

HETEROGENEOUS CPU+GPU COMPUTING

Ana Lucia Varbanescu – University of Amsterdam
a.l.varbanescu@uva.nl

Stijn Heldens – Twente University
mail@stijnh.nl

Significant contributions by: **Pieter Hijma** (UvA, NL),
Jie Shen (TUDelft, NL), **Basilio Fraguera** (A Coruna University, ESP)

PART I

Heterogeneous processing: pro's and con's

Hardware Performance metrics

- Clock frequency [GHz] = absolute hardware speed
 - Memories, CPUs, interconnects
- Operational speed [GFLOPs]
 - Instructions per cycle + frequency
- Memory bandwidth [GB/s]
 - differs a lot between different memories on chip
- Power [Watt]
- Derived metrics
 - FLOP/Byte, FLOP/Watt

Theoretical peak performance

$$\text{Peak} = \text{chips} * \text{cores} * \text{threads/core} * \text{vector_lanes} * \\ \text{FLOPs/cycle} * \text{clockFrequency}$$

- Some examples:

- Intel Core i7 CPU

- $2 \text{ chips} * 4 \text{ cores} * 4\text{-way vectors} * 2 \text{ FLOPs/cycle} * 2.4 \text{ GHz} = \mathbf{154 \text{ GFLOPs}}$

- NVIDIA GTX 580 GPU

- $1 \text{ chip} * 16 \text{ SMs} * 32 \text{ cores} * 2 \text{ FLOPs/cycle} * 1.544 \text{ GHz} = \mathbf{1581 \text{ GFLOPs}}$

Performance ratio (CPU:GPU): 1:10 !!!

DRAM Memory bandwidth

Bandwidth = memory bus frequency * bits per cycle
* bus width

- Memory clock != CPU clock!
- In bits, divide by 8 for GB/s
- Some Examples:
 - Intel Core i7 DDR3: $1.333 * 2 * 64 =$ **21 GB/s**
 - NVIDIA GTX 580 GDDR5: $1.002 * 4 * 384 =$ **192 GB/s**

Performance ratio (CPU:GPU): 1:8 !!!

Power

- Chip manufactures specify Thermal Design Power (TDP)
- We can measure dissipated power
 - Whole system
 - Typically (much) lower than TDP
- Power efficiency
 - FLOPS / Watt
- Examples (with theoretical peak and TDP)
 - Intel Core i7: $154 / 160 = \mathbf{1.0 \text{ GFLOPs/W}}$
 - NVIDIA GTX 580: $1581 / 244 = \mathbf{6.3 \text{ GFLOPs/W}}$
 - ATI HD 6970: $2703 / 250 = \mathbf{10.8 \text{ GFLOPs/W}}$

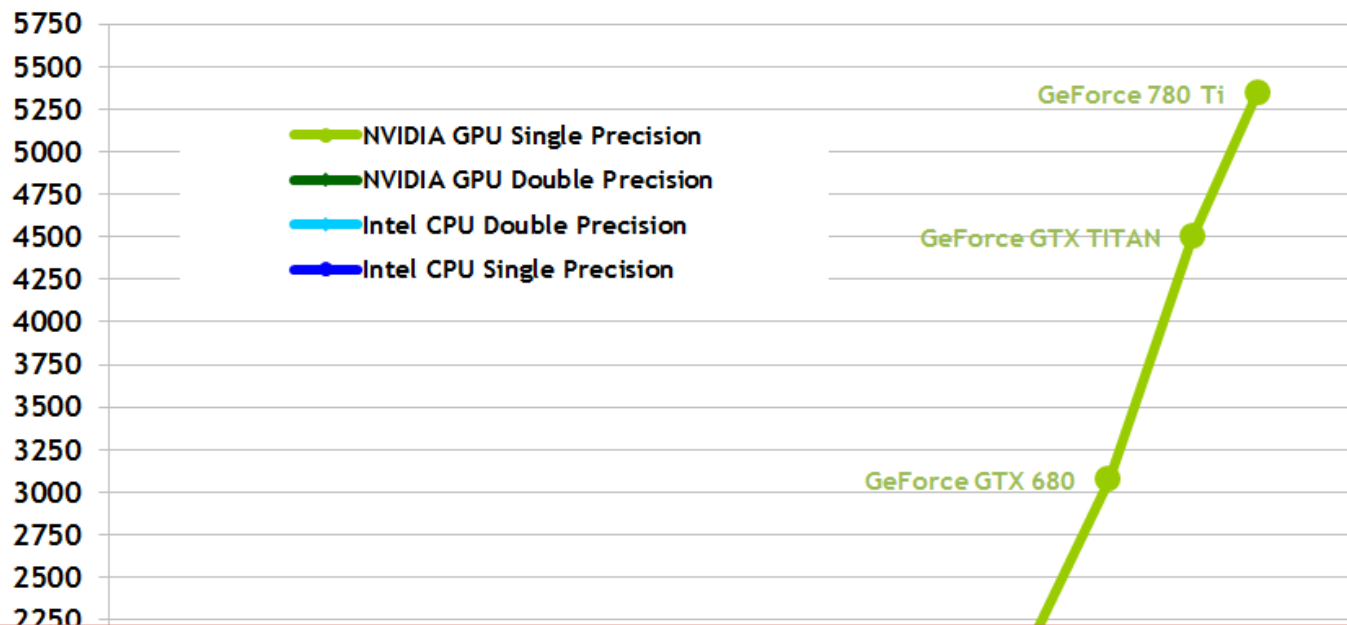
Summary

	Cores	Threads/ALUs	GFLOPS	Bandwidth
Sun Niagara 2	8	64	11.2	76
IBM BG/P	4	8	13.6	13.6
IBM Power 7	8	32	265	68
Intel Core i7	4	16	85	25.6
AMD Barcelona	4	8	37	21.4
AMD Istanbul	6	6	62.4	25.6
AMD Magny-Cours	12	12	125	25.6
Cell/B.E.	8	8	205	25.6
NVIDIA GTX 580	16	512	1581	192
NVIDIA GTX 680	8	1536	3090	192
AMD HD 6970	384	1536	2703	176
AMD HD 7970	32	2048	3789	264
Intel Xeon Phi 7120	61	240	2417	352

