# AMD IN HPC

**André Heidekrueger**
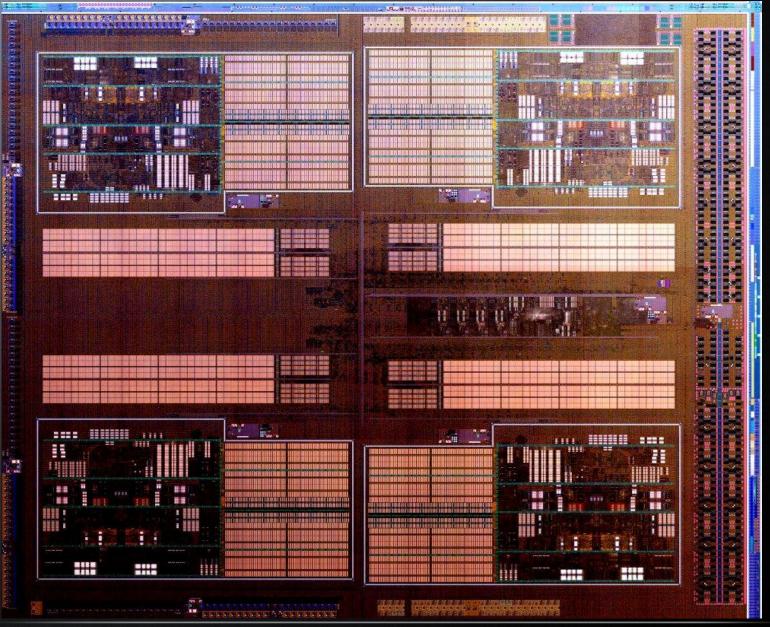**HPC Presales Engineer**
**AMD EMEA**

*November 2011*
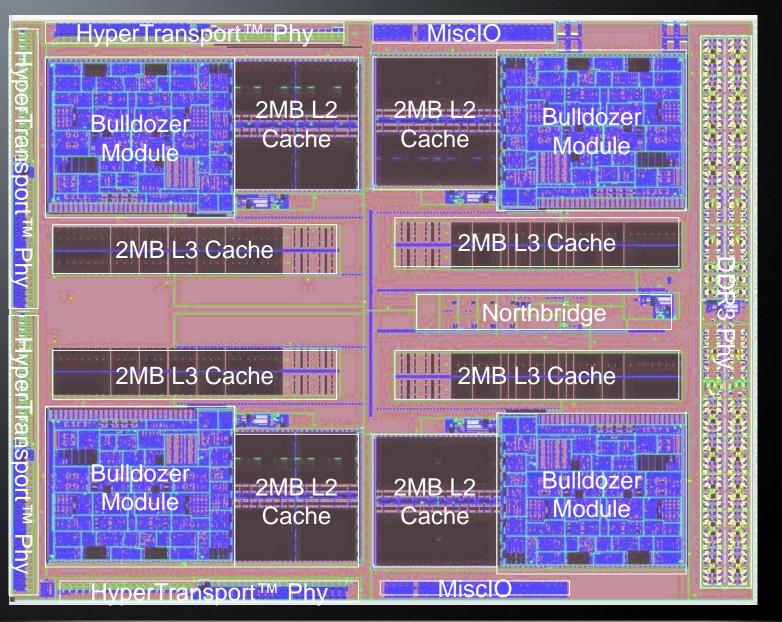
# AMD IN TOP 500 – JUNE 2011 LIST

- **22 of the 50** fastest supercomputers on the TOP500 list are using AMD's CPU and/or GPU technology.

- **66 systems** on the TOP500 are powered by AMD Opteron™ processors, and some new ones coming soon

# THE DIE | *Floorplan (315 mm$^2$)*

AMD

# COMPUTE UNIT - OPTIMIZED PERFORMANCE/WATT

## Leadership Multi-Threaded Micro-Architecture

### Full Performance From Each Core

- Dedicated execution units per core
  - Up to 128 Instructions/Core in flight/Cycle/Core
- No shared execution units as with SMT

### High Frequency / Low-Power Design

- Core Performance Boost
  - "Boosts" frequency of cores when available power allows
  - No idle core requirement
- Power efficiency enhancements
  - Significantly reduced leakage power
    - 32nm, Hi-K metal gate
  - More aggressive dynamic power mgt
    - Memory Power Capping
    - Northbridge P-States

### Virtualization Enhancements

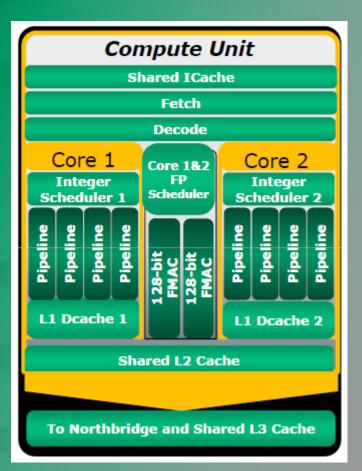- Faster switching between VMs
  - TLB Flush by ASID, VM State Caching, Decode Assists
- AMD-V extended migration support

