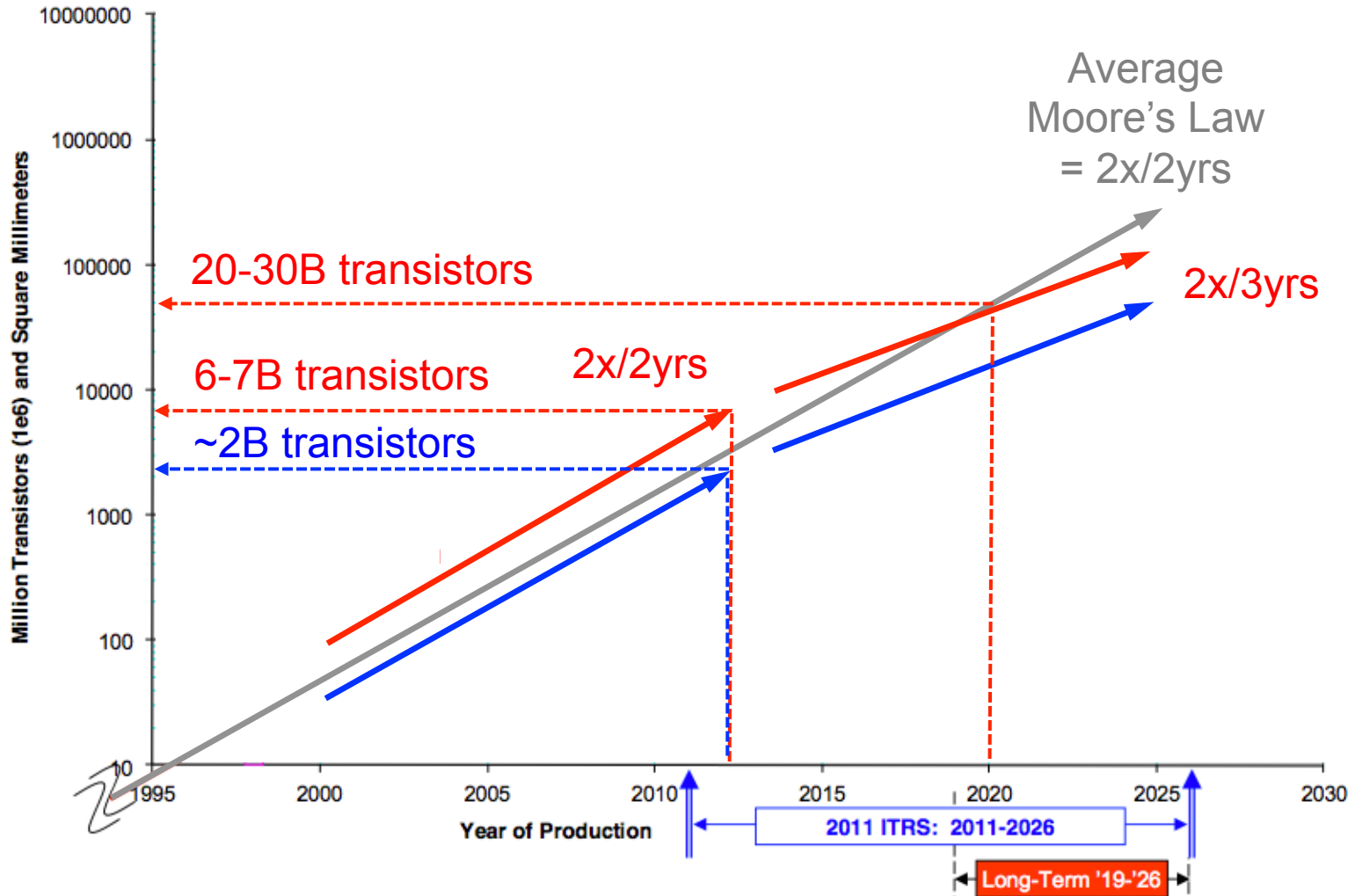# Trends in Heterogeneous Systems Architectures
# (and how they'll affect parallel programming models)

