



Hasso  
Plattner  
Institut

IT Systems Engineering | Universität Potsdam

# Code Mobilität: Adaptive Programmiermodelle für GPU-Computing

Frank Feinbube

Hasso Plattner Institute  
Operating Systems and Middleware  
Prof. Dr. Andreas Polze





# The beautiful new world of Hybrid Compute Environments?

4



# The beautiful new world of Hybrid Compute Environments?

5

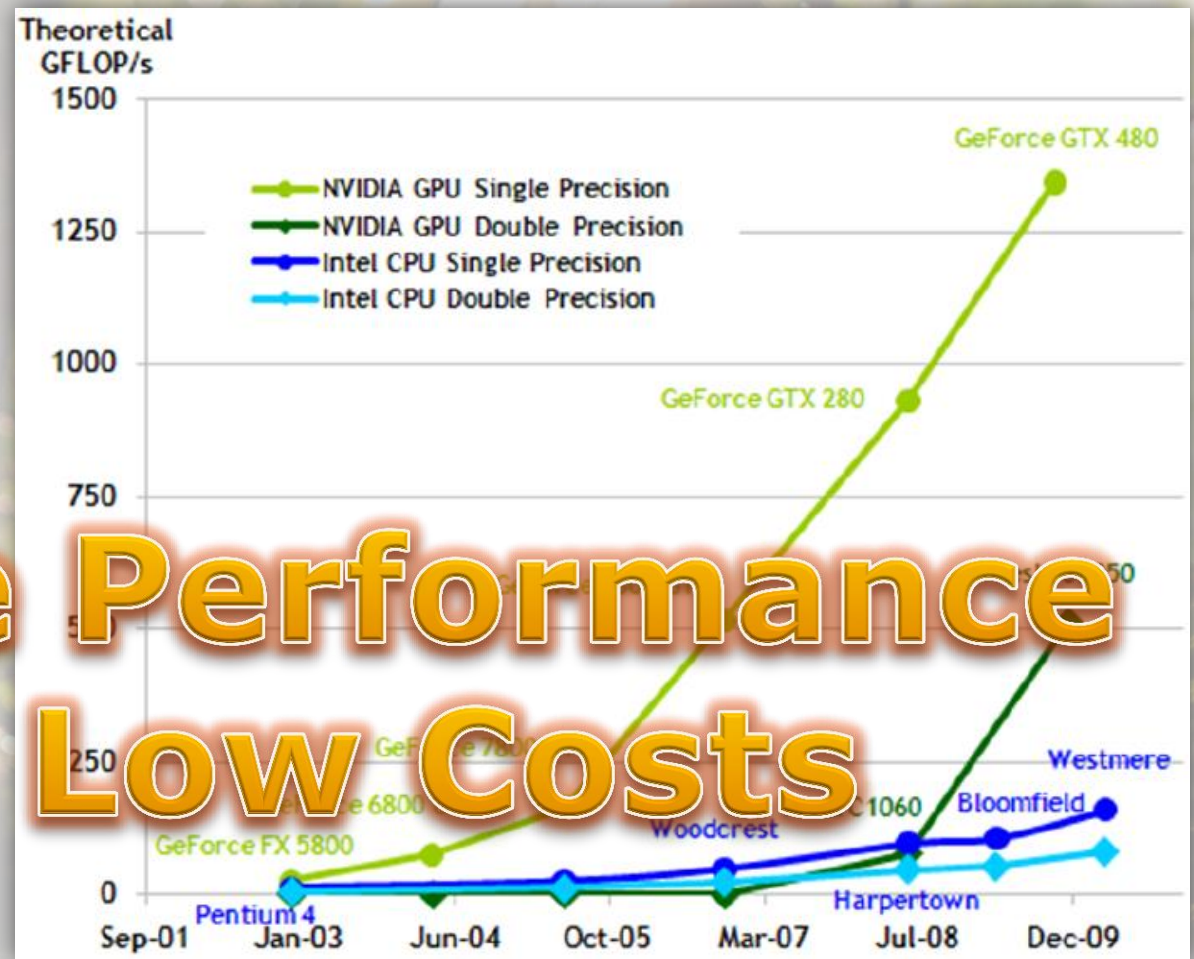
Medical, Finance  
Military, Research  
Video and Photo  
Energy, ....

**Variety of Successful Applications**

[http://www.nvidia.com/object/cuda\\_testimonials.html](http://www.nvidia.com/object/cuda_testimonials.html)

# The beautiful new world of Hybrid Compute Environments?

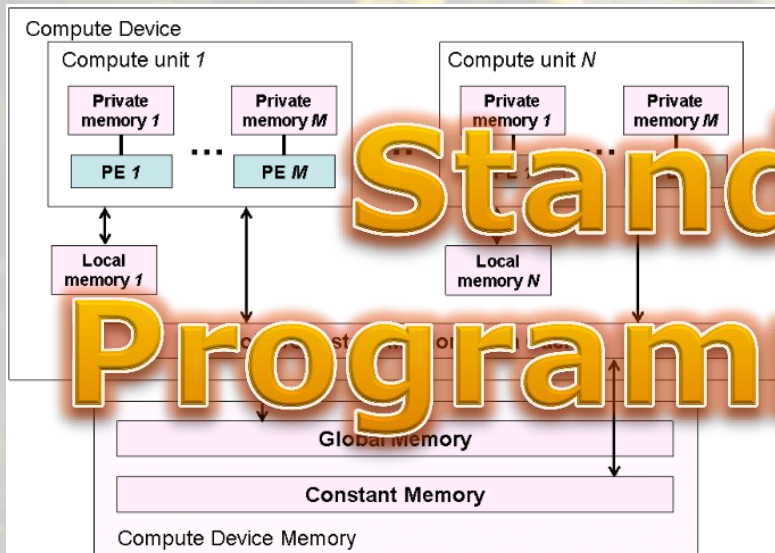
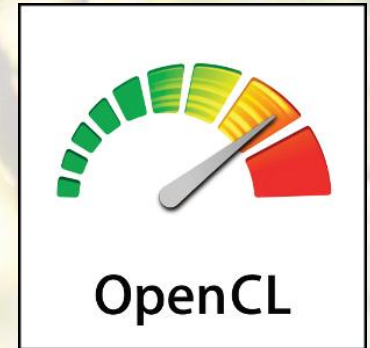
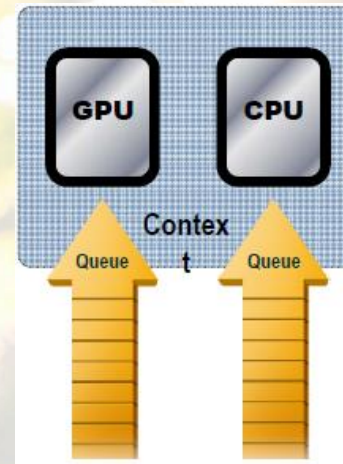
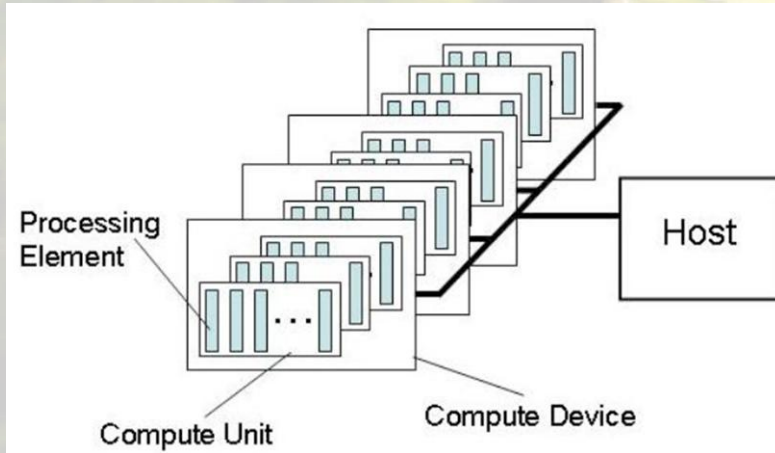
6



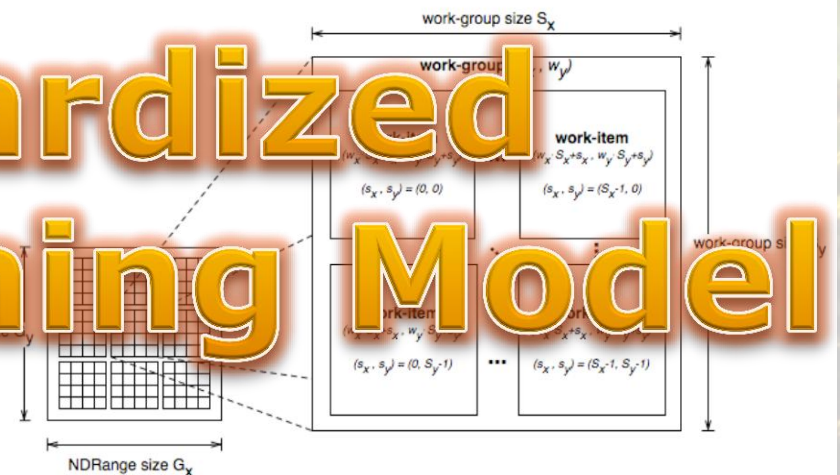
**Huge Performance + Low Costs**

# The beautiful new world of Hybrid Compute Environments?

7



# Standardized Programming Model



# History of GPU Computing

8

## Fixed Function Graphic Pipelines

- 1980s-1990s; configurable, not programmable; first APIs (DirectX, OpenGL); Vertex Processing

## Programmable Real-Time Graphics

- Since 2001: APIs for Vertex Shading, Pixel Shading and access to texture; DirectX9

## Unified Graphics and Computing Processors

- 2006: NVIDIAs G80; unified processors arrays; three programmable shading stages; DirectX10

## General Purpose GPU (GPGPU)

- compute problem as native graphic operations; algorithms as shaders; data in textures