### Shared Double-sized FPU

- Amortizes very powerful 256-bit unit across both cores
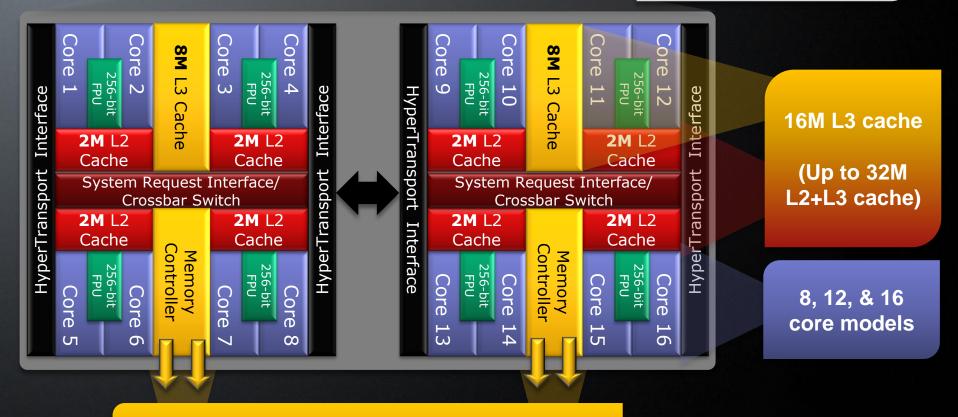  - Includes FMAC

### Improved IPC

- Micro-architecture and ISA enhancements
  - SSE4.1/4.2, AVX 1.0/1.1, SSSE3, AES, LWP, XOP, FMA4
  - Independent Core/FP Schedulers, Deeper Buffers, Larger Caches, Branch Fusion, Aggressive Pre-fetching, High Degree of Thread Concurrency Throughout

**Compute Unit**

| Shared ICache |
| Fetch |
| Decode |

**Core 1** — Integer Scheduler 1 — Pipeline / Pipeline / Pipeline / Pipeline — L1 Dcache 1

**Core 1&2 FP Scheduler** — 128-bit FMAC / 128-bit FMAC

**Core 2** — Integer Scheduler 2 — Pipeline / Pipeline / Pipeline / Pipeline — L1 Dcache 2

**Shared L2 Cache**

**To Northbridge and Shared L3 Cache**

AMD

# AMD OPTERON™ 6200 SERIES PROCESSOR ("INTERLAGOS")

**Multi- Chip Module (MCM) Package**

**Same platform as AMD Opteron™ 6100 Series processor.**

HyperTransport Interface

Core 1 · Core 2 · 256-bit FPU
**8M** L3 Cache
Core 3 · Core 4 · 256-bit FPU

**2M** L2 Cache · **2M** L2 Cache

System Request Interface/ Crossbar Switch

**2M** L2 Cache · **2M** L2 Cache

Core 5 · Core 6 · 256-bit FPU
Memory Controller
Core 7 · Core 8 · 256-bit FPU

HyperTransport Interface

HyperTransport Interface

Core 9 · Core 10 · 256-bit FPU
**8M** L3 Cache
Core 11 · Core 12 · 256-bit FPU

**2M** L2 Cache · **2M** L2 Cache

System Request Interface/ Crossbar Switch

**2M** L2 Cache · **2M** L2 Cache

Core 13 · Core 14 · 256-bit FPU
Memory Controller
Core 15 · Core 16 · 256-bit FPU

HyperTransport Interface

**16M L3 cache**

**(Up to 32M L2+L3 cache)**

**8, 12, & 16 core models**

**4 DDR3 memory channels supporting LRDIMM, ULV-DIMM, UDIMM, & RDIMM**

**AMD**

# A BROAD PORTFOLIO OF "INTERLAGOS" PRODUCTS



Four Core AMD
Opteron™ 6200
Series Processors

Huge memory
bandwidth per core

Eight Core AMD
Opteron™ 6200
Series Processors

Blend of cores and
memory with an
edge on memory
bandwidth

Twelve Core AMD
Opteron™ 6200
Series Processors

Blend of cores and
memory with a little
more computation
muscle

Sixteen Core AMD
Opteron™ 6200
Series Processors

The industry's only
16-core x86 processor
for massive thread
density

AMD

# FLEX FP: MORE FLEXIBLE TECHNICAL PROCESSING

*More performance and new instruction support*

| Fetch |
|---|
| Decode |

Integer Scheduler ← → Integer Scheduler

FP Scheduler

Integer Core | MMX | MMX | 128-bit FMAC | 128-bit FMAC | Integer Core

Shared L2 Cache

---

**Runs SSE and AVX simultaneously**

- No processing penalty for simultaneous execution

**Executes two SSE or AVX (128-bit) instructions simultaneously or one AVX (256-bit) instruction per Bulldozer module**

- Wider range of FP processing capabilities than competition

**Processes calculations in a single cycle using FMA4* and XOP instructions**

- Executes more instructions in fewer cycles than competition

**Uses dedicated floating point scheduler**

- No waiting on integer scheduler to run instructions
- Designed to be always available to schedule floating point operations

---

*FMAC can execute an FMA4 execution (a=b+c*d) in one cycle vs. 2 cycles that would be required for FMA3 or standard SSE floating point calculation.

