

GPUfs: Integrating a file system with GPUs

Mark Silberstein
(UT Austin/Technion)

Bryan Ford (Yale), Idit Keidar (Technion)
Emmett Witchel (UT Austin)

Traditional System Architecture

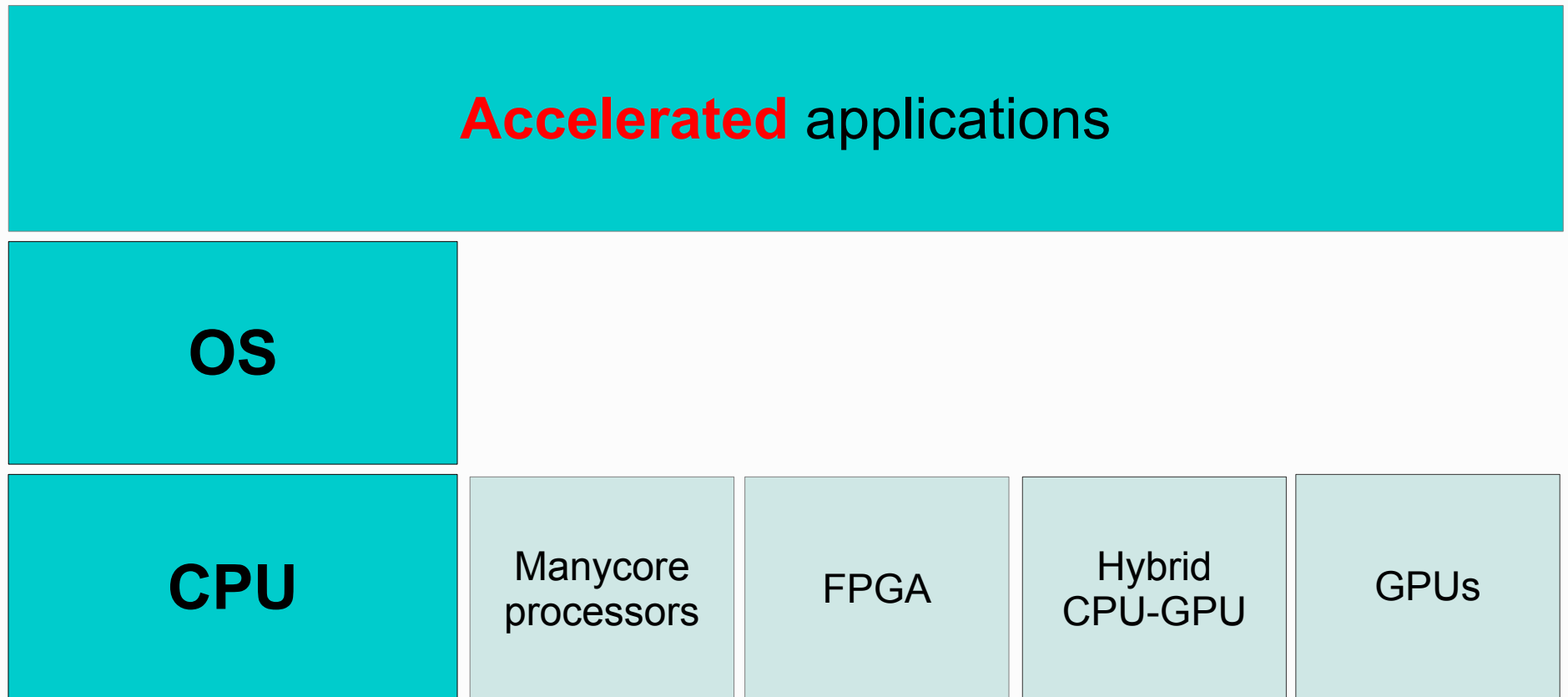


Applications

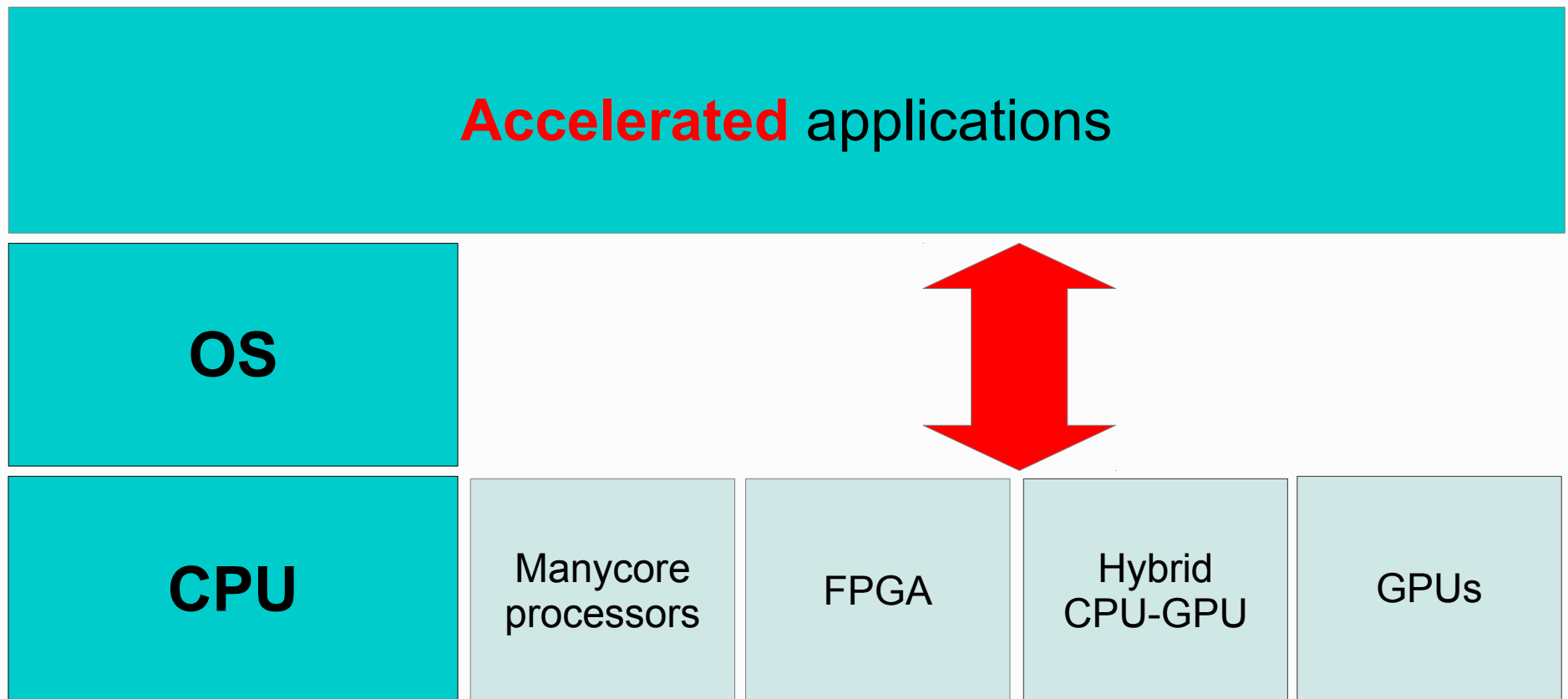
OS

CPU

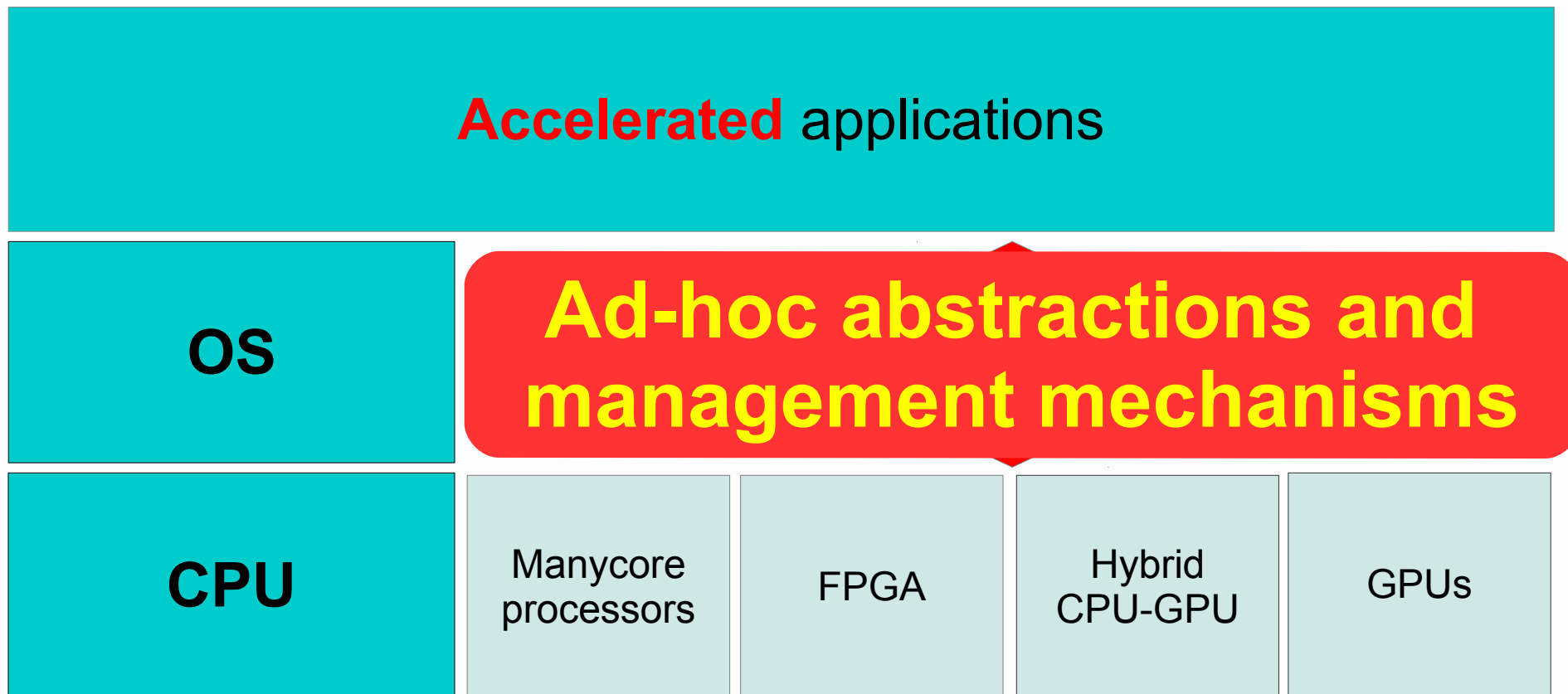
Modern System Architecture



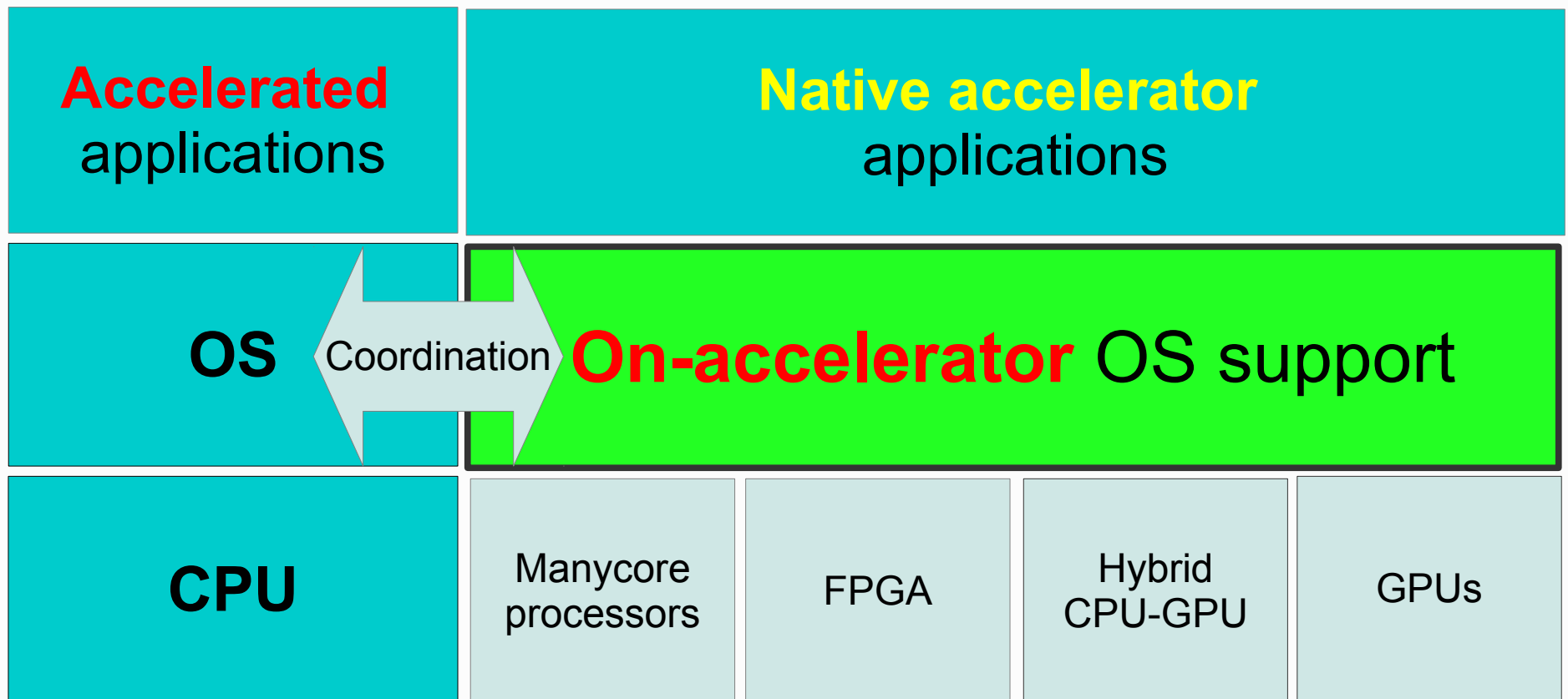
