdt("myBDT", "model.root");
auto y1 = bdt.Compute({1.0, ...});

auto x = TMVA::RTensor<float>(data, shape);
auto y2 = bdt.Compute(x);
```
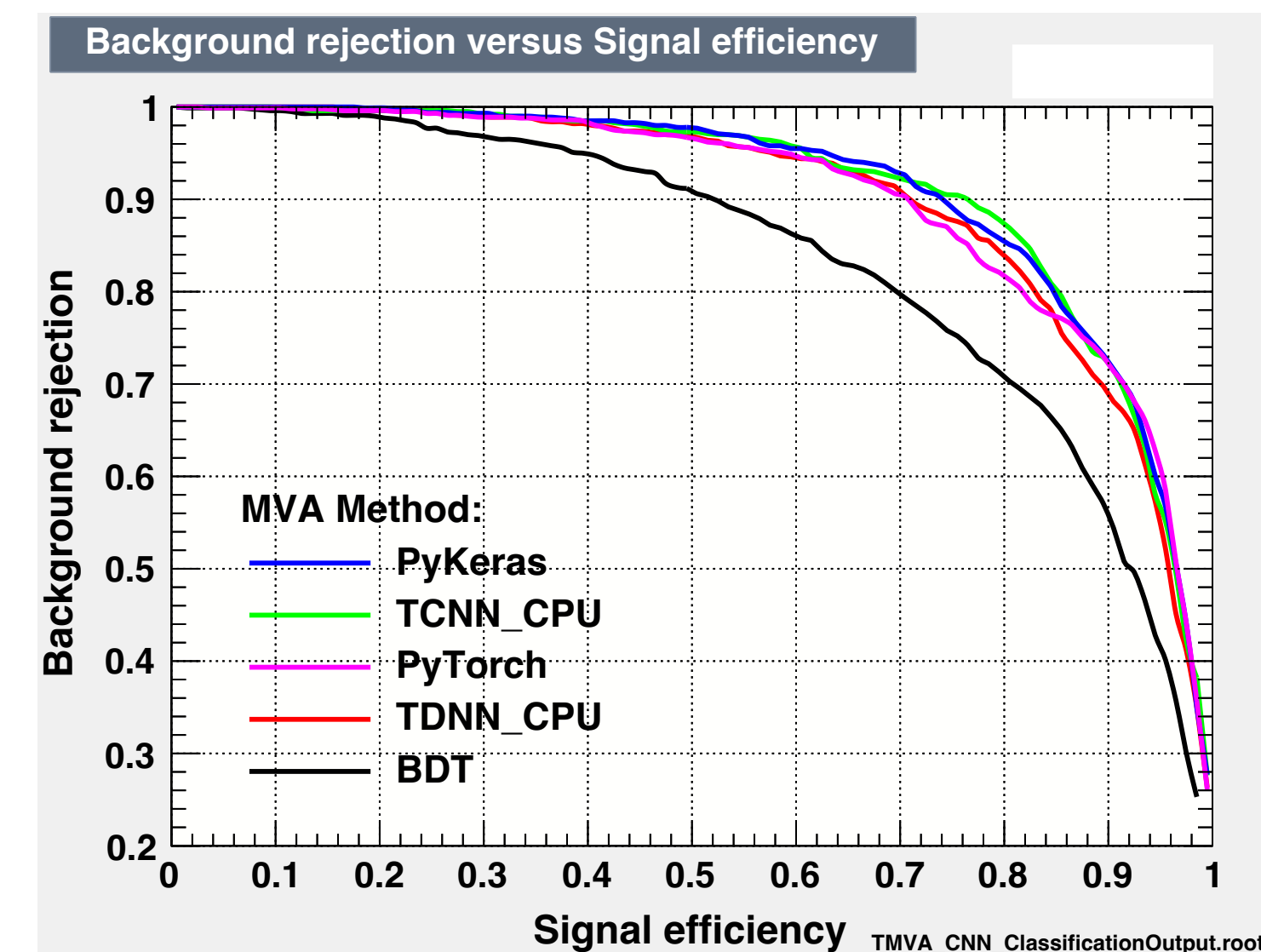
# Interoperability of TMVA

- Pythonization of TMVA interfaces

  - from string API to python keyword args

- Training of Python ML tools within TMVA workflow

  - interfaces to scikit-learn, Keras and PyTorch

- Working on generic data-loader for ML workflows

  - Generator doing batching and shuffling from ROOT files on the fly

  - Allows for efficient training of datasets larger than RAM sizes

  - Direct feeding of data from disk to GPU



Background rejection versus Signal efficiency

MVA Method:
- PyKeras
- TCNN_CPU
- PyTorch
- TDNN_CPU
- BDT

TMVA_CNN_ClassificationOutput.root

