# HSA Foundation

## Antonio Miele

Politecnico di Milano

POLITECNICO
MILANO 1863

# References

- This presentation is based on the material and slides published on the HSA foundation website:
  - http://www.hsafoundation.com/

# Heterogeneous processors have proliferated – make them better

- Heterogeneous SoCs have arrived and are a tremendous advance over previous platforms
- SoCs combine CPU cores, GPU cores and other accelerators, with high bandwidth access to memory
- How do we make them even better?
  - Easier to program
  - Easier to optimize
  - Easier to load balance
  - Higher performance
  - Lower power
- HSA unites accelerators architecturally
- Early focus on the GPU compute accelerator, but HSA will go well beyond the GPU

# HSA foundation

- Founded in June 2012
- Developing a new platform for heterogeneous systems
- www.hsafoundation.com
- Specifications under development in working groups to define the platform
- Membership consists of 43 companies and 16 universities
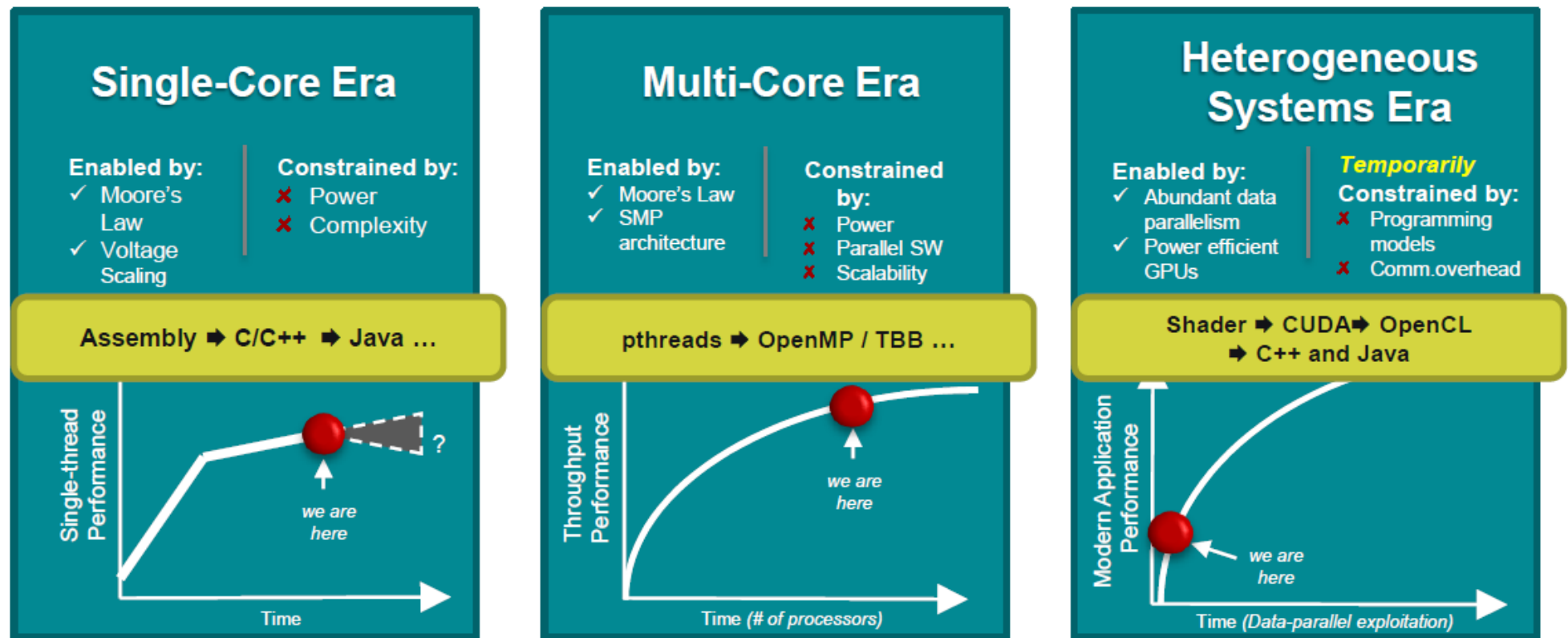- Adding 1-2 new members each month