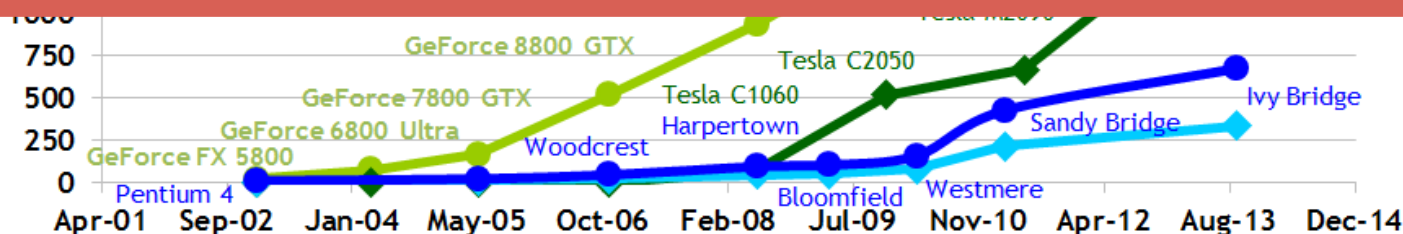
GPU vs. CPU performance

1 GFLOP = 10^9 ops

Theoretical GFLOP/s



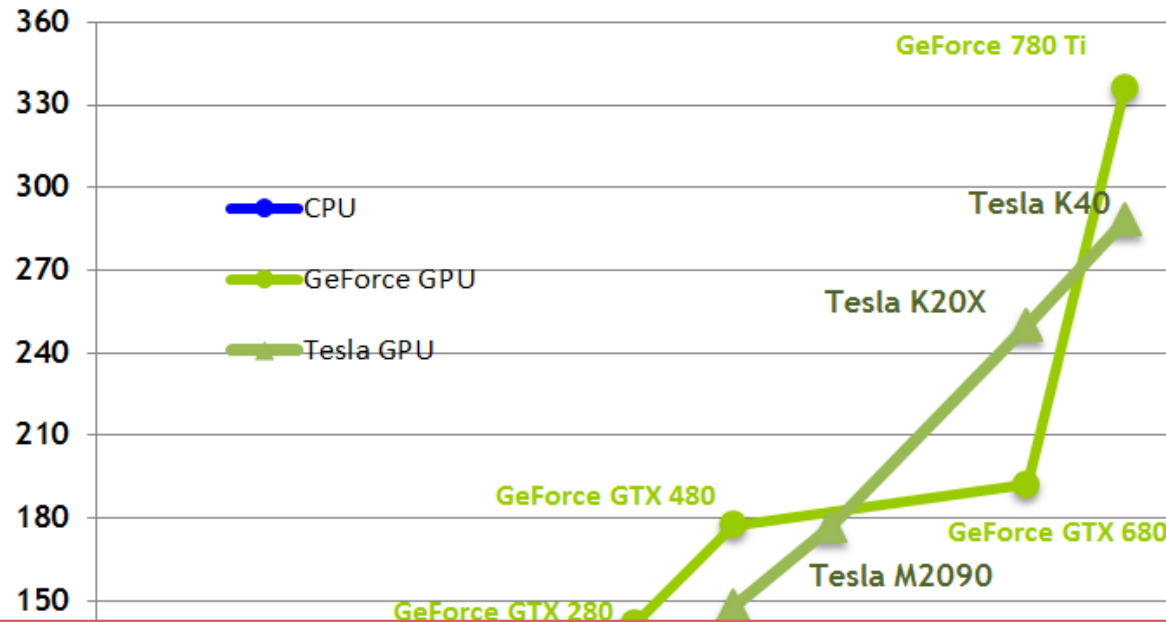
These are theoretical numbers! In practice, efficiency is much lower!



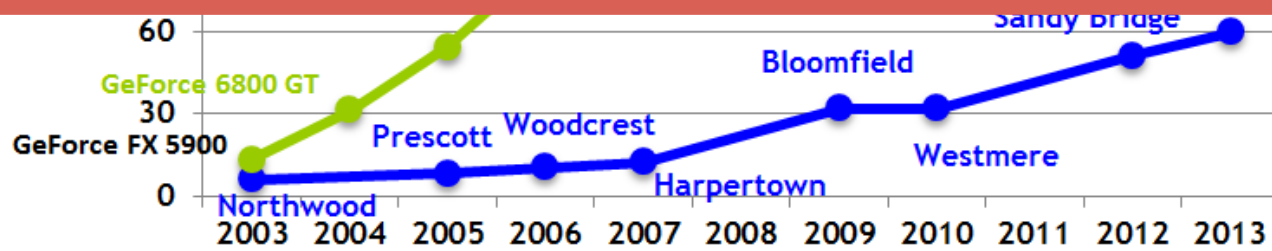
GPU vs. CPU performance

Theoretical GB/s

1 GB = 8×10^9 bits



These are theoretical numbers! In practice, efficiency is much lower!



Heterogeneity vs. Homogeneity

- Increase performance
 - Both devices work in parallel
 - Gain is much more than 10%
 - Decrease data communication
 - Which is often the bottleneck of the system
 - Different devices for different roles
- Increase flexibility and reliability
 - Choose one/all *PUs for execution
 - Fall-back solution when one *PU fails
- Increase power efficiency
- Cheaper per flop

