



# The Python Scientific Software Ecosystem: Past, Present and Future

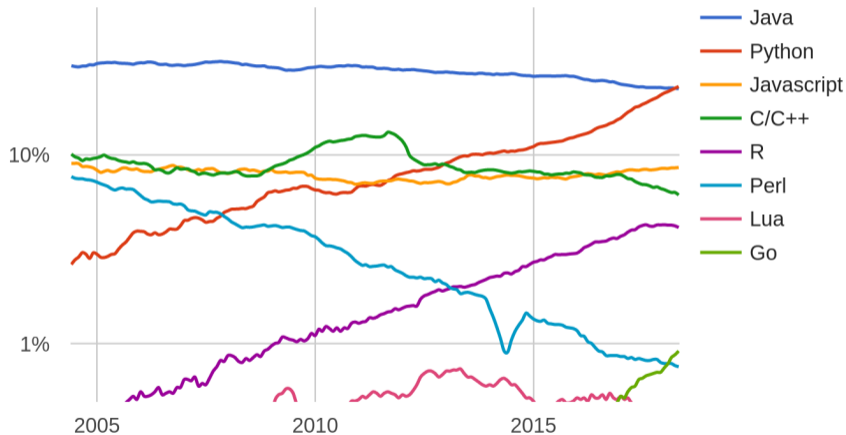
Jim Pivarski

Princeton University – DIANA-HEP

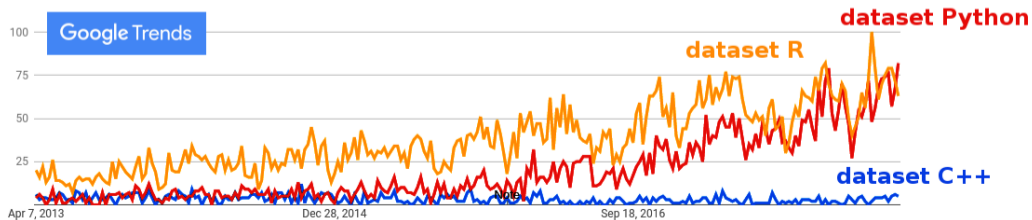
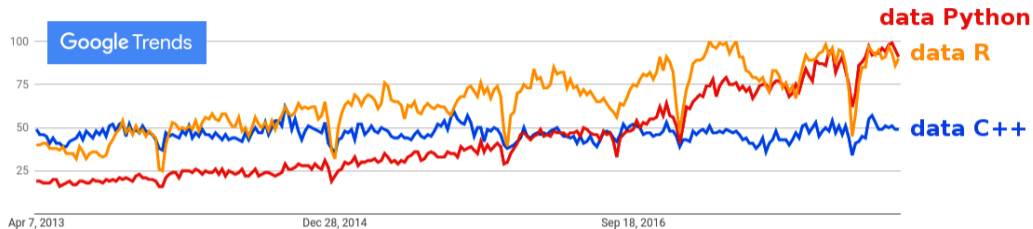
July 7, 2018



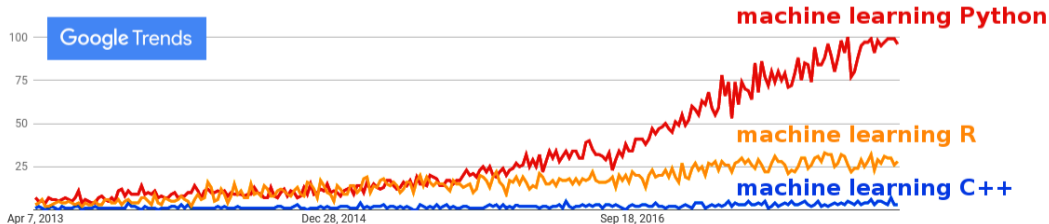
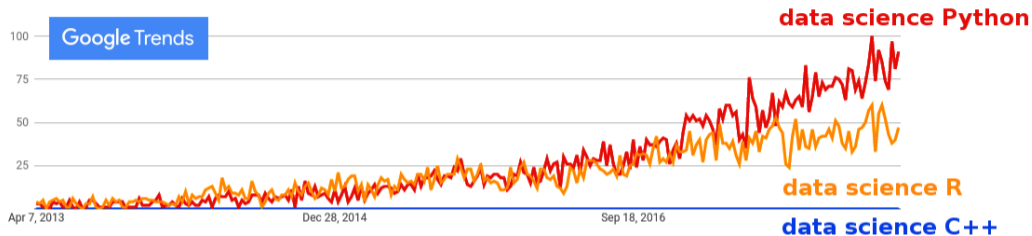
## PYPL Popularity of Programming Language



# Why we're here



# Why we're here





All of the machine learning libraries I could find either have a Python interface or are primarily/exclusively Python.



PYTORCH



Keras



GLUON



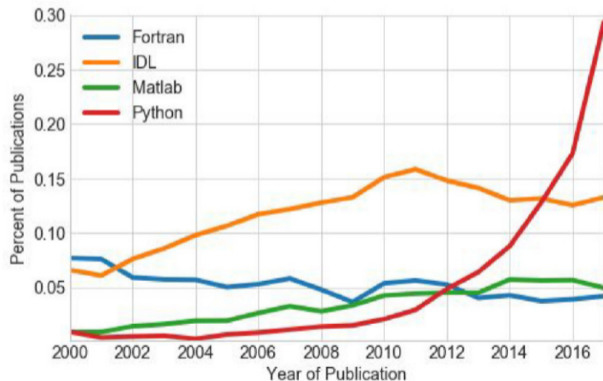
ONNX



dmlc  
XGBoost



## Mentions of Software in Astronomy Publications:



Compiled from NASA ADS ([code](#)).

Thanks to Juan Nunez-Iglesias,  
Thomas P. Robitaille, and Chris Beaumont.



# Python in HEP?



## Python in HEP?

- ▶ Collaboration frameworks like Athena and CMSSW are configured or driven by Python.





## Python in HEP?

- ▶ Collaboration frameworks like Athena and CMSSW are configured or driven by Python.
- ▶ Today, it is common for physicists to do their analyses in Python/PyROOT.  
(Half Python, half C++? Anyway, a lot more than in LHC Run I.)