Software-hardware gap is widening



Software-hardware gap is widening



On-accelerator OS support closes the programmability gap



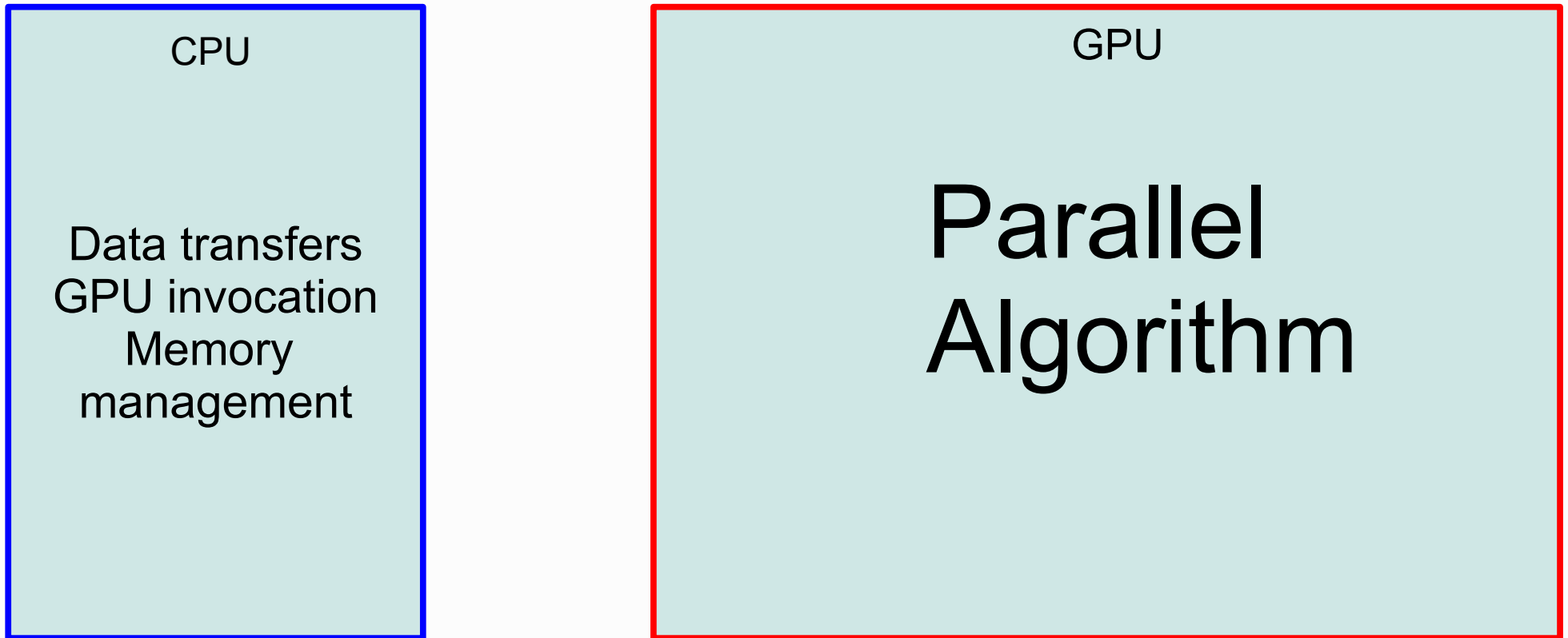
- **GPUfs: File I/O support for GPUs**

- Motivation
- Goals
- Understanding the hardware
- Design
- Implementation
- Evaluation

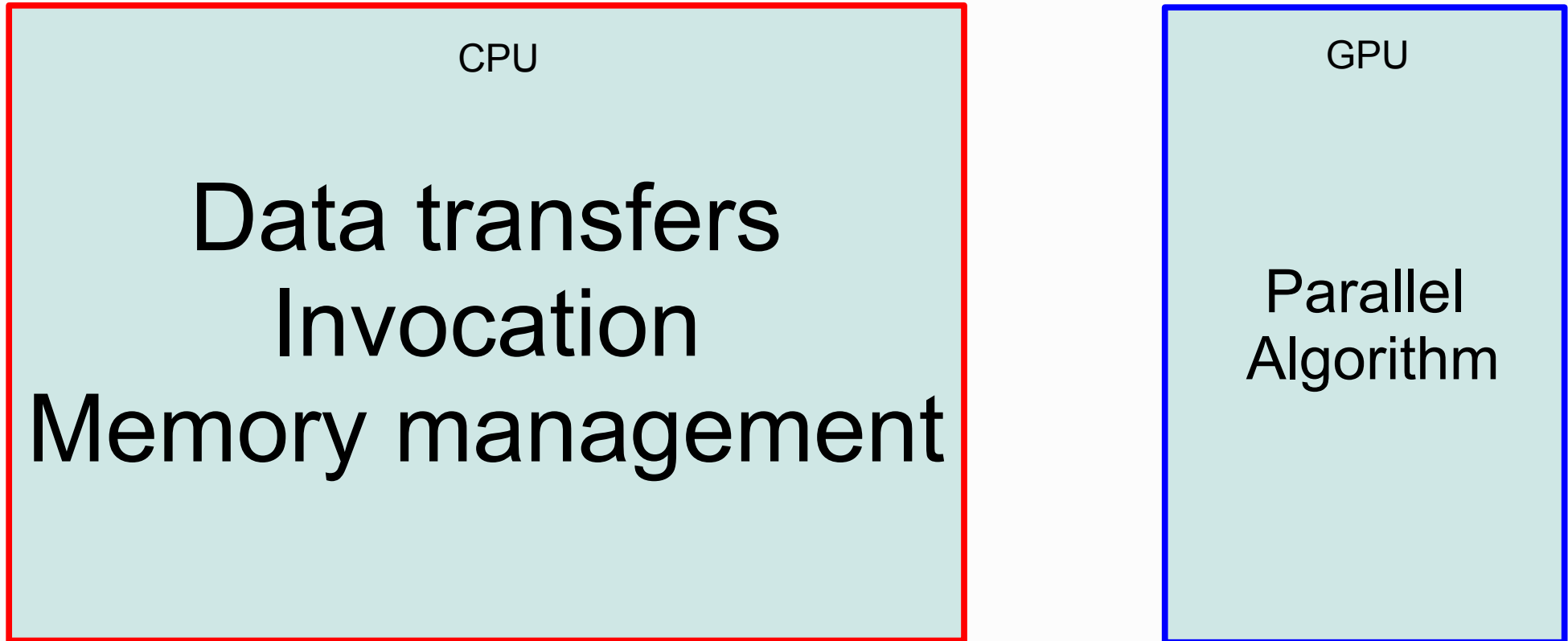


Building systems with GPUs is hard.
Why?