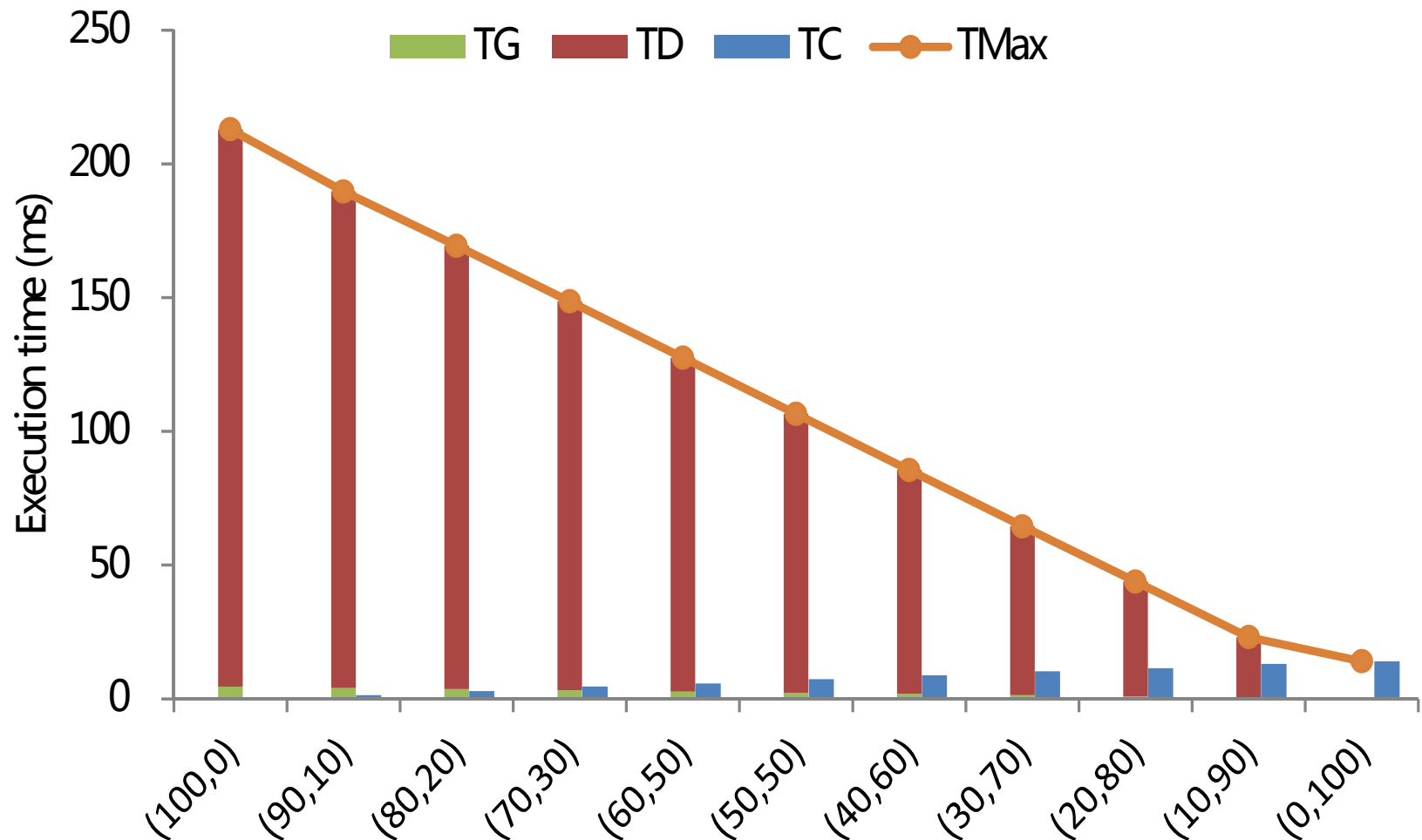
Example 1: dot product

- Dot product
 - Compute the dot product of 2 (1D) arrays
- Performance
 - T_G = execution time on GPU
 - T_C = execution time on CPU
 - T_D = data transfer time CPU-GPU
- GPU best or CPU best?

The diagram shows the dot product of two 1D arrays. The first array is $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$ and the second array is $\begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix}$. The result is $\begin{bmatrix} 58 \end{bmatrix}$. A yellow arrow labeled "Dot Product" points from the first array to the result.

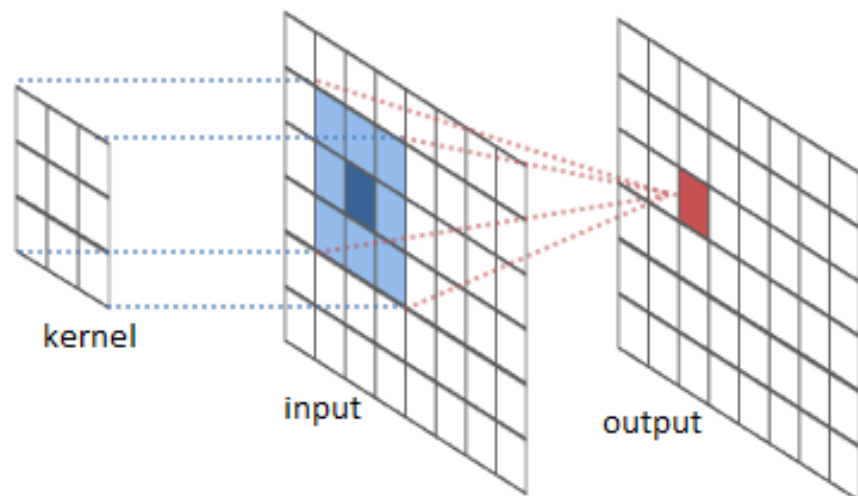
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix} = \begin{bmatrix} 58 \end{bmatrix}$$