AMD

# FLEX FP

| | AMD Opteron™ 6100 Series FPU | Intel Sandy Bridge | AMD Opteron™ 6200 Series FlexFP |
|---|---|---|---|
| Execute 128-bit FP | ✓ | ✓ | ✓ |
| Support SSSE3, SSE4.1, SSE4.2 | | ✓ | ✓ |
| Execute 128-bit AVX | | ✓ | ✓ |
| Execute 256-bit AVX | | ✓ | ✓ |
| Execute two 128-bit SSE or AVX ADD instructions in 1 cycle | | | ✓ |
| Execute two 128-bit SSE or AVX MUL instructions in 1 cycle | | | ✓ |
| Switch between SSE and AVX instructions without penalty | | | ✓ |
| Execute FMA operations (A=B+C*D) | | | ✓ |
| Supports XOP | | | ✓ |
| FLOPs per cycle (128-bit FP) | 48 | 32 | 64 |
| FLOPS per cycle (128-bit AVX) | - | 32 | 64 |
| FLOPS per cycle (256-bit AVX) | - | 64 | 64 |

FP Scheduler

128-bit FMAC   128-bit FMAC

128-Bit AVX   128-Bit AVX

FP Scheduler

128-bit FMAC   128-bit FMAC

256-Bit AVX

**Two 128-bit FMACs shared per module, allowing for dedicated 128-bit execution per core or shared 256-bit execution per module**

Sandybridge information from
http://software.intel.com/en-us/avx/

AMD

# NEW "BULLDOZER" INSTRUCTIONS

| Instructions | Applications/Use Cases |
|---|---|
| **SSSE3, SSE4.1, SSE4.2** (AMD and Intel) | • Video encoding and transcoding<br>• Biometrics algorithms<br>• Text-intensive applications |
| **AESNI PCLMULQDQ** (AMD and Intel) | • Application using AES encryption<br>• Secure network transactions<br>• Disk encryption (MSFT BitLocker)<br>• Database encryption (Oracle)<br>• Cloud security |
| **AVX** (AMD and Intel) | Floating point intensive applications:<br>• Signal processing / Seismic<br>• Multimedia<br>• Scientific simulations<br>• Financial analytics<br>• 3D modeling |
| **FMA4** (AMD Unique) | HPC applications |
| **XOP** (AMD Unique) | • Numeric applications<br>• Multimedia applications<br>• Algorithms used for audio/radio |

**Software that currently supports SSSE3, SSE4.1, SSE4.2, AESNI, and AVX should run on "Bulldozer."**

- No recompile of code if the software only checks ISA feature bits
- Recompile needed if software also checks for processor type

**For FMA4 or XOP, software will need to be written to call specific instructions or be compiled with a compiler that will automatically generate code that leverages these instructions.**

AMD

# COMPILER SUPPORT
## AMD Opteron™ 6200 Series Processors

| Compiler | Status | Support for SSSE3, SSE4.1/4.2, AVX (Intel and AMD) | Support for FMA4, XOP (AMD) | Options to auto generate code with new instruction support |
|---|---|---|---|---|
| **Microsoft Visual Studio 2010 SP1** | **Available** | Yes | Yes | No |
| **GCC 4.5, 4.6** | **Available** | Yes | Yes | Yes |
| **Open64 4.2.5** | **Available** | Yes | Yes | Yes |
| **PGI 11.6, 11.7** | **Available** | Yes | No | Yes (for SSSE3, SSE4.1/4.2, AVX) |
| **PGI 11.9** | **In Development** | Yes | Yes | Yes |
| **ICC 12\*** | **Available** | Yes (but ICC runtime fails on AMD processors) | No | Yes (but ICC runtime fails on AMD processors) |

*Intel has an –mAVX flag which is designed to run on any x86 processor; however, the ICC runtime makes assumptions about cache line sizes and other parameters that causes code not to run on AMD processors

AMD

# *ACML SUPPORT | A CLOSER LOOK*
## AMD Opteron™ 6200 Series Processors

| | **Linear Algebra** | **Fast Fourier Transforms (FFT)** | **Others** | **Compiler Support** |
|---|---|---|---|---|
| **ACML 5.0 (Aug 2011)** *Alpha available NOW* | • SGEMM (single precision)<br>• DGEMM (double precision)<br>• L1 BLAS | • Complex-to-Complex (C-C) single precision FFTs | • Random Number Generators<br>• AVX compiler switch for Fortran | • Alpha support for gcc 4.6 and Open64 4.2.5<br>• PGI 11.8 and ICC 12 added with ACML 5.0 production release<br>• Cray to begin deployment of ACML with their compiler with ACML 5.0 |
| **ACML 5.1 (Dec 2011)** | • CGEMM (complex single decision)<br>• ZGEMM (complex double precision) | • Real-to-complex (R-C) single precision FFTs<br>• Double precision C-C and R-C FFTs | | All compilers listed for ACLM 5.0 will be supported |

For additional information on ACML, go to:
http://developer.amd.com/libraries/acml/pages/default.aspx

AMD⤇

# *UP TO 33% MEMORY THROUGHPUT INCREASE\**

AMD Opteron™ 6200 Series Processors



- New redesigned Northbridge controller
- 1600 MHz DDR-3 support
- LR-DIMM support
- 1.25V LV-DDR3 support
- New memory power management features:
    - Aggressive power down
    - Partial channel power down
    - And memory power capping

*Based on measurements by AMD labs as of 8/9/11.  Comparison is AMD Opteron 6200 Series with DDR3-1600 vs. AMD Opteron 6100 Series with DDR3-1333. See substantiation section for config info.