## Python in HEP?

- ▶ Collaboration frameworks like Athena and CMSSW are configured or driven by Python.
- ▶ Today, it is common for physicists to do their analyses in Python/PyROOT. (Half Python, half C++? Anyway, a lot more than in LHC Run I.)
- ▶ Python is the most natural bridge to machine learning and other statistical software written outside of HEP.

# I was an early adopter (thesis workflow from 2006)



Earlier steps in  
Mathematica, Perl,  
GnuPlot, Emacs Lisp,  
PAW, and ROOT

## Fit function

```
DOUBLE PRECISION FUNCTION GMMF(RN,GAM,WSPREAD,HC)
  IMPLICIT NONE
  C IMPLICIT REAL (A-H,O-Z), INTEGER (I-N)
  C Observed shape of resonance peak, starting with Breit-Wigner
  C Normalized per unit W, W0=1
  C F_XF(x,s) convoluted with unit-area Breit-Wigner, and with beam
  resolution
  C unit-area Gaussian
  C If GAM=1 W0, the Breit-Wigner is replaced with a Breit function
  C RW=1 (GMMF2p1) if ((w-W0)^2<GAM^2/4) (pure integral over w)
  C Gaussian(w) = exp(-.5*(w-W0)^2/(GAM^2*(WSPREAD))) (unit int....)

  INTEGER MWIN, MWIN_L, MWIN_R
  REAL*8 ROOT2, PI, RTMDF1, RTMDFPI, PWA, TWY

  REAL*8 FKFT(MWIN) ! F_XF(x,s) (Mureev-Fadin eq.28)
  REAL*8 HA(W0:RN) ! WA-W0 of centroid WA bin
  REAL*8 DWA(W0:RN) ! bin width in WA-W0
  REAL*8 HT(0:MTWIN) ! 1/MTW bin center
  REAL*8 DM(0:MTWIN) ! result of BW-FK_F convolution
  REAL*8 BMSC(MWIN:MTWIN) ! result of BW-FK_F convolution

  DO 11 I=MTWIN-MWIN,MTWIN ! sum over W's
    H=HT(IABS(ITEM)) ! W-W0
    IF (ITEM.LT.0) H=-H
    SIG=(W-H)/((GAM/2)*WSPREAD) ! sqrt of exponent
    IF ((ITEM.EQ.MTWIN).AND.(SIG.LT.2.)) GO TO 12 ! running out of bins?
    IF (SIG.GT.3.) GO TO 13 ! Gaussian getting negligible?
    P=COS(SIG) ! exponent
    LFP.L.PWA*HT(I) THEN
      IF (DM(IABS(ITEM)).GT.1)*WSPREAD GO TO 12 ! binning too coarse
      GMMF=GMSC*(P+RTMDF1(IABS(ITEM))*RTMDFPI)
    ENDIF
  11 CONTINUE
  12 GMMF=SUM/(RTMDFPI*WSPREAD) ! convolution
  C 1/(sqrt(2*pi)*WSPREAD) normalizes Gaussian
  RETURN
END
```

## PyMinut

## SEAL-MINUT

```
include "Minut/Minut.h"
include "Minut/FunctionMinimum.h"
include "Minut/FCNBase.h"
include "Minut/FunctionCross.h"
include "Minut/MCross.h"
include "Minut/MinosError.h"

std::pair<double,double> MMinos::operator()(unsigned int par, unsigned
int mascal) const {
  MinosError merr = merr(par, mascal);
  return merr();
}

double MMinos::lower(unsigned int par, unsigned int mascal) const {
  MMinosParameterState upar = theMinimum.userState();
  double err = theMinimum.userState().error(par);
  MMinosCross sopt = local(par, mascal);
  double lower = sopt.isValid() ? -1.*err*(1.-sopt.value()) :
(sopt.atLimit() ? upar.parameter(par).lowerLimit() : upar.value(par));
  return lower;
}

// get_minut.cpp -1-hfs/cleo3/offline/re/current/other_sources/python
include/python_2_/-1/cdot/dm9/mccann/software/fit/minut/Minut-1_5_/
cdot/dm9/mccann/software/arc/minut/Minut-1_5_/src/*-o -shared -o
_minut.so
// get_minut.cpp -1-user/include/python_2_3_-1/root/src/Minut-1_5_2_/
root/src/Minut-1_5_2_/src/*-o -shared -o _minut.so

include <Python.h>
include "Minut/MinosParameters.h"
include "Minut/MMinograd.h"
include "Minut/MMinosize.h"
...
PyObject* dominos(PyObject *self, PyObject *args)
{
  // parameter list:
  PyObject *p_fcn; // objective function p_fcn
  int npar; // number of parameters in p_fcn
  double up; // 1 for chi^2, 0.5 for loglike
  PyObject *p_min; // the minimum you previously found

  if (!PyObject_IsTuple(args, "G1d51d00d0", &p_fcn, &npar, &p_min,
&mascal, &strategy, &p_double, &p_double, &p_arnum, &p_grad,
&p_checkrad)) {
    PyErr_SetString(PyExc_TypeError, "calling format must be: FCN(f),
npar(1), up(0), minimum(FunctionMinimum), mascal(1 or 0), strategy(1),
dlower(b), dupper(b), parnum(1), gradient(f or None), checkrad(b or
None)");
    return NULL;
  }
}
```



## Fitting script

```
# get_runs has been given a thorough look-over: it is correct (7 Oct
# 2005) (get_runs contains all corrections, from numbers of events to
# real, live cross-section.)

from minut import *
execfile("../home/mccann/antithesis/utilities/py2")
import gkwfit
import gkwftau
...
def dofittguss(h):
  def gauss(h, s, x): return exp(-.5**2/2./s**2)/(sqrt(2.*pi))/s
  def figauss(m,s):
    c = 0
    for x in h.data:
      c += -log(gauss(m, s, x))
    return c
  m = Minut(fitguss, start=[0., 1.], up=0.5)
  m.migrad()
  m.nstate(10, 11)
  err0 = m.minos_errors[0][1] + m.minos_errors[0][0]/2.
  err1 = m.minos_errors[1][1] + m.minos_errors[1][0]/2.
  return m.value[0], err0, m.value[1], err1, lambda c:
0.1*retmrae(m.h.data)*Gauss(m.value[0], m.value[1], x)
```