Example 1: dot product

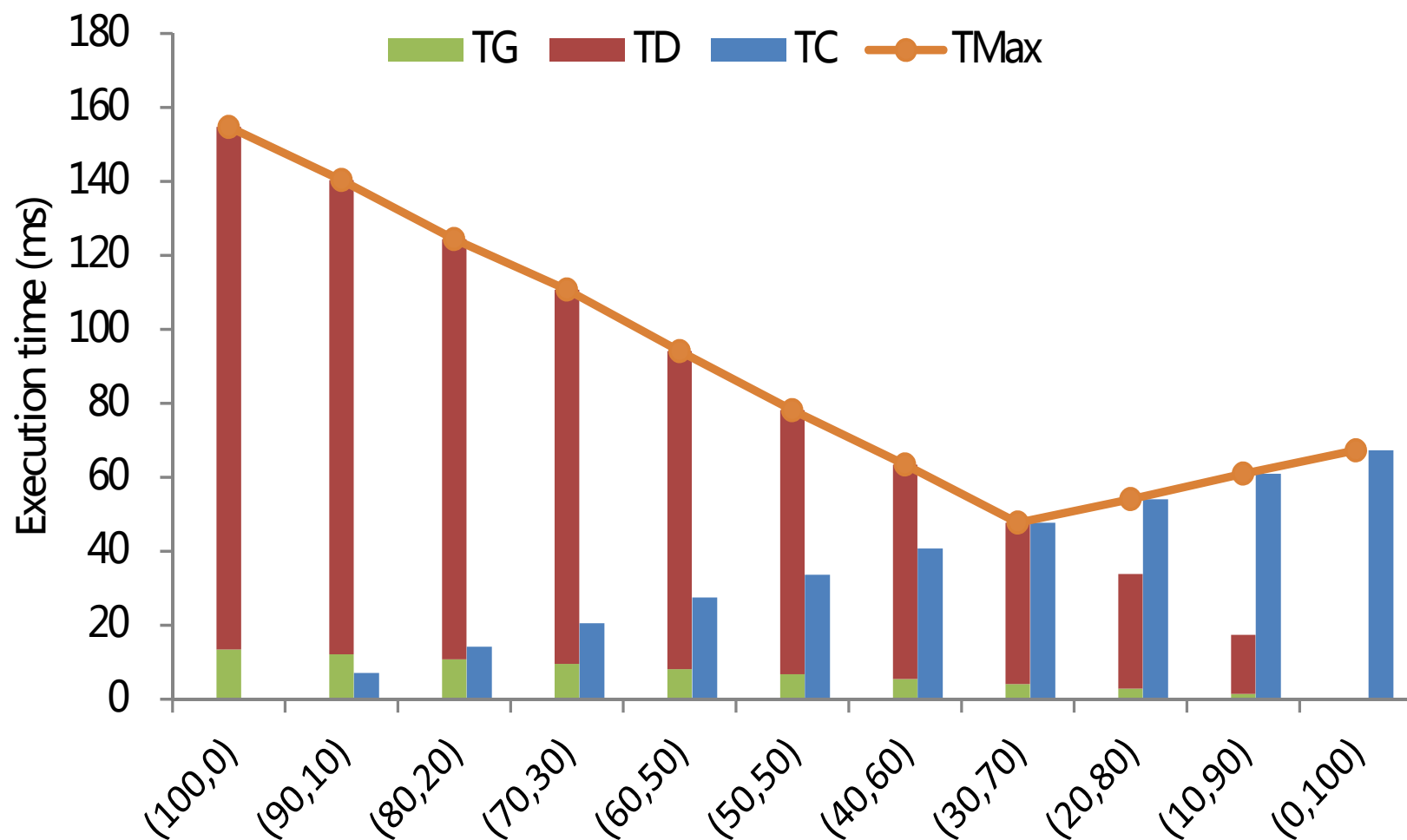


Example 2: separable convolution

- Separable convolution (CUDA SDK)
 - Apply a convolution filter (kernel) on a large image.
 - Separable kernel allows applying
 - Horizontal first
 - Vertical second
- Performance
 - T_G = execution time on GPU
 - T_C = execution time on CPU
 - T_D = data transfer time
- GPU best or CPU best?

