# Concurrency 2: Programming Heterogenous Multicore Systems
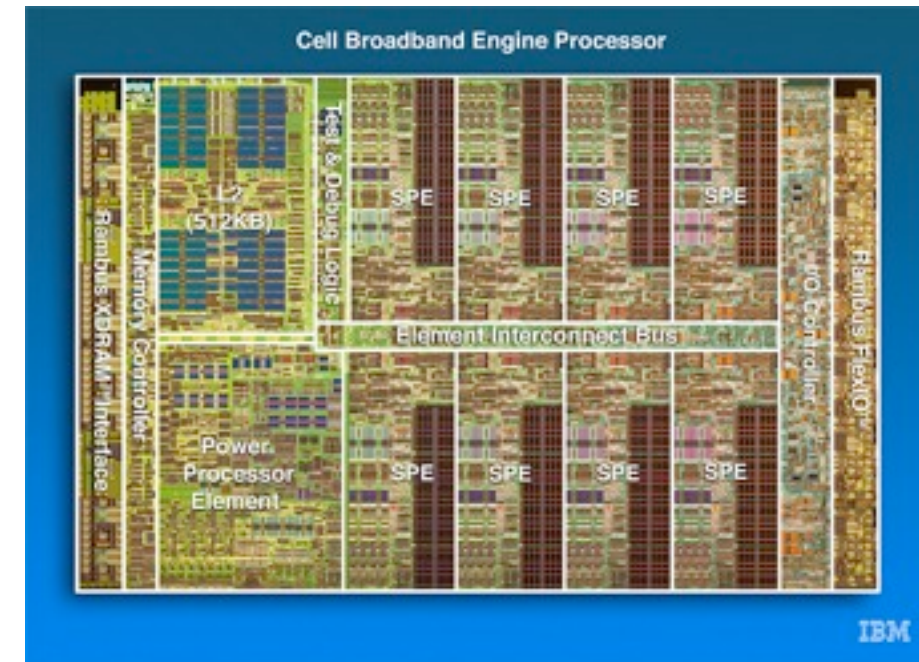
Advanced Operating Systems (M)
Lecture 17

# Lecture Outline

- Heterogeneous instruction set systems

- Programming models

  - Multi-kernel peer model

  - Offload to slave processor

  - Abstraction via virtual machine

- Discussion

  - Hybrid models – Microsoft's Accelerator framework

# Heterogeneous Instruction Set Systems

- Increasingly common for a single system to have cores running different instruction sets

  - CPU + GPU

  - CPU + offload of TCP, crypto, or multimedia functions

  - Cell processor with PPE + multiple SPE

- Desirable when different instruction sets have radically different performance characteristics

  - GPU hardware does simple SIMD-style computations at high speed, but performs very poorly for code with large numbers of conditional branches

  - A typical CPU is better suited for complex conditional code, but performs poorly with SIMD operations