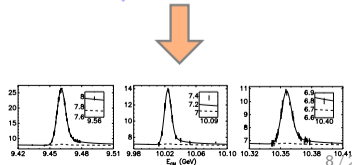
(pickle)

## GNU plotutils



## Plotting script

```
from math import *
import biggles, Numeric, cPickle as pickle
import gkwfit
import gkwftau
...
allth = pickle.load(file("../home/mccann/antithesis/novemberdata.p"))
ulrums = allth["ulrums"]
ulrums = allth["ulrums"]
ulrums = allth["ulrums"]
...
q = biggles.FramesPlot()
addData(q, [None], ulrums["High"], 0.)
addfunc(q, gkwftau, 10000, 10000.)
addfunc(q, shatunc, bkgrd, 10000., 10000., linestyle="dashed")
```



# Which got me involved in open source (PyMinuit is now “iminuit”)



[What's new?](#) | [Help](#) | [Directory](#) | [Sign in](#)



**pyminuit**

*Minuit numerical function minimization in Python*

[Project Home](#) [Downloads](#) [Wiki](#) [Issues](#) [Source](#)

## PyMinuit

*Minuit numerical function minimization in Python*

### Minuit

Minuit has been the standard package for minimizing general N-dimensional functions in high-energy physics since its introduction in 1972. It features a robust set of algorithms for optimizing the search, correcting mistakes, and measuring non-linear error bounds. It is the minimization engine used behind-the-scenes in most high-energy physics curve fitting applications.

**New:** more robust [installation instructions!](#)

### Python interface

PyMinuit is an extension module for Python that passes low-level Minuit functionality to Python functions. Interaction and data exploration is more user-friendly, in the sense that the user is protected from segmentation faults and index errors, parameters are referenced by their names, even in correlation matrices, and Python exceptions can be passed from the objective function during the minimization process. This extension module also makes it easier to calculate Minos errors and contour curves at an arbitrary number of sigmas from the minimum, and features a new N-dimensional scanning utility.

**License:** [GNU General Public License v2](#)

**Labels:** [python](#), [minuit](#), [optimization](#), [computation](#), [HEP](#), [math](#), [fitting](#), [physics](#)

**Featured Downloads:** [Show all](#)

↓ [Minuit-1.7.9.tar.gz](#)

↓ [pyminuit-1.0.2.tgz](#)

**Featured Wiki Pages:** [Show all](#)

[FunctionReference](#)

[GettingStartedGuide](#)

[HowToInstall](#)

**Links:**

- [Official Minuit homepage at CERN](#)
- [Minuit documentation through the ages](#)

**Project owners:** [Join project](#)  
[jpivarski](#)



But to be honest, Python isn't my idea of an ideal language.



But to be honest, Python isn't my idea of an ideal language.

I wish the scientific and data analysis community had adopted something more functional and statically typed. In particular, I wish there wasn't a dichotomy between statements and expressions. And `True == 1` is evil.



But to be honest, Python isn't my idea of an ideal language.

I wish the scientific and data analysis community had adopted something more functional and statically typed. In particular, I wish there wasn't a dichotomy between statements and expressions. And `True == 1` is evil.

HOWEVER, a heavy dash of consensus is worth a smattering of language features!



Did Python just arrive at the right time?

Or is it a better language for the job?





Did Python just arrive at the right time?

- ▶ Ruby, Lua numerical stacks were not ready before Python already had a foothold.
- ▶ Python was one of the first glue languages of the Linux/open source era.

Or is it a better language for the job?



## Did Python just arrive at the right time?

- ▶ Ruby, Lua numerical stacks were not ready before Python already had a foothold.
- ▶ Python was one of the first glue languages of the Linux/open source era.

## Or is it a better language for the job?

- ▶ Perl → Python
- ▶ “Tcl War”
- ▶ R → Python



## Did Python just arrive at the right time?

- ▶ Ruby, Lua numerical stacks were not ready before Python already had a foothold.
- ▶ Python was one of the first glue languages of the Linux/open source era.

## Or is it a better language for the job?

- ▶ Perl → Python
- ▶ “Tcl War”
- ▶ R → Python
- ▶ *An Empirical Investigation into Programming Language Syntax*, Andreas Stefik, Susanna Siebert.

# An Empirical Investigation into Programming Language Syntax



```

action Main
  number x = z(1, 100, 3)
end

action z(integer a, integer b,
integer c) returns number
  number d = 0.0
  number e = 0.0
  integer i = a
  repeat b - a times
    if i mod c = 0 then
      d = d + 1
    end
    else then
      e = e + 1
    end
    i = i + 1
  end
  if d > e then
    return d
  end
  else then
    return e
  end
end
    
```

(a) Quorum

```

$x = $z(1, 100, 3);

sub z{
  $a = $_[0];
  $b = $_[1];
  $c = $_[2];
  $d = 0.0;
  $e = 0.0;
  for ($i = $a; $i <= $b; $i++){
    if ($i % $c == 0) {
      $d = $d + 1;
    }
    else {
      $e = $e + 1;
    }
  }
  if ($d > $e) {
    $d;
  }
  else {
    $e;
  }
}
    
```

(b) Perl

```

^Main {
  ~ x \ z(1, 100, 3)
}

~ z(a \ b \ c) | ~ {
  ~ d \ 0.0
  ~ e \ 0.0
  @ i \ a
  # (b - a) {
  : i ; c ! 0 {
  : d \ d + 1
  }
  : e \ e + 1
  }
  i \ i + 1
  : d ~ e {
  ~ d
  }
  : e
  }
}
    
```