# HSA consortium

# HSA goals

- To enable power-efficient performance
- To improve programmability of heterogeneous processors
- To increase the portability of code across processors and platforms
- To increase the pervasiveness of heterogeneous solutions throughout the industry

# Paradigm shift

- Inflection in processor design and programming

# Key features of HSA

- **hUMA** – Heterogeneous Unified Memory Architecture
- **hQ** – Heterogeneous Queuing
- **HSAIL** – HSA Intermediate Language

# Key features of HSA

- <u>hUMA – Heterogeneous Unified Memory Architecture</u>

- hQ – Heterogeneous Queuing

- HSAIL – HSA Intermediate Language

# Legacy GPU compute

- Multiple memory pools
- Multiple address spaces
    - No pointer-based data structures
- Explicit data copying across PCIe
    - High latency
    - Low bandwidth
- High overhead dispatch

- Need lots of compute on GPU to amortize copy overhead
- Very limited GPU memory capacity
- Dual source development
- Proprietary environments
- Expert programmers only

# Existing APUs and SoCs

**Physical Integration**



APU = Accelerated Processing Unit (i.e., a SoC containing also a GPU)

- Physical integration of GPUs and CPUs
- Data copies on an internal bus
- Two memory pools remain
- Still queue through the OS
- Still requires expert programmers

- FPGAs and DSPs have the same issues

# Existing APUs and SoCs

- CPU and GPU still have separate memories for the programmer (different virtual memory spaces)
    1. CPU explicitly copies data to GPU memory
    2. GPU executes computation
    3. CPU explicitly copies results back to its own memory

# An HSA enabled SoC



- Unified Coherent Memory enables data sharing across all processors
  - Enabling the usage of pointers
  - Not explicit data transfer -> values move on demand
  - Pageable virtual addresses for GPUs -> no GPU capacity constraints
- Processors architected to operate cooperatively
- Designed to enable the application to run on different processors at different times

# Unified coherent memory



**Coherent Memory:**

Ensures CPU and GPU caches both see an up-to-date view of data

C P U

Cache

HW Coherency

GPU

Cache

Physical Memory

**Pageable memory:**

The GPU can seamlessly access virtual memory addresses that are not (yet) present in physical memory

Virtual Memory

**Entire memory space:**
Both CPU and GPU can access and allocate any location in the system's virtual memory space

# Unified coherent memory

- CPU and GPU have a unified virtual memory spaces
    1. CPU simply passes a pointer to GPU
    2. GPU executes computation
    3. CPU can read the results directly – no explicit copy need!



CPU / GPU Uniform Memory

# Unified coherent memory

# Unified coherent memory

# Unified coherent memory

# Unified coherent memory

# Unified coherent memory

# Unified coherent memory

# Unified coherent memory

# Unified coherent memory

# Unified coherent memory



**DATA POINTERS**

**HSA**

**AMD**

**SYSTEM MEMORY**

**GPU**

**KERNEL**

L   R

L   R   L   R

TREE

RESULT
BUFFER

# Unified coherent memory

# Unified coherent memory

# Unified coherent memory



**DATA POINTERS**

**HSA**

**SYSTEM MEMORY**

**AMD**

**GPU**

**KERNEL**

L  R

L  R  L  R

TREE

RESULT
BUFFER

# Unified coherent memory

# Unified coherent memory

# Key features of HSA

- **hUMA** – Heterogeneous Unified Memory Architecture
- **hQ** – Heterogeneous Queuing
- **HSAIL** – HSA Intermediate Language

# hQ: heterogeneous queuing

- Task queuing runtimes
  - Popular pattern for task and data parallel programming on Symmetric Multiprocessor (SMP) systems
  - Characterized by:
    - A work queue per core
    - Runtime library that divides large loops into tasks and distributes to queues
    - A work stealing scheduler that keeps system balanced
- HSA is designed to extend this pattern to run on heterogeneous systems

# hQ: heterogeneous queuing

- How compute dispatch operates today in the driver model

# hQ: heterogeneous queuing

- How compute dispatch improves under HSA
  - Application codes to the hardware
  - User mode queuing
  - Hardware scheduling
  - Low dispatch times

  - No Soft Queues
  - No User Mode Drivers
  - No Kernel Mode Transitions
  - No Overhead!

