

History and Adoption of Programming Languages in NHEP

Jim Pivarski

Princeton University – IRIS-HEP

February 8, 2202



This talk is a historical overview of widely used programming languages in NHEP, focusing on the motivations for each language's adoption.



This talk is a historical overview of widely used programming languages in NHEP, focusing on the motivations for each language's adoption.

There have always been physicists at the edge, trying out new languages, but most physicists only use one or two. The field as a whole changes slowly.



This talk is a historical overview of widely used programming languages in NHEP, focusing on the motivations for each language's adoption.

There have always been physicists at the edge, trying out new languages, but most physicists only use one or two. The field as a whole changes slowly.

Thesis: (1) change motivated more by “pain points” than incremental benefits,

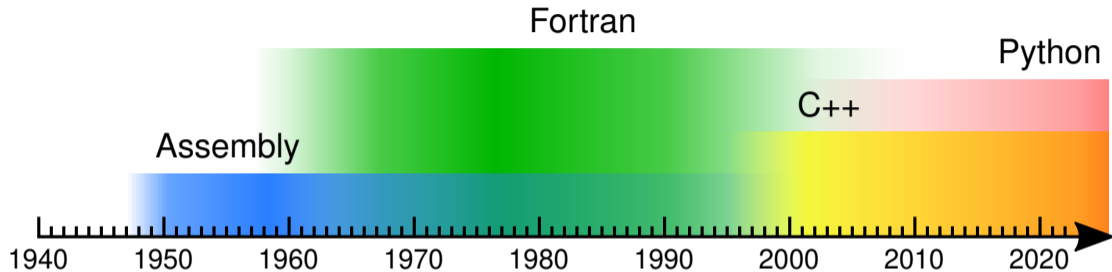


This talk is a historical overview of widely used programming languages in NHEP, focusing on the motivations for each language's adoption.

There have always been physicists at the edge, trying out new languages, but most physicists only use one or two. The field as a whole changes slowly.

Thesis: (1) change motivated more by “pain points” than incremental benefits,
(2) though each of these transitions happened in a unique way.

Three major transitions (so far)



Adoption of

Fortran: immediate; for syntax and portability; no infrastructure to replace

C++: long overdue; for data structures; replaced infrastructure in a burst

Python: slowly overtook its alternatives; for interactivity; different niche

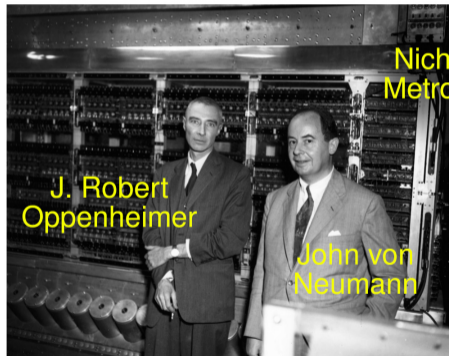


Part 1: Fortran

NHEP was an early adopter of digital computers



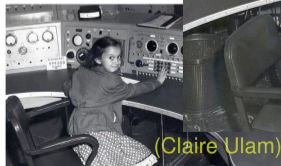
One of the very first applications was Monte Carlo (neutron transport).



Nicholas Metropolis



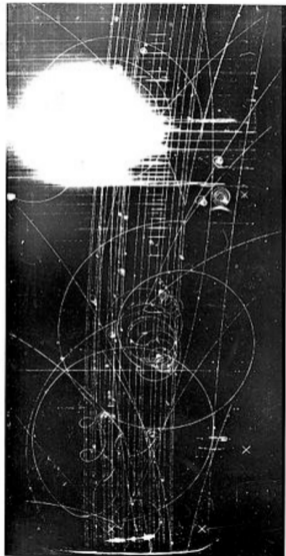
Stanislaw Ulam



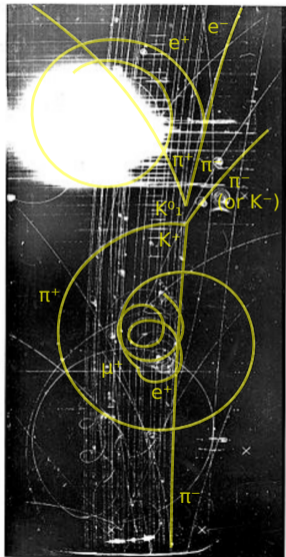
And so was data analysis

Luis Alvarez's group at the Bevatron: \$2M bubble chamber, \$0.2M IBM 650.





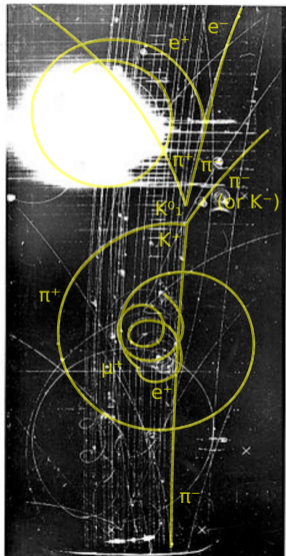
Unlabeled photos
come out of the
detector.



Unlabeled photos
come out of the
detector.

Labeling them
turns them into
quantities to
compute.

The problem was the same then as it is now



Unlabeled photos come out of the detector.

Labeling them turns them into quantities to compute.

THE MORE EVENTS THE BETTER!!!

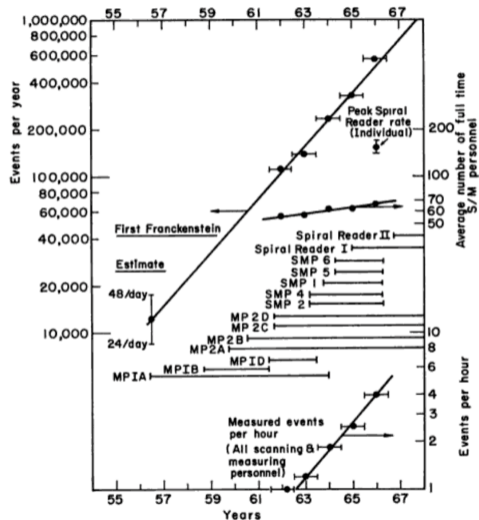
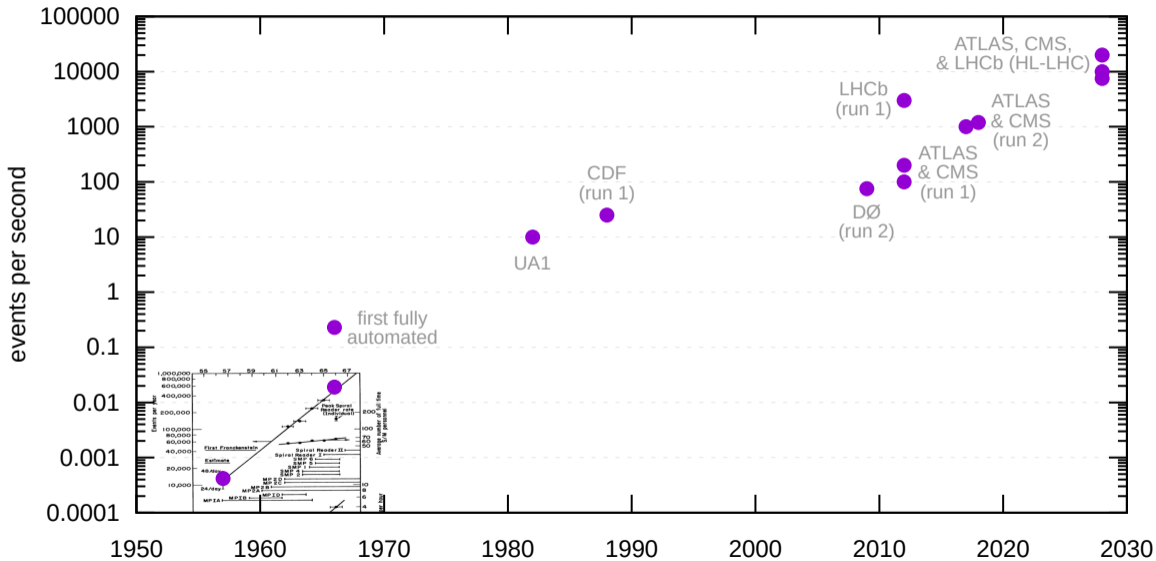


Fig. 9. Measuring Rates.

MUB 12506

The problem was the same then as it is now





So they invented special input devices to streamline data-entry.

Identifying tracks was beyond the capabilities of software

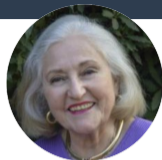


Madeleine (née Goldstein)
Isenberg, UCLA class of '65

"We scanners would review each frame of film, and per the brief instructions we had been given, looked for any 'unusual activity.'

"The scanner had to use both hands, a joystick in each, and turn them clockwise or anti-clockwise, to align a double crosshair cursor at several sequential positions on a track."

Identifying tracks was beyond the capabilities of software

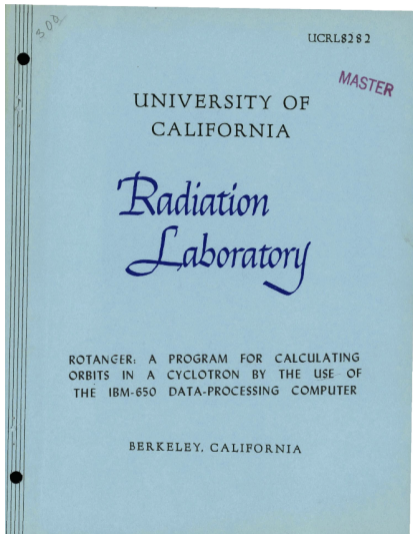


Madeleine (née Goldstein)
Isenberg, UCLA class of '65

“A quick but firm tap on the foot-pedal punched the coordinate values onto an IBM card that had been fed into the keypunch machine.

“The precious stack of IBM cards were passed to the physicists, who would then process the data in the existing IBM processors, using software that would calculate the best fit for these coordinates, and thereby mathematically simulate the curvature of the track.”

