ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich
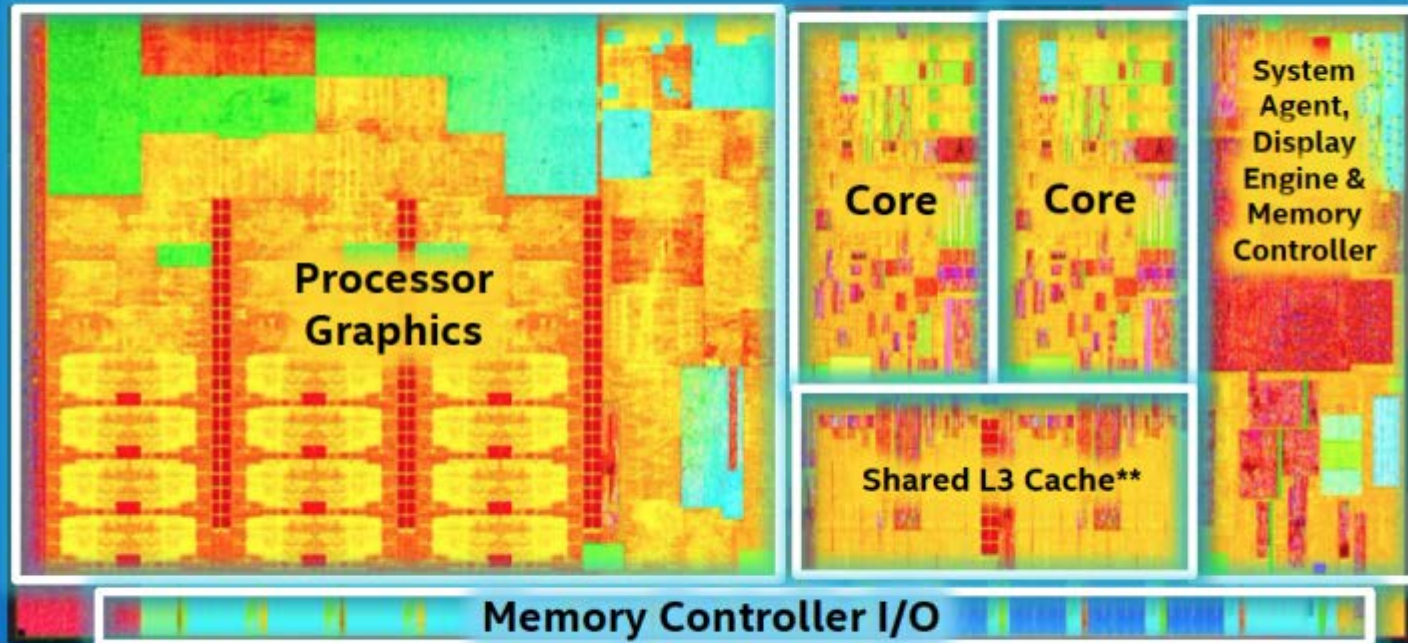
Department of Information Technology
and Electrical Engineering

# Heterogeneous Architecture

Luca Benini

lbenini@iis.ee.ethz.ch

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

IHEI
Department of Information Technology
and Electrical Engineering

# Intel's Broadwell



Intel® Core™ M Processor Die Map
14nm 2nd Generation Tri-Gate 3-D Transistors

Processor Graphics

Core

Core

System Agent, Display Engine & Memory Controller

Shared L3 Cache**

Memory Controller I/O

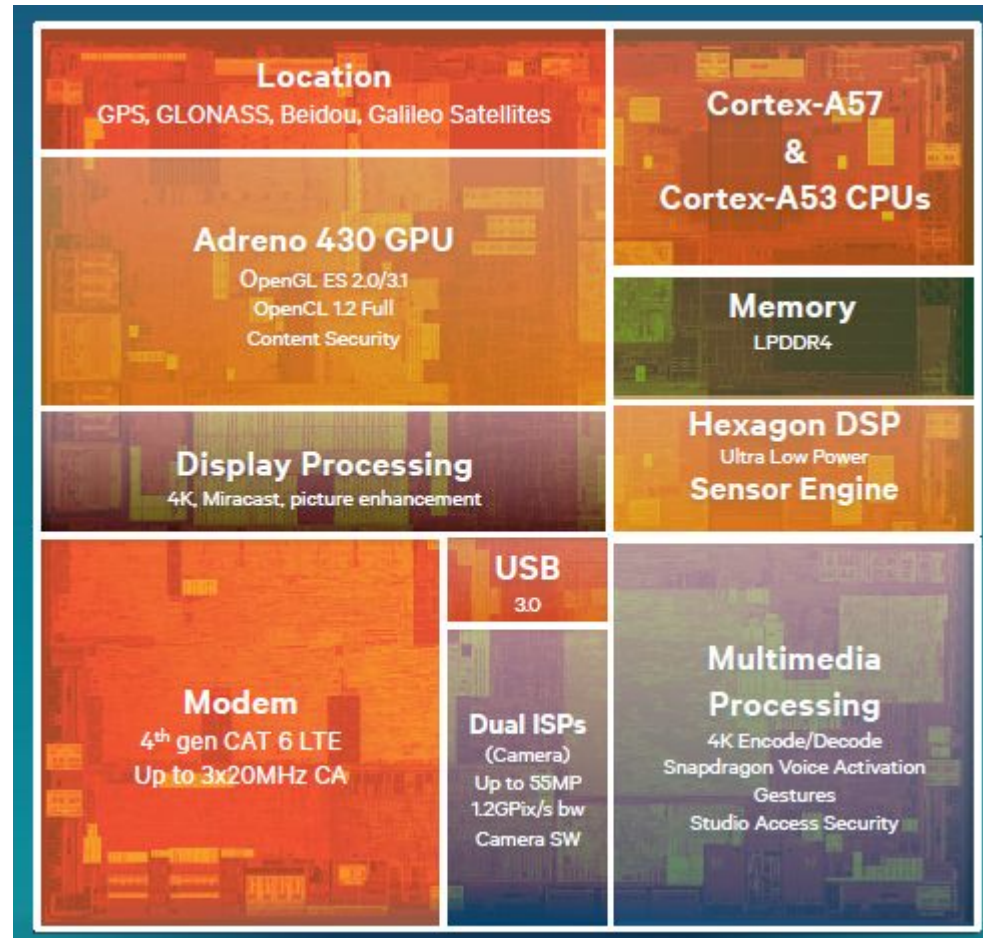Dual Core Die Shown Above | Transistor Count: 1.3 Billion | Die Size: 82mm²
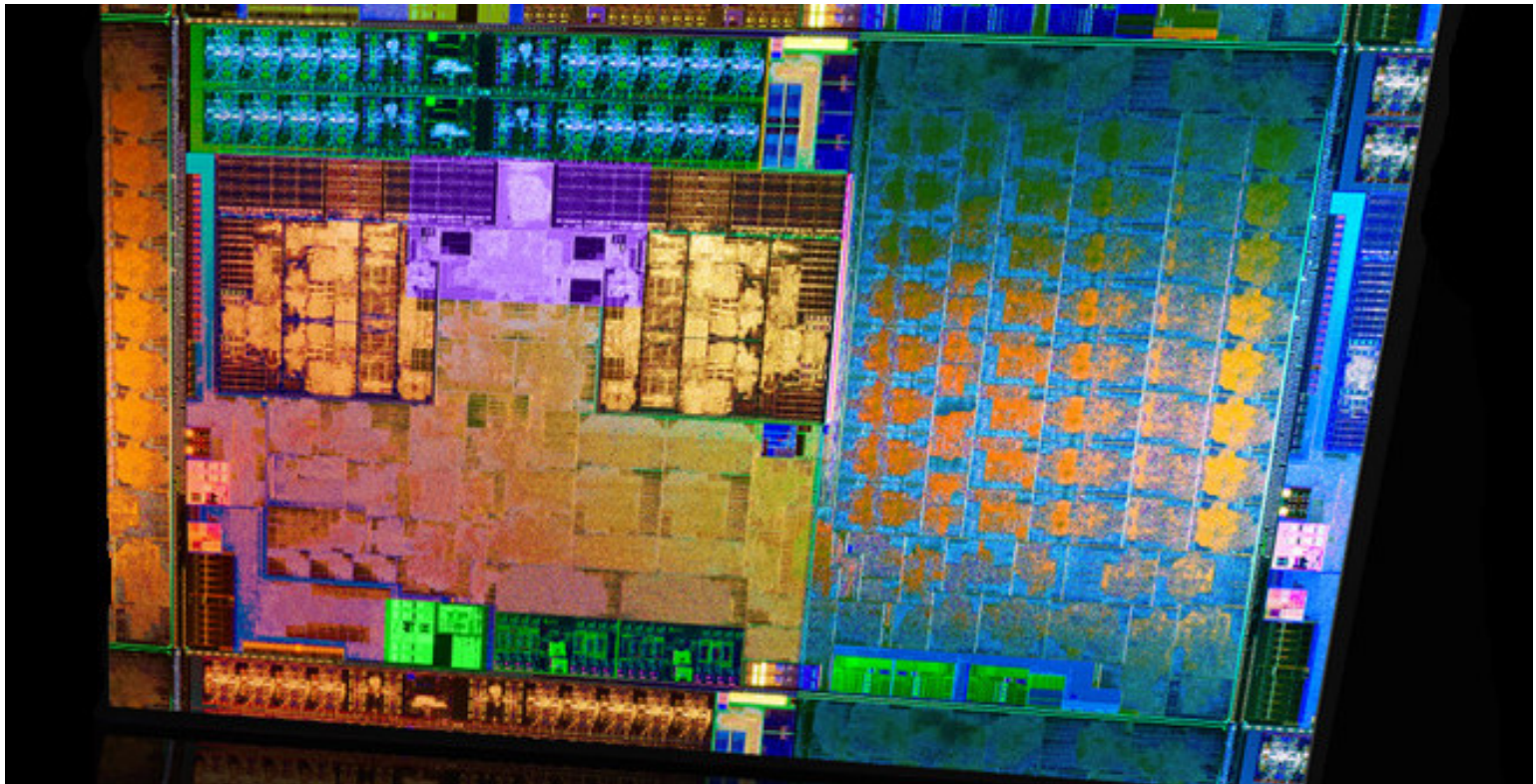4th Gen Core Processor ( Y series): .96B        4th Gen Core Processor ( Y series): 131mm²
** Cache is shared across both cores and processor graphics
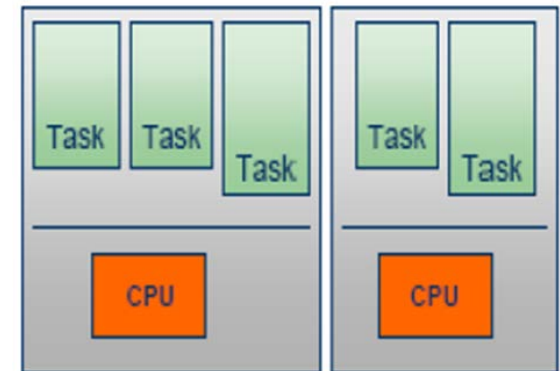
# Qualcomm's Snapdragon 810