Goal of GPU programming frameworks

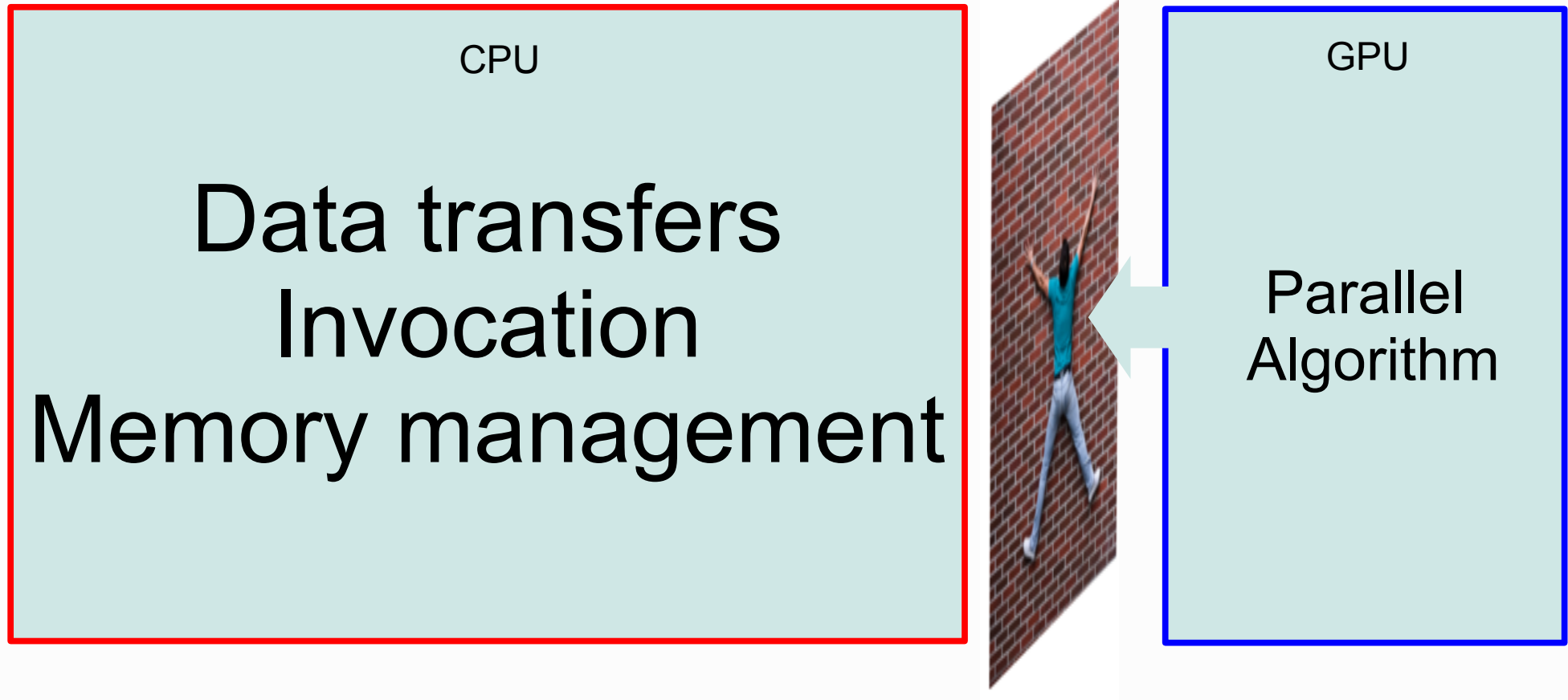


Headache for GPU programmers

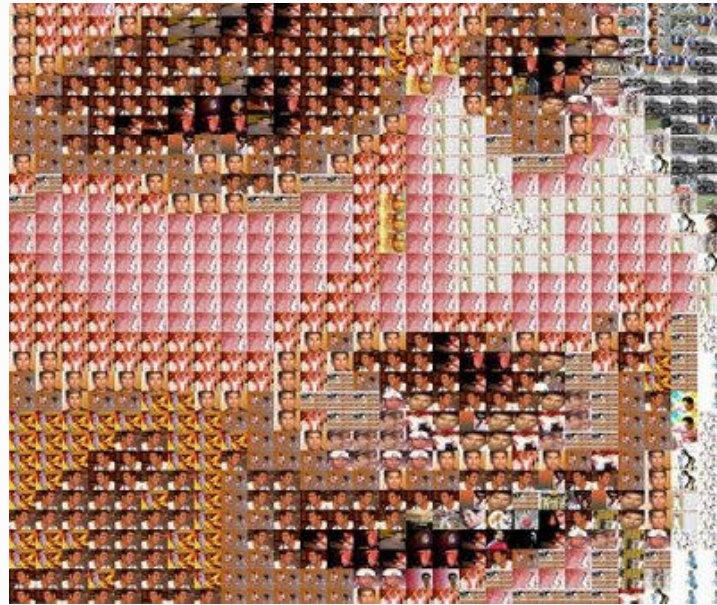


Half of the CUDA SDK 4.1 samples:
at least **9 CPU LOC per 1 GPU LOC**

GPU kernels are isolated



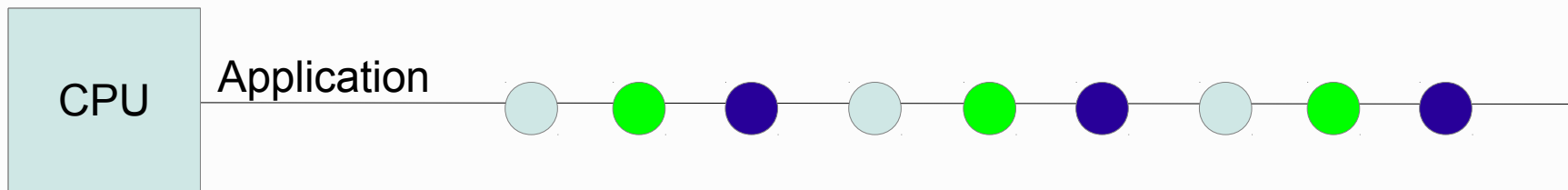
Example: accelerating photo collage



<http://www.codeproject.com/Articles/36347/Face-Collage>

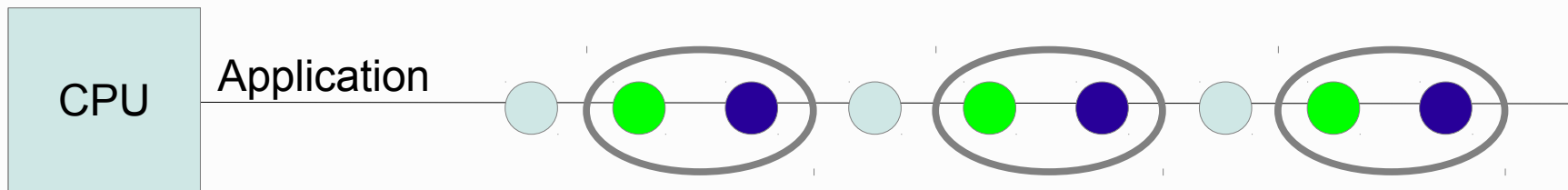
```
While(Unhappy()) {  
    Read_next_image_file()  
    Decide_placement()  
    Remove_outliers()  
}
```

CPU Implementation



```
While(Unhappy()) {  
    Read_next_image_file()  
    Decide_placement()  
    Remove_outliers()  
}
```

Offloading computations to GPU

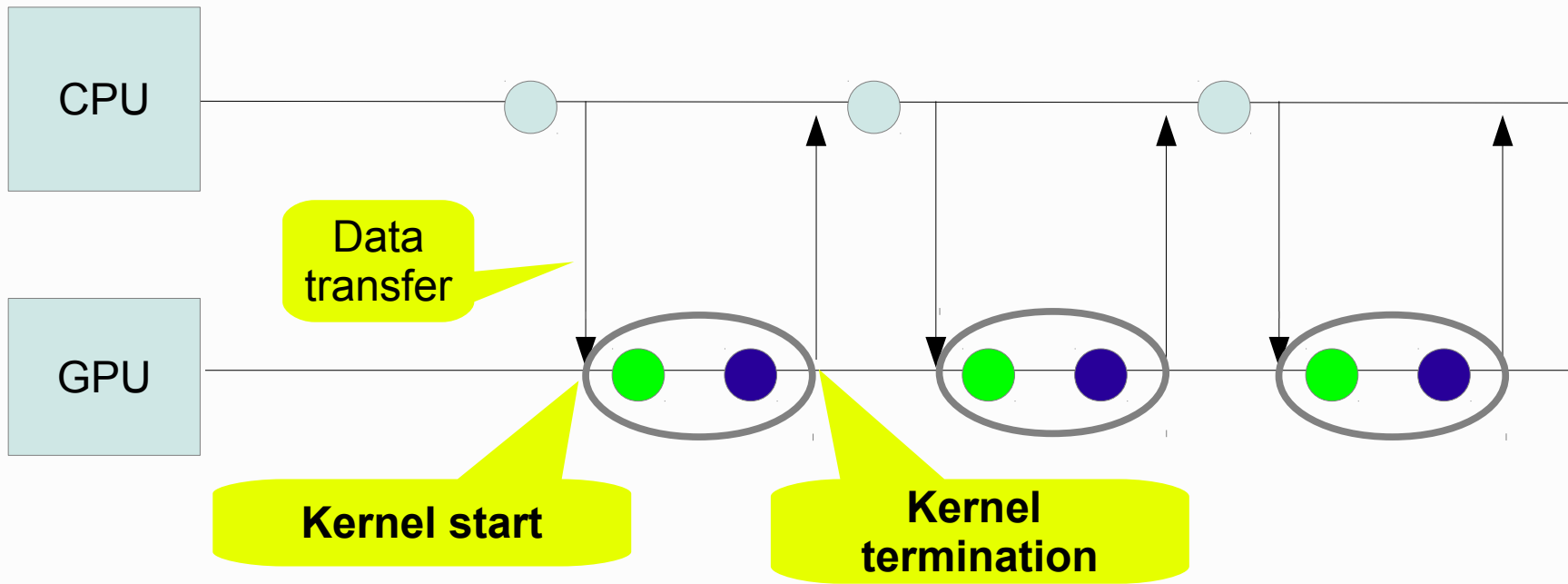


Move to GPU

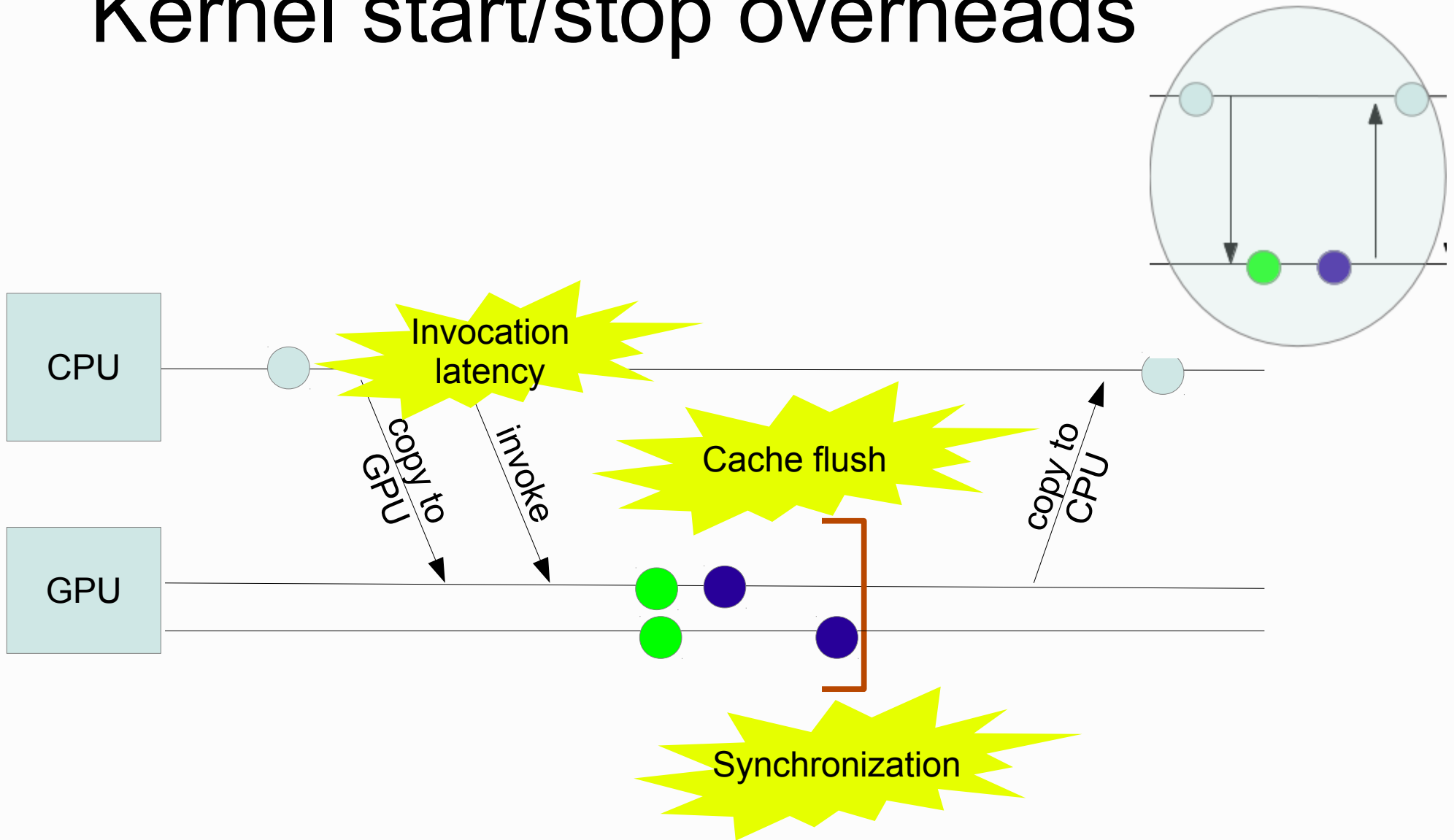
```
While(Unhappy()){  
    Read next image file()  
    Decide_placement()  
    Remove_outliers()  
}
```

