



Sören Metje

GPU Computing with Python

Table of contents

- 1** Introduction
- 2** GPU Computing
- 3** Python Frameworks
- 4** Summary
- 5** Appendix

Motivation

Why Python?

High-level programming

Compatible with many platforms and systems

Many high quality frameworks and libraries

Widely distributed in many different domains

Big community



Motivation

Why GPU Computing?

Reduce wall-clock time

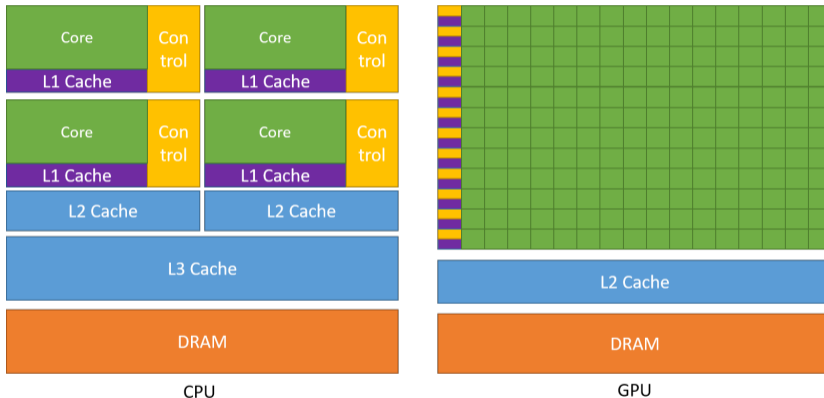
Achieve higher cost-efficiency

Source: NVIDIA - GPU-Accelerated Google Cloud [\[NVID\]](#)

Outline

- 1 Introduction
- 2 GPU Computing**
- 3 Python Frameworks
- 4 Summary
- 5 Appendix

GPU Architecture



Source: CUDA Toolkit Documentation [NVIa]

Use Cases

When to use GPU Computing?

Large data set available

Parallel processing possible

Use cases: Fluid dynamics, Image processing,
Deep learning, ...

When not to use GPU Computing?

Data set is too small

Data set is too big (exceeds GPU memory size)

Large amount of small sequential operations

Source: NVIDIA - GPU-Accelerated Google Cloud [\[NVID\]](#)

CUDA

Definition

NVIDIA CUDA (Compute Unified Device Architecture) is a parallel computing platform and programming model for general computing on GPUs.



Initial release: June 23, 2007

Gives access to the GPU's virtual instruction set

Enables execution of compute kernels

Accessible through frameworks, libraries, and compiler directives

Closed source

CUDA Compute Kernel

Definition

A compute kernel is a function compiled for accelerators (such as GPUs).

C++

```

    __global__ void add(float *A, float *B, float *C) {
        int i = threadIdx.x;
        C[i] = A[i] + B[i];
    }

```

Compute Kernel Limitations

Allowed operations: basic math operations, if / else, for / while loops

Can not explicitly return a value

Write results to passed array

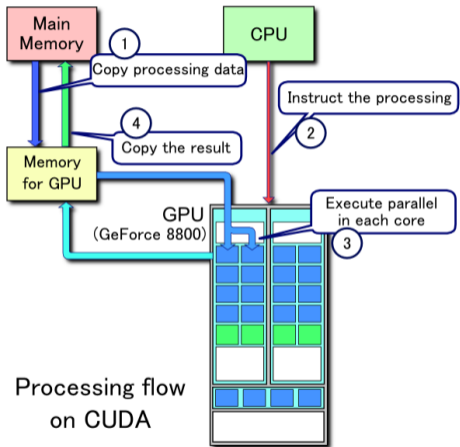
C++

```

    float * A, * B, * C;
    for (int i = 0; i < N; i++)
        C[i] = A[i] + B[i];
}

```

CUDA Processing Flow



Processing flow
on CUDA

Source: Wikipedia - CUDA [Wika]

AMD ROCm

Definition

AMD ROCm (Radeon Open Compute) is a software stack for GPU programming.