**Example of a possible ML workflow loading batches**
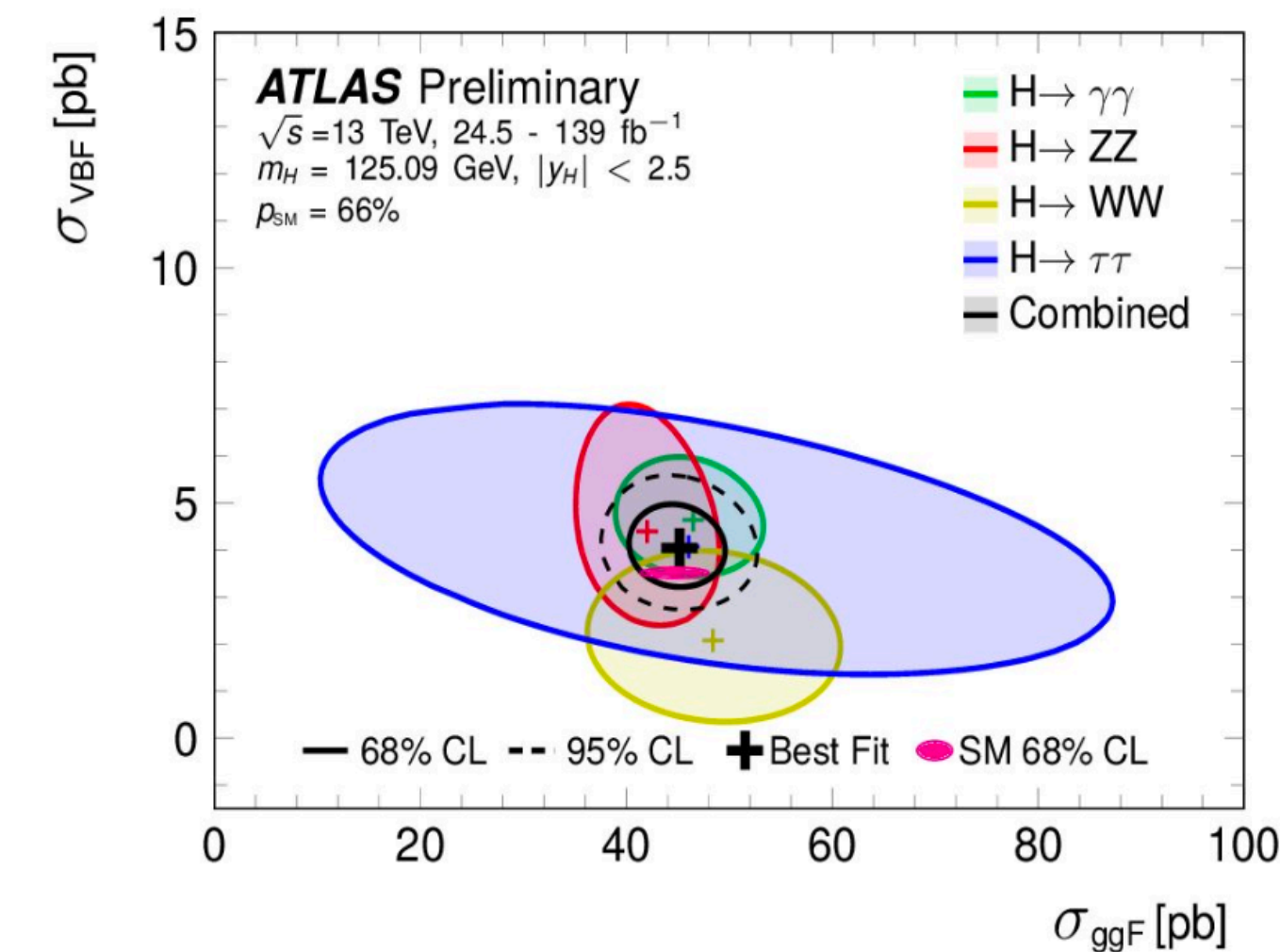
```
df = ROOT.RDataFrame("Events", "http://file.root")
generator = TMVA.BatchGenerator(df, cols, batchSize)
for step in gradientSteps:
    x = generator()
    model.fit(x)
```

# Statistical Modeling

- Physics analyses continuously generate increasingly complex likelihood models to describe their data
  - Higgs combination fits, EFT interpretations
    - O(1000) parameters
    - O(100) likelihood components
    - O(100) datasets



Recent ATLAS Higgs combination fit result [1]. This fit took about five hours to complete.

- Only one tool capable of handling such models: RooFit
- Recent development give possibility to bring down fitting time from hours to minutes
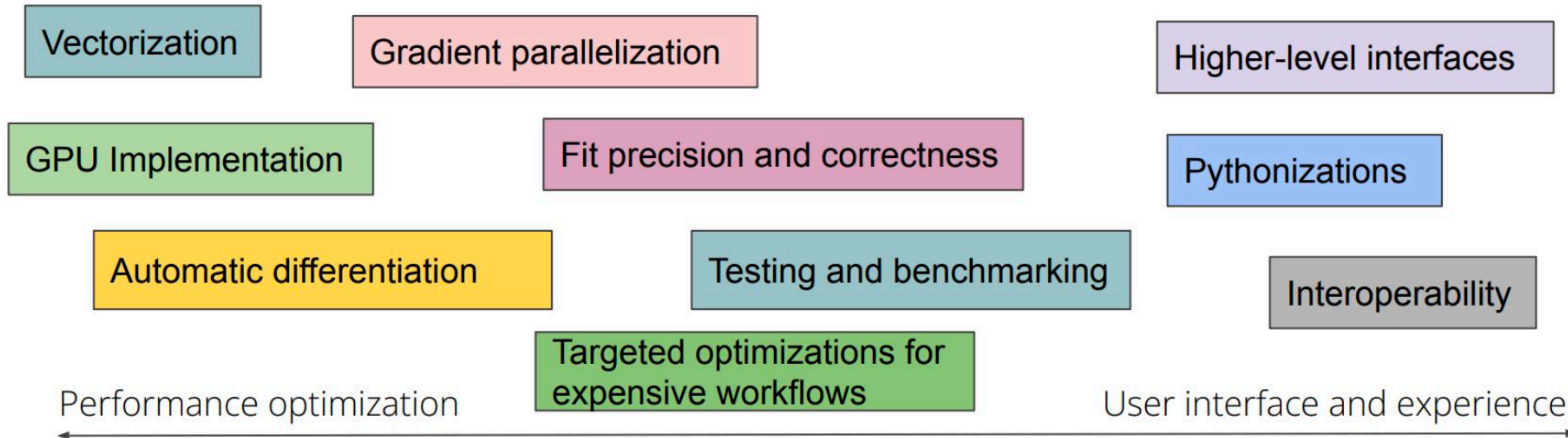  - from a work day to a coffe break!