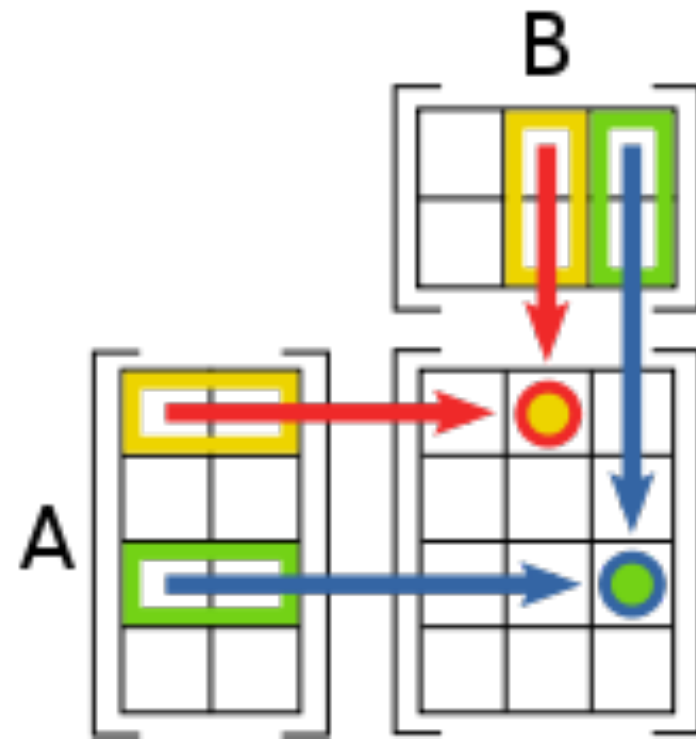


Example 2: separable convolution

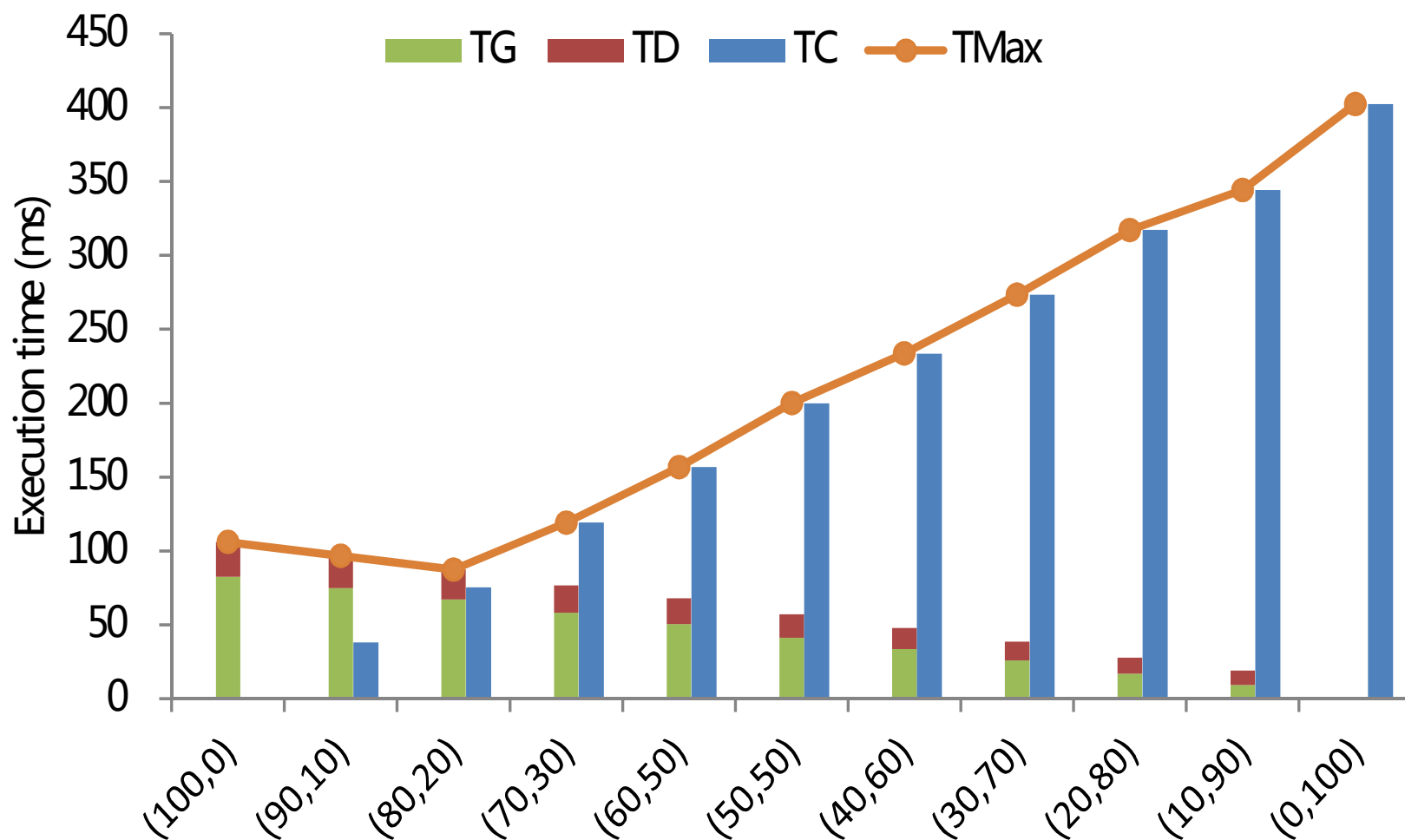


Example 3: matrix multiply

- Matrix multiply
 - Compute the product of 2 matrices
- Performance
 - T_G = execution time on GPU
 - T_C = execution time on CPU
 - T_D = data transfer time CPU-GPU
- GPU best or CPU best?