"NVIDIA CUDA for AMD GPUs"

Initial release: November 14, 2016

Available on [GitHub](#) (open-source)

Supported by:

- I PyTorch
- I TensorFlow
- I CuPy

Not as widely supported as NVIDIA CUDA

Gaining traction in the TOP500



TOP 500 List

Rank	System	Cores	Rmax (PFlop/s)	Rpeak (PFlop/s)	Power (kW)
1	Frontier - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE DOE/SC/Oak Ridge National Laboratory United States	8,730,112	1,102.00	1,685.65	21,100
2	Supercomputer Fugaku - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science Japan	7,630,848	442.01	537.21	29,899
3	LUMI - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE EuroHPC/CSC Finland	1,110,144	151.90	214.35	2,942
4	Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM DOE/SC/Oak Ridge National Laboratory United States	2,414,592	148.60	200.79	10,096

Source: TOP 500 List - June 2022 [\[TOP\]](#)

Outline

- 1 Introduction
- 2 GPU Computing
- 3 Python Frameworks**
- 4 Summary
- 5 Appendix

Python Frameworks for GPU Computing

Computing Frameworks

Numba

CuPy

Scikit-cuda

RAPIDS

Triton (presented by Dimitris Oikonomou on 2022-05-19)

...

Deep Learning Frameworks

PyTorch

TensorFlow

Keras

...

Python Frameworks for GPU Computing

Computing Frameworks

Numba

CuPy

Scikit-cuda

RAPIDS

Triton (presented by Dimitris Oikonomou on 2022-05-19)

...

Deep Learning Frameworks

PyTorch

TensorFlow

Keras

...

Numba

Definition

Numba is a just-in-time compiler for numerical functions in Python.

Translates Python functions to optimized machine code at runtime

Compiles code for CPU and GPU

Supports

- I NVIDIA GPUs: CUDA
- I ~~AMD GPUs: ROCm (deprecated)~~

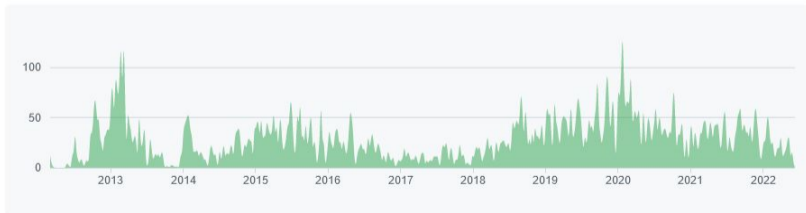
Available on [GitHub](#) (open-source)

Numba - Development

Mar 4, 2012 – Jun 13, 2022

Contributions: Commits ▾

Contributions to main, excluding merge commits and bot accounts



Source: GitHub - Numba [[Gita](#)]

Numba - Programming Approaches

2 Approaches for GPU Programming

Universal functions

CUDA Kernels

Numba - Universal Functions

Definition

A universal function (or ufunc for short) is a function that operates on arrays in an element-by-element fashion.

Python

```

    çÎ      òè      ò
    òè      ò      Î      Î      Î
    îò  Îîî      è
    ò
    Î      Î Î ò  Î      ò      Î
    ç      Î
    Îîî      è Î ç
  
```

Numba - Universal Functions on GPU

Python

```

    @jit
    def add(a, b):
        return a + b
  
```

- 1 Compile CUDA kernel
- 2 Allocate GPU memory
- 3 Copy data to the GPU
- 4 Executed CUDA kernel
- 5 Copy result back to the CPU
- 6 Return the result

Numba - CUDA Kernels

Python

```

    cI      e iI
      I

e iI
iò Iii  ò ò
    I      iò e iI  i      e iI  i  ò      ò i      òIi  i
      I  ò  I      I  ò      iò

    I I  ò      I  ò      I
    I I  ò      I  ò      I
      ò      ò

Iii  ò  ò      ç  è  ò  i  òIi  ò  ç  è
      ò
  
```

CuPy

Definition

CuPy is a NumPy/SciPy-compatible array library for GPU-accelerated computing.

"NumPy for GPU computing"

Provides various math operations

Supports

- I NVIDIA GPUs: CUDA
- I AMD GPUs: ROCm (experimental)

Available on [GitHub](#) (open-source)

CuPy - Development

Apr 12, 2015 – Jun 13, 2022

Contributions: Commits ▾

Contributions to master, excluding merge commits and bot accounts



Source: GitHub - CuPy [[Gitb](#)]

