Most parallel code runs on CPUs designed for scalar workloads
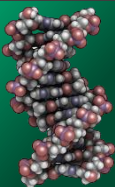
# WASTES POWER

AMD

# APPLICATION AREAS WITH ABUNDANT PARALLEL WORKLOADS

**Natural UI & Gestures**

Touch, gesture, and voice

**Biometric Recognition**

Secure, fast, accurate: face, voice, fingerprints

**Augmented Reality**

Superimpose graphics, audio, and other digital information as a virtual overlay

**Content Everywhere**

Content from any source to any display seamlessly

**Beyond HD Experiences**

Streaming media, new codecs, 3D, transcode, audio

**AV Content Management**

Searching, indexing and tagging of video & audio. multimedia data mining

AMD

# *ARCHAEOLOGY?*
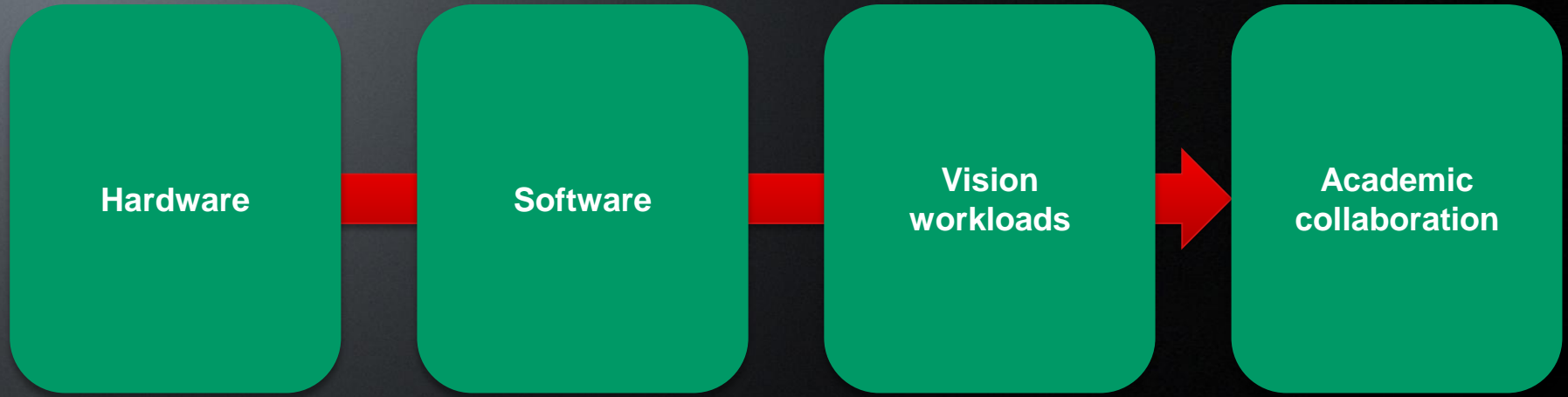


Tintagel Castle and Fountains Abbey
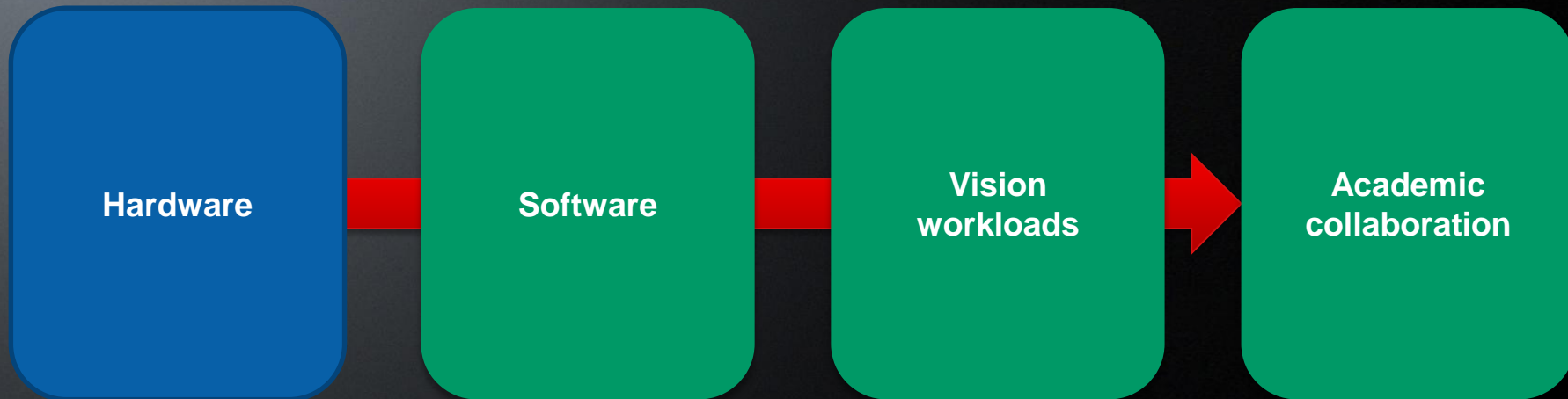Source: English Heritage



Knossos
Source: wikipedia

# *WHERE AM I GOING WITH THIS?*

**Hardware** → **Software** → **Vision workloads** → **Academic collaboration**

AMD

Hardware

Software

Vision workloads
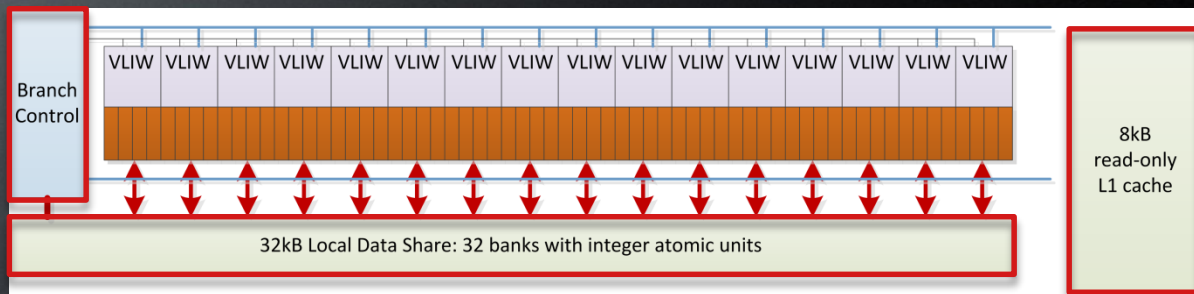
Academic collaboration

AMD GPU and APU Architecture
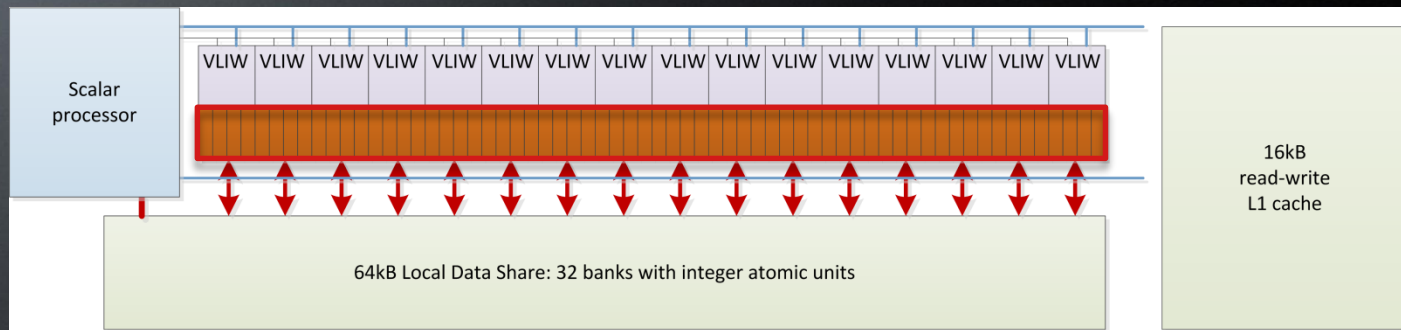The HD7970 and Graphics Core Next

