

# Future of User Analysis

Jim Pivarski

Princeton University – IRIS-HEP

October 1, 2020



# ~~Future of~~ Trends in User Analysis

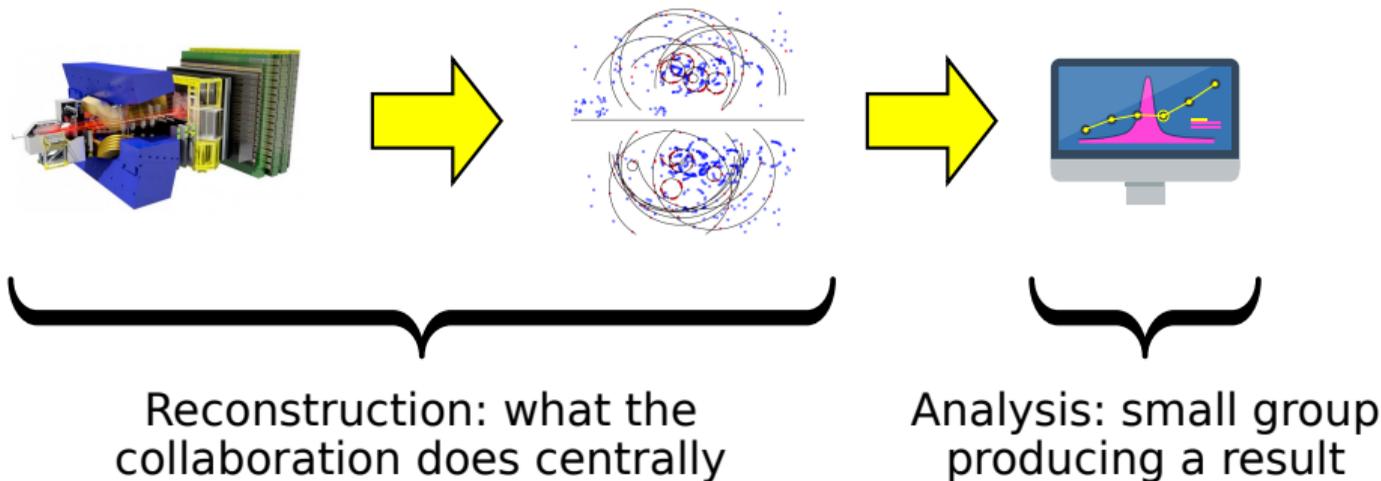
Jim Pivarski

Princeton University – IRIS-HEP

October 1, 2020



I'm aware that LHCb is moving important physics decisions earlier in the pipeline, but I think it's still meaningful to talk about “user analysis” as the last step.





*has happened*

*is happening*

*will happen*



Languages

Data formats

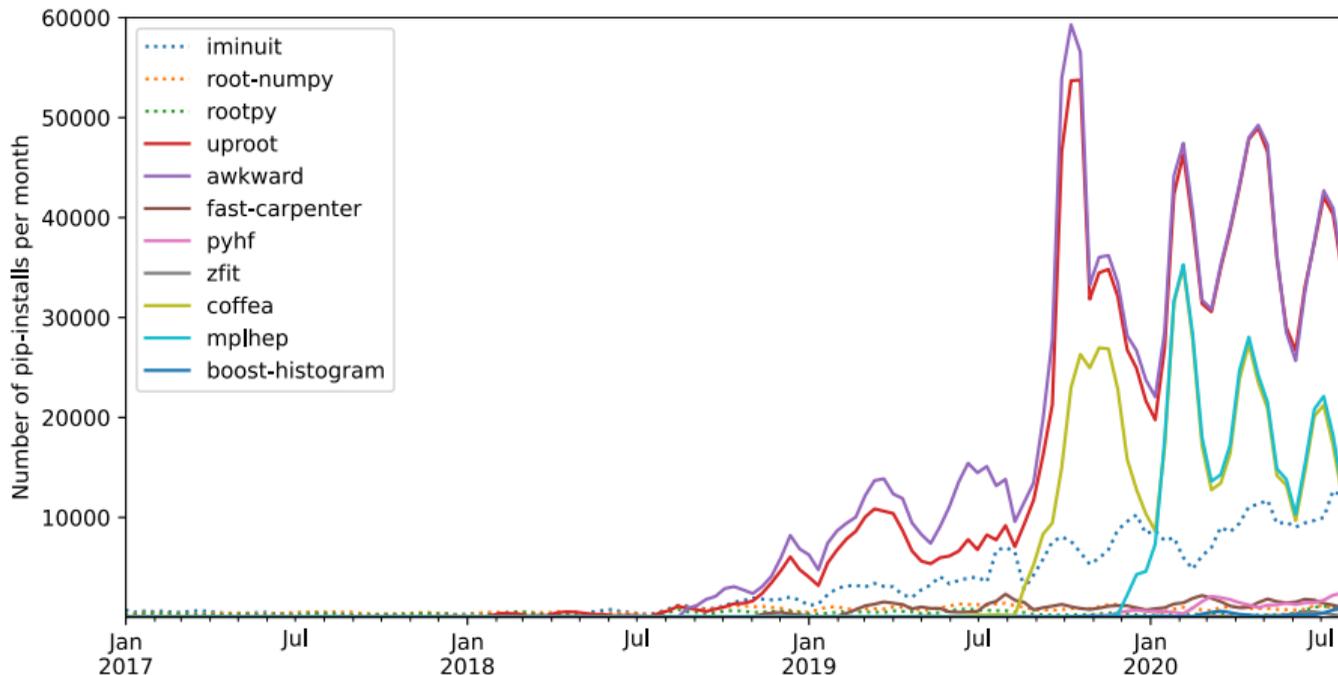
Libraries

Services



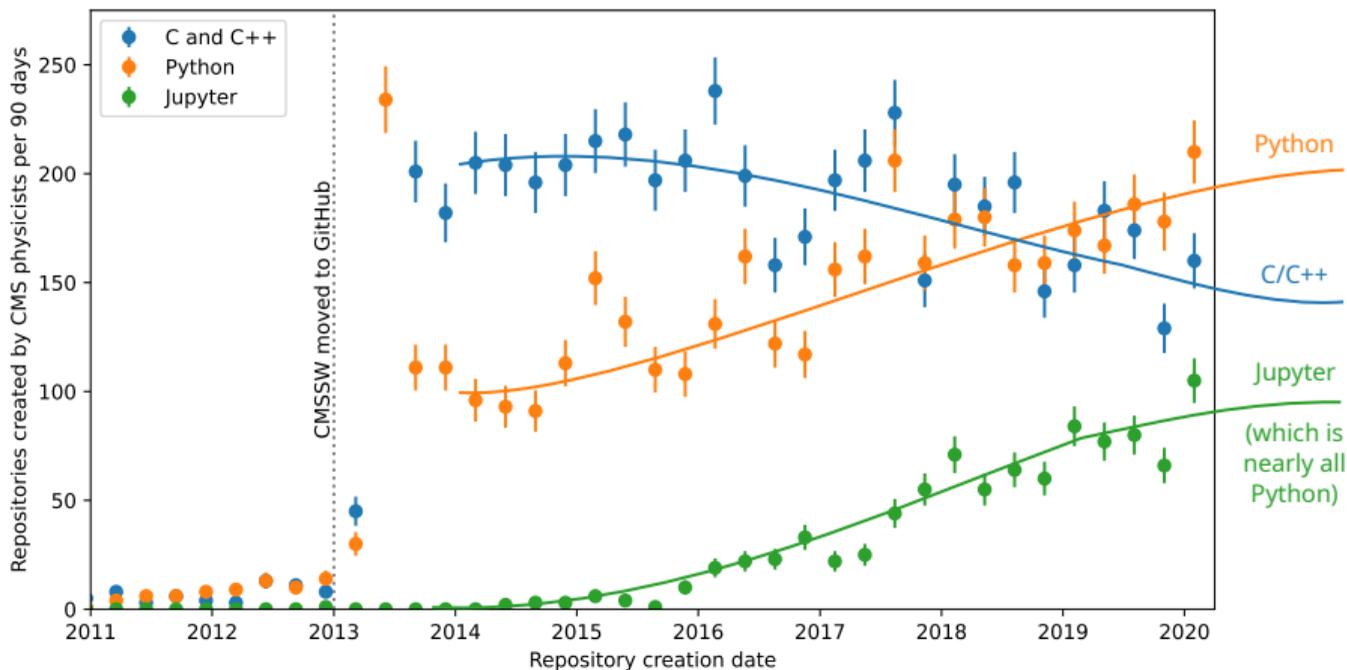
## Plot #1: pip-install statistics in the past three years

iminuit, root-numpy, and rootpy predate the rise, and therefore set the scale.





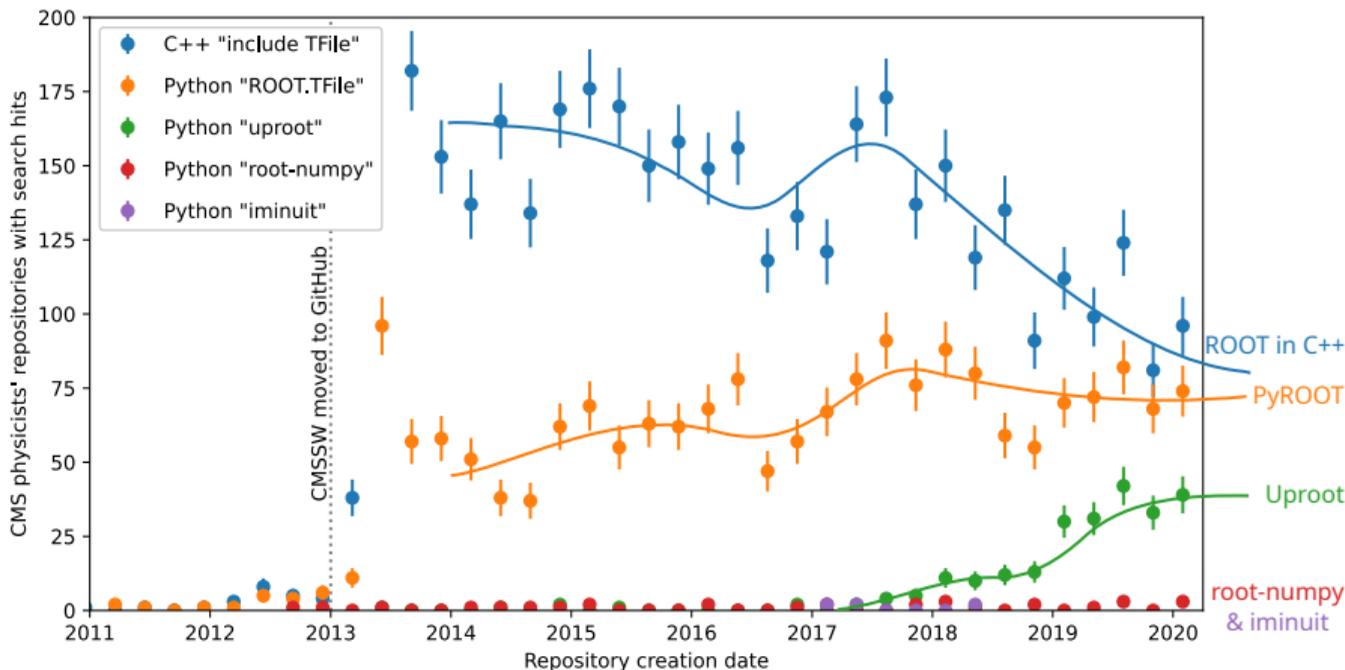
Plot #2: Language of non-fork GitHub repos by users who forked CMSSW