**AMD**

# *POWER EFFICIENCY OVERVIEW*

## Consistent Power and Thermals



AMD Opteron™ 6200 Series fits into the same general thermal range as previous generation

## Reduces processor power at idle by up to 46%*



### C6 power state

Shuts down clocks and power to idle cores

## Enables more control for IT



### TDP Power Cap

Set thermal design power (TDP) to meet power and workload demands for more flexible, denser deployments

## Intelligent Circuit Design



Minimizes the number of active transistors for lower power and better performance

*See processor power savings slide in substantiation section

**AMD**

**Core P-states specify multiple frequency and voltage points of operation**

– Higher frequency P-states deliver greater performance but require higher voltage and thus more power

– The hardware and operating system vary which P-state a core is in to deliver performance as needed, but use lower frequency P-states to save power whenever possible

**AMD Turbo CORE:  when the processor is below its power/thermal limits the frequency and voltage can be boosted above the normal maximum and stay there until it gets back to the power/thermal limits**



FREQUENCY

POWER

VOLTAGE

AMD

# *OVERVIEW OF P-STATES, APM, AND LIMITS*

Each P-state has a specific voltage and frequency

Higher performance

Lower power

C-state boost → **Pb0**

All cores boost → **Pb1**

**P2**

Example P-state limit (P3)

**P3**

**P4**

If the TDP limit is set below the P-state limit, APM controller dithers between P-states to maximize performance

**P5**

**P6**

**P7**

HW View

**P0**

**P1**

**P2**

**P3**

**P4**

**P5**

SW View

← Boost P-states

← Base P-state

Example TDP limit (higher resolution) set between P5 and P6

**AMD**

# AMD TURBO CORE TECHNOLOGY

**Base frequency with TDP headroom**

**All core boost activated (300 to 500MHz)**

**Max turbo activated (up to 1GHz, half cores)**

Mhz potential

Up to 500Mhz

Up to 1 Ghz

| | |
|---|---|
| **All Core Boost** | When there is TDP headroom in a given workload, AMD Turbo CORE technology is automatically activated and can increase clock speeds by 300 to 500 MHz across **all cores.** |
| **Max Turbo Boost** | When a lightly threaded workload sends half the Bulldozer modules into C6 sleep state but also requests max performance, AMD Turbo CORE technology can increase clock speeds by up to 1 GHz across **half the cores.** |

AMD

## Power Capping Power Thresholds[1]

Can allow the user to set the maximum processor power ceiling via BIOS[2] or APML[3].

### What's the Benefit?

- More control over power settings
- Flexibility to set power limits without capping CPU frequencies[4]

[1]Planned support in "Bulldozer" processors
[2]For platforms where TDP power capping feature is enabled in the system BIOS
[3]For platforms that have designed in APML platform support
[4]TDP power capping can still allow the processor to operate a maximum specified frequency

**AMD**

# REDUCING POWER LEAKAGE
## ENHANCED NEAR ZERO POWER CORE STATE WITH "C6"

## AMD Opteron 6100 & 4100 Series Processors

Single power plane, all cores powered at all times



VDD Core

Voltage is reduced but still applied to cores resulting in leakage / static power

## 'Interlagos' & 'Valencia' Processors

Single power plane, but each module can be turned on and off



VDD Core

Voltage is gated off to virtually remove all core static power / leakage

AMD

**Today**  **"Bulldozer"**

Decreasing Level of Power Consumption →

**Active** — All cores running workloads; core/module frequency can run independently to save power

**Idle** — No cores running workloads; core/module frequency reduced to 800MHz to save more power

L2  **Smart Fetch**  L2 — After a set idle time L2 cache is flushed to L3, allowing cores to 'sleep' to save power while maintaining MP coherency

**C6[1] (NEW!)** — On "Bulldozer" any idle module can independently enter 'C6', gating power up to 95% for considerable power savings; module state is saved to DRAM

**C1e** — In addition to reducing memory and I/O power, 'C6' further reduces core power on "Bulldozer" vs. current core in C1e

