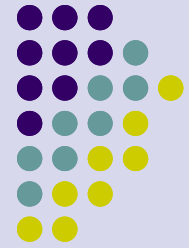


Prototyping and Developing GPU Accelerated Solutions with Python and CUDA

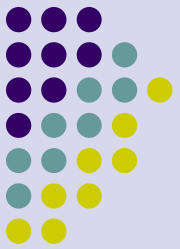


Luciano Martins

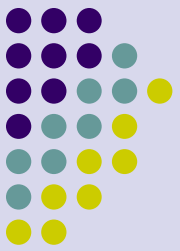
Principal Software Engineer
Oracle Corporation



Agenda

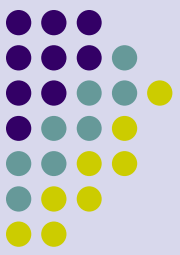


- Python introduction
- GPU programming
- Why python with GPU?
- Accelerating python
- Comparing codes with/without GPU support
- Summary



Python introduction

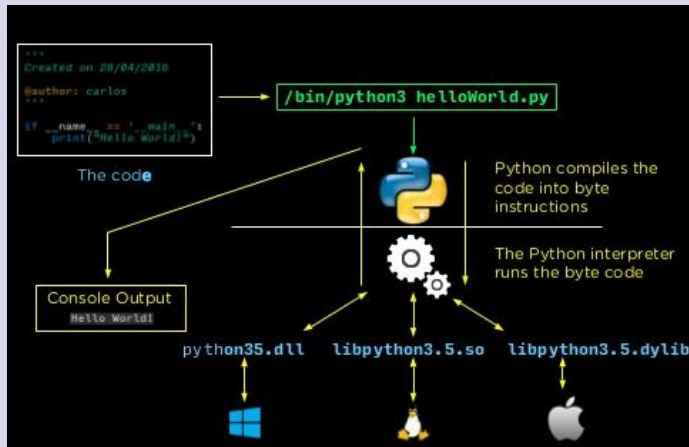
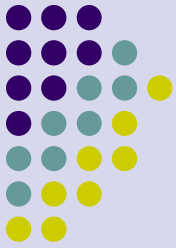
- created by Guido van Rossum in 1991
- "Zen of Python" which:
 - Beautiful is better than ugly
 - Explicit is better than implicit
 - Simple is better than complex
 - Complex is better than complicated
 - Readability counts
- interpreted language (CPython, JPython, ...)
- dynamically typed; based on objects



Python introduction

- small core structure:
 - ~30 keywords
 - ~80 built-in functions
- indention is a pretty serious thing
- a huge modules ecosystem
- binds to many different languages
- supports GPU acceleration via modules ←

Python introduction



Python Py

Python is a clear and powerful object-oriented programming language, comparable to Perl, Ruby, Scheme, or Java.