## Plot #3: Search for substrings in those repos: what are those users typing?

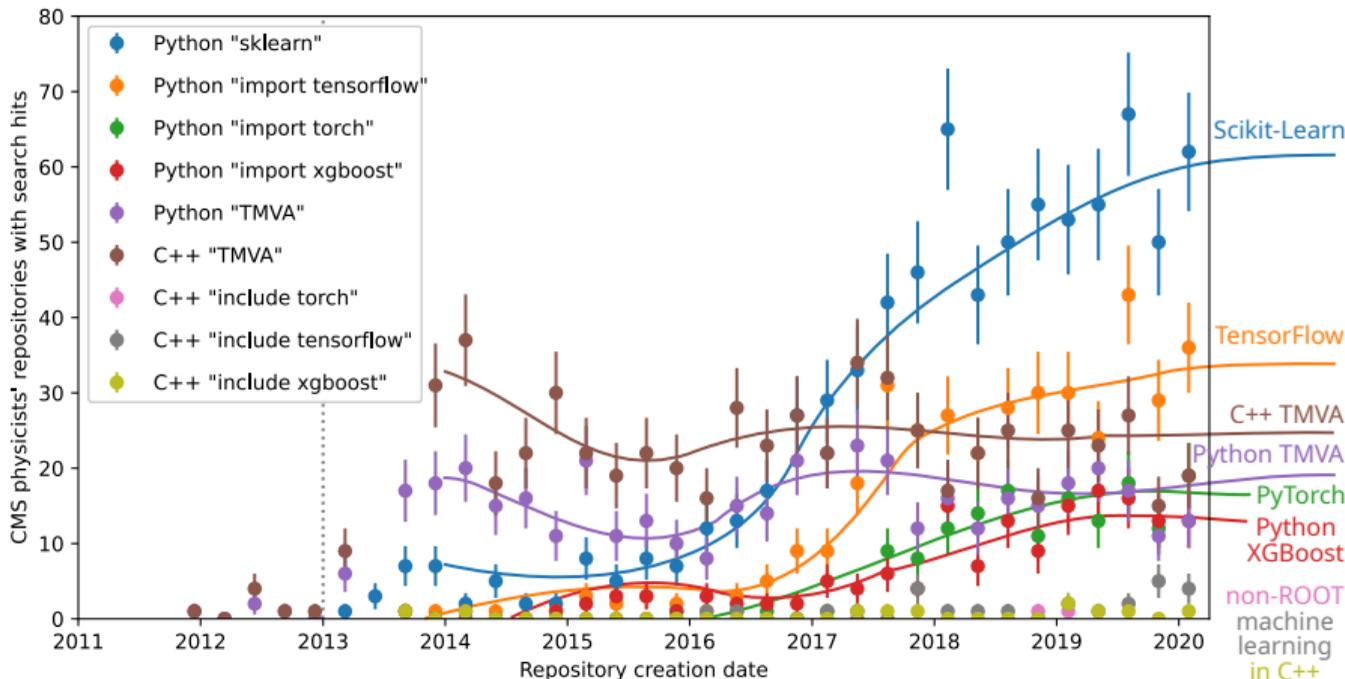
We assume "tfile" correlates to ROOT; "root" by itself (case insensitive) is too generic.





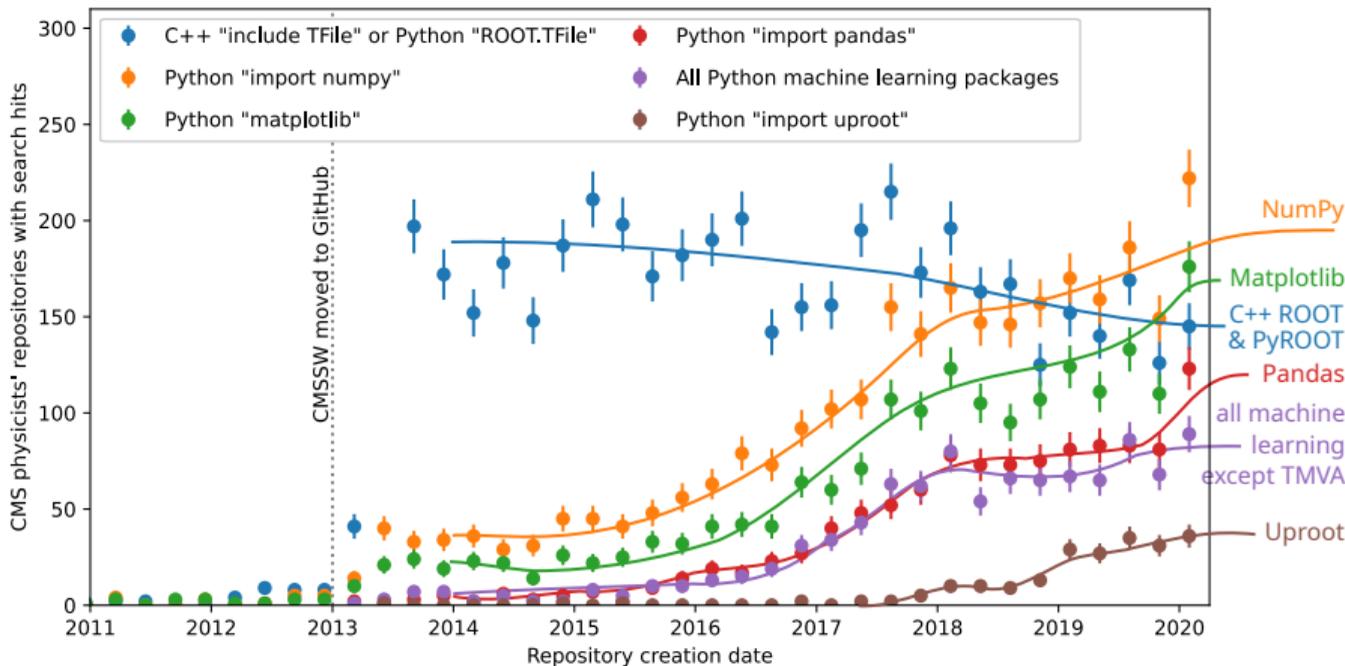
## Plot #3: What about machine learning?

**Surprise:** Scikit-Learn! **Not a surprise:** TMVA is the only significant C++ library.

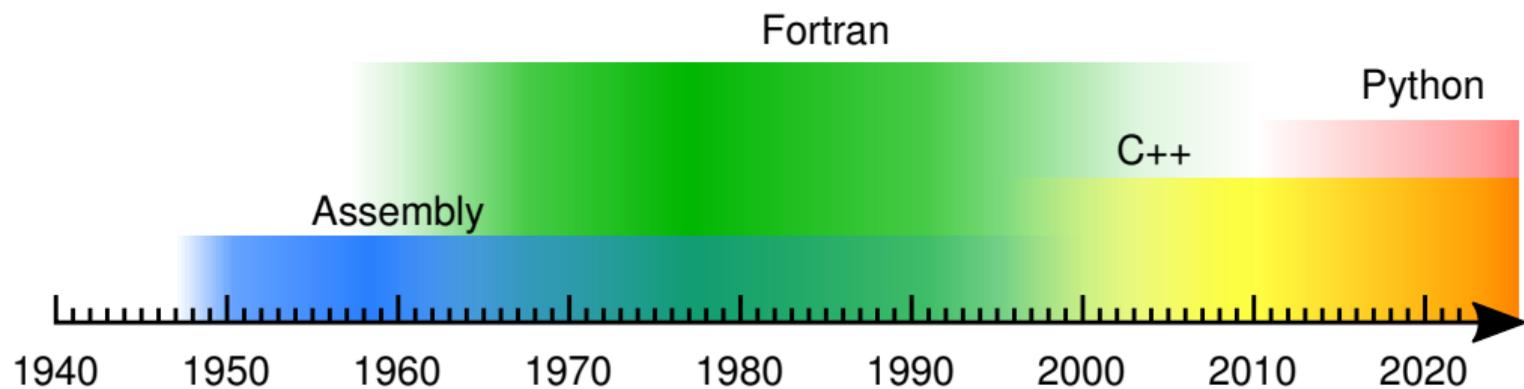




Plot #4: Was Python adoption driven by Uproot or machine learning (ML)?  
Twice as much basic analysis (NumPy/Matplotlib) than ML; trend predates Uproot.

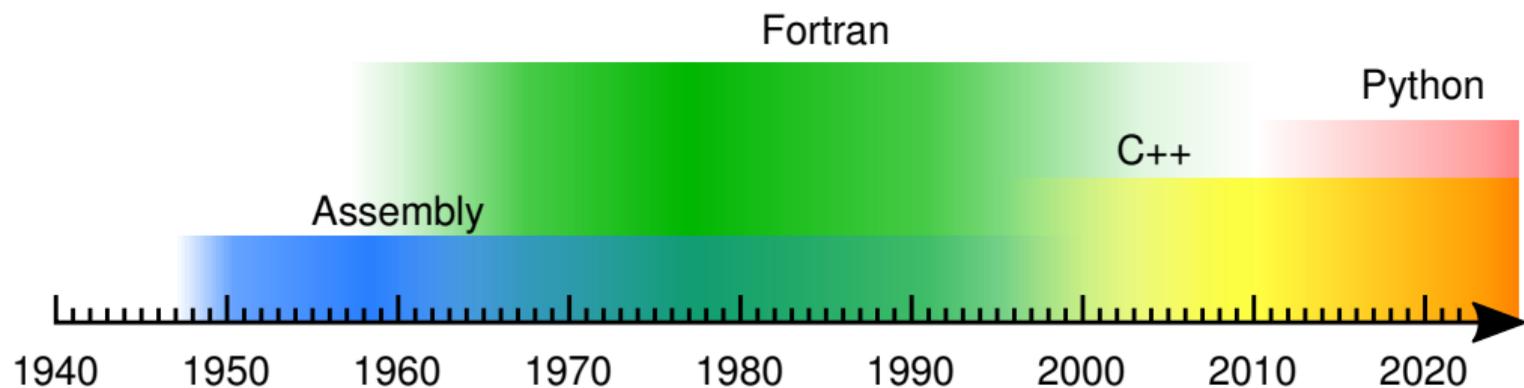


# Our community doesn't change or add languages often



C++ first appeared in 1985: 10–15 years before we adopted it

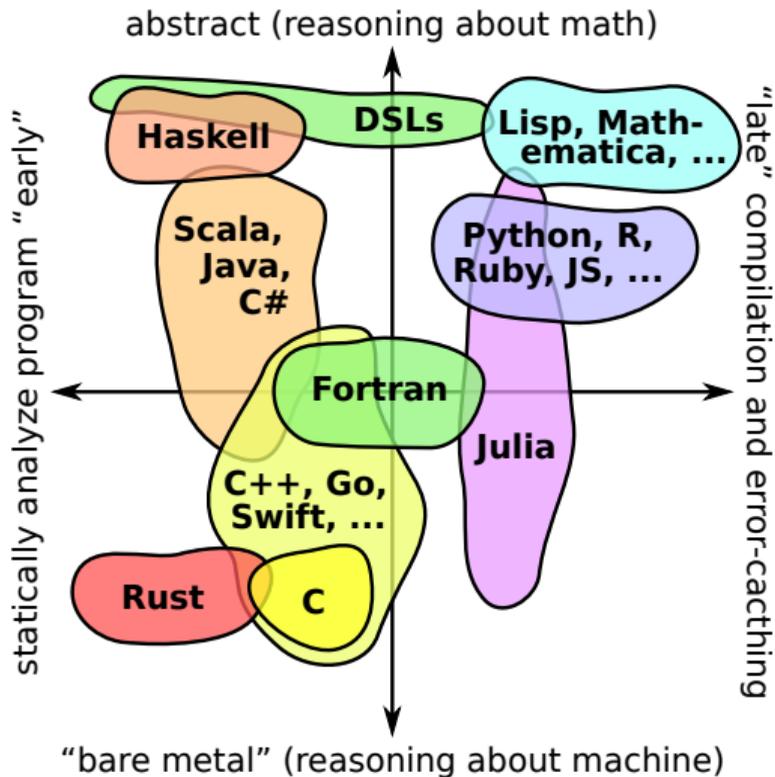
Python first appeared in 1990: 20–25 years before we adopted it