Offloading computations to GPU

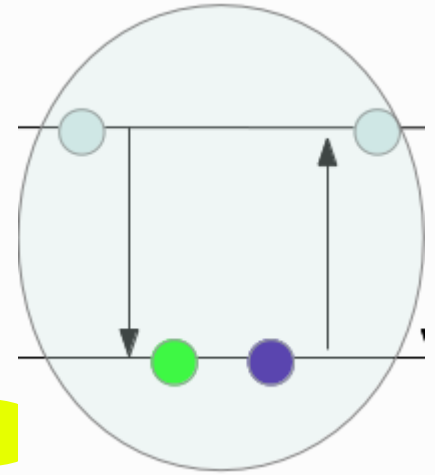
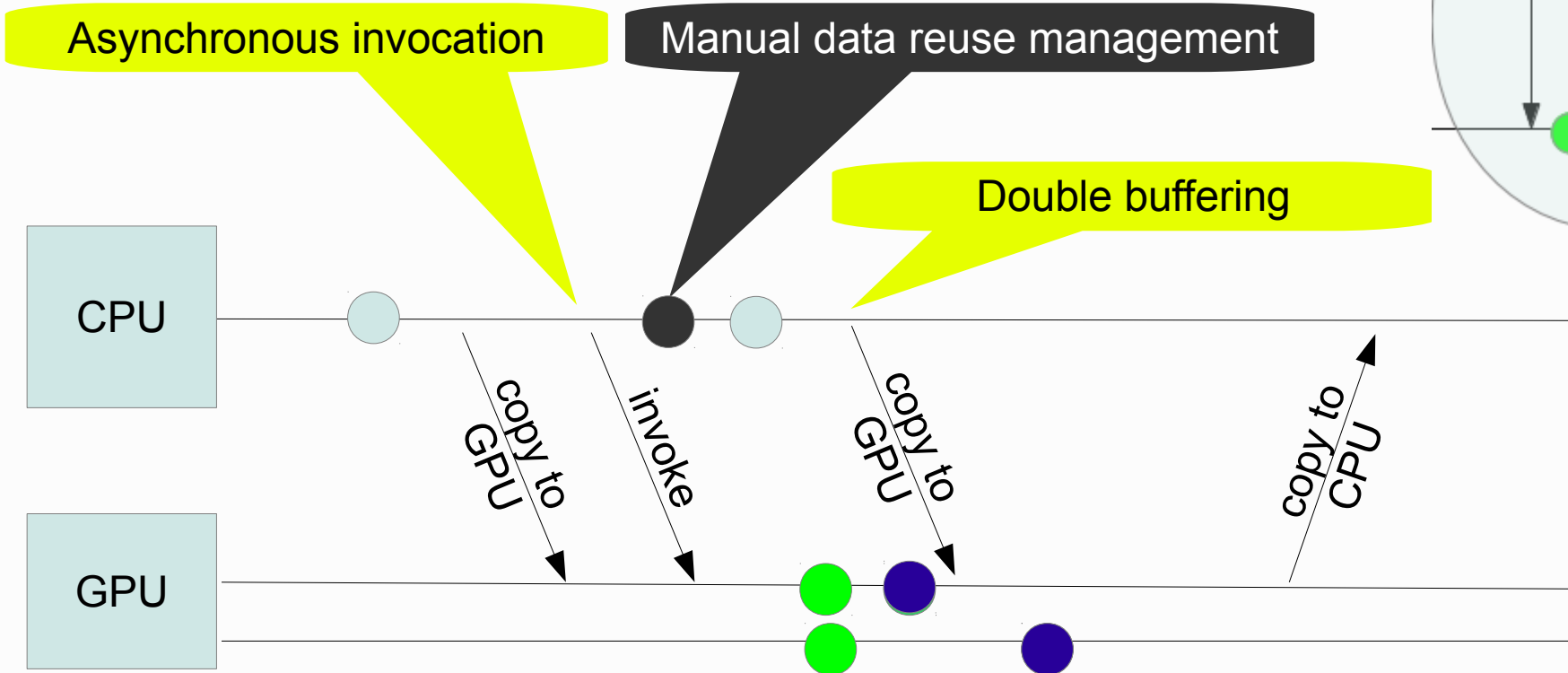
Co-processor programming model