At first, the software was written in assembly (example from 1958)



At first, the software was written in assembly (example from 1958)



500

UCRL-8282

constants in the expansion of the magnetic field are changed, and others. An interpolation routine to fix the regenerator position is included.

The program is in fixed point, and values are given with eight decimals. This sets an upper bound on the magnitude of numbers, which can be only 99.99999999. An overflow occurs a little before this value, about 10% in the computation of v^2 ; however, a rescaling could prevent this. The interval of each advance in the numerical integration h is variable and has been used at 0.1 radian for most orbits. The error ϵ between successive values of the variables as computed in the iterative integration procedure is also adjustable and has been used at .00000020 for most cases. The similar terms ϵ_0 and h_0 are used in the starting program and are set at one-half the values for ϵ and h .

At the end of each integration step, the square of the total velocity is computed. The constancy of this term is used to judge the accumulation of errors in the integration of the three equations.

COMPUTATIONS
Equations of Motion

The three-component equations obtained for the Lorentz equation are:

$$\frac{d^2 r}{dt^2} - r \left(\frac{d\theta}{dt} \right)^2 = - \frac{e}{m} r \frac{d\theta}{dt} B_z$$
$$r \frac{d^2 \theta}{dt^2} + 2 \frac{dr}{dt} \frac{d\theta}{dt} = - \frac{e}{m} \left(\frac{ds}{dt} B_r - \frac{dr}{dt} B_z \right)$$

and

$$\frac{d^2 s}{dt^2} = - \frac{c}{m} r \frac{d\theta}{dt} B_r$$

where we have assumed an azimuthally uniform field, i. e., $B_\theta = 0$, and note that $B_z = B_z(r, z)$ and $B_r = B_r(r, z)$.

The synchronous radius R corresponds to the beginning radius of the regenerator, and ρ is the orbit departure from this radius. The radius from the center of the cyclotron is $r = R + \rho$. The magnetic field values are normalized by dividing them by the field on the median plane at the synchronous orbit $B_z(R, 0)$.

Our equations then become:

$$\frac{d^2 \rho}{dt^2} - (R + \rho) \left(\frac{d\theta}{dt} \right)^2 = - \frac{e B_z(R, 0)}{m} (R + \rho) \frac{d\theta}{dt} \frac{B_r(r, z)}{B_z(R, 0)}$$

ROTANGER:
ORBITS IN
THE IBM-6

BE

At first, the software was written in assembly (example from 1958)



500

constants in the exponential interpolation routine

The program is This sets an upper bound of 99.99999999. An over-computation of v^2 ; h each advance in the t at 0.1 radian for most variables as computed and has been used at h_0 are used in the steps ϵ and h .

At the end of ϵ computed. The constant errors in the integration

The three-com

$$r \frac{d^2 \theta}{dt^2} + 2$$

and

ROTANGER: ORBITS IN THE IBM-6

where we have assumed note that $B_s = B_s(r, z)$

The synchronous regenerator, and p is the center of the cycle normalized by dividing synchronous orbit B_s

Our equations

$$\frac{d^2 \rho}{dt^2} - (R + \rho)$$

- 16 - UCRL-8282

Appendix 6.

Block Diagram of Rotanger

```

graph TD
    Start([Read Initial Conditions]) --> SubStart[Starting Part of Integration Subroutine]
    SubStart --> Deriv[Compute second Derivatives including B_s and B_r  
Compute third Derivative]
    Deriv --> IntSub[Integration Subroutine]
    IntSub --> EpsilonTest{epsilon test}
    EpsilonTest -- No --> Print1[Print]
    EpsilonTest -- Yes --> VTest{v test}
    VTest -- No --> Stop([Stop])
    VTest -- Yes --> RTest{r test}
    RTest -- No --> Stop
    RTest -- Yes --> ThetaTest{theta_1 test theta > theta_1}
    ThetaTest -- No --> Print2[Print]
    ThetaTest -- Yes --> Interp[Interpolate to theta_1  
Compute delta(rho) and delta(theta)  
Obtain new initial conditions]
    Interp --> Print3[Print]
    Print1 --> Advance[Advance tau by h]
    Print2 --> Advance
    Print3 --> Advance
    Advance --> SubStart
    
```

At first, the software was written in assembly (example from 1958)



ROTANGER:
 ORBITS IN
 THE IBM-6

300

constants in the expansion interpolation routine

The program is
 This sets an upper bound
 99.99999999. An over
 computation of v^2 ; h
 each advance in the t
 at 0.1 radian for most
 variables as compute
 and has been used at
 h_0 are used in the sta
 e and h .

At the end of e
 computed. The cons
 errors in the integra

The three-com

$$r = \frac{d^2 p}{dt^2} + 2$$

and

where we have assum
 note that $B_a = B_a(r, z)$

The synchronou
 regenerator, and p is
 the center of the cycl
 normalized by divid
 synchronous orbit B

Our equations

$$\frac{d^2 p}{dt^2} = -(R + p)$$

Block Diagram

```

        graph TD
            Start(( )) --> C1[Compute including]
            C1 --> C2[Compute]
            C2 --> C3[Compute including]
            C3 --> C4[Compute]
            C4 --> I[Interpolate]
            I --> C5[Compute]
            C5 --> O[Obtain]
            O --> P[Print]
            P --> End(( ))
    
```

-18-

UCRL-8282

Subprogram for calculating the z and θ coordinates of motion,
 including the computation of magnetic field components.

Label	Operation	Storage Location	Comment	Remarks
0313 0403123750	STZ			
1780 0010613781	XAR			" h "
1781 1181002750	SR			" $h \cdot R \cdot p$ "
1782 0100002753	STZ	0000		" p "
1783 0100040078	STZ			" z "
0088 1600350088	SL			" $h \cdot R$ "
0089 4000670484	SRZB			
0027 0918300000	LD			" z "
0028 2401000070	SR	0010		" z "
0070 0000090304	LD			" z "
0304 0401000310	STZ	0018		" z "
0311 0000340304	LD			" z "
0324 0401070400	STZ	0019		" z "
0400 0000300402	LD			" z "
0401 0401000400	STZ	0018		" z "
0404 0000300400	LD			" z "
0405 0401000400	STZ	100		" z "
0406 0000100407	LD			" z "
0407 0401000510	STZ	100		" z "
0500 0017000301	LD			" z "
0501 0401070500	STZ	100		" z "
0502 0010030503	LD			" z "
0503 0401000400	STZ	100		" z "
0403 0010030704	LD			" z "
1784 1000001700	STZ	0012		" z "
1785 1000001700	STZ	0012		" z "
1786 1000001700	STZ	0012		" z "
1787 1000001700	STZ	0012		" z "
1788 1000001700	STZ	0012		" z "
1789 1000001700	STZ	0012		" z "
1790 1000001700	STZ	0012		" z "
1791 1000001700	STZ	0012		" z "
1792 1000001700	STZ	0012		" z "
1793 1000001700	STZ	0012		" z "
1794 1000001700	STZ	0012		" z "
1795 1000001700	STZ	0012		" z "
1796 1000001700	STZ	0012		" z "
1797 1000001700	STZ	0012		" z "
1798 1000001700	STZ	0012		" z "
1799 1000001700	STZ	0012		" z "
1800 1000001700	STZ	0012		" z "
1801 1000001700	STZ	0012		" z "
1802 1000001700	STZ	0012		" z "
1803 1000001700	STZ	0012		" z "
1804 1000001700	STZ	0012		" z "
1805 1000001700	STZ	0012		" z "
1806 1000001700	STZ	0012		" z "
1807 1000001700	STZ	0012		" z "
1808 1000001700	STZ	0012		" z "
1809 1000001700	STZ	0012		" z "
1810 1000001700	STZ	0012		" z "
1811 1000001700	STZ	0012		" z "
1812 1000001700	STZ	0012		" z "
1813 1000001700	STZ	0012		" z "
1814 1000001700	STZ	0012		" z "
1815 1000001700	STZ	0012		" z "
1816 1000001700	STZ	0012		" z "
1817 1000001700	STZ	0012		" z "
1818 1000001700	STZ	0012		" z "
1819 1000001700	STZ	0012		" z "
1820 1000001700	STZ	0012		" z "
1821 1000001700	STZ	0012		" z "
1822 1000001700	STZ	0012		" z "
1823 1000001700	STZ	0012		" z "
1824 1000001700	STZ	0012		" z "
1825 1000001700	STZ	0012		" z "
1826 1000001700	STZ	0012		" z "
1827 1000001700	STZ	0012		" z "
1828 1000001700	STZ	0012		" z "
1829 1000001700	STZ	0012		" z "
1830 1000001700	STZ	0012		" z "
1831 1000001700	STZ	0012		" z "
1832 1000001700	STZ	0012		" z "
1833 1000001700	STZ	0012		" z "
1834 1000001700	STZ	0012		" z "
1835 1000001700	STZ	0012		" z "
1836 1000001700	STZ	0012		" z "
1837 1000001700	STZ	0012		" z "
1838 1000001700	STZ	0012		" z "
1839 1000001700	STZ	0012		" z "
1840 1000001700	STZ	0012		" z "
1841 1000001700	STZ	0012		" z "
1842 1000001700	STZ	0012		" z "
1843 1000001700	STZ	0012		" z "
1844 1000001700	STZ	0012		" z "
1845 1000001700	STZ	0012		" z "
1846 1000001700	STZ	0012		" z "
1847 1000001700	STZ	0012		" z "
1848 1000001700	STZ	0012		" z "
1849 1000001700	STZ	0012		" z "
1850 1000001700	STZ	0012		" z "
1851 1000001700	STZ	0012		" z "
1852 1000001700	STZ	0012		" z "
1853 1000001700	STZ	0012		" z "
1854 1000001700	STZ	0012		" z "
1855 1000001700	STZ	0012		" z "
1856 1000001700	STZ	0012		" z "
1857 1000001700	STZ	0012		" z "
1858 1000001700	STZ	0012		" z "
1859 1000001700	STZ	0012		" z "
1860 1000001700	STZ	0012		" z "
1861 1000001700	STZ	0012		" z "
1862 1000001700	STZ	0012		" z "
1863 1000001700	STZ	0012		" z "
1864 1000001700	STZ	0012		" z "
1865 1000001700	STZ	0012		" z "
1866 1000001700	STZ	0012		" z "
1867 1000001700	STZ	0012		" z "
1868 1000001700	STZ	0012		" z "
1869 1000001700	STZ	0012		" z "
1870 1000001700	STZ	0012		" z "
1871 1000001700	STZ	0012		" z "
1872 1000001700	STZ	0012		" z "
1873 1000001700	STZ	0012		" z "
1874 1000001700	STZ	0012		" z "

Remarks

Beats test to check proper
 field distribution in the orbit
 in first experiment. h is
 value at which trajectories
 are orthogonal. Printed
 value used was $h = h_0$.