(c) Randomo

## Abstract

Recent studies in the literature have shown that syntax remains a significant barrier to novice computer science students in the field. While this syntax barrier is known to exist, whether and how it varies across programming languages has not been carefully investigated. For this article, we conducted four empirical studies on programming language syntax as part of a larger analysis into the, so called, programming language wars. We first present two surveys conducted with students on the intuitiveness of syntax, which we used to garner formative clues on what words and symbols might be easy for novices to understand. We followed up with two studies on the accuracy rates of novices using a total of six programming languages: Ruby, Java, Perl, Python, Randomo, and Quorum. Randomo was designed by randomly choosing some keywords from the ASCII table (a metaphorical placebo). To our surprise, we found that languages using a more traditional C-style syntax (both Perl and Java) did not afford accuracy rates significantly higher than a language with randomly generated keywords, but that languages which deviate (Quorum, Python, and Ruby) did. These results, including the specifics of syntax that are particularly problematic for novices, may help teachers of introductory programming courses in choosing appropriate first languages and in helping students to overcome the challenges they face with syntax. [\[Less\]](#)

DOI: 10.1145/2534973

```

public static void main(String[] args) {
  double x = z(1, 100, 3);
}

public static double z(int a, int b,
int c) {
  double d = 0.0;
  double e = 0.0;
  for(int i = a; i < b - a; i++){
    if(i % c == 0){
      d = d + 1;
    }else{
      e = e + 1;
    }
  }
  if(d > e){
    return d;
  }else{
    return e;
  }
}
    
```

(d) Java

```

def z(a, b, c):
  d = 0.0
  e = 0.0
  for i in range(a, b):
    if i % c == 0:
      d = d + 1
    else:
      e = e + 1
  if d > e:
    return d
  else:
    return e
x = z(1, 100, 3)
    
```

(e) Python

```

def z(a, b, c)
  d = 0.0
  e = 0.0
  for i in a..b-a
    if i % c == 0
      d = d + 1
    else
      e = e + 1
  end
  if d > e
    return d
  else
    return e
  end
end
x = z(1, 100, 3)
    
```

(f) Ruby

Table XXIII.

A table showing the cross-task average and standard deviation for each language and the Tukey HSD comparisons. Statistically significant differences are in bold.

	Tukey HSD						
	Mean	Std. Dev.	Python	Quorum	Perl	Java	Randomo
Ruby	.558	.243	0.994	0.940	0.165	<b>0.044</b>	<b>0.003</b>
Python	.528	.246		0.999	0.412	0.141	<b>0.011</b>
Quorum	.508	.232			0.655	0.292	<b>0.033</b>
Perl	.429	.189				0.987	0.573
Java	.396	.194					0.922
Randomo	.344	.207					

# An Empirical Investigation into Programming Language Syntax



```
action Main
  number x = z(1, 100, 3)
end
```

```
$x = &z(1, 100, 3);
```

```
^Main {
  x \ z(1, 100, 3)
}
```

```
action z(int c, integer b)
  integer c)
  number d =
  number e =
  integer i =
  repeat b -
  if i mod c
    d = d +
  end
  else then
    e = e +
  end
  i = i + 1
end
if d > e then
  return d
end
else then
  return e
end
end
```

(a)

```
public stat
double x =
}

public stat
int c)
double d =
double e =
for(int i =
  if(i % c ==
    d = d +
  }else{
    e = e +
  }
}
if(d > e){
  return d;
}else{
  return e;
}
}
```

(d) Java

```
x = z(1, 100, 3)
```

(e) Python

```
end
```

```
x = z(1, 100, 3)
```

(f) Ruby

## Abstract

Recent studies in the literature have shown that syntax remains a significant barrier to novice computer science students in the field. While this syntax barrier is known to exist, whether and how it varies across programming languages has not been carefully investigated. For this article, we conducted four empirical studies on programming language syntax as part of a larger analysis into the, so called, programming language wars. We first present two surveys conducted with students on the intuitiveness of syntax, which we used to garner formative clues on what words and symbols might be easy for novices to understand. We followed up with two studies on the accuracy rates of novices using a total of six programming languages: Ruby, Java, Perl, Python, Randomo, and Quorum. Randomo was designed by randomly choosing some keywords from the ASCII table (a metaphorical placebo). To our surprise, we found that languages using a more traditional C-style syntax (both Perl and Java) did not afford accuracy rates significantly higher than a language with randomly generated keywords, but that languages which deviate (Quorum, Python, and Ruby) did. These results, including the specifics of syntax that are particularly problematic for novices, may help teachers of introductory programming courses in choosing appropriate first languages and in helping students to overcome the challenges they face with syntax. (Less)

DOI: [10.1145/2534973](https://doi.org/10.1145/2534973)

Java	.396	.194								0.922
Randomo	.344	.207								

ch language and bold.

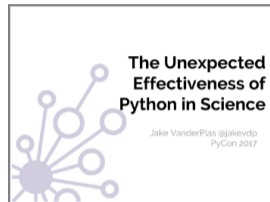
Java	Randomo
0.044	0.003
0.141	0.011
0.292	0.033
0.987	0.573



Python was good enough and first.

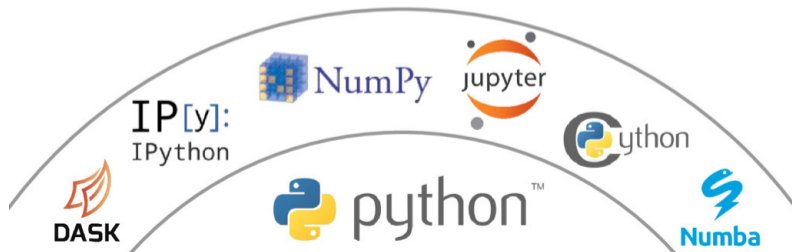
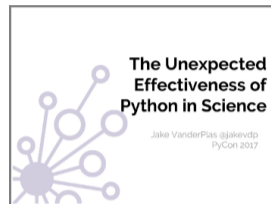


## Python's Scientific Stack





## Python's Scientific Stack

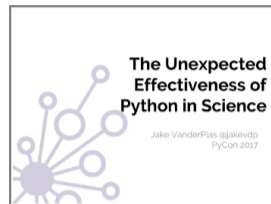
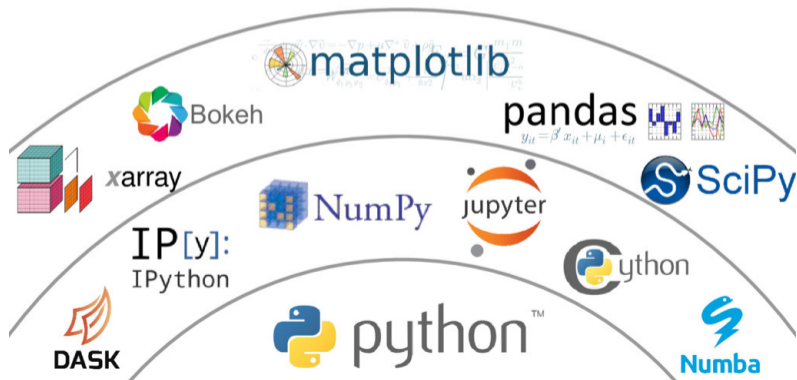