- CPU + GPU model is ubiquitous; others becoming more common



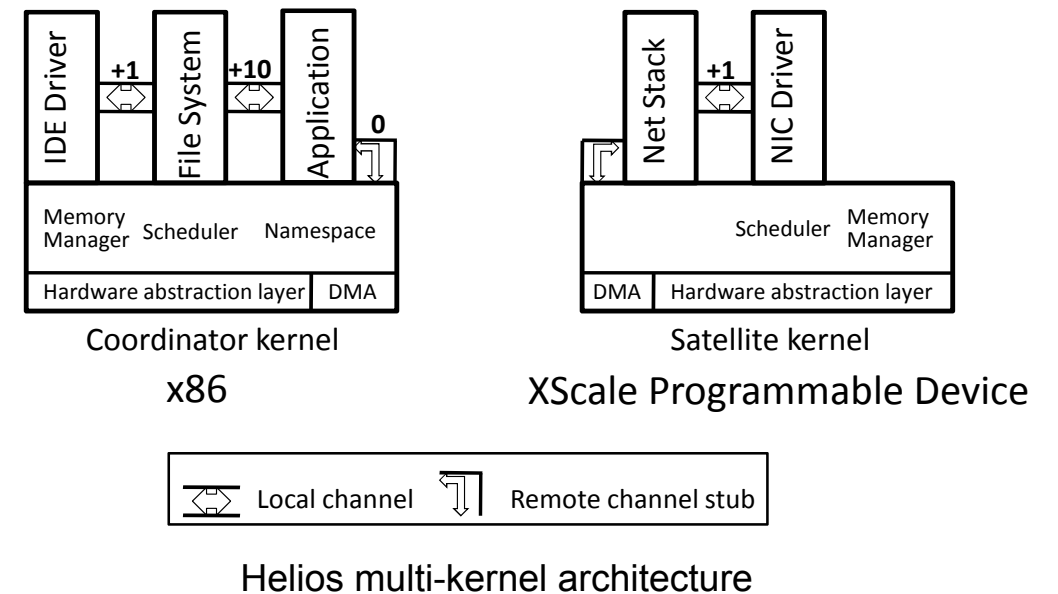Cell Broadband Engine Processor

# Programming Models

- How to program a heterogenous instruction set system?

  - If the cores have radically different characteristics, do they need different programming models and/or languages?

  - Should the cores be peers, or is a master/slave model appropriate?

- Three main alternatives have been explored

  - Multi-kernel model – heterogenous cores

  - Offload to slave processor

  - Abstraction via virtual machines

# Multi-kernel: Heterogenous Cores

- **If cores are full-featured, a multi-kernel model may be appropriate**

  - The multi-kernel model is a distributed system, with message passing – the underlying instruction set is unimportant, but the kernel needs to be recompiled for each architecture

  - Applications are either limited to a subset of the cores, require compilation as fat binaries, or use JIT compilation

    - May not be possible to effectively balance load across the system, due to limitations where certain processes can execute

    - Performance may suffer if related processes can't be co-located due to resource constraints

  - Heavy-weight approach, but offers considerable flexibility

  - Not widely implemented – systems with multiple full-featured cores generally use a homogenous instruction set

# Multi-kernel: Example – Helios

- A research prototype multi-kernel system designed to exploit heterogenous cores

- Multi-kernel extension to Singularity
  - Runs on x86 NUMA systems, and on x86 systems with offload to an ARM processor on a RAID card

- Based on a *satellite kernel* abstraction, allowing weak cores to delegate some work to more full-featured cores
  - All kernels export the same services and message-passing APIs, but some services are implemented by forwarding messages to other cores

  - Applications distributed as JIT compiled byte code; express *affinity* to other processes in metadata to allow dynamic load balancing across cores

  - Good performance on benchmarks, but these only considered a limited set of processes on the ARM core, with clear communication patterns and affinity – unclear how this will work in general with highly asymmetric cores
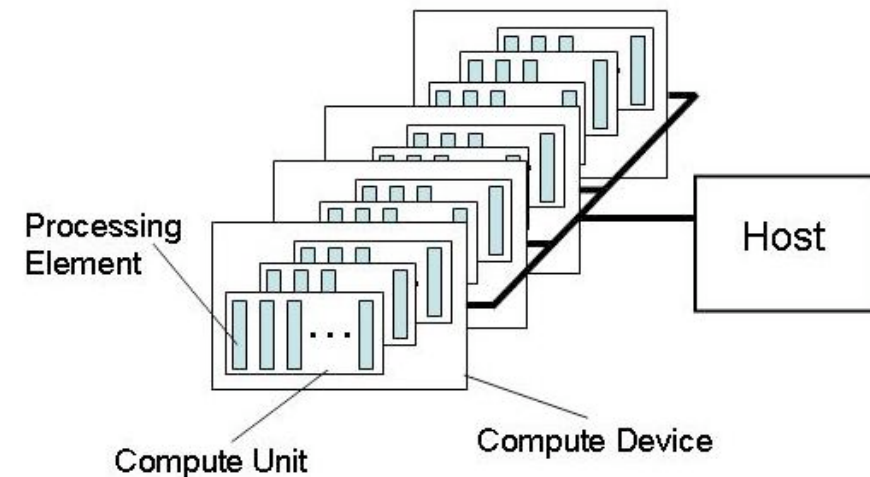