# hQ: heterogeneous queuing

- AQL (Architected Queueing Layer) enables any agent to enqueue tasks

# hQ: heterogeneous queuing

- AQL (Architected Queueing Layer) enables any agent to enqueue tasks
  - Single compute dispatch path for all hardware
  - No driver translation, direct access to hardware
  - Standard across vendors



- All agents can enqueue
  - Allowed also self-enqueuing
- Requires coherency and shared virtual memory

# hQ: heterogeneous queuing

- A work stealing scheduler that keeps system balanced

# Advantages of the queuing model

• Today's picture:

# Advantages of the queuing model



- The unified shared memory allows to share pointers among different processing elements thus avoiding explicit memory transfer requests

# Advantages of the queuing model



- Coherent caches remove the necessity to perform explicit synchronization operation

# Advantages of the queuing model



- The supported signaling mechanism enables asynchronous events between agents without involving the OS kernel

# Advantages of the queuing model



- Tasks are directly enqueued by the applications without using OS mechanisms

# Advantages of the queuing model

- HSA picture:

# Device side queuing

- Let's consider a tree traversal problem:
  - Every node in the tree is a job to be executed
  - We may not know at priory the size of the tree
  - Input parameters of a job may depend on parent execution



- Each node is a job
- Each job may generate some child jobs

# Device side queuing

- State-of-the-art solution:
  - The job has to communicate to the host the new jobs (possibly transmitting input data)
  - The host queues the child jobs on the device

- Each node is a job
- Each job may generate some child jobs

**Considerable memory traffic!**

# Device side queuing

- Device side queuing:
  - The <mark>job</mark> <mark>running on the device</mark> <mark>directly queues</mark> <mark>new jobs</mark> in the device/host queues



- Each node is a job
- Each job may generate some child jobs

# Device side queuing

- Benefits of device side queuing:
  - Enable more natural expression of nested parallelism necessary for applications with irregular or data-driven loop structures(i.e., breadth first search)
  - Remove of synchronization and communication with the host to launch new threads (remove expensive data transfer)
  - The finer granularities of parallelism is exposed to scheduler and load balancer

# Device side queuing

- OpenCL 2.0 supports device side queuing
  - Device-side command queues are out-of-order
  - Parent and child kernels execute asynchronously
  - Synchronization has to be explicitly managed by the programmer

# Summary on the queuing model

- ==User mode queuing== for ==low latency dispatch==
  - – ==Application dispatches directly==
  - – No OS or driver required in the dispatch path
- ==Architected Queuing Layer==
  - – Single compute dispatch path for all hardware
  - – No driver translation, direct to hardware
- ==Allows for dispatch to queue from any agent==
  - – CPU or GPU
- ==GPU self-enqueue enables lots of solutions==
  - – Recursion
  - – Tree traversal
  - – Wavefront reforming