# More significant: what has grown around Python



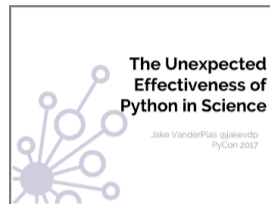
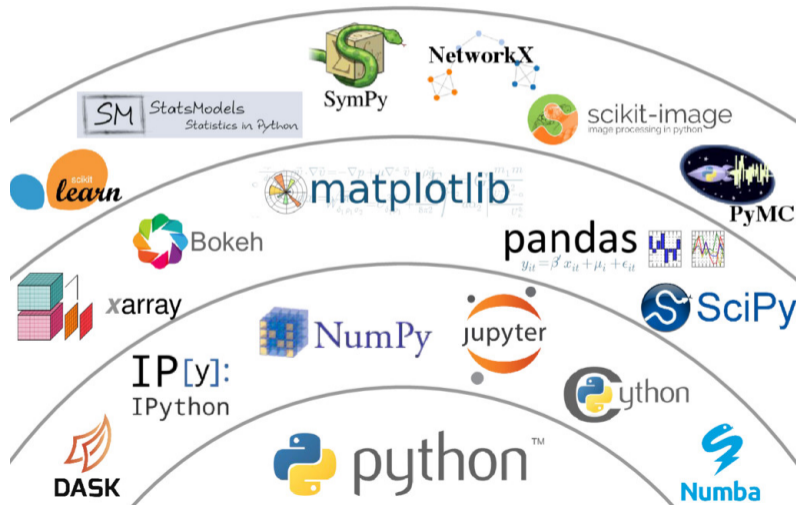
## Python's Scientific Stack



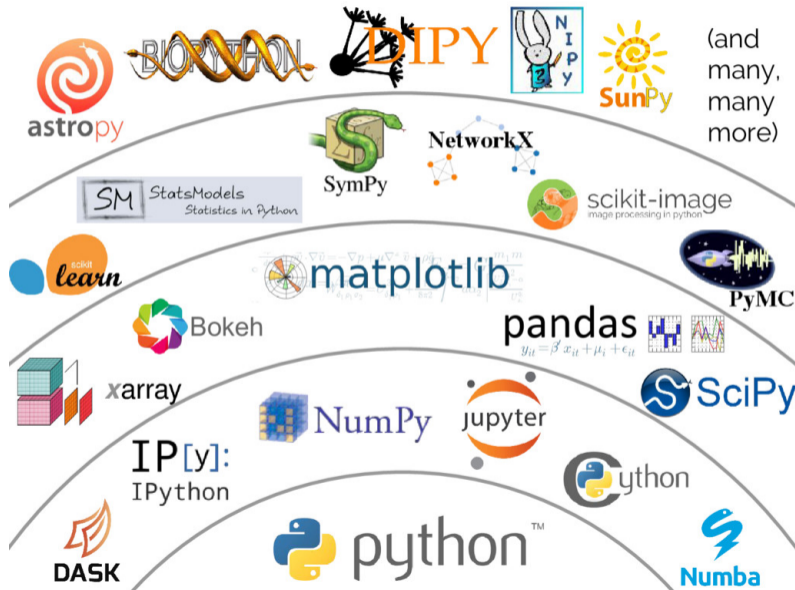
# More significant: what has grown around Python



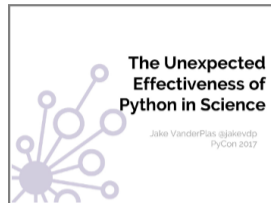
## Python's Scientific Stack



# More significant: what has grown around Python



(and many, many more)



# The key to ecosystem development was a common array library



- 1994 **Python** 1.0 released.
- 1995 First array package: **Numeric** (a.k.a. Numerical, Numerical Python, NumPy).
- 2001 Diverse scientific codebases merged into **SciPy**.
- 2003 **Matplotlib**
- 2003 Numeric was limited; **numarray** appeared as a competitor with more features (memory-mapped files, alignment, record arrays).
- 2005 Two packages were incompatible; could not integrate numarray-based code into SciPy. Travis Oliphant merged the codebases as **Numpy**.
- 2008 **Pandas**
- 2010 **Scikit-Learn**
- 2011 **AstroPy**
- 2012 **Anaconda**
- 2014 **Jupyter**
- 2015 **Keras**

# The key to ecosystem development was a common array library



- 1994 **Python** 1.0 released.
- 1995 First array package: **Numeric** (a.k.a. Numerical, Numerical Python, NumPy).
- 2001 Diverse scientific codebases merged into **SciPy**.
- 2003 **Matplotlib**
- 2003 Numeric was limited; **numarray** appeared as a competitor with more features (memory-mapped files, alignment, record arrays).
- 2005 Two packages were incompatible; could not integrate numarray-based code into SciPy. Travis Oliphant merged the codebases as **Numpy**.
- 2008 **Pandas**
- 2010 **Scikit-Learn**
- 2011 **AstroPy**
- 2012 **Anaconda**
- 2014 **Jupyter**
- 2015 **Keras**

The scientific Python ecosystem could have failed before it started if the Numeric-numarray split hadn't been resolved!

