# INTRODUÇÃO À PROGRAMAÇÃO PARA GPUs USANDO CUDA

Pedro Bruel

phrb@ime.usp.br

29 de Setembro de 2015

Instituto de Matemática e Estatística
Universidade de São Paulo

Os *pdf*s com as aulas e todo o código fonte usado nos exemplos estão no GitHub:

- github.com/phrb/aulas-gpu

# RECURSOS

Os *pdf*s com as aulas e todo o código fonte usado nos exemplos estão no GitHub:

- github.com/phrb/aulas-gpu

Outros recursos:

- CUDA Toolkit Documentation: docs.nvidia.com/cuda
- GPU Teaching Kit: syllabus.gputeachingkit.com
- iPython: ipython.org/notebook.html
- CUDA Toolkit: developer.nvidia.com/cuda-toolkit
- Anaconda: continuum.io/downloads

Os próximos *slides* foram adaptados do material disponível no GPU Teaching Kit:

- syllabus.gputeachingkit.com

GPU Teaching Kit

GPU Teaching Kit

Accelerated Computing

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

Lecture 2.4 – Introduction to CUDA C

Introduction to the CUDA Toolkit

# Objective

– To become familiar with some valuable tools and resources from the CUDA Toolkit
  – Compiler flags
  – Debuggers
  – Profilers

# GPU Programming Languages

| | |
|---|---|
| **Numerical analytics** ▷ | MATLAB, Mathematica, LabVIEW |
| **Fortran** ▷ | CUDA Fortran |
| **C** ▷ | CUDA C |
| **C++** ▷ | CUDA C++ |
| **Python** ▷ | PyCUDA, Copperhead, Numba, NumbaPro |
| **F#** ▷ | Alea.cuBase |

# NVCC Compiler

- NVIDIA provides a CUDA-C compiler
  - nvcc
- NVCC compiles device code then forwards code on to the host compiler (e.g. g++)
- Can be used to compile & link host only applications

# Compiler Flags

– Remember there are two compilers being used
  – NVCC: Device code
  – Host Compiler:  C/C++ code
– NVCC supports some host compiler flags
  – If flag is unsupported, use –Xcompiler to forward to host
    – e.g. –Xcompiler –fopenmp
– Debugging Flags
  – -g:  Include host debugging symbols
  – -G: Include device debugging symbols
  – -lineinfo:  Include line information with symbols

# Developer Tools - Debuggers



**NSIGHT**

**CUDA-GDB**

**CUDA MEMCHECK**

**NVIDIA Provided**

allinea DDT

TotalView®

**3rd Party**

https://developer.nvidia.com/debugging-solutions

# CUDA-MEMCHECK

– Memory debugging tool
  – No recompilation necessary
      %> cuda-memcheck ./exe
– Can detect the following errors
  – Memory leaks
  – Memory errors (OOB, misaligned access, illegal instruction, etc)
  – Race conditions
  – Illegal Barriers
  – Uninitialized Memory
– For line numbers use the following compiler flags:
  – -Xcompiler -rdynamic -lineinfo

http://docs.nvidia.com/cuda/cuda-memcheck

NVIDIA   ILLINOIS

# CUDA-GDB

– cuda-gdb is an extension of GDB
  – Provides seamless debugging of CUDA and CPU code
– Works on Linux and Macintosh
  – For a Windows debugger use NSIGHT Visual Studio Edition

http://docs.nvidia.com/cuda/cuda-gdb

# Developer Tools - Profilers



**NSIGHT**

**NVVP**

**NVPROF**

**NVIDIA Provided**

**TAU**

**VampirTrace**

**3rd Party**

https://developer.nvidia.com/performance-analysis-tools

# NVPROF

Command Line Profiler
– Compute time in each kernel
– Compute memory transfer time
– Collect metrics and events
– Support complex process hierarchy's
– Collect profiles for NVIDIA Visual Profiler
– No need to recompile