[1]Planned support in "Bulldozer" processors

AMD

# GENERATIONAL COMPARISONS

| | AMD Opteron™ 6100 Series Processors | AMD Opteron™ 6200 Series Processors |
|---|---|---|
| Cores | 8 or 12 core | 4, 8, 12 or 16 core |
| Cache (L2 per core / L3 per die) | 512KB / 6MB | 2MB (shared between 2 cores) / 8MB |
| Memory Channels and speed | four; up to 1333MHz | four; up to 1600MHz |
| Floating point capability | 128-bit FPU per core (FADD/FMUL) | 128-bit dedicated FMAC per core or 256-bit AVX shared between 2 cores |
| Integer Issues Per Cycle | 3 | 4 |
| Turbo CORE Technology | No | Yes (+500MHz with all cores active) |
| Power (ACP) | 65W, 80W, 105W | TBD (planned 65W, 80W, 105W) |
| New Instruction Sets | | SSSE3, SSE 4.1/4.2, AVX, AES, FMA4, XOP, PCLMULQDQ |
| Power Gating | AMD CoolCore™, C1E | AMD CoolCore™, C1E, C6 |
| Process / Die Size | 45nm SOI | 32nm SOI (smaller overall die size) |
| Performance | | 35% higher processing throughput* |

The above reflect current expectations regarding features and performance and is subject to change.

*Based on pre-production measurements by AMD labs as of 8/9/11 and comparison of top bin AMD Opteron™ 6200 Series processors to top bin AMD Opteron 6100 Series processors

AMD

# PERFORMANCE EFFICIENCY TODAY: AMD'S HPC PRODUCT PORTFOLIO

Energy efficient CPU and discrete GPU processors focused on addressing the most demanding HPC workloads

**Multi-core x86 Processors**

- **Outstanding Performance**
- **Superior Scalability**
- **Enhanced Power Efficiency**

**ATI FirePro™ Professional Graphics**

- **3D Accelerators For Visualization**
- **Full support for GPU computation with OpenCL**

**AMD FireStream™ GPU Accelerators**

- **Optimized for server integration**
- **Single-slot and dual-slot form factors**
- **Industry standard OpenCL SDK**

AMD **Accelerated** Parallel Processing TECHNOLOGY

AMD

# A CASE FOR SERVER FUSION - EXASCALE

- Current trajectory puts traditional x86 computing at just over 20Pflops by 2018

- A data center that could achieve an exaflop in 2018 using only x86 processors would consume over 3TW

- To achieve exascale capability by 2018, x86 performance would need to increase by 2x each year, starting in 2010



EXAFLOP

Heterogeneous compute required to bridge the gap

Homogeneous x86 compute hits a wall

AMD

# HIGH EFFICIENCY LINPACK IMPLEMENTATION ON AMD MAGNY COURS + AMD 5870 GPU

**System GFLOPS**



- GPU DPFP Peak: 544 GFLOPS
- GPU DGEMM kernel: 87% of Peak

2.5 GFLOPS/W

HPL code:

http://code.compeng.uni-frankfurt.de/

- Node DPFP Peak: 745.6 GFLOPS
- Linpack efficiency: 75.5% of Peak
- Linpack scaling across 4 nodes: 70% of Peak

**HIC** for **FAIR**
Helmholtz International Center

GOETHE UNIVERSITÄT
FRANKFURT AM MAIN

**FIAS** Frankfurt Institute for Advanced Studies

# WORLD'S FIRST HPC CLUSTER BASED ON FUSION APUS

Penguin Computing has successfully installed the world's first HPC cluster powered by AMD accelerated processing units (APUs) at Sandia National Labs in Albuquerque, New Mexico.

The Altus A2A00 system comprises 104 servers powered by A-series Fusion Llano APUs (one chip per server) with four x86 cores and 320/400 stream processors that are interconnected through a QDR Infiniband fabric. It delivers a theoretical peak performance of **59.6TFLOPs**. The Altus 2A00 was specifically designed by Penguin Computing, in partnership with AMD, to support the AMD Fusion APU architecture. It is the world's first Fusion APU system in a rack mountable chassis in a 2U form factor.

Penguin Altus 2A00
Leverage AMD's APU architecture for HPC

AMD

- APU: Fusion of CPU & GPU compute power within one processor

- High-bandwidth I/O

AMD

# GRAPHICS AND MEDIA PROCESSING EFFICIENCY IMPROVEMENTS

## 2010 IGP-based Platform



**CPU Chip**

~17 GB/sec  ~17 GB/sec

CPU Cores

UNB

MC

DDR3 DIMM Memory

~7 GB/sec

GPU   UVD

SB Functions

PCIe

*Graphics requires memory bandwidth to bring full capabilities to life*

**Bandwidth pinch points and latency hold back the GPU capabilities**

## 2011 APU-based Platform



**APU Chip**

CPU Cores

UVD

GPU

UNB / MC

DDR3 DIMM Memory

~27 GB/sec

~27 GB/sec

PCIe

- 3X bandwidth between GPU and memory
- Even the same sized GPU is substantially more effective in this configuration
- Eliminate latency and power associated with the extra chip crossing
- Substantially smaller physical foot print

AMD

# ACCELERATED DATA PARALLEL PROCESSING CAPABILITIES

- APU bandwidth enhancements not only improve traditional graphics, but also data parallel compute effectiveness

- Both GPUs cooperate on graphics and compute

- OpenCL 1.1 and DirectCompute compliant in both the APU and the optional discrete GPU

- AMD Fusion Experience Fund is helping to fuel the application developer community



**APU Chip**

CPU Cores

UVD

GPU

UNB / MC

DDR3 DIMM Memory

~27 GB/sec

~27 GB/sec

>100 GB/sec

**Discrete GPU Card/Chip**
(optional)

GPU

MC

UVD

Discrete GPU GDDR Memory

PCIe Gen2 (x16)
~7 GB/sec

AMD

# INTRODUCING OPENCL™

*The open standard for parallel programming across heterogeneous processors*

AMD

## IT'S A HETEROGENEOUS WORLD

- Heterogeneous computing
  - The new normal

- Many CPU's – 2, 4, 8, …

- Very many GPU processing elements – 100's

- Different vendors, configurations, architectures

- The multi-million dollar question
  - How do you avoid developing and maintaining different source code versions?