# Other necessary HW mechanisms

- Task preemption and context switching have to be supported by all computing resources (also GPUs)
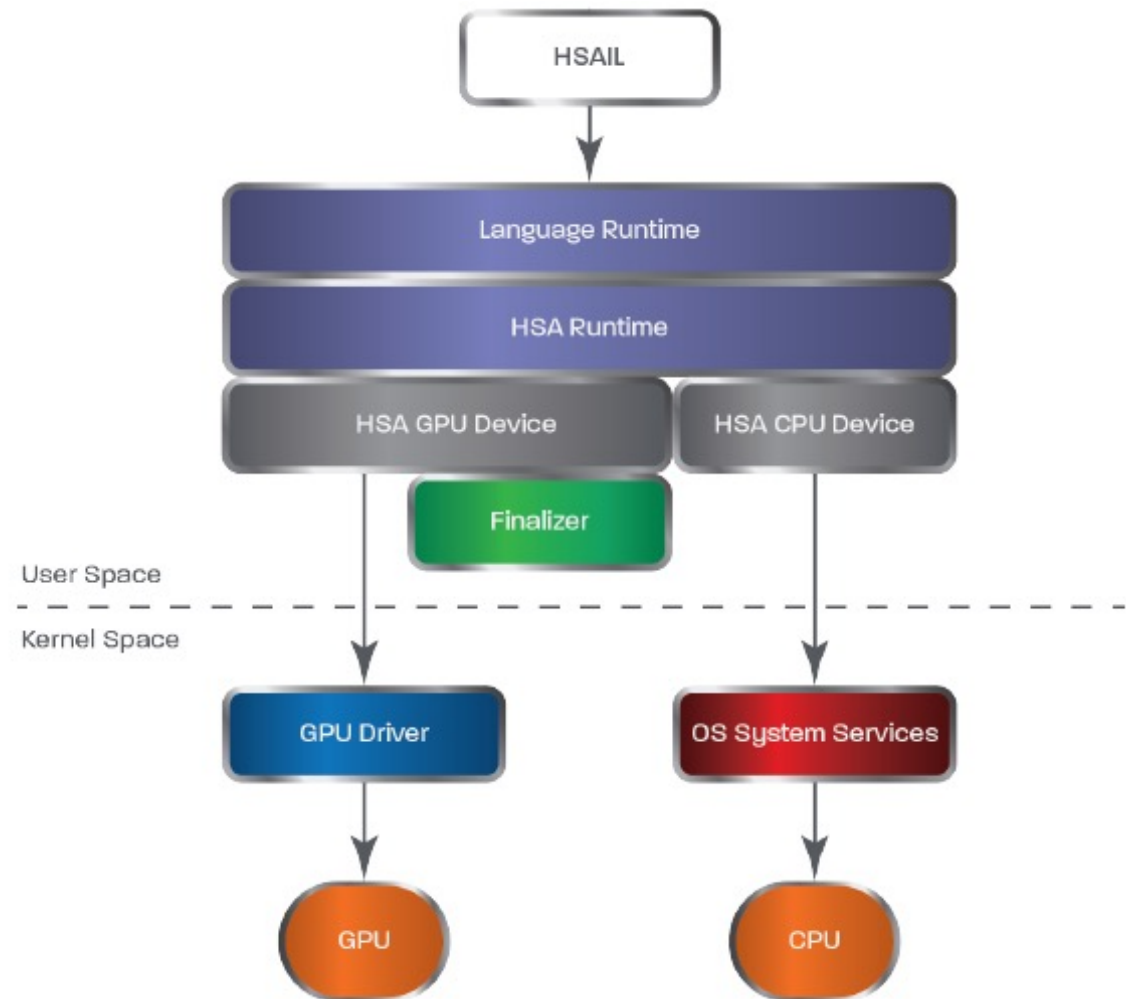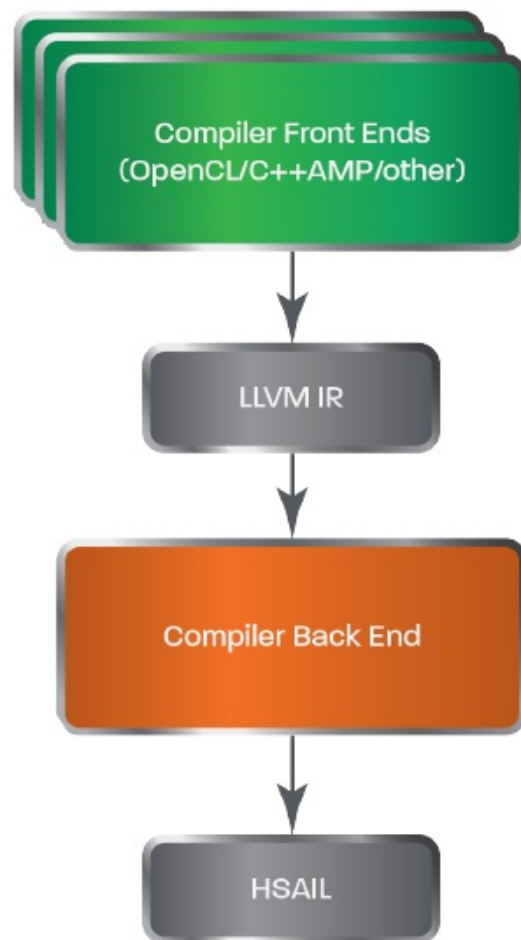
# Key features of HSA