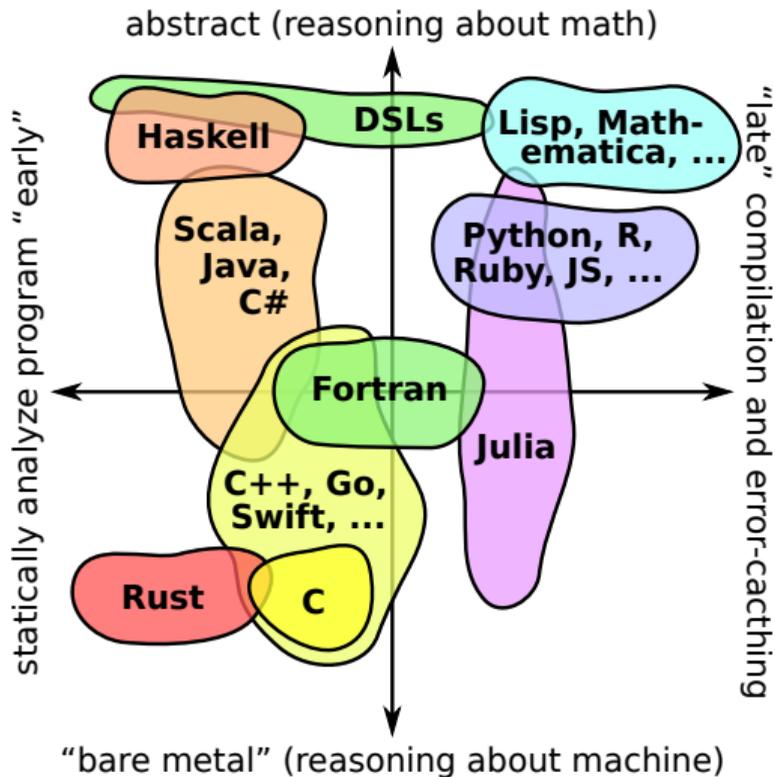
C++ first appeared in 1985: 10–15 years before we adopted it

Python first appeared in 1990: 20–25 years before we adopted it

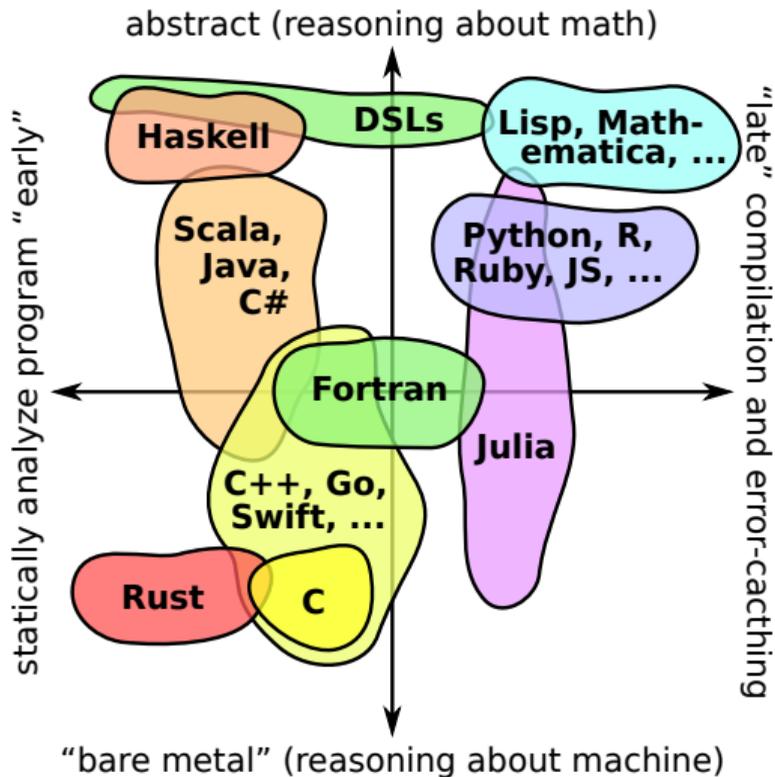
Julia first appeared in 2012: maybe we'll be using it by 2030



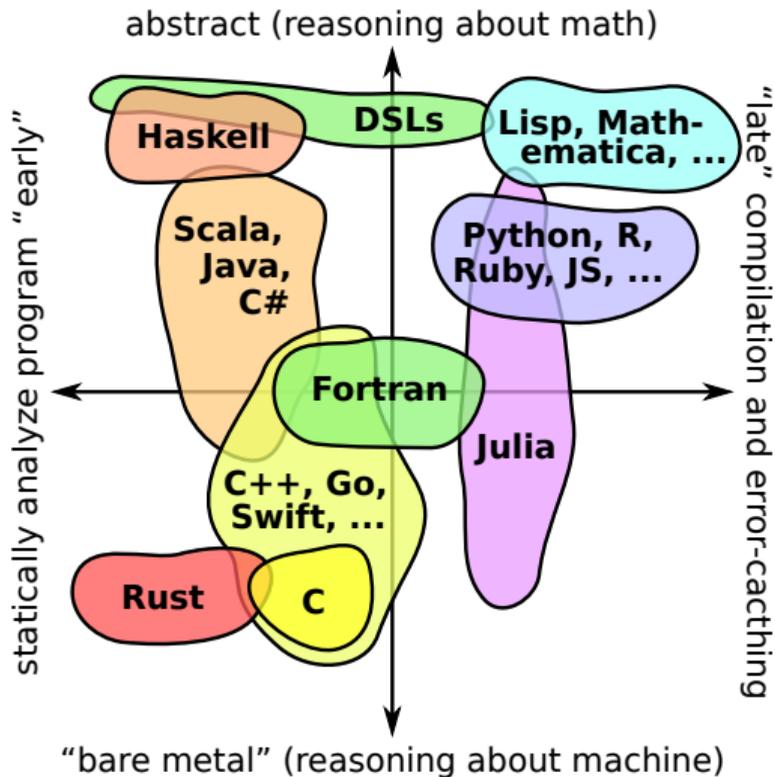
- ▶ Python, R, Ruby, Javascript, Lua, etc. are all relatively abstract, late error-catching glue languages.



- ▶ Python, R, Ruby, Javascript, Lua, etc. are all relatively abstract, late error-catching glue languages. Python isn't special: it's just where the best libraries self-gravitated.

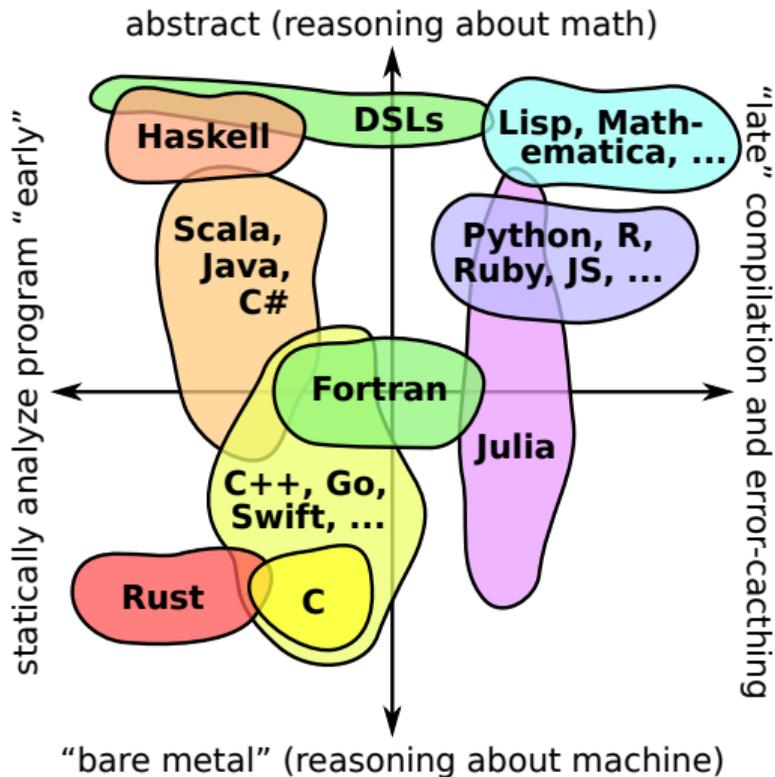


- ▶ Python, R, Ruby, Javascript, Lua, etc. are all relatively abstract, late error-catching glue languages. Python isn't special: it's just where the best libraries self-gravitated.
- ▶ Julia is different: everything is "just-in-time" (JIT) compiled, so it is flexible like Python, but fast, too.



- ▶ Python, R, Ruby, Javascript, Lua, etc. are all relatively abstract, late error-catching glue languages. Python isn't special: it's just where the best libraries self-gravitated.
- ▶ Julia is different: everything is "just-in-time" (JIT) compiled, so it is flexible like Python, but fast, too.
- ▶ Rust is a unique extreme. Triggers?

# An aside on future languages



- ▶ Python, R, Ruby, Javascript, Lua, etc. are all relatively abstract, late error-catching glue languages. Python isn't special: it's just where the best libraries self-gravitated.
- ▶ Julia is different: everything is "just-in-time" (JIT) compiled, so it is flexible like Python, but fast, too.
- ▶ Rust is a unique extreme. Triggers?
- ▶ DSLs (domain specific languages) have physics concepts baked in with opportunities for strict error checking and flexible backends.