Example 3: matrix multiply

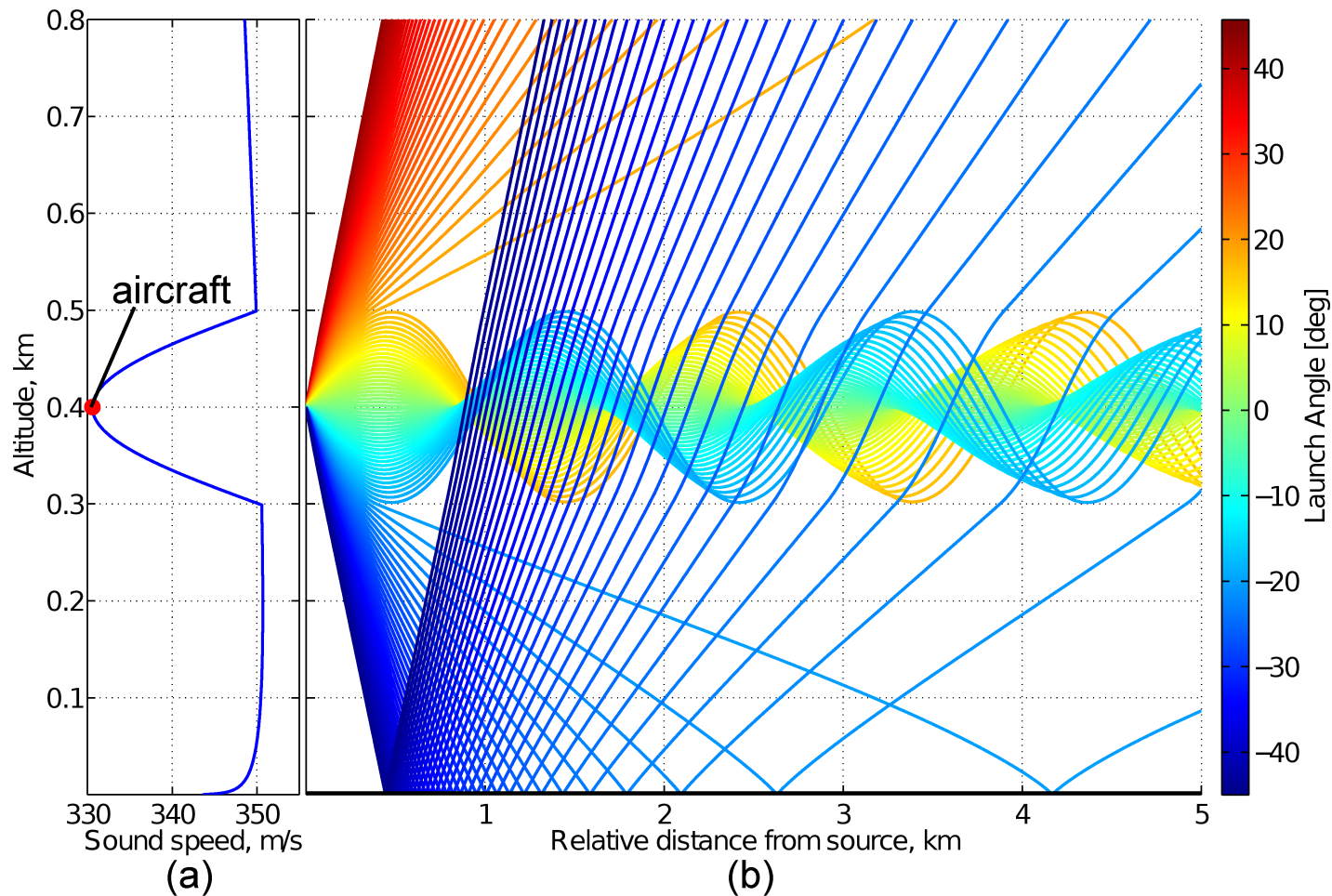




Example 4: Sound ray tracing



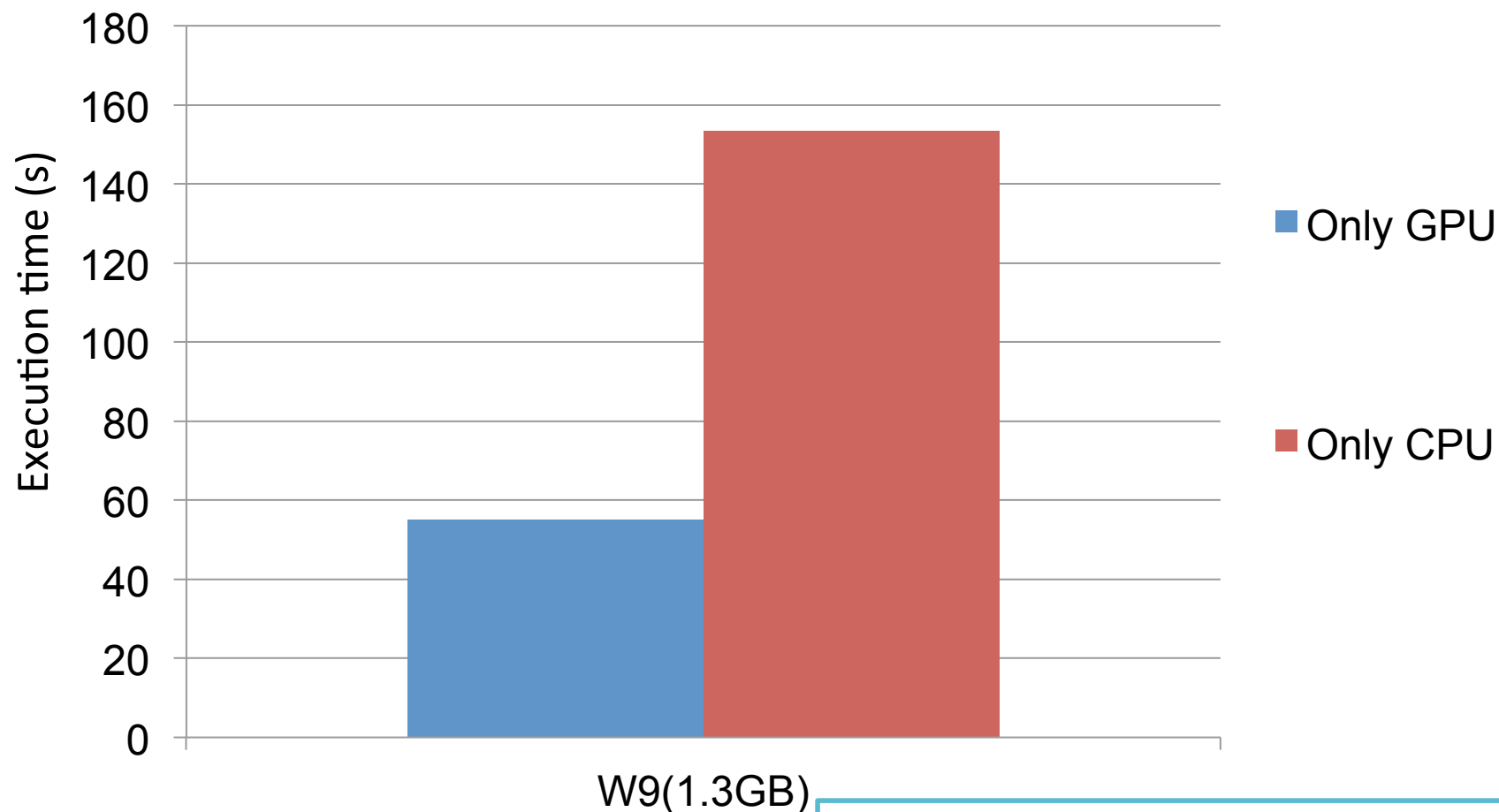
Example 4: Sound ray tracing



Which hardware?

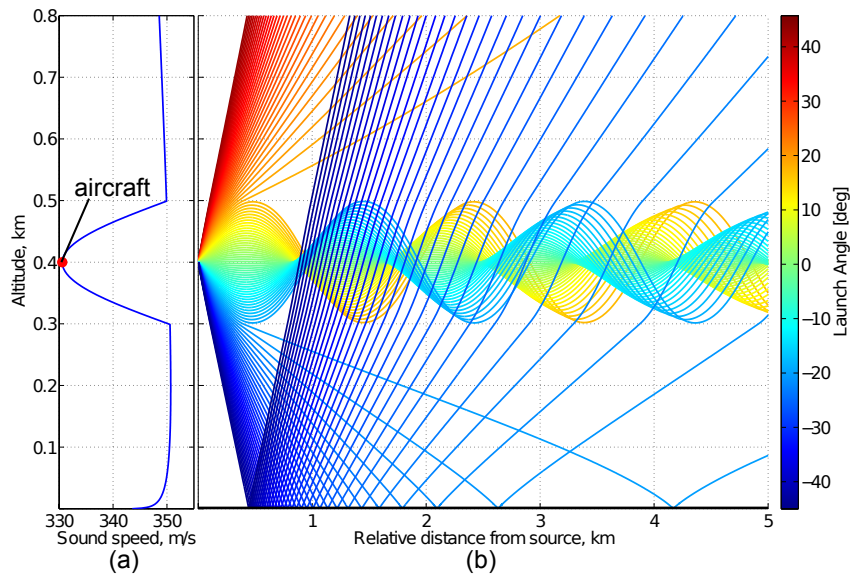
- Our application has ...
 - Massive data-parallelism ...
 - No data dependency between rays ...
 - Compute-intensive per ray ...
-
- ... clearly, this is a perfect GPU workload !!!