For coefficient in region $z > h_0$.

Revised in program of
 calculation of field
 components.

At first, the software was written in assembly (example from 1958)



I
C

ROTANGER:
ORBITS IN
THE IBM-6

B E

constants in the expansion interpolation routine

The program is
This sets an upper bound
99.99999999. An over
computation of \sqrt{r} ; h
each advance in the r
at 0.1 radian for most
variables as compute
and has been used at
h_g are used in the sta
e and h.

At the end of e
computed. The cons
errors in the integra

The three-com

$$r \frac{d^2 \theta}{dt^2} + 2$$

and

where we have assum
note that $B_s = B_s(r, z$

The synchronou
regenerator, and p is
the center of the cycl
normalized by divid
synchronous orbit B_s

Our equations

$$\frac{d^2 \theta}{dt^2} - (R + p)$$

Block Diagram

```

    graph TD
      A[Compute including] --> B[Compute]
      B --> C[Compute including]
      C --> D[Interpolate]
      D --> E[Compute]
      E --> F[Obtain m]
      F --> G[Print]
    
```

Address	Instruction	Operation	Storage Location	Element	Remarks
1775	3100001776	000			
1776	4000001777	RAT			
1777	1910001778	000			$X^2 + \frac{1}{2} Y^2 + Z^2$
1778	3100001779	000			
1779	2000001780	07L	000		$Z^2 + \frac{1}{2} Y^2 + X^2$
1780	4000001781	000			$+ \frac{1}{2} Y^2 + X^2 + Z^2$
1781	1100001782	000			R_2 / R_0
1782	5100001783	000	000		$R_2 / R_0 (Z^2 / R_0, R_0, R_0)$
1783	4000001784	000			End computation of axial component of magnetic field.
1784	1901001785	000			
1785	6000001805	000			
1805	1900001786	000			
1786	3100001787	000			
1787	0000001788	000	000L		
1788	4001017189	000			
1789	1901017190	000			
1790	1000001791	000			
1791	4000001792	000			
1792	3100001793	000			
1793	2100001794	000			
1794	1001017195	000			
1795	6000001796	000			
1796	1910171797	000			
1797	3100001798	000			
1798	2000001799	000	000		
1799	4010118000	000			
1800	1901001801	000			
1801	0000001802	000	000L		
1802	4001018003	000			
1803	4400001804	0000			
1804	3100001805	000			
1805	1601071807	000			
1807	4000001808	000			
1808	3100001809	000			
1809	4400001810	0000			
1810	1601001811	000			
1811	2000071812	000	000L		
1812	4010071813	000			
1813	1910171814	000			
1814	3100001815	000			
1815	0000001816	000			
1816	4000001817	000			
1817	1900071818	000			
1818	3100001819	000			
1819	1000001820	000			
1820	0000121840	000			

11 / 50

At first, the software was written in assembly (example from 1958)



300

T

C

B

constants in the exponential interpolation routine

The program is This sets an upper bound of 99.99999999. An over-computation of \sqrt{r} is each advance in the t at 0.1 radian for most variables as computed and has been used at h_0 are used in the steps and h .

At the end of each computed. The constant errors in the integration

The three-com

$$r = \frac{d^2 \theta}{dt^2} + 2$$

and

where we have assumed that $B_s = B_0(r, z)$

The synchronous regenerator, and p is the center of the cycle normalized by dividing synchronous orbit B_0

Our equations

$$\frac{d^2 \theta}{dt^2} = (R + p)$$

Block Diagram

```

    graph TD
      Start(( )) --> Loop[Compute including Compute]
      Loop --> Interpolate[Interpolate]
      Interpolate --> Compute[Compute]
      Compute --> Obtain[Obtain new]
      Obtain --> Print[Print]
      Print --> Loop
    
```

Instruction	Address	Description
0313	0403123750	1775 3100001776
1780	0410413761	1774 0000011777
1781	1111002752	1775 1910471770
1782	0100023753	1778 3100001778
1783	0210410074	1779 2000031779
0088	1640350099	1780 0000011780
0099	4000270484	1781 1100031781
0027	0918300068	1782 3100001782
0040	2401000076	1783 0001001783
0070	0900070004	1784 1901001784
0304	0401000010	1785 0000001805
0311	0000340304	1805 1900001786
0324	0401070400	1786 3100001787
0400	0000300400	1787 0000001788
0440	0401000000	1788 000101789
0484	0000030400	1789 190101790
0485	0401000000	1790 1800001791
0486	0000100487	1791 0000001792
0487	0401000010	1792 1900001793
0510	0017000010	1793 3100001794
0521	0401070500	1794 190101795
0522	0018000010	1795 0000001796
0523	0401000000	1796 1910471797
0443	0010000010	1797 3100001798
1784	1900001799	1798 2000001799
1783	3100001800	1799 0010411800
1736	1901071707	1800 1901001801
1747	0018001802	1807 3100001801
1750	1900001709	1801 0000001802
1759	3100001803	1800 001041803
1760	1901001801	1803 0400001804
1761	0000001804	1809 3100001808
1762	1900001705	1806 1601071807
1763	3100001805	1807 0000001808
1764	1901001809	1808 3100001809
1765	0000011766	1809 0400001810
1766	0018011767	1810 1601001811
1767	1901001768	1811 0000071810
1768	3100001812	1810 0010471813
1801	0000001802	1813 1910471804
1809	0001001800	1824 3100001825
1770	3000001771	1805 0000001814
1771	0400001806	1814 1910471813
1804	0000001770	1815 3100001816
1778	1901071773	1816 0000001817
1773	0000001774	1817 1900001816
1774	1910471775	1818 3100001819
		1819 1800001820
		1820 0000121840

-20- UCRL-8282

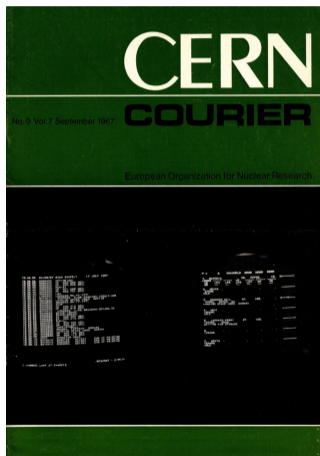
Instruction	Address	Description	Operation	Storage Location	Element	Remarks
1840	0001001841		R0			
1841	0000001840		SR			
1842	0000001843		SR			
1843	0410413043		R0	003		
1844	0000181043		SR			
1845	1910411846		SR			
1846	3100001848		SR			
1848	0000001849		SR			
1849	1910411850		SR			
1850	3100001851		SR			
1851	0010471853		SR	100		
1852	0010411853		R0			
1853	1100001854		SR			
1854	0100104855		SR	004		
1855	1910411856		SR			
1856	3100001857		SR			
1857	0000001858		SR			
1858	1910411859		SR			
1859	3100001860		SR			
1860	0010471861		SR	100		
1861	0000141865		SR			
1862	1010411862		R0			
1863	1910411863		SR			
1864	0000001866		SR			
1865	0000121867		SR	003		
1866	0018011868		R0			
1867	0000001869		SR			
1868	1900001870		SR			
1869	0000001871		SR			
1870	1900121873		SR			
1871	3100001874		SR			
1872	0000001875		SR			
1873	0010471876		SR	100		
1874	0000001877		SR			
1875	0010400180		SR			

Description location 0111 refers to integrating routine for the three simultaneous equation.

11/50



September 1967



Fortran
mentioned
3 times

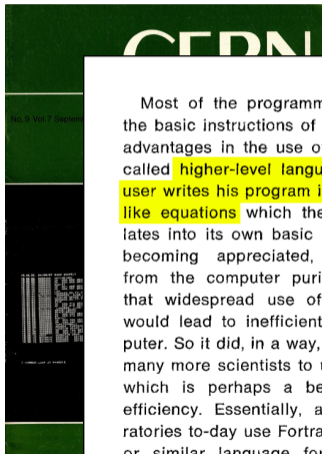
March 1972



Fortran
mentioned
18 times

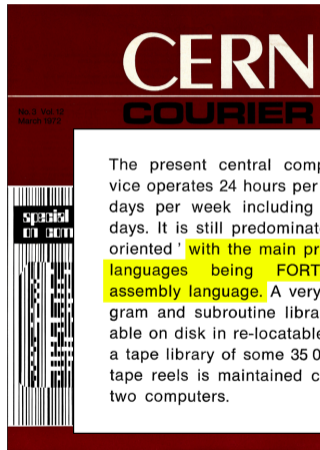


September 1967



Most of the programming was done in the basic instructions of the computer, but advantages in the use of Autocode (a so-called higher-level language in which the user writes his program in text and algebra-like equations which the computer translates into its own basic instructions) were becoming appreciated, amidst protests from the computer purists who believed that widespread use of such languages would lead to inefficient use of the computer. So it did, in a way, but it also allowed many more scientists to use the computer, which is perhaps a better criterion for efficiency. Essentially, all research laboratories to-day use Fortran, Algol, Autocode or similar language for scientific work.

March 1972



The present central computing service operates 24 hours per day, seven days per week including most holidays. It is still predominately 'batch-oriented' with the main programming languages being FORTRAN and assembly language. A very large program and subroutine library is available on disk in re-locatable form, and a tape library of some 35 000 labelled tape reels is maintained close to the two computers.

an
tioned
mes



NHEP adopted Fortran for data analysis in the first years after its release in 1956.