AMD

# THE SIMD CORE

- The SIMD unit on the old Radeon architectures had a branch control but full scalar execution was performed globally
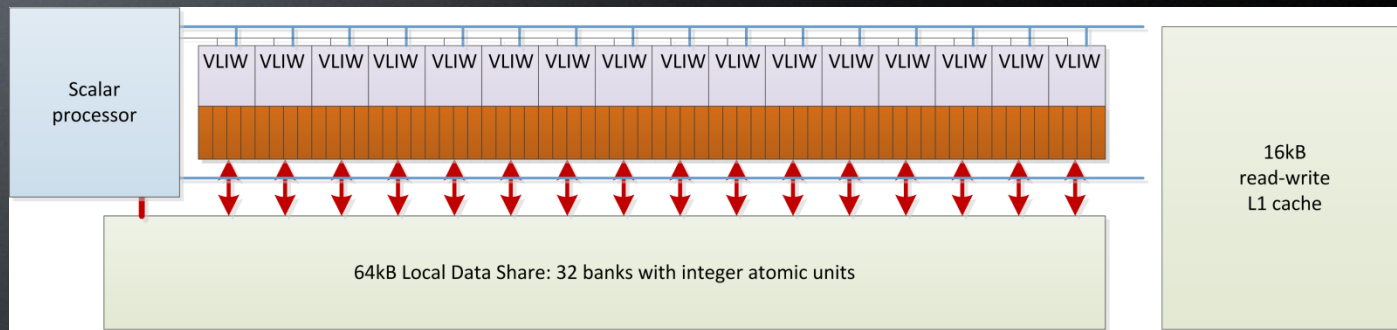
# THE SIMD CORE

- On the HD7970 we have a full scalar processor and the L1 cache and LDS have been doubled in size
- Then let us consider the VLIW ALUs

# THE SIMD CORE

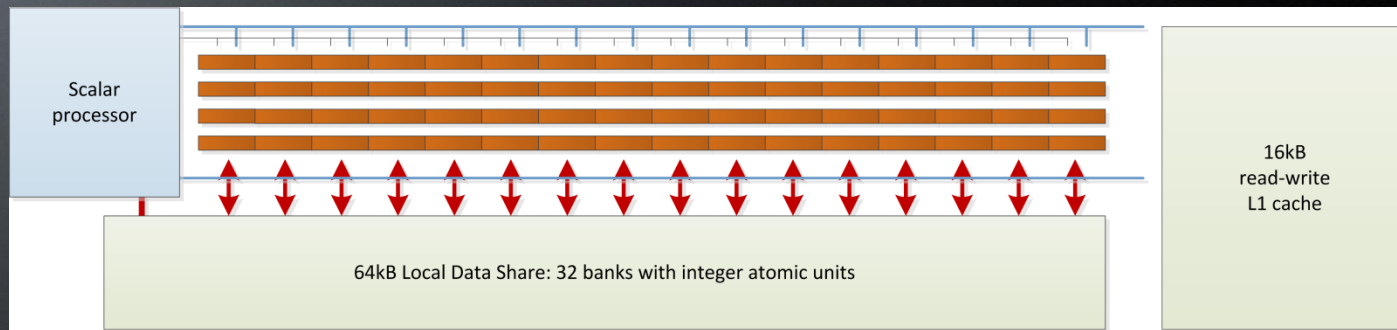- Remember we could view the architecture two ways:
  - An array of VLIW units
  - A VLIW cluster of vector units

# THE SIMD CORE

- Now that we have a scalar processor we can dynamically schedule instructions rather than relying on the compiler
- No VLIW!



- The heart of Graphics Core Next:
  - A scalar processor plus four 16-wide vector units
  - Each lane of the vector, and hence each IL work item, is now scalar

# THE SIMD CORE



- The scalar core manages a large number of threads
  - Each thread requires its set of vector registers
  - Significant register state for both scalar and vector storage
  - 10 waves per SIMD, 40 waves per CU (core), 2560 work items per CU, 81920 work items on the HD7970

# FAMILIAR?

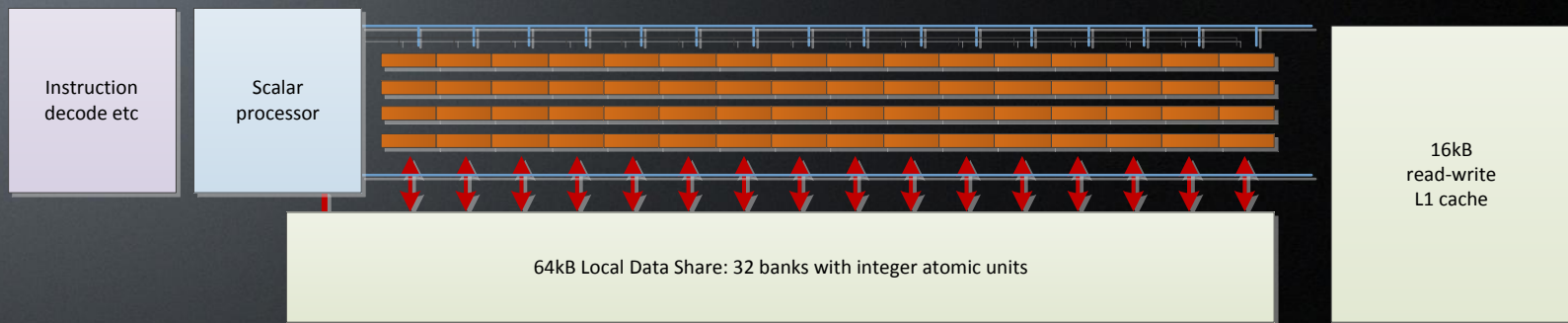- If we add the frontend of the core…

"Graphics Core Next" core

| Instruction decode etc | Scalar processor | | 16kB read-write L1 cache |
|---|---|---|---|

64kB Local Data Share: 32 banks with integer atomic units

"Barcelona" core

| Instruction decode etc | Scalar units (3 scalar ALUs, branch control etc) | FADD SSE<br>FMUL SSE<br>FMISC SSE | 64kB read-write L1 cache | 512kB read-write L2 cache |
|---|---|---|---|---|

AMD

# AMD RADEON HD7970 - GLOBALLY

- Two command processors
  - Capable of processing two command queues concurrently

- Full read/write L1 data caches

- SIMD cores grouped in fours
  - Scalar data and instruction cache per cluster
  - L1, LDS and scalar processor per core