# Numpy is high-level, array-at-a-time math



```
>>> import numpy
>>> a = numpy.arange(12)
>>> a
array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11])
>>> a.shape = (3, 4)
>>> a
array([[ 0,  1,  2,  3],
       [ 4,  5,  6,  7],
       [ 8,  9, 10, 11]])
>>> a.sum(axis=0)
array([12, 15, 18, 21])
>>> a.min(axis=1)
array([0, 4, 8])
>>> a**2
array([[ 0,  1,  4,  9],
       [16, 25, 36, 49],
       [64, 81, 100, 121]])
>>> numpy.sqrt(a)
array([[0.          ,  1.          ,  1.41421356,  1.73205081],
       [2.          ,  2.23606798,  2.44948974,  2.64575131],
       [2.82842712,  3.          ,  3.16227766,  3.31662479]])
```



# Numpy is also a low-level way to poke raw bytes



```
>>> import numpy
>>> hello = b"Hello, world!"

>>> # Python strings are immutable
>>> hello[4:8] = "?????"
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
TypeError: 'bytes' object does not support item assignment

>>> # any buffer may be cast as an array
>>> a = numpy.frombuffer(hello, dtype=numpy.uint8)
>>> a
array([ 72, 101, 108, 108, 111,  44,  32, 119, 111, 114, 108, 100,  33],
      dtype=uint8)

>>> # and changed (with possibly disastrous consequences)
>>> a.flags.writeable = True
>>> a[4:8] = [69, 86, 73, 76]
>>> hello
b'HellEVILorld!'
```

# Numpy is also a low-level way to poke raw bytes



```
>>> import ctypes
```

```
>>> # any *pointer* may be cast as an array
```

```
>>> ptr = ctypes.cast(id(hello), ctypes.POINTER(ctypes.c_uint8))
```

```
>>> ptr.__array_interface__ = {
```

```
...     "version": 3,
```

```
...     "typestr": numpy.ctypeslib._dtype(ctypes(ptr.contents)).str,
```

```
...     "data": (ctypes.addressof(ptr.contents), False),
```

```
...     "shape": (100,)      # how many bytes do you want to access?
```

```
... }
```

```
>>> b = numpy.array(ptr, copy=False)
```

```
>>> b
```

```
array([  2,   0,   0,   0,   0,   0,   0,   0, 224, 136, 151,   0,   0,
         0,   0,   0,  13,   0,   0,   0,   0,   0,   0,  47,  49,
        110, 120,  37, 235,  98,  94,  72, 101, 108, 108,  69,  86,  73,
         76, 111, 114, 108, 100,  33,   0,   0,   0,   1,   0,   0,   0,
         0,   0,   0,   0, 224, 136, 151,   0,   0,   0,   0,   0,  12,
         0,   0,   0,   0,   0,   0,   0, 255, 255, 255, 255, 255, 255,
        255, 255,   0,   1,  18,   1,   3,   1,  22,   1,  13,   1,   5,
         1,   0, 105,   0,   0,   1,   0,   0,   0], dtype=uint8)
```

```
>>> "".join(map(chr, b[32:45]))
```

```
'HellEVIWorld!'
```



NumPy





Although you can write Python `for` loops over Numpy arrays, you don't reap the benefit unless you express your calculation in Numpy ufuncs (universal functions).

```
pz = numpy.empty(len(pt))
for i in range(len(pt)):
    pz[i] = pt[i]*numpy.sinh(eta[i])
```

$\mathcal{O}(N)$  Python bytecode instructions, type-checks, interpreter locks.

VS `pz = pt * numpy.sinh(eta)`

$\mathcal{O}(1)$  Python bytecode instructions, type-checks, interpreter locks.

$\mathcal{O}(N)$  statically typed, probably vectorized native bytecode operations on contiguous memory.



Although you can write Python `for` loops over Numpy arrays, you don't reap the benefit unless you express your calculation in Numpy ufuncs (universal functions).

```
pz = numpy.empty(len(pt))
for i in range(len(pt)):
    pz[i] = pt[i]*numpy.sinh(eta[i])
```

VS `pz = pt * numpy.sinh(eta)`

$\mathcal{O}(N)$  Python bytecode instructions, type-checks, interpreter locks.

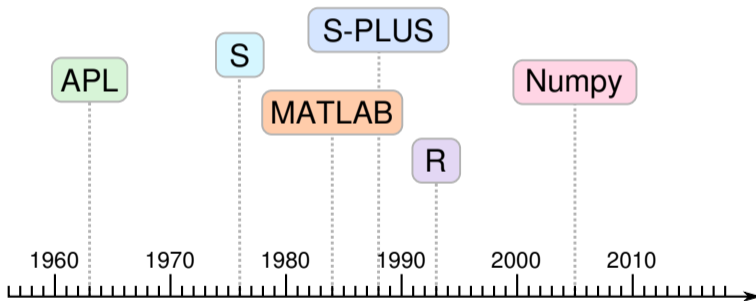
$\mathcal{O}(1)$  Python bytecode instructions, type-checks, interpreter locks.

$\mathcal{O}(N)$  statically typed, probably vectorized native bytecode operations on contiguous memory.

In other words, a Single (Python) Instruction on Multiple Data.

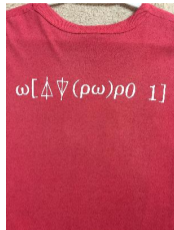


**APL**, “A Programming Language” introduced the idea of single commands having sweeping effects across large arrays.



All members of the APL family are intended for interactive data analysis. Numpy, however, is a library in a general-purpose language, not a language in itself.

APL pioneered conciseness;  
discovered the mistake of being too concise.



Conway's Game of Life was one line of code:

```
life ← {↑1 ω∇.∧3 4 = +/,-1 0 1○.Θ-1 0 1○.Φ ⊂ ω}
```

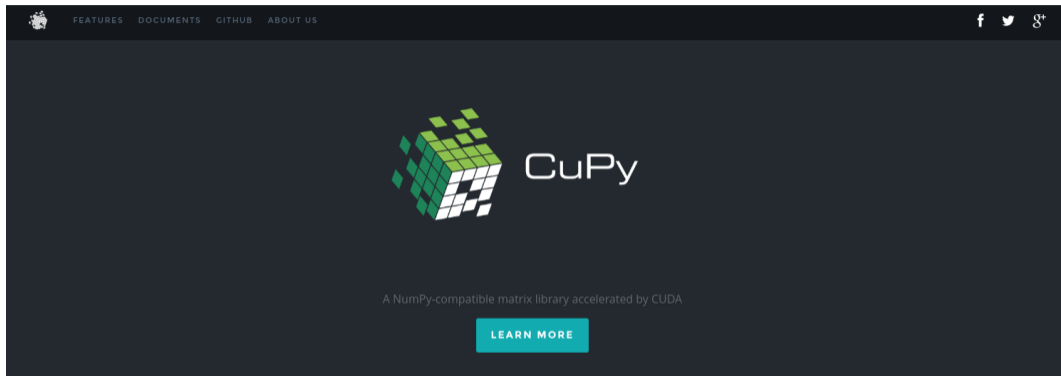
“Map” was implicit, “reduce” was a slash, functions were symbols. For example:

APL	Numpy
$m \leftarrow +/(3 + \iota 4)$	<code>m = (numpy.arange(4) + 3).sum()</code>



As an array abstraction, Numpy presents a high-level way for users to think about vectorization.