# RooFit Evolution

RooFit is evolving on many different areas:



Vectorization

Gradient parallelization

Higher-level interfaces

GPU Implementation

Fit precision and correctness

Pythonizations

Automatic differentiation

Testing and benchmarking

Interoperability

Targeted optimizations for expensive workflows

Performance optimization → User interface and experience

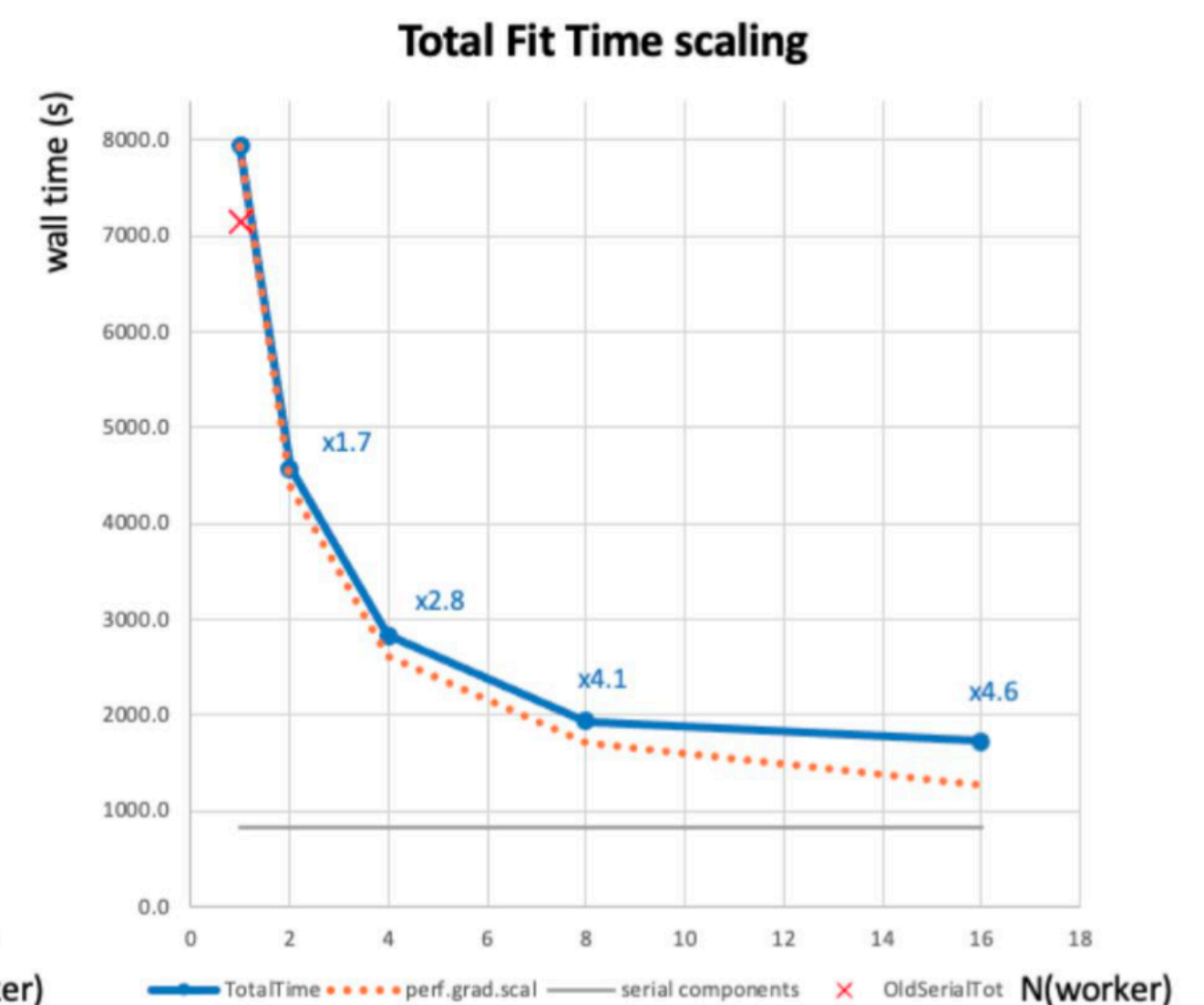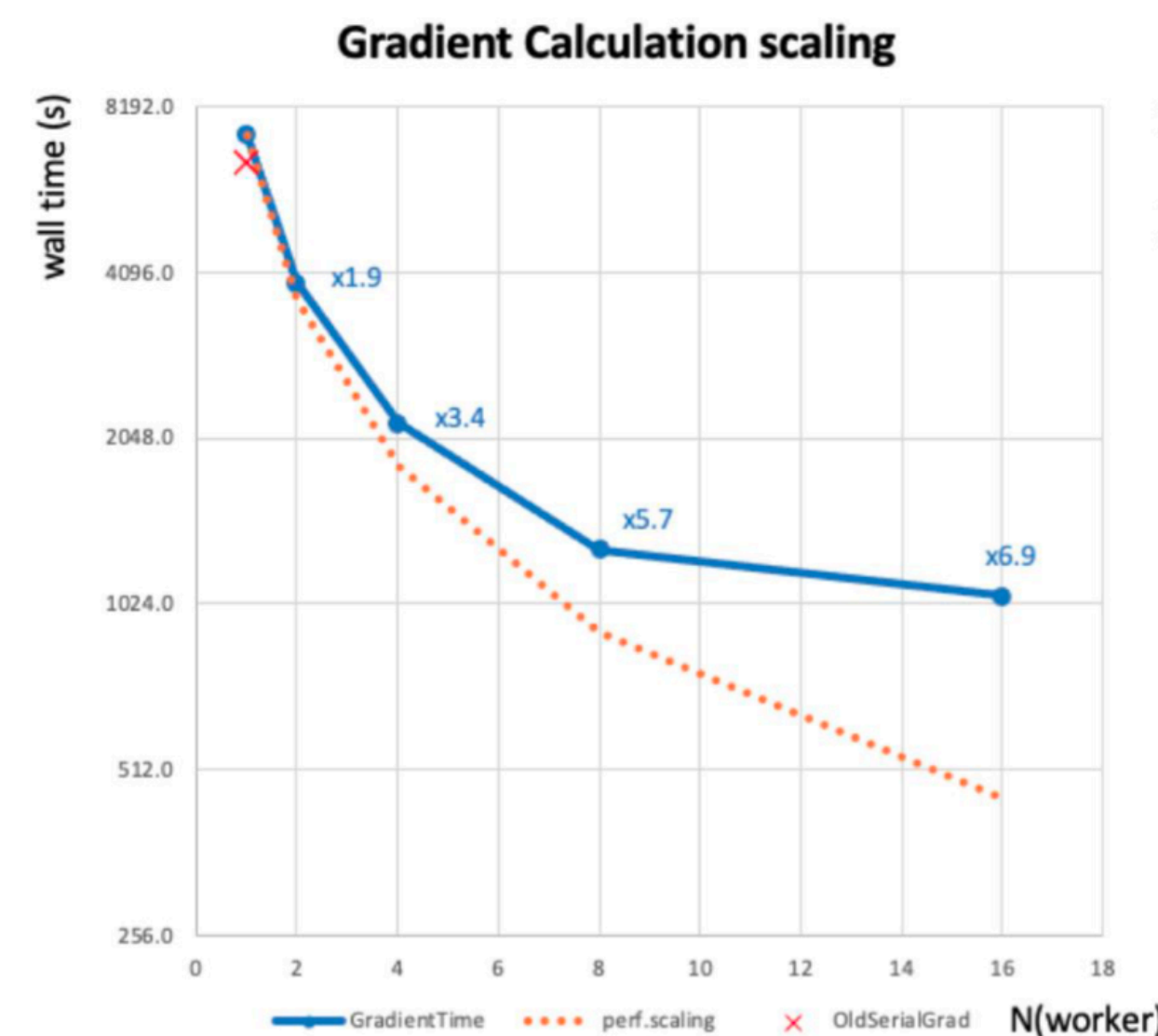# Gradient parallelization