# AMD Bristol Ridge

# Heterogeneous Multicores

Deszo Sima (Univ. Budapest)

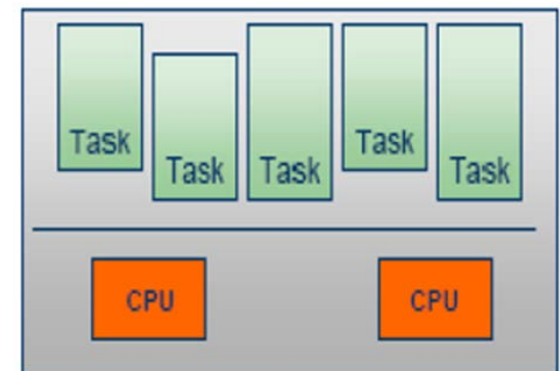# Vocabulary in the multi-era

- ### AMP (Asymmetric MP)
  - Each processor has local memory
  - Tasks statically allocated to one processor
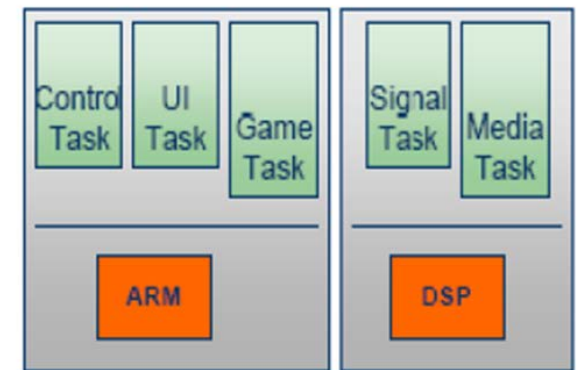
- ### SMP (Symmetric MP)
  - Processors share memory
  - Tasks dynamically scheduled to any processor
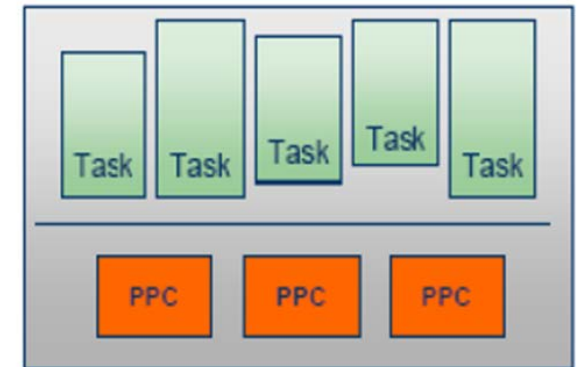
# Vocabulary in the multi-era

- **Heterogeneous:**
  - Specialization among processors
  - Often different instruction sets
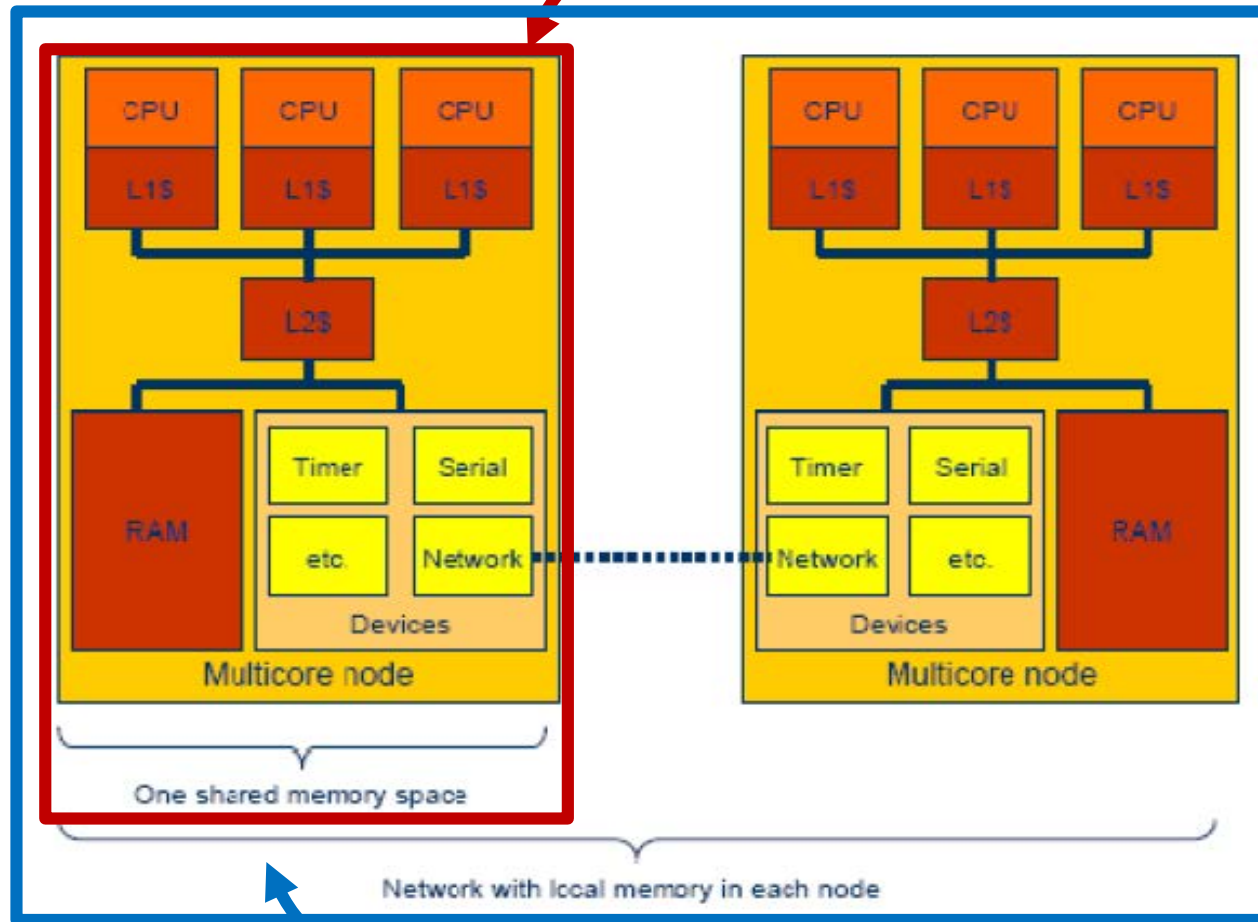  - Usually AMP design