Votes
4.05K

Fans
3.35K

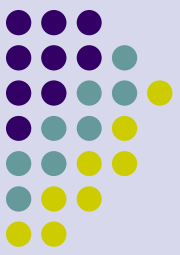
Stacks
3.71K

Integrations
38

Jobs New
3.05K

COMPANIES USING PYTHON

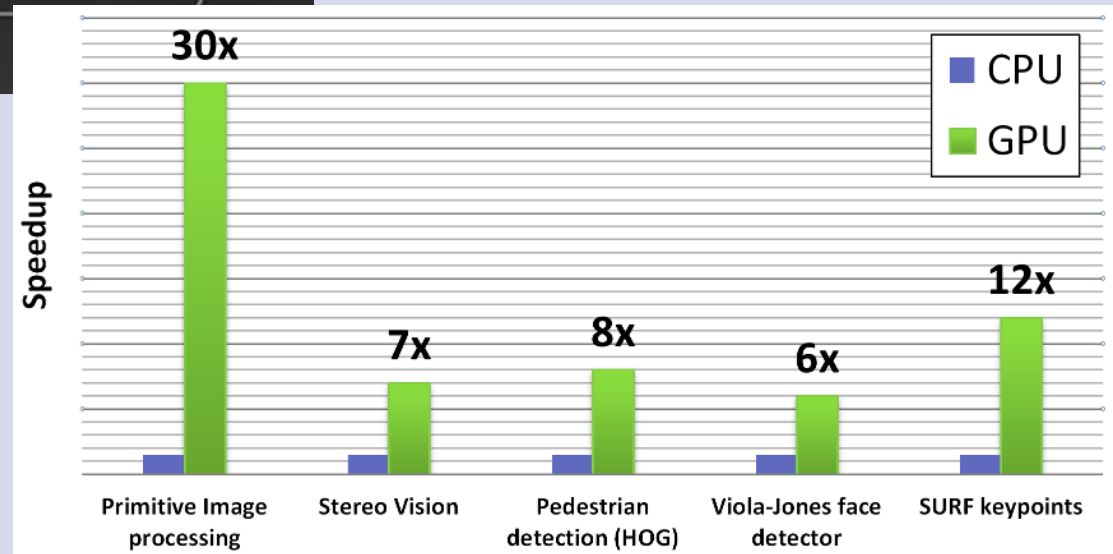
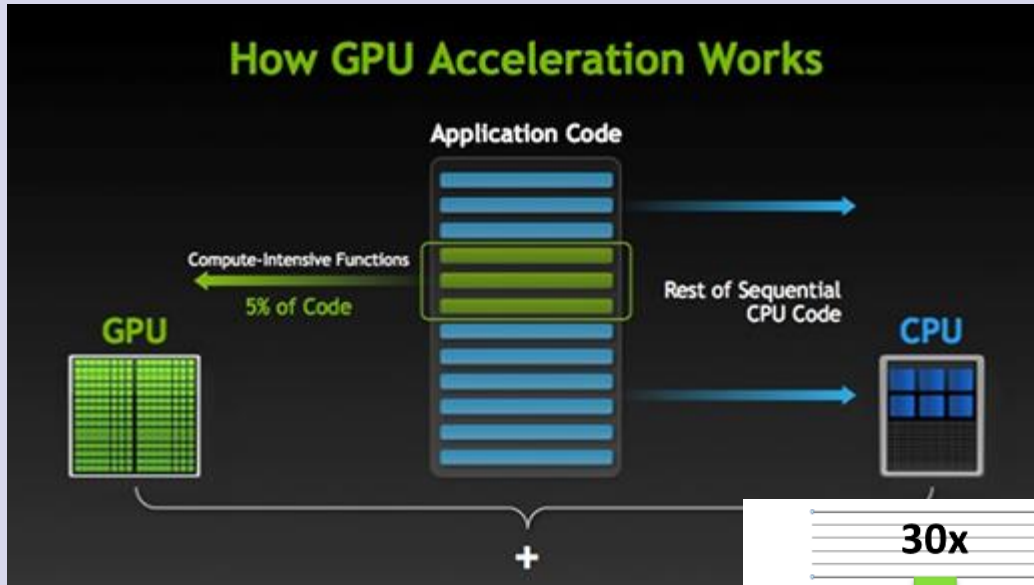
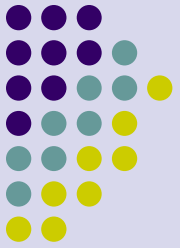
Language Rank	Types	Spectrum Ranking
1. Python		100.0
2. C		99.7
3. Java		99.5
4. C++		97.1
5. C#		87.7
6. R		87.7
7. JavaScript		85.6
8. PHP		81.2
9. Go		75.1
10. Swift		73.7

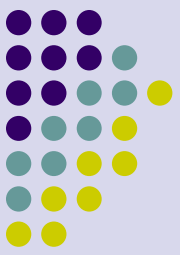


GPU programming

- “the use of a graphics processing unit (GPU) together with a CPU to accelerate deep learning, analytics, and engineering applications” (NVIDIA)
- most common GPU accelerated operations:
 - large vector/matrix operations (BLAS)
 - speech recognition
 - computer vision
 - way more

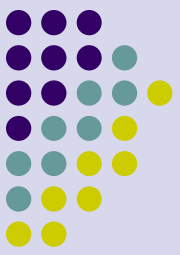
GPU programming





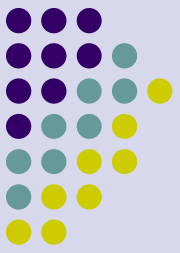
Why python with GPU?