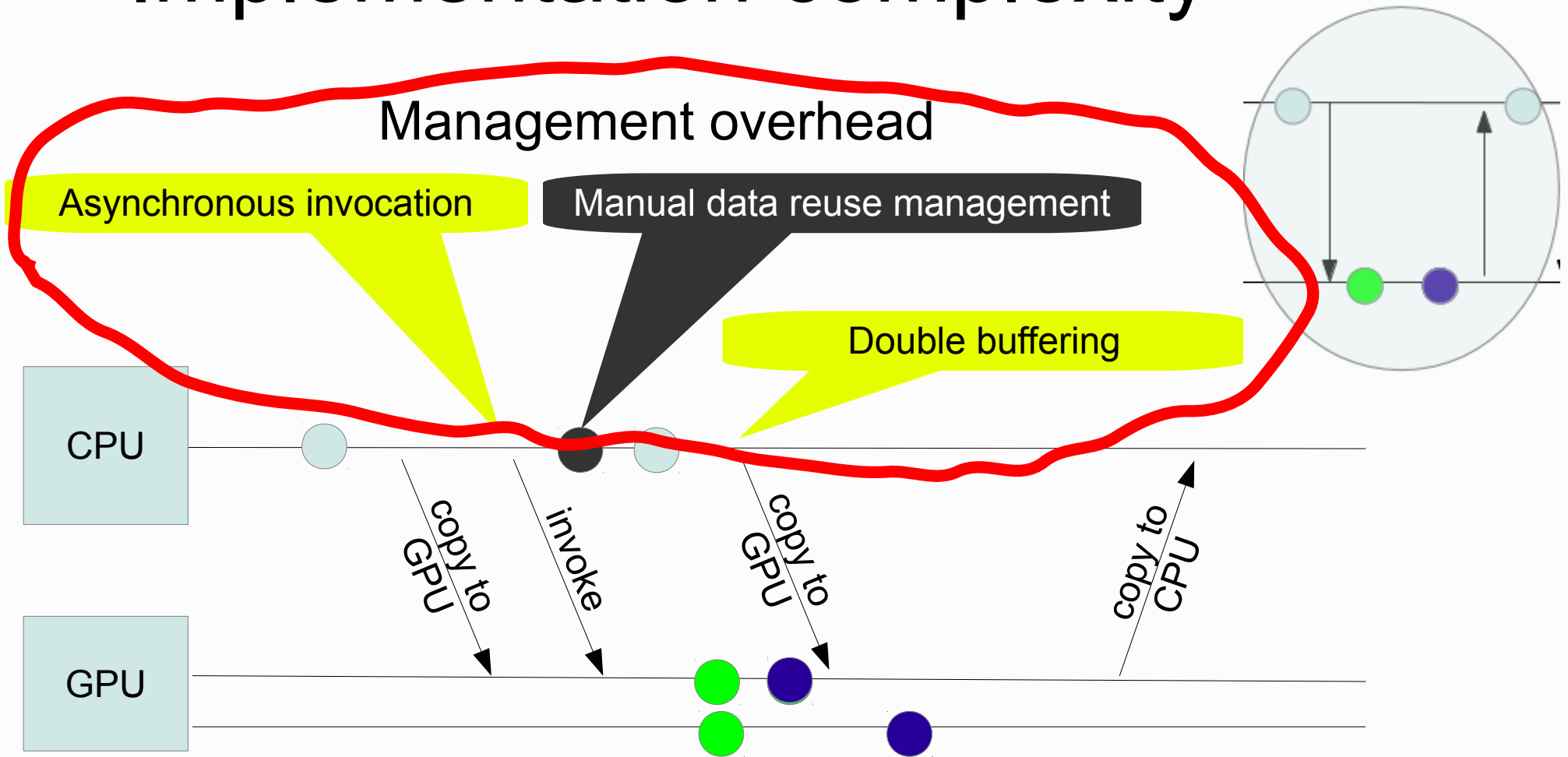
Kernel start/stop overheads



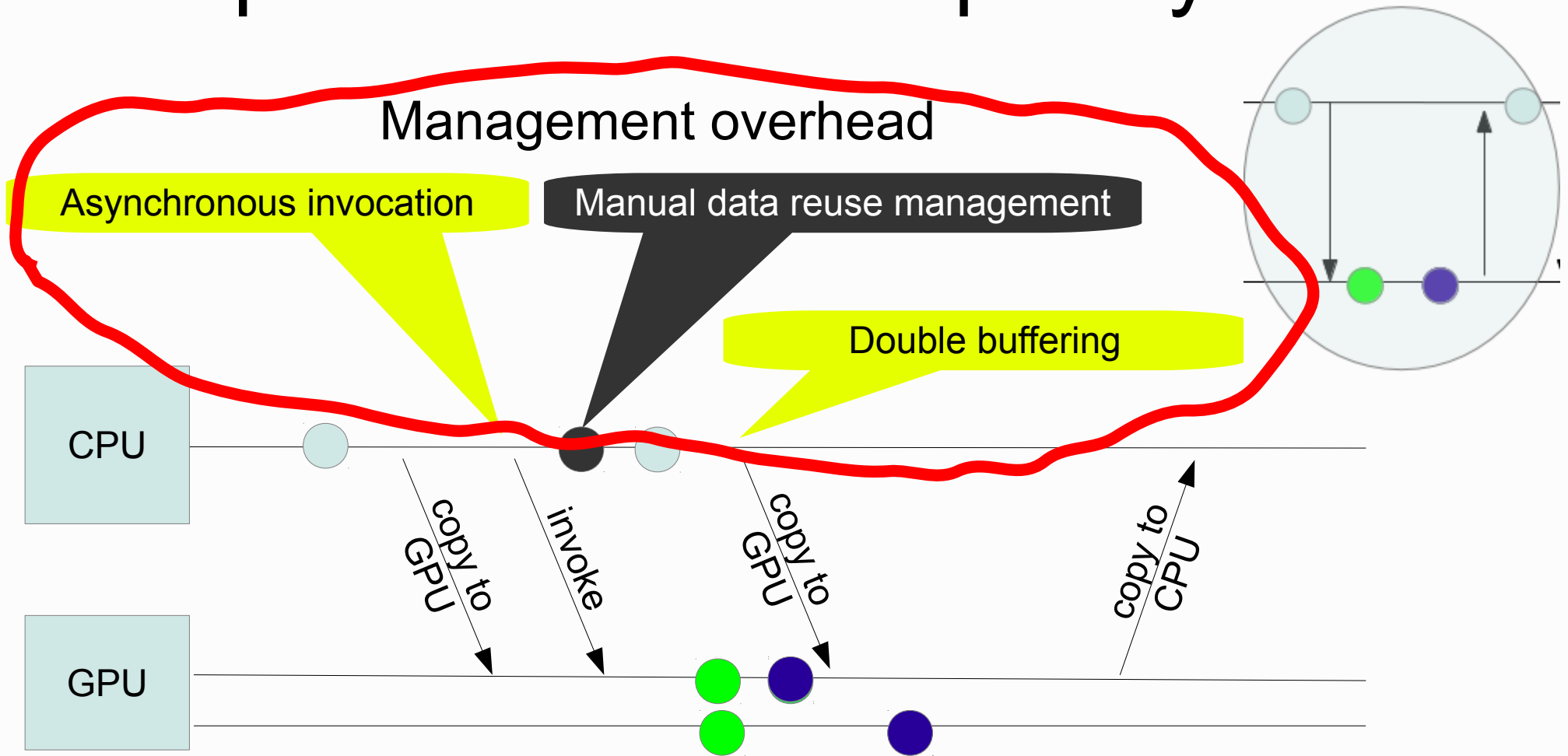
Hiding the overheads



Implementation complexity



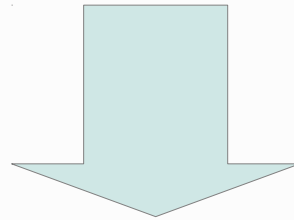
Implementation complexity



Why do we need to deal with low-level system details?

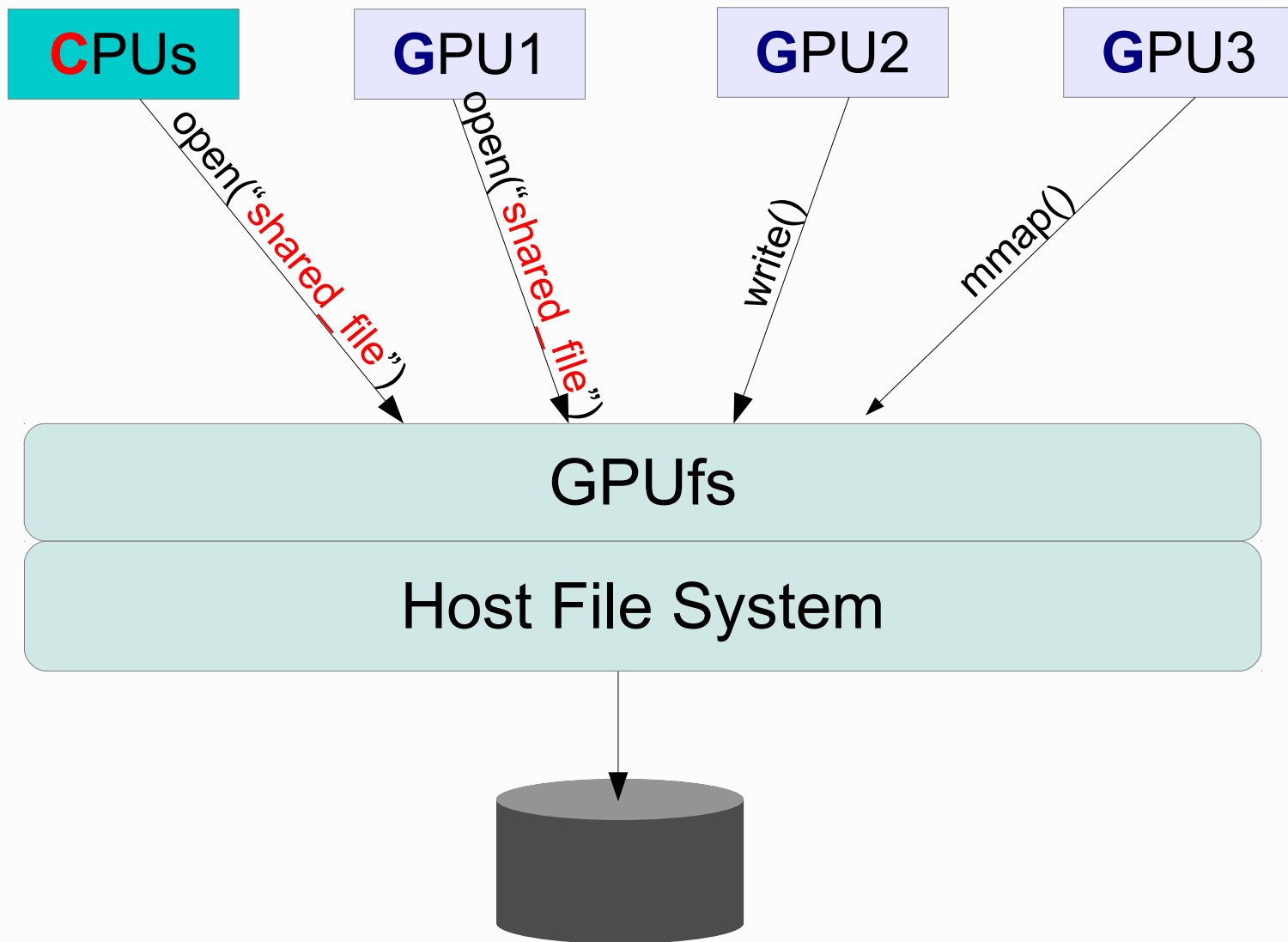
The reason is....