Coordinator kernel
x86

Satellite kernel
XScale Programmable Device

Helios multi-kernel architecture

# Slave Processor

- The system has a master CPU, plus one or more slave processors running a different instruction set

  - The programming model for the slave processors is different to the master

  - Slave processors often too limited to run a full kernel and general-purpose programming language; code for the slave processors written in a special language, and compiled separately (e.g., OpenCL or CUDA)

- Slave processors don't run independently – they're resources invoked by the master CPU

  - Graphics processing
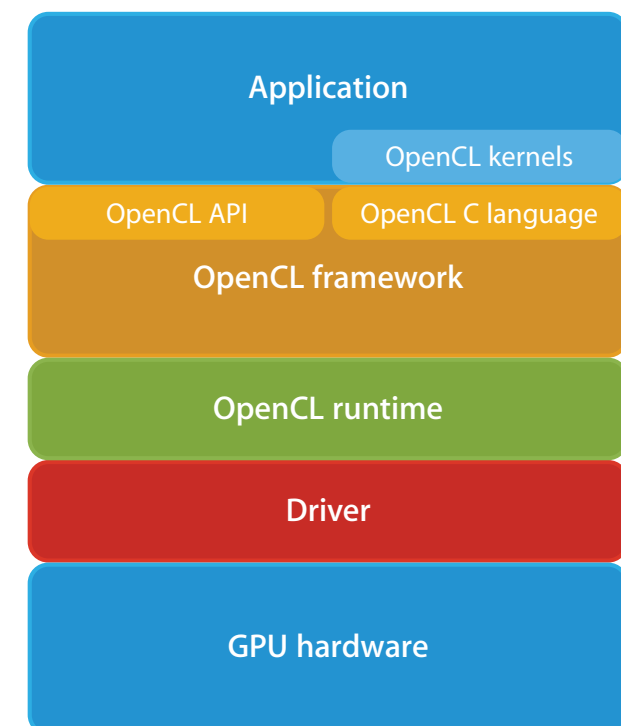
  - TCP or crypto offload

# Slave Processor: Example – OpenCL

- Slave processors are GPU devices – large grid of compute cores designed for SIMD-style array processing

  - 512 cores available on modern GPUs

  - Weak support for conditional branches – the model is that each core runs the same code on different data

- Slave processor code written in OpenCL C

  - Extended subset of C: http://www.khronos.org/opencl/

  - Adds built-in vector, 2D, and 3D image types

  - Adds pointer qualifiers to reference host and GPU memory; use of pointers restricted since memory is not shared between host and device

    - Must explicitly copy inputs and outputs to/from slave device

  - Very restricted standard library

  - Defines the concept of a *kernel function* that can be JIT compiled and executed on a device, and runtime support code to allow device management

  - Runtime creates massive numbers of threads, each running the kernel on different parts of the data



[Source: The OpenCL specification, v1.0]

**Figure 3.1**: *Platform model … one host plus one or more compute devices each with one or more compute units each with one or more processing elements.*