- **Homogeneous:**
  - All processors have the same instruction set
  - Processors can run any task
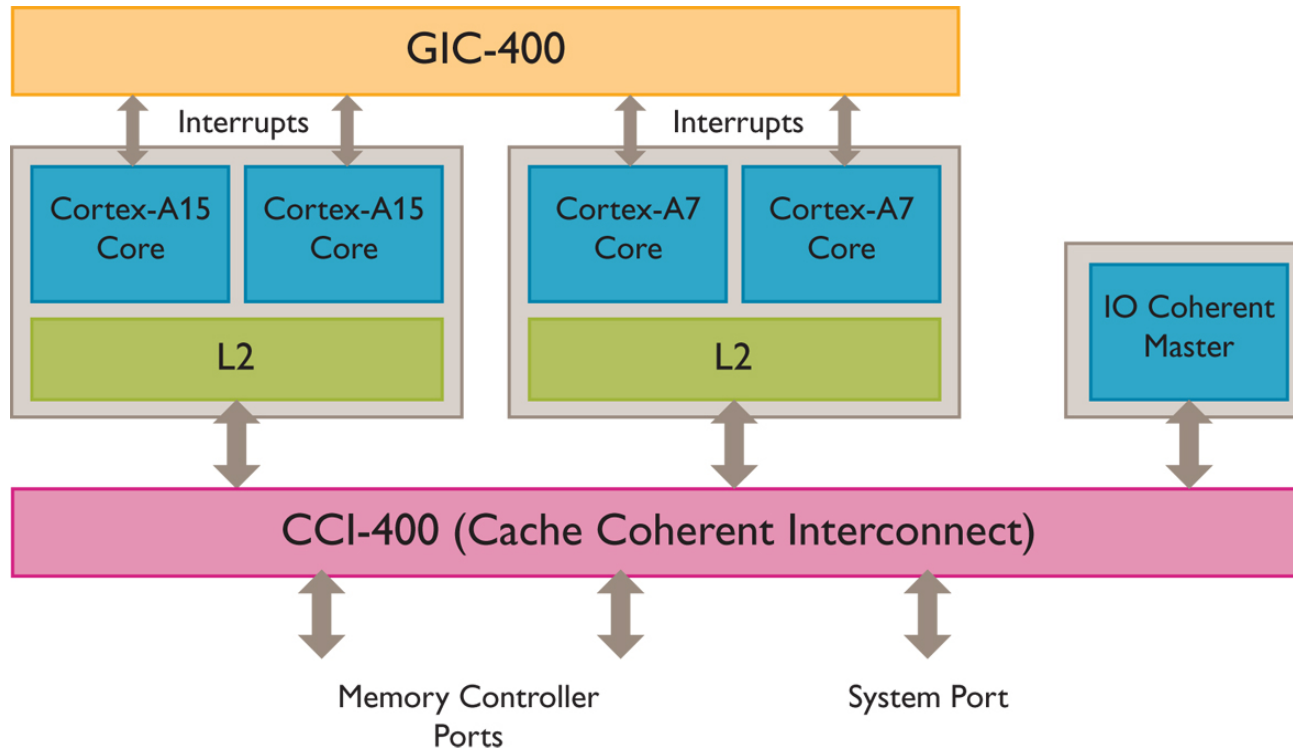  - Usually SMP design
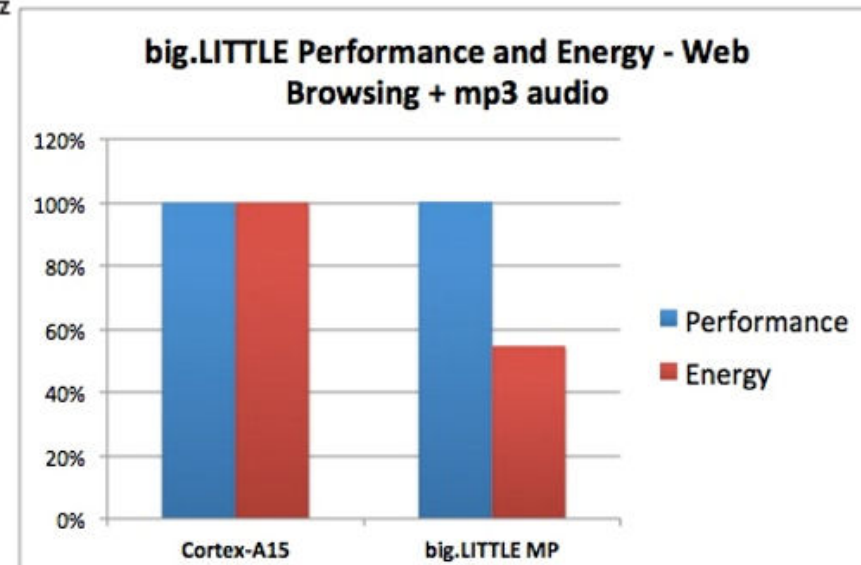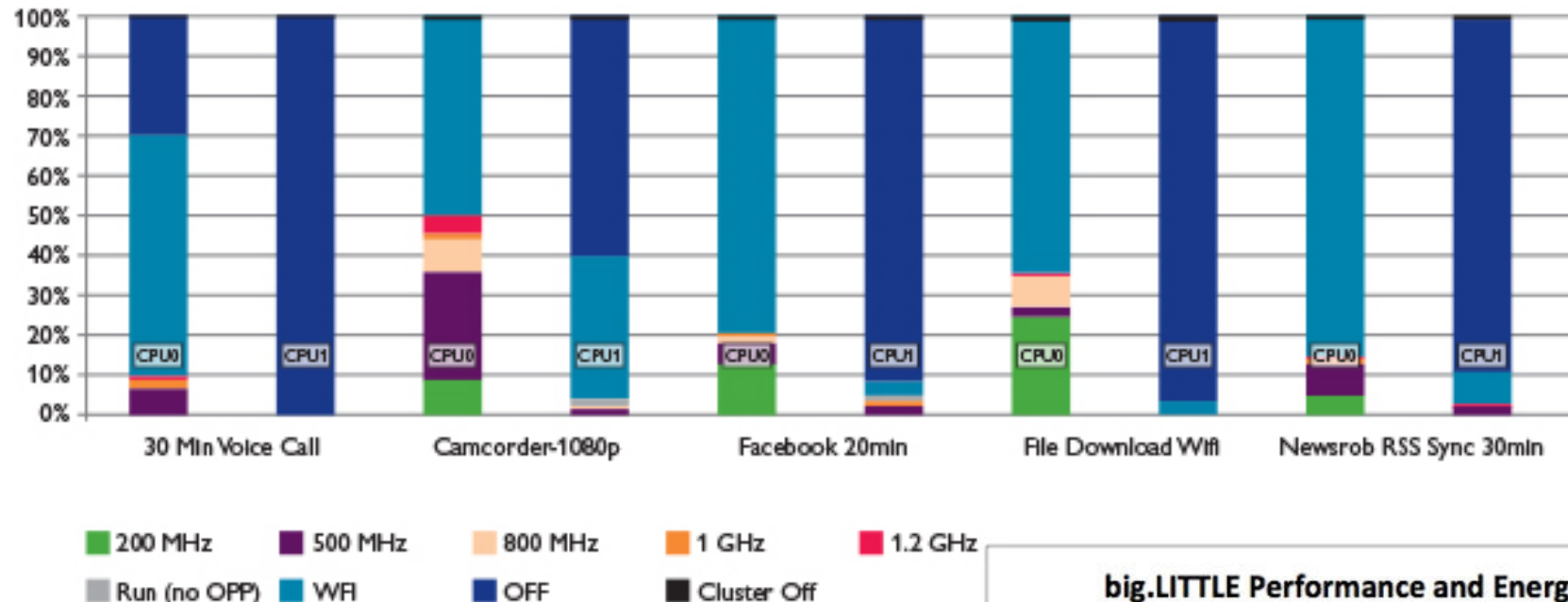
# Future many-cores



**Locally homogeneous**

**Globally heterogeneous**

# ARM big.LITTLE



- DVFS
- Task Migration
- SW Overhead
- Memory coherency

# ARM big.LITTLE Energy Saving

# 4.3 Principle of ARM's big.LITTLE technology (3)

## Exclusive/inclusive use of the clusters



**Exclusive/inclusive use of the clusters**

**Exclusive use of the clusters**

Clusters are used exclusively, i.e. at a time one of the clusters is in use as shown below for the cluster migration model (to be discussed later)

**Inclusive use of the clusters**

Clusters are used inclusively, i.e. at a time both clusters can be used partly or entirely

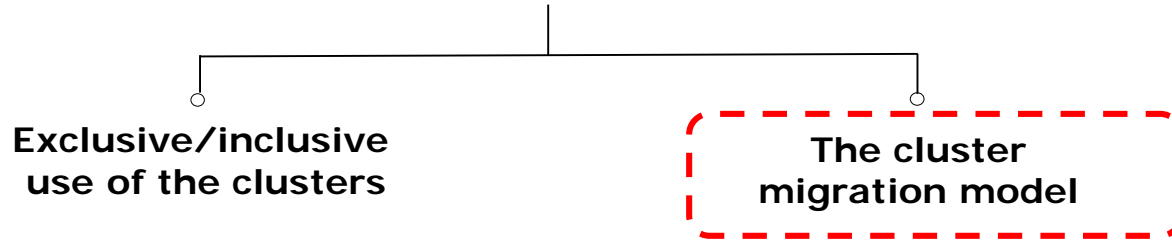Usage models of synchronous adaptive SMPs in the n + n configuration

Usage models of synchronous adaptive SMPs in the n+n configuration

Exclusive/inclusive
use of the clusters

The cluster
migration model

# 4.3 Principle of ARM's big.LITTLE technology (4)

The cluster migration model [5]

The cluster migration model

Exclusive use
of the clusters

Inclusive use
of the clusters

Cluster migration

Core migration

Core migration



*big.LITTLE processing
with cluster migration*

*big.LITTLE processing
with core migration*

*big.LITTLE MP*

# 4.3 Principle of ARM's big.LITTLE technology (5)

Big.LITTLE processing with cluster migration [5]

- There are two core clusters, the LITTLE core cluster and the big core cluster.
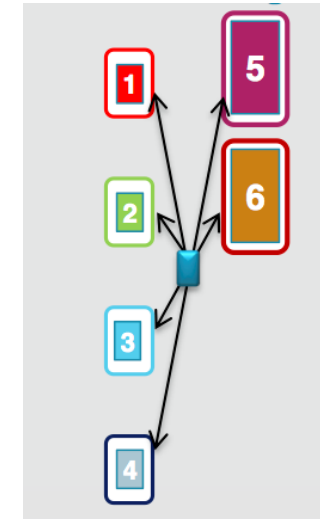- Tasks run on either the LITTLE or the big core cluster, so only one core cluster is active at any time (except a short interval during a cluster switch).
- Low workloads, such as background synch tasks, audio or video playback run typically on the LITTLE core cluster.
- If the workload becomes higher than the max performance of the LITTLE core cluster the workload will be migrated to the big core cluster and vice versa.

**Cluster Migration**

# 4.3 Principle of ARM's big.LITTLE technology (6)

Cluster switches [6]

- Cluster selection is driven by OS power management.

- OS (e.g. the Linux cpufreq routine) samples the load for all cores in the cluster and selects an operating point for the cluster.

- It switches clusters at terminal points of the current clusters DVFS curve, as illustrated in the next Figure.

Power/performance curve during cluster switching [7]



- A switch from the low power cluster to the high performance cluster is an extension of the DVFS strategy.
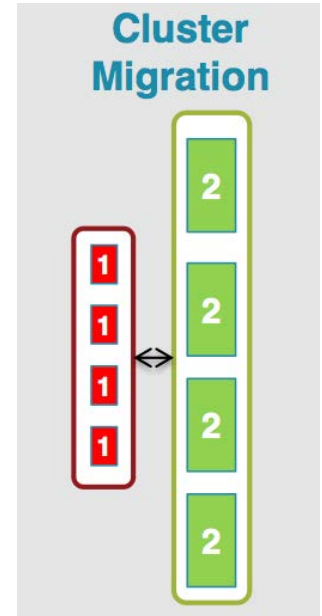- A cluster switch lasts about 30 kcycles.

# 4.3 Principle of ARM's big.LITTLE technology (8)

Big.LITTLE processing with core migration [5], [8]

- There are two core clusters, the LITTLE core cluster and the big core cluster.
- Cores are grouped into pairs of one big core and one LITTLE core.
  The LITTLE  and the big core of a group are used exclusively.
- Each LITTLE core can switch to its big counterpart if it meets a higher load than its max. performance and vice versa.
- Each core switch is independent  from the others.

# 4.3 Principle of ARM's big.LITTLE technology (9)

Core switches [6]

- Core selection in any core pair is performed by OS power management.
- The DVFS algorithm monitors the core load.