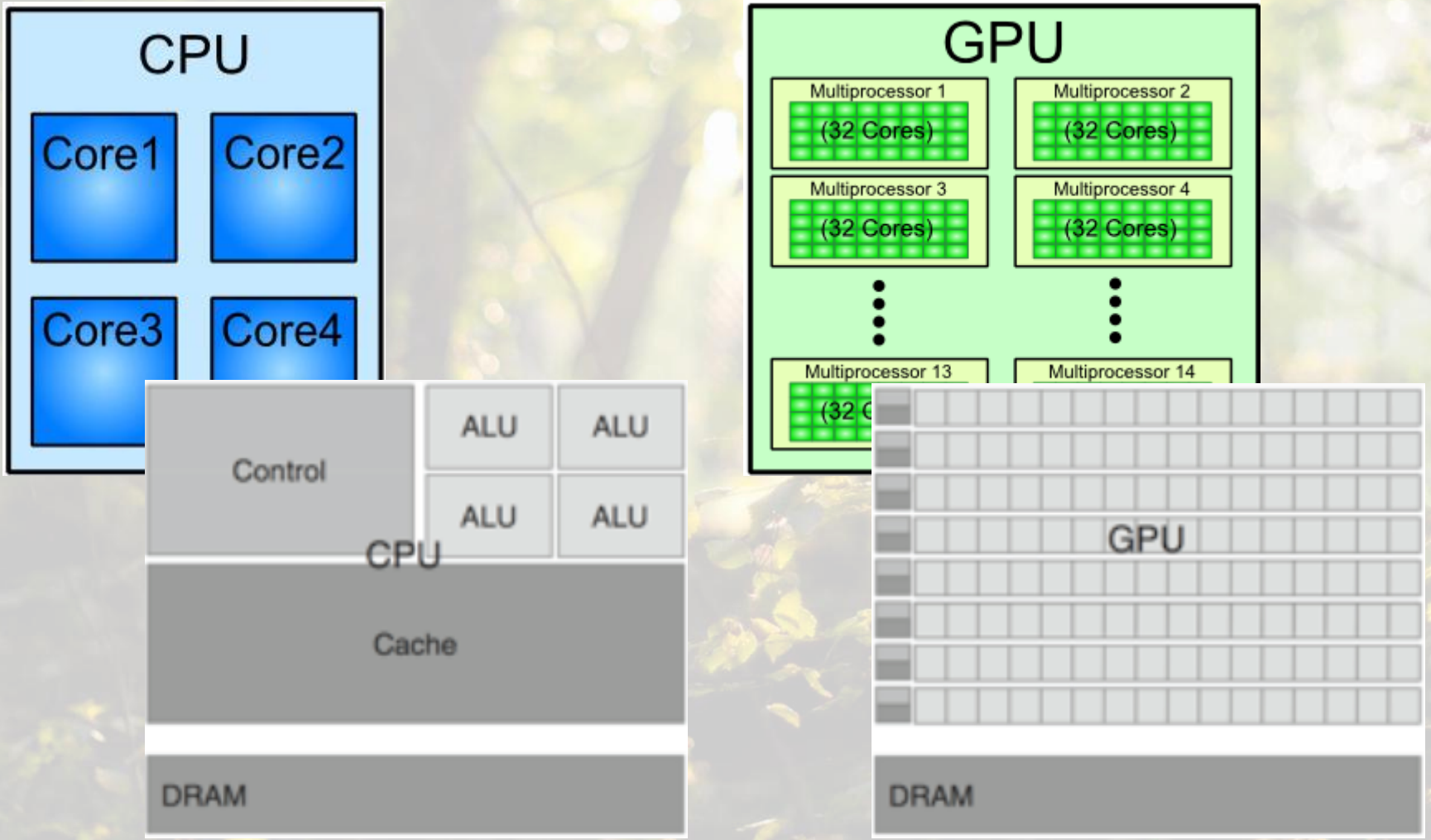
## GPU Computing

- Programming CUDA; shaders programmable; load and store instructions; barriers; atomics



# CPU vs. GPU Architecture

9



# GF100

Host Interface

GigaThread Engine

Memory Controller

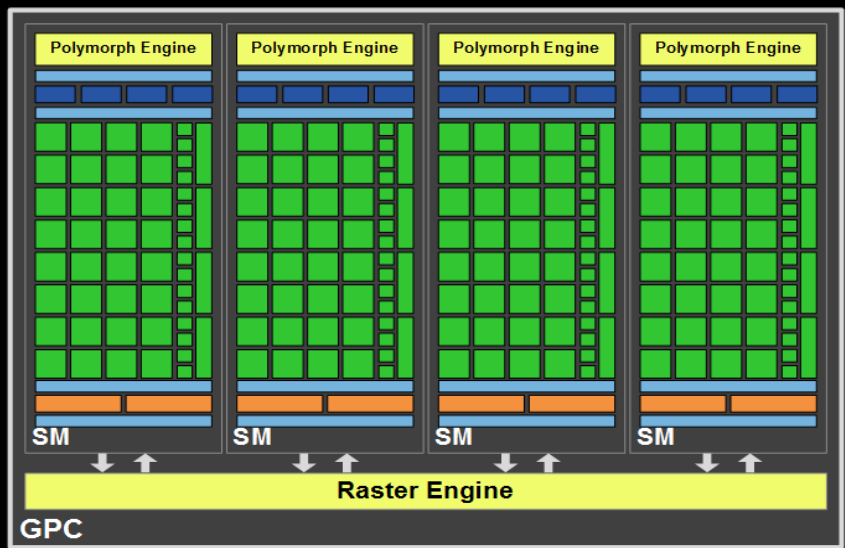
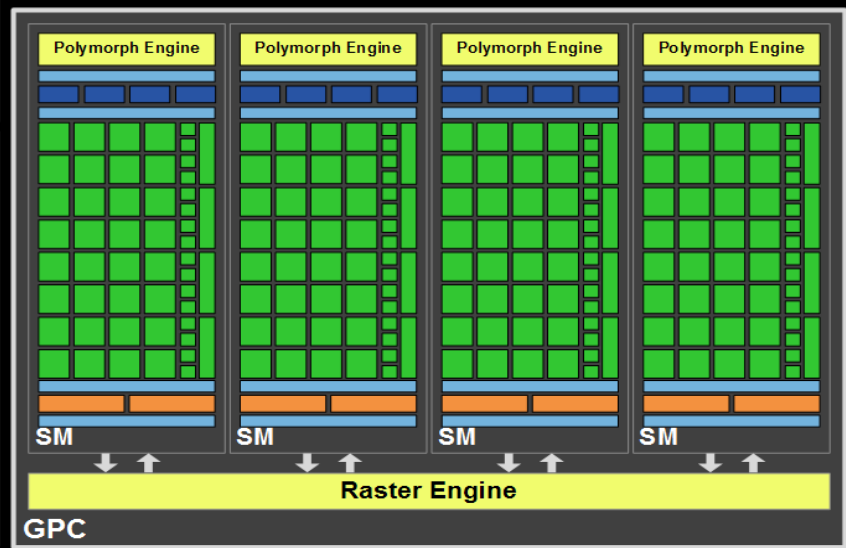
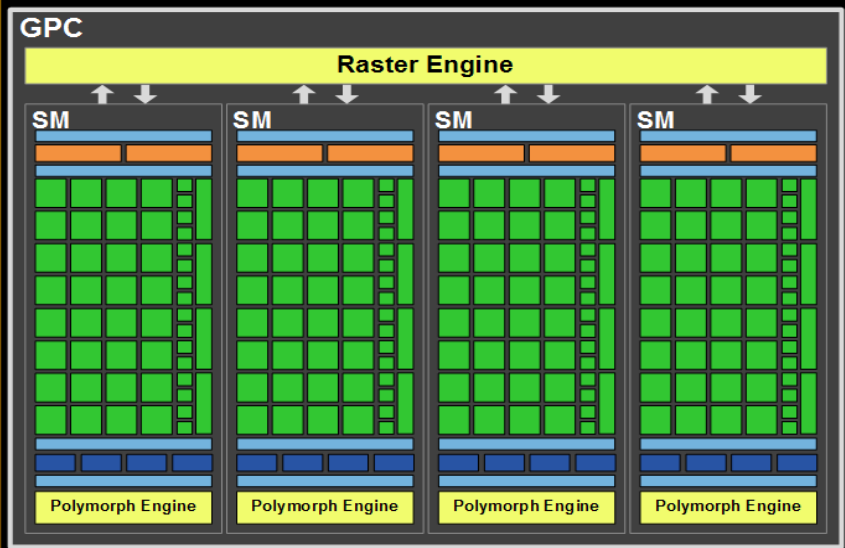
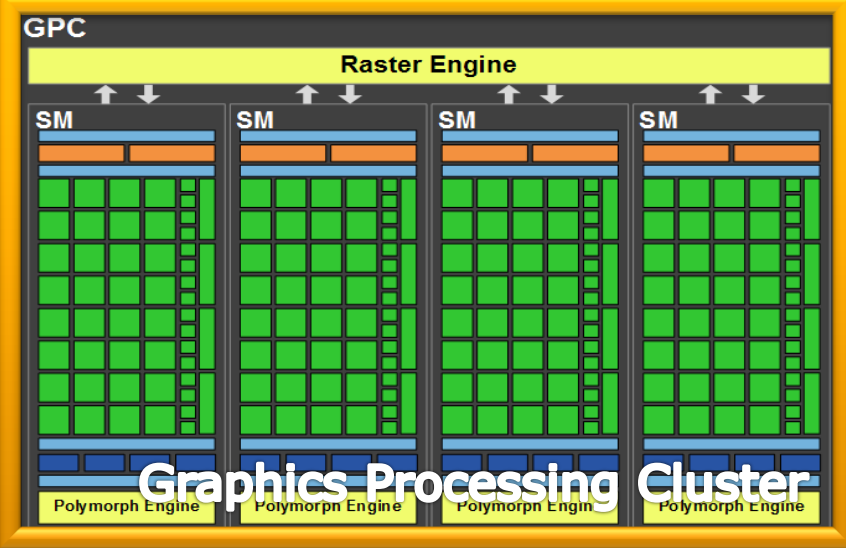
Memory Controller