The OpenCL architecture



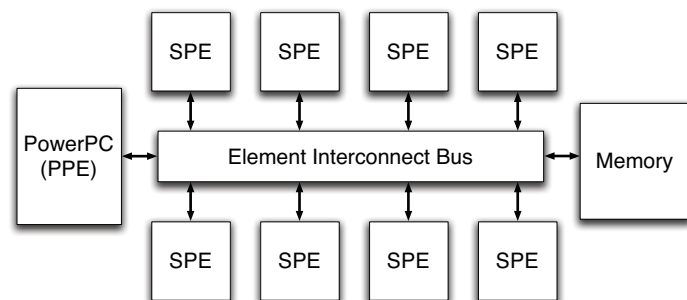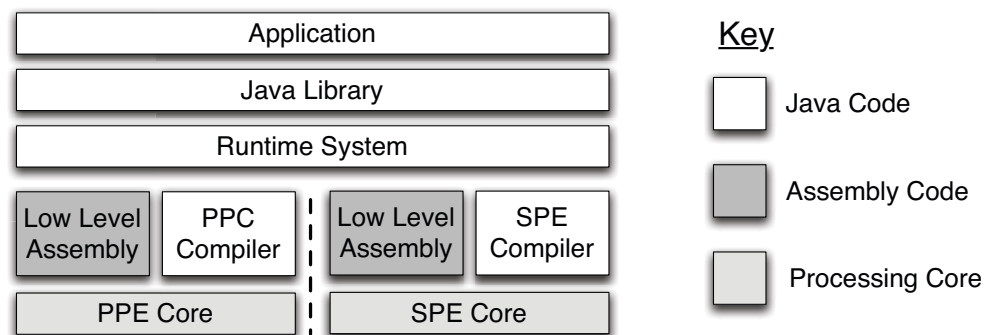[Source: Apple, OpenCL technology brief]

8

# Slave Processor: Example – OpenCL

- ## OpenCL runtime on host system manages offload of OpenCL code to slave devices

  - JIT compilation; host driver code controls exactly what OpenCL functions execute, and when

  - OpenCL devices do not run an OS – they're dumb devices, managed by a device driver

- ## Low-level API and programming model

  - High conceptual burden to use

  - Cannot run general purpose code; programming and communications model is too restricted

  - Does not easily integrate with host applications – too much boilerplate
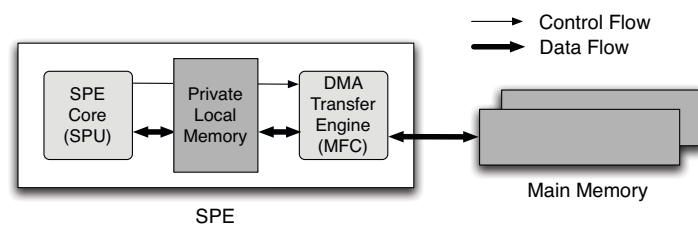
# VM Abstraction

- Use of separate languages and compilation stages for slave processors is complex and error prone

  - OpenCL and CUDA have complex and poorly-defined semantics

  - High cognitive overhead on programmers; difficult to develop and debug

- Alternative: write in a high-level language targeting a virtual machine; let the VM handle the offload

  - E.g., a Java virtual machine that can JIT compile for different cores

  - Pushes complexity onto the VM – simple for application programmer

  - But, high-level languages often not a good fit for slave processors

    - e.g., the JVM has no natural means to express SIMD-style array processing operations, and encourages conditional execution, imperative code, and mutable state – the opposite of what is needed for good GPU code

    - But, a language optimised for GPU processing would perform poorly on a general-purpose CPU, with a small number of cores optimised for imperative code

# VM Abstraction: Example – Hera-JVM



Key
- Java Code
- Assembly Code
- Processing Core



(a) The architecture of the Cell processor.



(b) An SPE core's memory subsystem.

- A JVM for the Cell processor, than can offload methods from PPE to SPE cores

  - JIT compilation; methods compiled for appropriate core based on runtime code placement algorithm

  - Data caching: SPE memory is not cache coherent; data cached on SPE when method starts; cache flushed at synchronisation points, following Java memory model

  - Methods copied to SPE memory in their entirety; migration onto the SPE causes an entire method, and any methods it calls, to run on the SPE

  - Garbage collector understands both architectures, and the caches on the SPEs

  - Hard to decide which methods to migrate to SPE:

    - Explicit annotations (@RunOnSPECore, @RunOnPPECore) work, but place high overhead on programmer

    - Behaviour hints (@ArithmeticCode, @ObjectAccessCode, @LargeWorkingSet) allow the JVM runtime to automatically migrate methods to the SPEs, but are suboptimal

    - Optimal solution is an open problem

# Discussion

- ## Offload to slave processor model is common

  - Hard for programmer, but gives good performance

  - Main kernel treats the GPU as a resource, that can be claimed by a process, and managed as any other resource