- Parallelize at gradient calculation level
  - N parameters: ~O(2N) function evaluations for computing numerical derivatives
  - Line search: serial part ~O(3) function evaluation
- Dynamic load balancing over workers through random work stealing algorithm
  - complexity of derivative calculation varies
- Designed to have maximum speed impact of complex fits with many parameters

- Wall time decrease in Higgs combination fit:
  - from 2h12m26s → 28m52s (~4/5 times faster)
- Result validated: all parameters agreed with serial fit within 1% of their uncertainties !
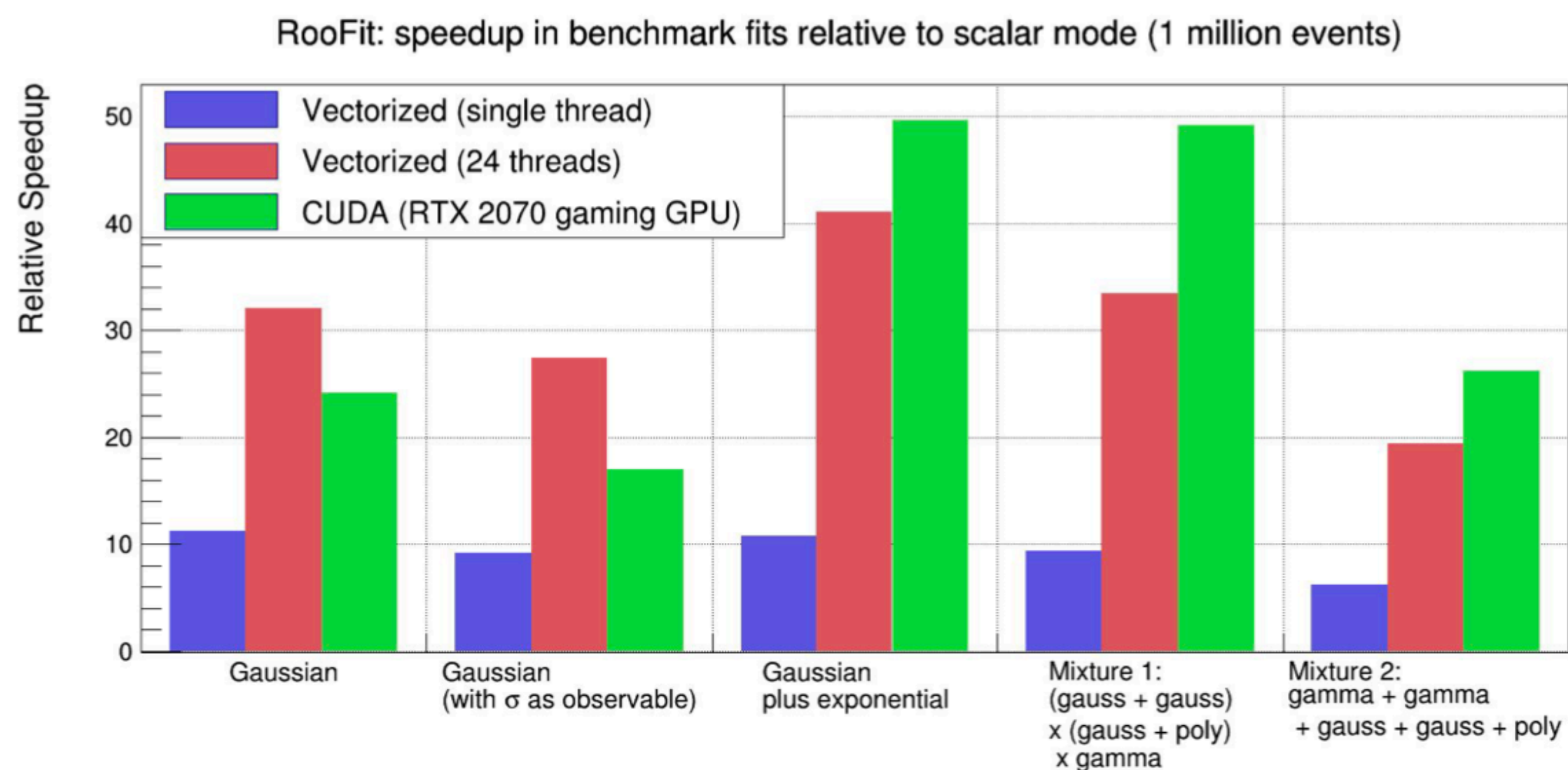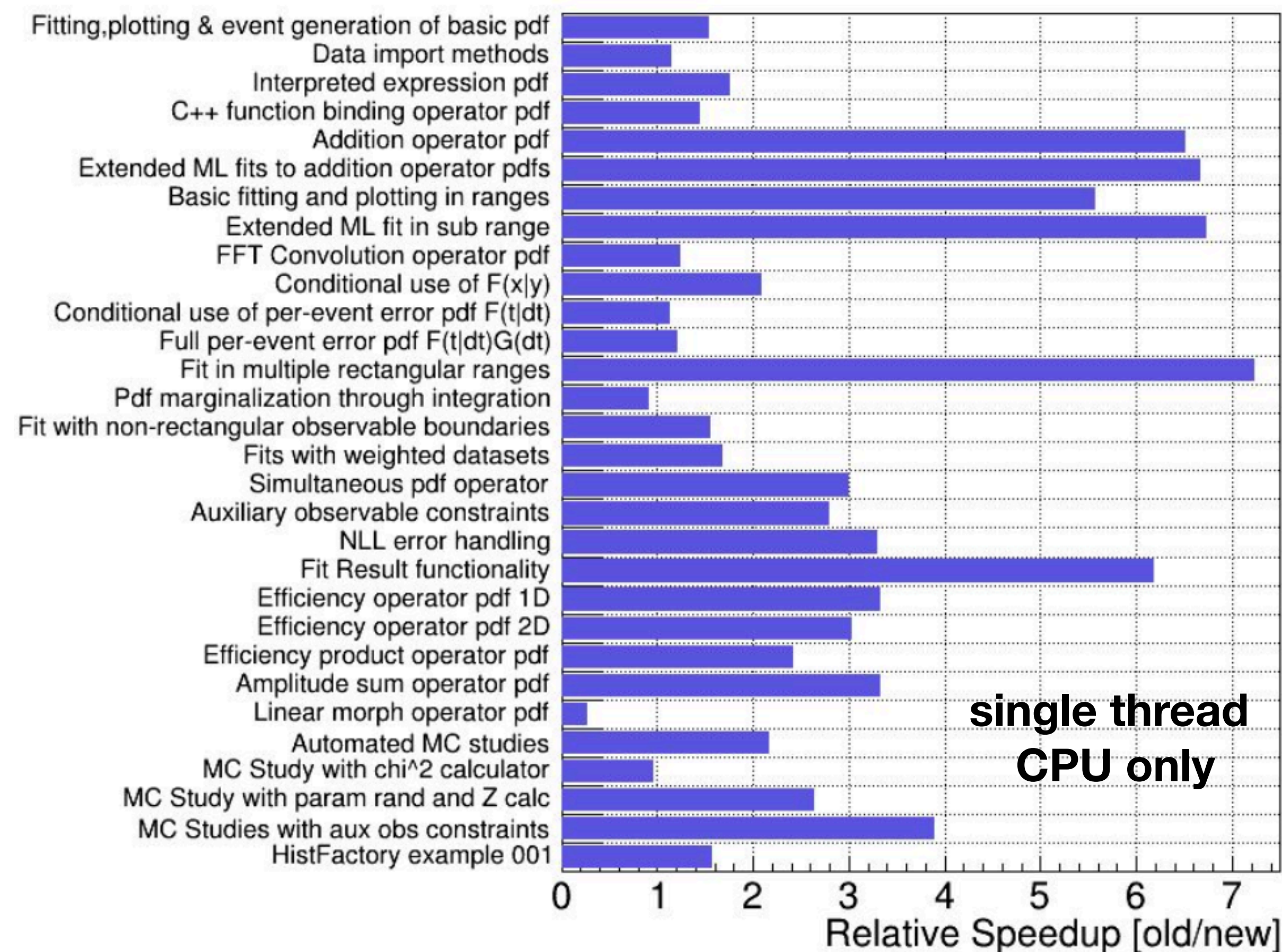
# Vectorisation and GPU Computation

- Batch evaluation of RooFit models:
  - allow code vectorisation
  - optionally allow multi-threading or GPU computation (with CUDA)
- possible due to a re-structure of RooFit computational graph.