  When a LITTLE core cannot service the actual load, a switch to its big counterpart is initiated and the LITTLE core is turned off and vice versa.

**big.LITTLE MP processing with core migration** [8],[5]

- The OS scheduler has all cores of both clusters at its disposal and can activate all cores at any time.
- Tasks can run or be moved between the LITTLE cores and the big cores as decided by the scheduler.
- big.LITTLE MP termed also as Heterogeneous Multiprocessing (HMP).

Use of the big.LITTLE technology in recent mobile processors

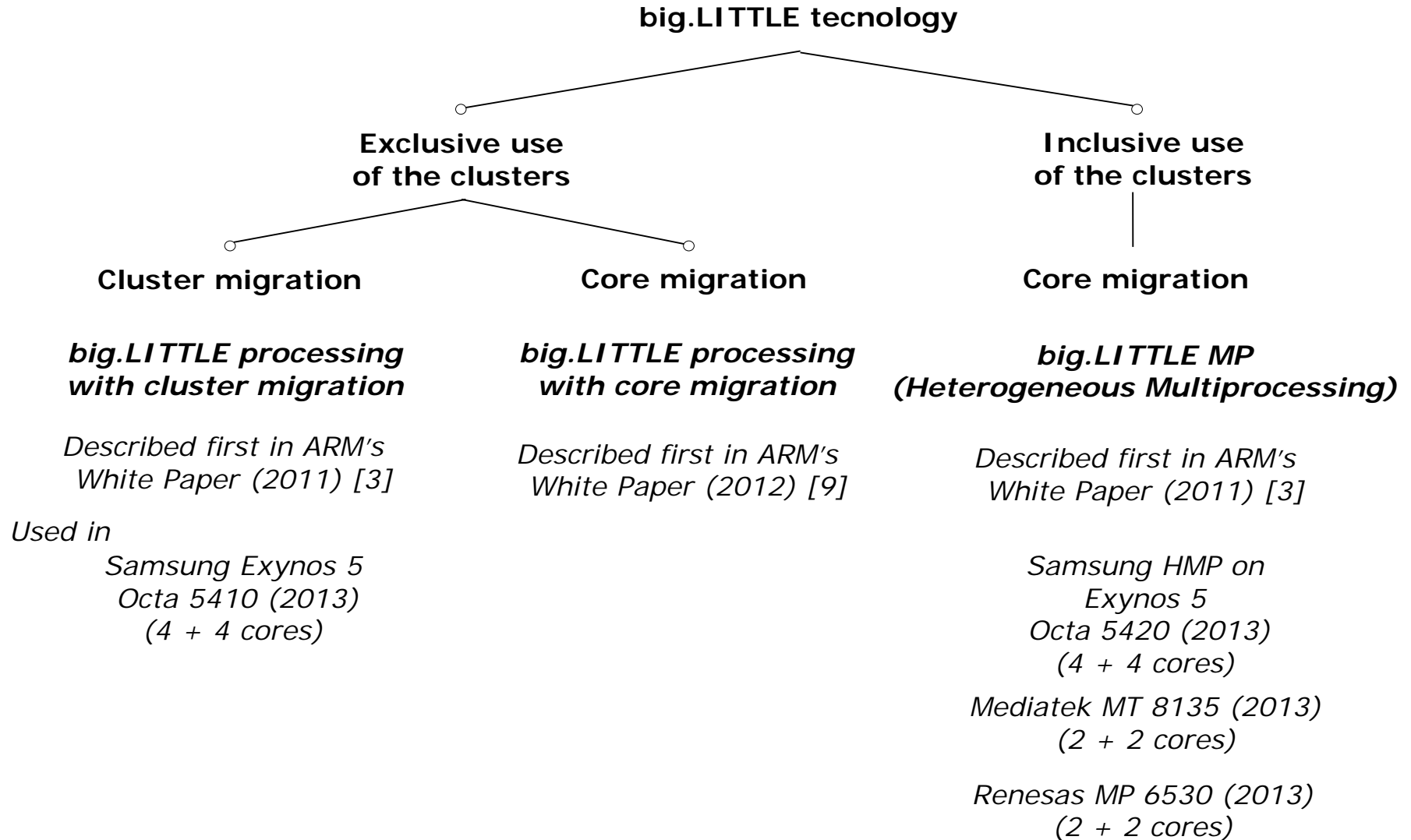**big.LITTLE tecnology**

**Exclusive use
of the clusters**

**Inclusive use
of the clusters**

**Cluster migration**

**Core migration**

**Core migration**

***big.LITTLE processing
with cluster migration***

***big.LITTLE processing
with core migration***

***big.LITTLE MP
(Heterogeneous Multiprocessing)***

*Described first in ARM's
White Paper (2011) [3]*

*Described first in ARM's
White Paper (2012) [9]*

*Described first in ARM's
White Paper (2011) [3]*

*Used in*

*Samsung Exynos 5
Octa 5410 (2013)
(4 + 4 cores)*

*Samsung HMP on
Exynos 5
Octa 5420 (2013)
(4 + 4 cores)*

*Mediatek MT 8135 (2013)
(2 + 2 cores)*

*Renesas MP 6530 (2013)
(2 + 2 cores)*

# CPU+GPU integration

Deszo Sima (Univ. Budapest)

## 1. Introduction to architectural integration of the CPU and GPU

**Aim of architectural integration of the CPU and GPU**

**Accelerated graphics processing**

- Accelerating graphics processing
  by the higher bandwidth
  of connecting the CPU and the GPU.
- Cost reduction

**Heterogeneous processing**

- Accelerating HPC by the GPU,
  i.e. by the large number of FP units
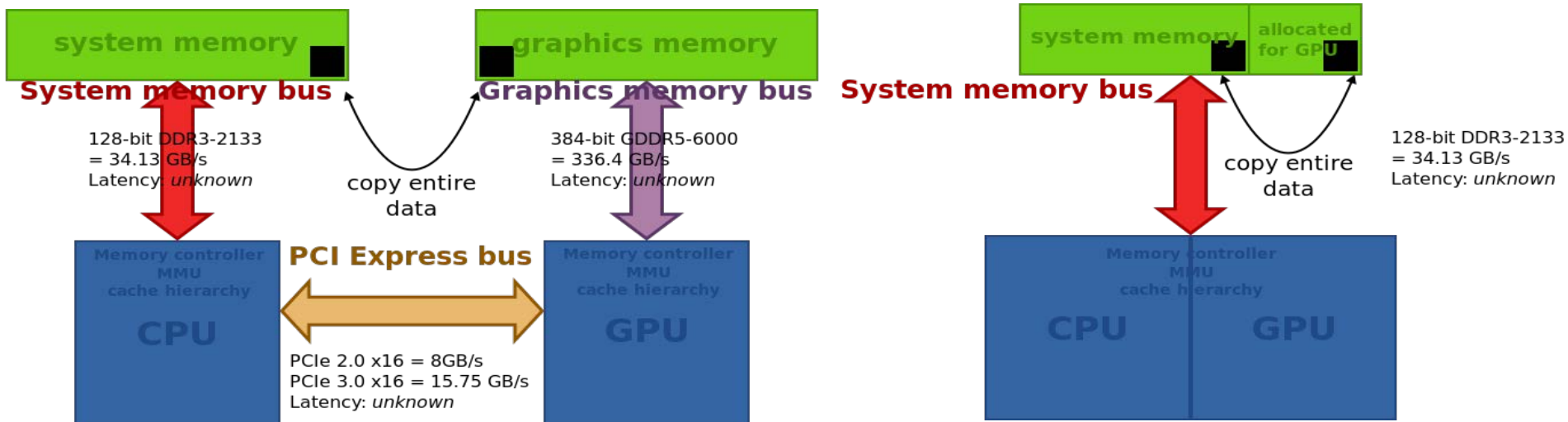  available in a GPU.

Needs HLL support (CUDA, OPenCL)