# NVIDIA's Visual Profiler (NVVP)
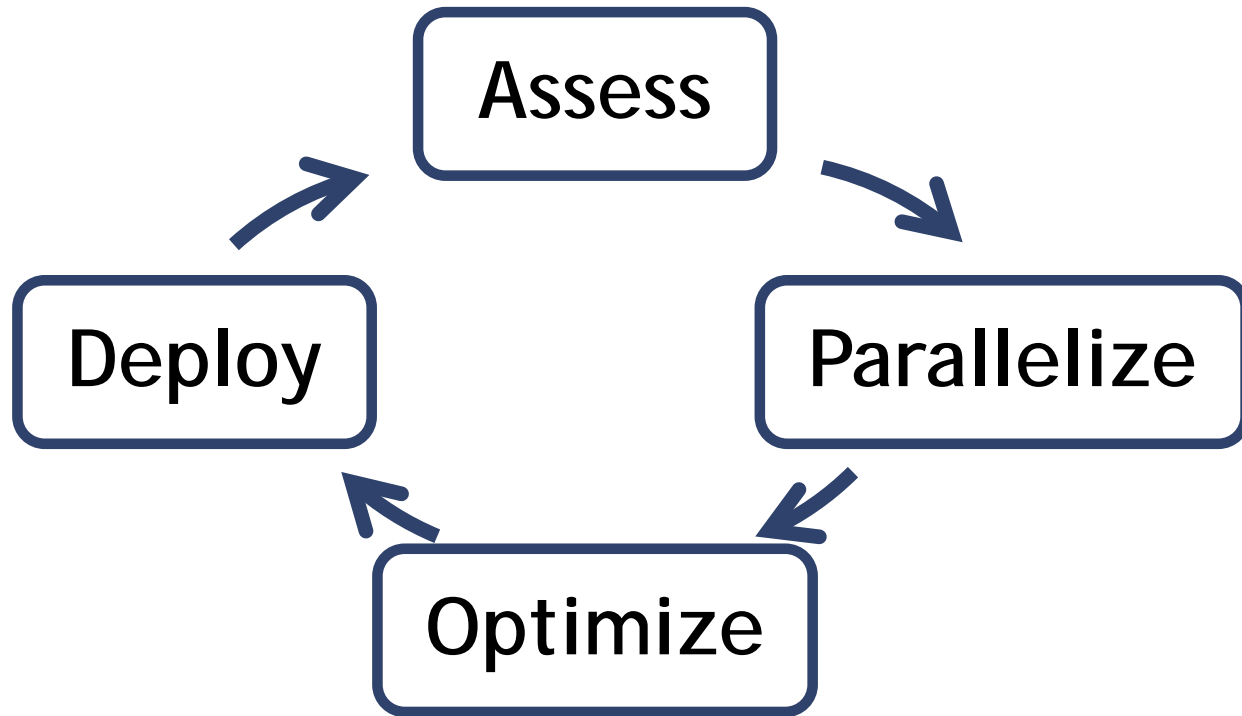
## Timeline



## Guided System
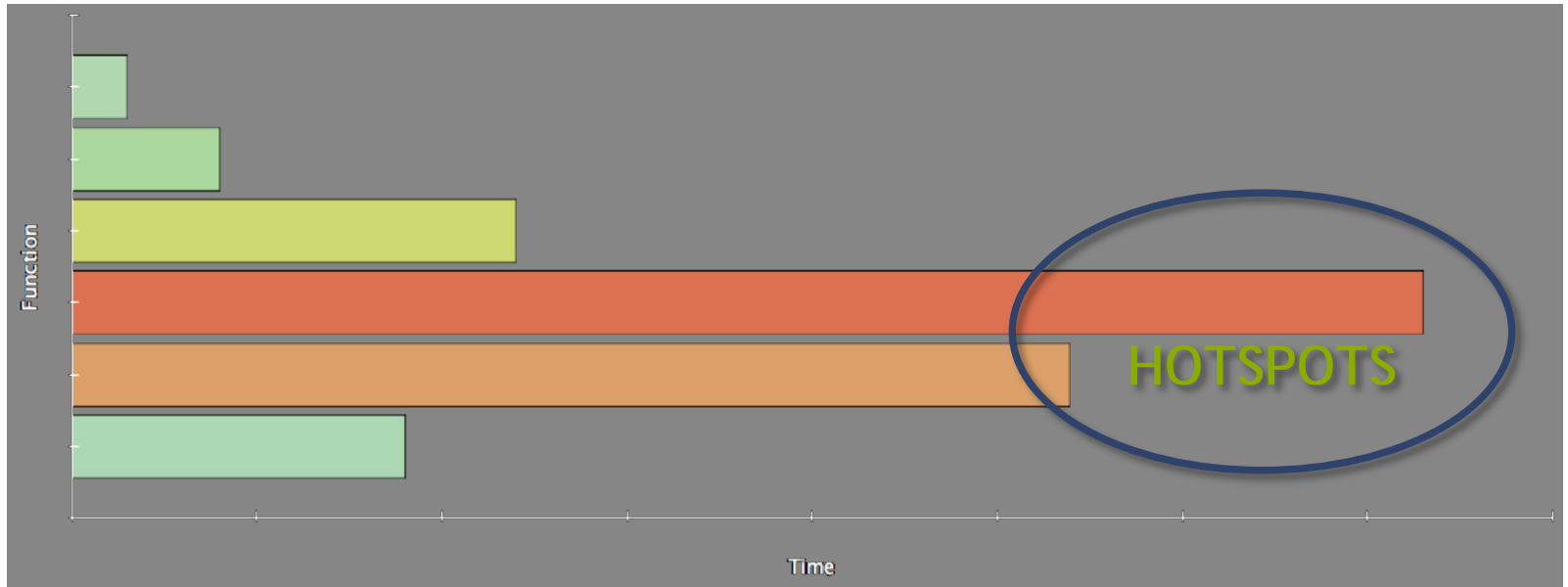


## Analysis

# Profiler Summary

- Many profile tools are available
- NVIDIA Provided
  - NVPROF:  Command Line
  - NVVP:  Visual profiler
  - NSIGHT: IDE (Visual Studio and Eclipse)
- 3rd Party
  - TAU
  - VAMPIR

# Optimization

# Assess



– Profile the code, find the hotspot(s)