Memory Controller

Memory Controller

Memory Controller

Memory Controller



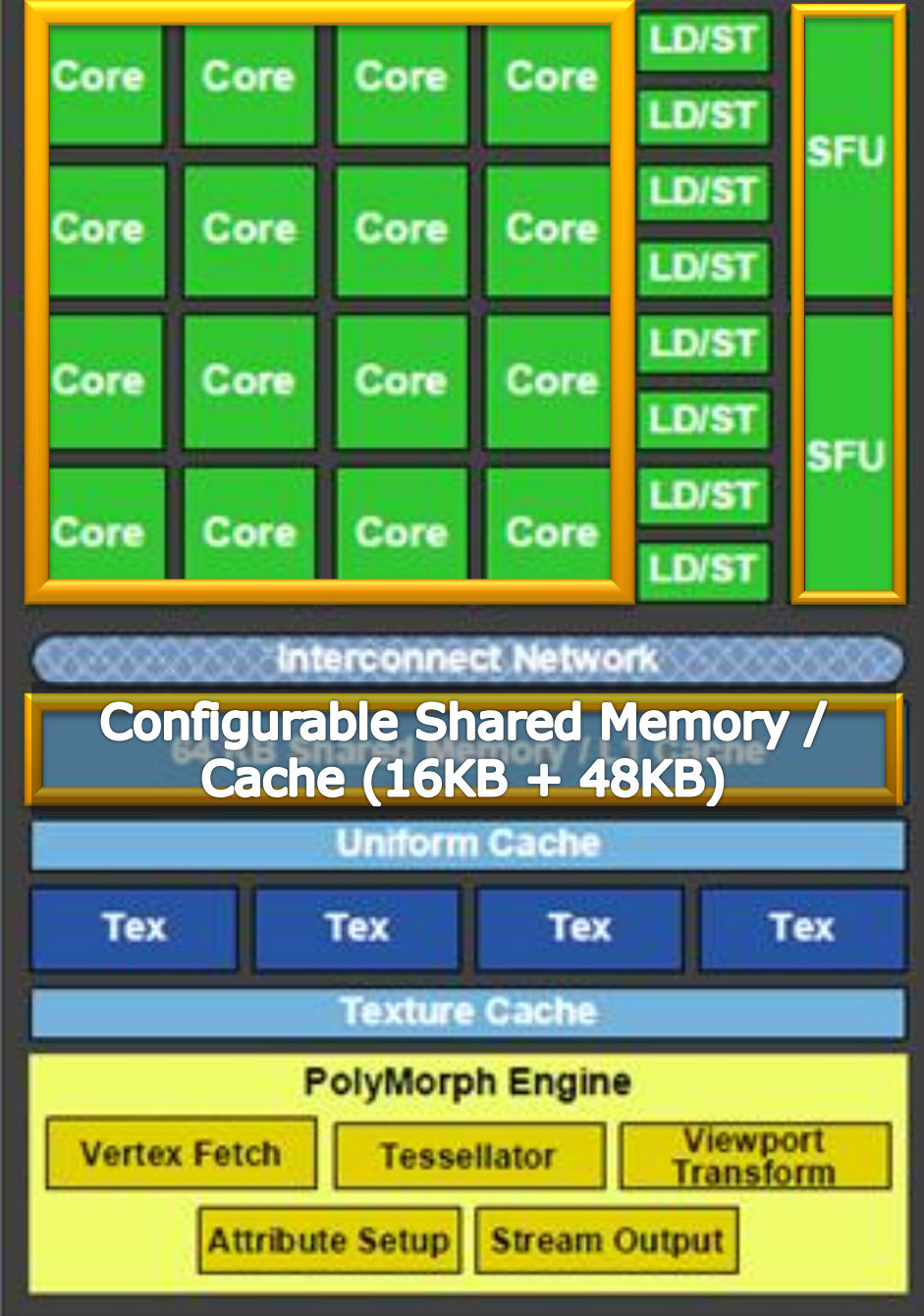
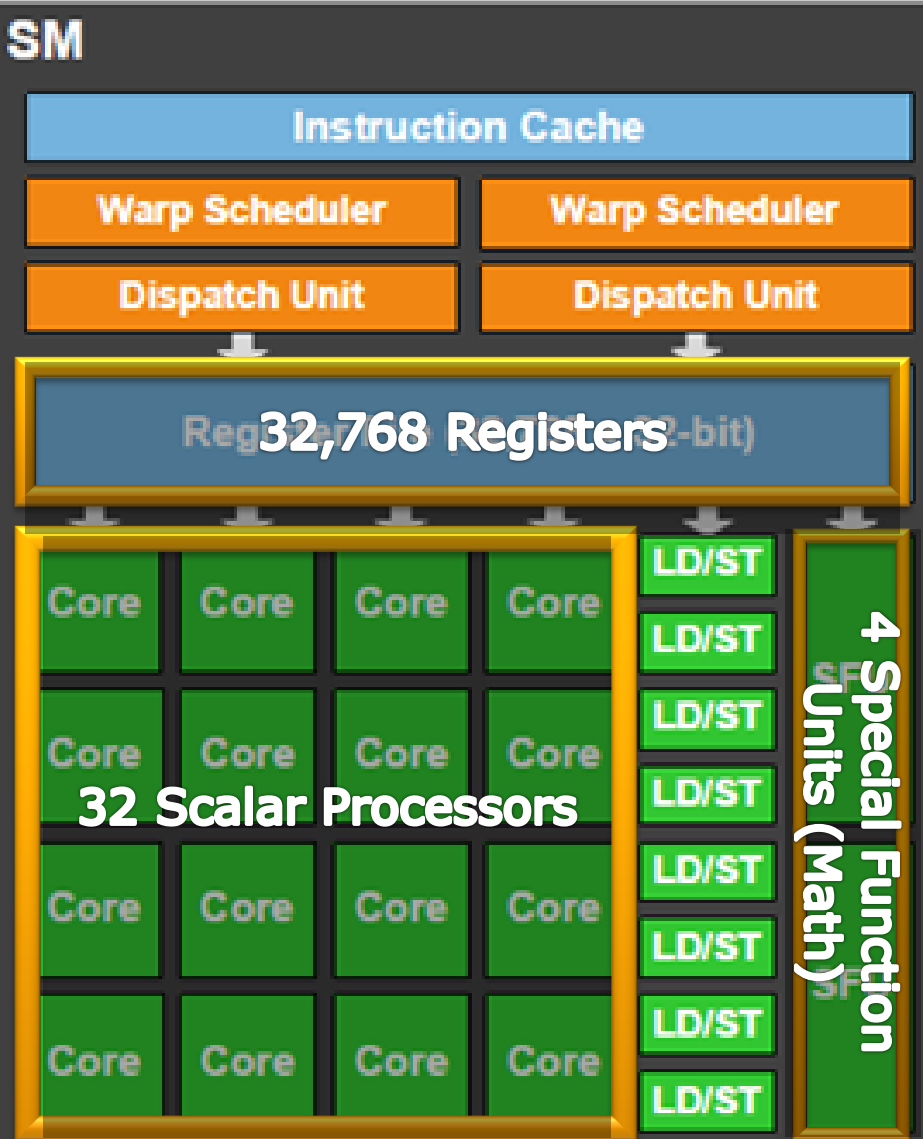
# GF100



Streaming Multiprocessor (SM)

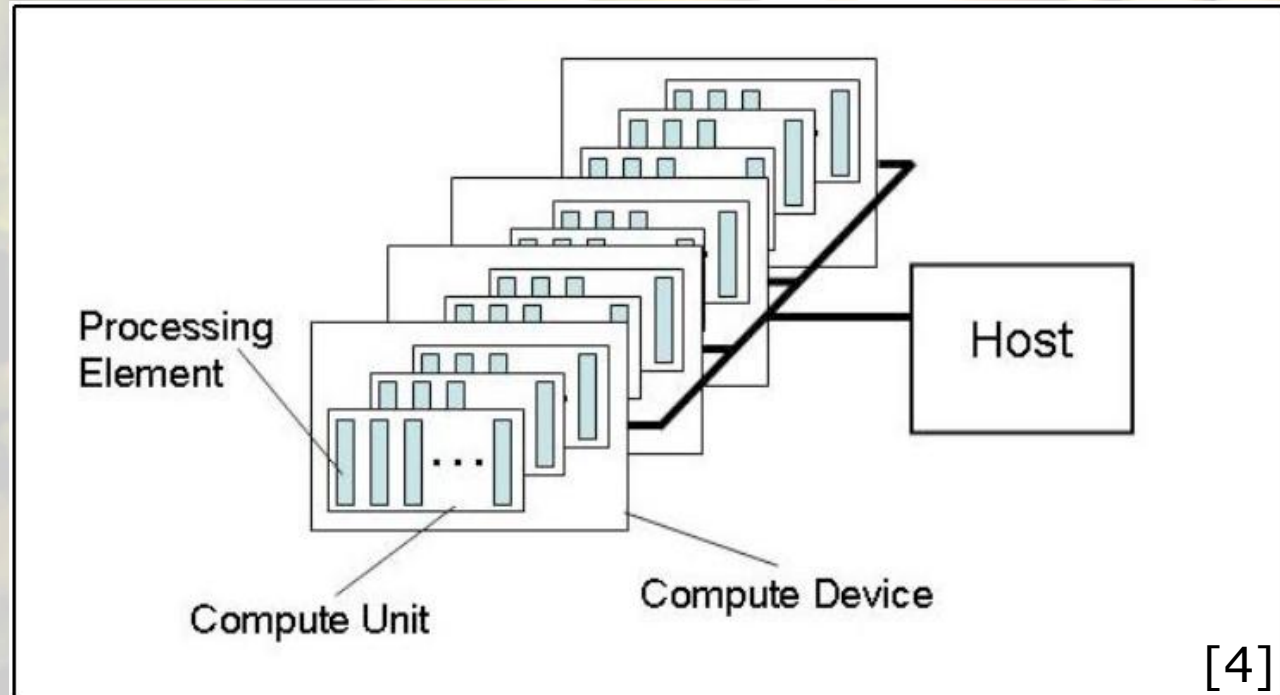
Scalar Processors (SP) = Cores

# GF100



# OpenCL Platform Model

13



- OpenCL exposes CPUs, GPUs, and other Accelerators as “devices”
- Each “device” contains one or more “compute units”, i.e. cores, SMs,...
- Each “compute unit” contains one or more SIMD “processing elements”

# OpenCL Execution Model

14

An OpenCL kernel is executed by an array of work items.

- All work items run the same code (SPMD)
- Each work item has an index that it uses to compute memory addresses and make control decisions



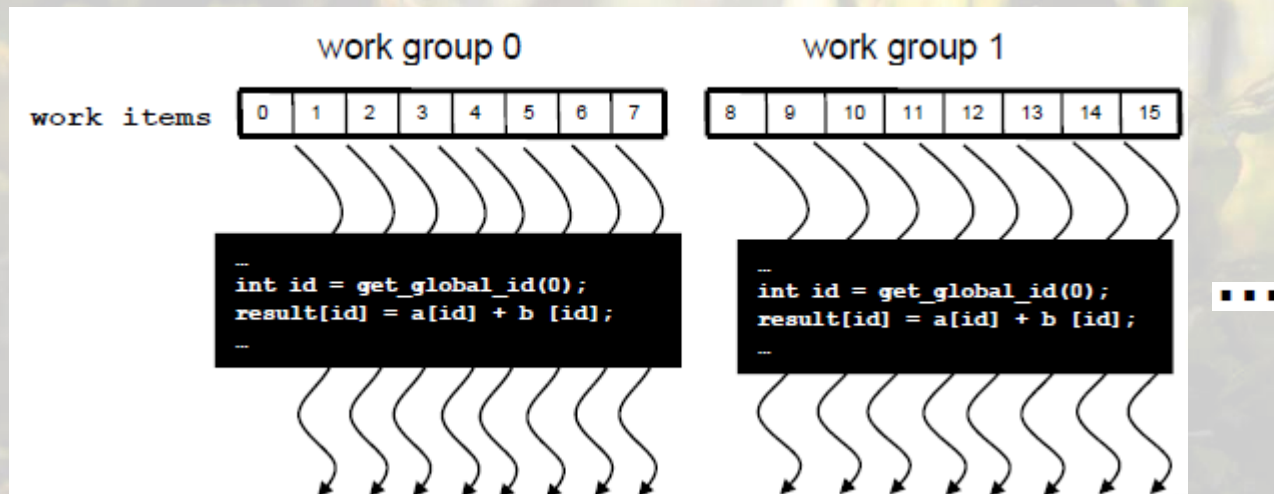
[1]