GPUs are peer-processors

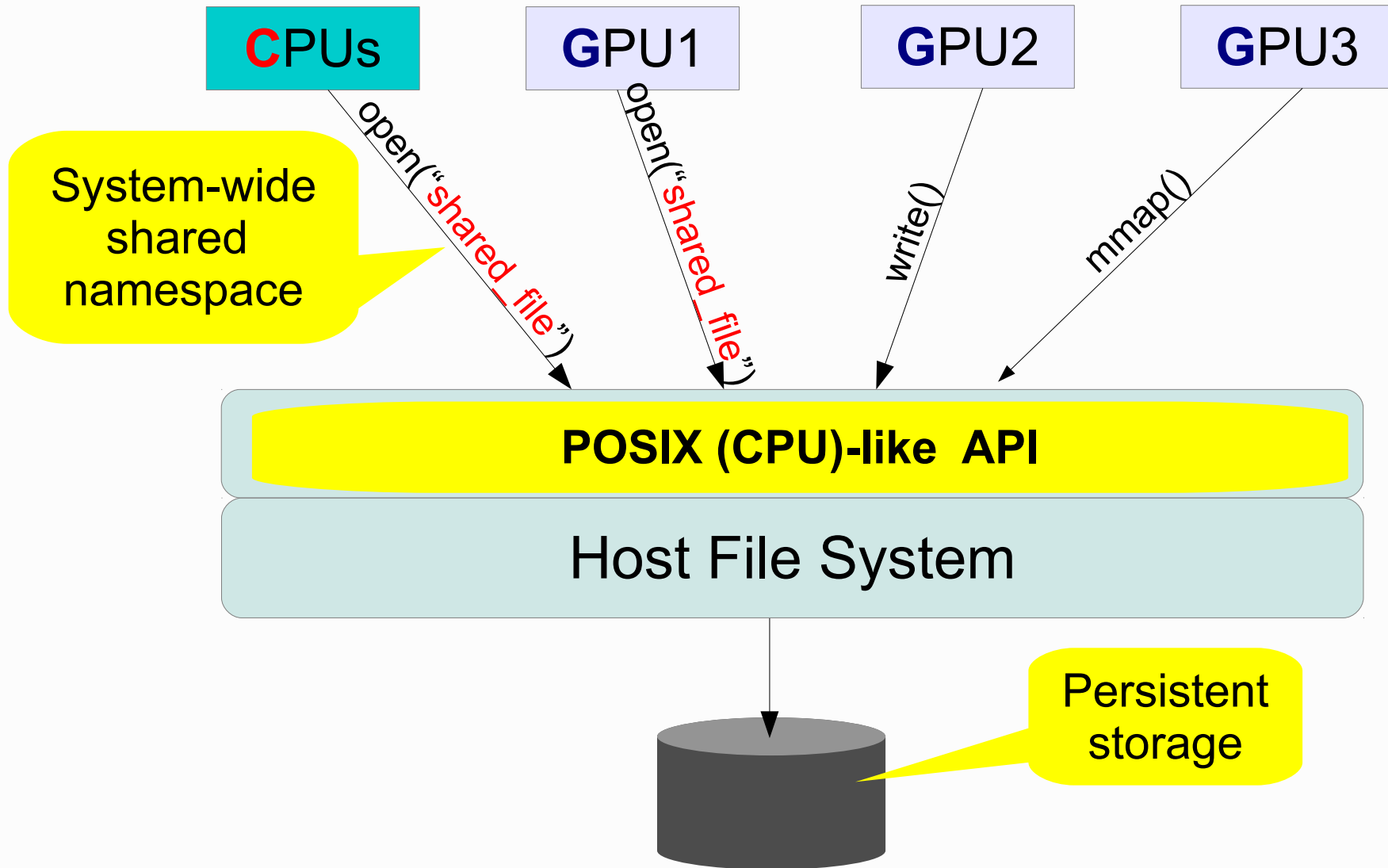


They need I/O OS services

GPUfs: application view

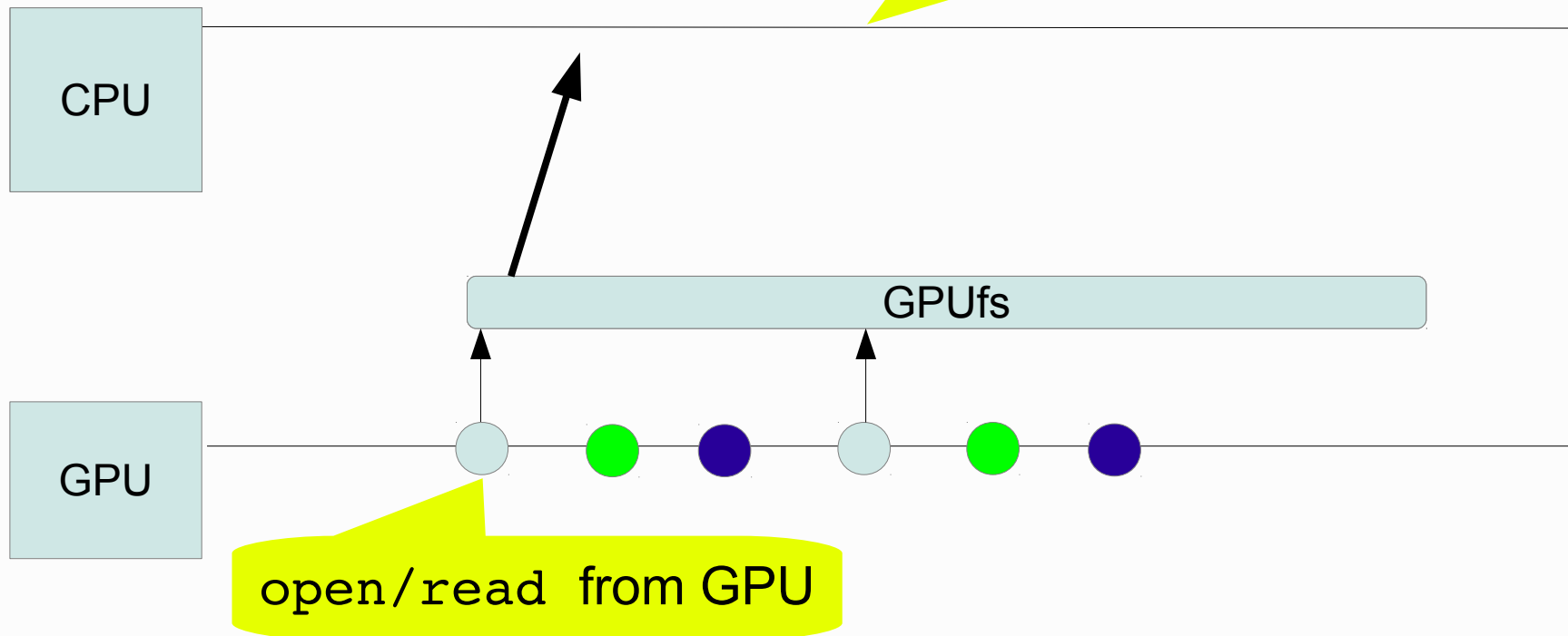


GPUfs: application view

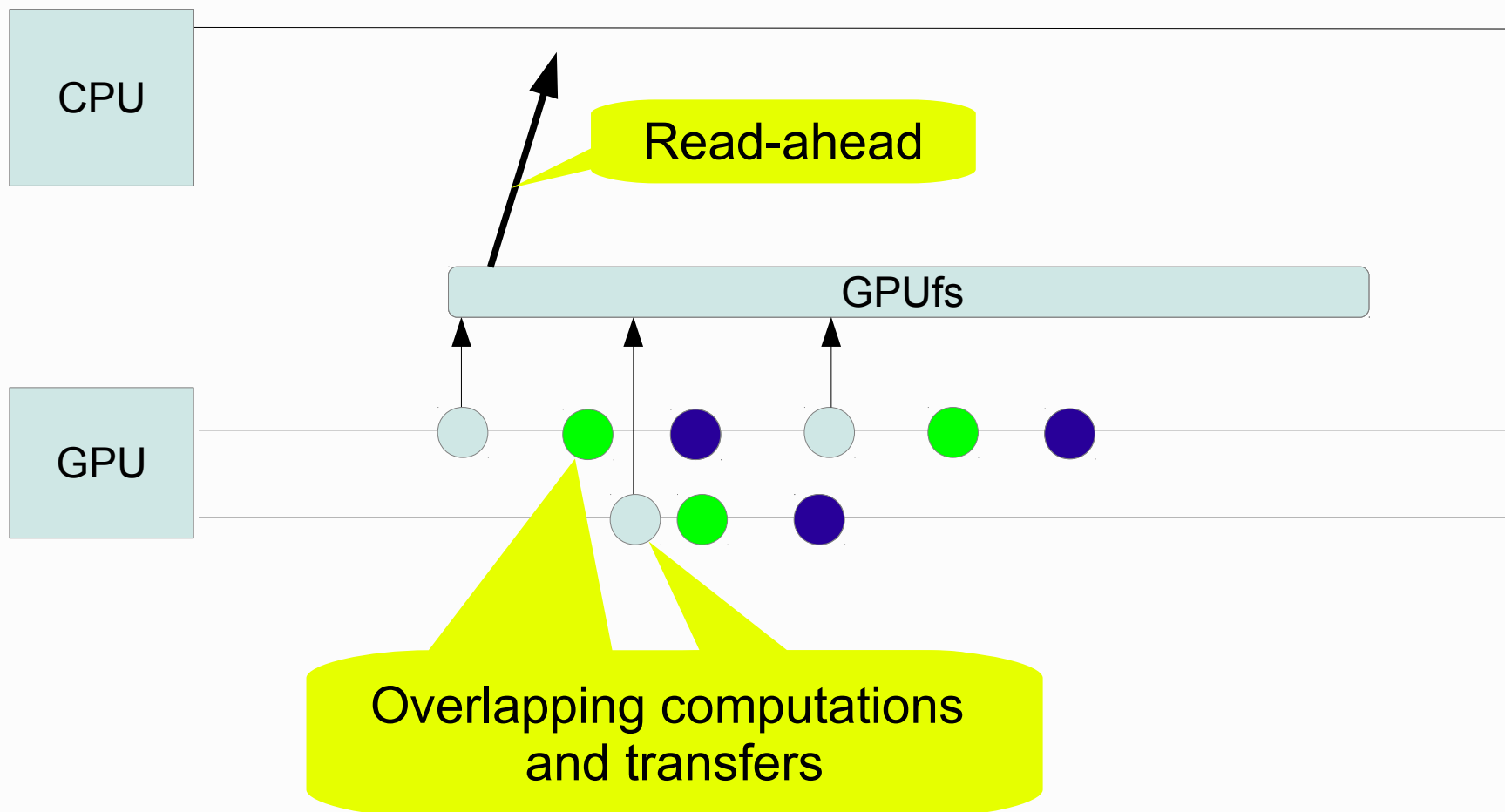


Accelerating collage app with GPUfs

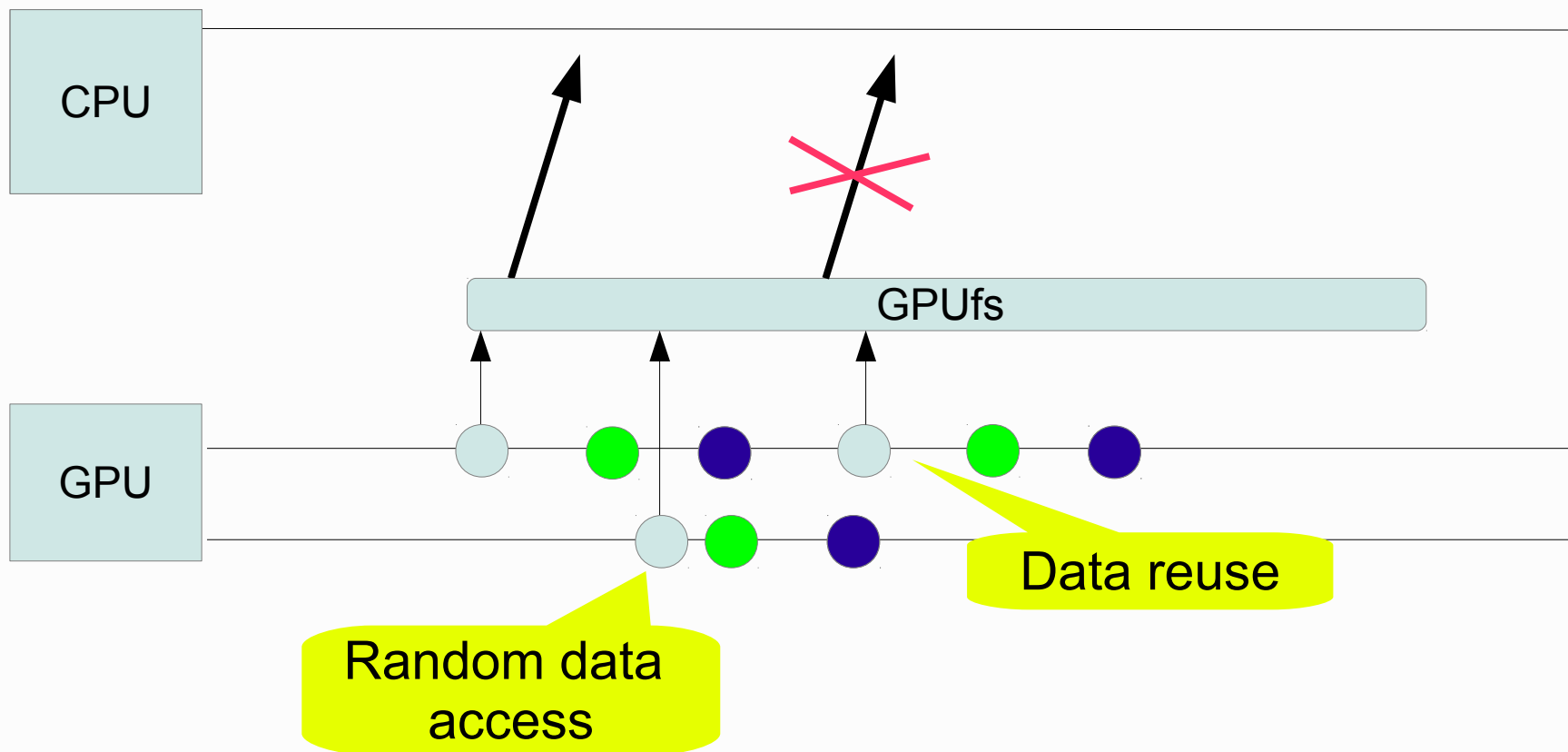
No CPU
management code



Accelerating collage app with GPUfs



Accelerating collage app with GPUfs



Challenge

GPU \neq CPU

Massive parallelism

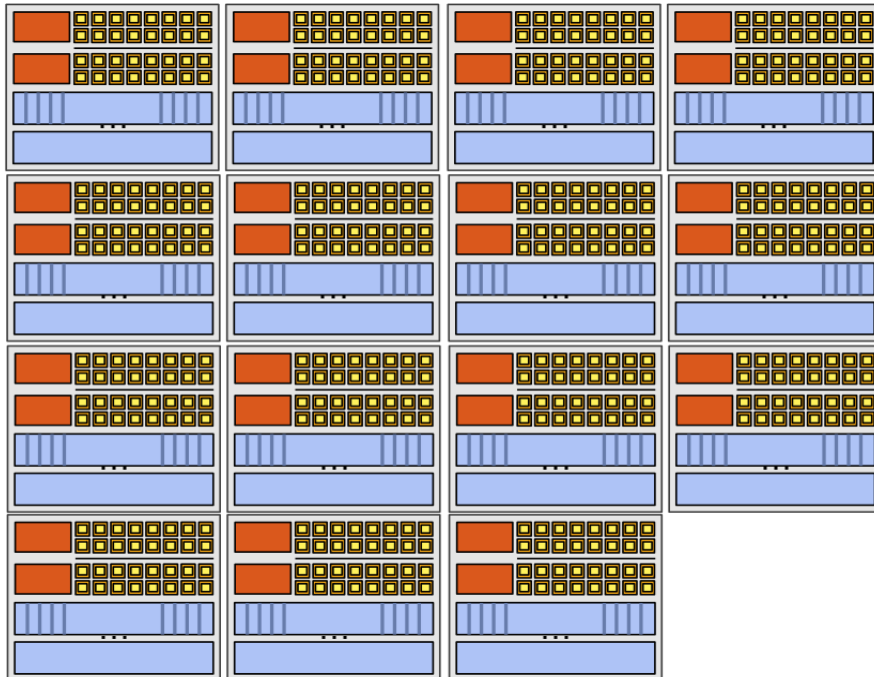
Parallelism is essential for performance in deeply multi-threaded wide-vector hardware