**Implementation alternatives of graphics memory** [1]

**Implementation of the graphics memory**

**Discrete graphics memory**

**Unified system memory**
**(Shared memory)**
**(Unified Virtual Memory)**
**(Unified Memory Architecture)**

*Implementation example*



system memory

graphics memory

**System memory bus**

**Graphics memory bus**

128-bit DDR3-2133
= 34.13 GB/s
Latency: *unknown*

384-bit GDDR5-6000
= 336.4 GB/s
Latency: *unknown*

copy entire
data

**PCI Express bus**

Memory controller
MMU
cache hierarchy

**CPU**

Memory controller
MMU
cache hierarchy

**GPU**

PCIe 2.0 x16 = 8GB/s
PCIe 3.0 x16 = 15.75 GB/s
Latency: *unknown*

system memory | allocated for GPU

**System memory bus**

128-bit DDR3-2133
= 34.13 GB/s
Latency: *unknown*

copy entire
data

Memory controller
MMU
cache hierarchy

**CPU**        **GPU**

*Typ. use*

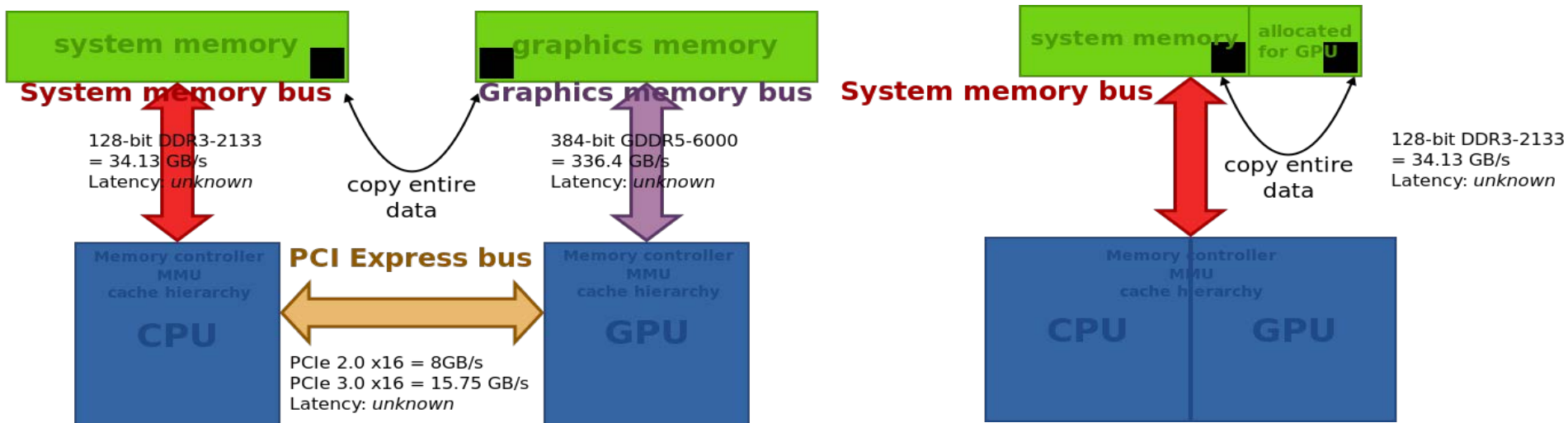*Graphics cards*

On-die integrated graphics

**Key benefit/drawback of the Unified system memory** [1]

**Implementation of the graphics memory**

**Discrete graphics memory**

**Unified system memory**
**(Shared memory)**
**(Unified Virtual Memory)**
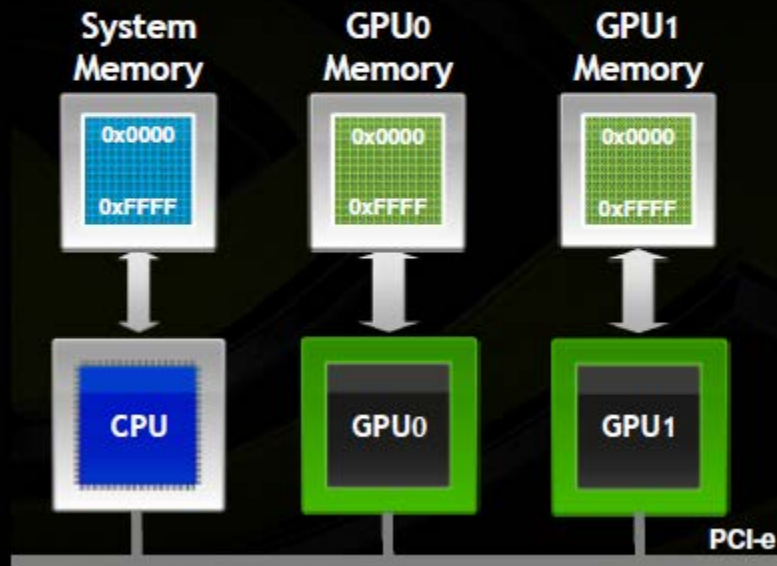**(Unified Memory Architecture)**

*Implementation example*

| system memory | | graphics memory | | system memory | allocated for GPU |

**System memory bus**

128-bit DDR3-2133
= 34.13 GB/s
Latency: *unknown*

copy entire
data

**Graphics memory bus**

384-bit GDDR5-6000
= 336.4 GB/s
Latency: *unknown*

**System memory bus**

copy entire
data

128-bit DDR3-2133
= 34.13 GB/s
Latency: *unknown*

Memory controller
MMU
cache hierarchy

**CPU**

**PCI Express bus**

Memory controller
MMU
cache hierarchy

**GPU**

PCIe 2.0 x16 = 8GB/s
PCIe 3.0 x16 = 15.75 GB/s
Latency: *unknown*

Memory controller
MMU
cache hierarchy

**CPU**     **GPU**

It eliminates the need for an extra graphics memory controller and bus but has the constraints of a reduced memory bandwidth.
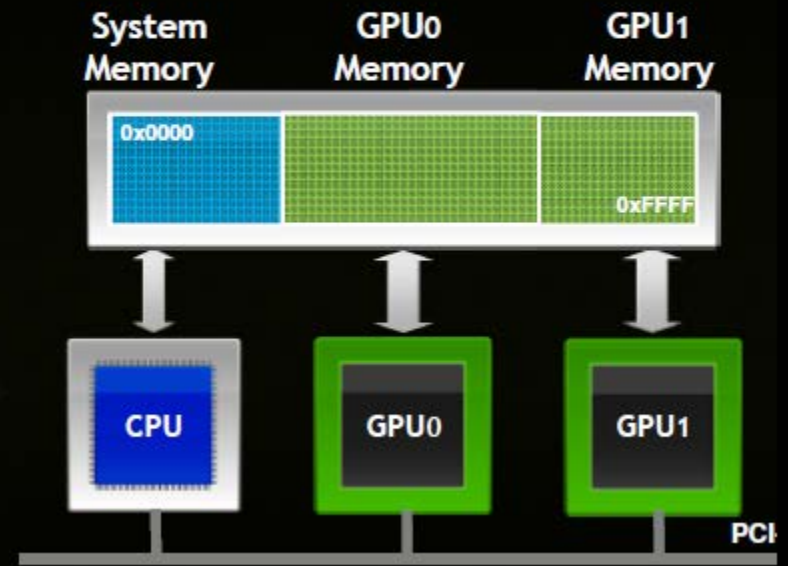
**Address spaces for Discrete and Unified System Memory (USM)** [2]



No USM: *Multiple Memory Spaces* — USM : *Single Address Space*

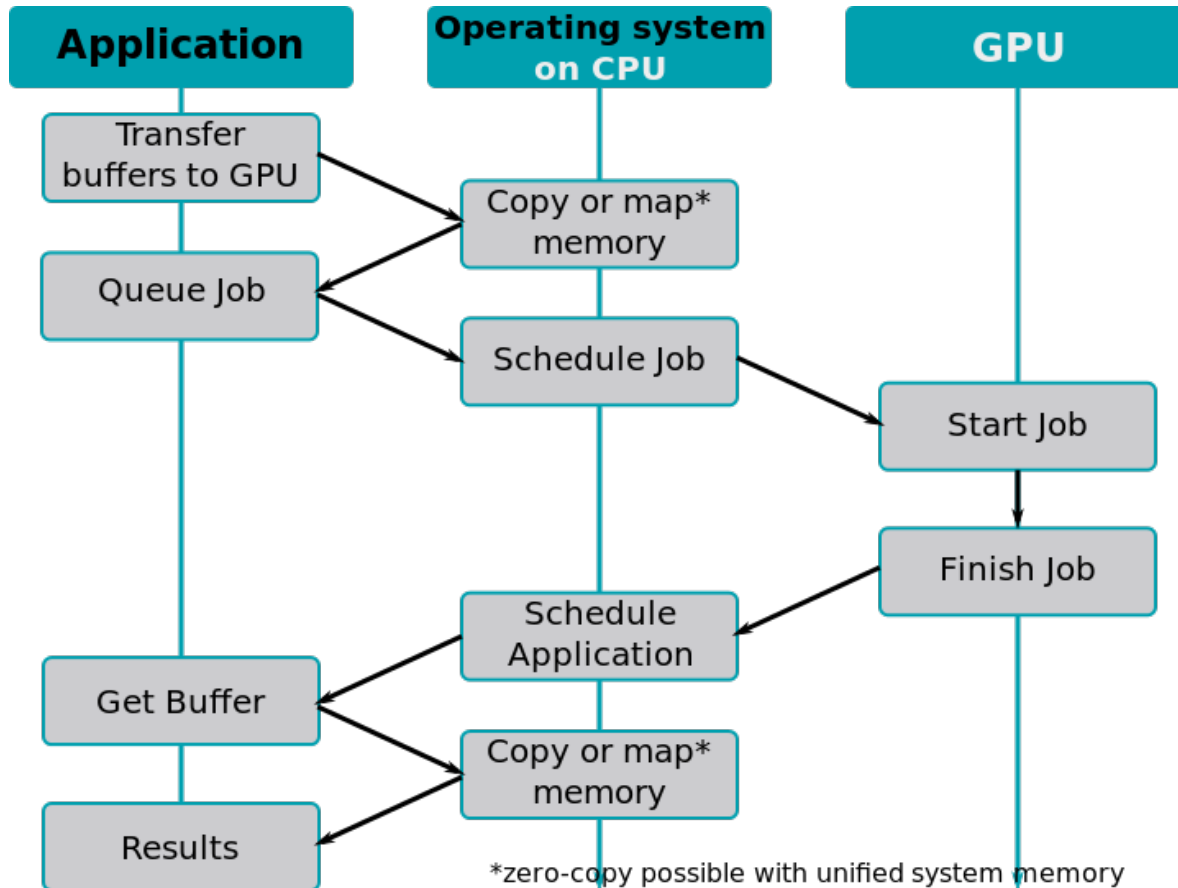**Main difficulty of heterogeneous computing**

- GPUs operate  from their own address spaces, thus without appropriate hardware and software
    support data to be processed by GPUs need to be loaded into their address space and
    the results need to be sent back to the host.

- Transferring data may be avoided with unified system memory and suitable software support
    (CUDA 4.0 or OpenCL 1.2 or higher) assuming appropriate hardware.

   Then data transfer will be substituted by address mapping (called zero copy).

**Tasks to be performed for copying data/mapping address spaces for GPUs** [1]



*zero-copy possible with unified system memory

# 2. AMD's approach to support heterogeneous computing

**Aim of HSA-1** [4]



**EFFECTIVE COMPUTE OFFLOAD IS MADE EASY BY HSA**

APP Accelerated Software Applications

Accelerated Processing Unit

Graphics Workloads

Data Parallel Workloads

Serial and Task Parallel Workloads

DDR3 Controller

GMC

Channel

UNB

L2 Cache    L2 Cache

GPU

Dual    Dual
Core x86    Core x86
Module    Module

PCIe    PCIe    Display PLL    DP / HDMI    Display Controller

AMD HD Media Accelerator

AMD Fusion[12] DEVELOPER SUMMIT

**Aim of HSA-2** [5]

**The entire HSA solution stack** [6]



HSAIL: HSA Intermediate Language

**Establishment of the HAS Foundation**

- In 6/2012 AMD, ARM, TI Qualcomm, Samsung and several other leading semiconductor firms established the Heterogeneous System Architecture (HSA) Foundation.