- Up to 32 cores / compute units

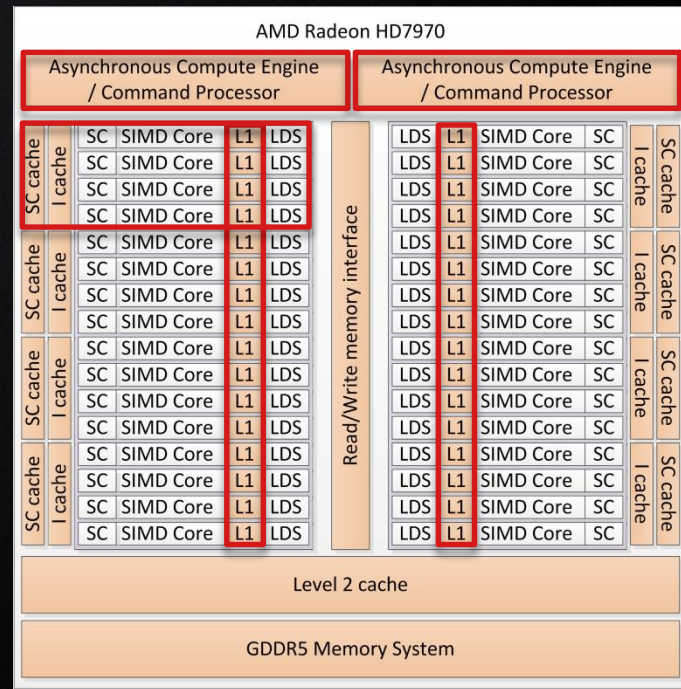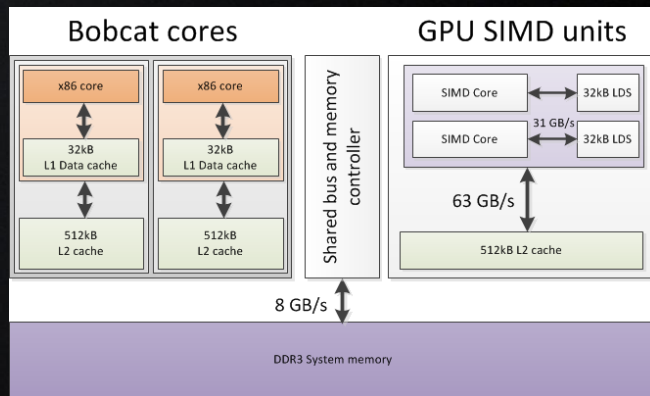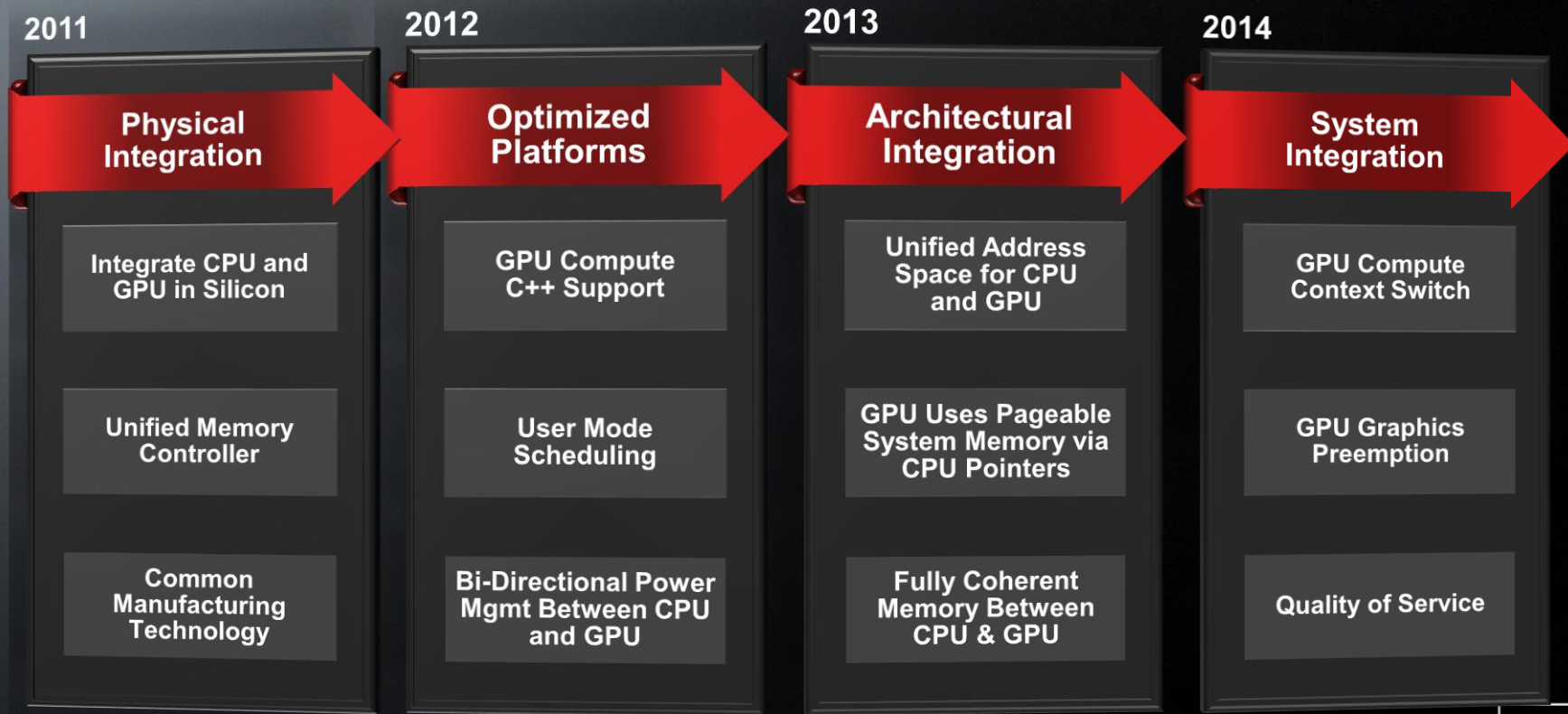# APU: ACCELERATED PROCESSING UNIT

- The APU has arrived and it is a great advance over previous platforms
- Combines scalar processing on CPU with parallel processing on the GPU and high bandwidth access to memory
- How do we make it even better going forward?
  - Easier to program
  - Easier to optimize
  - Easier to load balance
  - Higher performance
  - Lower power



APU E350 (2011)

AMD

# HETEROGENEOUS SYSTEM ARCHITECTURE ROADMAP

## 2011
### Physical Integration
- Integrate CPU and GPU in Silicon
- Unified Memory Controller
- Common Manufacturing Technology

## 2012
### Optimized Platforms
- GPU Compute C++ Support
- User Mode Scheduling
- Bi-Directional Power Mgmt Between CPU and GPU

## 2013
### Architectural Integration
- Unified Address Space for CPU and GPU
- GPU Uses Pageable System Memory via CPU Pointers
- Fully Coherent Memory Between CPU & GPU

## 2014
### System Integration
- GPU Compute Context Switch
- GPU Graphics Preemption
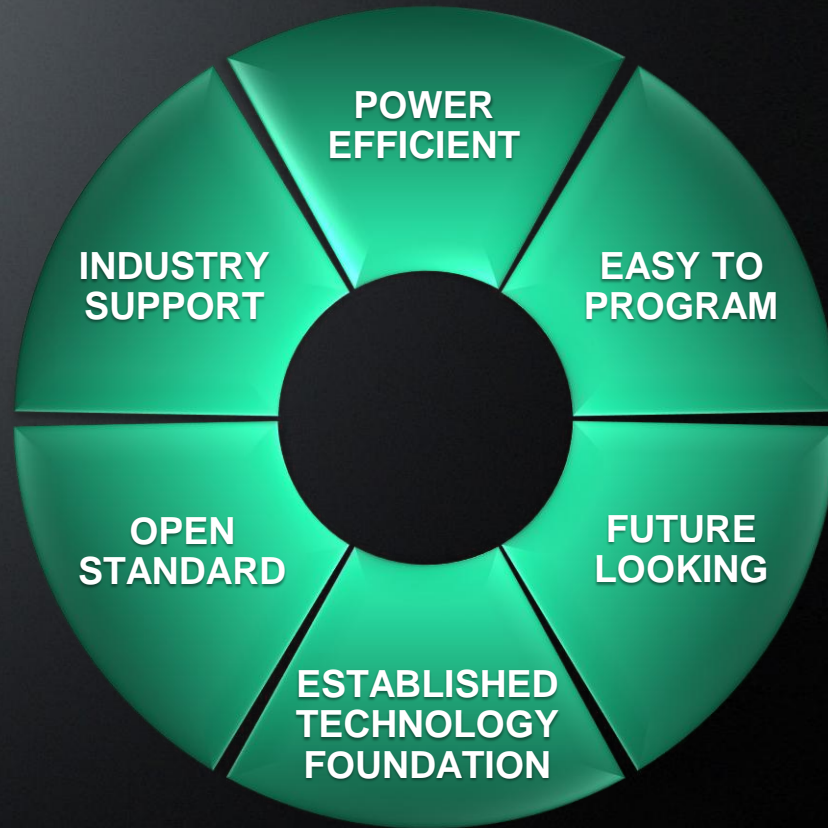- Quality of Service