NHEP adopted Fortran for data analysis in the first years after its release in 1956.

“Algebra-like equations” (i.e. **FOR**mula-**TRAN**slation) was an obvious benefit.



NHEP adopted Fortran for data analysis in the first years after its release in 1956.

“Algebra-like equations” (i.e. **FOR**mula-**TRAN**slation) was an obvious benefit.

Also portability: assembly programs only work on a single model of computer.



NHEP adopted Fortran for data analysis in the first years after its release in 1956.

“Algebra-like equations” (i.e. **FOR**mula-**TRAN**slation) was an obvious benefit.

Also portability: assembly programs only work on a single model of computer.

Why Fortran and not something else (e.g. ALGOL)?

It was IBM's product, and most labs were buying IBM computers.

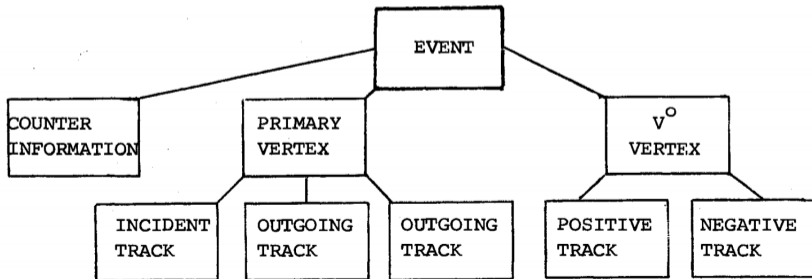


Part 2: C++



Example: High Energy Physics events are made up of vertices, every vertex has tracks associated to it. Also, to each event is associated a bank of information concerning electronic counter information to be used later. Assume the event to be a two-prong with an associated V^0 .

The pictorial graph for this event information is then

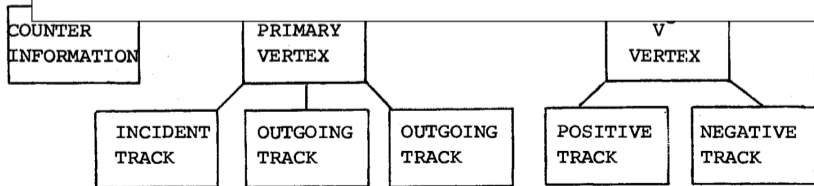


R.K. Böck, *Initiation to Hydra* (1974)



Example: High Energy Physics events are made up of vertices, every vertex has tracks associated to it. Also, to each event is associated a bank of information concerning electronic counter information to be used later. Assume the event to be a two-prong with an associated V^0 .

This wasn't a part of Fortran until 1991.
Physicists created *libraries* for tree-like data.



R.K. Böck, *Initiation to Hydra* (1974)



P. Lebrun and A. Kreymer, *High Level Language Memory Management on Parallel Architectures* (1989)

Years ago, the need for "pointer based" FORTRAN packages, such as HYDRA or BOS, later on ZBOOK [1], ZEBRA [2] or YBOS, became a necessity to efficiently manage the user heap space on a tight fixed size physical memory, in order to support large applications, such as code management systems or histogram packages.[1] Later on, these memory management systems were used not only to allocate dynamically space in a fixed size heap space, but also as a tool to organize and manage sensibly a complicate set of structures describing a detector, or an elementary particle collision. These packages are the essential building blocks for HEP data bases. They were designed to run mainly on single CPU systems and ignored entirely the existence of virtual memory available within FORTRAN through - system dependent ! - system calls.

Fortran lacked an essential feature for NHEP



inside a program. Logical relations between banks are expressed by including the address of one bank in the link-table of another bank. For example, all tracks of a vertex-point are linked together by each track pointing to the next. Such a *data-structure* contains not only the numeric information but also logical information about the object it describes.

The program modularity is achieved by organizing the program into processors each having a well defined task. This task is *entirely* describable as a transformation applied to a data-structure in the dynamic store : some banks provide the input data to the processor and some contain the desired results. For a given application, a steering program is written to coordinate the operations of the processors needed. Any processor consists of at least one FORTRAN subroutine, its operation being invoked by transferring control to this

subroutine. As a matter of internal organization, the processor may be divided up into the primary and several secondary subroutines. The programming of a processor has to observe certain conventions in order to be compatible with the HYDRA system. Precisely these conventions, which are the same through the whole program (indeed through all HYDRA programs) are responsible for the easy documentation and the good readability of the program.

The processors are supported by the HYDRA system. Its services are requested with CALL statements much like the services of the FORTRAN system which are part of the definition of the basic language. In this sense, the HYDRA system is an extension of the FORTRAN language to provide - primarily - dynamic memory management facilities. Some languages contain these facilities in their basic definition, but the HYDRA-

FORTTRAN combination has two important advantages — the execution speed is that of a normal FORTRAN program, with very little overhead for the HYDRA system, and FORTRAN is a commonly accepted language.

Because of the need for machine independence (so that the same programs can be used on a variety of computers) the processors for the new bubble chamber program, as well as the HYDRA system packages, have been written in ANSI FORTRAN which is the internationally accepted minimum requirement expected from anybody's compiler.

The bubble chamber programs of the HYDRA form will come into operation in 1972. They should help to tear down the walls that have sometimes threatened, on the data handling side, to separate physicists from computer specialists, or bubble chamber groups from each other and from physicists using other techniques.

R.K. Böck and J. Zoll, *Central Computers in the Analysis*, CERN Courier No. 3, Vol. 12, March 1972



COBOL (1959)

```
01 Point.  
   05 x      pic 9(3).  
   05 y      pic 9(3).
```

Simula (1962)

```
Class Point (x, y);  
  Integer x, y;  
  ! define attributes...  
  ! define methods...  
End of Point;
```

PL/I (1964)

```
define structure  
  1 point,  
  2 x integer,  
  2 y integer;
```

ALGOL-68 (1968)

```
MODE POINT = STRUCT(  
  INT x,  
  INT y  
);
```

Pascal (1970)

```
type Point = record  
  x, y: integer;  
end;
```

C (1972)

```
typedef struct Point {  
  int x;  
  int y;  
} Point;
```



Paul Kunz, *Physics Analysis Tools* (1991)

The C programming language is much better than FORTRAN for both data structures and configuration control. Shown in Figure 4 are some segments of C code that one might use in dealing with a track entity. Note the expressive power of the language in that access to variables is by full name. Also, the C language deals directly with the dynamic memory allocation of such structures since the memory allocation functions are part of its standard library. Finally, there's nothing lost in using a symbolic debugger because structures are part of the language, thus known to the existing debugger.

Clearly, C is much better at handling data structures than FORTRAN plus some additional package. Although many large collaborations have discussed abandoning FORTRAN, none has done so (yet). One reason they stayed with FORTRAN is reluctance to learn a new language, which is ironic since in each case they had to learn a big system to complement FORTRAN.

```
struct track {
    float px;
    float py;
    float pz;
};
... ..
struct track **mctrack;
... ..
px = mctrack[it]->px;
```

Figure 4. Segments of C code.



Paul Kunz, *Physics Analysis Tools* (1991)

The C programming language is much better than FORTRAN for both data structures and configuration control. Shown in Figure 4 are some segments of C code that one might use in dealing with a track entity. Note the expressive power of the language in that access to variables is by full name. Also, the C language deals directly with the dynamic memory allocation of such structures since the memory allocation functions are part of its standard library. Finally, there's nothing lost in using a symbolic debugger because structures are part of the language, thus known to the existing debugger.

Clearly, C is much better at handling data structures than FORTRAN plus some additional package. Although many large collaborations have discussed abandoning FORTRAN, none has done so (yet). One reason they stayed with FORTRAN is reluctance to learn a new language, which is ironic since in each case they had to learn a big system to complement FORTRAN.

```
struct track {
    float px;
    float py;
    float pz;
};
... ..
struct track **mctrack;
... ..
px = mctrack[it]->px;
```

Figure 4. Segments of C code.

In the first 15 years of CHEP (1985–2000), similar suggestions were made for ALGOL, PL/I, Pascal, Ada, Eiffel, Objective C, Java, and of course C++.



Fortran-90 (1991)

```
type Point
  integer :: x
  integer :: y
end type Point
```



René Brun, *Technologies, Collaborations and Languages: 20 Years of HEP Computing* (2012)

Fortran-90 (1991)

```
type Point
  integer :: x
  integer :: y
end type Point
```

One major stumbling block in the move to FORTRAN 90 was the question of Input/Output. With ZEBRA, we had a simple way to describe data structures (banks) built out of basic types (typically integers and floats). Because FORTRAN 90 supported derived data types, it was theoretically possible to implement the most complex data structures that we used to model with ZEBRA. In particular ZEBRA was able to write and read these data structures from machine independent files.

The need for introspection to deal with derived data types doomed the efforts to move to FORTRAN 90. It was going to be as hard (if not more) with C++, but we did not know this at the time.

Using FORTRAN 90, it appeared pretty hard to make a general implementation equivalent to ZEBRA without parsing the data type description in the FORTRAN 90 modules. In fact, we encountered the same problem later with C++, but we naively ignored at that time how much work it was to implement a data dictionary or reflection system describing at run time the type of objects. Mike Metcalf was aware of the problem and we reported this to a special session of the FORTRAN committee at CERN in 1992. As most members in the committee had no experience with this problem and thought that this was a database problem and not a language problem, the enhancements that we were expecting in the language did not happen.



So actually, NHEP needs:

1. rich data structures in the programming language,
2. serialization of those structures to/from disk,
3. read compatibility for old data versions (schema evolution),
4. mapping between persistent data and language's structures (which can be implemented with type-introspection).



So actually, NHEP needs:

1. rich data structures in the programming language,
2. serialization of those structures to/from disk,
3. read compatibility for old data versions (schema evolution),
4. mapping between persistent data and language's structures (which can be implemented with type-introspection).

HYDRA/ZEBRA/etc. approached these as *a single* problem.