CuPy - Math Functions

Python

\hat{I}

$\hat{I} \hat{I}^T$

$\hat{I} \hat{I} \hat{I}$

\hat{I}

\hat{I}

$\hat{I} \hat{I}$

$\hat{I} \hat{I}$

$\hat{I} \hat{I}$

$j \otimes \mathcal{E}$

\hat{I}^T

$\hat{I}^T \otimes \mathcal{E}$

CuPy - CUDA Kernels

Python

```

    è  Î  è

Îîî  è  “Î  aò  ò
ò  ò  ,  çÎ  î  Îîî  è  è
      î  ç  è  2  ç  è  Sí  òÎîSí
      î  î  î

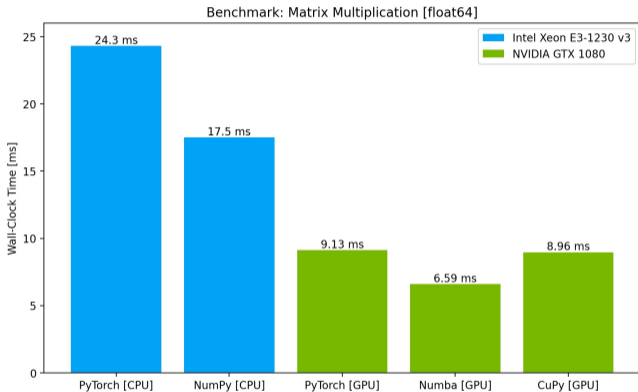
Îîî
è  Î  Î  ò  î  ò
è  Î  Î  ò  î  ò
è  ò  î  ò

Îîî  ç  è  ò  î  òÎî  ò  ç  è
      Î  Î
  
```

Benchmark - Matrix Multiplication

Code for benchmarking is available on [GWDG GitLab](#)

Benchmark - Matrix Multiplication Float64



Wall-clock time in milliseconds for computing result of fast matrix multiplication. Mean of 100 loops and 4 runs.
Multiplied a 1024x2048 matrix and a 2048x512 matrix of float64.
Measured time does not count in time to copy matrices to GPU and result matrix back from GPU.

Code for benchmarking is available on [GWDG GitLab](#)

Benchmark - Matrix Multiplication

Code for benchmarking is available on [GWDG GitLab](#)

Benchmark - Matrix Multiplication

Code for benchmarking is available on [GWDG GitLab](#)

Outline

- 1 Introduction
- 2 GPU Computing
- 3 Python Frameworks
- 4 Summary**
- 5 Appendix

New Trends

What's new?

GPUs get more powerful

More complex models and computations

New high-level Python frameworks provide features for various use cases

ROCm (open-source) is gaining traction

Summary

GPUs

- I achieve massive data parallelism
- I reduce wall-clock time
- I increase cost efficiency

Summary

GPUs

- I achieve massive data parallelism
- I reduce wall-clock time
- I increase cost efficiency

Numba Advantages

- I Enables implementing own universal functions & kernels in Python running on GPU

Summary

GPUs

- I achieve massive data parallelism
- I reduce wall-clock time
- I increase cost efficiency

Numba Advantages

- I Enables implementing own universal functions & kernels in Python running on GPU

CuPy Advantages

- I Easily move existing NumPy code towards GPU computing
- I Directly use NumPy-style array operations and execute them on GPU
- I Great starting-point to learn about GPU computing

References I

[Avimanyu Bandyopadhyay](#). *Hands-On GPU Computing with Python*. Packt Publishing Ltd, 2019. ISBN: 9781789341072.

[NVIDIA](#). *Programming Guide :: CUDA Toolkit Documentation*. URL:

[https://docs.nvidia.com/cuda/cuda-toolkit/docs/index.html](#).

[NVIDIA](#). *PTX ISA :: CUDA Toolkit Documentation*. URL:

[https://docs.nvidia.com/cuda/cuda-toolkit/docs/ptx/index.html](#).

[NVIDIA](#). *CUDA Zone - Library of Resources | NVIDIA Developer*. URL:

[https://developer.nvidia.com/cuda-zone](#).

[NVIDIA](#). *GPU-Accelerated Google Cloud*. URL:

[https://cloud.google.com/accelerator/gpu-overview](#).

[Wikipedia](#). *CUDA*. URL:

[https://en.wikipedia.org/wiki/CUDA](#).

[Wikipedia](#). *ROCm*. URL:

[https://en.wikipedia.org/wiki/ROCm](#).

[Wikipedia](#). *Compute kernel*. URL:

[https://en.wikipedia.org/wiki/Compute_kernel](#).

[Microsoft](#). *VISC: Virtual Instruction Set Computing*. URL:

[https://www.microsoft.com/en-us/research/publication/virtual-instruction-set-computing/](#).