CPU

System Memory

Fusion GPU …

GPU Memory

Discrete GPU …

AMD

- Industry Standard
- Open Standard
- Cross Platform
- Multi-Vendor

- Royalty Free
- Broad ISV Support

**CPUs**
**Multiple cores driving performance increases**

**Emerging Intersection**

**GPUs**
**Increasingly general purpose data-parallel computing**

**OpenCL**

**Multi-processor programming – e.g. OpenMP**

**Heterogeneous Computing**

**Graphics APIs and Shading Languages**

**OpenCL™ is a programming framework for heterogeneous compute resources**

Source Khronos

AMD

- Initial proposal made by Apple
  - A broad diversity of industry perspectives
  - Processor vendors, application developers, system OEMs, tool vendors, …



**Source Khronos**

Context

Programs

Kernels

Memory Objects

Command Queues

```
__kernel void
dp_mul(__global const float *a,
       __global const float *b,
       __global float *c)
{
    int id = get_global_id(0);
    c[id] = a[id] * b[id];
}
```

dp_mul
CPU program
binary

dp_mul
GPU program
binary

dp_mul

arg [0] value

arg [1] value

arg [2] value

Images

Buffers

In Order Queue

Out Order Queue

**Compile code**

**Create data & arguments**

**Send to execution**

```
// create the OpenCL context on a GPU device
cl_context = clCreateContextFromType(0,
    CL_DEVICE_TYPE_GPU, NULL, NULL, NULL);


// get the list of GPU devices associated with context
clGetContextInfo(context, CL_CONTEXT_DEVICES, 0,
                                    NULL, &cb);

devices = malloc(cb);

clGetContextInfo(context, CL_CONTEXT_DEVICES, cb,
    devices, NULL);


// create a command-queue
cmd_queue = clCreateCommandQueue(context, devices[0],
    0, NULL);


memobjs[0] = clCreateBuffer(context,CL_MEM_WRITE_ONLY,

    sizeof(cl_char)*strlen("Hello World", NULL,
                                        NULL);

// create the program
program = clCreateProgramWithSource(context, 1,
    &program_source, NULL, NULL);
```

```
// build the program
err = clBuildProgram(program, 0, NULL, NULL, NULL,
                                            NULL);


// create the kernel
kernel = clCreateKernel(program, "vec_add", NULL);


// set the args values
err  = clSetKernelArg(kernel, 0, (void *) &memobjs[0],
                                sizeof(cl_mem));


// set work-item dimensions
global_work_size[0] = strlen("Hello World");;


// execute kernel
err = clEnqueueNDRangeKernel(cmd_queue, kernel, 1,
    NULL, global_work_size, NULL, 0, NULL, NULL);


// read output array
err = clEnqueueReadBuffer(cmd_queue, memobjs[0],
    CL_TRUE, 0, strlen("Hello World") *sizeof(cl_char),
    dst, 0, NULL, NULL);
```

AMD

```
// create the OpenCL context on a GPU device
cl_context = clCreateContextFromType(0,
```

**Define platform and queues**

```
// get the list of GPU devices associated with context
clGetContextInfo(context, CL_CONTEXT_DEVICES, 0,
                                      NULL, &cb);

devices = malloc(cb);

clGetContextInfo(context, CL_CONTEXT_DEVICES, cb,
    devices, NULL);


// create a command-queue
cmd_queue = clCreateCommandQueue(context, devices[0],
    0, NULL);
```

**Build  the program**

```
err                                     , NULL, NULL,
NULL);
```

**Create and setup kernel**

```
k                                       LL);

// set the args values
err  = clSetKernelArg(kernel, 0, (void *) &memobjs[0],
                              sizeof(cl_mem));

// set work-item dimensions
global_work_size[0] = n;

// execute kernel
err = clEnqueueNDRangeKernel(cmd_queue, kernel, 1,
    NULL, global_work_size, NULL, 0, NULL, NULL);
```

**Define Memory objects**

```
                                        READ_ONLY
    CL_MEM_COPY_HOST_PTR, sizeof(cl_char)*strlen("Hello
    World"), srcA, NULL);}
// create the program
program = clCreateProgramWithSource(context, 1,
    &program_source, NULL, NULL);
```

```
// read output array
err = cl                                , CL_TRUE,
    0, n*
```