# Work Groups: Scalable Cooperation

15

Divide monolithic work item array into work groups

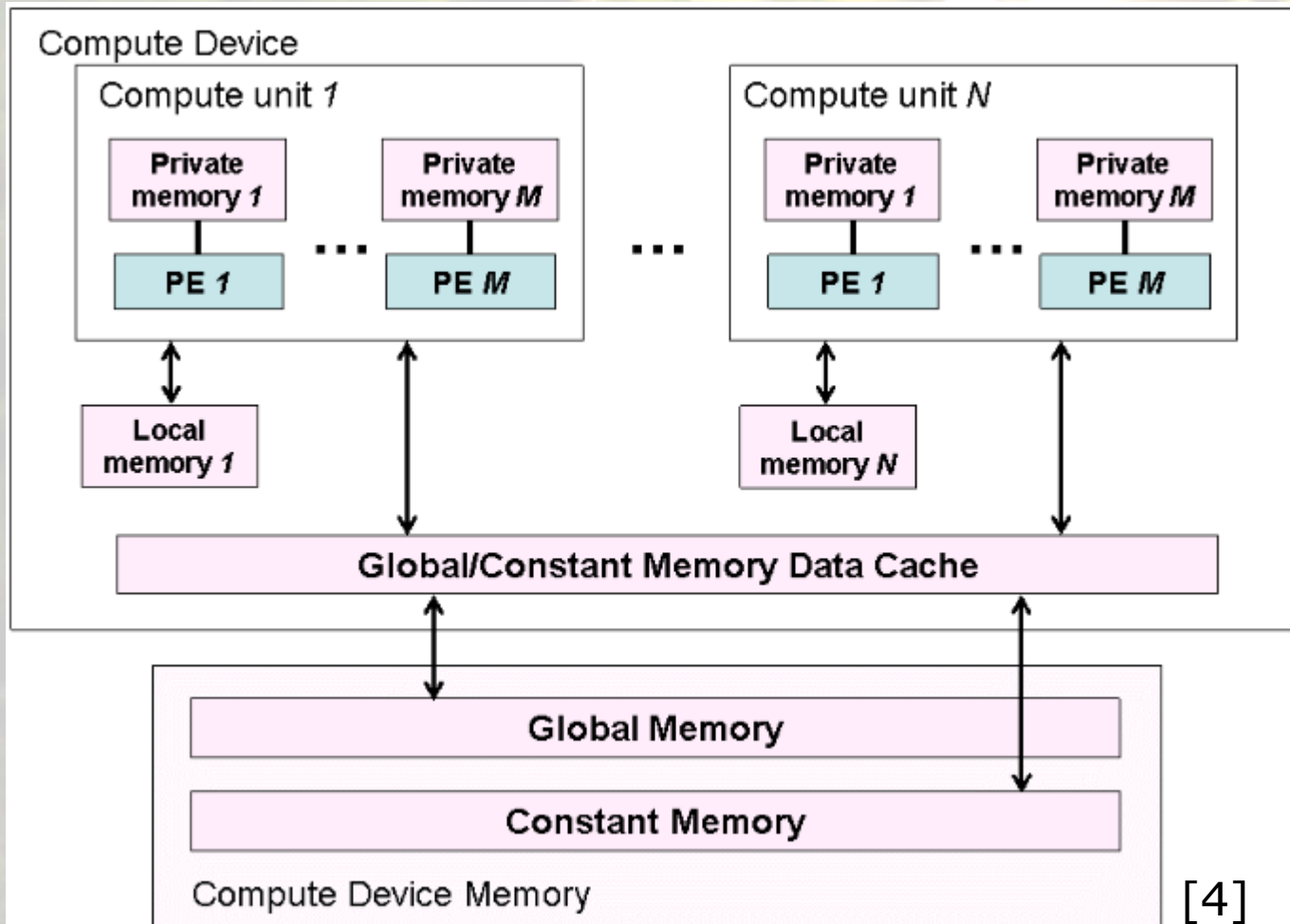
- Work items within a work group cooperate via **shared memory, atomic operations** and **barrier synchronization**
- Work items in different work groups cannot cooperate



[1]