Vectorization is key to using GPUs and modern CPUs efficiently.

A screenshot of the CuPy website. The header has a small logo on the left and navigation links for 'FEATURES', 'DOCUMENTS', 'GITHUB', and 'ABOUT US' in the center. On the right are social media icons for Facebook, Twitter, and Google+. The main content area features the CuPy logo, which is a 3D cube of green and white squares, followed by the text 'CuPy'. Below this is the tagline 'A NumPy-compatible matrix library accelerated by CUDA' and a teal button with the text 'LEARN MORE'.

## HIGH PERFORMANCE WITH CUDA

CuPy is an open-source matrix library accelerated with NVIDIA CUDA. It also uses CUDA-related libraries including cuBLAS, cuDNN, cuRand, cuSolver, cuSPARSE, cuFFT and NCCL to make full use of the GPU architecture.



## PROJECTS

QuantStack developers contribute to a number of open-source projects, including jupyter, xtensor, bqplot, conda-forge and many others.

### High-Performance-Computing



C++ tensor algebra library



BLAS extension to xtensor



C++ wrappers for SIMD intrinsics  
and optimized math  
implementations



# Can Numpy deal with HEP data?



$N$ -dimensional arrays of values are great for image processing, signal processing, PDEs, but HEP data have nested structure with variable-length contents:

muons		
$p_T$	phi	eta
31.1	-0.481	0.882
9.76	-1.24	0.924
8.18	-0.119	0.923

VS

mu1	mu1	mu1	mu2	mu2	mu2
$p_T$	phi	eta	$p_T$	phi	eta
31.1	-0.481	0.882	9.76	-0.124	0.924
5.27	1.246	-0.991	n/a	n/a	n/a
4.72	-0.207	0.953	n/a	n/a	n/a
8.59	-1.754	-0.264	8.714	0.185	0.629



# Can Numpy deal with HEP data?



$N$ -dimensional arrays of values are great for image processing, signal processing, PDEs, but HEP data have nested structure with variable-length contents:

muons		
$p_T$	phi	eta
31.1	-0.481	0.882
9.76	-0.124	0.924
8.18	-0.119	0.923

VS

mu1	mu1	mu1	mu2	mu2	mu2
$p_T$	phi	eta	$p_T$	phi	eta
31.1	-0.481	0.882	9.76	-0.124	0.924
5.27	1.246	-0.991	n/a	n/a	n/a
4.72	-0.207	0.953	n/a	n/a	n/a
8.59	-1.754	-0.264	8.714	0.185	0.629

```
[ [Muon (31.1, -0.481, 0.882), Muon (9.76, -0.124, 0.924), Muon (8.18, -0.119, 0.923) ],  
  [Muon (5.27, 1.246, -0.991) ],  
  [Muon (4.72, -0.207, 0.953) ],  
  [Muon (8.59, -1.754, -0.264), Muon (8.714, 0.185, 0.629) ], ... ]
```

# Can Numpy deal with HEP data?



$N$ -dimensional arrays of values are great for image processing, signal processing, PDEs, but HEP data have nested structure with variable-length contents:

muons		
$p_T$	phi	eta
31.1	-0.481	0.882
9.76	-0.124	0.924
8.18	-0.119	0.923

VS

mu1	mu1	mu1	mu2	mu2	mu2
$p_T$	phi	eta	$p_T$	phi	eta
31.1	-0.481	0.882	9.76	-0.124	0.924
5.27	1.246	-0.991	n/a	n/a	n/a
4.72	-0.207	0.953	n/a	n/a	n/a
8.59	-1.754	-0.264	8.714	0.185	0.629

```
[ [Muon (31.1, -0.481, 0.882), Muon (9.76, -0.124, 0.924), Muon (8.18, -0.119, 0.923) ],  
  [Muon (5.27, 1.246, -0.991) ],  
  [Muon (4.72, -0.207, 0.953) ],  
  [Muon (8.59, -1.754, -0.264), Muon (8.714, 0.185, 0.629) ], ... ]
```

I've been working on ways to represent arbitrary physics objects as vectorizable arrays:

offsets	0,	3,	4,	5,	7		
$p_T$	31.1,	9.76,	8.18,	5.27,	4.72,	8.59,	8.714
phi	-0.481,	-0.123,	-0.119,	1.246,	-0.207,	-1.754,	0.185
eta	0.882,	0.924,	0.923,	-0.991,	0.953,	-0.264,	0.629

# Extending Numpy's concept of "broadcasting"



In Numpy, arithmetic is applied element-wise; scalars are duplicated to fit:

```
>>> MET = numpy.array([(10.2, -0.480), (34.1, 1.251), (26.5, -0.22), (19.0, -1.75)],  
... dtype=[("E", float), ("phi", float)])  
>>> MET["E"] * 1.1  
array([11.22, 37.51, 29.15, 20.9 ])
```

# Extending Numpy's concept of "broadcasting"



In Numpy, arithmetic is applied element-wise; scalars are duplicated to fit:

```
>>> MET = numpy.array([(10.2, -0.480), (34.1, 1.251), (26.5, -0.22), (19.0, -1.75)],
... dtype=[("E", float), ("phi", float)])
>>> MET["E"] * 1.1
array([11.22, 37.51, 29.15, 20.9 ])
```