AMD

# HETEROGENEOUS SYSTEM ARCHITECTURE
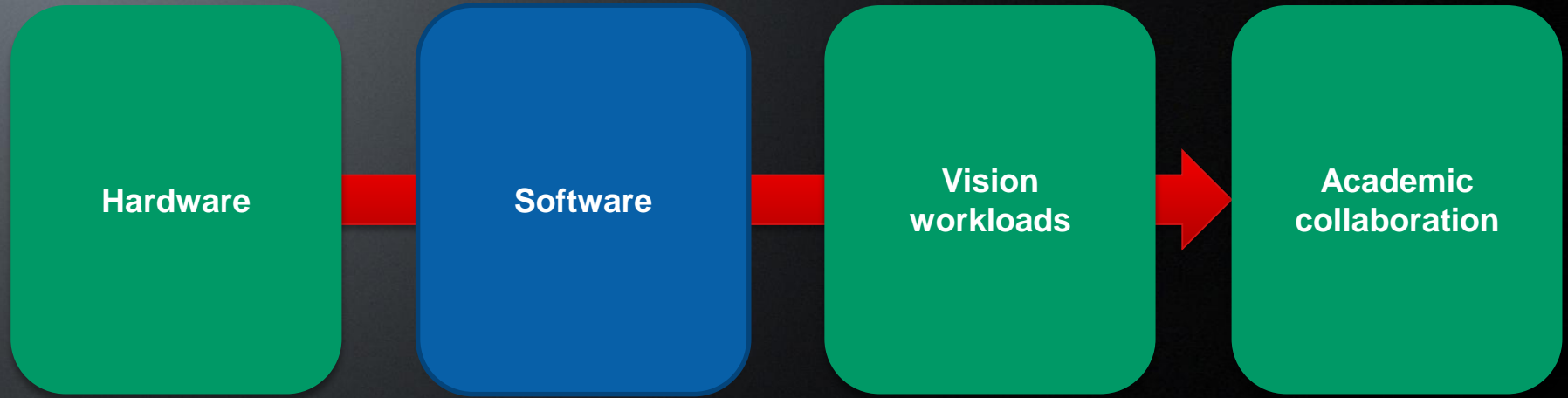Brings All the Processors in a System into Unified Coherent Memory

# *HETEROGENEOUS SYSTEM ARCHITECTURE – AN OPEN PLATFORM*

- Open Architecture, published specifications
  - HSAIL virtual ISA
  - HSA memory model
  - HSA system architecture

- ISA agnostic for both CPU and GPU

- HSA Foundation formed in June 2012

- Inviting academic partners to join us, in all areas
  - Hardware design
  - Operating Systems
  - Tools and Middleware
  - Applications

Hardware

Software
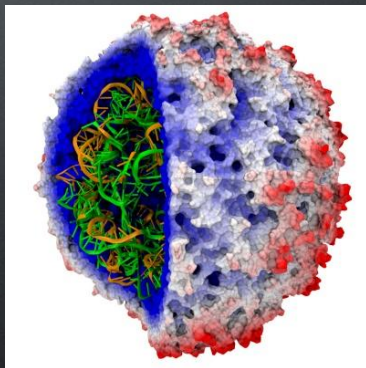
Vision workloads

Academic collaboration

AMD Software Infrastructure

AMD

# *OPENCL™ TODAY*

- OpenCL
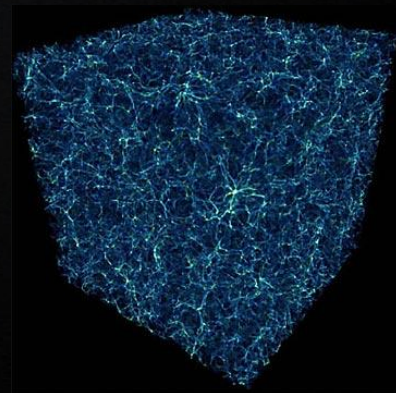  - Open development platform for multi-vendor heterogeneous architectures
  - Broad industry support: Created by architects from AMD, Apple, IBM, Intel, Nvidia, Sony, etc.
  - It provides C API and an extended subset of C99 kernel language.
  - Excellent performance, but programming can be longwinded and difficult

Molecular Dynamics
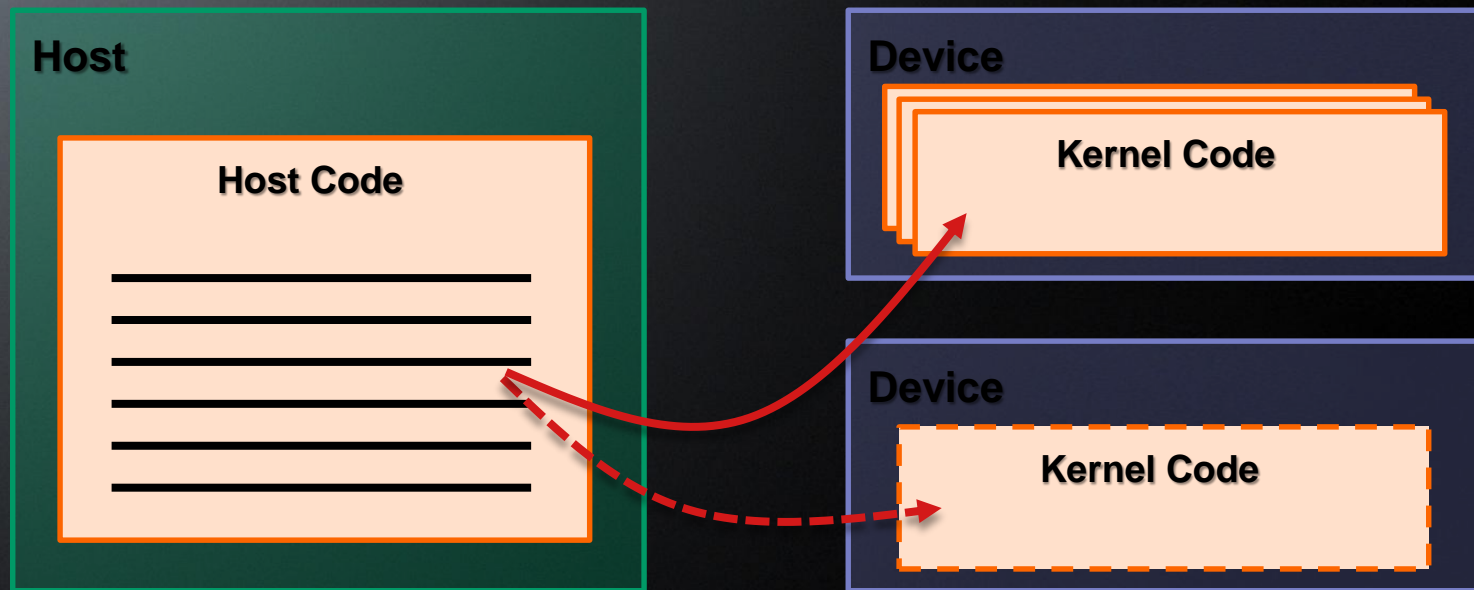BUDE, NMAD, Folding@Home