## LHADA/CutLang $\rightarrow$ ADL: syntax

### Event selection regions

# preselection region

region **preselection**

select **size**(AK4jets) >= 3

select **size**(AK8jets) >= 1

select **MR** > 800

select **Rsq** > 0.08

# control region for tt+jets

region **ttjetsCR**

select **preselection**

select **size**(leptonsVeto) == 0

select **size**(Wjets) >= 1

select **dphimegajets** < 2.8

select **MT** > 100

select **size**(bjetsLoose) >= 1

bin **MR** [] 800 1000 and **Rsq** [] 0.08 0.1

bin **MR** [] 800 1000 and **Rsq** [] 0.1 0.2

bin ...



## F.A.S.T. F.A.S.T analysis tools based on YAML: Syntax

```
BasicVars:
  variables:
    - Muon_Pt: "sqrt(Muon_Px ** 2 + Muon_Py ** 2)"
    - IsoMuon_Idx: (Muon_Iso / Muon_Pt) < 0.10
      # This next variable will create a single
      # number for each event, using a set of inputs
      # whose length varies for each event
    - NIsoMuon:
      formula: IsoMuon_Idx
      reduce: count_nonzero
    - HasTwoMuons: NIsoMuon >= 2
      # Capture first muon's Pt, padded
      # with NaNs if NMuon < 1
    - Muon_lead_Pt: {reduce: 0, formula: Muon_Pt}
      # Capture second muon's Pt, padded
      # with NaNs if NMuon < 2
    - Muon_sublead_Pt: {reduce: 1, formula: Muon_Pt}
```

```
DiMu_controlRegion:
  weights: {nominal: weight}
  selection:
    All:
      - {reduce: 0, formula: Muon_pt > 30}
      - leadJet_pt > 100
    - All:
      - DiMuon_mass > 60
      - DiMuon_mass < 120
    - Any:
      - nCleanedJet == 1
      - DiJet_mass < 500
      - DiJet_deta < 2
```

```
DiMuonMass:
  dataset_col: true
  binning:
    - in: DiMuon_Mass
      out: dimu_mass
      bins: {low: 60, high: 120, nbins: 60}
  weights: {weighted: EventWeight}
```



## NAIL Natural Analysis Implementation Language: syntax

```
flow.Define("Muon_p4", "@p4v(Muon)")
flow.Define("Electron_p4", "@p4v(Electron)")
flow.Define("Electron_pid", "Electron_pt*0+11")
flow.Define("Muon_pid", "Muon_pt*0+13")
flow.MergeCollections("Lepton", ["Muon", "Electron"])
flow.Define("Lepton_index", "Range(nLepton)")
flow.Distinct("LPair", "Lepton")
flow.Selection("twoLeptons", "nLepton>=2")
flow.Define("isOSSF", "LPair0_charge != LPair1_charge && LPair0_pid == LPair1_pid", requires=["twoLeptons"])
flow.Selection("hasOSSF", "Sum(isOSSF) > 0")
flow.TakePair("Z", "Lepton", "LPair", "Argmax(-abs(MemberMap((LPair0_p4+LPair1_p4), M()) *isOSSF-91.2))", requires=["hasOSSF"])
flow.Selection("threeLeptons", "nLepton>=3", requires=["twoLeptons", "hasOSSF"])
flow.SubCollection("ResidualLeptons", "Lepton", sel="Lepton_index != LPair0[Z_indices] && Lepton_index != LPair1[Z_indices]",
flow.ObjectAt("ResidualLepton", "ResidualLeptons", "Argmax(ResidualLeptons_pt)")

histosPerSelection={
  "threeLeptons" : ["METplusLepton_Mt"]
}

flow.binningRules = [
  (".*Mt.*", "100,0,1000")
]
```



## hep\_tables and dataframe\_expressions: syntax

Associating electrons with corresponding GEN particles:

```
mc_part = df.TruthParticles('TruthParticles')
mc_ele = mc_part[(mc_part.pdgId == 11) | (mc_part.pdgId == -11)]

eles = df.Electrons('Electrons')

def good_e(e):
    'Good electron particle'
    return (e.ptgev > 20) & (abs(e.eta) < 1.4)

good_eles = eles[good_e]
good_mc_ele = mc_ele[good_e]
```

Key line: for each good electron, select only the very near good MC electrons.

Filter by all that are near by for each electron we are "looking" at!

Build the data model to make life easy...

```
def associate_particles(source, pick_from):
    ...
    Associate each particle from source with a close by one from the partic

    Args:
        source          The particles we want to start from
        pick_from       For each particle from source, we'll find a close
        name            Naming we can use when we extend the data model.

    Returns:
        with_assoc      The source particles that had a close by match
        ...

    def dr(p1, p2):
        'short hand for calculating DR between two particles.'
        return DeltaR(p1.eta(), p1.phi(), p2.eta(), p2.phi())

    def very_near(picks, p):
        'Return all particles in picks that are DR less than 0.1 from p'
        return picks[lambdas ps: dr(ps, p) < 0.1]

    source[f'all'] = lambda source_p: very_near(pick_from, source_p)

    source[f'has_match'] = lambda e: e.all.Count() > 0
    with_assoc = source[source.has_match]
    with_assoc['mc'] = lambda e: e.all.First()

    return with_assoc

matched = associate_particles(good_eles, good_mc_ele)
```

Note the lambda capture!!!



## PartiQL (demo) → AwkwardQL (real implementation)

```
# "For events with at least three leptons (electrons or muons) and a same-flavor  
# opposite-sign lepton pair, find the same-flavor opposite-sign lepton pair with a  
# mass closest to 91.2 GeV; make a histogram of the pT of the leading other lepton."  
leptons = electrons union muons  
cut count(leptons) >= 3 named "three_leptons" {  
  Z = electrons as (lep1, lep2) union muons as (lep1, lep2)  
    where lep1.charge != lep2.charge  
    min by abs(mass(lep1, lep2) - 91.2)  
  third = leptons except [Z.lep1, Z.lep2] max by pt  
  hist third.pt by regular(100, 0, 250) named "third_pt"  
}
```



*has happened*

*is happening*

*will happen*



Languages

Data formats

Libraries

Services

# Language choice does not determine performance!



```
import numpy
```

```
def run(height, width, maxiterations=20):
```

```
    y, x = numpy.ogrid[-1:0:height*1j, -1.5:0:width*1j]
```

```
    c = x + y*1j
```

```
    fractal = numpy.full(c.shape, maxiterations,  
                        dtype=numpy.int32)
```

```
    for h in range(height):
```

```
        for w in range(width):
```

```
            z = c[h, w]
```

```
            for i in range(maxiterations):
```

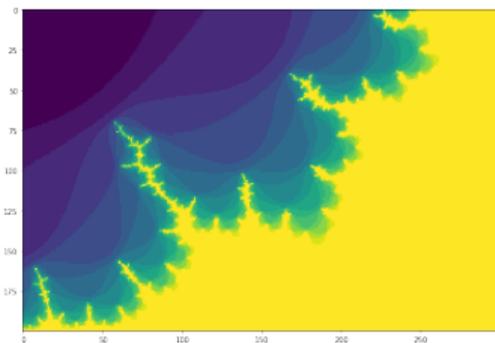
```
                z = z**2 + c[h, w]
```

```
                if abs(z) > 2:
```

```
                    fractal[h, w] = i
```

```
                    break
```

```
    return fractal
```



```
        # for each pixel (h, w)...
```

```
        # iterate at most 20 times
```

```
        # applying  $z \rightarrow z^2 + c$ 
```

```
        # if it diverges ( $|z| > 2$ )
```

```
        # color the plane with the iteration number
```

```
        # we're done, no need to keep iterating
```

# Language choice does not determine performance!



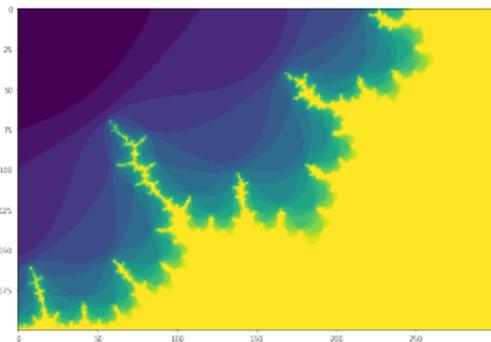
```
import numpy, numba
@numba.jit
def run(height, width, maxiterations=20):
    y, x = numpy.ogrid[-1:0:height*1j, -1.5:0:width*1j]
    c = x + y*1j
    fractal = numpy.full(c.shape, maxiterations,
                        dtype=numpy.int32)

    for h in range(height):
        for w in range(width):
            z = c[h, w]
            for i in range(maxiterations):
                z = z**2 + c[h, w]
                if abs(z) > 2:
                    fractal[h, w] = i
                    break

    return fractal
```

*# for each pixel (h, w)...*

*# iterate at most 20 times*  
*# applying  $z \rightarrow z^2 + c$*   
*# if it diverges ( $|z| > 2$ )*  
*# color the plane with the iteration number*  
*# we're done, no need to keep iterating*



Now **50× faster**, about as fast as C code (-O3).



The screenshot shows the Numba website homepage. At the top, there is a navigation bar with a lightning bolt icon on the left and several menu items: "Learn Numba in 5 minutes", "Documentation" (with a dropdown arrow), "Install", "Examples", "Talks/Tutorials", and "Community" (with a dropdown arrow). The main content area features the Numba logo, which consists of a blue lightning bolt icon followed by the word "Numba" in a large, bold, blue sans-serif font. Below the logo, the text reads "Numba makes Python code fast". Further down, a paragraph states: "Numba is an open source JIT compiler that translates a subset of Python and NumPy code into fast machine code." At the bottom of this section, there are two buttons: a white button with a blue border and text "Learn More", and a dark grey button with white text "Try Numba »".