  - from event-by event to batch node evaluation



RooFit: speedup in benchmark fits relative to scalar mode (1 million events)

RooFit/HistFactory stress tests: speedup of NLL minimization by using BatchMode

single thread CPU only

➡ Additional speed-up obtained in Batch mode by better CPU caching and vectorisation !

# RooFit Pythonizations

- RooFit is more pythonic in latest ROOT (6.26):
  - use Python keyword arguments instead of RooFit command arguments
  - pass around Python sets or lists or dictionaries instead of RooFit collections
  - converters from/to Numpy/Pandas and RooFit data classes
- All Pythonizations are documented in the reference guide

Example code from a RooFit tutorial

```python
# Create background pdf poly(x)*poly(y)*poly(z)
px = ROOT.RooPolynomial("px", "px", x, [-0.1, 0.004])
py = ROOT.RooPolynomial("py", "py", y, [0.1, -0.004])
pz = ROOT.RooPolynomial("pz", "pz", z)
bkg = ROOT.RooProdPdf("bkg", "bkg", [px, py, pz])

# Create composite pdf sig+bkg
fsig = ROOT.RooRealVar("fsig", "signal fraction",
                        0.1, 0., 1.)
model = ROOT.RooAddPdf("model", "model",
                        [sig, bkg], [fsig])

data = model.generate((x, y, z), 20000)

# Make plain projection of data and pdf on x observable
frame = x.frame(Title="Projection on X", Bins=40)
data.plotOn(frame)
```
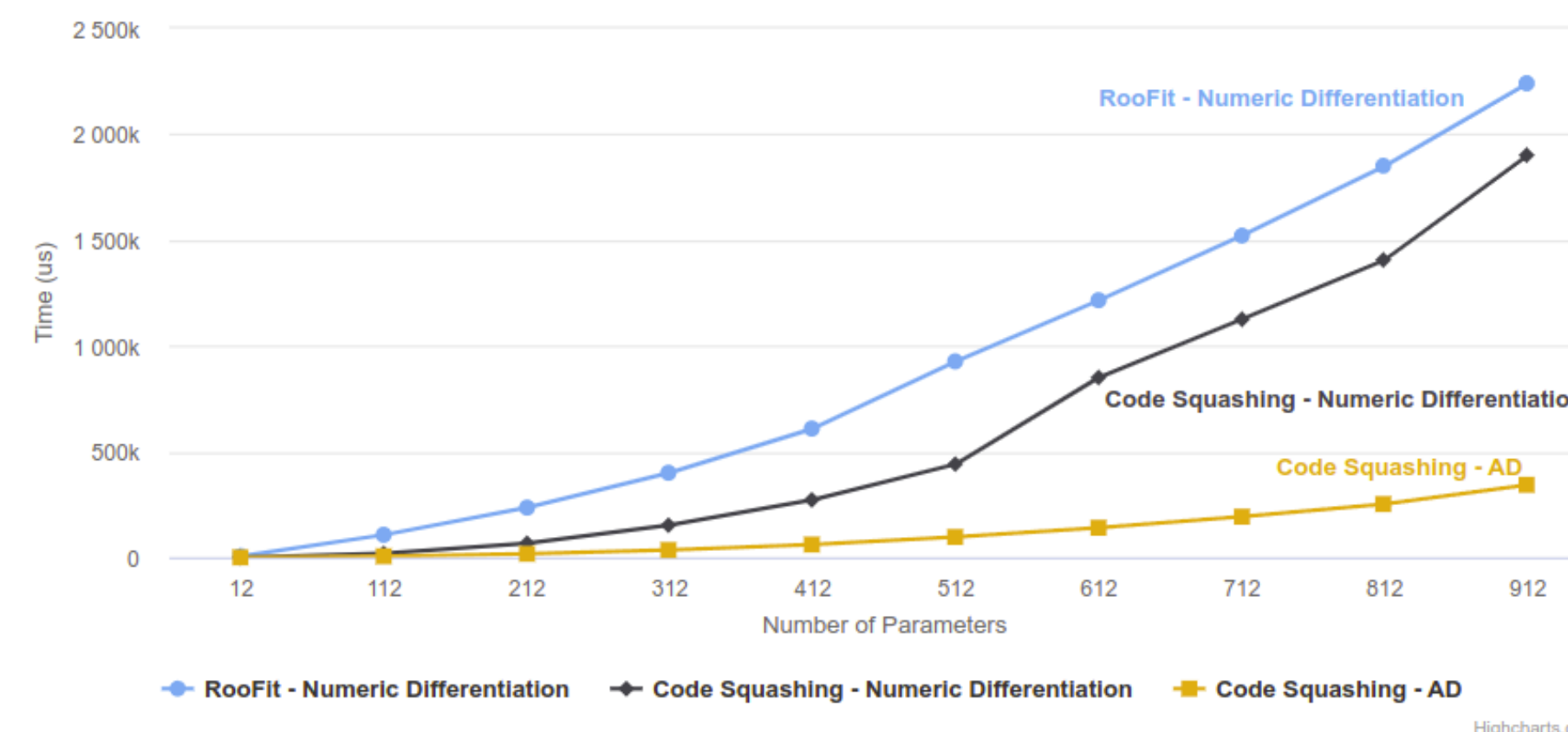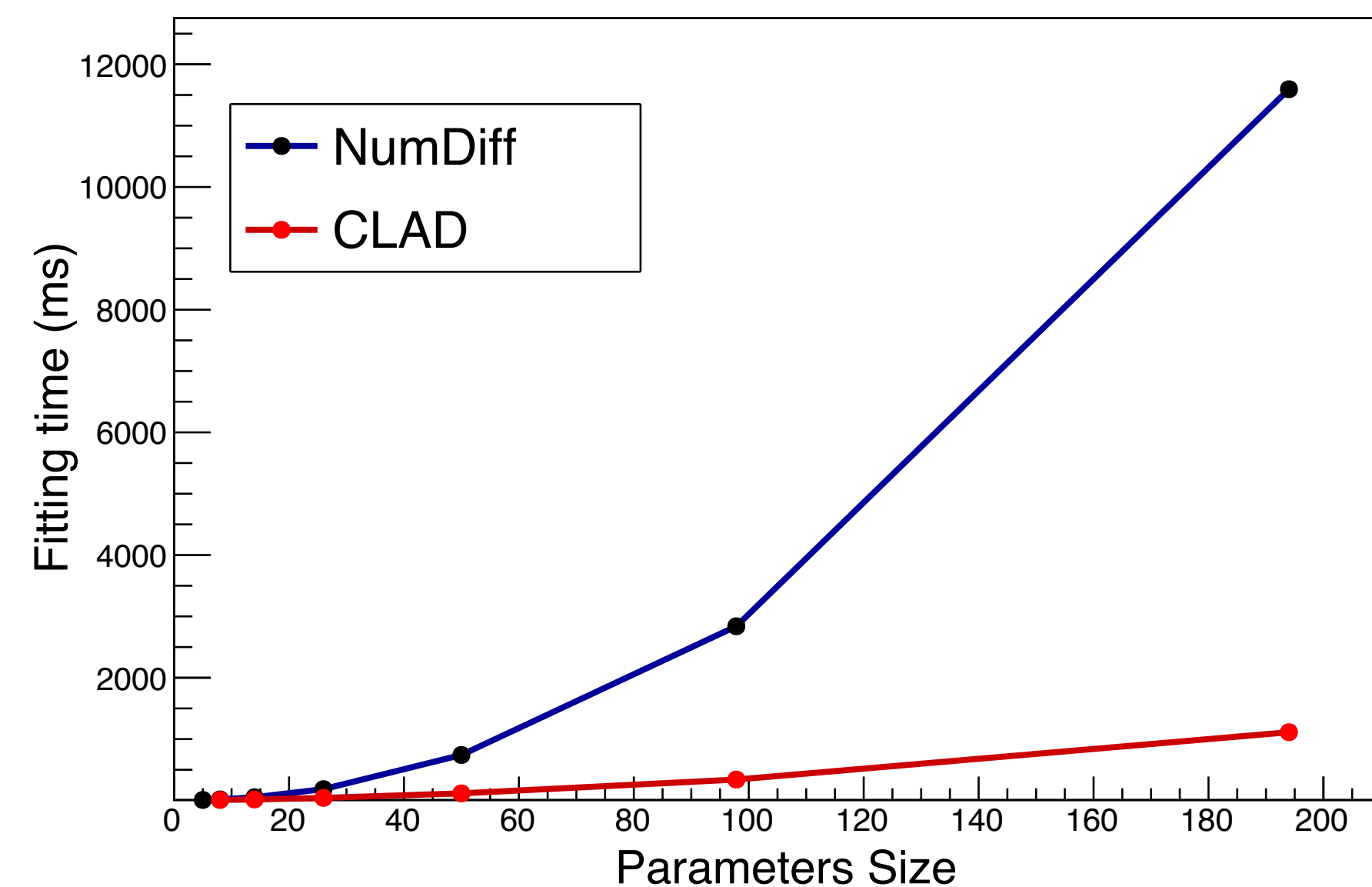
# Auto-differentiation (AD)

- AD available in ROOT **TFormula** class for fitting

- Large speedup can be obtained with respect to using numerical differentiation in case of large number of parameters
  - ~ 10x speedup for 100 parameters fit

- Also AD for Hessian (covariance matrix)
  - working in adding support in Minuit