Ray Tracing
LuxRenderer

Search for missing matter in the Universe

# HOST AND DEVICE MODEL

- OpenCL™ uses a host and device model
  - Host code dispatches kernel code to the devices

# *EXAMPLE – VECTOR ADDITION (HOST PROGRAM)*

```c
// 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);
```

Define platform and queues

```c
// allocate the buffer memory objects
memobjs[0] = clCreateBuffer(context, CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR, sizeof(cl_float)*n, srcA, NULL);
memobjs[1] = clCreateBuffer(context, CL_MEM_READ_ONLY | CL_MEM_
memobjs[2] = clCreateBuffer(context, CL_MEM_WRITE_ONLY, sizeof(
```

```c
// create the program
program = clCreateProgramWithSource(context, 1, &program_source
```

**Same boiler-plate code across
virtually every OpenCL host code**

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

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

// set the args values
err  = clSetKernelArg(kernel, 0, sizeof(cl_mem), (void *) &memobjs[0]);
err |= clSetKernelArg(kernel, 1, sizeof(cl_mem), (void *) &memobjs[1]);
err |= clSetKernelArg(kernel, 2, sizeof(cl_mem), (void *) &memobjs[2]);

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

Create and setup kernel

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

Execute the kernel

```c
// read output array
err = clEnqueueReadBuffer(cmd_queue, memobjs[2], CL_TRUE, 0, n*sizeof(cl_float), dst, 0, NULL, NULL);
```

Read results on the host

AMD

# OPENCL™ C++ HOST API CODE FOR VECTOR ADD

```
std::function<Event (const EnqueueArgs&, Buffer, Buffer, Buffer)> vadd =
            make_kernel<Buffer, Buffer, Buffer>(Program(program_source), "vadd");


memobj[0] = Buffer (CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR, sizeof(float) * n, srcA);
memobj[1] = Buffer (CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR, sizeof(float) * n, srcB);
memobj[2] = Buffer (CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR, sizeof(float)


vadd(EnqueueArgs(NDRange(n)), memobj[0], memobj[1], memobj[2]));


enqueueReadBuffer(memobj[2], CL_TRUE, sizeof(float) * n, dest);
```

**Program automatically created and compiled**

**Defaults, no need to reference context, command queue**

**No clRelease*XXX* cleanup code required**

AMD

# *OPENCL C++ INTERFACE AND KERNEL LANGUAGE EXTENSION*

- Khronos has defined a common C++ header file containing a high level interface to OpenCL
  - Common defaults for the platform and command-queue
  - Simplify basic API by parameterizing through template types
  - Maintain object lifetimes through constructors and destructors
  - Support function-like kernel dispatch

- OpenCL C++ kernel language extension
  - Support static C++ language features
  - Support templating and overloading
  - Innovative integration with OpenCL C address spaces

- Make it easier to create applications for heterogeneous platforms.

AMD

# CODE SIZE REDUCTION

- Sometimes substantial reduction in code size:

| Application | C lines | C++ lines | Reduction |
|---|---|---|---|
| Vector addition | 268 | 140 | 47.7% |
| Pi computation | 306 | 166 | 45.8% |
| Ocean simulation | 1386 | 533 | 61.5% |
| Particle simulation | 733 | 601 | 18.0% |
| Radix sort | 627 | 593 | 5.4% |

- Depending on the complexity, you may not see this in real use cases. But for beginners, the simplification of the program makes a difference.
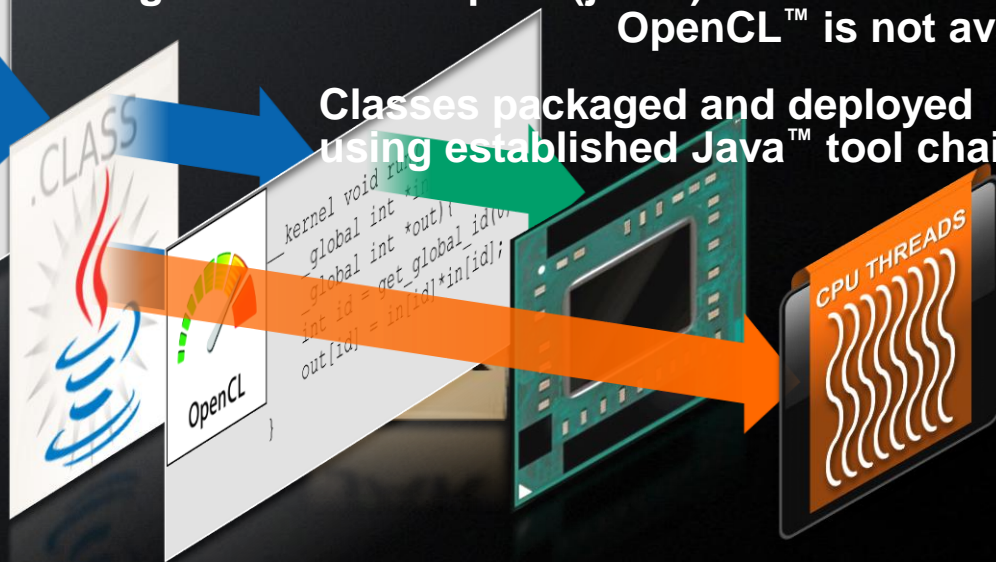
AMD

# JAVA ENABLEMENT BY APARAPI

*Aparapi = Runtime capable of converting Java™ bytecode to OpenCL™*
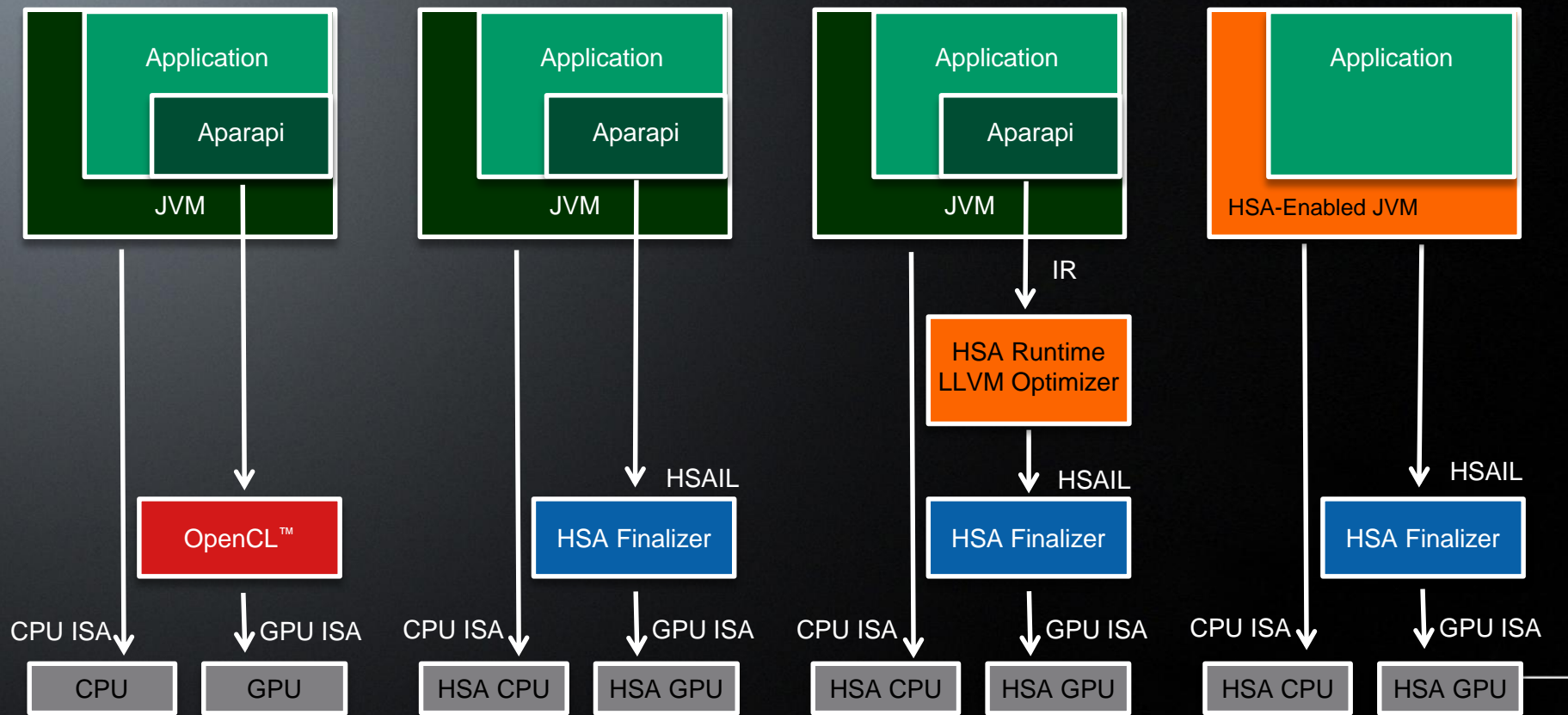