- **hUMA** – Heterogeneous Unified Memory Architecture
- **hQ** – Heterogeneous Queuing
- <u>HSAIL – HSA Intermediate Language</u>

# HSA intermediate layer (HSAIL)

- A portable "virtual ISA" for vendor-independent compilation and distribution
  - Like Java bytecodes for GPUs

- Low-level IR, close to machine ISA level
  - Most optimizations (including register allocation) performed before HSAIL

- Generated by a high-level compiler (LLVM, gcc, Java VM, etc.)
  - Application binaries may ship with embedded HSAIL

- Compiled down to target ISA by a vendor-specific "finalizer"
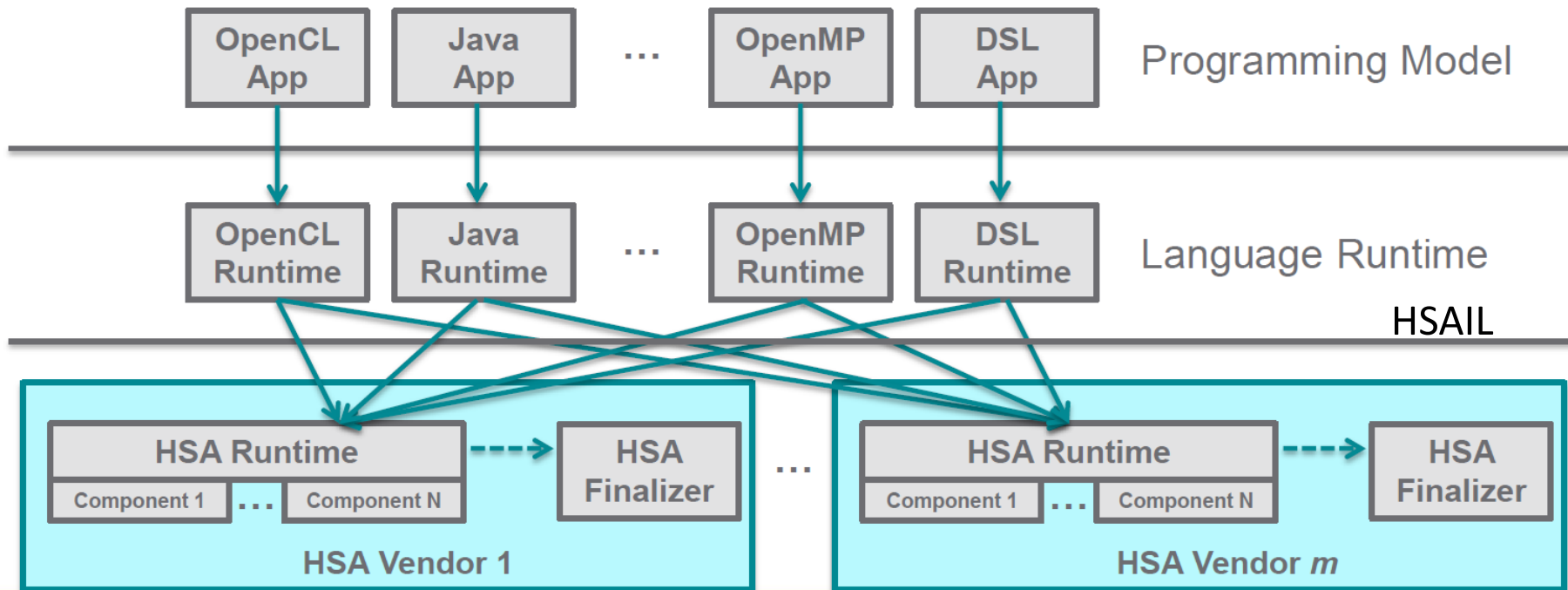  - Finalizer may execute at run time, install time, or build time