# OpenCL Memory Architecture

16



**Private**  
Per work-item

**Local**  
Shared within  
a workgroup

**Global/  
Constant**  
Visible to  
all workgroups

**Host Memory**  
On the CPU



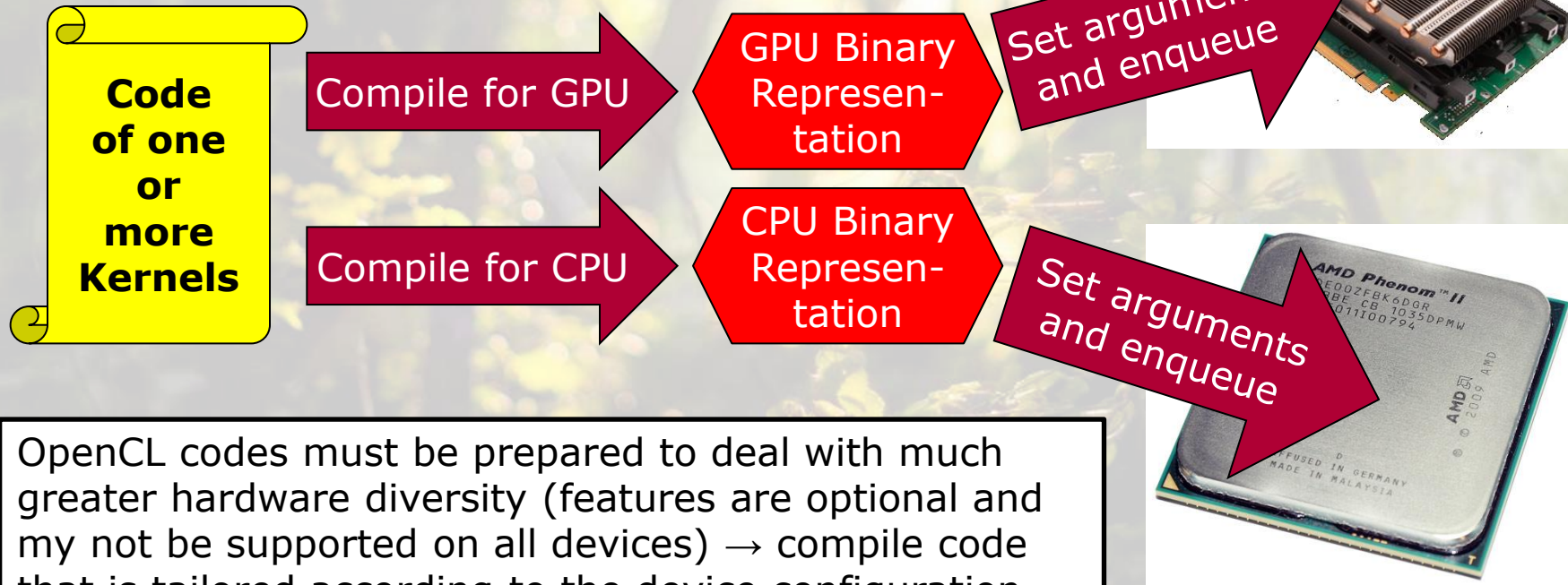
# Building and Executing OpenCL Code

17

## Kernel

## Program

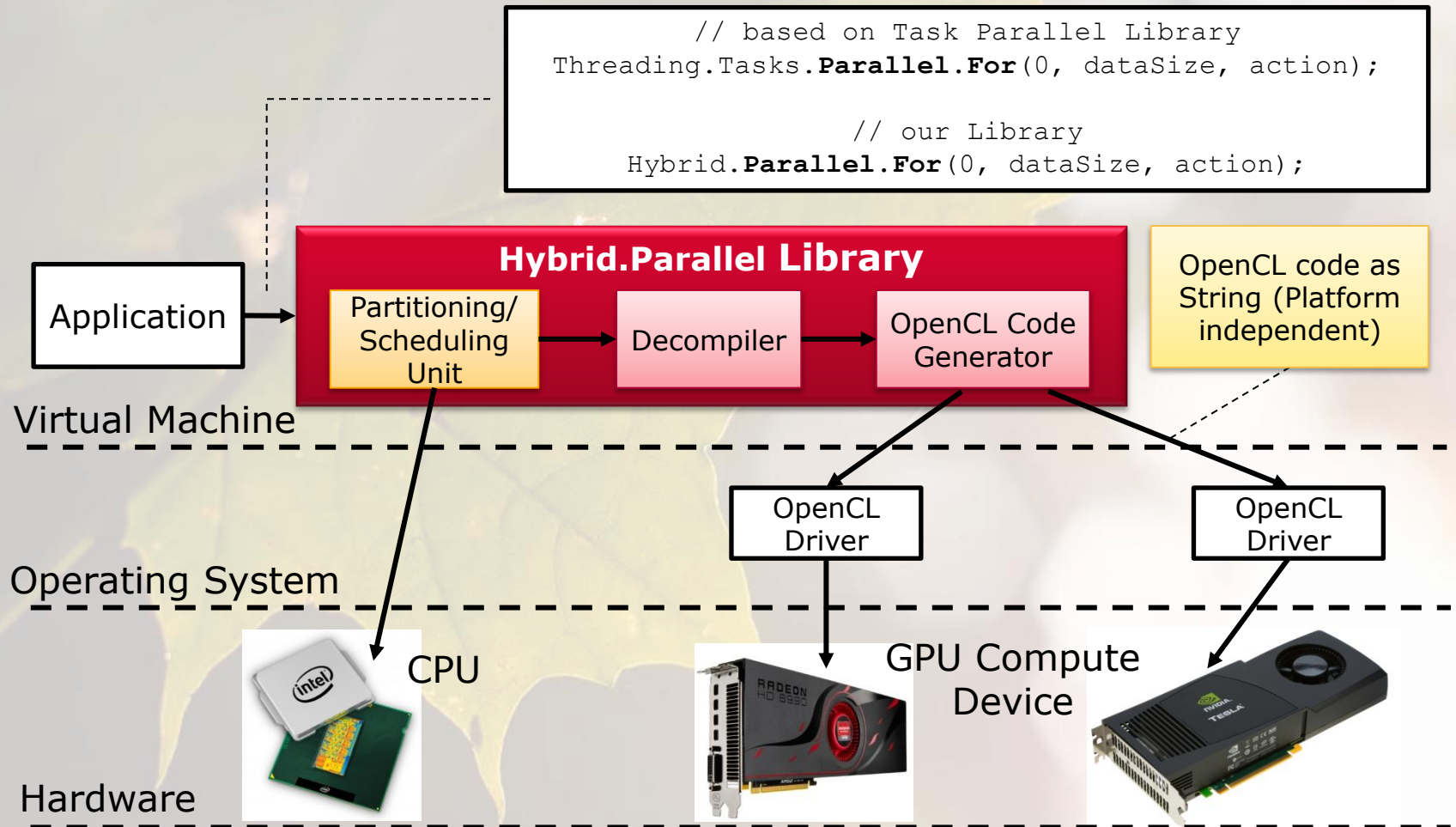
## Device

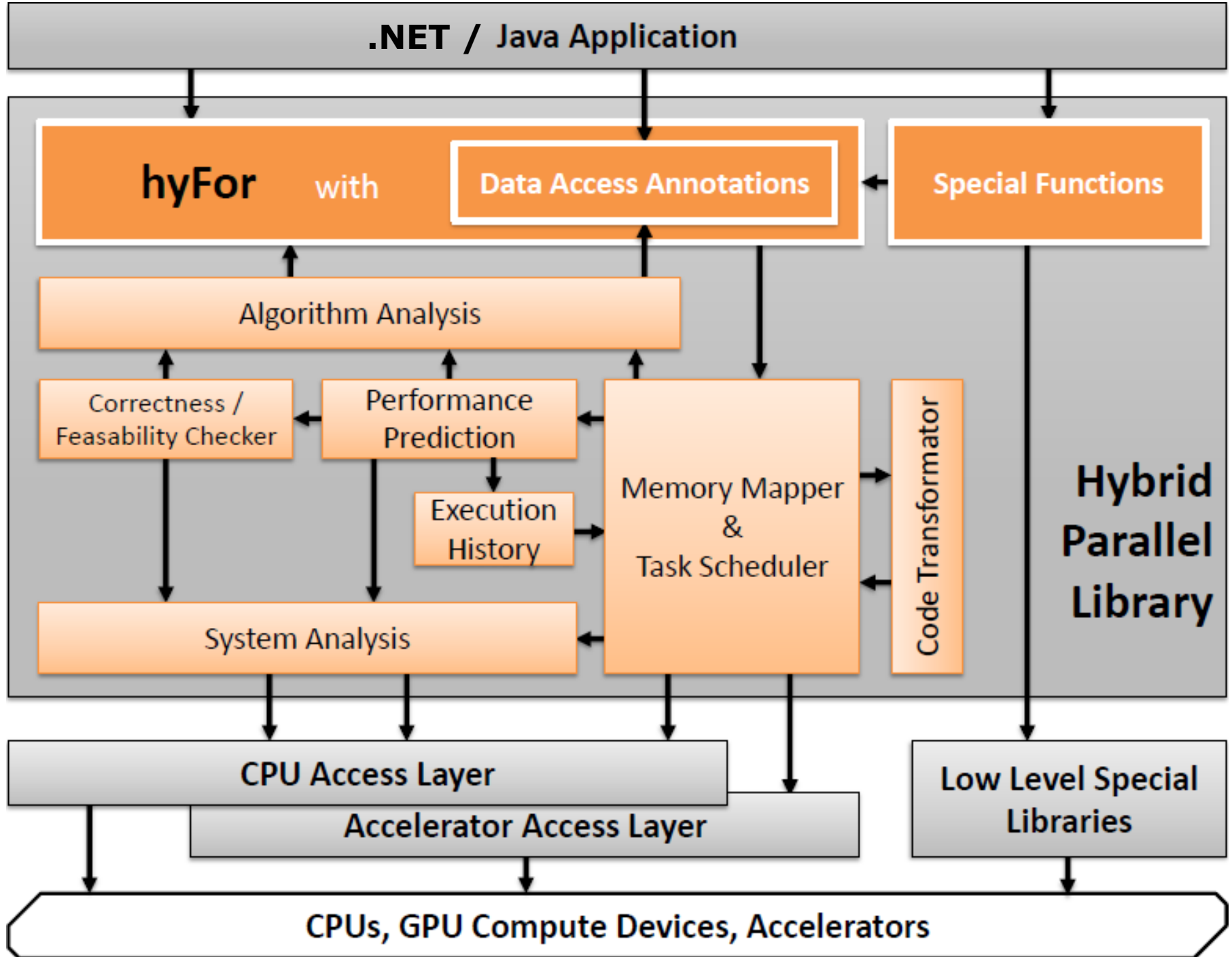


OpenCL codes must be prepared to deal with much greater hardware diversity (features are optional and may not be supported on all devices) → compile code that is tailored according to the device configuration

# Hybrid Computing for High-Level Programmers: The Hybrid Parallel Project

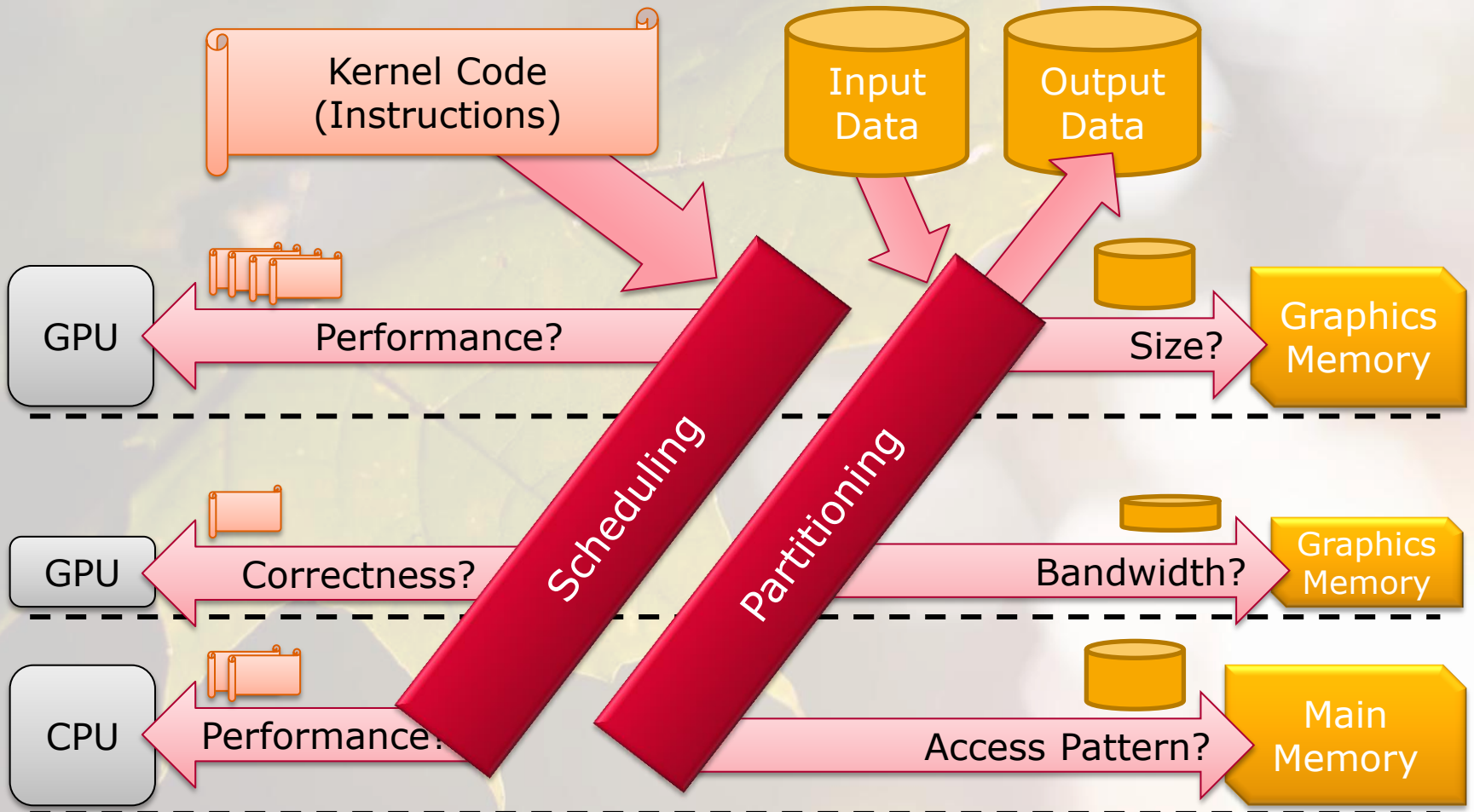
18





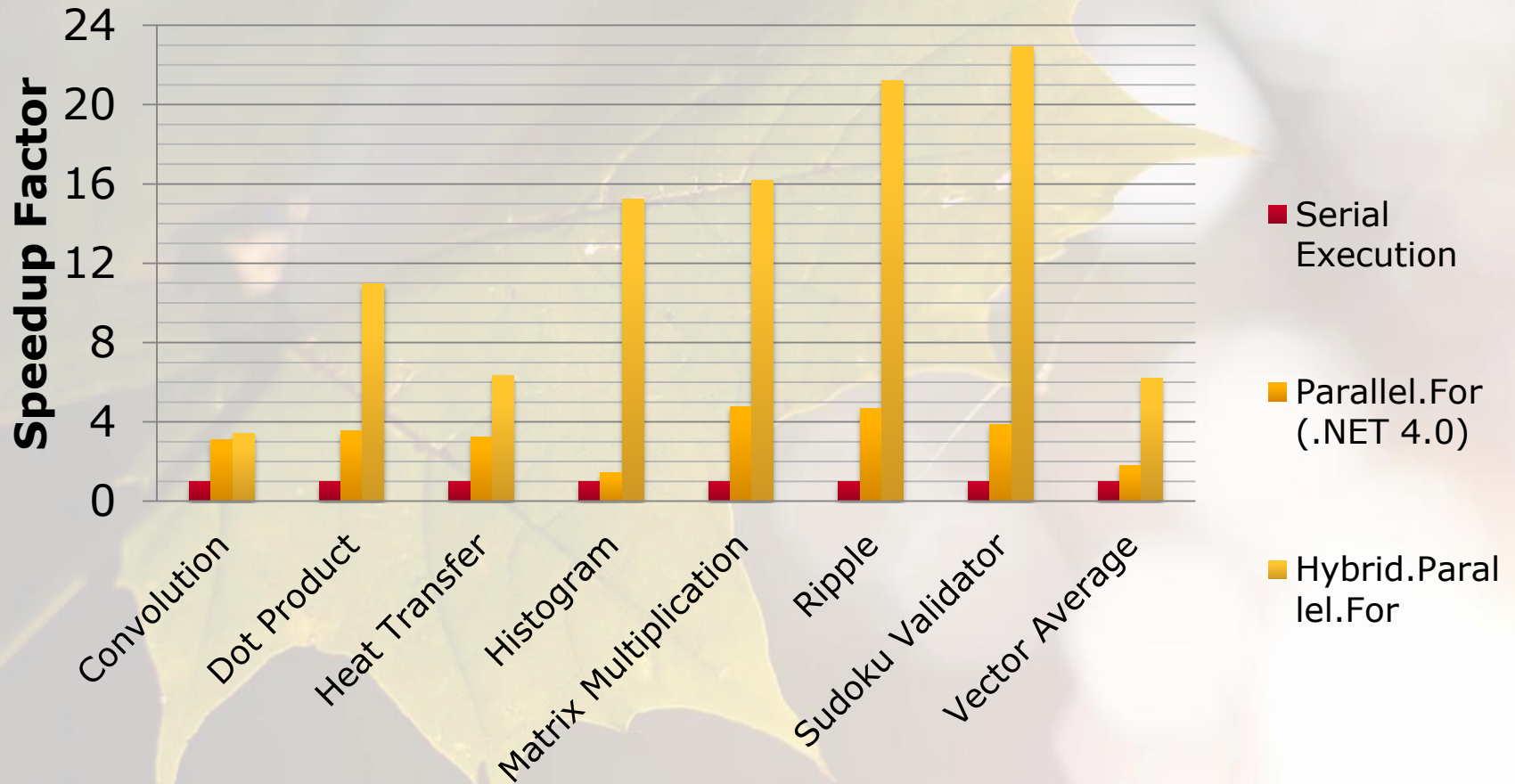
# Hybrid Computing for High-Level Programmers: The Hybrid Parallel Project

