

## THE MISSING **PYTHON INTRODUCTION** FOR SCIENTISTS.

**TAMAS GAL** – ERLANGEN CENTRE FOR ASTROPARTICLE PHYSICS (ECAP)



Friedrich-Alexander-Universität **Erlangen-Nürnberg** 



ERLANGEN CENTR FOR ASTROPARTICI F PHYSICS



## WHO AM I?

- Tamás Gál, born 1985 in Debrecen (Hungary)
- Astroparticle physicist at the Erlangen Centre for Astroparticle Physics (ECAP) working on the KM3NeT neutrino detector experiment and open science/data
- Sysadmin (DevOps) at ECAP (including the ECAP and KM3NeT IT • services)
- **Programming** background: •
  - Coding enthusiast since ~1993
  - First real application written in **Amiga Basic** (toilet manager, tons of GOTOs;)
  - Mostly Julia, Python, Rust, JavaScript and C/C++ for work
  - Haskell for fun
  - Earlier also Obj-C, Java, Perl, PHP, Delphi, MATLAB, whatsoever...
- Started with **Python** around **1998** (to replace Perl/Shell)
- Editor: Vim for ~25 years but switched to Emacs (EVIL mode) in 2020 •
- Other: ADV motorbikes, climbing, electronics, modular synths, DIY...





### DISCLAIMER

The following presentation contains oversimplifications and blatant omissions due to time constraints. The viewer may also experience well-timed product placements.





## THE PYTHON PROGRAMMING LANGUAGE

- Interpreted high-level general-purpose programming language
- Dynamically-typed and garbage-collected
- "batteries included"
- necessary

• Object-oriented, procedural (imperative), functional, structured, reflective



• Tries to avoid premature optimisation: move time-critical functions to extension modules written in "faster" languages (like C or Fortran) when



### >>> import this

The Zen of Python, by Tim Peters Beautiful is better than ugly. Explicit is better than implicit. Simple is better than complex. Complex is better than complicated. Flat is better than nested. Sparse is better than dense. Readability counts.

Special cases aren't special enough to break the rules. Although practicality beats purity.

Errors should never pass silently.

Unless explicitly silenced.

In the face of ambiguity, refuse the temptation to guess. There should be one-- and preferably only one --obvious way to do it. Although that way may not be obvious at first unless you're Dutch. Now is better than never.

Although never is often better than \*right\* now. If the implementation is hard to explain, it's a bad idea. If the implementation is easy to explain, it may be a good idea. Namespaces are one honking great idea -- let's do more of those

### THE ZEN OF PYTHON







TIOBE Programming Community Index

Source: www.tiobe.com



**Python** is the most popular language (according to TIOBE)! ... and has beaten Java and C++! Btw. Julia is #33 ;)

## POPULAR LANGUAGES

### Most loved languages <u>https://survey.stackoverflow.co/2022</u>





### YOUR JOURNEY THROUGH PYTHON? (JUST A VERY ROUGH GUESS, NOT A MEAN GAME)

Raise your hand and keep it up until you answer a question with "no" (means you're out).

- Have you ever launched the Python interpreter?
- Wrote for/while-loops or if/else statements?
- ...your own functions?
- ...classes?
- ...list/dict/set comprehensions?
- Do you know what a generator is?
- Have you ever implemented a decorator?
- ...a metaclass?
- ...a C-extension?
- Do you know and can you explain the output of the following line for Python?

### print(5 is 7 - 2, 300 is 302 - 2)

Advanced

I am the BDFL of Python





**Python 2.7**: True, False **Python 3.6:** True, False Python 3.7: True, True Python 3.8: True, True, and warnings ... Python 3.9: True, True, and warnings ... Python 3.10: True, True, and warnings ...

### **ANSWER TO** print(5 is 7 - 2, 300 is 302 - 2)



### This picture makes none



PyObject\* PyLong\_FromLong(long v) Return value: New reference. Return a new PyLongObject object from v, or NULL on failure.

The current implementation keeps an array of integer objects for all integers between -5 and 256, when you create an int in that range you actually just get back a reference to the existing object.

"is" is an operator which checks if two objects are identical: "x is y" is true iff x and y are pointing to the same object.

In Python 3.7+ the constant folding is moved from the peephole optimiser to the new AST optimiser, which effectively avoids the extra allocation. (https://github.com/python/cpython/commit/7ea143ae795a9fd57eaccf490d316bdc13ee9065)





## WHY IS PYTHON SO POPULAR (FOR SCIENCE)?

- Ease of use scientists don't know/want how to program
- Readable code source code is more often read than written
- Interactive workflow
- Lots of scientific libraries (and machine learning is everywhere)
- Batteries included: tons of (built-in) useful supplementary functionalities
- General purpose language so that scientists can focus on a single language to rule them all...
- ...can they?





### PERFORMANCE OF LANGUAGES

Microbenchmarks from <a href="https://julialang.org/benchmarks/">https://julialang.org/benchmarks/</a>







**TO UNDERSTAND THE PERFORMANCE ISSUES** 

## CPYTHON

- Python (in contrast to other languages like C, Julia) has no formal specification
- The Python Language Reference is written in English: <u>https://docs.python.org/3.9/reference</u>
- CPython is the reference implementation which contains implementation details which are not part of the language (e.g. GC with reference counting)
- Python is defined partly by the Python Language Reference and its main implementation **CPython**
- There are several other implementations: PyPy, Jython (stuck at Python 2.7), IronPython (2.7) and 3.4), etc.
- From now on, with "**Python**" we refer to **CPython**



## FROM SOURCE TO RUNTIME



Recommended reading (Victor Skvortsov - Python behind the scenes): https://tenthousandmeters.com/blog/python-behind-the-scenes-1-how-the-cpython-vm-works/