**Developer creates Java™ source**

**Source compiled to class files (bytecode) using standard compiler (javac)**

**For execution on any OpenCL™ 1.1+ capable device OR execute via a thread pool if OpenCL™ is not available**

**Classes packaged and deployed using established Java™ tool chain**

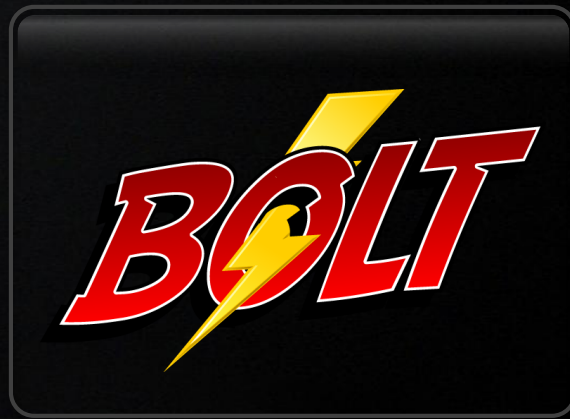# JAVA AND APARAPI HSA ENABLEMENT ROADMAP

| Application | Application | Application | Application |
|---|---|---|---|
| Aparapi | Aparapi | Aparapi | |
| JVM | JVM | JVM | HSA-Enabled JVM |

IR

HSA Runtime
LLVM Optimizer

| OpenCL™ | HSAIL | HSAIL | HSAIL |
|---|---|---|---|
| | HSA Finalizer | HSA Finalizer | HSA Finalizer |

CPU ISA — GPU ISA    CPU ISA — GPU ISA    CPU ISA — GPU ISA    CPU ISA — GPU ISA

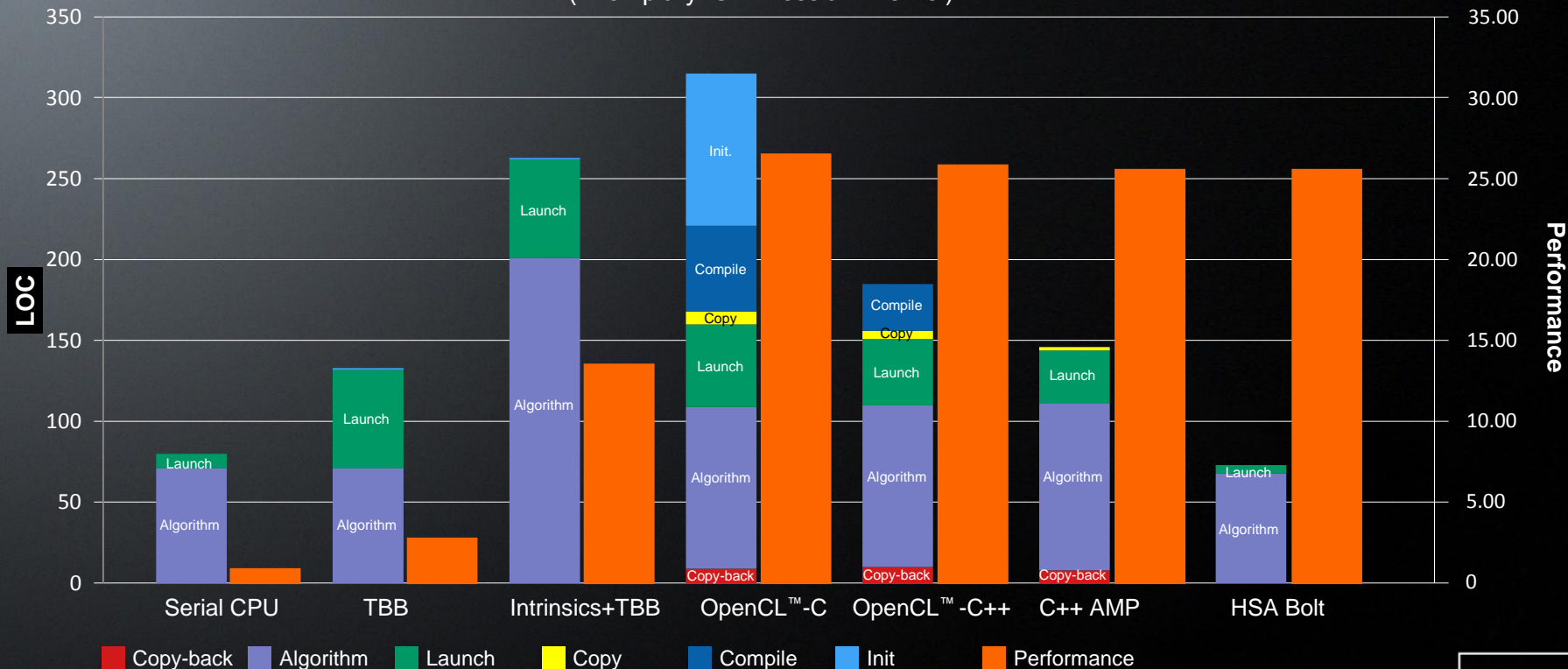| CPU | GPU | HSA CPU | HSA GPU | HSA CPU | HSA GPU | HSA CPU | HSA GPU |

AMD

# *INTRODUCING HSA BOLT – PARALLEL PRIMITIVES LIBRARY FOR HSA*

- Easily leverage the inherent power efficiency of GPU computing

  - Common routines such as scan, sort, reduce, transform
  - More advanced routines like heterogeneous pipelines
  - Bolt library works with OpenCL or C++ AMP

- Enjoy the unique advantages of the HSA platform

  - Move the computation, not the data
  - Use appropriate computation resource (CPU or GPU or others)

- Finally a single source code base for the CPU and GPU!

  - Scientists can focus on innovations!

- Bolt preview available in AMD APP SDK 2.8.