- Integrating AD in RooFit
  - started with HistFactory models (binned likelihood fits, many parameters, few integrals) : ~ 5x speedup

*from G. Singh, AD in ROOT @2nd MODE Workshop, Sep 2022*  34

# Summary

- Future analyses in HEP and Nuclear Physics will use larger data sets and be more complex (e.g. larger parameter spaces)
- Analysis tools for physicists will need to be
  - efficient and performant
  - with simple interfaces
  - robust, reliable and supported
  - allow interoperability (e.g. with Python ecosystem)
- ROOT is rapidly evolving in all its components to fulfil these requirements
  - thanks to a large effort from all the team and user contributions
- Need your help for continuous feedback and support

# References

- [root.cern](root.cern)

- [https://cern.ch/forum](https://cern.ch/forum)
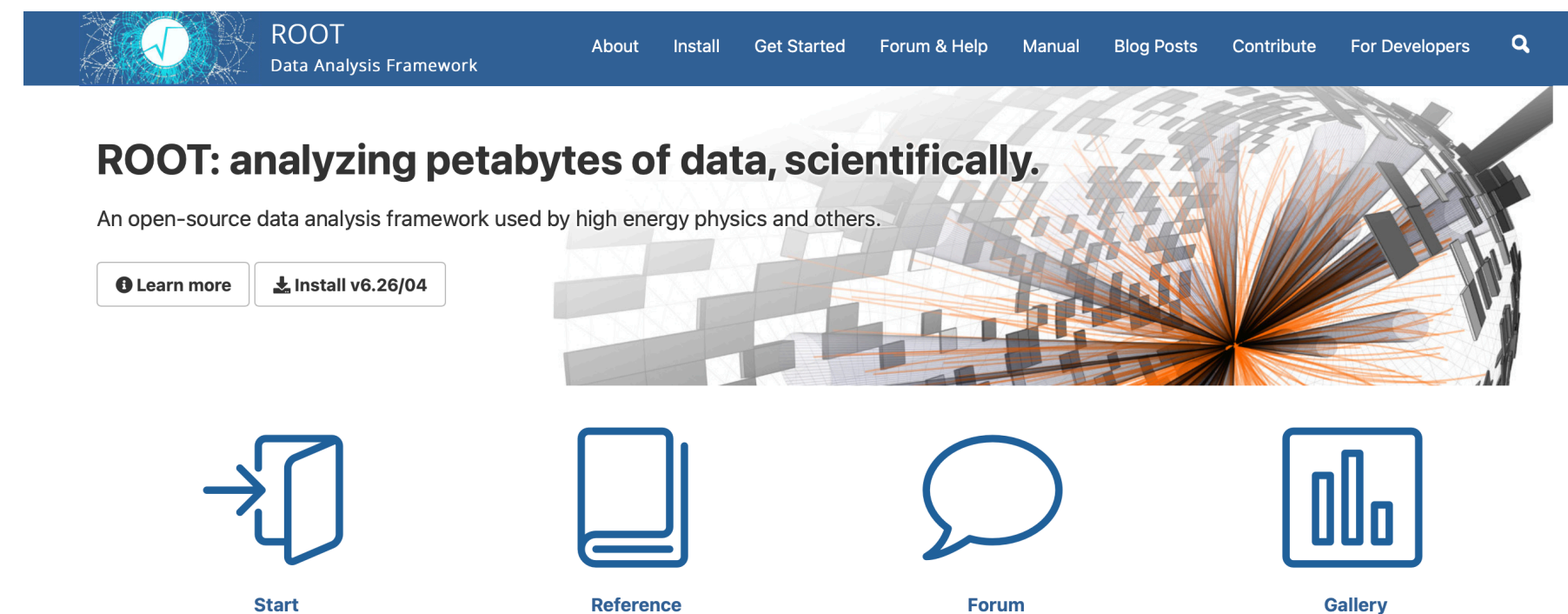
- [https://github.com/root-project](https://github.com/root-project)


- For more information of current developments see presentations at latest ROOT Users Workshop

- see also the ICHEP 22 presentations in the Computing session

# Backup Slides

# RDF Benchmarks

Fully compiled C++ RDataFrame

Coffea 0.7.12 (using chunksize=2**19)

| query, | 1x data (s), | 10x data (s) | query, | 1x data (s), | 10x data (s) |
|--------|--------------|--------------|--------|--------------|--------------|
| Q1 | 0.37 | 1.50 | Q1 | 1.40 | 4.24 |
| Q2 | 0.46 | 3.70 | Q2 | 1.51 | 5.76 |
| Q3 | 0.73 | 6.23 | Q3 | 1.81 | 7.96 |
| Q4 | 0.65 | 5.92 | Q4 | 1.65 | 6.58 |
| Q5 | 0.84 | 7.45 | Q5 | 2.41 | 12.43 |
| Q6 | 3.08 | 27.99 | Q6 | 13.89 | 124.59 |
| Q7 | 2.56 | 22.27 | Q7 | 4.19 | 29.12 |
| Q8 | 1.17 | 10.22 | Q8 | 3.27 | 17.70 |