- Interpreted languages has the reputation of being slow for high performance needs
- Python needs assistance for those tasks
- Keep the best of both scenarios:
 - Quick development and prototyping with python
 - Use high processing power and speed of GPU
- Can deliver quick results for complex projects
- Gives a business decision choice at the end



Accelerating python

- GPU + python projects are arising every day
- Accelerated code may be pure python or adding C code
- Focusing here on the following modules
 - PyCUDA
 - Numba
 - cudamat
 - cupy
 - scikit-cuda

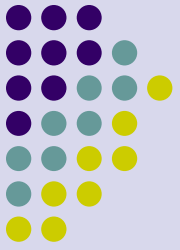


Accelerating python – PyCUDA

- A python wrapper to CUDA API
- Requires C programming knowledge (kernel)
- Gives speed to python – near zero wrapping
- Compiles the CUDA code copy to GPU
- CUDA errors translated to python exceptions
- Easy installation

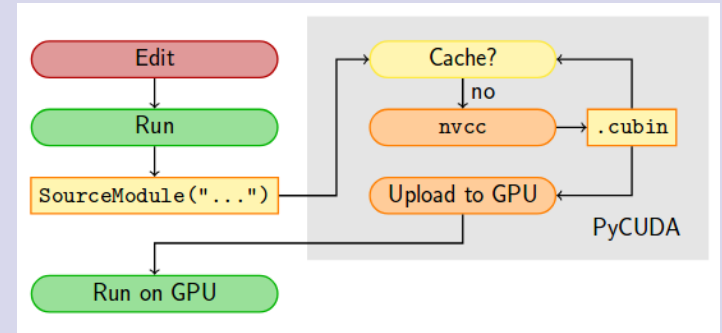
```
root@hell:~# pip3 install pycuda
Collecting pycuda
  Downloading pycuda-2017.1.1.tar.gz (1.6MB)
    100% |████████████████████████████████████████| 1.6MB 963kB/s
```

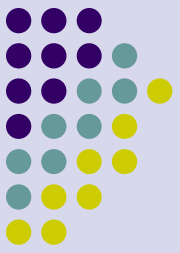
```
root@hell:~# pip3 show pycuda
Name: pycuda
Version: 2017.1.1
Summary: Python wrapper for Nvidia CUDA
Home-page: http://mathematician.de/software/pycuda
Author: Andreas Kloeckner
Author-email: inform@tiker.net
```