NVIDIA Fermi*
23,000
active threads



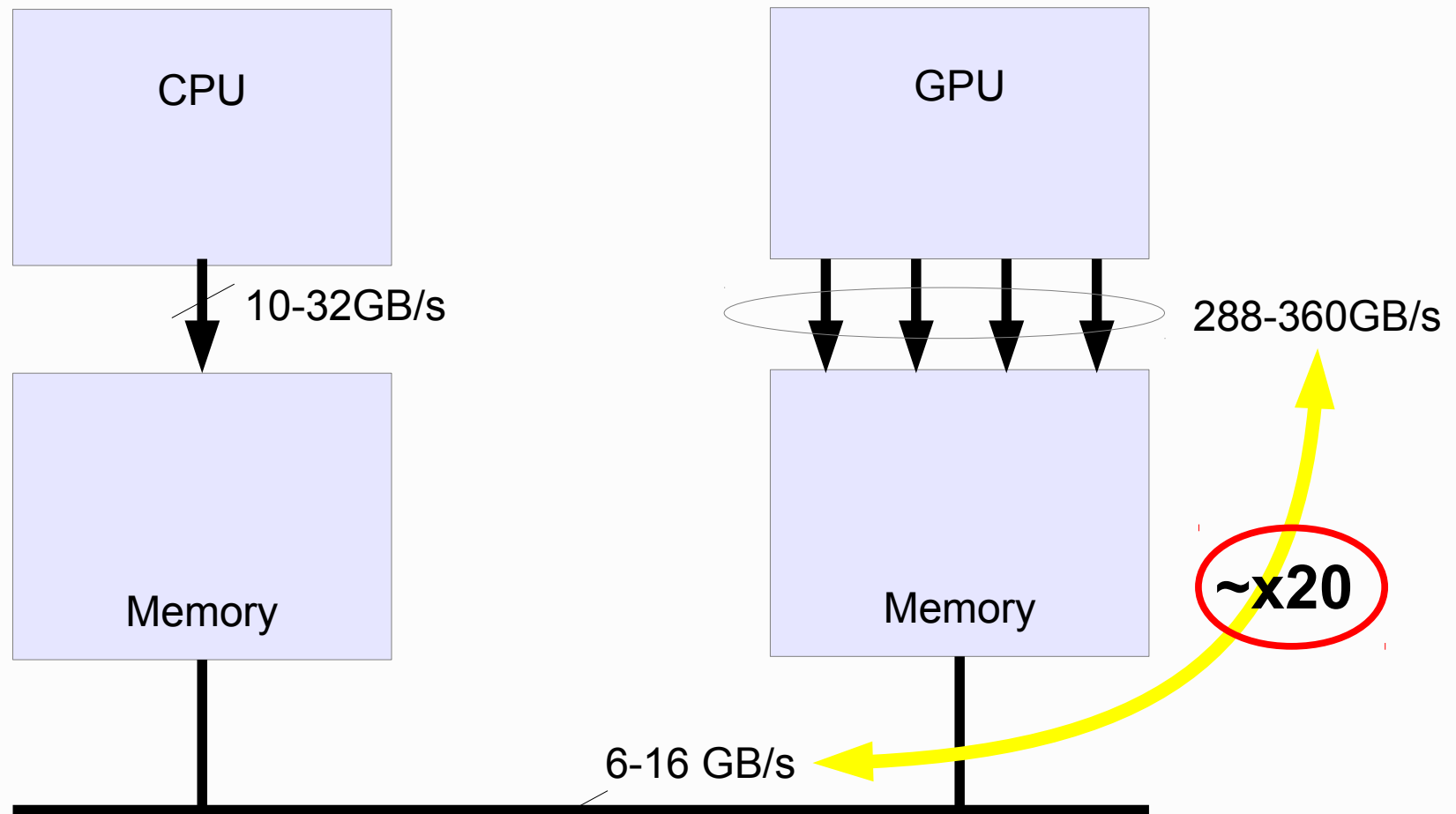
AMD HD5870*
31,000
active threads



From M. Houston/A. Lefohn/K. Fatahalian – A trip through the architecture of modern GPUs*

Heterogeneous memory

GPUs inherently impose high bandwidth demands on memory

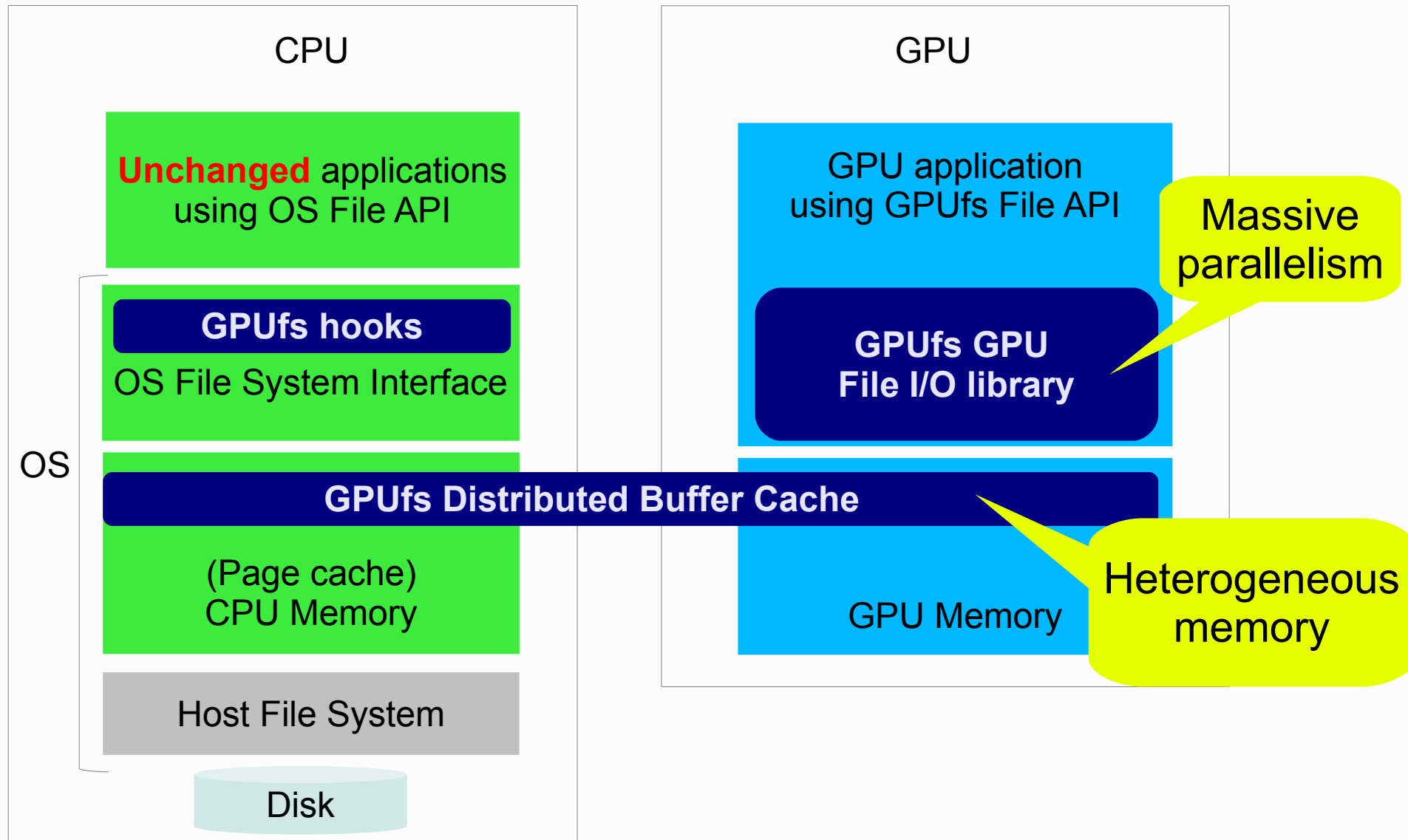


How to build an FS layer
on this hardware?

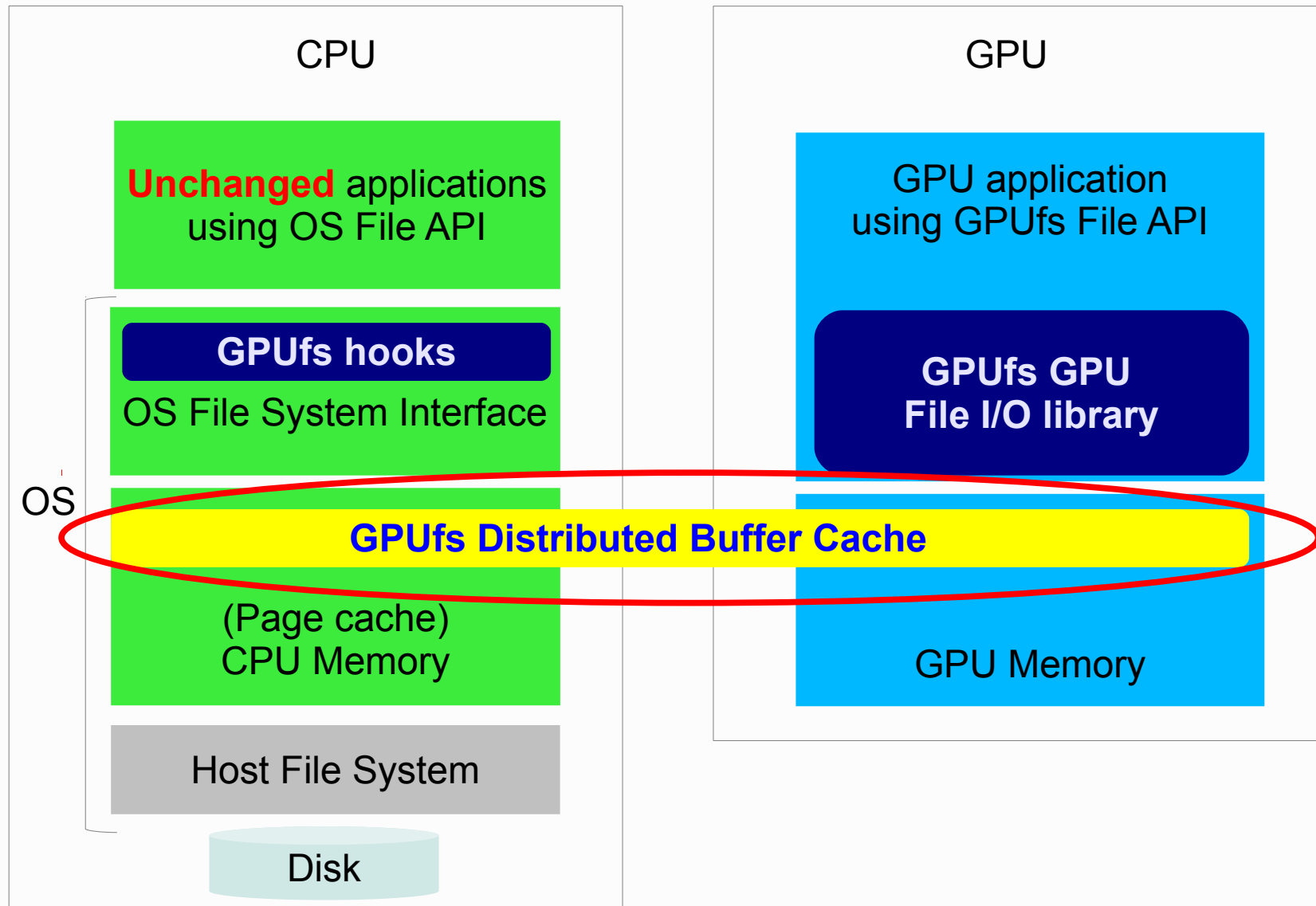
GPUfs: principled **redesign** of the **whole** file system stack