So actually, NHEP needs:

1. rich data structures in the programming language,
2. serialization of those structures to/from disk,
3. read compatibility for old data versions (schema evolution),
4. mapping between persistent data and language's structures (which can be implemented with type-introspection).

HYDRA/ZEBRA/etc. approached these as *a single* problem.

In Java, (1–2) is a language problem, (3–4) is for databases.



So actually, NHEP needs:

1. rich data structures in the programming language,
2. serialization of those structures to/from disk,
3. read compatibility for old data versions (schema evolution),
4. mapping between persistent data and language's structures (which can be implemented with type-introspection).

HYDRA/ZEBRA/etc. approached these as *a single* problem.

In Java, (1–2) is a language problem, (3–4) is for databases.

Object databases focus on (4)...

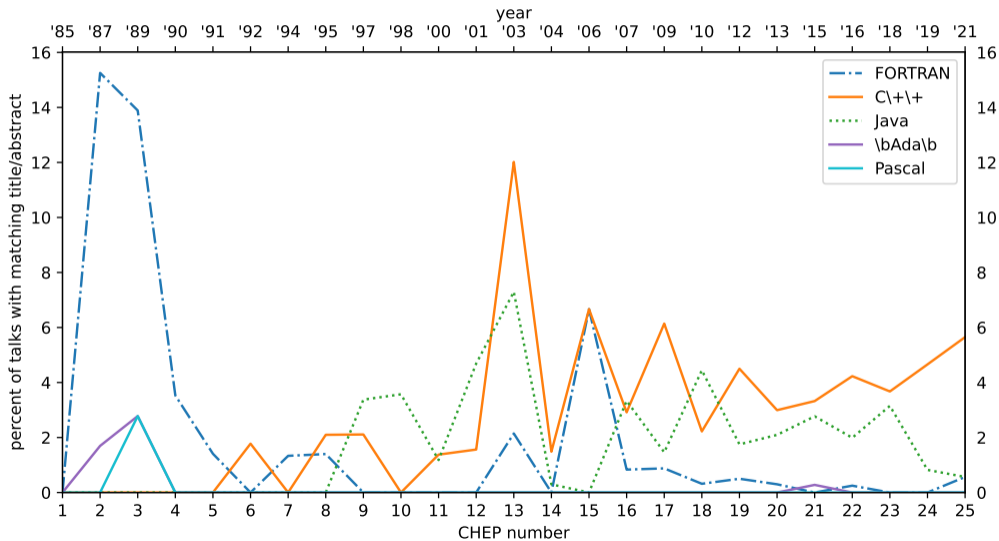


So it wasn't just a matter of adding another language to the mix.

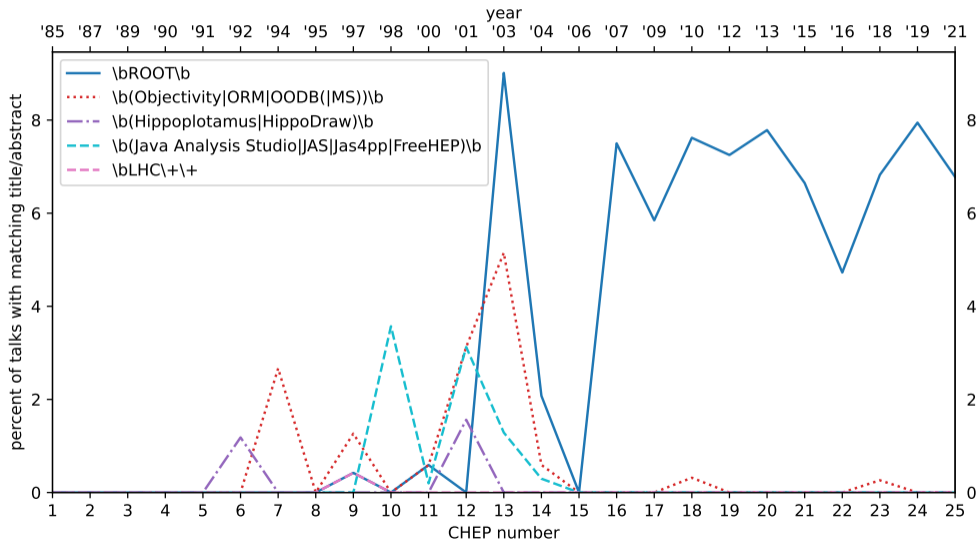
The whole infrastructure, including I/O, had to change.

You only want to do that once!

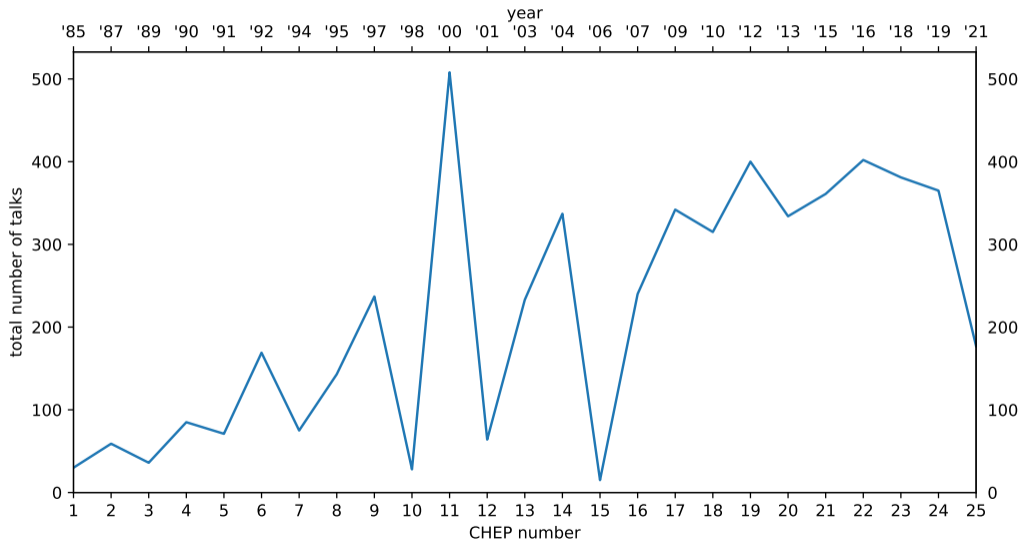
Programming languages in CHEP title/abstract regex matches



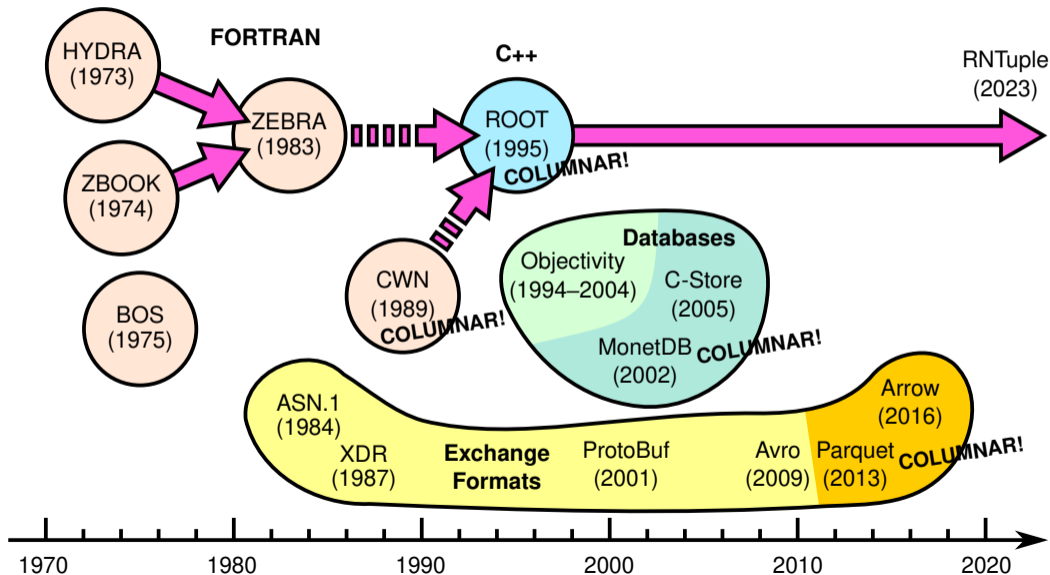
Data infrastructure in CHEP title/abstract regex matches



Total number of talks (denominator): dips are small CHEPs



ROOT I/O is custom, but was more advanced than alternatives





Dremel: Interactive Analysis of Web-Scale Datasets (Google, 2010)

storage and reduce CPU cost due to cheaper compression. Column stores have been adopted for analyzing relational data [1] but to the best of our knowledge have not been extended to nested data models. The columnar storage format that we present is supported by many data processing tools at Google, including MR, Sawzall [20], and FlumeJava [7].

In this paper we make the following contributions:

- We describe a novel columnar storage format for nested data. We present algorithms for dissecting nested records into columns and reassembling them (Section 4).

Columnar, nested data was a ROOT feature 15 years earlier.



NHEP adoption of C++ (or similar) was held back by the fact that we had unique infrastructure.



NHEP adoption of C++ (or similar) was held back by the fact that we had unique infrastructure.

It's important to have data structures in the language, but also on disk.



NHEP adoption of C++ (or similar) was held back by the fact that we had unique infrastructure.

It's important to have data structures in the language, but also on disk.

Our solutions were both *more advanced* and *less modular* than others.



NHEP adoption of C++ (or similar) was held back by the fact that we had unique infrastructure.

It's important to have data structures in the language, but also on disk.

Our solutions were both *more advanced* and *less modular* than others.

Many options considered in late 1990's; C++ with ROOT I/O became dominant.



Part 3: Python





NHEP has a history of custom solutions for interactivity:
SPEAKEASY, Minuit, PAW, KUIP, CINT, Cling...



NHEP has a history of custom solutions for interactivity:
SPEAKEASY, Minuit, PAW, KUIP, CINT, Cling...

But the number of industry solutions is also vast.