### PYTHON BYTECODE print(5 is 7 - 2, 300 is 302 - 2)

		import	dis;	<pre>dis.dis(compile('print(5</pre>
9	נח	1		0 LOAD_NAME
•	S (			2 LOAD_CONST
$\mathbf{c}$				4 LOAD_CONST
	g			6 COMPARE_OP
				8 LOAD_CONST
$\bigcirc$				10 LOAD_CONST
	ົດ			12 COMPARE_OP
	n			14 CALL_FUNCTION
	Γ			16 PRINT_EXPR
	$\vdash$			18 LOAD_CONST
				20 RETURN VALUE
		import	dis;	<pre>dis.dis(compile('print(5</pre>
7		<b>import</b>	dis;	<pre>dis.dis(compile('print(5</pre>
. 7	Je	<b>import</b> 1	dis;	<pre>dis.dis(compile('print(5</pre>
3.7	rue	<b>import</b> 1	dis;	<pre>dis.dis(compile('print(5</pre>
<b>3</b> .7	True	<b>import</b> 1	dis;	<pre>dis.dis(compile('print(5</pre>
on 3.7	True	<b>import</b> 1	dis;	<pre>dis.dis(compile('print(5</pre>
10n 3.7	, True	<b>import</b> 1	dis;	<pre>dis.dis(compile('print(5       0 LOAD_NAME       2 LOAD_CONST       4 LOAD_CONST       6 COMPARE_OP       8 LOAD_CONST       10 LOAD_CONST</pre>
thon 3.7	le, True	import 1	dis;	<pre>dis.dis(compile('print(5</pre>
vthon 3.7	rue, True	import 1	dis;	<pre>dis.dis(compile('print(5</pre>
Python 3.7	True, True	import (	dis;	<pre>dis.dis(compile('print(5         0 LOAD_NAME         2 LOAD_CONST         4 LOAD_CONST         6 COMPARE_OP         8 LOAD_CONST         10 LOAD_CONST         10 LOAD_CONST         12 COMPARE_OP         14 CALL_FUNCTION         16 PRINT_EXPR</pre>

20 RETURN\_VALUE



5 (None)





## THE TYPE OF A PyObject

"An object has a 'type' that determines what it represents and what kind of data it contains. An object's type is fixed when it is created. Types themselves are represented as objects. The type itself has a type pointer pointing to the object representing the type 'type', which contains a pointer to itself!" — object.h



## DATA IN PYTHON

### • Every piece of data is a PyObject

### >>> **dir**(42)

['\_\_abs\_\_', '\_\_add\_\_', '\_\_and\_\_', '\_\_bool\_\_', '\_\_ceil\_\_', '\_\_class\_\_', '\_\_delattr\_\_', '\_\_dir\_\_', '\_\_divmod\_\_', '\_\_doc\_\_', '\_\_eq\_\_', '\_\_float\_\_', '\_\_floor\_\_', '\_\_floordiv\_\_', '\_\_format\_\_', '\_\_ge\_\_', '\_\_getattribute\_\_', '\_\_getnewargs\_\_', '\_\_gt\_\_', '\_\_hash\_\_', '\_\_index\_\_', '\_\_init\_\_', '\_\_init\_subclass\_\_', '\_\_int\_\_', '\_\_invert\_\_', '\_\_le\_\_', '\_\_lshift\_\_', '\_\_lt\_\_', '\_\_mod\_\_', '\_\_mul\_\_', '\_\_ne\_\_', '\_\_neg\_\_', '\_\_new\_\_', '\_\_or\_\_', '\_\_pos\_\_', '\_\_pow\_\_', '\_\_radd\_\_', '\_\_rand\_\_', '\_\_rdivmod\_\_', '\_\_reduce\_\_', '\_\_reduce\_ex\_\_', '\_\_repr\_\_', '\_\_rfloordiv\_\_', '\_\_rlshift\_\_', '\_\_rmod\_\_', '\_\_rmul\_\_', '\_\_ror\_\_', '\_\_round\_\_', '\_\_rpow\_\_', '\_\_rrshift\_\_', '\_\_rshift\_\_', '\_\_rsub\_\_', '\_\_rtruediv\_\_', '\_\_rxor\_\_', '\_\_setattr\_\_', '\_\_sizeof\_\_', '\_\_str\_\_', '\_\_sub\_\_', '\_\_subclasshook\_\_', '\_\_truediv\_\_', '\_\_trunc\_\_', '\_\_xor\_\_', 'bit\_length', 'conjugate', 'denominator', 'from\_bytes', 'imag', 'numerator', 'real', 'to\_bytes']





### DATA IN PYTHON



- Lower limit for the size in bytes (on a 64bit system), oversimplified: 8+8+8 (list) + 3\*(8+8+8+8+8) (integer) = **168**

• "Technically" it's only 24 bytes of information if we see it as an array of integers



### YOUR BEST FRIEND AND WORST ENEMY: GIL - Global Interpreter Lock

- The **GIL** prevents parallel execution of (Python) bytecode
- trouble with race conditions and deadlocks)
- Even though Python has real threads, they never execute (byte)code at the same time
- Threads perform pretty badly on **CPU bound** tasks
- They do a great job speeding up **I/O heavy** tasks

It is a very simple solution to memory safety (Python uses reference counting, which can cause

**Context switching** between threads creates overhead (the user cannot control thread-priority)



### THREADS AND CPU BOUND TASKS

```
single thread:
```

```
>>> N = 100_{000}_{000}
>>> def count(n):
        while n != 0:
. . .
       n -= 1
. . .
. . .
>>> %time count(N)
CPU times: user 2,97 s, sys: 6.28 ms, total: 2.98 s
Wall time: 2.98 s
```

### two threads:

```
>>> from threading import Thread
>>> def count_threaded(n):
       t1 = Thread(target=count, args=(N/2,))
. . .
       t2 = Thread(target=count, args=(N/2,))
. . .
       t1.start()
. . .
       t2.start()
. . .
    t1.join()
. . .
       t2.join()
. . .
. . .
>>> %time count_threaded(N)
CPU times: user 3.18 s, sys: 15.3 ms, total: 3.19 s
Wall time: 3.22 s
```

This is probably not really what you expected...



### THREADS FIGHTING FOR THE GIL



By David M Beazley @PyCON'2010: http://dabeaz.com/GIL/gilvis

### OS X: 4 threads on 1 CPU (Python 2.6)



### THREADS FIGHTING FOR THE GIL

OS X: 4 threads on 4 CPUs (Python 2.6)



By David M Beazley @PyCON'2010: http://dabeaz.com/GIL/gilvis



## **OBJECT, NO REAL PARALLEL EXECUTION OF CODE...**

HOW COULD PYTHON EVER COMPETE WITH ALL THOSE SUPER FAST C/C++/FORTRAN SOFTWARE?