## GOOD

```
>>> @nb.njit
... def monte_carlo_pi(nsamples):
...     total = 0
...     for i in range(nsamples):
...         x = random.random()
...         y = random.random()
...         if x**2 + y**2 < 1:
...             total += 1
...     return 4.0 * total / nsamples
...
>>> monte_carlo_pi(int(1e9))
3.14156838
```

## BAD

```
>>> @nb.njit
... def monte_carlo_pi(nsamples):
...     total = "hello"
...     for i in range(nsamples):
...         if total == "hello":
...             total = 0
...         x = random.random()
...         y = random.random()
...         if x**2 + y**2 < 1:
...             total += 1
...     return 4.0 * total / nsamples
...
>>> monte_carlo_pi(int(1e9))
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
... (big traceback) ...
Cannot unify Literal[str](hello) and
  Literal[int](0) for 'total.5',
  defined at <stdin> (4)
```



## GOOD

```
>>> @nb.njit
... def monte_carlo_pi(nsamples):
...     total = 0
...     for i in range(nsamples):
...         x = random.random()
...         y = random.random()
...         if x**2 + y**2 < 1:
...             total += 1
...     return 4.0 * total / nsamples
...
>>> monte_carlo_pi(int(1e9))
3.14156838
```

## BAD

```
>>> @nb.njit
... def monte_carlo_pi(nsamples):
...     total = "hello"
...     for i in range(nsamples):
...         if total == "hello":
...             total = 0
...         x = random.random()
...         y = random.random()
...         if x**2 + y**2 < 1:
...             total += 1
...     return 4.0 * total / nsamples
...
>>> monte_carlo_pi(int(1e9))
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
... (big traceback) ...
Cannot unify Literal[str](hello) and
  Literal[int](0) for 'total.5',
  defined at <stdin> (4)
```

# PyPy JIT-compiles *all* of Python; Julia defines less dynamism



## PyPy

A fast, [compliant](#) alternative implementation of Python

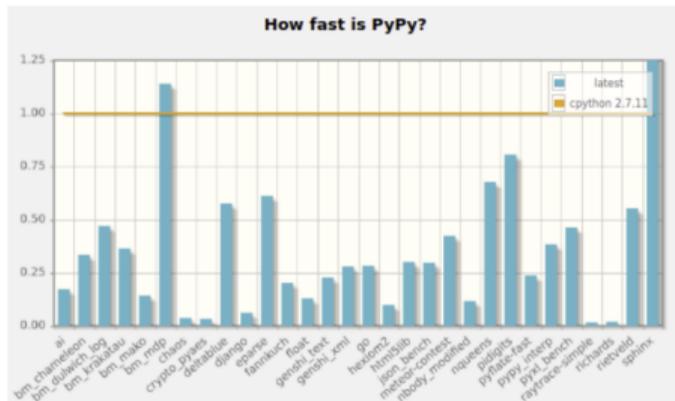
[Download PyPy](#)



[What is PyPy ?](#)

[Documentation \(external link\)](#)

On average, PyPy is **4.2 times faster** than CPython



PyPy trunk (with JIT) benchmark times normalized to CPython. Smaller is better. Based on the geometric average of all benchmarks



## The Julia Programming Language

[Download v1.5.2](#)

[Documentation](#)

★ Star 29,557

### Julia in a Nutshell

#### Fast

Julia was designed from the beginning for [high performance](#). Julia programs compile to efficient native code for [multiple platforms](#) via LLVM.

#### Reproducible

[Reproducible environments](#) make it possible to recreate the same Julia environment every time, across platforms, with [pre-built binaries](#).

#### General

Julia provides [asynchronous I/O](#), [metaprogramming](#), [debugging](#), [logging](#), [profiling](#), [a package manager](#), and more. One can build entire [Applications](#) and [Microservices](#) in Julia.

#### Dynamic

Julia is [dynamically typed](#), feels like a scripting language, and has good support for [interactive use](#).

#### Composable

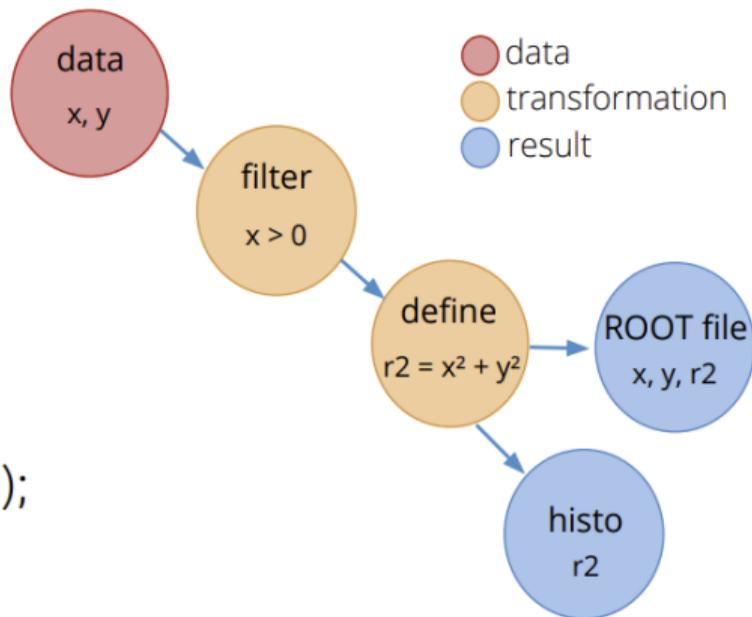
Julia uses [multiple dispatch](#) as a paradigm, making it easy to express many object-oriented and [functional programming](#) patterns. The talk on the [Unreasonable Effectiveness of Multiple Dispatch](#) explains why it works so well.

#### Open source

Julia is an open source project with over 1,000 contributors. It is made available under the [MIT license](#). The [source code](#) is available on GitHub.

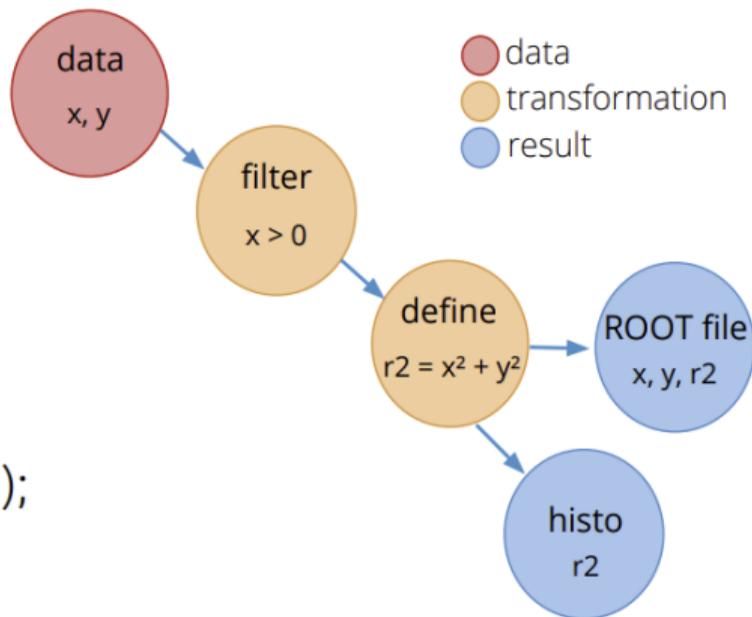


```
ROOT::RDataFrame df(dataset);  
auto df2 = df.Filter("x > 0")  
    .Define("r2", "x*x + y*y");  
auto rHist = df2.Histo1D("r2");  
df2.Snapshot("newtree", "newfile.root");
```





```
ROOT::RDataFrame df(dataset);  
auto df2 = df.Filter("x > 0")  
    .Define("r2", "x*x + y*y");  
auto rHist = df2.Histo1D("r2");  
df2.Snapshot("newtree", "newfile.root");
```

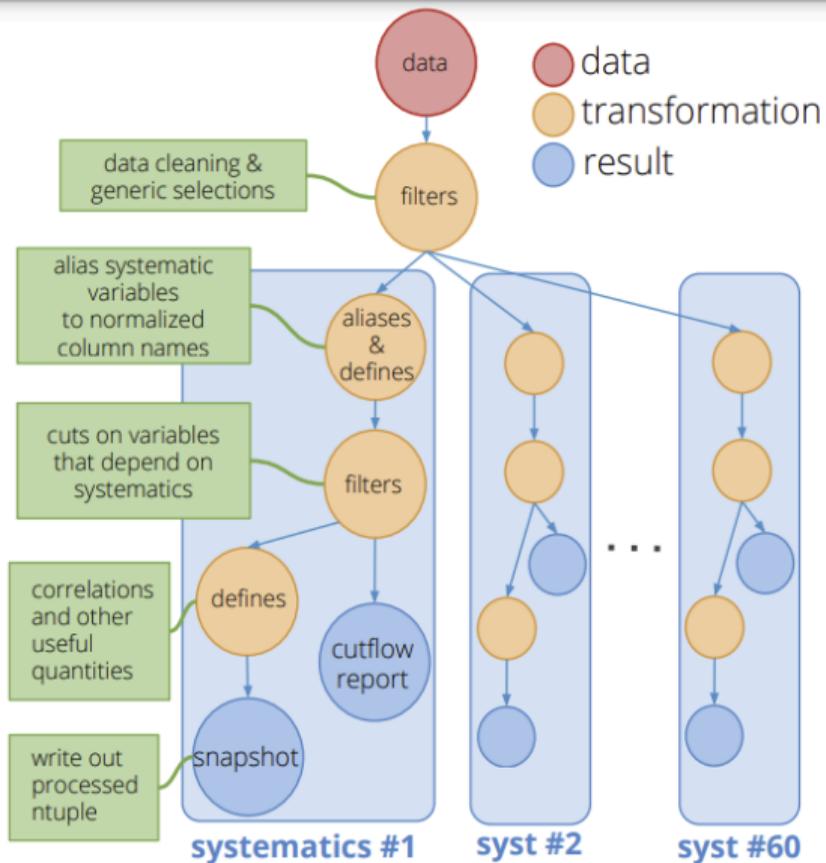


Without an explicit **for** loop, tasks without a dependency arrow can be parallelized or distributed across a cluster.

(see Enrico Guiraud @ CHEP 2018)



# Case study: ATLAS SUSY ntuple $\rightarrow$ ntuple



Local ntuple  $\rightarrow$  ntuple processing, MC data is processed to add quantities relevant for publication

- $\rightarrow$  program's main reads similarly to this graph
- $\rightarrow$  the large blue boxes represent one single function that applies the same operations to an RDF variable and is re-used for all different systematics
- $\rightarrow$  cuts, calculations and writing of the 60 output trees all happen in the same multi-thread event loop



# Case study

# tuple → ntuple

LoopSUSYFrame graph  
Thursday, 26 June 2018 23:44

dots

Filter

(cleaning and generic selections)

... x ~ 60

Define with syst names

Define with syst

Filter on variables that depend on syst

... This block is mirrored for every syst

CutFlow Report

use both of Define

Snapshot

data

filters

data cleaning & generic selections

alias systematic variables to normalized column names

aliases & defines

cuts on variables that depend on systematics

filters

correlations and other useful quantities

defines

cutflow report

write out processed ntuple

snapshot

systematics #1

syst #2

able processing, MC data is quantities relevant for publication

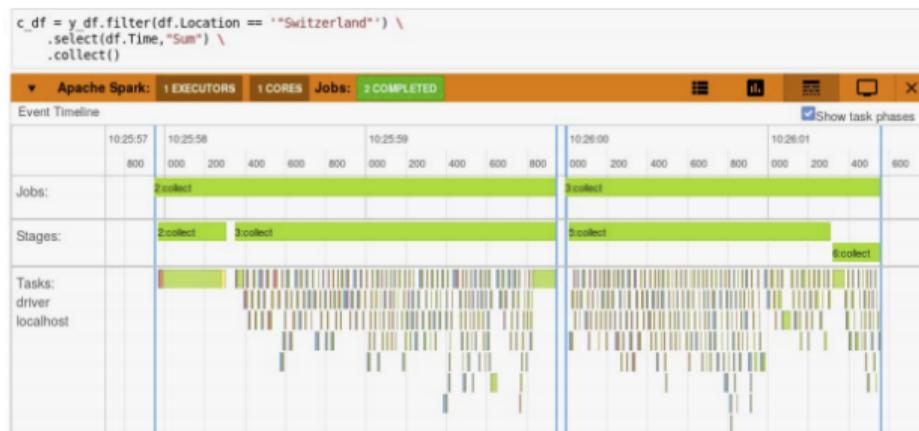
similarly to this graph

present one single function operations to an RDF variable different systematics

ing of the 60 output trees multi-thread event loop

# Distributed Monitoring

- ▶ **Bridge the gap** between interactive computing and distributed data processing
- ▶ Automatically appears when a Spark job is submitted from a cell
- ▶ Progress bars, task timeline, resource utilisation