# HSA intermediate layer (HSAIL)

- HSA compilation stack
- HSA runtime stack

# HSA software stack

- HSA supports many languages

# Specifications and software

# HSA architecture V1

- GPU compute C++ support
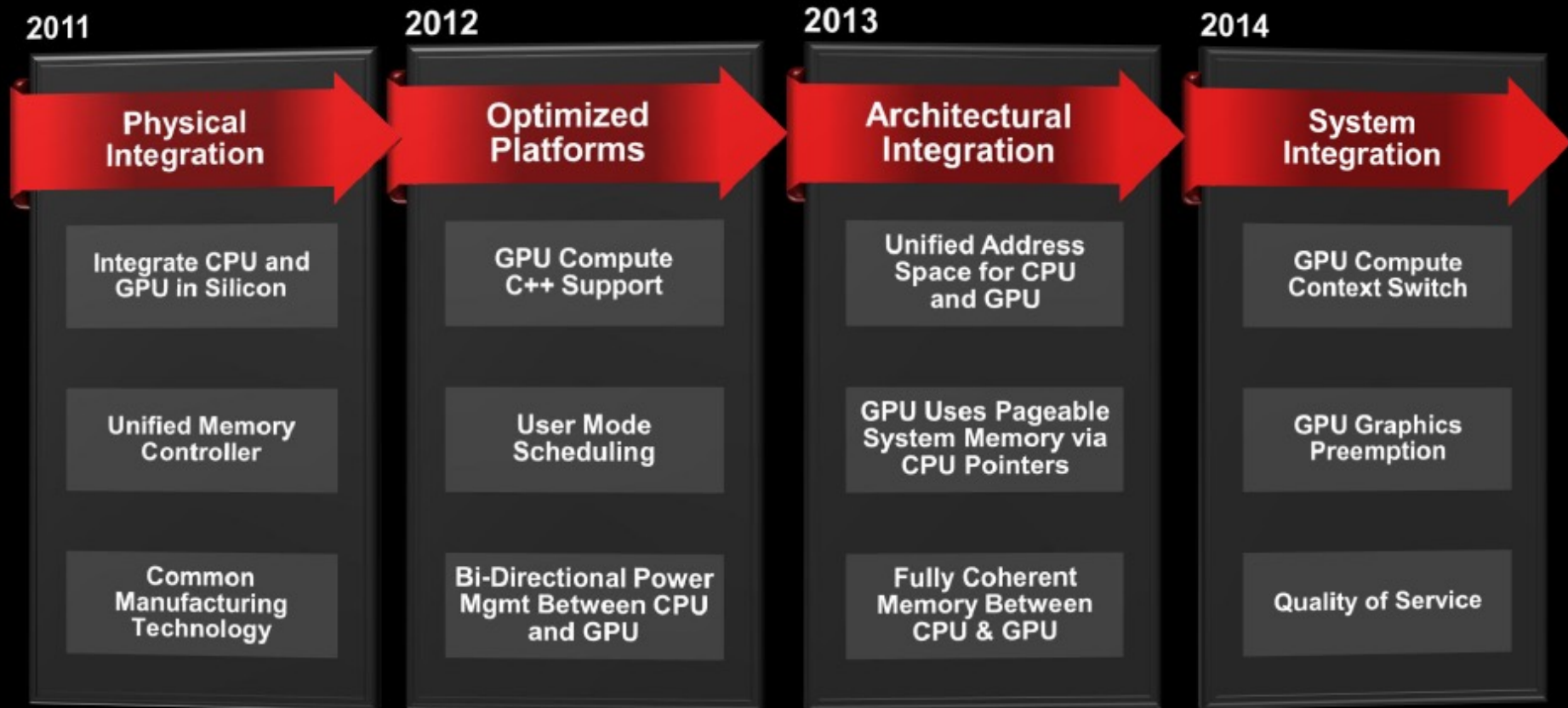- User Mode Scheduling
- Fully coherent memory between CPU & GPU
- GPU uses pageable system memory via CPU pointers
- GPU graphics pre-emption
- GPU compute context switch