20



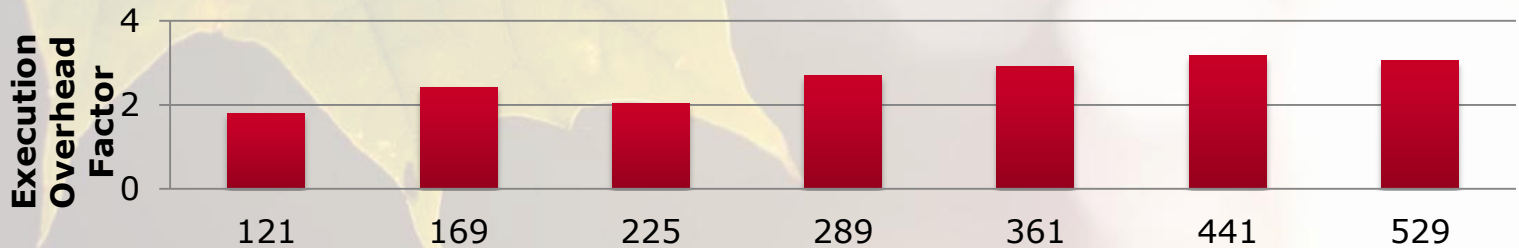
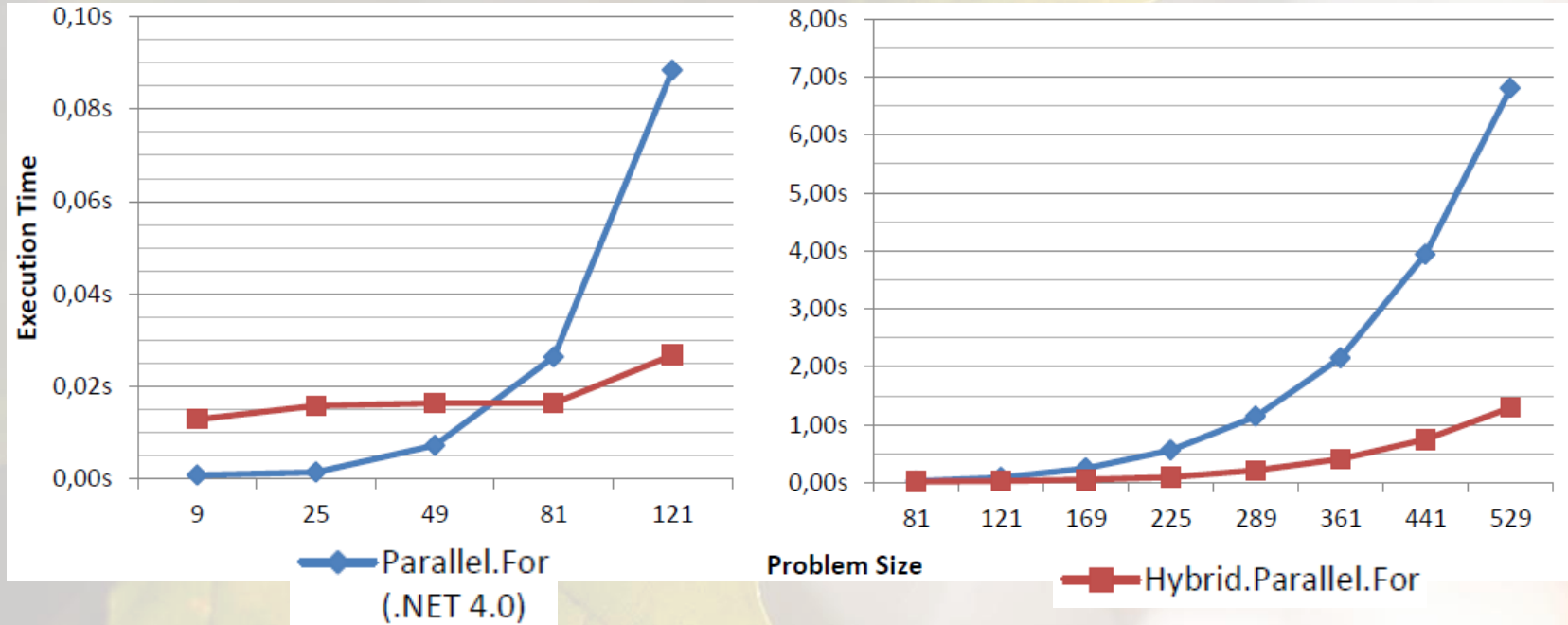
# Hybrid Parallel For .NET Evaluation

21



# Hybrid Parallel For .NET Evaluation

22



## „Operating Systems must support GPU abstractions“

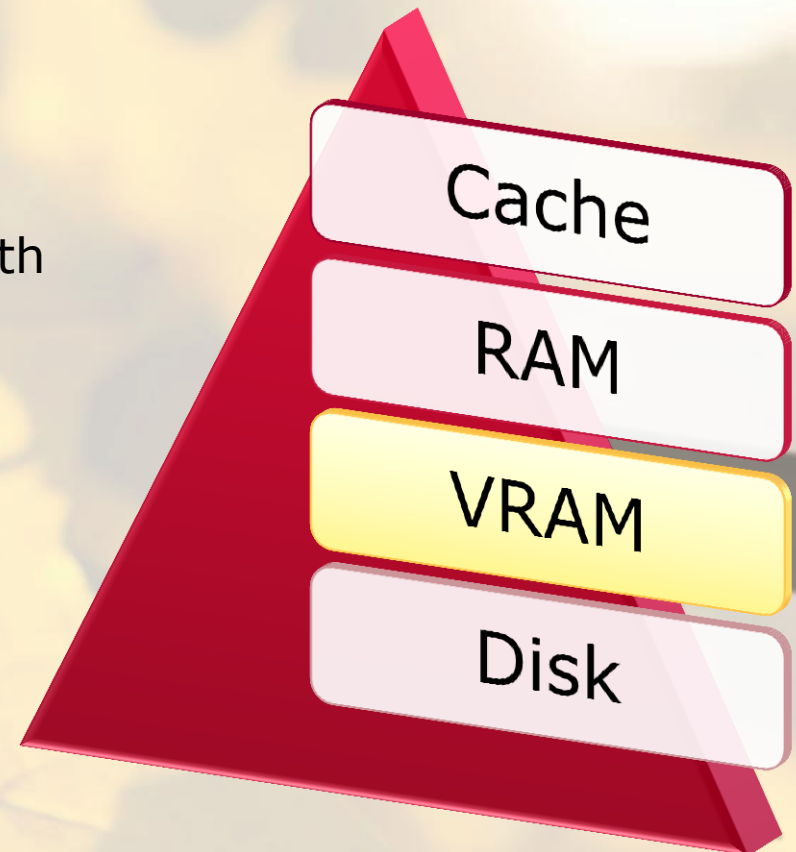
– *Rossbach, Currey, Witchel [HotOS11]*

- **PTask** API by Rossbach et al. (based on CUDA)
  - Workload as acyclic graphs of ptasks → executed on GPU
  - Proof-of-concept: encryption / decryption for Linux EncFS
- **GPUStore** by Sun et al. (based on CUDA)
  - Specially tailored for storage systems
  - GPUs for encryption and RAID functionality
- **Gdev** framework by Kato et al. (based on Direct Rendering Int.)
  - Reverse engineering of graphics card driver internals

# GPU Page Cache

24

- Use Global Memory as paging backing store
  - Typical desktop systems with VRAM of 1+ GB
  - VRAM unused most of the time
- Decrease paging operation latencies