- **Relaxed FS API semantics** for parallelism
- **Relaxed FS consistency** for heterogeneous memory
- **GPU-specific implementation** of synchronization primitives, lock-free data structures, memory allocation,

GPUfs high-level design



GPUfs high-level design



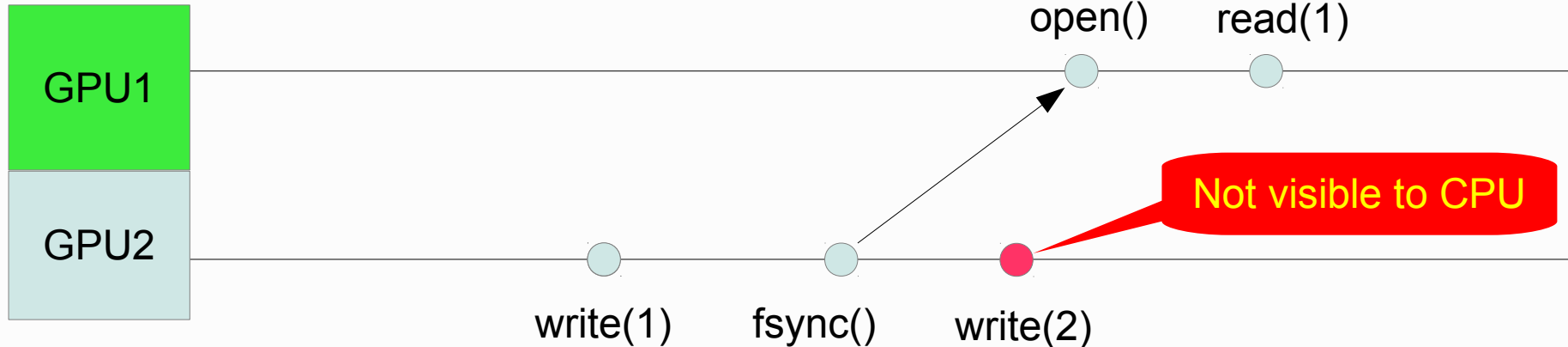
Buffer cache semantics

Local or Distributed file system
data consistency?

GPUfs buffer cache

Weak data consistency model

- close(sync)-to-open semantics (AFS)



Remote-to-Local memory performance ratio is similar to a **distributed system**

On-GPU File I/O API

In the paper

<code>open/close</code>	→	<code>gopen/gclose</code>
<code>read/write</code>	→	<code>gread/gwrite</code>
<code>mmap/munmap</code>	→	<code>gmmap/gmunmap</code>
<code>fsync/msync</code>	→	<code>gfsync/gmsync</code>
<code>ftrunc</code>	→	<code>gftrunc</code>

Changes in the semantics are crucial

Implementation bits

In the paper

- Paging support
- Dynamic data structures and memory allocators
- Lock-free radix tree
- Inter-processor communications (IPC)
- Hybrid H/W-S/W barriers
- Consistency module in the OS kernel

~1,5K GPU LOC, ~600 CPU LOC

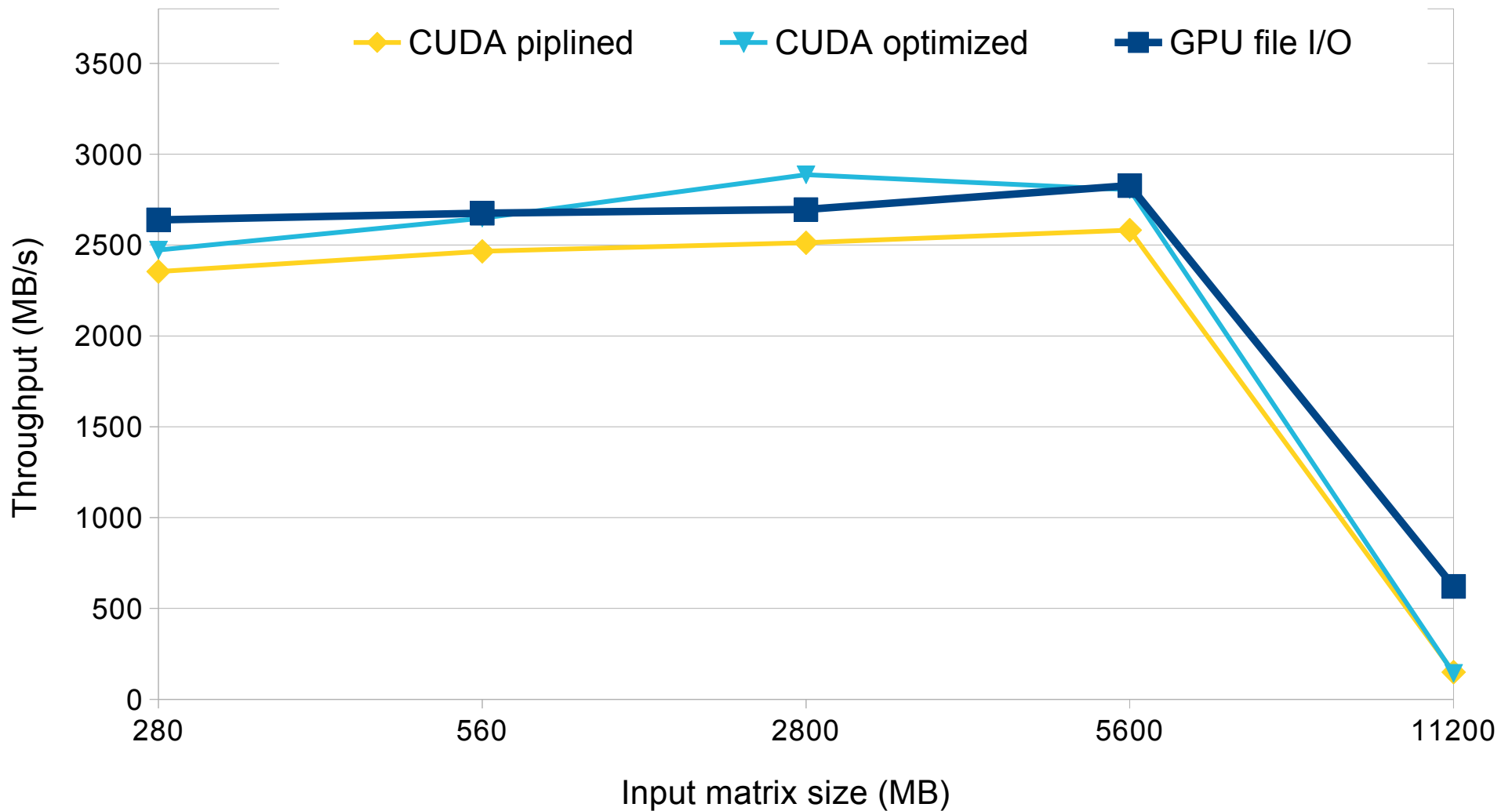
Evaluation

All benchmarks are
written as a GPU
kernel:

**no CPU-side
development**

Matrix-vector product (Inputs/Outputs in files)

Vector 1x128K elements, Page size = 2MB, GPU=TESLA C2075



Word frequency count in text

- Count frequency of modern English words in the works of Shakespeare, and in the Linux kernel source tree

English dictionary: 58,000 words

Challenges

Dynamic working set

Small files

Lots of file I/O (33,000 files, 1-5KB each)

Unpredictable output size

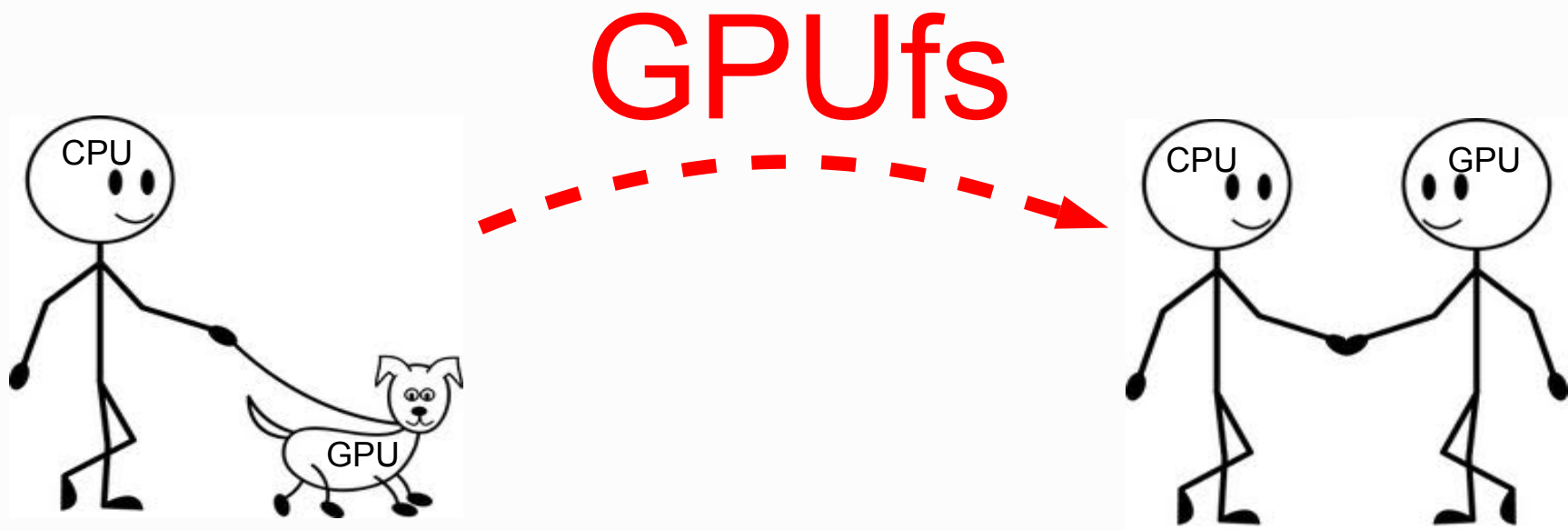
Results

	8CPUs	GPU-vanilla	GPU-GPUfs
Linux source 33,000 files, 524MB	6h	50m (7.2X)	53m (6.8X)
Shakespeare 1 file, 6MB	292s	40s (7.3X)	40s (7.3X)

Results

	8CPUs	GPU-vanilla	GPU-GPUs
Linux source 33,000 files, 524MB	6h	50m (7.2X) 8% overhead	53m (6.8X)
Shakespeare 1 file, 6MB	292s	40s (7.3X)	40s (7.3X)
Unbounded input/output size support	✓	✗	✓

GPUfs is the first system to provide native access to host OS services from GPU programs



Code is available for download at:

<https://sites.google.com/site/silbersteinmark/Home/gpufs>

<http://goo.gl/ofJ6J>