**OK, HUGE OVERHEAD FOR EVERY SINGLE** 

### **C-EXTENSIONS AND INTERFACES TO "FAST" LANGUAGES LIKE** C/C++/FORTRAN!

**DO THE HEAVY STUFF IN THE BACKGROUND.** 

# THOSE CAN RELEASE THE GIL AND

### A DUMB SPEED COMPARISON Calculating the mean of 1000000 randomly generated numbers.

### pure Python

```
>>> def mean(numbers):
... return sum(numbers) / len(numbers)
...
>>> numbers = list(range(1_000_000))
>>> %timeit mean(numbers)
2.98 ms ± 21.7 µs per loop (mean ± std. dev. of 7 runs, 100 loops each)
```

### Numba (~3x faster)

```
>>> @nb.njit
... def numba_mean(numbers):
        s = 0
. . .
        N = len(numbers)
. . .
        for i in range(N):
. . .
            s += numbers[i]
. . .
        return s/N
. . .
. . .
>>> numbers = np.random.random(1_000_000)
>>> %timeit numba_mean(numbers)
1.03 ms ± 35.2 μs per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

### NumPy (~16x faster)

>>> import numpy as np

>>> numbers = np.random.random(1\_000\_000)

>>> %timeit np.mean(numbers) 190 μs ± 1.35 μs per loop (mean ± std. dev. of 7 runs, 10,000 loops each)

### Julia (~16x faster)

```
julia> numbers = rand(1_000_000);
julia> using BenchmarkTools
julia> @benchmark mean($numbers)
BenchmarkTools.Trial: 10000 samples with 1 evaluation.
Range (min ...max): 182.292 μs ...230.167 μs GC (min ...max): 0.00% ...0.00%
Time (median): 183.166 μs GC (median): 0.00%
Time (mean ± σ): 183.927 μs ± 2.457 μs GC (mean ± σ): 0.00% ± 0.00%
182 μs Histogram: log(frequency) by time 195 μs <
Memory estimate: 0 bytes, allocs estimate: 0.
```



### CRAZY LLVM COMPILER OPTIMISATIONS Summing up consecutive numbers from 0 to N=100,000,000

### pure Python

>>> def simple\_sum(N): s = 0 . . . for i in range(1, N+1): . . . s += i . . . return s . . . . . . >>> %timeit simple\_sum(N) 3.02 s ± 56.6 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

### **Numba** (~34000000x faster)

```
>>> @nb.njit
... def simple_sum(N):
       s = 0
       for i in range(1, N+1):
. . .
         s += i
. . .
        return s
. . .
. . .
>>> %timeit simple_sum(N)
85.6 ns ± 0.602 ns per loop (mean ± std. dev. of 7 runs, 10,000,000 loops each)
```

### $\Rightarrow$ N\*(N+1)/2 = 5000000050000000

NumPy (~75x faster)

>>> np\_numbers = np.arange(1, N+1, dtype=np.int64)

>>> %timeit np.sum(np\_numbers) 38 ms ± 27.8 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)

### **Julia** (~150000000x faster)

ulia> function simple_sum(N)	julia>	0code_na .sectio	ative debuginfo=:none si on      TEXT,text,	mple_sum(1_000_000 regular,pure_inst
5 - 0		.build	_version macos, 13, 0	
for i ∈ 1:N		.globl	_julia_simple_sum_1025	; Beg
s += i		.p2ali	gn <mark>2</mark>	
end	_julia_	_simple_s	sum_1025:	; @julia_simple_
		.cfi_s	tartproc	
S	; %bb.0	9:		; %top
end		cmp	x0, #1	; =1
<pre>imple_sum (generic function with 1 method)</pre>		b.lt	LBB0_2	
	; %bb.1	1:		; %L18.preheade
ulies Obenehment, simple sum(N)		bic	x8, x0, x0, asr #63	
ulla> wbenchmark simple_sum(N)		sub	x9, x8, #1	; =1
enchmarkTools.Trial: 10000 samples with 997 evaluations.		sub	x10, x8, #2	; =2
Range (minmax): 19.015 ns489.468 ns 🕴 GC (minmax): 0.00% .		mul	x11, x9, x10	
Time (median): 19 267 ns GC (median): 0 00%		umulh	x9, x9, x10	
$T_{\text{ineq}} = \left( \frac{1}{1000} - \frac{1}{1000} \right) = \frac{1}{1000} = \frac{1}{1000$		extr	X9, X9, X11, #1	
$11me (mean \pm \sigma)$ : 20.046 ns ± 6.731 ns ; 60 (mean ± $\sigma$ ): 0.46% ±	±	add	X8, X9, X8, ISI #1	1
		sub	x0, x8, #1	; =
	1000.0	ret		
	LDD0_2:		×0 <b>#0</b>	
		mov	x0, #0	
19 ns Histogram: log(frequency) by time 23.9 ns <		ret	ndanaa	
		.CTI_e	nuproc	
Memory estimate: <mark>16 bytes</mark> , allocs estimate: 1.				
······································				





# *"WHY WE CAN'T HAVE NICE THINGS"*

**AKA VECTORISE FOR YOUR LIFE!** 

## CLASSES TO STRUCTURE CODE

- Classes and cleverly designed class hierarchies make it easier to understand code, structure and architecture
- It's tempting to take "object-oriented" seriously and create a detailed model of your application
- The more fine-grained it gets, the bigger the impact on performance
- Imagine a neutrino detector which consists of PMTs to measure photons, each triggered signal represented by an instance of the "Hit"-class:

```
>>> class Hit:
        def __init__(self, pmt_id, time, tot, triggermask):
. . .
             self.pmt_id = pmt_id
• • •
             self.time = time
. . .
             self.tot = tot
• • •
             self.triggermask = triggermask
• • •
• • •
```



## CLASSES TO STRUCTURE CODE

- Now imagine that a KM3NeT detector (the experiment I work on) produces
- this class due to the large overhead