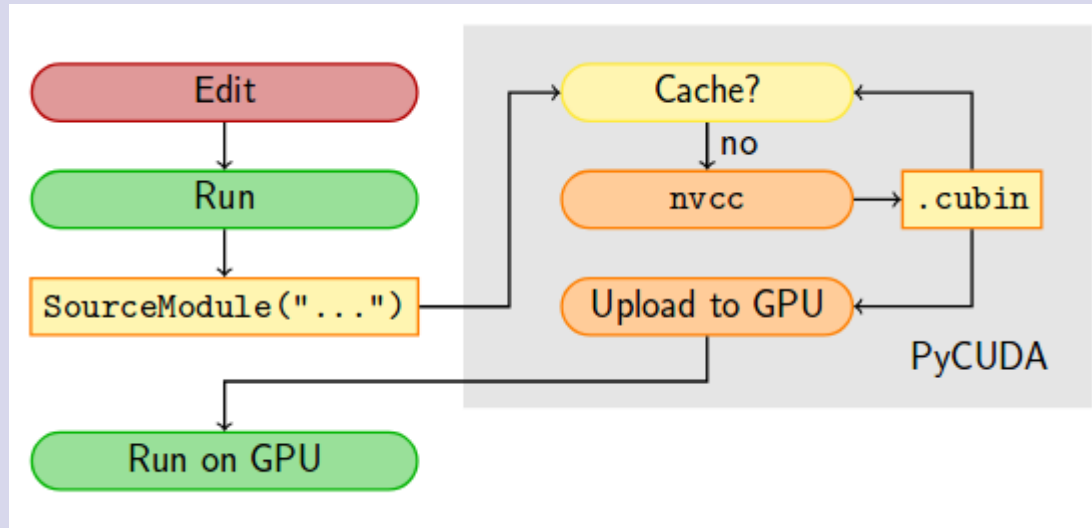
Accelerating python – PyCUDA

```
Editor - /Downloads/GTC2018/GTC_sample.py
GTC_sample.py* x
1#!/usr/bin/env python
2
3# importing modules
4import pycuda.driver as cuda
5import pycuda.autoinit
6from pycuda.compiler import SourceModule
7import numpy
8
9# creating a numpy array
10a = numpy.random.randn(4,4)
11
12# maing this array as single precision format
13a = a.astype(numpy.float32)
14
15# allocating GPU memory to store the a array
16a_gpu = cuda.mem_alloc(a.nbytes)
17
18# copy the array to GPU memory
19cuda.memcpy_htod(a_gpu, a)
20
21# declaring the kernel to be run on the GPU
22# using PyCUDA as C wrapper
23mod = SourceModule("""
24__global__ void doublify(float *a)
25{
26    int idx = threadIdx.x + threadIdx.y*4;
27    a[idx] *= 2;
28}
29""")
30
31# importing to the python realm the 'doublify' function
32# from the kernel and passing the a_gpu as argument
33func = mod.get_function("doublify")
34func(a_gpu, block=(4,4,1))
35
36# fetching the data back from the GPU and displaying it
37# also showing the original a
38a_doubled = numpy.empty_like(a)
39cuda.memcpy_dtoh(a_doubled, a_gpu)
40print(a_doubled)
41print(a)
```

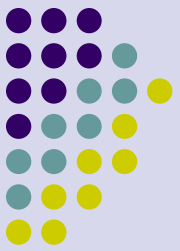




Accelerating python – PyCUDA

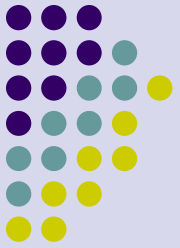


Accelerating python – Numba



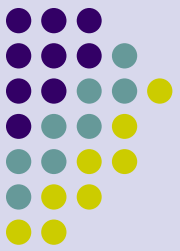
- high performance functions written in Python
- On-the-fly code generation
- Native code generation for the CPU and GPU
- Integration with the Python scientific stack
- Take advantage of Python decorators
- No need to write C code
- Code translation done using LLVM compiler

Accelerating python – Numba

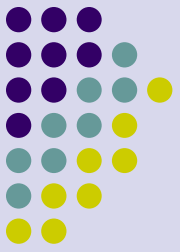


```
Spyder (Python 3.5)
Editor - /Downloads/GTC2018/GTC_sample.py
GTC_sample.py
1 #!/usr/bin/env python
2
3 import numpy as np
4 from numba import guvectorize
5
6 @guvectorize(['void(float64[:,], intp[:,], float64[:,])'], '(n),(n)->(n)')
7 def move_mean(a, window_arr, out):
8     window_width = window_arr[0]
9     asum = 0.0
10    count = 0
11    for i in range(window_width):
12        asum += a[i]
13        count += 1
14        out[i] = asum / count
15    for i in range(window_width, len(a)):
16        asum += a[i] - a[i - window_width]
17        out[i] = asum / count
18
19 arr = np.arange(20, dtype=np.float64).reshape(2, 10)
20 print(arr)
21 print(move_mean(arr, 3))
22
```

