# ICNPG 2023

Python y un poquito de Julia

# GPUs y python

- Numba
- Cupy
- RAPIDS
- PyCUDA
- Etc...

- **Interpreted languages:**

  - **Python, Ruby, or JavaScript, etc.**

  - **Execute code directly from the source code without the need for a compilation step.**

  - **often have a simpler syntax and are easier to learn, read, and write.**

  - **better support for dynamic typing and dynamic memory management.**

  - **slower than compiled languages because they must interpret the code at runtime.**

- Compiled languages:

  - C, C++,  Fortran, Julia, etc.

  - require a compilation step to translate the source code into machine code.

  - tend to have a more complex syntax and require more experience to learn, read, and write.

  - They often have better support for static typing and memory management.

# Numba

https://numba.pydata.org/

## Accelerate Python Functions

Numba translates Python functions to optimized machine code at runtime using the industry-standard LLVM compiler library. Numba-compiled numerical algorithms in Python can approach the speeds of C or FORTRAN.

You don't need to replace the Python interpreter, run a separate compilation step, or even have a C/C++ compiler installed. Just apply one of the Numba decorators to your Python function, and Numba does the rest.

```python
from numba import njit
import random

@njit
def monte_carlo_pi(nsamples):
    acc = 0
    for i in range(nsamples):
        x = random.random()
        y = random.random()
        if (x ** 2 + y ** 2) < 1.0:
            acc += 1
    return 4.0 * acc / nsamples
```

# Numba

https://numba.pydata.org/

## Portable Compilation

Ship high performance Python applications without the headache of binary compilation and packaging. Your source code remains pure Python while Numba handles the compilation at runtime. We test Numba continuously in more than 200 different platform configurations.

Numba supports Intel and AMD x86, POWER8/9, and ARM CPUs (including Apple M1), NVIDIA GPUs, Python 3.7-3.10, as well as Windows/macOS/Linux. Precompiled Numba binaries for most systems are available as conda packages and pip-installable wheels.

## GPU Acceleration

NVIDIA CUDA

https://numba.readthedocs.io/en/stable/cuda/index.html

With support for NVIDIA CUDA, Numba lets you write parallel GPU algorithms entirely from Python.

# Numba

https://numba.pydata.org/

https://numba.readthedocs.io/en/stable/cuda/index.html

## GPU Acceleration

**NVIDIA CUDA**

With support for NVIDIA CUDA, Numba lets you write parallel GPU algorithms entirely from Python.

```
1    import numpy as np
2    from numba import cuda
```

```
1    @cuda.jit
2    def f(a, b, c):
3        # like threadIdx.x + (blockIdx.x * blockDim.x)
4        tid = cuda.grid(1)
5        size = len(c)
6
7        if tid < size:
8            c[tid] = a[tid] + b[tid]
```

```
1    N = 100000
2    a = cuda.to_device(np.random.random(N))
3    b = cuda.to_device(np.random.random(N))
4    c = cuda.device_array_like(a)
```

NumPy/SciPy-compatible Array Library for GPU-accelerated Computing with Python

https://cupy.dev/

**HIGHLY COMPATIBLE WITH NUMPY & SCIPY**

```
>>> import cupy as cp
>>> x = cp.arange(6).reshape(2, 3).astype('f')
>>> x
array([[ 0.,  1.,  2.],
       [ 3.,  4.,  5.]], dtype=float32)
>>> x.sum(axis=1)
array([  3.,  12.], dtype=float32)
```

```
>>> x = cp.arange(6, dtype='f').reshape(2, 3)
>>> y = cp.arange(3, dtype='f')
>>> kernel = cp.ElementwiseKernel(
...     'float32 x, float32 y', 'float32 z',
...     '''
...     if (x - 2 > y) {
...         z = x * y;
...     } else {
...         z = x + y;
...     }
...     ''', 'my_kernel')
>>> kernel(x, y)
array([[ 0.,  2.,  4.],
       [ 0.,  4., 10.]], dtype=float32)
```

The N-dimensional array (`ndarray`)

Universal functions (`cupy.ufunc`)

Routines (NumPy)

Routines (SciPy)

CuPy-specific functions

Low-level CUDA support

Custom kernels

Distributed

Environment variables

Comparison Table

Python Array API Support

# RAPIDS
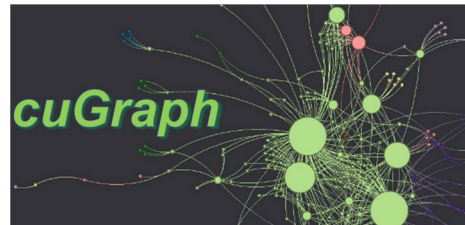## GPU Accelerated Data Science

### ⊞ FASTER PANDAS WITH CUDF

cuDF is a near drop in replacement to pandas for most use cases and has greatly improved performance.