```
>>> class Hit:
         def __init__(self, pmt_id, ti
. . .
              self.pmt_id = pmt_id
. . .
              self.time = time
. . .
             self.tot = tot
• • •
              self.triggermask = trigge
. . .
• • •
```

several thousands of hits per event, with a trigger rate of more than 200Hz

• It's certainly a bad idea to have an array of thousands of hits with instances of

		>>> hits		
		<pre>[<mainhit< pre=""></mainhit<></pre>	at	0x10846e430>,
	trinnerm	<mainhit< td=""><td>at</td><td>0x108ba2100&gt;,</td></mainhit<>	at	0x108ba2100>,
me. tot.		<mainhit< td=""><td>at</td><td>0x108422490&gt;,</td></mainhit<>	at	0x108422490>,
	or ±9901 m	<mainhit< td=""><td>at</td><td>0x1085d66a0&gt;,</td></mainhit<>	at	0x1085d66a0>,
		<mainhit< td=""><td>at</td><td>0x1085d6b80&gt;,</td></mainhit<>	at	0x1085d6b80>,
		<mainhit< td=""><td>at</td><td>0x1085d61c0&gt;,</td></mainhit<>	at	0x1085d61c0>,
		<mainhit< td=""><td>at</td><td>0x108c366a0&gt;,</td></mainhit<>	at	0x108c366a0>,
rmask		<mainhit< td=""><td>at</td><td>0x108c36850&gt;,</td></mainhit<>	at	0x108c36850>,
		<mainhit< td=""><td>at</td><td>0x108c36460&gt;,</td></mainhit<>	at	0x108c36460>,
		<mainhit< td=""><td>at</td><td>0x108c36070&gt;,</td></mainhit<>	at	0x108c36070>,
		<mainhit< td=""><td>at</td><td>0x108c362b0&gt;,</td></mainhit<>	at	0x108c362b0>,



## **STRUCT OF ARRAYS** VS ARRAY OF STRUCTS

- arrays of the same size, one for each field

```
>>> hits = np.array([(1, 2, 3, 4), (5, 6, 7, 8), (9, 10, 11, 12)], dtype=hit_dtype).view(np.recarray)
>>> hits.time
array([ 2., 6., 10.])
>>> hits.tot
array([ 3, 7, 11], dtype=uint8)
>>> hits[2].time
10.0
```

• Instead of having an array of hits, use a single instance with multiple

NumPy offers the "recarray" functionality for convenient attribute access

>>> hit\_dtype = np.dtype([("pmt\_id", np.uint32), ("time", np.float64), ("tot", np.uint8), ("triggermask", np.uint64)])



## **STRUCT OF ARRAYS VS ARRAY OF STRUCTS**

- field
- NumPy offers the "**recarray**" functionality for convenient attribute access
- Similar to the way "**Panda DataFrames**" work (columns)

```
>>> hits = np.array([(1, 2, 3, 4), (5, 6, 7, 8), (9, 10, 11, 12)], dtype=hit_dtype).view(np.recarray)
>>> hits.time
array([ 2., 6., 10.])
                                              ......
>>> hits.tot
                                             $(SIGNATURES)
array([ 3, 7, 11], dtype=uint8)
>>> hits[2].time
                                              ......
10.0
```

• Instead of having an array of hits, use a single instance with multiple arrays of the same size, one for each

>>> hit\_dtype = np.dtype([("pmt\_id", np.uint32), ("time", np.float64), ("tot", np.uint8), ("triggermask", np.uint64)])



Return only triggered hits.

triggered(hits::Vector{T}) where {T<:AbstractHit} = filter(h->h.trigger\_mask > 0, hits)





"NUMERIC" IN 1995, "NUMPY" IN 2006 (JUST KICKED IN WHEN I STARTED STUDYING PHYSICS)

## NUMPY

### NumPy is the fundamental package for scientific computing with Python.

- gives us a powerful N-dimensional array object: ndarray
- broadcasting functions
- tools for integrating C/C++ and Fortran
- most of the scientific libraries build upon NumPy

linear algebra, Fourier transform and random number capabilities



## NUMPY: ndarray

>>> a



- >>> a = np.arange(6)
- array([0, 1, 2, 3, 4, 5])

- >>> a.dtype dtype('int64') >>> a.ndim >>> a.shape (6,)
- Contiguous array in memory with a fixed type, no pointer madness!
  - C/Fortran compatible memory layout,
    - so they can be passed to those
      - without any further efforts.



### NUMPY: ARRAY OPERATIONS AND ufuncs

a * 23				
array([	0,	23,	46,	69
a**a				
array([	1,	1,	4,	27

### a **ufunc**, which can operate both on scalars and arrays (element-wise)

np.exp	<mark>(</mark> a)	
array([	1. ,	2.718
	54.59815003,	148.413
	54.55015005,	140.41



easy and intuitive element-wise operations

7.3890561 , 20.08553692, 328183, 31591 ])



## **RESHAPING ARRAYS**

a = np.arange(6)а



No rearrangement of the elements in memory but setting the iterator limits internally!

a.reshape(2, 3)

array([[0, 1, 2], [3, 4, 5]])



## **RESHAPING ARRAYS IS CHEAP**



>>> %timeit a.reshape(100, 5\_000, 20) 131 ns ± 1.86 ns per loop (mean ± std. dev. of 7 runs, 10,000,000 loops each)



# OK, ALL FINE, NUMPY TO RULE THEM ALL?

## "HIDDEN" ALLOCATIONS

```
import resource
import sys
import numpy as np
```

```
def peak_memory_usage():
    """Return peak memory usage in MB""""
    mem = resource.getrusage(resource.RUSAGE_SELF).ru_maxrss
    factor_mb = 1 / 1024
    if sys.platform == "darwin":
        factor_mb = 1 / (1024 * 1024)
        return mem * factor_mb
```

```
def main():
    print(f"All libraries loaded: {peak_memory_usage()} MB")
```

```
N = 1_{000}_{000}
```

```
a = np.random.rand(N)
print(f"'a' allocated: {peak_memory_usage()} MB")
```

```
b = np.random.random(N)
print(f"'b' allocated: {peak_memory_usage()} MB")
```