Results [1]

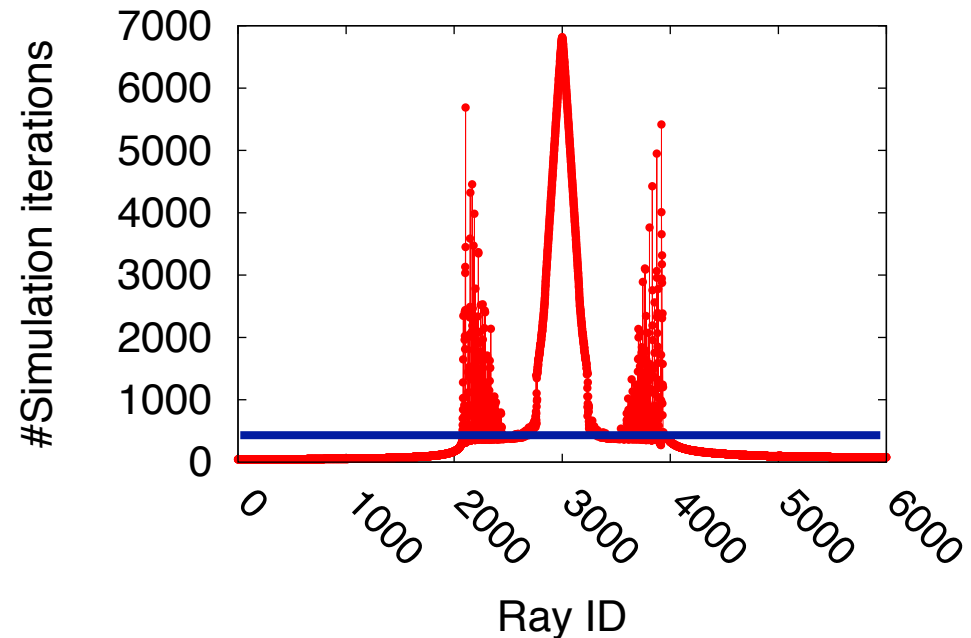


Only 2.2x performance improvement!
We expected 100x ...

Workload profile

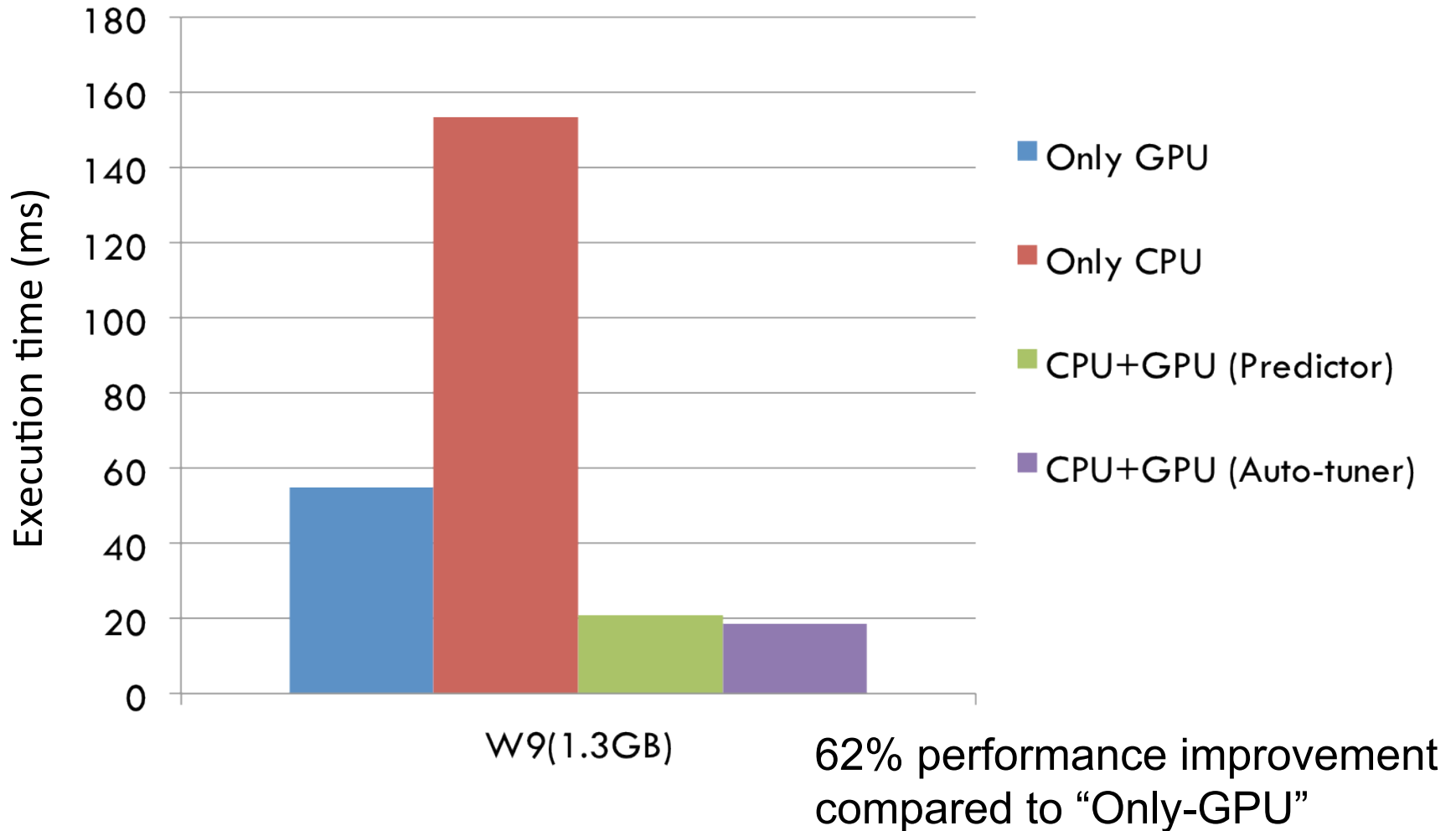


Peak
Processing iterations: ~ 7000



Bottom
Processing iterations: ~ 500

Results [2]



So ...