- It is a non-profit consortium that aims at defining and promoting open standards for heterogeneous computing.
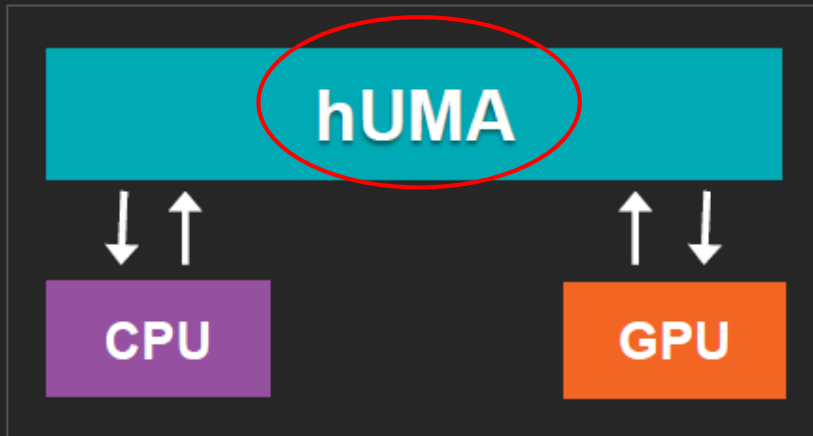
# 2.2 Key hardware enhancements needed to implement HSA

**2.2 Key hardware enhancements needed to implement HSA** [5]



**EQUAL ACCESS TO ENTIRE MEMORY**

hUMA

CPU    GPU

◢ First time ever: GPU and CPU have uniform visibility into entire memory space (up to 32 GB)

hUMA:  heterogeneous UMA
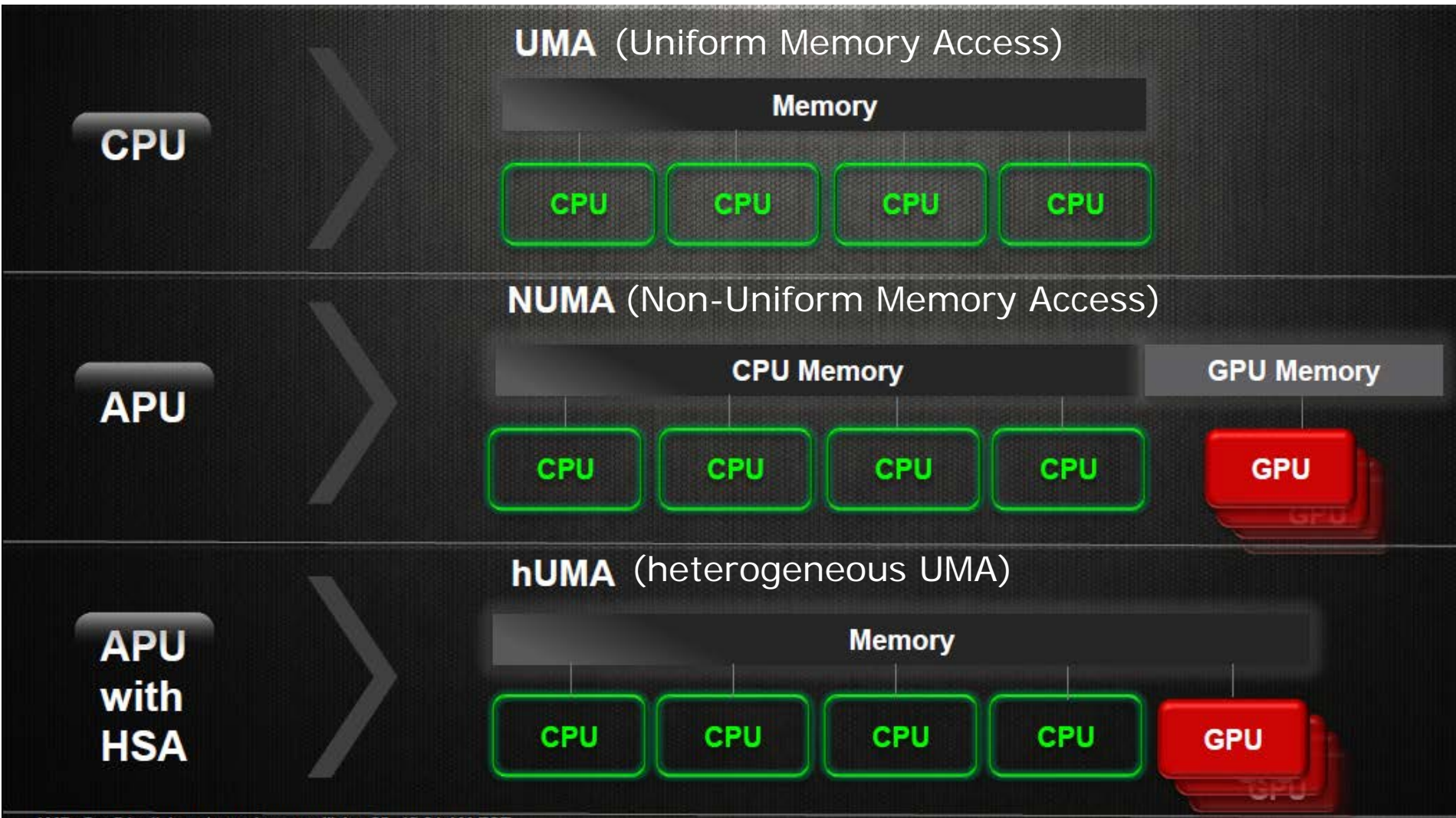UMA:    Uniform Memory Access

**EQUAL FLEXIBILITY TO DISPATCH**

hQ

CPU    GPU

◢ Heterogeneous queuing (hQ) defines how processors interact equally

◢ GPU and CPU have equal flexibility to create/dispatch work

**2.2.1 hUMA (Heterogeneous Uniform Memory Access)**

**The hUMA (heterogenous UMA) memory access scheme** [3]

**GPU co-processing of data structures without hUMA** [3]



- CPU explicitly copies data to GPU memory
- GPU completes computation
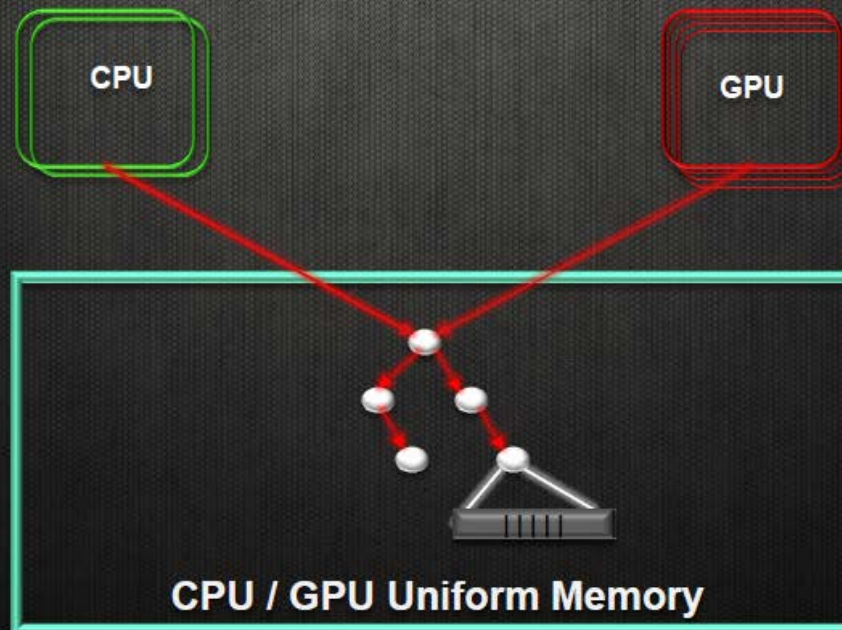- CPU explicitly copies result back to CPU memory

CPU

GPU

CPU Memory

GPU Memory

Only the data array can be copied since GPU cannot follow embedded data-structure links

*A Pointer is a named variable that holds a memory address. It makes it easy to reference data or code segments by a name and eliminates the need for the developer to know the actual address in memory. Pointers can be manipulated by the same expressions used to operate on any other variable