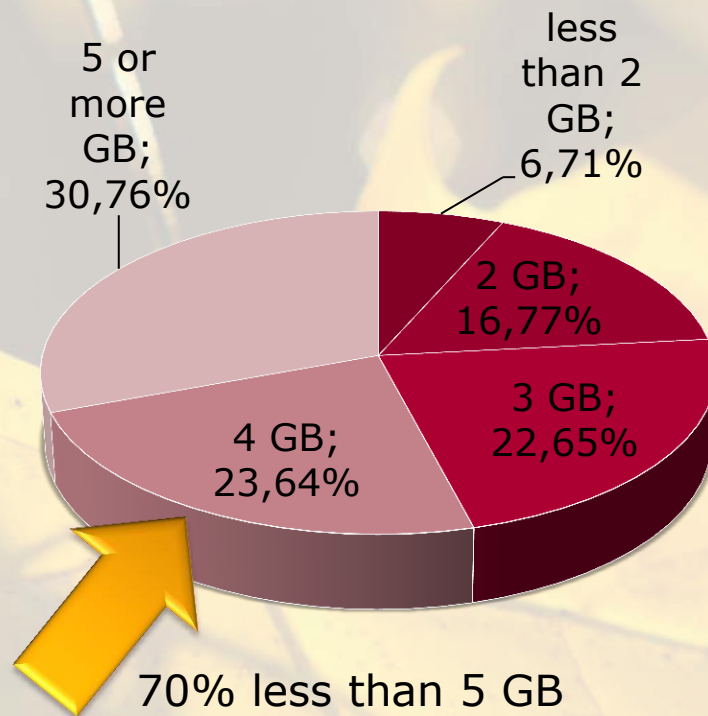




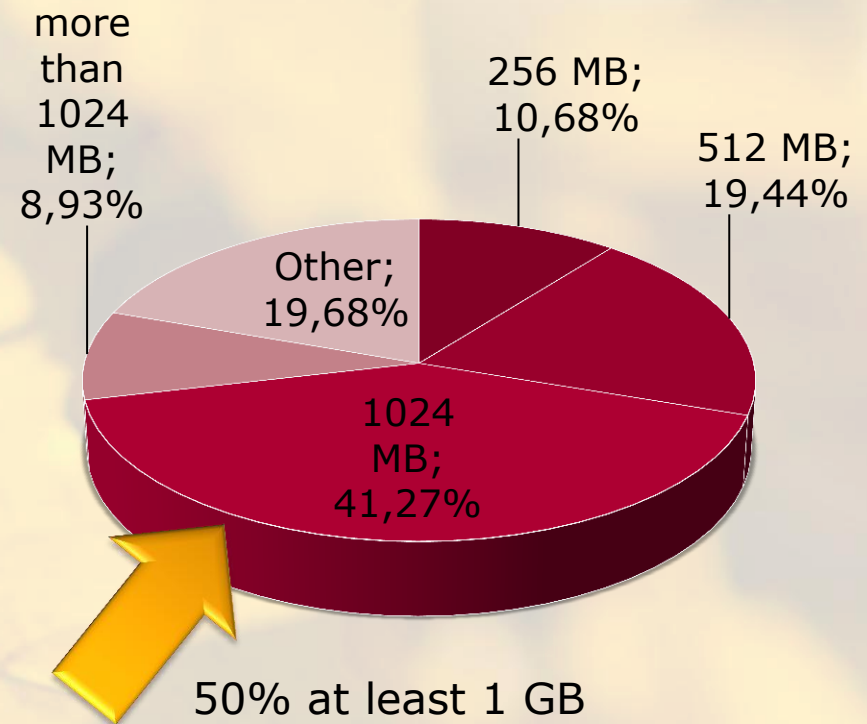
# Hardware Survey

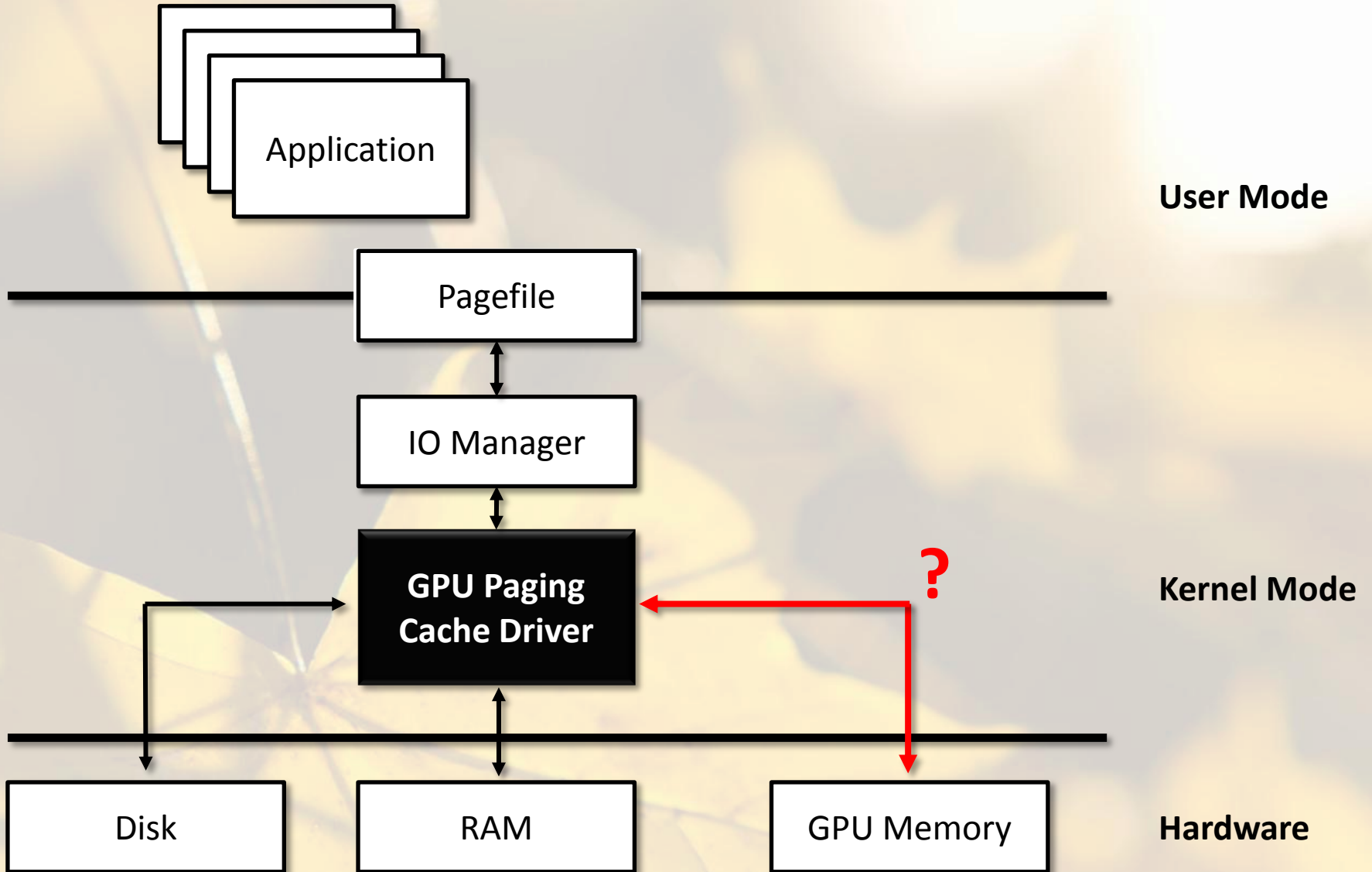
25

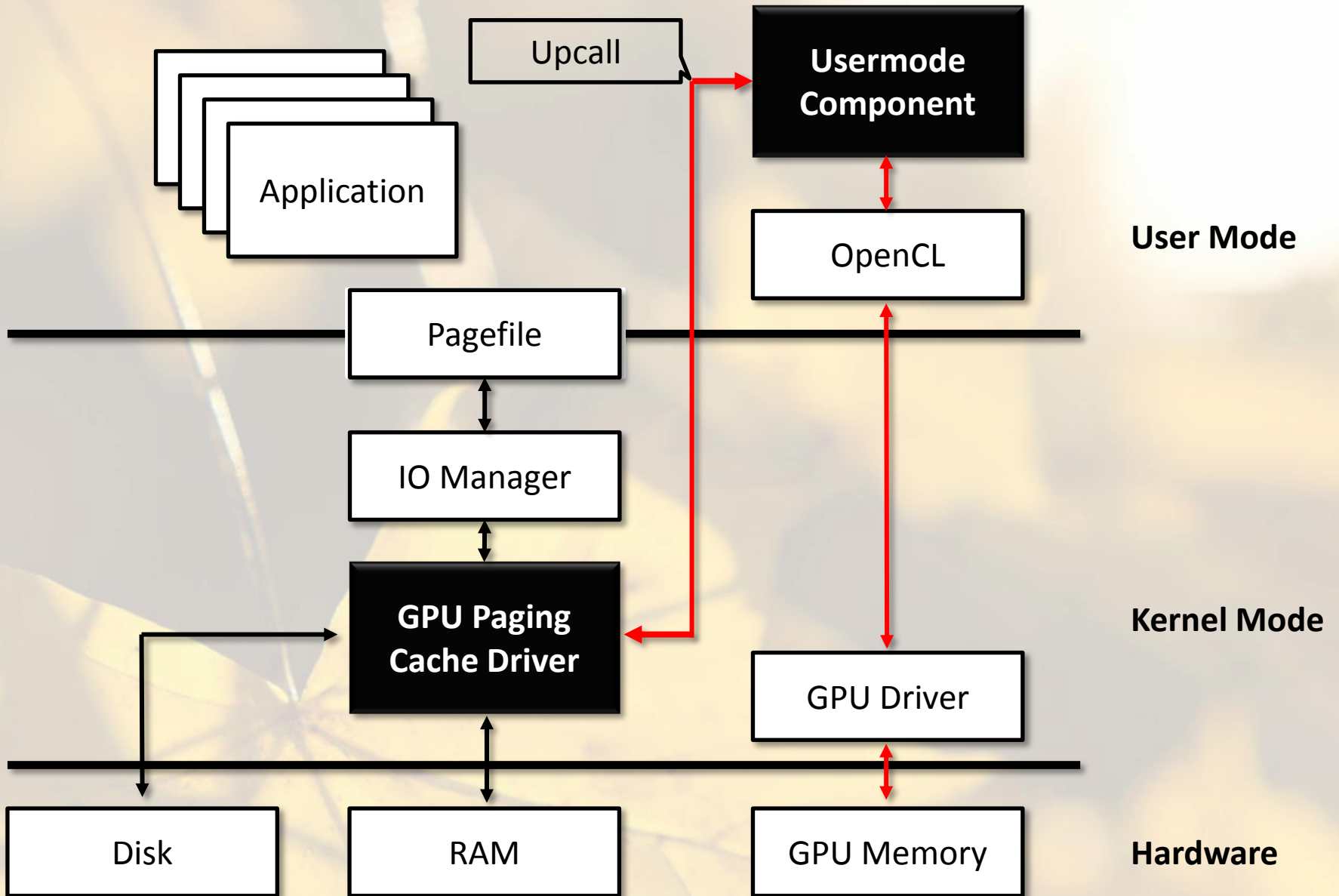
**System RAM Distribution**  
(Steam HW Survey May 2012)