```
c = 2*a + 3*b
print(f"'c' calculated: {peak_memory_usage()} MB")
```

```
if __name__ == "__main__":
    main()
```

tamasgal@silentbox:tmp/missing\_python\_intro
00:16:38 > python3 hidden\_allocations.py
All libraries loaded: 28.578125 MB
'a' allocated: 36.28125 MB
'b' allocated: 43.921875 MB

'c' calculated: 67.59375 MB

### 1000000 float64s are 8 MB

+8 MB after "a" is allocated

+8 MB after "b" is allocated

+24 MB = 8+8+8 MB after "c" is calculated

## c = 2\*a + 3\*b



## AVOIDING UNNEEDED ALLOCATIONS

```
import resource
import sys
import numpy as np
def peak_memory_usage():
    """Return peak memory usage in MB"""
    mem = resource.getrusage(resource.RUSAGE_SELF).ru_maxrss
    factor_mb = 1 / 1024
    if sys.platform == "darwin":
        factor_mb = 1 / (1024 * 1024)
    return mem * factor_mb
def main():
    print(f"All libraries loaded: {peak_memory_usage()} MB")
   N = 1_{000}_{000}
    a = np.random.rand(N)
    print(f"'a' allocated: {peak_memory_usage()} MB")
    b = np.random.random(N)
    print(f"'b' allocated: {peak_memory_usage()} MB")
    np.multiply(a, 2, out=a)
    print(f"'2*a' calculated: {peak_memory_usage()} MB")
    np.multiply(b, 3, out=b)
    print(f"'3*b' calculated: {peak_memory_usage()} MB")
    c = np.add(a, b)
    print(f"'c' calculated: {peak_memory_usage()} MB")
if __name__ == "__main__":
```

main()

tamasgal@silentbox:tmp/missing\_python\_intro
14:53:23 > python3 numpy\_out.py
All libraries loaded: 28.703125 MB
'a' allocated: 36.421875 MB
'b' allocated: 44.0625 MB
'2\*a' calculated: 44.0625 MB
'3\*b' calculated: 44.0625 MB
'c' calculated: 51.703125 MB

### Reusing a and b with "out="

a = a + b vs. np.add(a, b, out=a)

Only a single 8 MB allocation is needed for c (mutating the a and b!)



## MEMORY AND PERFORMANCE PROBLEMS

- The code may work for small sample files
- Quickly escalates the peak memory usage if not handled with care
- Users blindly perform chains of transformations
- Trigger redundant loops instead using kernels and a single loop
- Unnecessary memory allocations for temporary data
- Pandas e.g. uses NumPy behind the scenes and is a constant source of scaling issues (hello Pandas concat/append/join/etc.)



## MASKING/SLICING IS THE ROOT OF ALL EVIL

- Masks are often used in NumPy to reduce data
- If Python loops were fast, we could simply iterate through hits (array of structs) and select them based on a condition => minimal memory footprint, single loop
- Potential extra loops and allocations which might blow the memory
- Instead, masks are tempting due do the nice and easy syntactic sugar

selected hits = hits[mask]





## WRITE YOUR CODE SO THAT IT ITERATES THROUGH A CONFIGURABLE SIZE OF CHUNKS!

### (FOR EXAMPLE X EVENTS INSTEAD OF FILE-BY-FILE)



**ROUTINES FOR THE FAST EVALUATION OF ARRAY EXPRESSIONS ELEMENT-WISE BY USING A VECTOR-BASED VIRTUAL MACHINE.** 

## NUMEXPR USAGE EXAMPLE

import numpy as np import numexpr as ne

a = np.arange(5)b = np.linspace(0, 2, 5)

ne.evaluate("a\*\*2 + 3\*b")

array([ 0., 2.5, 7., 13.5, 22.])



## NUMEXPR SPEED-UP

NumPy: 2 \* a \* \* 3 - 4 \* a \* \* 5 + 6 \* np.log(a)

Numexpr with 4 threads: ne.set\_num\_threads(4)

ne.evaluate("2 \* a\*\*3 - 4 \* a\*\*5 + 6 \* log(a)")~10x faster (with 4 threads) ...minimal memory footprint

a = np.random.random(1000000)

 $82.4 \text{ ms} \pm 1.88 \text{ ms} \text{ per loop}$ 

7.85 ms  $\pm$  103 µs per loop



### **NUMEXPR - SUPPORTED OPERATORS**

- Logical operators: &, , ~
- Comparison operators: <, <=, ==, !=, >=, >
- Unary arithmetic operators: -
- Binary arithmetic operators:

### +, -, \*, /, \*\*, %, <<, >>



### **NUMEXPR - SUPPORTED FUNCTIONS**

- where(bool, number1, number2): number -- number1 if the bool condition is true, number2 otherwise.
- {**sin**,**cos**,**tan**}(float|complex): float|complex -- trigonometric sine, cosine or tangent.
- {arcsin, arccos, arctan}(float|complex): float|complex -- trigonometric inverse sine, cosine or tangent.
- arctan2(float1, float2): float -- trigonometric inverse tangent of float1/float2.
- {**sinh**,**cosh**,**tanh**}(float|complex): float|complex -- hyperbolic sine, cosine or tangent.
- {arcsinh,arccosh,arctanh}(float|complex): float|complex -- hyperbolic inverse sine, cosine or tangent.
- {log,log10,log1p}(float|complex): float|complex -- natural, base-10 and log(1+x) logarithms.
- {**exp**,**expm1**}(float|complex): float|complex -- exponential and exponential minus one.
- **sqrt**(float|complex): float|complex -- square root.
- **abs**(float|complex): float|complex -- absolute value.
- **conj**(complex): complex -- conjugate value.
- {**real**,**imag**}(complex): float -- real or imaginary part of complex.
- **complex**(float, float): complex -- complex from real and imaginary parts.
- contains(str, str): bool -- returns True for every string in `op1` that contains `op2`.
- **sum**(number, axis=None): Sum of array elements over a given axis. Negative axis are not supported.
- prod(number, axis=None): Product of array elements over a given axis. Negative axis are not supported.



## AVOIDING EXTRA ALLOCATIONS WITH NUMEXPR