Interactive mode languages [\[edit\]](#)

Interactive mode languages act as a kind of shell: expressions or statements can be entered one at a time, and the result of their evaluation is seen immediately. The interactive mode is also termed a [read-eval-print loop](#) (REPL).

- [APL](#)
- [BASIC](#) (some dialects)
- [Clojure](#)
- [Common Lisp](#)
- [Dart](#) (with Observatory or Dartium's developer tools)
- [ECMAScript](#)
 - [ActionScript](#)
 - [ECMAScript for XML](#)
 - [JavaScript](#)
 - [JScript](#)
 - [Source](#)
- [Erlang](#)
- [Elixir](#) (with iex)
- [F#](#)
- [Fril](#)
- [GAUSS](#)
- [Groovy](#)
- [Haskell](#) (with the GHCi or Hugs interpreter)
- [IDL](#)
- [J](#)
- [Java](#) (since version 9)
- [Julia](#)
- [Lua](#)
- [MUMPS](#) (an ANSI standard general purpose language)
- [Maple](#)
- [Mathematica](#) (Wolfram language)
- [MATLAB](#)
- [ML](#)
- [OCaml](#)
- [Perl](#)
- [PHP](#)
- [Pike](#)
- [PostScript](#)
- [Prolog](#)
- [Python](#)
- [PROSE](#)
- [R](#)
- [REBOL](#)
- [Rexx](#)
- [Ruby](#) (with IRB)
- [Scala](#)
- [Scheme](#)
- [Smalltalk](#) (anywhere in a Smalltalk environment)
- [S-Lang](#) (with the S-Lang shell, slsh)
- [Speakeasy](#)
- [Swift](#)
- [Tcl](#) (with the Tcl shell, tclsh)
- [Unix shell](#)
- [Windows PowerShell](#) (.NET-based CLI)
- [Visual FoxPro](#)

In recent years, the industry has consolidated on Python

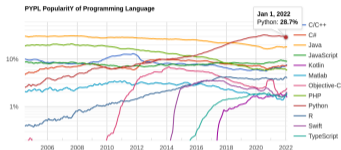


Python is currently leading every “most popular programming language” index.

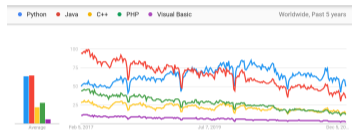
Tiobe

Jan 2022	Jan 2021	Change	Programming Language	Rankings	Change
1	3	▲	Python	13.58%	+1.86%
2	1	▼	C	12.44%	-1.94%
3	2	▼	Java	10.60%	-1.30%
4	4		C++	8.29%	+0.73%
5	5		C#	5.60%	+1.73%
6	6		Visual Basic	4.74%	+0.80%
7	7		JavaScript	2.09%	-0.11%
8	11	▲	Assembly language	1.85%	+0.23%

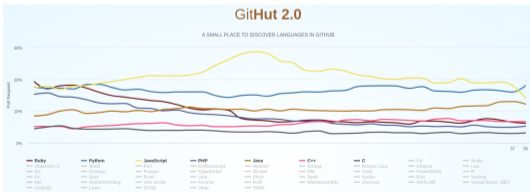
PYPL



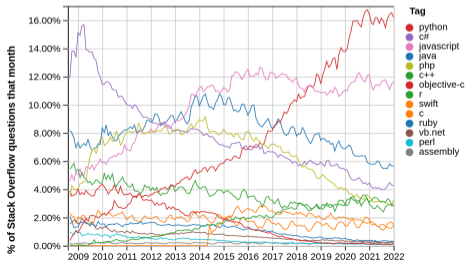
Google Trends



GitHub

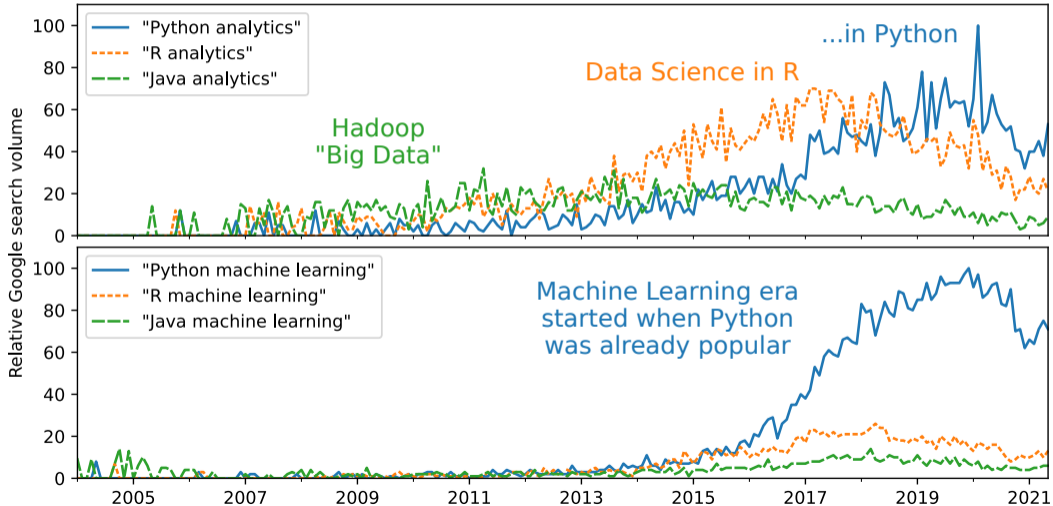


StackOverflow

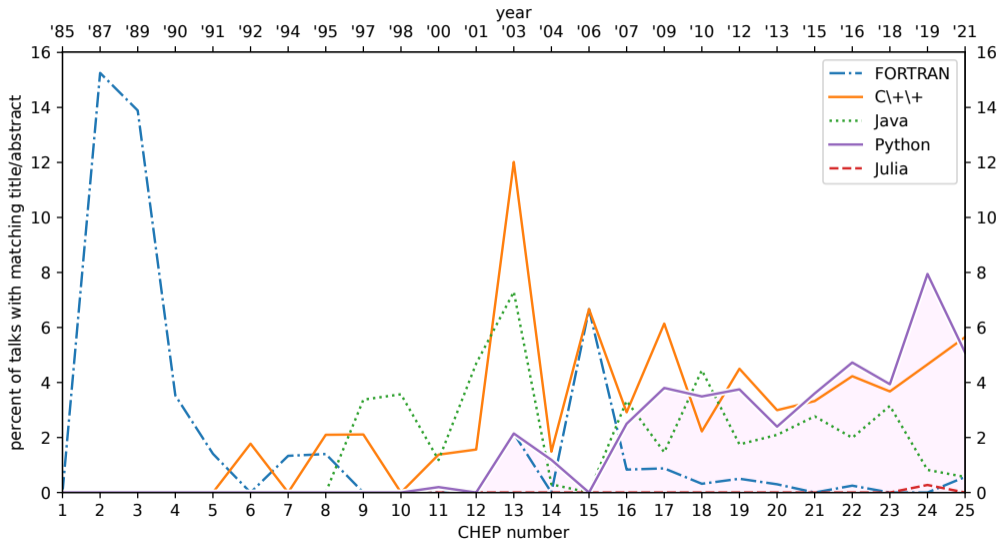




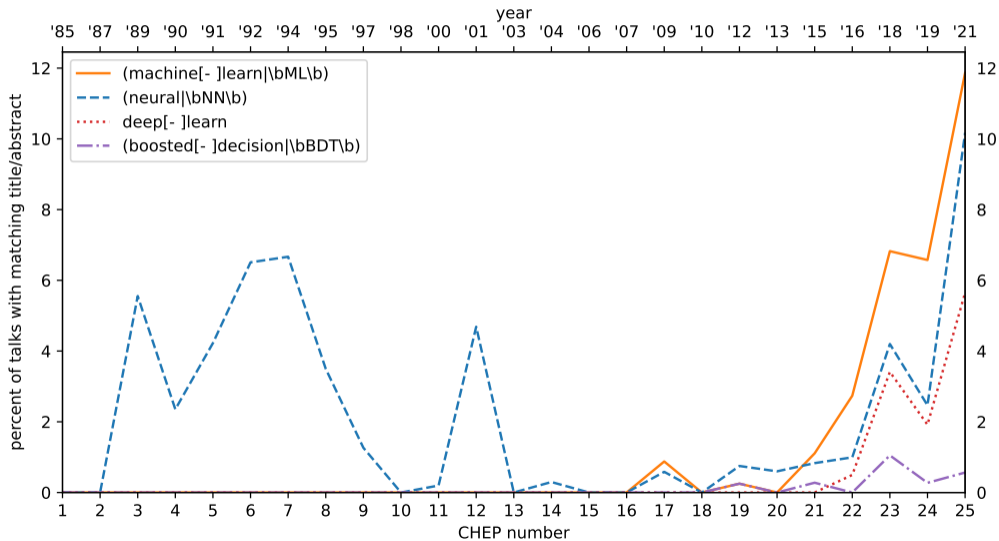
Correlated words in Google searches (Google Trends).



Python use has also been rising—steadily—in NHEP



Since *before* the return of machine learning





Emerging Standard ? Python as “Software Glue”

■ Clear trend towards Python

- ❖ Used by: ATLAS (Athena), CMS, D0, LHCb (Gaudi), SND, ...
- ❖ Used by: Lizard/Anaphe, HippoDraw, JAS (Jython)...
- ❖ Architecturally, scripting is “just another service”
- ❖ ROOT is the exception to the “Python rule”
 - CINT interpreter plays a central role
 - Developers and users seem happy

■ Python is popular with developers...

- ❖ Rapid prototyping; gluing together code
- ❖ (Almost) auto-generation of wrappers (SWIG)

■ ...but acceptance by users not yet proven

- ❖ Another language to learn, syntax, ...

*“Summary of Track 2: Data Analysis and Visualisation”
Lucas Taylor, Northeastern U. CHEP 01, Beijing, 3-7 S*

Note: PyROOT introduced in 2004 (v4.00/04).