Accelerating python – cudamat

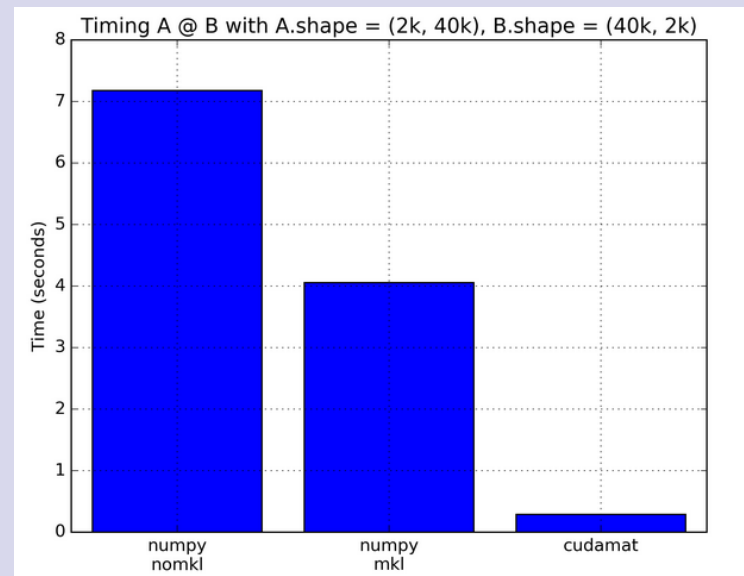


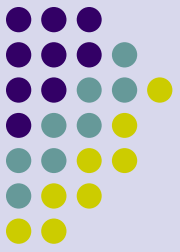
- provides a CUDA-based python matrix class
- Primary goal: easy dense matrix manipulation
- Useful to perform matrix ops on GPU
- Perform many matrix operations
 - multiplication and transpose
 - Elementwise addition, subtraction, multiplication, and division
 - Elementwise application of exp, log, pow, sqrt
 - Summation, maximum and minimum along rows or columns



Accelerating python – cudamat

```
Editor - /Downloads/GTC2018/GTC_sample.py
GTC_sample.py x
1 #!/usr/bin/env python
2
3 import cudamat as cm
4 import numpy as np
5
6 cm.cuda_set_device (0)
7 cm.init()
8
9 # Create some matrices on the GPU.
10 A = cm.CUDAMatrix(np.random.randn(10, 20))
11 B = cm.CUDAMatrix(np.random.randn(10, 20))
12 T = cm.CUDAMatrix(np.random.randn(10, 20))
13
14 # Add A and B and store the result in Tx
15 A.add(B, target = T)
16
17 # Multiply A by 50
18 A.mult(50)
19
```





Accelerating python – cupy

- an implementation of NumPy-compatible multi-dimensional array on CUDA
- Useful to perform matrix ops on GPU
- CuPy is faster than NumPy in many ways

Table 1: Matrix operation performance

Size	NumPy [ms]	CuPy [ms]
10^4	0.03	0.58
10^5	0.20	0.97
10^6	2.00	1.84
10^7	55.55	12.48
10^8	517.17	84.73

```
1 >>> import time, numpy, cupy
2 >>>
3 >>> size = 10 ** 8
4 >>> def test(xp):
5 ...     return xp.arange(size).reshape(1000, -1).T * 2
6 >>>
7 >>> for xp in [numpy, cupy]:
8 ...     test(xp) # Avoid first call overhead
9 ...     # Synchronize CPU and GPU for benchmark
10 ...     cupy.cuda.runtime.deviceSynchronize()
11 ...     t1 = time.time()
12 ...     test(xp)
13 ...     cupy.cuda.runtime.deviceSynchronize()
14 ...     t2 = time.time()
15 ...     print(xp.__name__, t2 - t1)
16 ('numpy', 0.5105748176574707)
17 ('cupy', 0.08418107032775879)
```