```
import resource
import sys
import numexpr as ne
import numpy as np
def peak_memory_usage():
    """Return peak memory usage in MB"""
   mem = resource.getrusage(resource.RUSAGE_SELF).ru_maxrss
   factor_mb = 1 / 1024
   if sys.platform == "darwin":
        factor_mb = 1 / (1024 * 1024)
    return mem * factor_mb
def main():
   print(f"All libraries loaded: {peak_memory_usage()} MB")
   N = 1_{000}_{000}
   a = np.random.rand(N)
   print(f"'a' allocated: {peak_memory_usage()} MB")
   b = np.random.random(N)
   print(f"'b' allocated: {peak_memory_usage()} MB")
   c = ne.evaluate("2*a + 3*b")
   print(f"'c' calculated: {peak_memory_usage()} MB")
if __name__ == "__main__":
   main()
```

tamasgal@silentbox:tmp/missing\_python\_intro
00:37:25 > python3 numexpr\_example.py
All libraries loaded: 29.40625 MB
'a' allocated: 37.109375 MB
'b' allocated: 44.75 MB
'c' calculated: 52.59375 MB

## Only a single **8 MB** allocation is needed **for c**





**BY JIM PIVARSKI (SCIKIT-HEP)** 

### AWKARD ARRAY MOTIVATION

- NumPy arrays are rectangular tables or tensors: cannot express variable-length structures
- Tree-like data (very common in HEP) is difficult to express with NumPy arrays -- in an efficient way
- Speed and performance are crucial
- Easy to use and interactive interfaces for commonly used operations like cuts and aggregations



## AWKARD ARRAY

- Written in Python and C++
- Has Numba support to take it to the next level!
- Supports arbitrary tree representations with as many jagged/ragged structures as you need
- Offers lots of functions to work with ragged/ jagged data

```
In [1]: import awkward as ak
In [2]: arr = ak.Array([[1,2,3], [4,8], [6,7,8,9]])
In [3]: arr
Out[3]: <Array [[1, 2, 3], [4, 5], [6, 7, 8, 9]] type='3 * var * int64'>
In [4]: arr[:,0]
Out[4]: <Array [1, 4, 6] type='3 * int64'>
In [5]: ak.mean(arr, axis=1)
Out[5]: <Array [2, 4.5, 7.5] type='3 * ?float64'>
```



## AWKARD ARRAY

- All kinds of multiply nested structures are understood and "type stable" •
- Allows fancy indexing over multiple nested elements •
- **Database-like** operations
- Feels like Pandas, with awkwardly jagged arrays •

In [1]:	import <mark>awkward</mark> as <mark>ak</mark>
In [2]:	arr = ak.Array([{"pos_x": [3, 45, 65], "pos_y
In [3]:	arr
Out[3]:	<array 45,="" 65,="" 6]<="" [3,="" [5,="" [{pos_x:="" pos_y:="" td=""></array>
In [4]:	arr.pos_x
Out[4]:	<Array [[3, 45, 65], [1, 3]] type='2 * var *</td>



### AWKWARD ARRAY

### • A very nice introduction by Jim himself (just search for "awkward array" on YouTube)

### But Numpy doesn't have anything for unequal-length lists



Some libraries\* can represent arrays of unequal-length arrays, known as "jagged" or "ragged" arrays. \*Apache Arrow, XND, TensorFlow, Zarr (genetics), and ROOT (particle physics)  $_{17/43}$ 

https://www.youtube.com/watch?v=2NxWpU7NArk



collection 2



derived collection









## WHEN THINGS ARE GETTING NORECOMPLICATED THAN



2\*A + 3\*B



**THE JIT (LLVM) COMPILER FOR PYTHON** 



## NUMBA

Numba is a **compiler** for Python array and numerical functions that gives you the power to speed up code written directly in Python.

- uses **LLVM** to boil down pure Python code to **JIT optimised machine code**
- only accelerates selected functions decorated by yourself
- native code generation for CPU (default) and GPU
- integration with the Python scientific software stack (thanks to NumPy)
- runs side by side with regular Python code or third-party C extensions and libraries
- great **CUDA** support
- N-core scalability by releasing the GIL (beware: no protection from race conditions!)
- create NumPy ufuncs with the @[gu]vectorize decorator(s)
- unfortunately: limited support of data structures



## FROM SOURCE TO RUNTIME





## NUMBA JIT-EXAMPLE

numbers = np.arange(1000000).reshape(2500, 400)

```
def sum2d(arr):
   M, N = arr.shape
    result = 0.0
    for i in range(M):
        for j in range(N):
            result += arr[i,j]
    return result
```

 $289 \text{ ms} \pm 3.02 \text{ ms} \text{ per loop}$ 