- note that these benchmarks are not representative of large analysis workloads
- see also this ACAT talk by Nick Smith

Benchmark from github.com/nsmith-/coffea-benchmarks
Setup: AMD EPYC 7702P, using 48 physical cores, data read from filesystem cache

# Coming Soon in RDF

- Performance improvements (e.g. bulk processing, ROOT PoW 2022)

- Collection aggregations (muon_{pt,eta,phi} → muons) being discussed

- Simpler Pythonic interfaces (less C++ strings in Python code), PoW 2022

- Allow default values for missing branches, PoW 2022, GitHub issue

# RooFit Interoperability with NumPy

- New **converters** between **NumPy** arrays/**Pandas** dataframes and **RooDataSet/RooDataHist**:

  - `RooDataSet.from_numpy()`
  - `RooDataSet.to_numpy()`
  - `RooDataSet.from_pandas()`
  - `RooDataSet.to_pandas()`
  - `RooDataHist.from_numpy()`
  - `RooDataHist.to_numpy()`

- New **RooRealVar.bins()** function to get RooFit bin boundaries as NumPy array

```python
from ROOT import RooRealVar, RooCategory, RooGaussian

x = RooRealVar("x", "x", 0, 10)
cat = RooCategory("cat", "cat",
                  {"minus": -1, "plus": +1})
mean = RooRealVar("mean", "mean",
                  5, 0, 10)
sigma = RooRealVar("sigma", "sigma",
                   2, 0.1, 10)
gauss = RooGaussian("gauss", "gauss",
                    x, mean, sigma)


data = gauss.generate((x, cat), 100)
df = data.to_pandas()
```

|     | x        | cat |
| --- | -------- | --- |
| 0   | 6.997865 | -1  |
| 1   | 7.211196 | -1  |
| 2   | 3.198248 | 1   |
| 3   | 5.015824 | 1   |
| 4   | 7.782388 | 1   |
| ... | ...      | ... |
| 95  | 6.878027 | -1  |
| 96  | 0.475900 | 1   |
| 97  | 4.451101 | -1  |
| 98  | 3.481015 | -1  |
| 99  | 4.010105 | -1  |

100 rows × 2 columns