- There are **very** few GPU-only applications
 - CPU – GPU communication bottleneck.
 - Increasing performance of CPUs
- A part of the computation can be done by the CPU.
 - How to program an application to enable this?
 - Which part?

Main challenges: **programming** and **workload partitioning!**

PART III

Challenge 2: Programming

Programming models (PMs)

- Heterogeneous computing = a mix of different processors on the same platform.
- Programming
 - Single programming model (unified)
 - OpenCL is a popular choice
 - Mix of programming models
 - One(/several?) for CPUs – OpenMP
 - One(/several?) for GPUs – CUDA

Low level



OpenCL

High level

OpenACC
Directives for Accelerators

OpenMP 4.0

Heterogeneous Programming Library



OmpSs

Heterogeneous computing PMs

- CUDA + OpenMP/TBB
 - Typical combination for NVIDIA GPUs
 - Individual development per *PU
 - Glue code can be challenging
- OpenCL
 - Functional portability => can be used as a unified model
 - Performance portability via code specialization
- HPL
 - Library on top of OpenCL, to automate code specialization

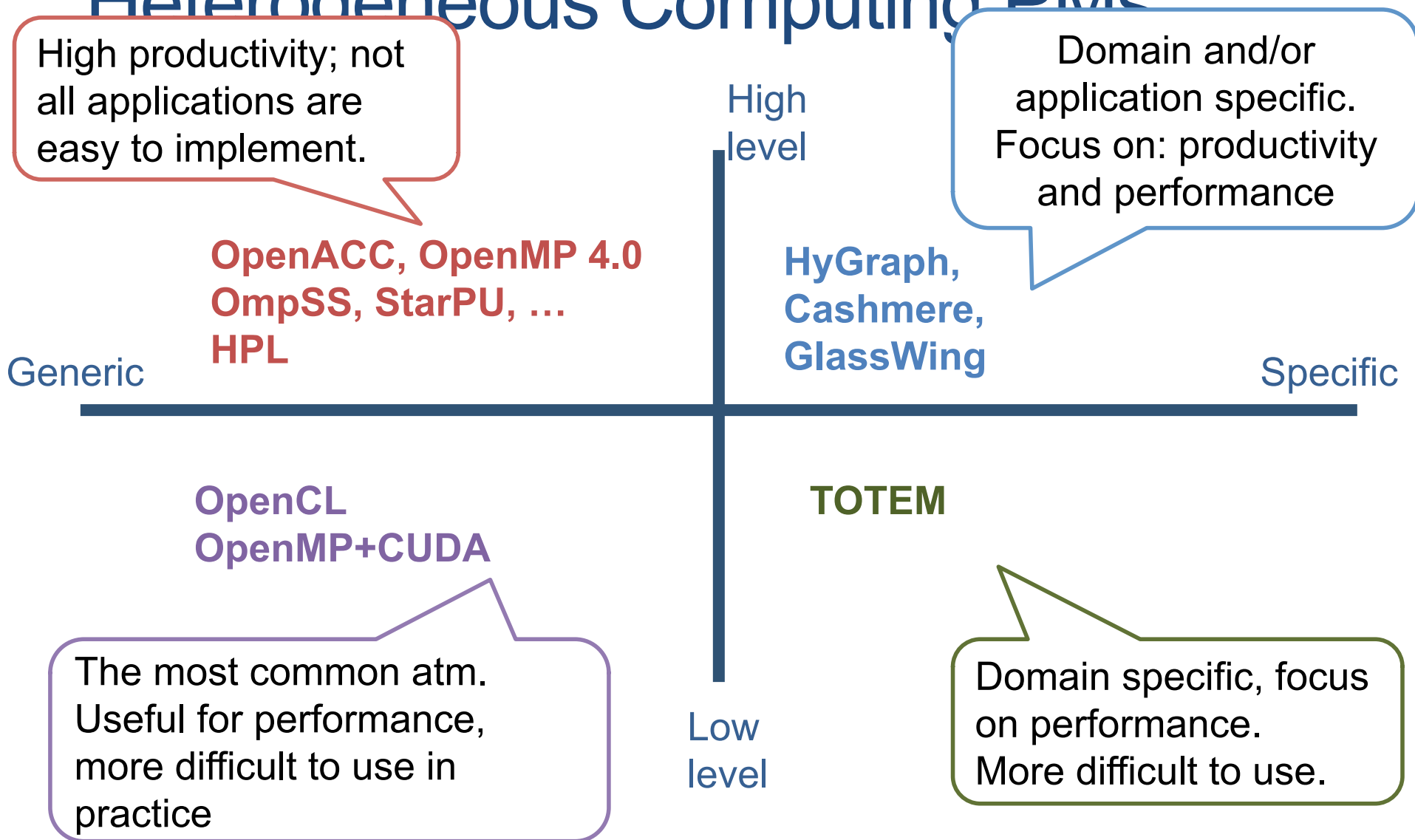
Heterogeneous computing PMs

- StarPU
 - Special API for coding
 - Runtime system for scheduling
- OmpSS
 - C + OpenCL/CUDA kernels
 - Runtime system for scheduling and communication optimization

Heterogeneous computing PMs

- Cashmere
 - Dedicated to Divide-and-conquer solutions
 - OpenCL backend.
- GlassWing
 - Dedicated to MapReduce applications
- TOTEM
 - Graph processing
 - CUDA+Multi-threading
- HyGraph
 - Graph processing
 - Based on CUDA+OpenMP

Heterogeneous Computing DMs



End of part III

- Questions ?

PART III

Workload partitioning models