```
@nb.jit
def sum2d_jit(arr):
    M, N = arr.shape
    result = 0.0
    for i in range(M):
        for j in range(N):
            result += arr[i,j]
    return result
```

 $2.13 \text{ ms} \pm 42.6 \mu \text{s}$  per loop

### ~135x faster, with a single line of code



## NUMBA VECTORIZE-EXAMPLE

a = np.arange(1000000, dtype='f8')b = np.arange(1000000, dtype='f8') + 23

NumPy: np.abs(a - b) / (np.abs(a) + np.abs(b))23 ms  $\pm$  845 µs per loop

### Numba avectorize:

anb.vectorize def nb\_rel\_diff(a, b): **return** abs(a - b) / (abs(a) + abs(b))

rel\_diff(a, b)

3.56 ms  $\pm$  43.2 µs per loop





## **PYTHON IS YOUR EVERYDAY HAMMER**



## **IF YOU ARE HOLDING A HAMMER, EVERYTHING LOOKS LIKE A NAIL**



## CHOOSE YOUR TOOLS WITH CARE

- Python is a powerful language and it can be used for many different tasks
- However: it's easy to write code with horrible performance
- It's even easier to write code which simply does not scale
- In contrary to Python's expressiveness, high-performant Python code is often neither nicely readable, nor easily maintainable
- High-level Python APIs can be very useful but usually act as a barrier
- Low-level development requires multiple languages and technology stacks, which increase complexity greatly
- You have to be more than a Python expert to write high-performant "Python code"
- Keep all this in mind and give Julia a go for scientific computing



## PERSONAL REMARK

- I am a big fan of **Julia** (surprise!)
- The Julia language is built for scientific computing, it "feels like Python and runs like C".
- No fear of writing for-loops
- Define and use your own types and type hierarchy to create expressive code
- Interactive prototyping, just like in Python with a REPL
- Even naively written code is often very close to optimal performance
- You can still use your **Jupyter-notebooks** (fun fact: "Ju" in Jupyter stands for Julia)
- Easy package management and deployment (you can easily set up your own package registry)
- **Reproducible** environments are a built-in feature (crucial for reproducible science)







https://discourse.julialang.org/t/help-testing-julia-tree-sitter-mode-in-emacs/928

### My own Emacs setup

### 30 struct NoRecoTrack<:AbstractRecoTrack end

ruct SinglebuPa			Formational Italia					
at64	×J Fil	le Edit Selection View Go Run I	lerminal Help	run_reconstruction	on.ji - workspace (Workspace	e) - Visual Studio Code		] 🗖 []] 0%
	ر <mark>ل</mark> م	JULIA ····	run_reconstruction.jl	reconstruction.jl	👶 analyze_results.jl м	♣ plo ▷ ∽ 🖽	···· ♣ gaussian_model.jl	
	````	V WORKSPACE ····	CEIReconstruction > scri	pts > iodopyridine > 👶 run	_reconstruction.jl >		cn.stce = 200 radicese = [1, 7, 20] row.usete.strothym	
	Q	> [@] ans DataFrame	88 hydrogen	= 1			mar uturaliant order = 3 licet = 5.28845400772312 weights = [1, 0.5, 0.28]	
		> 😚 attribute_carbon_pairs attribut	89 carbonA	= 2			Component 1 Component 2 Component 3 Component 4 Component 5 EVR = 0.32 EVR = 0.18 EVR = 0.13 EVR = 0.1 EVR = 0.073	Component 6 0
Ua 104		# carbonA 2	91 nitrogen	= 4				
		# carbonB 3	92 iodine =	5				1 1
Fit	æ	> [ ] carbons2d Matrix{Float64} 🖬	93 04 forkin	avec(momenta 3)				
ingleDUParams		<pre>&gt; [] charges Matrix{int32} with 5×1 &gt; [] dims UnitBande{Int64} with 3 el</pre>	95 mome	nta[:, :, k] = reframe	e(momenta[:,	Vine, rogen)		
itial::SingleD	Ē	> [] elems Vector{Element} with 5 e	96 end					
at64		> [ ] folders Vector{String} with 62 e	97 08 full mom	onto - conv(momonto)				· · ·
ed hits::Vecto	Ш	> [ ] full_momenta Array{Float64, 3	99 end 3×5×111	.280 Array{Float64, 3}			-200 0 200 -200 0 200 -200 0 200 -200 0 200 -200 0 200 -200 0 200 - x x x x x x x	-200 0 200 -200 x
:Model		> [e] group HDF5.Group	100					
	A	> [] guess_params Vector{Float64}	101 include( <u>"cle</u>	anup_data.jl") 🔽				
		> [ ] guessΣ Matrix{Float64} with 21	102 103 # Compute th	e guess by fitting a (	Gaussian mixture to th	ne carbons and		
		> [@] h5file HDF5.File	104 # a simple G	ausian to iodine and r	nitrogen			
		# hydrogen 1	105 begin	d - hcat(momenta[1,2	carbonA :1 mor	Cherenkov Photon	Yield ar $\times$ $\equiv$ ORCA4 Single DU RBR-MC $\times$ $\blacktriangleright$	ORCA4 S
ES)	~	# i 7	107 C gmm =	GMM(5, carbons2d  > pe	ermutedims ; kind	+ % 6 6	▶ ■ C → Code ∨	
	0	> [ ] ii UnitRange{Int64} with 2 elem	108				detection unit	
he prefit algo	00	# iodine 5	$109  guess\Sigma =$	Matrix{Float64}(I, 3	n_atoms, 3n_atoms			
		> [ ] joint_proba Vector{Float64} wit	110 guessµ =			[41]: duplot(	df, :χ²; label="\\$\\chi^2\\$")	
refit(hits::Ve	~		112 for (k,	C) in enumerate(sim_ca	arbons)	18		
ngth(hits)		> DOCUMENTATION	113 i =	C.index			······	
0		V PLOT NAVIGATOR	114 11 – 115 gues	$s\Sigma[ii, ii] = covars(C)$	gmm)[k]	16 15	1	60
		Plot 43 11:45:54 AM	116 gues	<pre>sµ[ii] = means(C_gmm)</pre>	[k, :]	14		
LU.U, U.U, U.U		Plot 42 11:45:54 AM	PROBLEMS 53 OUTPU	T DEBUG CONSOLE TERM	IINAL	13 $12$		
Lmax(I, (n.tot		Plot 41 11:45:54 AM				11 동 10		40 ~