# LINES-OF-CODE AND PERFORMANCE FOR DIFFERENT PROGRAMMING MODELS



(Exemplary ISV "Hessian" Kernel)

AMD A10-5800K APU with Radeon™ HD Graphics – CPU: 4 cores, 3800MHz (4200MHz Turbo); GPU: AMD Radeon HD 7660D, 6 compute units, 800MHz; 4GB RAM.
Software – Windows 7 Professional SP1 (64-bit OS); AMD OpenCL™ 1.2 AMD-APP (937.2); Microsoft Visual Studio 11 Beta

# OPENCL MODULE FOR OPENCV LIBRARY
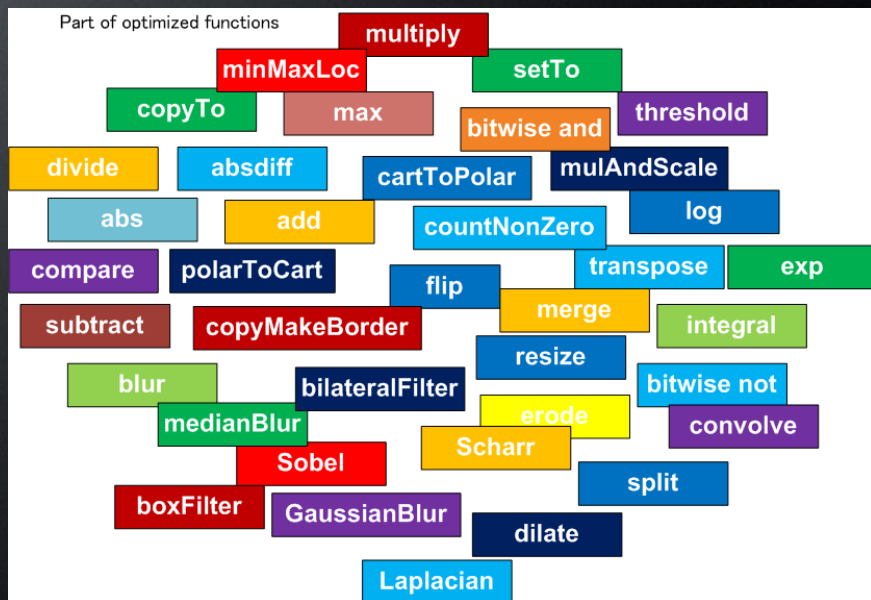
**OpenCV**

- Based on OpenCV 2.4, first OCL module released in Nov 2012.

- Implement and maintain an OpenCL module that is optimised for AMD GPU and APU platforms, including utility functions, low-level vision primitives, and high-level algorithms.

- Support any OpenCL 1.1 compatible device, tested on AMD's, Intel's and NVIDIA's GPU.

- Designed as a host-level API and device-level kernels.

- Requirement: OpenCL SDK, AMD FFT and BLAS library

- No prior knowledge of OpenCL required.

AMD

# OPENCL MODULE FOR OPENCV LIBRARY

- Over 80 kernels are provided with OCL module in 2.4.9 release.
- OCL module can run on NVIDIA, Intel and AMD GPUs without any modification.
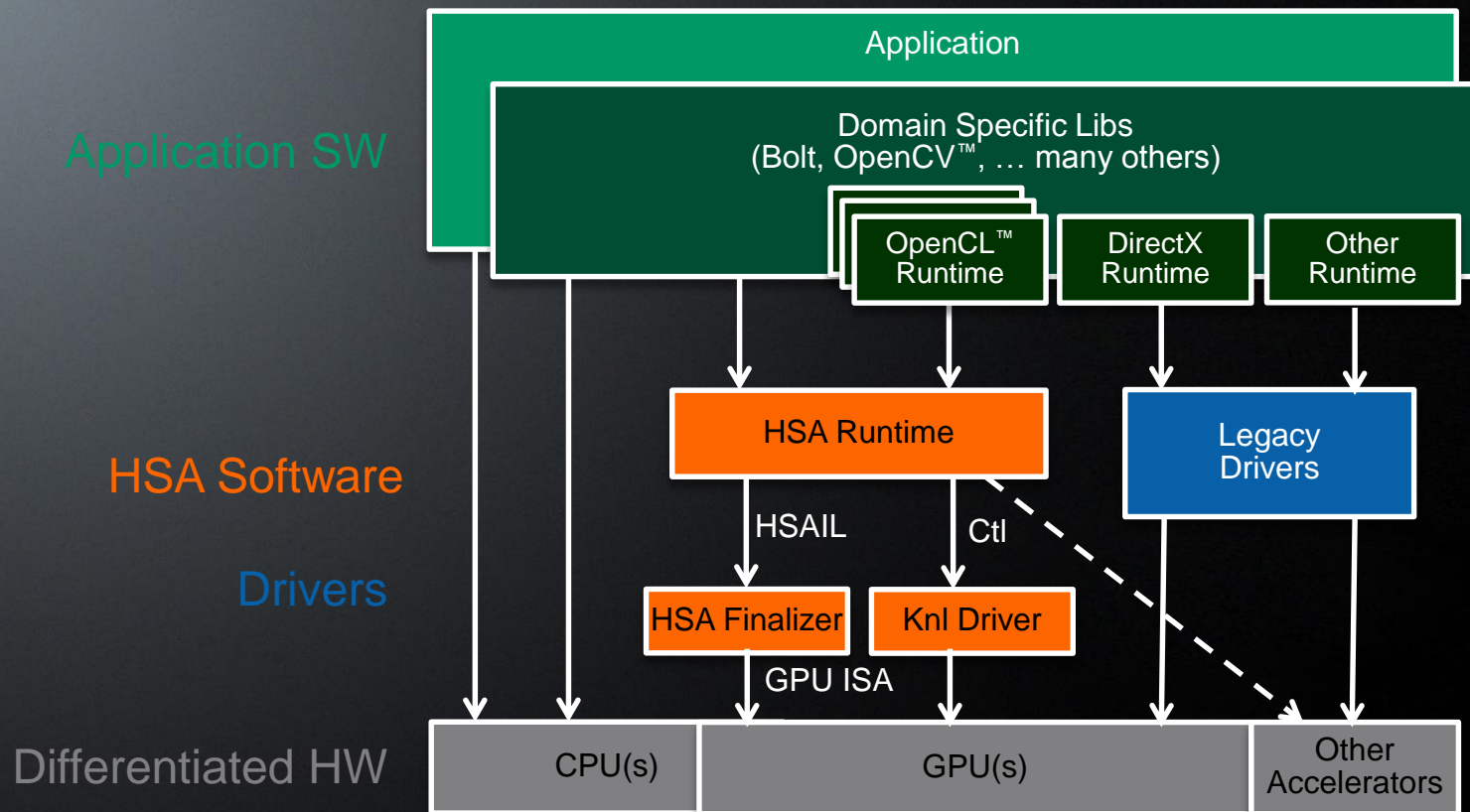


Part of optimized functions

multiply, minMaxLoc, setTo, copyTo, max, bitwise and, threshold, divide, absdiff, cartToPolar, mulAndScale, abs, add, countNonZero, log, compare, polarToCart, flip, transpose, exp, subtract, copyMakeBorder, merge, integral, resize, blur, bilateralFilter, bitwise not, medianBlur, erode, convolve, Scharr, Sobel, split, boxFilter, GaussianBlur, dilate, Laplacian