Job ID	Job Name	Status	Stages	Tasks	Submission Time	Duration
2	reduce	COMPLETED	2/2	48 / 48	5 minutes ago	3s
Stage Id	Stage Name	Status	Tasks	Submission Time	Duration	
5	reduce	COMPLETED	32 / 32	5 minutes ago	2s	
4	coalesce	COMPLETED	16 / 16	5 minutes ago	0s	
3	foreach	COMPLETED	1/1 (1 skipped)	32 / 32	5 minutes ago	1m:20s
Stage Id	Stage Name	Status	Tasks	Submission Time	Duration	
6	coalesce	SKIPPED	0 / 0	Unknown	-	
7	foreach	COMPLETED	32 / 32	5 minutes ago	1m:20s	



Google Summer of Code



## C++

```
d.Filter([](double t) { return t > 0.; }, {"theta"})  
  .Snapshot<vector<float>>("t", "f.root", {"pt_x"});
```

---

## C++ with cling's just-in-time compilation

```
d.Filter("theta > 0").Snapshot("t", "f.root", "pt_x");
```

---

## PyROOT, automatically generated Python bindings

```
d.Filter("theta > 0").Snapshot("t", "f.root", "pt_x")
```



Events are distributed, but event code is C++ on `ROOT::VecOps::RVec<Particle>`.

```
getPt_code = ""
using namespace ROOT::VecOps;
RVec<double> getPt(const RVec<FourVector> &tracks) {
    auto pt = [](const FourVector &v) { return v.pt(); };
    return Map(tracks, pt);
}
""
ROOT.gInterpreter.Declare(getPt_code)

getPtWeights_code = ""
using namespace ROOT::VecOps;
RVec<double> getPtWeights(const RVec<FourVector> &tracks) {
    auto ptWeight = [](const FourVector &v) { return 1. / v.Pt(); };
    return Map(tracks, ptWeight);
};
""
ROOT.gInterpreter.Declare(getPtWeights_code)

augmented_d = (d.Define("tracks_n", "(int)tracks.size()")
               .Filter("tracks_n > 2")
               .Define("tracks_pts", "getPt(tracks)")
               .Define("tracks_pts_weights", "getPtWeights(tracks)"))
```





# Awkward Array

```
>>> import awkward1 as ak
>>> import numpy as np
>>>
>>> @ak.mixin_class(ak.behavior)
... class Lorentz:
...     @property
...     def pt(self):
...         return np.sqrt(self.px**2 + self.py**2)
...
>>> array = ak.Array([{"px": 1, "py": 1, "pz": 1, "E": 1},
...                   {"px": 2, "py": 2, "pz": 2, "E": 2},
...                   {"px": 3, "py": 3, "pz": 3, "E": 3},
...                   {"px": 4, "py": 4, "pz": 4, "E": 4},
...                   {"px": 5, "py": 5, "pz": 5, "E": 5}],
...                   with_name="Lorentz")
...
>>> array[-2]
<LorentzRecord {px: 4, py: 4, pz: 4, E: 4} type='Lorentz["px": int64, "py": int64, "pz": int64, "E": int64]'>
```



# Awkward Array

```
>>> import awkward1 as ak
>>> import numpy as np
>>>
>>> @ak.mixin_class(ak.behavior)
... class Lorentz:
...     @property
...     def pt(self):
...         return np.sqrt(self.px**2 + self.py**2)
...
>>> array = ak.Array([{"px": 1, "py": 1, "pz": 1, "E": 1},
...                   {"px": 2, "py": 2, "pz": 2, "E": 2},
...                   {"px": 3, "py": 3, "pz": 3, "E": 3},
...                   {"px": 4, "py": 4, "pz": 4, "E": 4},
...                   {"px": 5, "py": 5, "pz": 5, "E": 5}],
...                   with_name="Lorentz")
...
>>> array[-2]
<LorentzRecord {px: 4, py: 4, pz: 4, E: 4} type='Lorentz["px": int64, "py": int64... '>
>>> array * 10
<Array [{px: 10, py: 10, pz: 10, ... E: 50}] type='5 * {"px": int64, "py": int64... '>
```



# Awkward Array

```
>>> import awkward1 as ak
>>> import numpy as np
>>>
>>> @ak.mixin_class(ak.behavior)
... class Lorentz:
...     @property
...     def pt(self):
...         return np.sqrt(self.px**2 + self.py**2)
...
>>> array = ak.Array([{"px": 1, "py": 1, "pz": 1, "E": 1},
...                   {"px": 2, "py": 2, "pz": 2, "E": 2},
...                   {"px": 3, "py": 3, "pz": 3, "E": 3},
...                   {"px": 4, "py": 4, "pz": 4, "E": 4},
...                   {"px": 5, "py": 5, "pz": 5, "E": 5}],
...                   with_name="Lorentz")
...
>>> array[-2]
<LorentzRecord {px: 4, py: 4, pz: 4, E: 4} type='Lorentz["px": int64, "py": int64... '>
>>> array * 10
<Array [{px: 10, py: 10, pz: 10, ... E: 50}] type='5 * {"px": int64, "py": int64... '>
>>> array.pt
<Array [1.41, 2.83, 4.24, 5.66, 7.07] type='5 * float64'>
```



# Awkward Array

```
>>> import awkward1 as ak
>>> import numpy as np
>>>
>>> @ak.mixin_class(ak.behavior)
... class Lorentz:
...     @property
...     def pt(self):
...         return np.sqrt(self.px**2 + self.py**2)
...
>>> array = ak.Array([[{"px": 1, "py": 1, "pz": 1, "E": 1},
...                    {"px": 2, "py": 2, "pz": 2, "E": 2},
...                    {"px": 3, "py": 3, "pz": 3, "E": 3}],
...                   [],
...                   [{"px": 4, "py": 4, "pz": 4, "E": 4},
...                    {"px": 5, "py": 5, "pz": 5, "E": 5}]],
...                  with_name="Lorentz")
...
>>> array[0, -1]
<LorentzRecord {px: 3, py: 3, pz: 3, E: 3} type='Lorentz["px": int64, "py": int64, "pz": int64, "E": int64]'>
>>> array * 10
<Array [[{px: 10, py: 10, ... E: 50}]] type='3 * var * {"px": int64, "py": int64, "pz": int64, "E": int64}'>
>>> array.pt
<Array [[1.41, 2.83, 4.24], ... [5.66, 7.07]] type='3 * var * float64'>
```



# Awkward Array

```
>>> import awkward1 as ak
>>> import numpy as np
>>>
>>> @ak.mixin_class(ak.behavior)
... class Lorentz:
...     @property
...     def pt(self):
...         return np.sqrt(self.px**2 + self.py**2)
...
>>> array = ak.Array([[[{"px": 1, "py": 1, "pz": 1, "E": 1},
...                     {"px": 2, "py": 2, "pz": 2, "E": 2}],
...                   [{"px": 3, "py": 3, "pz": 3, "E": 3}]],
...                   [],
...                   [{"px": 4, "py": 4, "pz": 4, "E": 4}],
...                   [],
...                   [{"px": 5, "py": 5, "pz": 5, "E": 5}]]],
...                   with_name="Lorentz")
>>> array[-1, 2, 0]
<LorentzRecord {px: 5, py: 5, pz: 5, E: 5} type='Lorentz["px": int64, "py": int64, "pz": int64, "E": int64]'>
>>> array * 10
<Array [ [ [ [px: 10, py: 10, ... E: 50] ] ] ] type='3 * var * var * {"px": int64, "py": int64, "pz": int64, "E": int64}'>
>>> array.pt
<Array [ [ [1.41, 2.83], [4.24, ... [], [7.07] ] ] ] type='3 * var * var * float64'>
```



```
array = ak.Array([\n    [{"x": 1.1, "y": [1]}, {"x": 2.2, "y": [1, 2]}, {"x": 3.3, "y": [1, 2, 3]}],\n    [],\n    [{"x": 4.4, "y": [1, 2, 3, 4]}, {"x": 5.5, "y": [1, 2, 3, 4, 5]}]\n])
```

NumPy-like expression

```
output = np.square(array["y", ..., 1:])
```

```
[ \n    [[], [4], [4, 9]], \n    [], \n    [[4, 9, 16], [4, 9, 16, 25]] \n]
```

**4.6 seconds to run (2 GB footprint)**

equivalent Python

```
output = []\nfor sublist in python_objects:\n    tmp1 = []\n    for record in sublist:\n        tmp2 = []\n        for number in record["y"][1:]:\n            tmp2.append(np.square(number))\n        tmp1.append(tmp2)\n    output.append(tmp1)
```

**138 seconds to run (22 GB footprint)**

(single-threaded on a 2.2 GHz processor with a dataset 10 million times larger than the one shown)



$$\begin{array}{c} [[\square, \square, \square], [], [\square, \square]] \\ + \\ [ \quad \square, \quad \square, \quad \square ] \\ \hline [[\square, \square, \square], [], [\square, \square]] \end{array}$$

```
>>> jagged = ak.Array([[1, 2, 3], [], [4, 5]])
>>> flat   = ak.Array([ 100, 200, 300])
>>> jagged + flat
<Array [[101, 102, 103], [], [304, 305]] type='3 * var * int64'>
```



## Cartesian product

$[[\square, \square, \square], [], [\square]]$

$\otimes$

$[[\square, \square], [\square], [\square, \square]]$

---

$[[\square, \square, \square, \square, \square, \square], [], [\square, \square]]$

```
>>> x = [[1, 2, 3], [], [4]]
>>> y = [["a", "b"], ["c"], ["d", "e"]]
>>> ak.cartesian([x, y], axis=1)
[[ (1, "a"), (1, "b"),
  (2, "a"), (2, "b"),
  (3, "a"), (3, "b")],
 [],
 [(4, "d"), (4, "e")]]
```

## Combinations without replacement

$\binom{[\square, \square, \square, \square]}{2}, \binom{[\square]}{2}, \binom{[\square, \square]}{2}$

---

$[[\square, \square, \square, \square, \square, \square], [], [\square, \square]]$

```
>>> x = [[1, 2, 3, 4], [], [5, 6]]
>>> ak.combinations(x, 2, axis=1)
[[ (1, 2), (1, 3), (1, 4),
  (2, 3), (2, 4), (3, 4)],
 [],
 [(5, 6)]]
```

# Like NumPy, but structured: **reduction**



```
array = ak.Array(  
    [  
        [ 2, 3, 5],  
        [ ],  
        [None, 7],  
        [ 11 ]  
    ]  
)
```

[ 2, 3, 5]	→	30
[ ]	→	1
[None, 7]	→	7
[ 11 ]	→	11

22      21      5

```
>>> ak.prod(array, axis=0)  
<Array [22, 21, 5] type='3 * int64'>  
>>> ak.prod(array, axis=1)  
<Array [30, 1, 7, 11] type='4 * int64'>
```