### ⊞ FASTER SCIKIT-LEARN WITH CUML

cuML brings huge speedups to ML modeling with an API that matches scikit-learn.

### ⊘ FASTER NETWORKX WITH CUGRAPH

cuGraph makes migration from networkX easy, accelerates graph analytics, and allows scaling far beyond existing tools.



- ▸ Dataframe processing with **cuDF** (similar API to pandas)
- ▸ Machine learning with **cuML** (similar API to scikit-learn)
- ▸ Graph processing with **cuGraph** (similar API to networkX)
- ▸ Spatial analytics with **cuSpatial** (similar API to geoPandas)
- ▸ Image processing with **cuCIM** (similar API to scikit-image)
- ▸ Signal processing **cuSignal** (similar API to scipy.signal)
- ▸ Seamless cross-filtered dashboards with **cuxfilter**
- ▸ Low level compute primitives with **RAFT**
- ▸ Apache Spark acceleration with **Spark RAPIDS**

```python
import cudf, requests
from io import StringIO

url = "https://github.com/plotly/datasets/raw/master/tips.csv"
content = requests.get(url).content.decode('utf-8')

tips_df = cudf.read_csv(StringIO(content))
tips_df['tip_percentage'] = tips_df['tip'] / tips_df['total_bill'] * 100

# display average tip by dining party size
print(tips_df.groupby('size').tip_percentage.mean())
```

# Ejercicio

```python
#@title numpy
import numpy as np
import time

n = 10000000

# create x and y vectors
x = np.arange(1, n+1, dtype=np.float32)
y = np.arange(n, 0, -1, dtype=np.float32)

start_time = time.time()

# create z vector
z = x + y

# normalize z by its Euclidean norm
norm = np.linalg.norm(z)
#z_norm = z / norm

# sum the elements of the final result
result = np.sum(z)/norm

# print the result and the time taken
print("Result:", result)
print("Time taken:", (time.time() - start_time) * 1000000, "us")

tnumpy = (time.time() - start_time) * 1000000
#print(x,y,z,z_norm)
```

$$S = \frac{\sum_{i=0}^{n}(x_i + y_i)}{\sqrt{\sum_{i=0}^{n}(x_i + y_i)^2}}$$

# Ejercicio: hacerlo en CUPY y NUMBA

**Tasa de contagio**

$$\frac{dS}{dt} = -\beta SI,$$

$$\frac{dI}{dt} = \beta SI - \gamma I,$$

$$\frac{dR}{dt} = \gamma I.$$

*Método de Euler*

$$S(t + \Delta t) = S(t) + \Delta t[-\beta I(t)S(t)]$$

$$I(t + \Delta t) = I(t) + \Delta t[\beta I(t)S(t) - \gamma I(t)]$$

$$R(t + \Delta t) = R(t) + \Delta t[\gamma I(t)]$$

**Suceptible, Infectado, Recuperado**

Julia is a good choice for numerical computing and scientific computing tasks that require **high performance and parallelism**, while Python is a good choice for general-purpose programming, web development, and data analysis tasks that require a large ecosystem of libraries and tools.

## Julia in a Nutshell

### Fast

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

### Dynamic

Julia is dynamically typed, feels like a scripting language, and has good support for interactive use.

### Reproducible

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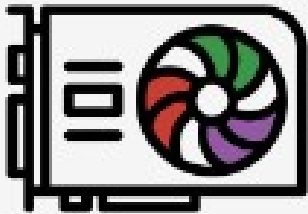
### Composable

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

### General

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

### Open source

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

# CUDA programming in Julia

The CUDA.jl package is the main entrypoint for programming NVIDIA GPUs in Julia. The package makes it possible to do so at various abstraction levels, from easy-to-use arrays down to hand-written kernels using low-level CUDA APIs.

CUDA.jl