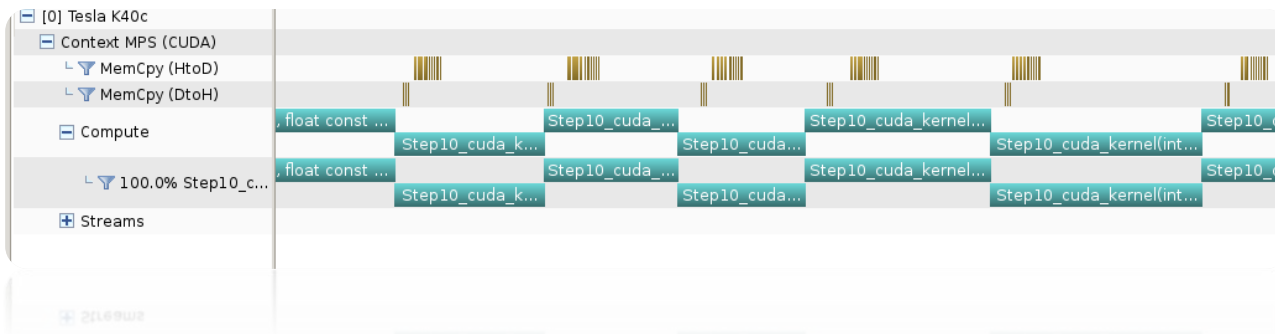– Focus your attention where it will give the most benefit

# Parallelize

**Applications**

| Libraries | Compiler Directives | Programming Languages |
|-----------|--------------------|-----------------------|

# Optimize

## Timeline



## Guided System



## Analysis

Obrigado!

# INTRODUÇÃO À PROGRAMAÇÃO PARA GPUs USANDO CUDA

Pedro Bruel

phrb@ime.usp.br

29 de Setembro de 2015

Instituto de Matemática e Estatística
Universidade de São Paulo