12

**GPU co-processing of data structures with hUMA** [3]

- CPU simply passes a pointer to GPU
- GPU completes computation
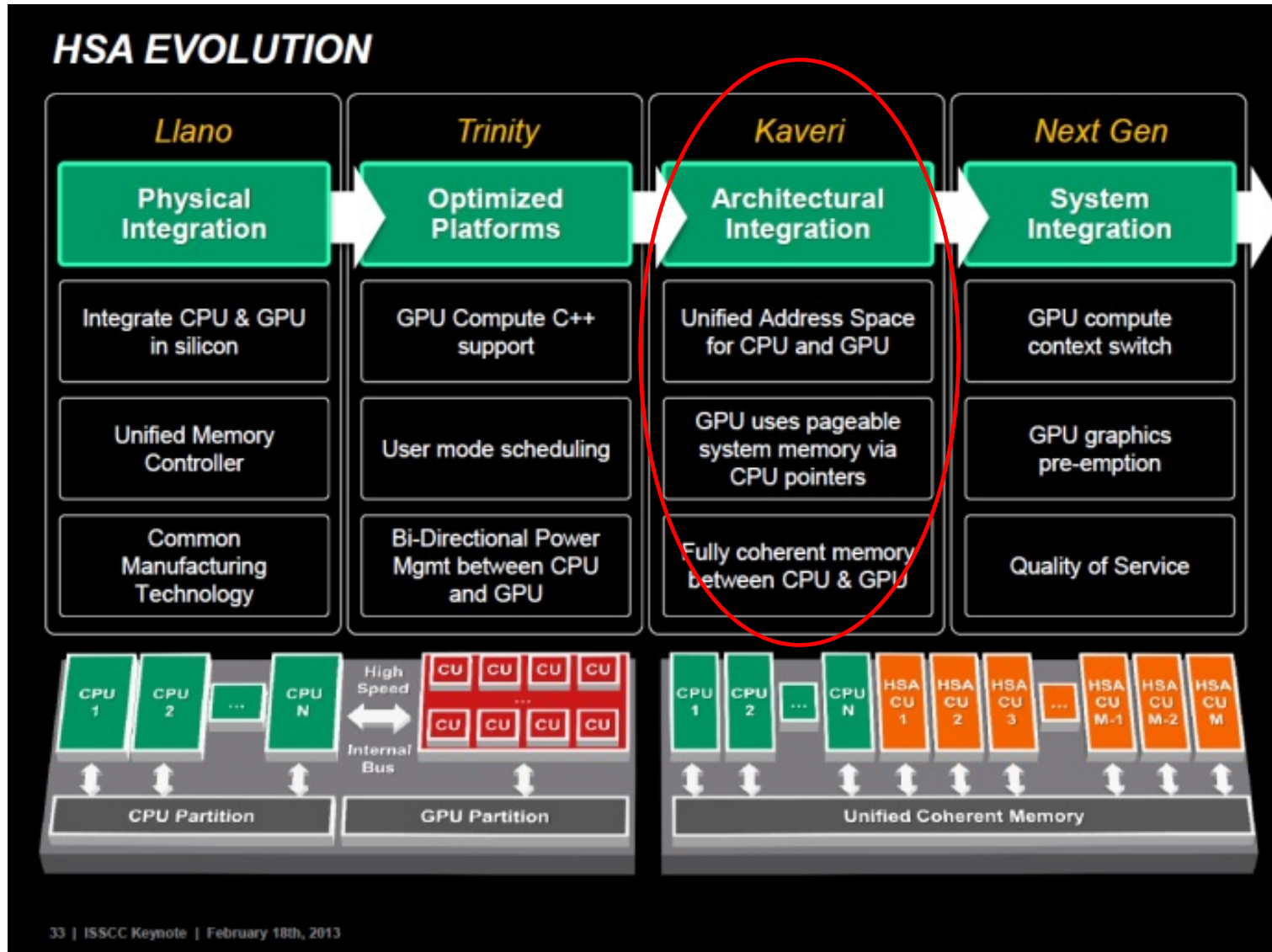- CPU can read the result directly – no copying needed!

CPU

GPU

CPU can pass a pointer to entire data structure since the GPU can now follow embedded links

**CPU / GPU Uniform Memory**

*A Pointer is a named variable that holds a memory address. It makes it easy to reference data or code segments by a name and eliminates the need for the developer to know the actual address in memory. Pointers can be manipulated by the same expressions used to operate on any other variable

**The evolution path to HSA in AMD's subsequent Family 15h based APU lines** [7]

**Key requirements of hUMA-1** [3]



**BI-DIRECTIONAL COHERENT MEMORY**
*Any updates made by one processing element will be seen by all other processing elements – GPU or CPU*
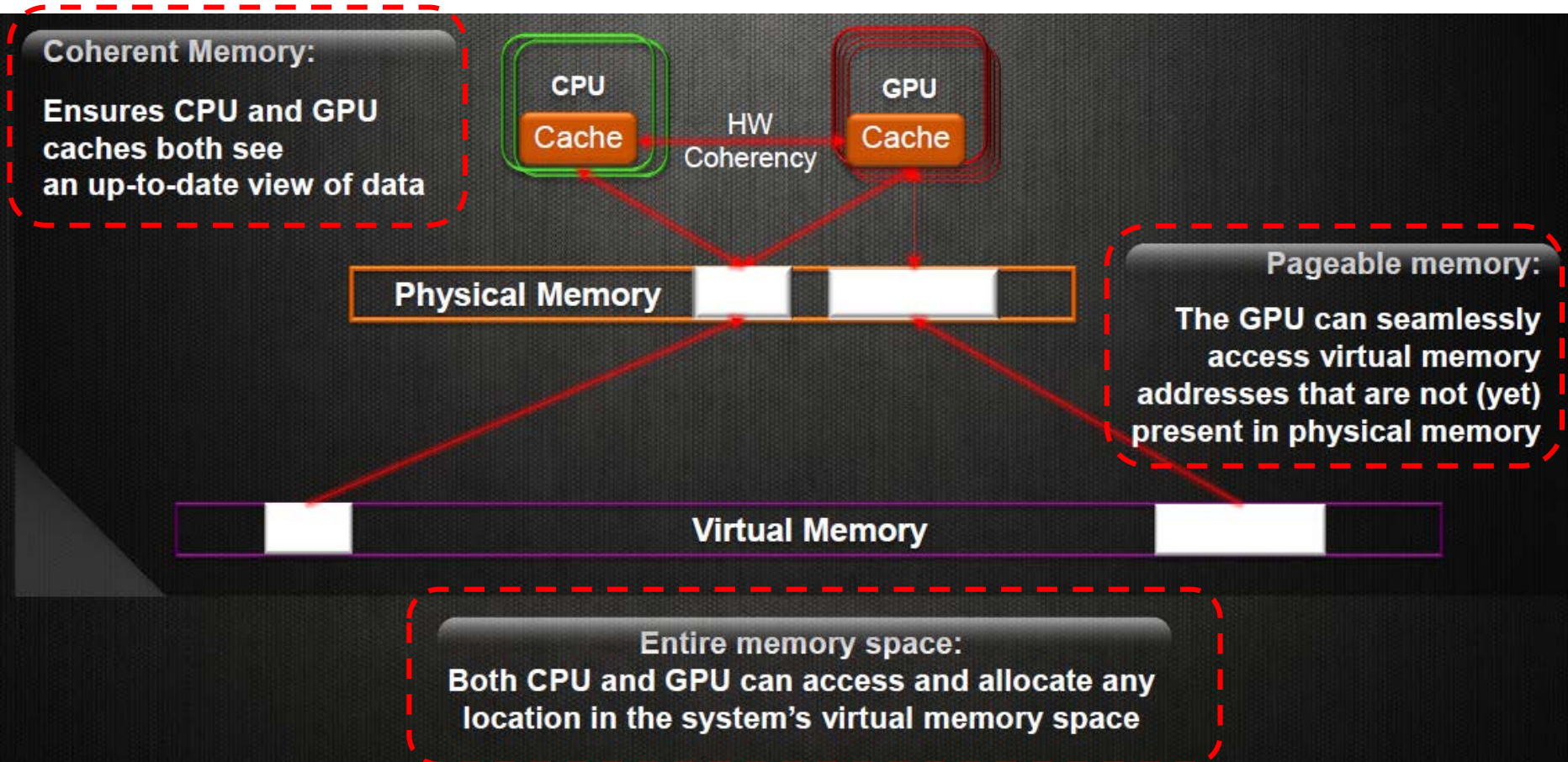
**PAGEABLE MEMORY**
*GPU can take page faults, and is no longer restricted to page locked memory*

**ENTIRE MEMORY SPACE**
*CPU and GPU processes can dynamically allocate memory from the entire memory space*

**Key requirements of hUMA-2** [3]



Coherent Memory:

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

CPU
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

**Hardware support of memory management in Kaveri** [8]

**2.2.2 hQ (heterogeneous Queuing)**

**Traditional management of application tasks queues** [9]

**Application task management with heterogeneous queuing (hQ)** [9]

**Main features of hQ** [9]

- Heterogeneous queuing (hQ) is symmetrical.

    It allows both the CPU and the GPU to generate tasks for themselves and for each other.

- Work is specified in a standard packet format that will be supported by all HSA-compatible hardware, so there's no need for the software to use vendor-specific code.

- Applications can put packets directly into the task queues  that will be accessed by the hardware.

- Each application can have multiple task queues, and a virtualization layer allows HSA hardware to see all the queues.
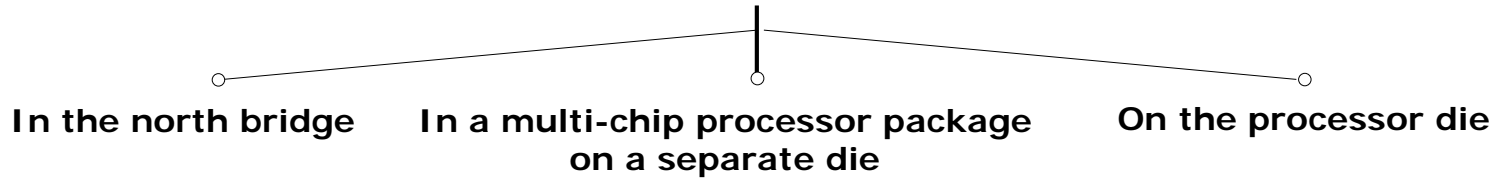
# 3. Intel's approach to support heterogeneous computing
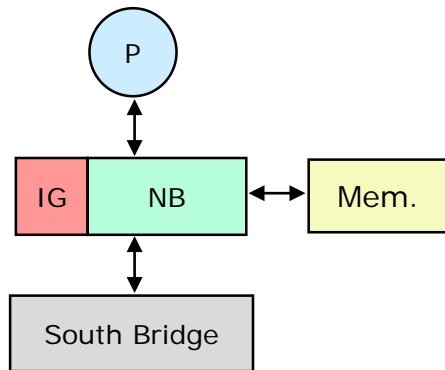
## 3.1 Intel's implementations of integrated graphics

Intel has a long history of integrated graphics, implemented first in the north bridge in the 8xx chipset (1999), as indicated below.
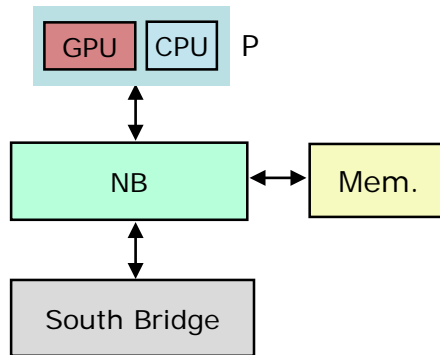
**Implementation of integrated graphics**

**In the north bridge**    **In a multi-chip processor package on a separate die**    **On the processor die**
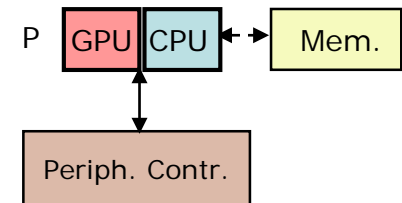
Both the CPU and the GPU are on separate dies
and are mounted into a single package