= [n.tot for n - [b multiplic		Plot 40 11:45:52 AM	[ Info: Initializi	ng GMM. 5 Gaussians di	ag covariance 2 (	° j	······ • • • • • • • • • • • • • • • •	$\times^{40}$ $\times$
- [n.multiplic]		Plot 39 11:45:52 AM	K-means converged	with 21 iterations (ob	jv = 1.620860257	87		
nhounds for k		Plot 37 11:45:50 AM	Info: K-means wi	th 5000 data points us s per parameter	ing 21 iteration	$\begin{array}{c} 6\\ 5\end{array}$		20
if i == k		Plot 36 11:45:50 AM	ERROR: UndefVarErr	pr: n_atom not defined		4		20
il/one/fit i		Plot 35 11:45:49 AM	Some of the types	have been truncated in	the stacktrace i	$\frac{3}{2}$		
A.JI/SPC/IIC.J.		Plot 34 11:45:48 AM	in the stack trace	, evaluate `TruncatedS	tacktraces.VERB0	1		
		Plot 33 11:45:47 AM	Ctacktraca			2	detection unit.	
	8	Plot 32 11:45:47 AM	[1] top-level sco	be	•			
	~~~	Plot 31 11:45:45 AM	<pre>@ ~/Documents/C</pre>	EIReconstruction/scrip	ts/iodopyridine/_	[17]: using (	SV	
	201	Plot 29 11:45:44 AM	∘julia>∏			using 1	ables	
	્ર ૪	master* $\bigcirc 0 \downarrow 1^{\uparrow} \gg \otimes 0 \land 0 \bigcirc 53$	Git Graph Julia env: CEIRed	construction INSERT		using F	ProgressMeter	-
						using N	Ierca	
	D	rovided by Bon	oît <b>Dichar</b> c			using S	StatsBase	
		I UVILLEU DY DEIT		4		using (	)ptim	
						using L	andauDistribution	JUP
						using L	aTeXStrings	
_								
						using F	PrettyTables	
						struct	TimePecMinimiser <: Eunction	
						tin	neresiduals::Vector{Float64}	
						bir	IS	
						fur	<pre>iction TimeResMinimiser(timeresiduals::\     limit = Int(round(At ( 2))</pre>	Vector{T}
							bins = -limit:limit	
		200 C					$t = filter(x -> abs(x) < \Delta t / 2, timeres$	siduals)
							<pre>h = fit(Histogram, t, bins)</pre>	
						0.00	new(n.weignts, bins[1:end-1])	
						end		
0.2/1								
261						functio	on (t::TimeResMinimiser)(p)	
						LA,	Lµ, Lơ, GA, Gµ, Gơ, Offset = p Gơ < $0$      $\sigma < 4$    offset < $0$    $GA =$	
10/11						One of r	ny lunyter prototypir	
<u>13/11</u>								19 30



## "WITH GREAT POWER COMES GREAT RESPONSIBILITY"

### - UNCLE BEN

## A FINAL WORD ON GREEN CODING/COMPUTING

- In HPC, we can easily waste energy with inefficient code
- It's just a matter of a few keystrokes to laun thousands of computing jobs
- It's also your responsibility to learn how to resources efficiently and with care

×	source	secs	mem	gz	cpu secs	cpu load					Но	w n	าล
1.0	Classic Fortran #3	0.71	10,984	638	2.84	100% 98% 98% 100%	te 100	ust	S	+++++++++++++++++++++++++++++++++++++++	<u>llia</u>	Ē	an
1.0	Rust #7	0.72	11,128	932	2.84	100% 100% 100% 100%	fas	1 R	Ű	ġ	JL	Ζ.	ort
1.0	Rust #4	0.72	11,124	817	2.84	100% 100% 100% 98%	کم 100 الم	ust	S	÷ U	ulia	÷ U	U.U.
1.0	Rust #5	0.72	11,124	1055	2.84	100% 100% 100% 100%			0 D	+	Ē	Ш	ass
1.0	Chapel #2	0.73	10,948	335	2.89	100% 98% 98% 98%	d Se	ust	U	+6	ulia	Ζ.	Ο
1.5	Julia #4	1.09	185,772	429	3.67	99% 78% 79% 79%	s 10	~	g	+	Ť.	U H	an
1.8	Julia #2	1.26	192,016	370	4.14	76% 76% 76% 98%	n ela	tust	9	Ŭ	ulia	T	9 1 1 9
2.0	<b>Go</b> #4	1.42	11,244	548	5.67	99% 100% 99% 99%	Iran	- <del>"</del>	-	+	T,	, <b>₽</b> ,	<u> </u>
2.0	Swift #3	1.43	11,344	601	5.68	100% 100% 99% 99%	prog	<b>H</b>	Ĺ	Ē	Ŧ	Ê.	Ţ
2.0	<b>C</b> gcc #3	1.44	11,392	463	5.70	100% 99% 99% 99%		ber	nchm	ark	s ga	me	

https://benchmarksgame-team.pages.debian.net

Table 4. Normalized global results for Energy, Time, and Memory

	Total							
ab		Energy		Time				
CN	(c) C	1.00	(c) C	1.00	(c) Pascal			
	(c) Rust	1.03	(c) Rust	1.04	(c) Go			
	(c) C++	1.34	(c) C++	1.56	(c) C			
	(c) Ada	1.70	(c) Ada	1.85	(c) Fortran			
	(v) Java	1.98	(v) Java	1.89	(c) C++			
	(c) Pascal	2.14	(c) Chapel	2.14	(c) Ada			
	(c) Chapel	2.18	(c) Go	2.83	(c) Rust			
use these	(v) Lisp	2.27	(c) Pascal	3.02	(v) Lisp			
	(c) Ocaml	2.40	(c) Ocaml	3.09	(c) Haskell			
	(c) Fortran	2.52	(v) C#	3.14	(i) PHP			
	(c) Swift	2.79	(v) Lisp	3.40	(c) Swift			
	(c) Haskell	3.10	(c) Haskell	3.55	(i) Python			
	(v) C#	3.14	(c) Swift	4.20	(c) Ocaml			
	(c) Go	3.23	(c) Fortran	4.20	(v) C#			
any times slower?	(i) Dart	3.83	(v) F#	6.30	(i) Hack			
	(v) F#	4.13	(i) JavaScript	6.52	(v) Racket			
A NE NE NE	(i) JavaScript	4.45	(i) Dart	6.67	(i) Ruby			
Chan S Ch	(v) Racket	7.91	(v) Racket	11.27	(c) Chapel			
- 12 ee va	(i) TypeScript	21.50	(i) Hack	26.99	(v) F#			
L L L L L L L L L L L L L L L L L L L	(i) Hack	24.02	(i) PHP	27.64	(i) JavaScript			
	(i) PHP	29.30	(v) Erlang	36.71	(i) TypeScript			
A B B A B B A B B A B B A B B A B B B A B B B A B B B B B B B B B B B B B B B B B B B B	(v) Erlang	42.23	(i) Jruby	43.44	(v) Java			
	(i) Lua	45.98	(i) TypeScript	46.20	(i) Perl			
	(i) Jruby	46.54	(i) Ruby	59.34	(i) Lua			
╮╩┐凵匚╘┝┥┝┥┍┑┟┦║	(i) Duly	69.91	(i) Perl	65.79	(v) Erlang			
╢┥╦╎╢┋┰╞╝┰┋╽	(i) Python	75.88	(i) Python	71.90	(i) Dart			
	() Derl	79.58	(i) Lua	00.01	(i) Jruby			
01 Mar 2023 u64q					-			

Source: Energy Efficiency across Programming Languages, SLE'17 68



# THANKS

"PEOPLE ARE VERY OPEN-MINDED ABOUT NEW THINGS — AS LONG AS THEY'RE EXACTLY LIKE THE OLD ONES."

- CHARLES F. KETTERING