You can mix fast, explicit **for** loops with NumPy-like slicing, all in Python.

```
>>> @nb.jit
... def make_cut(electrons, jets, builder):
...     for event_electrons, event_jets in zip(electrons, jets):
...         builder.begin_list()
...         for electron in event_electrons:
...             keep = False
...             for jet in event_jets:
...                 if abs(electron.x - jet.x) < 0.45: # calculate delta-R here
...                     keep = True
...                     break
...             builder.append(keep)
...         builder.end_list()
...     return builder
...
>>> cut = make_cut(electrons, jets, ak.ArrayBuilder()).snapshot()
>>> cut
<Array [[True, False, True], ... [True, False]] type='3 * var * bool'>
>>> electrons[cut].tolist()
[[{'x': 1, 'y': 1, 'z': 1, 't': 1}, {'x': 3, 'y': 3, 'z': 3, 't': 3}],
 [],
 [{'x': 4, 'y': 4, 'z': 4, 't': 4}]]
```

# vector: an Awkward 2D, 3D, Lorentz vector library in development



scikit-hep / vector

Unwatch

10

Star

2

Fork

1

Code

Issues 4

Pull requests

Actions

Projects

Wiki

Security

Insights

Settings

master

3 branches

0 tags

Go to file

Add file

Code



henryiii docs: update README

5f9e43e on Aug 11 41 commits

.github	ci: update steps	2 months ago
docs	style: add and fix styling (#15)	2 months ago
notebooks	refactor: always include T	2 months ago
src/vector	refactor: always include T	2 months ago
tests	refactor: cleanup and work on op.add	2 months ago
.gitignore	style: add and fix styling (#15)	2 months ago
.pre-commit-config.yaml	style: add and fix styling (#15)	2 months ago
.readthedocs.yml	Adding initial placeholder for documentation	7 months ago
LICENSE	Adding standard header	7 months ago
README.md	docs: update README	2 months ago

## About

Vector classes and utilities

[vector.readthedocs.io](https://vector.readthedocs.io)

scikit-hep

vector

Readme

BSD-3-Clause License

## Releases

No releases published

[Create a new release](#)

## Contributors 6



# First release of **hist**, a Pythonic histogram library, on Tuesday



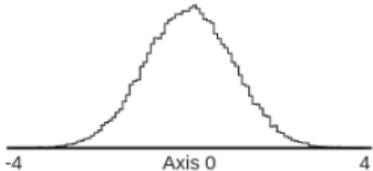
File Edit View Run Kernel Tabs Settings Help

hist-plots.ipynb Python 3

```
[1]: import uproot4
file = uproot4.open("hepdata-example.root")
```

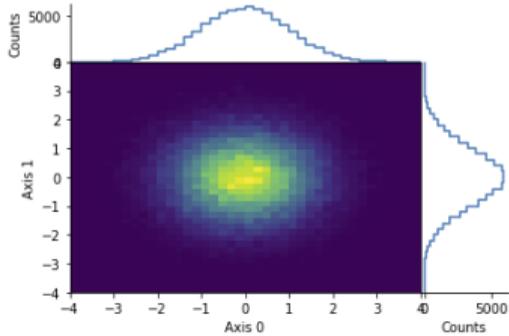
```
[2]: file["hpx"].to_hist()
```

[2]:



Regular(100, -4, 4, label='Axis 0')  
Double()  $\Sigma$  74994.0 (75000.0 with flow)

```
[3]: file["hpxpy"].to_hist().plot2d_full();
```



# hepunits: common units in HEP and **particle**: PDG data



```
File Edit View Run Kernel Tabs Settings Help
hepunits-particle.ipynb x
+ ✂ 📄 📄 ▶ ■ ⌂ ⏪ Code ⌄ ⌚ git Python 3 ○

[1]: import particle
      from hepunits.units import cm

      import IPython.display

      # Find all strange baryons with c*tau > 1 cm
      for x in particle.Particle.findall(lambda p: p.pdgid.is_baryon and p.pdgid.has_strange and p.ctau > 1 * cm):

          IPython.display.display(IPython.display.Latex("$" + x.latex_name + "$"))
          print(repr(x), end="\n\n")

Σ-
<Particle: name="Sigma-", pdgid=3112, mass=1197.45 ± 0.03 MeV>

Σ+
<Particle: name="Sigma~+", pdgid=-3112, mass=1197.45 ± 0.03 MeV>

Λ
<Particle: name="Lambda", pdgid=3122, mass=1115.683 ± 0.006 MeV>

Λ̄
<Particle: name="Lambda~", pdgid=-3122, mass=1115.683 ± 0.006 MeV>

Σ+
<Particle: name="Sigma+", pdgid=3222, mass=1189.37 ± 0.07 MeV>

Σ-
<Particle: name="Sigma~-", pdgid=-3222, mass=1189.37 ± 0.07 MeV>
```

# Scikit-HEP: a growing ecosystem of interoperating tools



Particle

zfit

uproot

iminuit

Boost  
Histogram

Scikit  
HEP

pyf  
differentiable  
Likelihoods

Decay  
Language

Coofi  
CUDA/OpenMP  
Fitting Framework  
for C++ & Python

Awkward  
Array

mplhep

hepstats



*has happened*

*is happening*

*will happen*



Languages

Data formats

Libraries

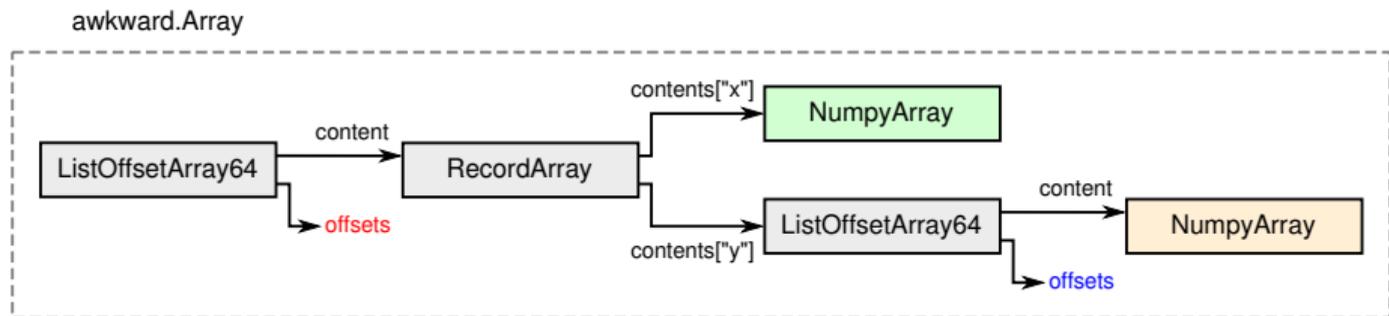
Services



# A **columnar** form can be iterated or vectorized

```
array = ak.Array([\n    [{"x": 1, "y": [11]},\n     {"x": 4, "y": [12, 22]},\n     {"x": 9, "y": [13, 23, 33]}],\n    [],\n    [{"x": 16, "y": [14, 24, 34, 44]}]\n])
```

```
outer offsets: 0, 3, 3, 5\ncontent for x: 1, 4, 9, 16\noffsets for y: 0, 1, 3, 6, 10\ncontent for y: 11, 12, 22, 13, 23, 33, 14, 24, 34, 44
```



# ROOT's TTree is columnar to one level (only intended for iteration)



```
tree = TTree([
    [{"x": 1, "y": [11]},
     {"x": 4, "y": [12, 22]},
     {"x": 9, "y": [13, 23, 33]}],
    []],
    [{"x": 16, "y": [14, 24, 34, 44]}]
])
```

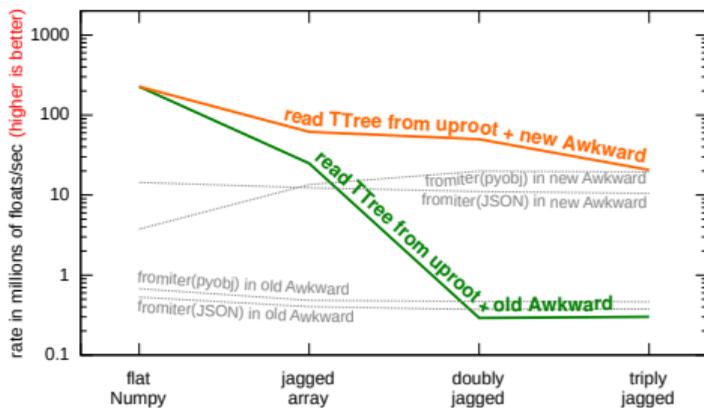
0  
x 1 y(1) 11 x 4 y(2) 12 22 x 9 y(3) 13 23 33      15 15      x 16 y(4) 14 24 34 44

# ROOT's TTree is columnar to one level (only intended for iteration)



```
tree = TTree([\n    [{"x": 1, "y": [11]},\n     {"x": 4, "y": [12, 22]},\n     {"x": 9, "y": [13, 23, 33]}],\n    [],\n    [{"x": 16, "y": [14, 24, 34, 44]}]\n])
```

0  
x 1 y(1) 11 x 4 y(2) 12 22 x 9 y(3) 13 23 33 x 16 y(4) 14 24 34 44



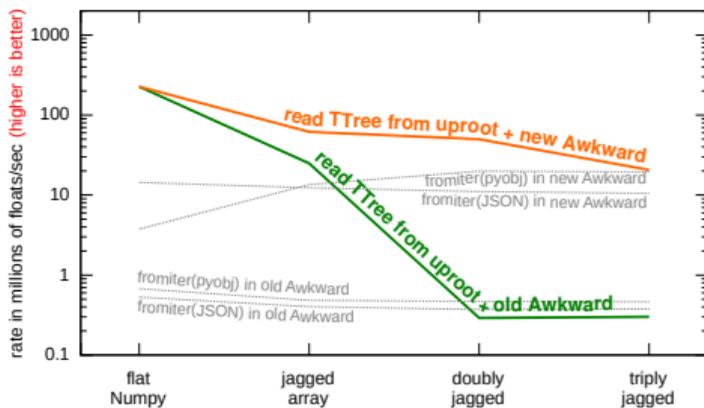
This is why Uproot reads **float** and **vector<float>** data much faster than **vector<vector<float>>** or higher.

# ROOT's TTree is columnar to one level (only intended for iteration)



```
tree = TTree([\n    [{"x": 1, "y": [11]},\n     {"x": 4, "y": [12, 22]},\n     {"x": 9, "y": [13, 23, 33]}],\n    [],\n    [{"x": 16, "y": [14, 24, 34, 44]}]\n])
```

0  
x 1 y(1) 11 x 4 y(2) 12 22 x 9 y(3) 13 23 33 x 16 y(4) 14 24 34 44



This is why Uproot reads **float** and **vector<float>** data much faster than **vector<vector<float>>** or higher.