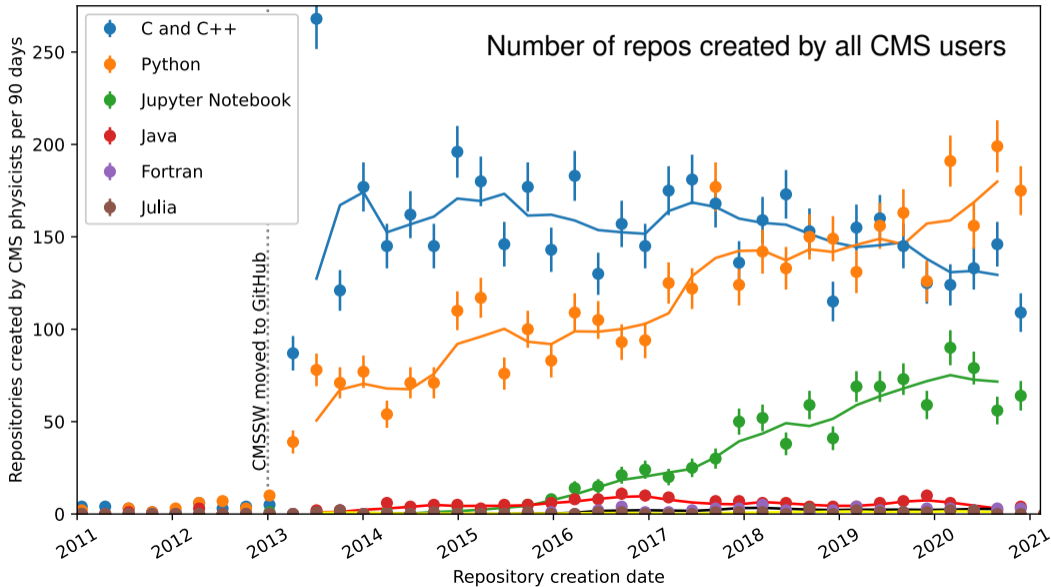
Stephan Lammel, *Computing models of major HEP experiments: DØ and CDF*, 1997

DØ has made the decision to move all large software projects to C++. Their framework approach has a set of modules that execute sequentially, each having a specific task. The glue that holds the individual software packages together will be an interpreted script system. The main task of this framework is to “guide” data between the various modules/packages ... prototype framework based on the Python scripting language has been developed and is ready for use.

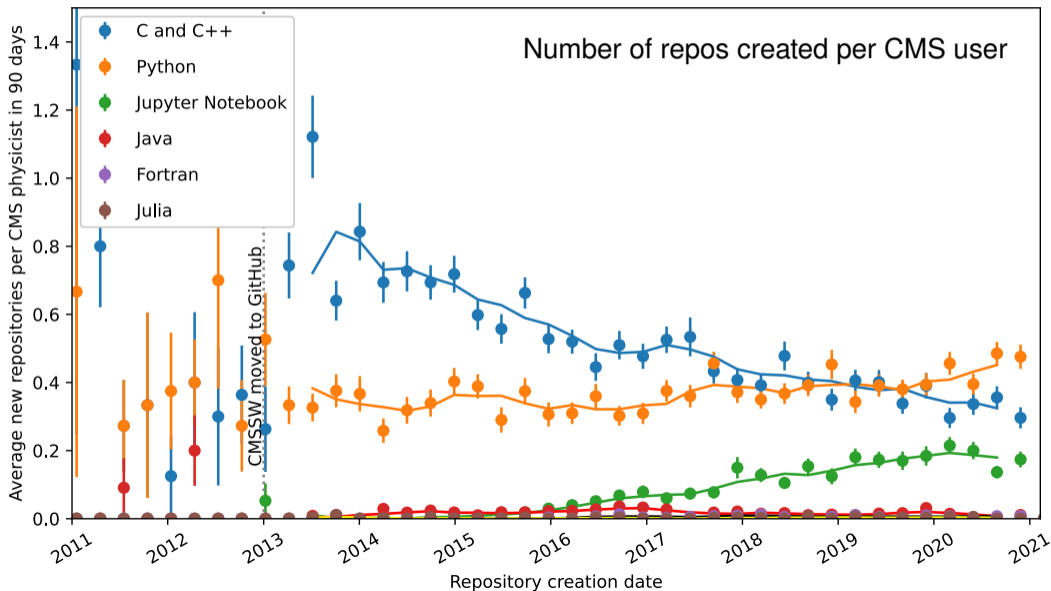
Jeff Templon, *Python as an Integration Language*, SPAG-1998-02, 1998

Most of the code would be written in C or C++, but the integration would be done through Python. This enables the uninitiated to make simple modifications to the analysis which were perhaps not thought of by the authors; all the neophyte needs to know is how the interfaces work. On the other hand, it will force the code authors to make the analysis subsystems independent of each other (one of the big problems with the current code), and will encourage rigorous testing of subunits.

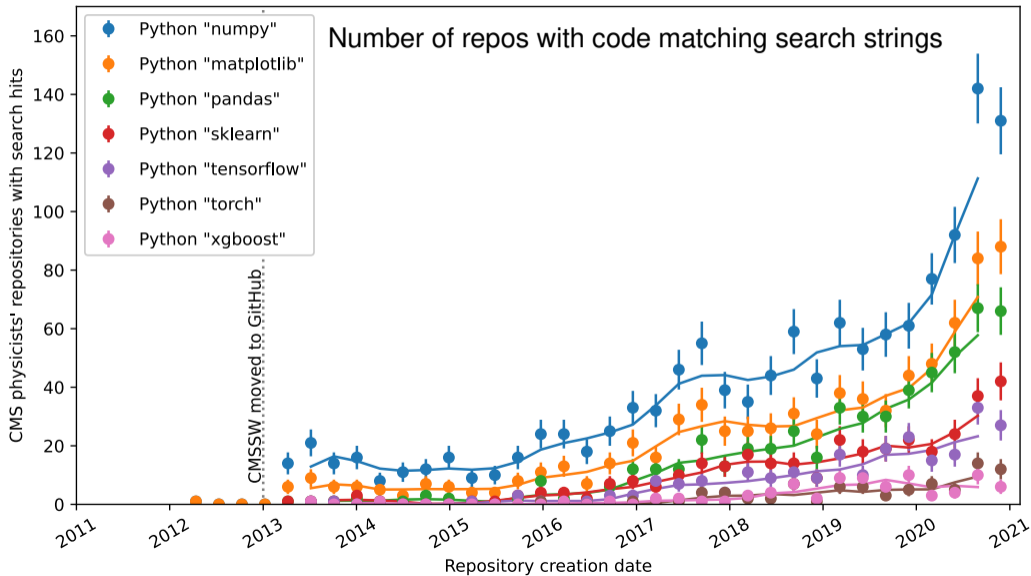
Trends in code written by CMS users (in GitHub)



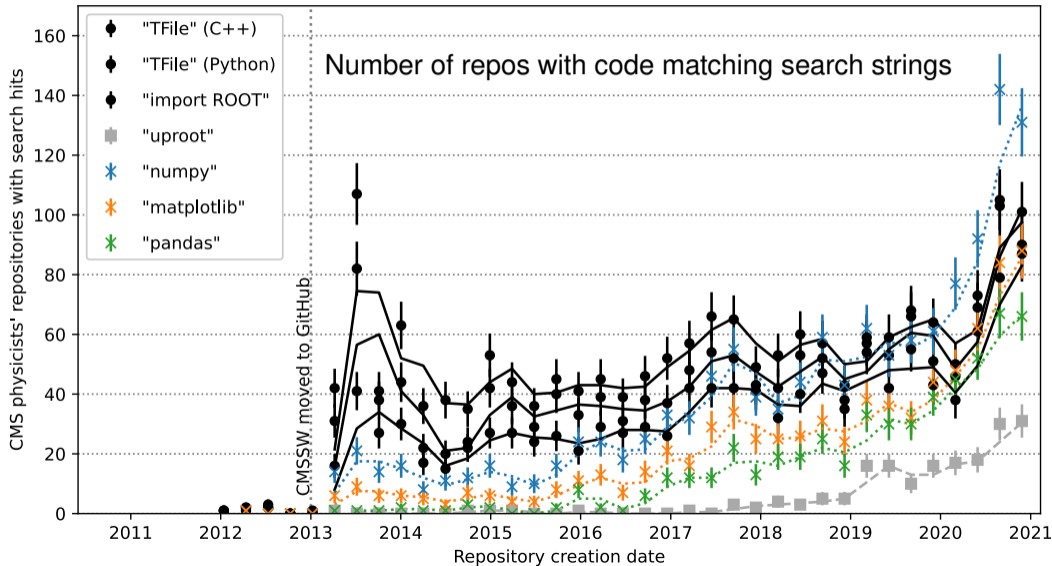
Trends in code written by CMS users (in GitHub)



Trends in code written by CMS users (in GitHub)



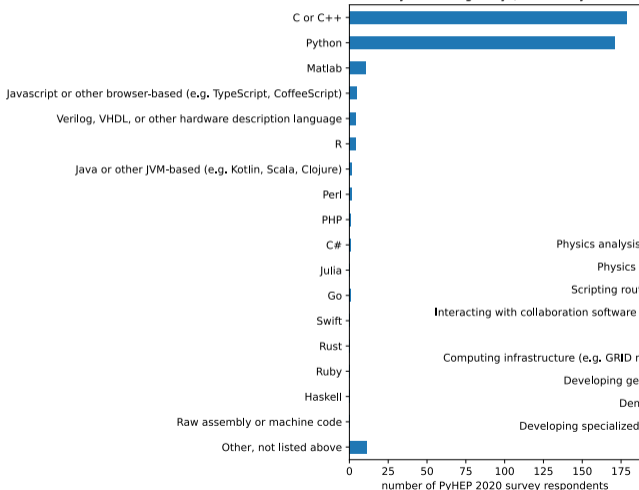
Trends in code written by CMS users (in GitHub)



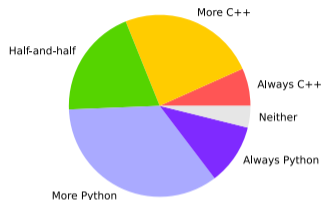
Survey responses from PyHEP 2020 attendees ($N = 406$)



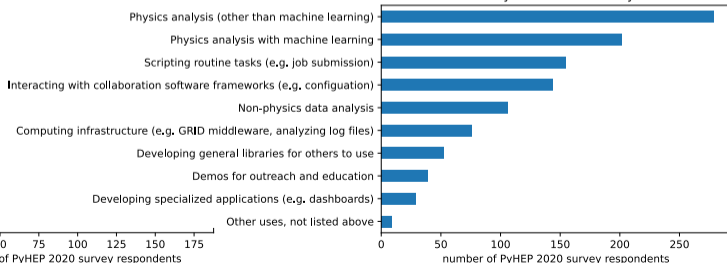
"Which do you use regularly (> 10% of your work)?"