**Simon McIntosh-Smith  simonm@cs.bris.ac.uk**
**Head of Microelectronics Research**
**University of Bristol, UK**

University of
BRISTOL

# Moore's Law today



2011 ITRS - Functions/chip and Chip Size

Average Moore's Law = 2x/2yrs

20-30B transistors

6-7B transistors

~2B transistors

2x/2yrs

2x/3yrs

2011 ITRS: 2011-2026

Long-Term '19-'26

Year of Production

Million Transistors (1e6) and Square Millimeters

# 🔥 Herb Sutter's new outlook

http://herbsutter.com/welcome-to-the-jungle/

"In the twilight of Moore's Law, the transitions to multicore processors, GPU computing, and HaaS cloud computing are not separate trends, but aspects of a single trend – mainstream computers from desktops to 'smartphones' are being permanently transformed into heterogeneous supercomputer clusters. Henceforth, a single compute-intensive application will **need to harness different kinds of cores, in immense numbers**, to get its job done."

"The free lunch is over.
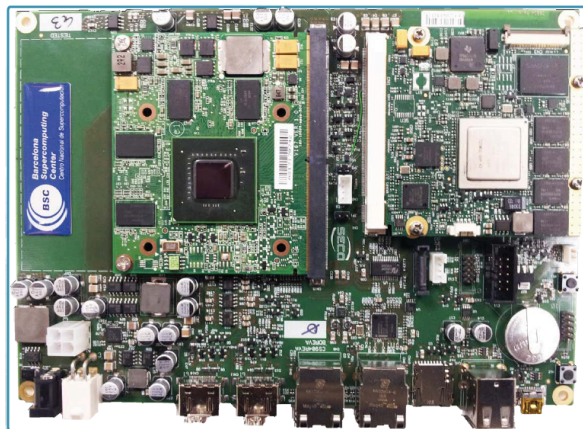Now welcome to the *hardware jungle*."

University of BRISTOL

# 🔥 Four causes of heterogeneity

- Multiple types of programmable core
  - CPU (lightweight, heavyweight)
  - GPU
  - Others (accelerators, …)
- Interconnect asymmetry
- Memory hierarchies
- Software (OS, middleware, tools, …)

University of BRISTOL

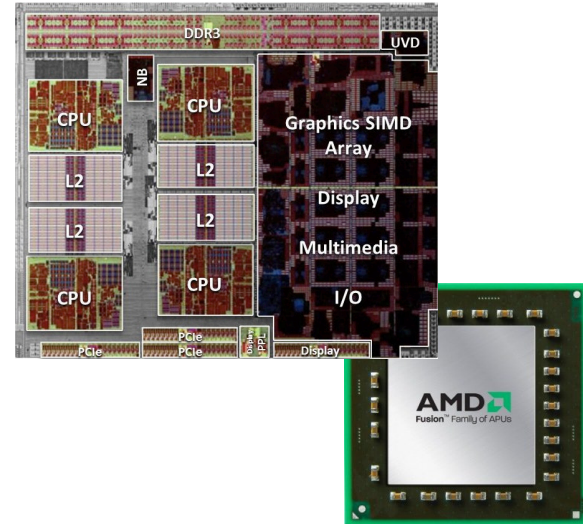# Heterogeneous Systems

AMD Llano Fusion APUs

Intel MIC

FP7 Mont Blanc ARM + GPU

NVIDIA Tegra, Project Denver

University of BRISTOL

# 🍂 Heterogeneity is mainstream

Quad-core ARM Cortex A9 CPU
Quad-core SGX543MP4+ Imagination GPU

Dual-core ARM 1.4GHz, ARMv7s CPU
Triple-core SGX554MP4 Imagination GPU

## Most tablets and smartphones are already powered by heterogeneous processors.

University of
BRISTOL

# Current limitations

- Disjoint view of memory spaces between CPUs and GPUs

- Hard partition between "host" and "devices" in programming models

- Dynamically varying nested parallelism almost impossible to support

- Large overheads in scheduling heterogeneous, parallel tasks

University of BRISTOL

# The emerging Heterogeneous System Architecture (HSA) standard

# Current HSA members

**Founders**

AMD · ARM · Imagination

MEDIATEK · QUALCOMM · SAMSUNG · TEXAS INSTRUMENTS

**Promoters**

LG Electronics

**Supporters**

Arteris · codeplay · FABRIC ENGINE · MULTICORE WARE

**Contributors**

apical · CEVA · DMP · MARVELL · ST life.augmented

Sony Mobile · ST ERICSSON · SONICS · symbio · tensilica · VIVANTE

**Academic**

NTHU Programming Language Lab · University of Illinois Computer Science · THE UNIVERSITY of EDINBURGH School of Informatics · University of BRISTOL