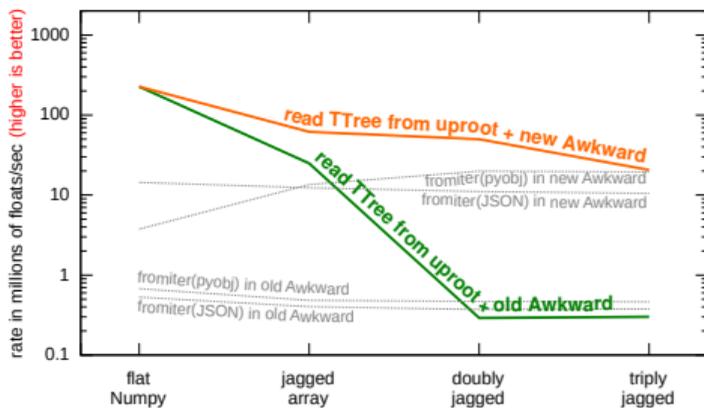
TTree interleaves data of more than one level of nesting, which can't be unraveled by NumPy.

# ROOT's TTree is columnar to one level (only intended for iteration)



```
tree = TTree([\n    [{"x": 1, "y": [11]},\n     {"x": 4, "y": [12, 22]},\n     {"x": 9, "y": [13, 23, 33]}],\n    [],\n    [{"x": 16, "y": [14, 24, 34, 44]}]\n])
```

0  
x 1 y(1) 11 x 4 y(2) 12 22 x 9 y(3) 13 23 33 x 16 y(4) 14 24 34 44

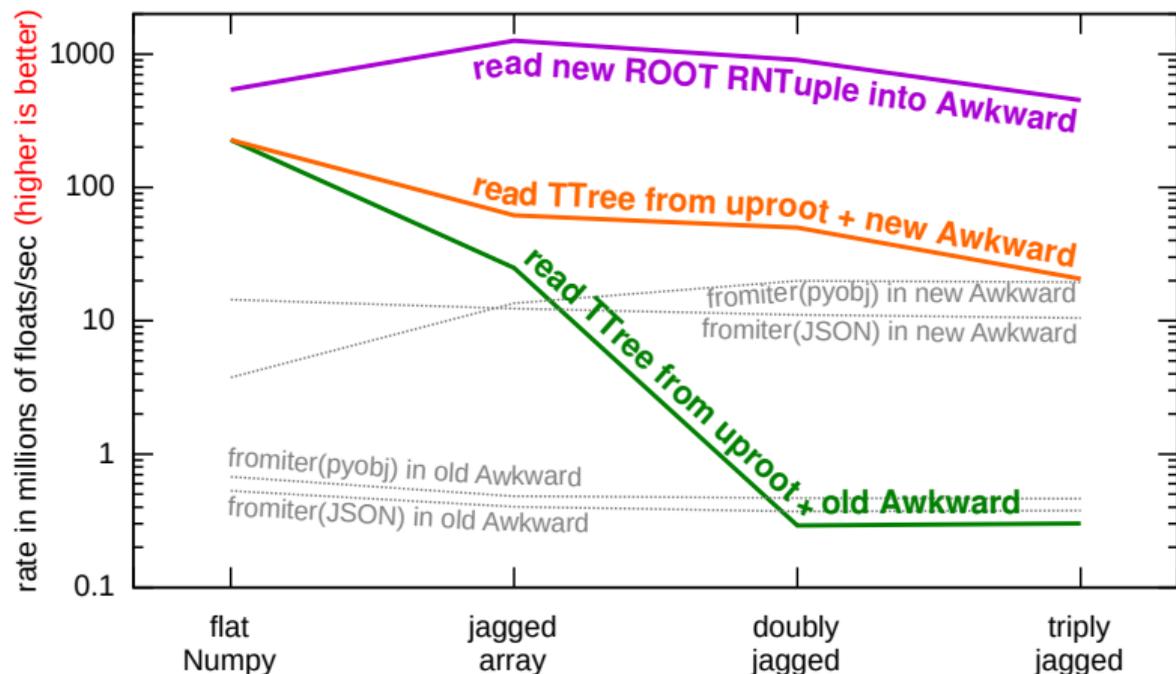


This is why Uproot reads **float** and `vector<float>` data much faster than `vector<vector<float>>` or higher.

TTree interleaves data of more than one level of nesting, which can't be unraveled by NumPy.

(But some performance can be recovered by moving the **for** loops from Python into C++.)

# RNTuple: the ROOT 7 replacement for TTree!



**RNTuple** is a redesign of TTree from scratch, using columnar data structures at all levels of depth.

Not having to deal with interleaved data at all is better than moving loops from Python to C++.



The first two layers are compatible/convertible to Awkward Array.

### Event iteration

Reading and writing in event loops and through RDataFrame  
RNTupleDataSource, RNTupleView, RNTupleReader/Writer

### Logical layer / C++ objects

Mapping of C++ types onto columns  
e.g. `std::vector<float>`  $\mapsto$  index column and a value column  
RField, RNTupleModel, REntry

### Primitives layer / simple types

“Columns” containing elements of fundamental types (float, int, ...) grouped into (compressed) pages and clusters  
RColumn, RColumnElement, RPage

### Storage layer / byte ranges

RPageStorage, RCluster, RNTupleDescriptor

### Approximate translation between TTree and RNTuple classes:

TTree	≈	RNTupleReader
		RNTupleWriter
TTreeReader	≈	RNTupleView
TBranch	≈	RField
TBasket	≈	RPage
TTreeCache	≈	RClusterPool

- ▶ Jakob Blomer and his students are developing RNTuple in ROOT.
- ▶ Oksana Shadura is developing the Python interface.

# Columnar data formats in industry: on disk, in memory, on GPUs



Apache Software Foundation / Apache Parquet

Apache Parquet is a [columnar storage](#) format available to any project in the Hadoop ecosystem, regardless of the choice of data processing framework, data model or programming language.

## Parquet Videos (more [presentations](#))



## Tweets by @ApacheParquet



Overview FAQ Blog Get Arrow Documentation Community ASF Links

# APACHE ARROW

A cross-language development platform for in-memory analytics

Star 6,241 Follow @ApacheArrow 7,504 followers

## What is Arrow?

### Format

Apache Arrow defines a language-independent columnar memory format for flat and hierarchical data, organized for efficient analytic operations on modern hardware like CPUs and GPUs. The Arrow memory format also supports zero-copy reads for lightning-fast data access without serialization overhead.

[Learn more about the design](#) or [read the specification](#).

### Libraries

Arrow's libraries implement the format and provide building blocks for a range of use cases, including high performance analytics. Many popular projects use Arrow to ship columnar data efficiently or as the basis for analytic engines.

Libraries are available for C, C++, C#, Go, Java, JavaScript, MATLAB, Python, R, Ruby, and

## RAPIDS

# Open GPU Data Science

GET STARTED

## GPU DATA SCIENCE

### ① ACCELERATED DATA SCIENCE

The RAPIDS suite of open source software libraries gives you the freedom to execute end-to-end data science and analytics pipelines entirely on GPUs. [Learn about RAPIDS >>](#)

### ⚙️ SCALE OUT ON GPUS

Seamlessly scale from GPU workstations to multi-GPU servers and multi-node clusters with Dask. [Learn about Dask >>](#)

### 🐍 PYTHON INTEGRATION

Accelerate your Python data science toolchain with minimal code changes and no new tools to learn. [Learn about our libraries >>](#)

### 🎯 TOP MODEL ACCURACY

Increase machine learning model accuracy by iterating on models faster and deploying them more frequently. [Learn about RAPIDS for model optimization >>](#)

# Columnar data formats in industry: on disk, in memory, on GPUs



Apache Software Foundation / Apache Parquet

Apache Parquet is a [columnar storage](#) format available to any project in the Hadoop ecosystem, regardless of the choice of data processing framework, data model or programming language.

## Parquet Videos (more [presentations](#))



## Tweets by @ApacheParquet



## What is Arrow?

### Format

Apache Arrow defines a language-independent columnar memory format for flat and hierarchical data, organized for efficient analytic operations on modern hardware like CPUs and GPUs. The Arrow memory format also supports zero-copy reads for lightning-fast data access without serialization overhead.

[Learn more about the design](#) or [read the specification](#).

### Libraries

Arrow's libraries implement the format and provide building blocks for a range of use cases, including high performance analytics. Many popular projects use Arrow to ship columnar data efficiently or as the basis for analytic engines.

Libraries are available for C, C++, C#, Go, Java, JavaScript, MATLAB, Python, R, Ruby, and



## GPU DATA SCIENCE

### Ⓜ ACCELERATED DATA SCIENCE

The RAPIDS suite of open source software libraries gives you the freedom to execute end-to-end data science and analytics pipelines entirely on GPUs. [Learn about RAPIDS >>](#)

### 🔗 SCALE OUT ON GPUS

Seamlessly scale from GPU workstations to multi-GPU servers and multi-node clusters with Dask. [Learn about Dask >>](#)

### 🐍 PYTHON INTEGRATION

Accelerate your Python data science toolchain with minimal code changes and no new tools to learn.

### 🎯 TOP MODEL ACCURACY

Increase machine learning model accuracy by iterating on models faster and deploying them more frequently. [Learn about RAPIDS for model optimization >>](#)

Since RNTuple is fully columnar, it will be easier to convert.



*has happened*

*is happening*

*will happen*

Languages

Data formats

Libraries

Services



As datasets get too big to ntuple, we'll have to move the pre-aggregation steps of analysis to a centralized service.



As datasets get too big to ntupleize, we'll have to move the pre-aggregation steps of analysis to a centralized service.

Need to reproduce the low-latency experience of a local ntuple on a shared system.

- ▶ It will be harder because it will be more in demand at peak times.
- ▶ It will be easier because popular data can be in shared caches.



As datasets get too big to ntupleize, we'll have to move the pre-aggregation steps of analysis to a centralized service.

Need to reproduce the low-latency experience of a local ntuple on a shared system.

- ▶ It will be harder because it will be more in demand at peak times.
- ▶ It will be easier because popular data can be in shared caches.

Strong bias toward reusing existing technologies: Kubernetes, Docker, Kafka, Spark, Dask, and perhaps Ray.



As datasets get too big to ntupleize, we'll have to move the pre-aggregation steps of analysis to a centralized service.

Need to reproduce the low-latency experience of a local ntuple on a shared system.

- ▶ It will be harder because it will be more in demand at peak times.
- ▶ It will be easier because popular data can be in shared caches.

Strong bias toward reusing existing technologies: Kubernetes, Docker, Kafka, Spark, Dask, and perhaps Ray.

Query services are now at the demo stage!

- ▶ ServiceX: <https://youtu.be/NEgqGRdt360>
- ▶ CoffeaCasa: <https://youtu.be/CDIFd1gDbSc>



- ▶ The transition to more-Python-than-C++ has already happened.
- ▶ DSLs are also being explored, but they're at a much earlier stage.
- ▶ Just-in-time compilation is important: Numba, PyPy, Julia, and RDataFrame.
- ▶ So is “removing the **for** loop.”
- ▶ Awkward Array has implicit loops at all levels for NumPy-like manipulation, and also provides an onramp into Numba for conventional code.
- ▶ The Scikit-HEP ecosystem is booming!
- ▶ Columnar internal structures are essential: ROOT is developing RNTuple.
- ▶ Columnar data formats are popular in industry, too. We must interoperate!
- ▶ Query services are at the demo stage.