"How often do you use Python relative to C or C++?"



"What are your main uses of Python?"





NHEP adoption of Python started long before machine learning and columnar analysis trends.



NHEP adoption of Python started long before machine learning and columnar analysis trends.

Interactive/fluid programming has always been a need, and has traditionally been met by a variety of alternatives.



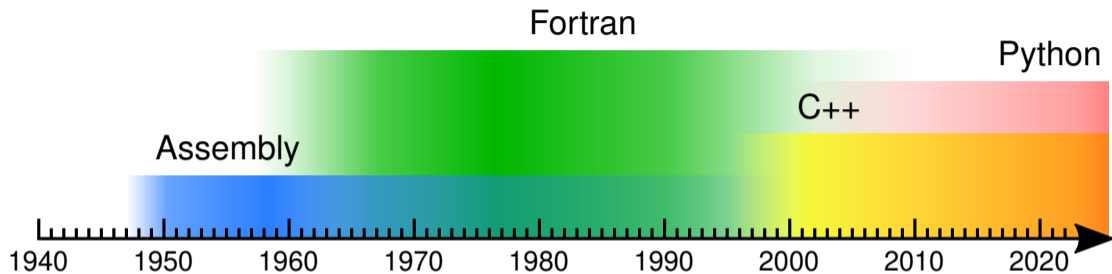
NHEP adoption of Python started long before machine learning and columnar analysis trends.

Interactive/fluid programming has always been a need, and has traditionally been met by a variety of alternatives.

What's new is the consolidation on one language, Python, and an increase in how much analysis logic can be “driven” from the interactive language.



Concluding conclusions



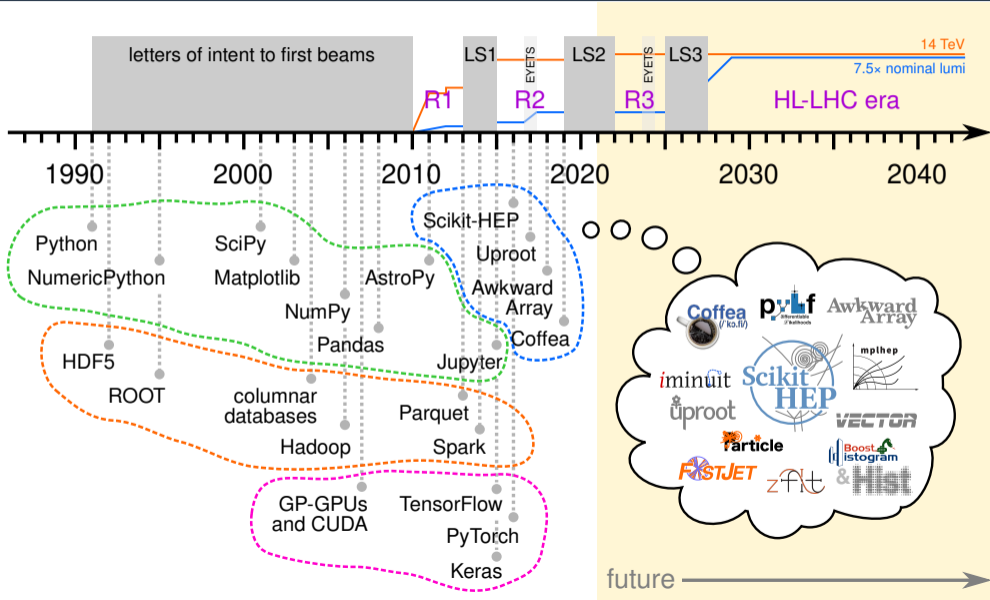
Adoption of

- Fortran: immediate; for syntax and portability; no infrastructure to replace
- C++: long overdue; for data structures; replaced infrastructure in a burst
- Python: slowly overtook its alternatives; for interactivity; different niche
- Julia? slowly mix in among the Python and C++ until it's all that's left?

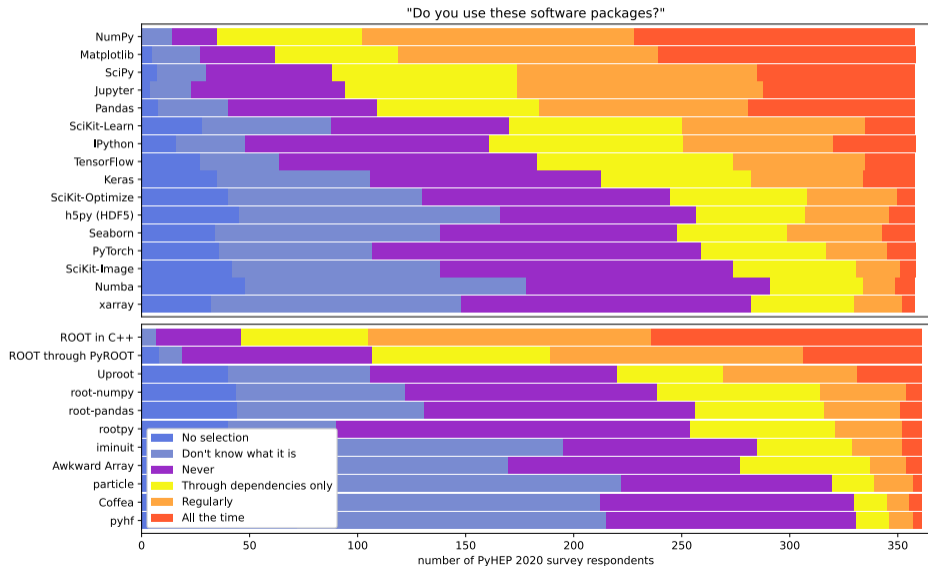


Backup

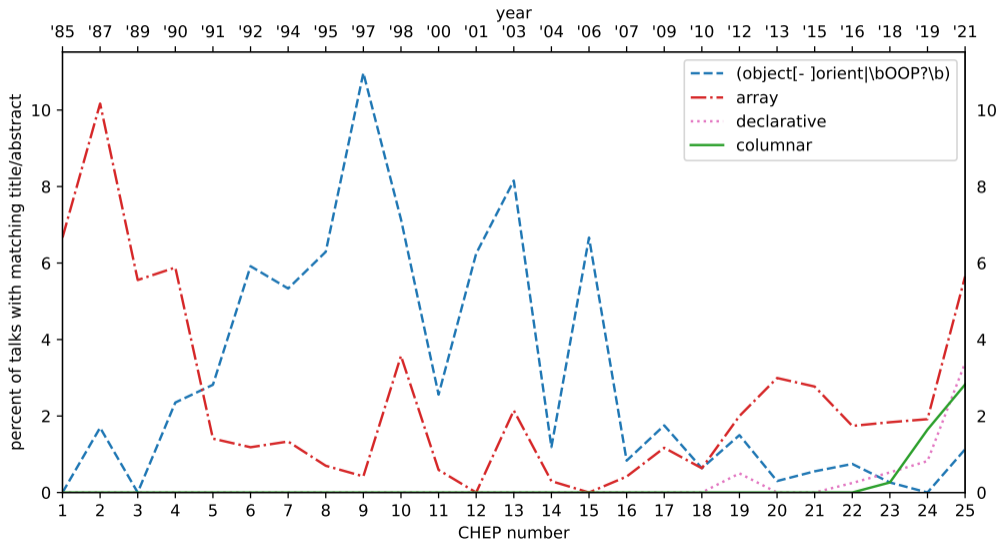
Python ecosystem has had as much time to evolve as the LHC

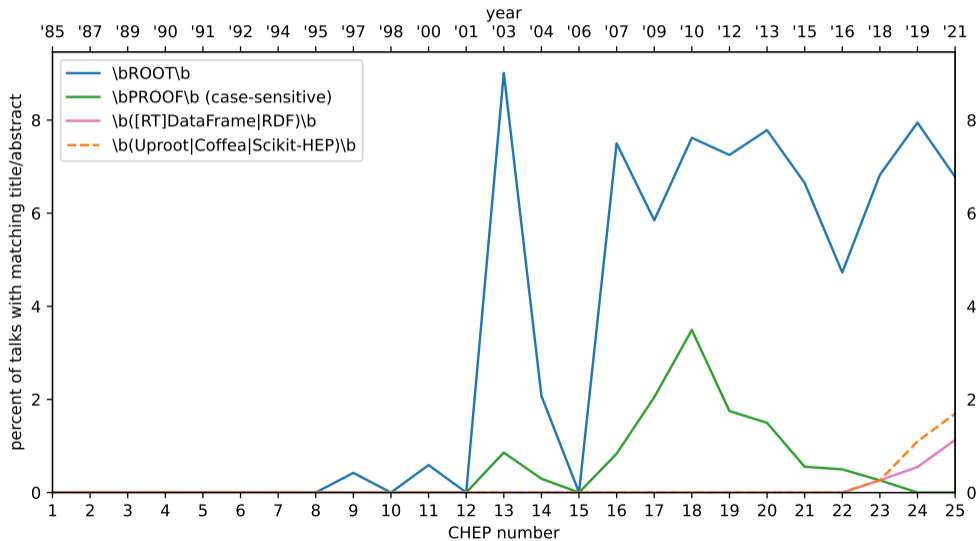


Widespread familiarity with data science tools (PyHEP survey)



More CHEP history





More CHEP history