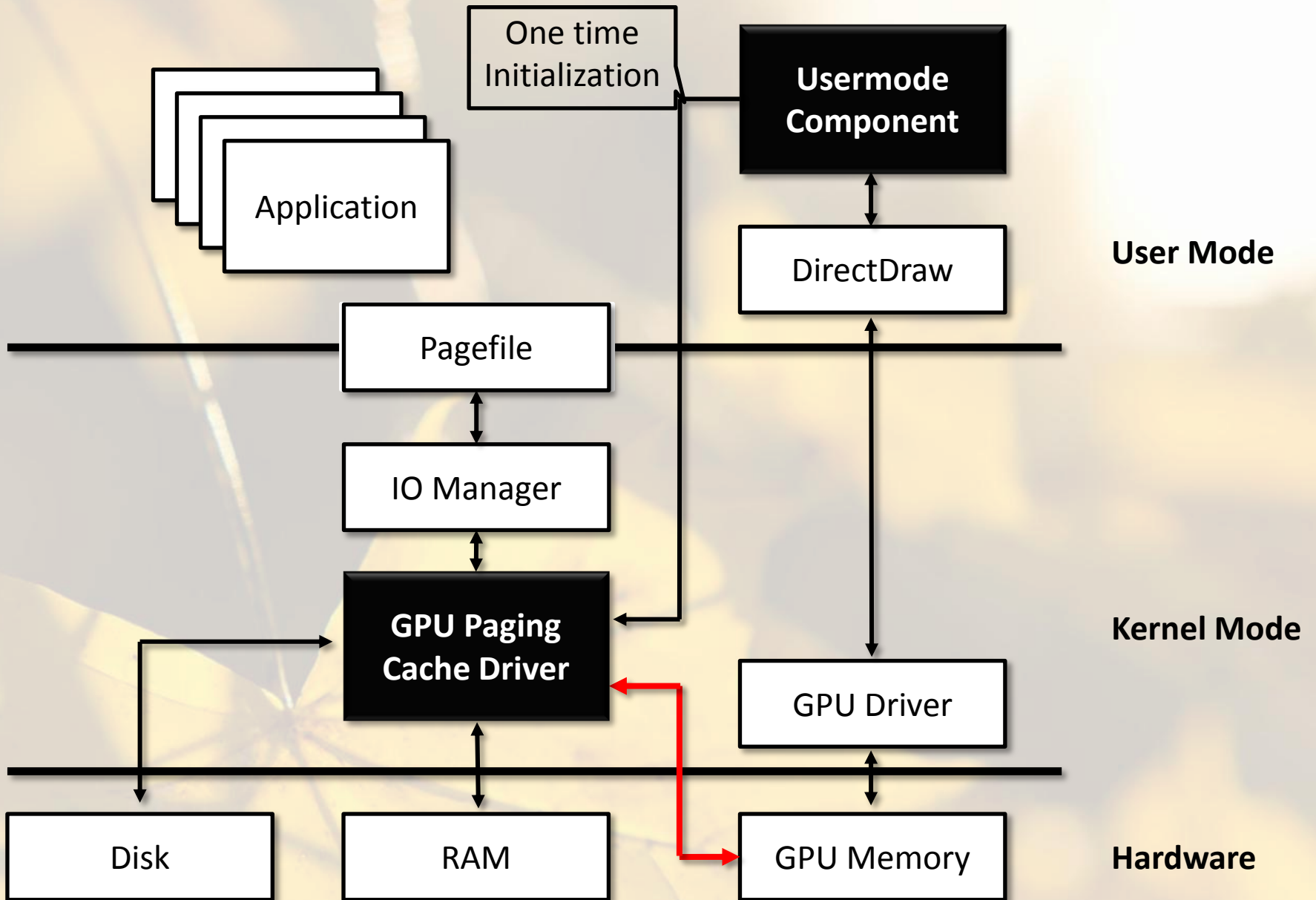


**VRAM Distribution**  
(Steam HW Survey May 2012)



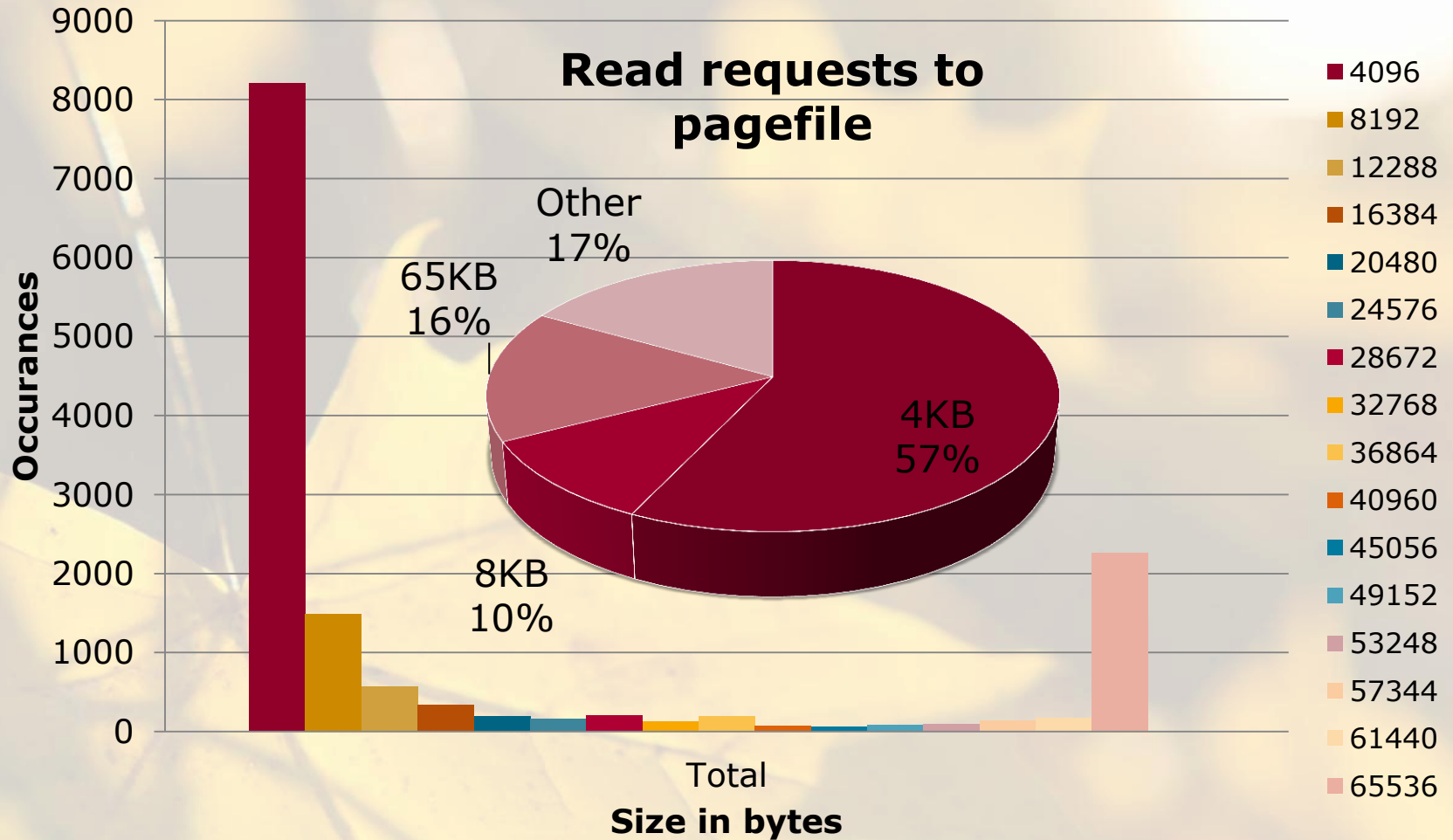






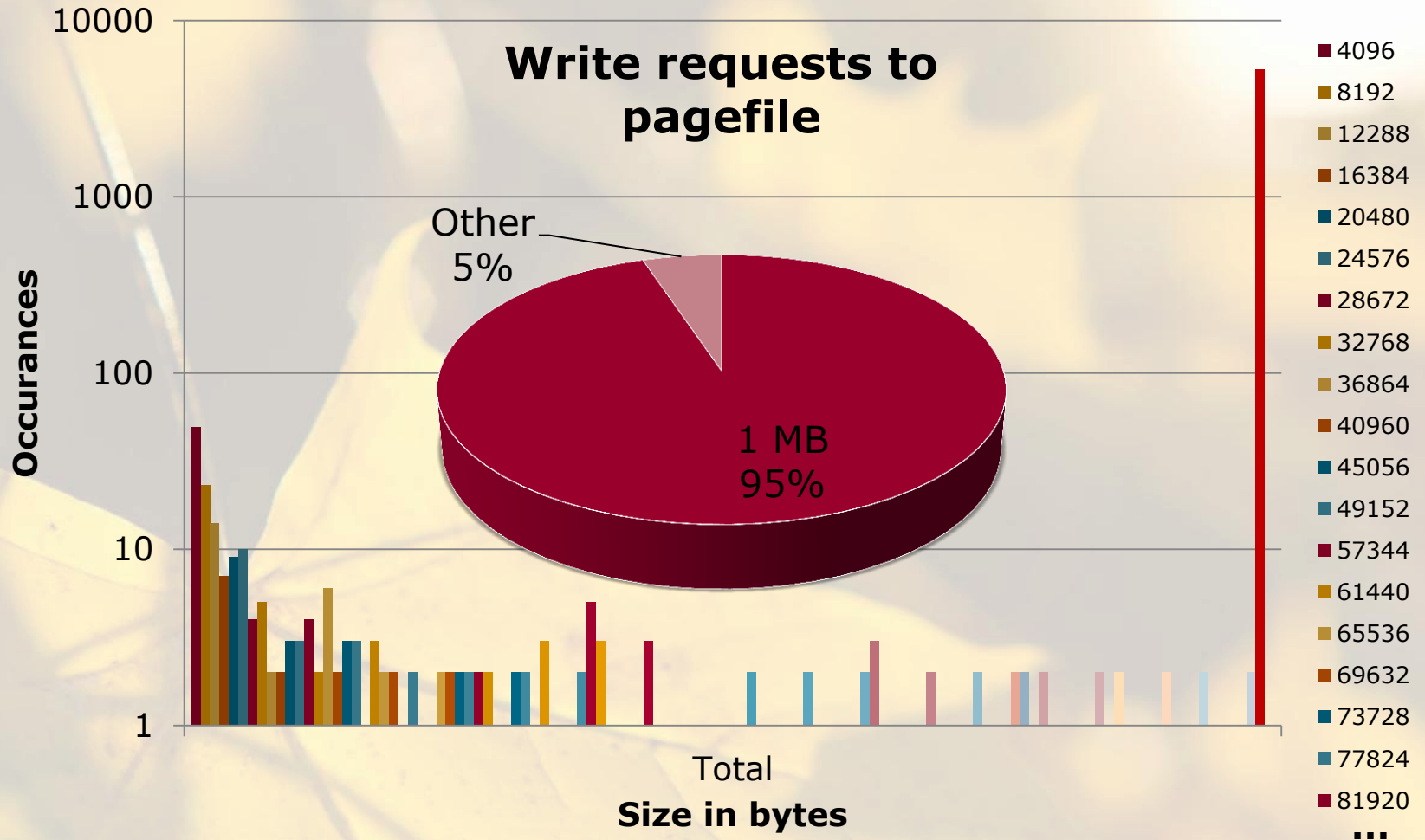
# Paging Request Read

29



# Paging Request Write

30



# Performance Metric

31

$$bandwidth = \frac{1}{RTT}$$

$$RTT = RTT_{Read} + RTT_{Write}$$

$$RTT_{Read} = \sum_{i=4KB}^{1MB} \frac{Probability_{Read}(i)}{Bandwidth_{Read}(i)}$$

$$RTT_{Write} = \sum_{i=4KB}^{1MB} \frac{Probability_{Write}(i)}{Bandwidth_{Write}(i)}$$

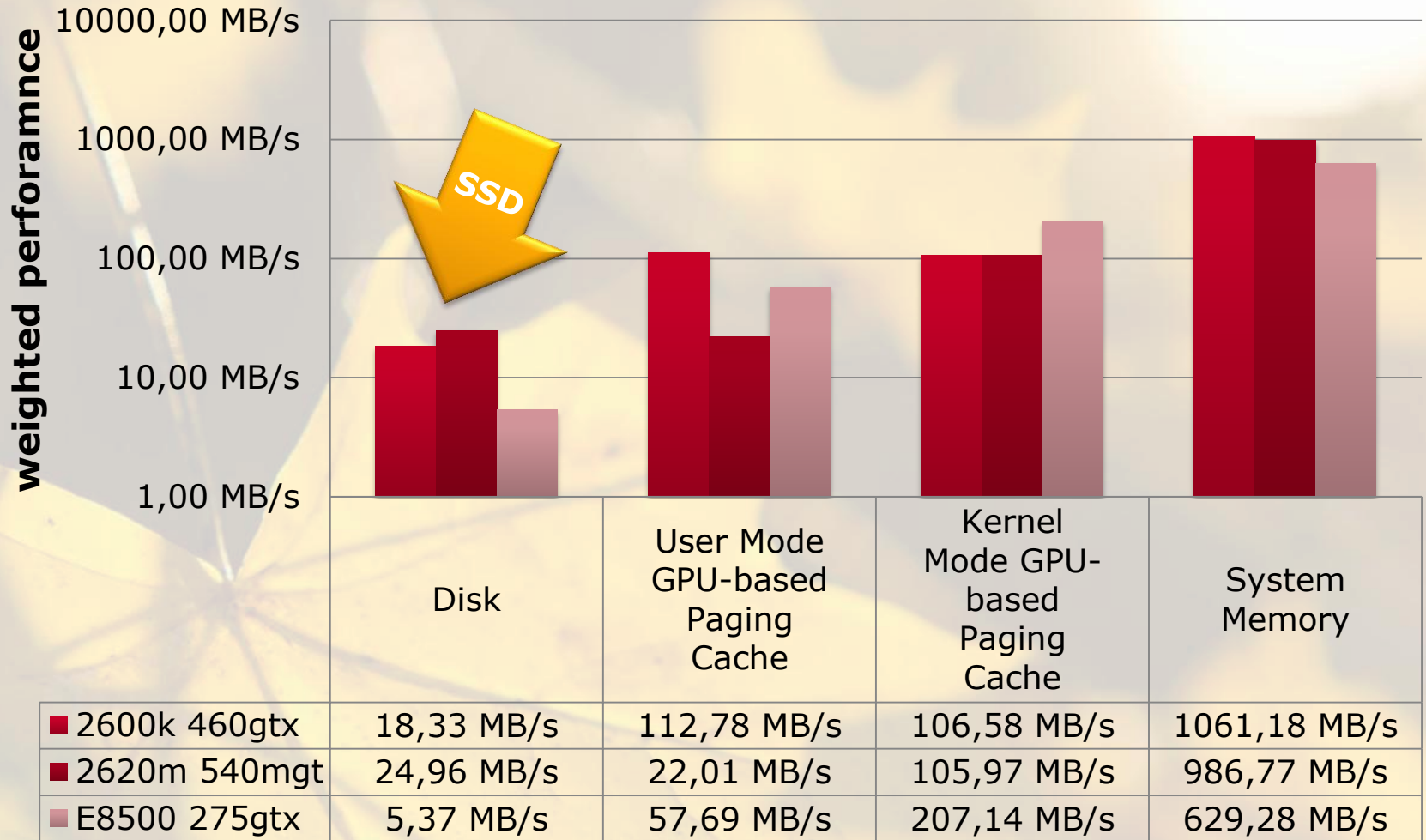
$$bandwidth = \frac{1}{\sum_{i=4KB}^{1MB} \frac{Probability_{Read}(i)}{Bandwidth_{Read}(i)} + \sum_{i=4KB}^{1MB} \frac{Probability_{Write}(i)}{Bandwidth_{Write}(i)}}$$

# Tested Hardware

32

ID	CPU	GPU	Memory	Disk	Operating System
2600k 460gtx	Intel Core i7-2600K 3,7 Ghz	Nvidia GeForce GTX 460 1023 MB GDDR5	8 GB DDR3	Intel Rapid Storage 64 GB SSD	Windows 7 Ultimate SP1 64 bit
E8500 275gtx	Intel Core 2 Duo E8500 3,17 Ghz	Nvidia GeForce GTX 275 896 MB GDDR3	8 GB DDR2	1 TB HDD	Windows 7 Ultimate SP1 64 bit
2620m 540mgt	Intel Core i7-2620M 2,7 Ghz	Nvidia GeForce GT 540M 3 GB GDDR3	8 GB DDR3	240 GB SSD	Windows 7 Ultimate SP1 64 bit





## „Operating Systems must support GPU abstractions“

- GPU as yet another processor
  - Resource abstraction, scheduling, security
- Transparent data transfer and execution on all types of accelerators