**Execute the kernel**

**Read results on the host**

**Create the program**

AMD

# AMD RADEON™ HD 6900 SERIES

- Dual graphics engines

- VLIW4 core architecture

- Fast 256-bit GDDR5 memory interface
  - Up to 5.5 Gbps

### HD 6970

- 1536 Stream Processors

- 2.7 TFLOPs SP

- 683 GFLOPs DP

- VLIW4 thread processors

  - 4-way co-issue

  - All stream processing units have
    equal capabilities

    - Special functions (transcendentals) occupy
      3 of 4 issue slots

**Stream Processing Units**

**FP ops per clock**

4 32-bit MAD
2 64-bit MUL or ADD
1 64-bit MAD or FMA
1 Special Function

**Integer ops per clock**

4 24-bit MUL, ADD or MAD
2 32-bit ADD
1 32-bit MUL

Branch Unit

General Purpose Registers

AMD

**C function**

```
for (int i = 0; i < 24; i++)
{
        Y[i] = a*X[i] + Y[i];
}
```

- Serial execution, one iteration after the other

AMD

# *EXPOSING PARALLELISM*

**C function**

```
for (int i = 0; i < 24; i++)
{
        Y[i] = a*X[i] + Y[i];
}
```

**OpenCL kernel**

```
__kernel void
saxpy(const __global float * X,
          __global float * Y,
              const float a)
{
        uint i = get_global_id(0);
        Y[i] = a* X[i] + Y[i];
}
```

- Serial execution, one iteration after the other

- Parallel execution, multiple iterations at the same time

AMD

# WORK ITEM

- Think of work item as a parallel "thread" of execution

**Work items**

**1 saxpy operation per iteration**
**=**
**1 saxpy operation per work item**

```
0   1   2   ...                        10  11                                      22  23
```

```
for (int i = 0; i < 24; i++)                    {

{                                                 uint i = get_global_id(0);

        Y[i] = a*X[i] + Y[i];                     Y[i] = a* X[i] + Y[i];

}                                               }
```

AMD

## WORK GROUP

**Divide the execution domain into groups**

**Can exchange data and synchronize inside a group**

**Work items**

0   1   2   ...

0   1   2   ...

0   1   2   ...

`get_local_id(0)`

**Work groups**

AMD

Private
(per work item)

Local
(per work group)

Memory
consistent only
at barriers!

Global
(visible to all)

# MAPPING WORK-GROUPS ON GPUS

- Work-groups with multiple wavefronts → all wavefronts in same work-group scheduled on same SIMD

- Work-group size should be integral multiple of wavefront size

- (Advanced) Tip: When possible, pass work-group size at OpenCL kernel compile time

Can synchronize
within work group

?

Executes on
the same
SIMD

Cannot synchronize
across work groups

# MAPPING WORK-GROUPS ON CPUS

Can synchronize
within work group

Executes on
the same core

Cannot synchronize
across work groups

AMD

# PERFORMANCE DEVELOPER TOOLS

## DEVELOPER TOOLS SUITE

# AMD APP SDK - OVERVIEW

- OpenCL 1.1

- Windows 7, Linux
  - Apple OpenCL support thru MacOS
- Hardware Support
  - Fusion APUs
  - Discrete GPUs (Evergreen  N Islands + …)
  - X86 (SSE3, SSE4. XOP, FMA4)

- Khronos (khr) OpenCL extensions
  - Atomics
  - Images
  - GL Sharing
  - D3D10 Sharing
  - Bytes / shorts
  - Device fission (CPU)
  - FP64 (double precision) (CPU/GPU - Cypress)

- AMD (amd) OpenCL extensions
  - FP64 (double precision)
  - Media Ops (SAD, Pack, Unpack,…)
  - Printf
  - Popcnt
- OpenVideo Decode UVD (Windows 7)
- Multi-GPU (Windows 7)
- FFT and BLAS-3 Libraries
- Binary Image Format / Offline Compile
- Developer Tool Suite (Later slides)
  - Debugger, Profiler, Optimizers

- Graphics Driver releases
  - OpenCL run-time available by default in Catalyst driver builds (Windows)
  - Silent install (Windows)

AMD

# gDEBugger V5.8

- OpenCL and OpenGL API level Debugger, Profiler and Memory Analyzer
  - Access internal system information
  - Find bugs
  - Optimize compute performance
  - Minimize memory usage

- Windows, Linux and MacOS

- Proven mature product (6+ yrs)
- Widely deployed
  - 10k-s of users
  - Multiple industries

AMD

# GDEBUGGER 6.0 (VISUAL STUDIO EXTENSION)

- Visual Studio 2010 extension
  - Native look and feel
  - Find bugs
  - Shortens development time
  - Minimize memory usage

- Includes all gDEBugger's capabilities
  - OpenCL and OpenGL
  - Trace OpenCL API calls
  - Single step through OpenCL Kernels
  - Insert breakpoints into OpenCL kernels
  - View Kernel local variables, buffers & images

- Run-time OpenCL and DirectCompute™ profiler
- Identify bottlenecks and optimization opportunities
- Access GPU hardware performance counters for AMD GPUs
- Visualize timeline and API trace for an OpenCL program
- Display IL and ISA (kernal disassembly)
- Visual Studio 2008/2010 Plugin
- Command line version for Linux