Implementations about
1999 – 2009 such as
Intel's 828xx chipset (1999)
and subsequent implementations

Intel's Havendale (DT) and
Auburndale (M)
(scheduled for 1H/2009
but cancelled)

Arrandale (DT, 1/2010) and
Clarkdale (M, 1/2010)

Intel's Sandy Bridge (1/2011) and
subsequent processors

AMD's Swift (scheduled for 2009
but canceled)

*AMD's Bobcat-based APUs (M, 1/2011)*
*Llano APUs (DT, 6/2011)*

*and subsequent processors*

**Unified Memory Architecture (UMA)**

Some early graphics cards of other vendors and essentially all (in the north bridge) integrated graphics designs of Intel already made use of the UMA design as early as in the second half of the 1990s, as the next Figure shows for the integrated graphics of the Intel 810 chipset (1999)
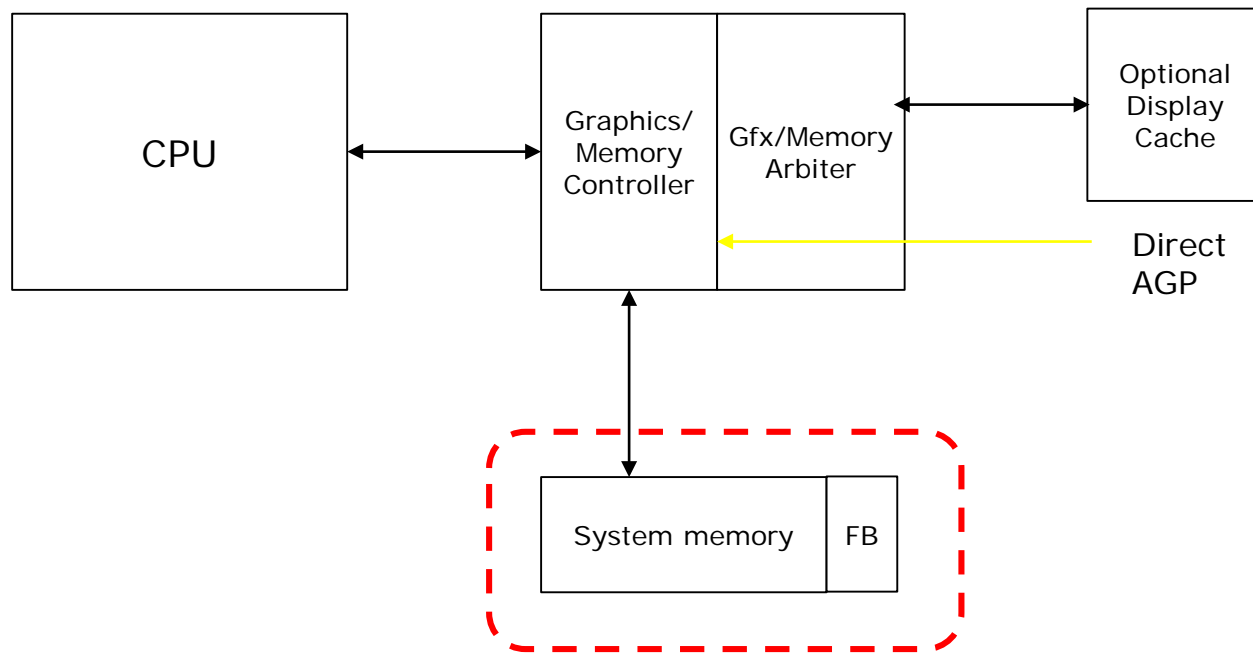


Figure: The UMA design of an early (in the north bridge) integrated graphics [12]

Key benefit/drawback of the UMA design

The UMA design eliminates the need for an extra graphics memory controller and bus but has the constraints of a reduced memory bandwidth.

**Intel's on-die integrated graphics designs**

Subsequently, we discuss only those on-die integrated CPU-GPU solutions that support also HPC i.e. have OpenCL support, as follows.
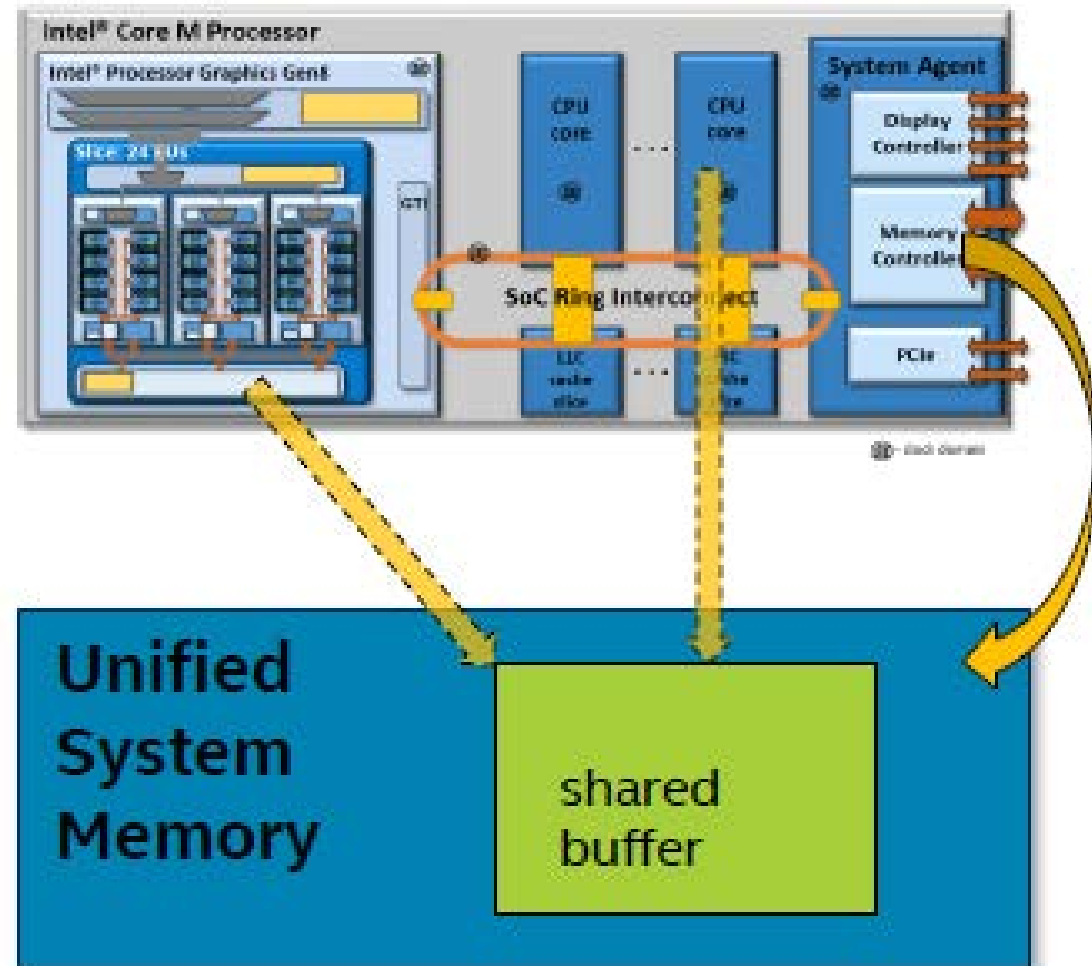
| OpenCL support | Processor Graphics Architecture | Example Processor Product Family Name | Example Graphics SKUs | Codename |
|---|---|---|---|---|
| OpenCL 2.0 | Intel® Processor Graphics **Gen8** | Intel® Core™ M Processor | Intel® HD Graphics 5300 | Broadwell |
| OpenCL 1.2 | Intel Processor Graphics **Gen7.5** | 4th Generation Intel Core Processor | Intel® Iris™ Pro 5200<br>Intel® Iris™ 5100<br>Intel HD Graphics 4600 | Haswell |
| OpenCL 1.2 | Intel Processor Graphics **Gen7** | Intel® Atom™ Processor<br>Intel® Celeron® Processor | Intel HD Graphics | Bay Trail |
| OpenCL 1.2 | Intel Processor Graphics **Gen7** | 3rd Generation Intel Core Processor | Intel HD Graphics 4000 | Ivy Bridge |
| No OpenCL | Intel Processor Graphics **Gen6** | 2nd Generation Intel Core Processor | Intel HD Graphics 3000 | Sandy Bridge |

Table: Intel's processors with on-die integrated GPUs [13]

**Key benefit of OpenCL 1.2 supported  Shared Physical Memory, illustrated on an example**
[13]

• "Zero Copy" CPU & Graphics data sharing

• Enabled by buffer allocation flags in OpenCL™,
  DirectX*, etc.

**Shared Virtual Memory - based on the Broadwell architecture and supported by OpenCL 2.0** [13]

- Significant feature, new in Gen8

- Seamless sharing of pointer rich data-structures in a shared virtual address space

- Hardware-supported byte-level CPU & GPU coherency

- Spec'd Intel® VT-d IOMMU features enable heterogeneous virtual memory, page tables, cache snooping protocols...

- Facilitated by OpenCL™ 2.0 Shared Virtual Memory:

**Key OpenCL 2.0 features** [14]

- **Shared Virtual Memory**
  - Host and device kernels can directly share complex, pointer-containing data structures such as trees and linked lists, providing significant programming flexibility and eliminating costly data transfers between host and devices

- **Dynamic Parallelism**
  - Device kernels can enqueue kernels to the same device with no host interaction, enabling flexible work scheduling paradigms and avoiding the need to transfer execution control and data between the device and host, often significantly offloading host processor bottlenecks

# 3.2 Intel's first implementation of Shared Virtual Memory
## The Core M processor

**3.2 Intel's first implementation of Shared Virtual Memory – The Core M processor**

- The Core M processor targets  tablets and 2-in-1 devices
- Based on the 14 nm Broadwell architecture
- Includes Gen8 graphics
- It is a SOC (System on Chip) design
- Announced: 8/2014

.

**Die layout of the Core M** [13]

**Block diagram of the Core M** [13]

## Compute architecture of Core M [13]