University of BRISTOL

# HSA overview

- The HSA Foundation launched mid 2012
- HSA is a new, ***open*** architecture specification
  - HSAIL virtual (parallel) instruction set
  - HSA memory model
  - HSA dispatcher and run-time

- Provides an optimised platform architecture for heterogeneous programming models such as OpenCL, C++AMP, et al

University of BRISTOL

# HSA overview

# Enabling more efficient heterogeneous programming

- Unified virtual address space for all cores
  - CPU and GPU
  - Enables PGAS-style distributed arrays
- Hardware queues per core with lightweight user mode task dispatch
  - Enables GPU context switching, preemption, efficient heterogeneous scheduling
- First class barrier objects
  - Aids parallel program composability

University of BRISTOL

# HSA Intermediate Layer (HSAIL)

- Virtual ISA for parallel programs
- Similar to LLVM IR and OpenCL SPIR
- *Finalised* to specific ISA by a JIT compiler
- Make late decisions on which core should run a task
- HSAIL features:
  - Explicitly parallel
  - Support for exceptions, virtual functions and other high-level features
  - Syscall methods (I/O, printf etc.)
  - Debugging support

University of BRISTOL

# HSA memory model

- Compatible with C++11, OpenCL, Java and .NET memory models

- Relaxed consistency

- Designed to support both managed language (such as Java) and unmanaged languages (such as C)

- Will make it much easier to develop 3rd party compilers for a wide range of heterogeneous products

  - E.g. Fortran, C++, C++AMP, Java et al

University of BRISTOL

# HSA dispatch

- HSA designed to enable heterogeneous task queuing
  - A work queue per core (CPU, GPU, …)
  - Distribution of work into queues
  - Load balancing by work stealing
- Any core can schedule work for any other, including itself
- Significant reduction in overhead of scheduling work for a core

University of
BRISTOL

# Today's Command and Dispatch Flow

# Today's Command and Dispatch Flow

# Today's Command and Dispatch Flow

# HSA enabled dispatch

# HSA roadmap from AMD



**HETEROGENEOUS SYSTEM ARCHITECTURE ROADMAP**

| 2011 Physical Integration | 2012 Optimized Platforms | 2013 Architectural Integration | 2014 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 |

Fusion¹² AMD DEVELOPER SUMMIT

University of BRISTOL

**20**

# Open Source software stack for HSA

A Linux execution and compilation stack will be open-sourced by AMD

- Jump start the ecosystem
- Allow a single shared implementation where appropriate
- Enable university research in all areas

| Component Name | Purpose |
| --- | --- |
| HSA Bolt Library | Enable understanding and debug |
| OpenCL HSAIL Code Generator | Enable research |
| LLVM Contributions | Industry and academic collaboration |
| HSA Assembler | Enable understanding and debug |
| HSA Runtime | Standardize on a single runtime |
| HSA Finalizer | Enable research and debug |
| HSA Kernel Driver | For inclusion in Linux distros |

University of BRISTOL

# HSA should enable nested parallel programs like this

Support for multiple algorithms, even within a single application

Task farms, pipeline, data parallelism, …



University of BRISTOL

# Conclusions

- Heterogeneity is an increasingly important trend
- The market is finally starting to create and adopt the necessary open standards
  - Proprietary models likely to start declining now
  - Don't get locked into any one vendor!
- Parallel programming models are likely to (re)proliferate
- HSA should enable much more dynamically heterogeneous nested parallel programs and programming models

University of BRISTOL

# www.cs.bris.ac.uk/Research/Micro