- ## Abstraction via virtual machine conceptually clean

  - In principle, allows transparent offload of work from main processor to subordinate processors such as GPUs

  - Difficult in practice: applications written without account for the different processor types and capabilities, and don't aid the runtime; insufficient information for the runtime to effectively offload work – likely inefficient

# Hybrid Virtual Machine/Slave Processor

- Hybrid model: wrap a device-specific programming model in the virtual machine, alongside a general purpose language

    - E.g., for GPU offload, add a data parallel array datatype, then JIT compile operations on those arrays to execute on the GPU cores

- Explicit model of device-specific operations, and control over when they execute

- Virtual machine hides low-level details

# Example: Accelerator

- Extension to C# to provide data-parallel arrays with GPU offload

  - Support operations such as conversion to/from standard arrays, element-wise arithmetic, reductions, transformations, and matrix algebra

  - Data parallel arrays are lazy, and don't compute their value until converted back to a standard array

  - Lazy evaluation helps efficiency: runtime JIT compiles all operations on a single data parallel array at once, and passes to the GPGPU for execution as a single block

- Similar model to OpenCL, except the complexity of managing the GPU is pushed onto the VM

  - Programming model is very similar, and there is similar control over when code executes on the GPU

```
static float[,] Blur(float[,] array, float[] kernel) {
  float[,] result;
  DFPA parallelArray = new DFPA(array);

  FPA resultX = new FPA(0f, parallelArray.Shape);
  for (int i = 0; i < kernel.Length; i++) {
    int[] shiftDir = new int[] { 0, i};
    resultX += PA.Shift(parallelArray, shiftDir) * kernel[i];
  }

  FPA resultY = new FPA(0f, parallelArray.Shape);
  for (int i = 0; i < kernel.Length; i++) {
    int[] shiftDir = new int[] { i, 0 };
    resultY += PA.Shift(resultX, shiftDir) * kernel[i];
  }
  PA.ToArray(resultY, out result);
  parallelArray.Dispose();
  return result;
}
```
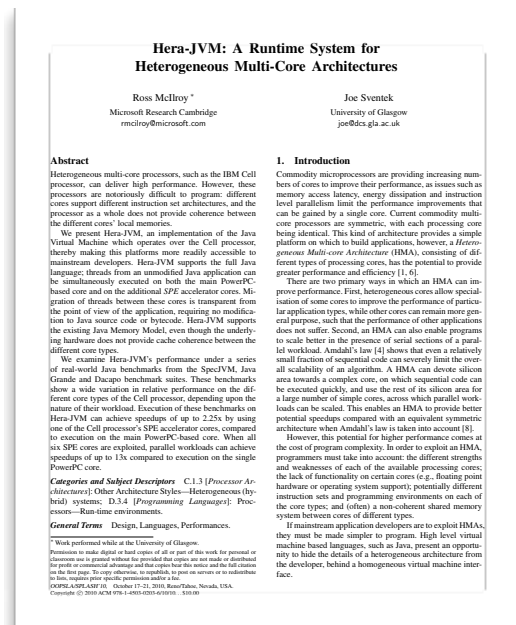
D. Tarditi, S. Puri, and J. Oglesby. Accelerator: using data parallelism to program GPUs for general-purpose use. In Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), San Jose, CA, USA, October 2006. ACM.

# Discussion and Further Reading

- D. Tarditi, S. Puri, and J. Oglesby. Accelerator: using data parallelism to program GPUs for general-purpose use. Proc. International Conference on Architectural Support for Programming Languages and Operating Systems, San Jose, CA, USA, October 2006. DOI: 10.1145/1168857.1168898

- R. McIlroy and J. Sventek, Hera-JVM: A Runtime System for Heterogeneous Multi-Core Architectures, Proc. ACM SIGPLAN Conference on Object-Oriented Programming Systems, Languages, and Applications (OOPSLA), Reno, Nevada, October 2010. DOI: 10.1145/1869459.1869478

- Both to be discussed in tutorial tomorrow (conceptual level – no need to consider performance results)