- Version of GPU ShaderAnalyzer tool targeting Accelerated Parallel Processing
- Statically analyze of OpenCL Kernels for AMD Radeon GPUs.
- Display compiled Kernels as IL or GPU ISA
- Display statistics from the AMD Kernel Compiler
- Estimate Kernel performance
- Includes support for previous 12 versions of Catalyst Driver
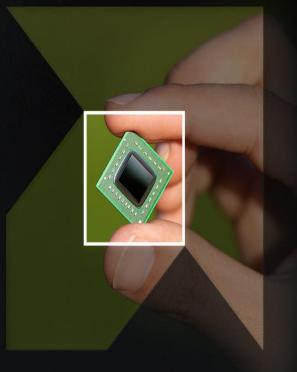
# *CODEANALYST PERFORMANCE ANALYZER*



- Profiling suite to identify, investigate and tune code performance on AMD platforms

- Find time critical hot-spots
  - C, C++, Fortran
  - Java, .Net managed cod

-

- Diagnose performance issues
  - Time based profiling
  - Event based sampling
  - Instruction based sampling

- Identify thread-affinity and core utilization problems

- Windows and Linux platforms

# *LINUX*

# *AMD APP OPENCL FOR LINUX*



- SDK 2.5
  - Linux OpenCL SDK available
  - gDEBugger
  - LLVM optimization passes
  - APP Profiler
  - Code Analyst

- OpenCL
  - Open standard
  - Closed implementation

- Multicoreware
  - Parallel Path Analyzer
  - Task Manager
  - GMAC



redhat.
v6, 5.5

SuSe
v11.3

ubuntu
v10.04

**AMD**

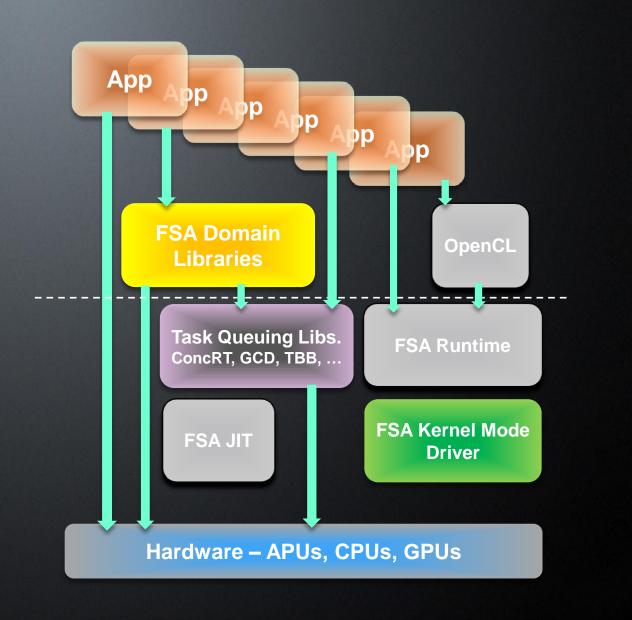| Physical Integration | Optimized Platforms | Architectural Integration | System Integration |
|---|---|---|---|
| Integrate CPU & GPU in silicon | GPU Compute C++ support | Unified Address Space for CPU and GPU | GPU compute context switch |
| Unified Memory Controller | User mode schedulng | GPU uses pageable system memory via CPU pointers | GPU graphics pre-emption |
| | | | Quality of Service |
| Common Manufacturing Technology | Bi-Directional Power Mgmt between CPU and GPU | Fully coherent memory between CPU & GPU | Extend to Discrete GPU |

AMD

# OPENCL DRIVER STACK (NOW)

App
App
App
App
App

**Domain Libraries**

**OpenCL, DirectX Runtimes,
User Mode Drivers**
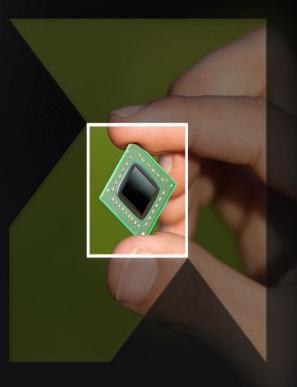
**Graphics Kernel Mode Driver**

**Hardware – APUs, CPUs, GPUs**

AMD user mode component

AMD kernel mode component

All other SW contributed by
AMD or 3rd parties

QUESTIONS?

# Disclaimer & Attribution

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions and typographical errors.

The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. There is no obligation to update or otherwise correct or revise this information. However, we reserve the right to revise this information and to make changes from time to time to the content hereof without obligation to notify any person of such revisions or changes.

NO REPRESENTATIONS OR WARRANTIES ARE MADE WITH RESPECT TO THE CONTENTS HEREOF AND NO RESPONSIBILITY IS ASSUMED FOR ANY INACCURACIES, ERRORS OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION.

ALL IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR ANY PARTICULAR PURPOSE ARE EXPRESSLY DISCLAIMED.  IN NO EVENT WILL ANY LIABILITY TO ANY PERSON BE INCURRED FOR ANY DIRECT, INDIRECT, SPECIAL OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

AMD◢