We could also apply operations element-wise if their nested structure is the same:

```
>>> muons = awkward.fromiter(
...     [[Muon(31.1, -0.481, 0.882), Muon(9.76, -0.124, 0.924), Muon(8.18, -0.119, 0.923)],
...     [Muon(5.27, 1.246, -0.991)], [Muon(4.72, -0.207, 0.953)],
...     [Muon(8.59, -1.754, -0.264), Muon(8.714, 0.185, 0.629)]]
>>> muons["pt"]
<JaggedArray [[31.1  9.76  8.18] [5.27] [4.72] [8.59  8.714]] at 7d5022ab3f90>
>>> muons["pt"] * numpy.sinh(muons["eta"])
<JaggedArray [[31.12755703 10.35740718  8.66877254] [-6.12037182] [5.21063432]
               [-2.29419425  5.84974849]] at 7d50223d77d0>
```

# Extending Numpy's concept of "broadcasting"



In Numpy, arithmetic is applied element-wise; scalars are duplicated to fit:

```
>>> MET = numpy.array([(10.2, -0.480), (34.1, 1.251), (26.5, -0.22), (19.0, -1.75)],
... dtype=[("E", float), ("phi", float)])
>>> MET["E"] * 1.1
array([11.22, 37.51, 29.15, 20.9 ])
```

We could also apply operations element-wise if their nested structure is the same:

```
>>> muons = awkward.fromiter(
...     [[Muon(31.1, -0.481, 0.882), Muon(9.76, -0.124, 0.924), Muon(8.18, -0.119, 0.923)],
...     [Muon(5.27, 1.246, -0.991)], [Muon(4.72, -0.207, 0.953)],
...     [Muon(8.59, -1.754, -0.264), Muon(8.714, 0.185, 0.629)]]
>>> muons["pt"]
<JaggedArray [[31.1  9.76  8.18] [5.27] [4.72] [8.59  8.714]] at 7d5022ab3f90>
>>> muons["pt"] * numpy.sinh(muons["eta"])
<JaggedArray [[31.12755703 10.35740718  8.66877254] [-6.12037182] [5.21063432]
               [-2.29419425  5.84974849]] at 7d50223d77d0>
```

One-per-event scalars could also be broadcast down to multi-per-event jagged arrays:

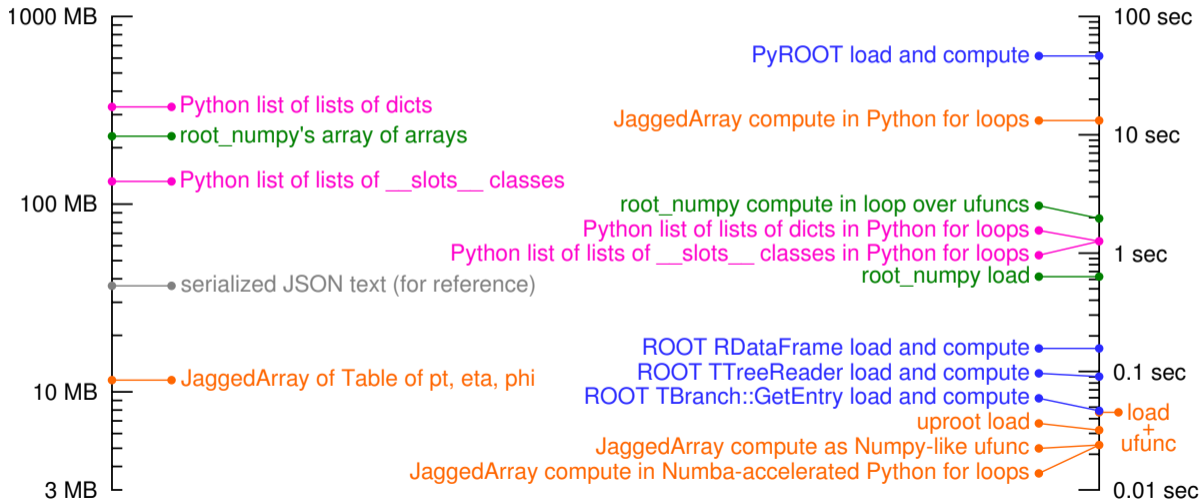
```
>>> muons["phi"] - MET["phi"]
<JaggedArray [[-0.001  0.356  0.361] [-0.005] [0.013] [-0.004  1.935]] at 7d50223d7790>
```

# Big performance gain, even without writing C code



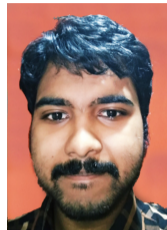
RAM memory

time to complete





Jaydeep Nandi, a CERN/HSF Google Summer of Code student, is investigating vectorized algorithms to replace for-loop manipulations.



Explode to event-wise pairs (using only hardware-SIMD operations):

```
>>> pairs = muons.pairs()
>>> pairs
<JaggedArray [[<Pair 0> <Pair 1> <Pair 2>] [] [] [<Pair 3>]]>
>>> pairs[0][0].tolist()
{"_0": {"pt": 31.1, "phi": -0.481, "eta": 0.882, "pz": 31.128},
 "_1": {"pt": 9.76, "phi": -0.123, "eta": 0.924, "pz": 10.358}}
```

Now we can do such things as compute invariant masses without loops:

```
>>> pt1, eta1, phi1 = pairs["_0"]["pt"], pairs["_0"]["eta"], pairs["_0"]["phi"]
>>> pt2, eta2, phi2 = pairs["_1"]["pt"], pairs["_1"]["eta"], pairs["_1"]["phi"]
>>> mass = numpy.sqrt(2*pt1*pt2*(numpy.cosh(eta1 - eta2) - numpy.cos(phi1 - phi2)))
```

Also considering problems like “minimize per event” and “match gen/reco candidates” in hardware-SIMD operations, exposed in a Numpy-like interface.



- ▶ Python is a popular language, even in sciences where performance is critical.
- ▶ It has good features for readability and is easy to learn, particularly by scientists whose primary interest is not programming.
- ▶ It came early enough to build up a numerical ecosystem.
- ▶ The community was almost fractured by the Numeric-numarray split.
- ▶ Modern Numpy is a fixed API that can be swapped out for GPU and SIMD implementations.
- ▶ Numpy's existing API is flexible enough to represent complex data structures.
- ▶ May be a road toward vectorized HEP analysis, in the guise of interactive array primitives.