References II

[GitHub. Numba. URL:](#) [https://github.com/numba/numba](#).

[Anaconda. Numba: A High Performance Python Compiler. URL:](#) [https://docs.continuum.io/numba/](#).

[Anaconda. Numba documentation. URL:](#) [https://numba.pydata.org/numba-doc/](#).

[Villoro. Villoro - numba. URL:](#) [https://villoro.github.io/numba/](#).

[GitHub. CuPy. URL:](#) [https://github.com/cupy/cupy](#).

[Preferred Networks, Inc. and Preferred Infrastructure, Inc. CuPy NumPy & SciPy for GPU. URL:](#) [https://preferred.ai/cupy/](#).

[NumPy. NumPy documentation. URL:](#) [https://numpy.org/doc/](#).

[Mindfire Solutions. Python: 7 Important Reasons Why You Should Use Python. URL:](#) [https://www.mindfiresolutions.com/python-7-important-reasons-why-you-should-use-python/](#).

[Marko Alešić. CPU Vs. GPU: A Comprehensive Overview. URL:](#) [https://www.markoalešić.com/cpu-vs-gpu-a-comprehensive-overview/](#).

References III

[TOP500](#). *TOP 500 List*. URL: [https://www.top500.org/](#).

[Amazon Web Services](#). *Amazon EC2 P3 instances*. URL: [https://aws.amazon.com/ec2/instance-types/p3/](#).

[Google](#). *Google Colab*. URL: [https://colab.research.google.com/](#).

GPU as a Service

Instance Size	GPUs - Tesla V100	GPU Peer to Peer	GPU Memory (GB)	vCPUs	Memory (GB)	Network Bandwidth	EBS Bandwidth	On-Demand Price/hr*	1-yr Reserved Instance Effective Hourly*	3-yr Reserved Instance Effective Hourly*
p3.2xlarge	1	N/A	16	8	61	Up to 10 Gbps	1.5 Gbps	\$3.06	\$1.99	\$1.05
p3.8xlarge	4	NVLink	64	32	244	10 Gbps	7 Gbps	\$12.24	\$7.96	\$4.19
p3.16xlarge	8	NVLink	128	64	488	25 Gbps	14 Gbps	\$24.48	\$15.91	\$8.39
p3dn.24xlarge	8	NVLink	256	96	768	100 Gbps	19 Gbps	\$31.218	\$18.30	\$9.64

Example Pricing: Amazon EC2 P3 instances. Source: AWS [[Ama](#)]

Advantages

- Scalable resources
- On-demand pricing

Disadvantages

- Expense for large periods of time
- Data confidentiality not guaranteed depending on vendor

GPU as a Service

Willkommen bei Colaboratory

Datei Bearbeiten Anzeige Einfügen Laufzeit Tools Hilfe Änderungen können nicht gespeichert werden

Inhalt

- Erste Schritte
- Data Science
- Maschinelles Lernen
- Weitere Ressourcen
 - Vorgestellte Beispiele
- Bereich

+ Code + Text In Google Drive kopieren Verbinden Bearbeiten

Data Science

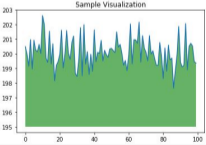
Mit Colab können Sie das volle Potenzial der beliebten Python-Bibliotheken nutzen, um Daten zu analysieren und visualisieren. Die Codezelle unten verwendet **NumPy**, um zufällige Daten zu erstellen, und **Matplotlib**, um sie zu visualisieren. Klicken Sie zum Bearbeiten des Codes einfach auf die Zelle und beginnen Sie mit der Bearbeitung.

```
[ ] import numpy as np
from matplotlib import pyplot as plt

ys = 200 + np.random.randn(100)
x = [x for x in range(len(ys))]

plt.plot(x, ys, '-')
plt.fill_between(x, ys, 195, where=(ys > 195), facecolor='g', alpha=0.6)

plt.title("Sample Visualization")
plt.show()
```



Google Colab Jupyter Notebooks. Source: Google [Goo]

NVIDIA PTX

Definition

NVIDIA PTX (Parallel Thread Execution) is a low-level parallel thread execution virtual machine and instruction set architecture (ISA)

- Exposes the GPU as a data-parallel computing device

- Interprets compiled code (analog to Java byte code interpreted by JVM)

Advantage

- Achieves portability of source code among multiple NVIDIA GPUs