# AMD roadmaps
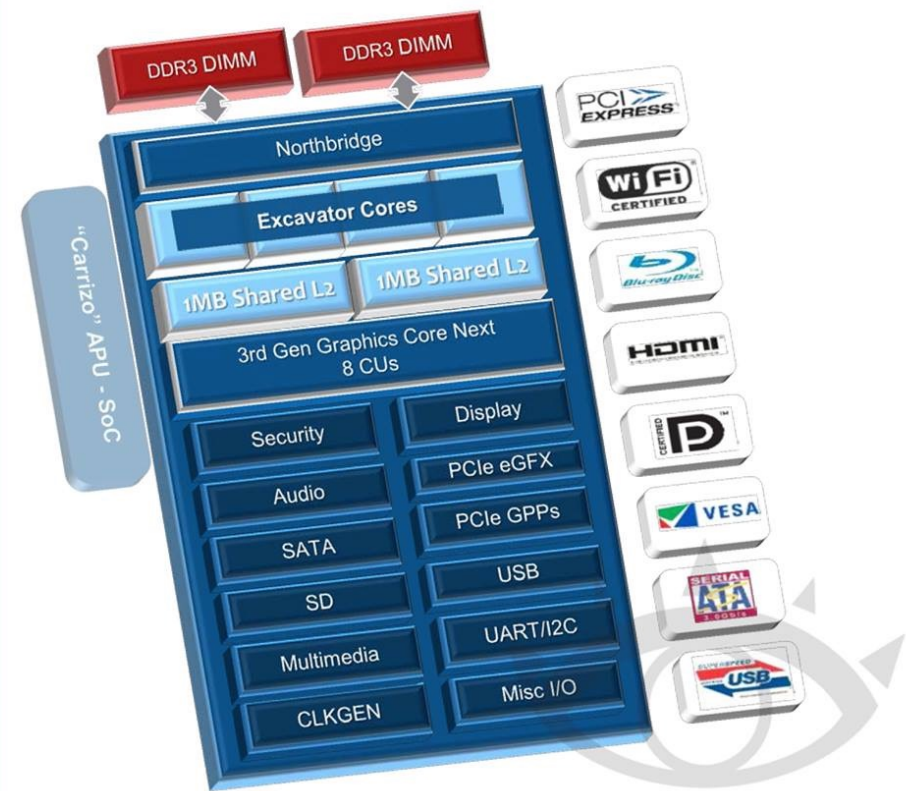


**HETEROGENEOUS SYSTEM ARCHITECTURE ROADMAP**

| 2011 | 2012 | 2013 | 2014 |
|---|---|---|---|
| **Physical Integration** | **Optimized Platforms** | **Architectural Integration** | **System Integration** |
| Integrate CPU and GPU in Silicon | GPU Compute C++ Support | Unified Address Space for CPU and GPU | GPU Compute Context Switch |
| Unified Memory Controller | User Mode Scheduling | GPU Uses Pageable System Memory via CPU Pointers | GPU Graphics Preemption |
| Common Manufacturing Technology | Bi-Directional Power Mgmt Between CPU and GPU | Fully Coherent Memory Between CPU & GPU | Quality of Service |

8 | XLDB - Stanford | Sept. 12, 2012

AMD Fusion¹² DEVELOPER SUMMIT

# AMD roadmaps