## OTHER LIBRARIES AND TOOLS

- APPML, contains FFT and BLAS functions, primarily targeting AMD GPUS and APUs.
  - User specifies problem parameters through library API
  - Kernel generator creates tailored OpenCL kernels
  - Kernels are dispatched for execution
- APP Kernel Analyzer
  - static analysis tool to compile, analyse and disassemble OCL kernel for GPUs
- APP profiler
  - performance analysis tool that gathers data from OCL runtime and GPUs during the execution.
- Code Analyst
- gDEBugger, an OpenCL and OpenGL debugger and memory analyser
- GPUPerfAPI
  - a library can be integrated directly into your application for accessing GPU performance counters.
- CodeXL, offers GPU debugging, CPU and GPU profiling, static OpenCL kernel analysis capabilities.
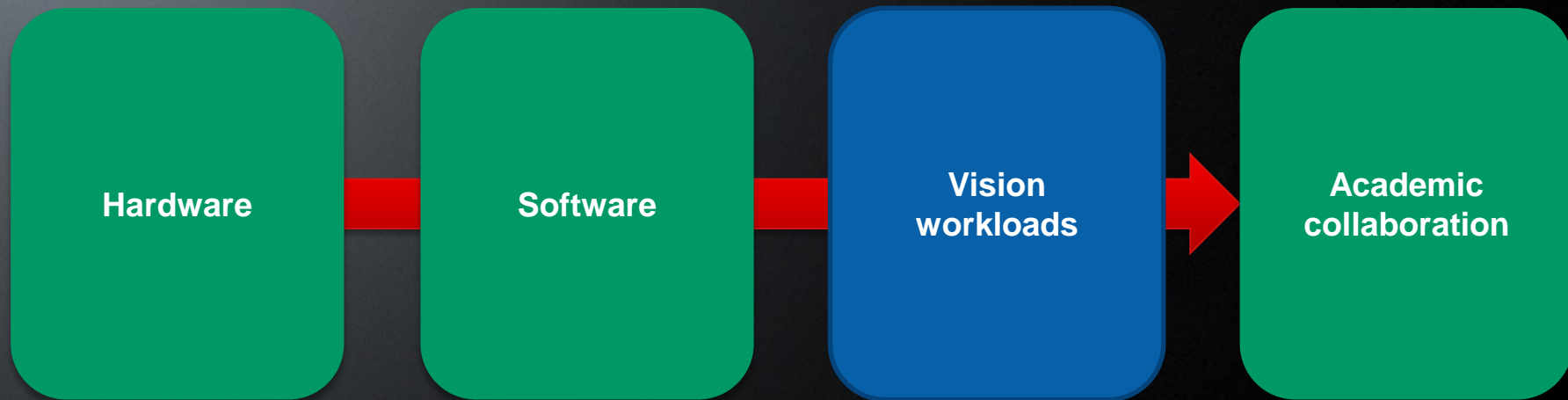
AMD

# HSA SOLUTION STACK

**Application SW**

**HSA Software**

**Drivers**

**Differentiated HW**

Application

Domain Specific Libs
(Bolt, OpenCV™, … many others)

OpenCL™ Runtime

DirectX Runtime

Other Runtime

HSA Runtime

Legacy Drivers

HSAIL

Ctl

HSA Finalizer

Knl Driver

GPU ISA

CPU(s)

GPU(s)

Other Accelerators

AMD

# AMD'S OPEN SOURCE COMMITMENT TO HSA

- We will open source our linux execution and compilation stack
  - Jump start the ecosystem
  - Allow a single shared implementation where appropriate
  - Enable university research in all areas

| Component Name | AMD Specific | Rationale |
| --- | --- | --- |
| HSA Bolt Library | No | Enable understanding and debug |
| OpenCL HSAIL Code Generator | No | Enable research |
| LLVM Contributions | No | Industry and academic collaboration |
| HSA Assembler | No | Enable understanding and debug |
| HSA Runtime | No | Standardize on a single runtime |
| HSA Finalizer | Yes | Enable research and debug |
| HSA Kernel Driver | Yes | For inclusion in linux distros |

AMD

**Hardware** → **Software** → **Vision workloads** → **Academic collaboration**

Example ISV workload: Haar Face Detection
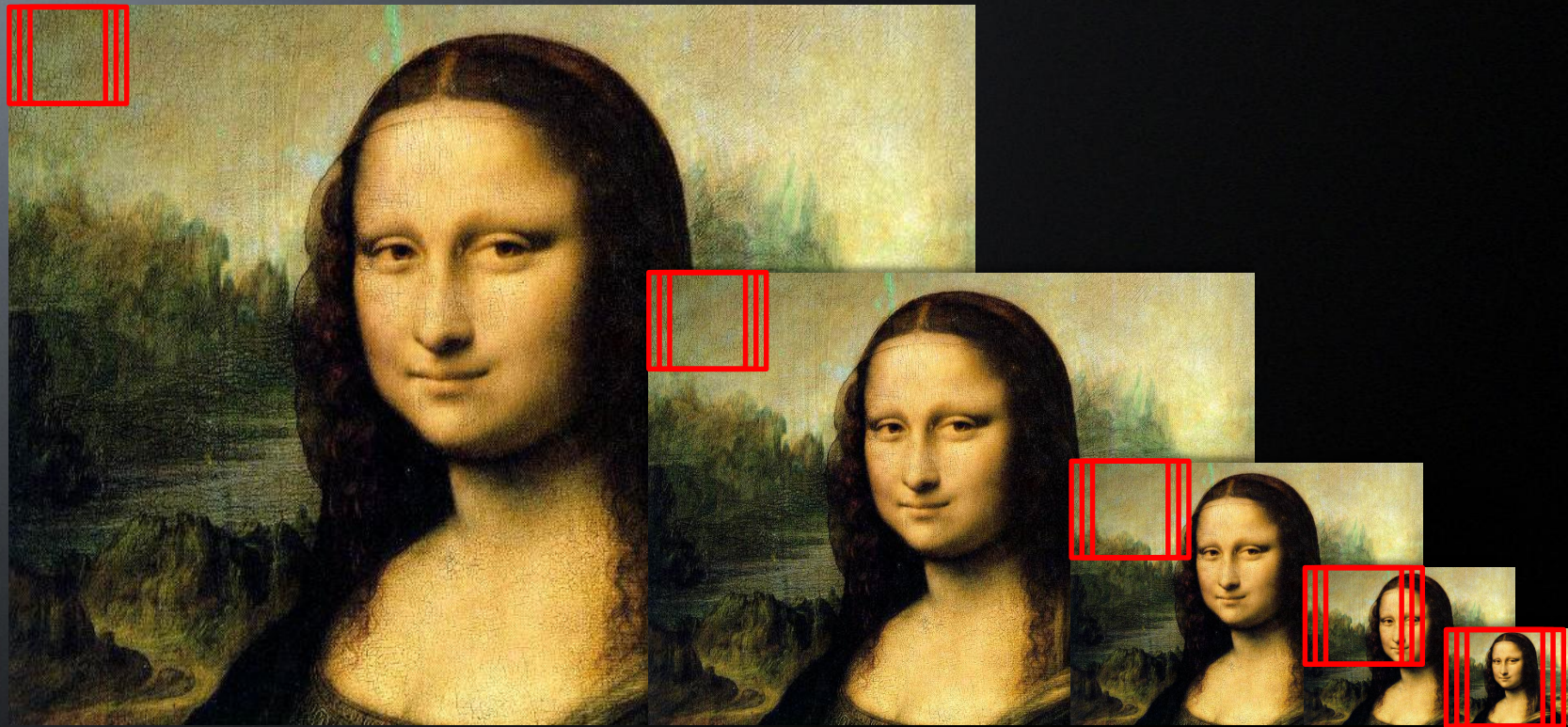Cornerstone Technology for Computer Vision

AMD

**Quick HD Calculations**

Search square = 21 x 21

Pixels = 1920 x 1080 = 2,073,600

Search squares = 1900 x 1060 = ~2 Million

# LOOKING FOR DIFFERENT SIZE FACES – BY SCALING THE VIDEO FRAME



**More HD Calculations**
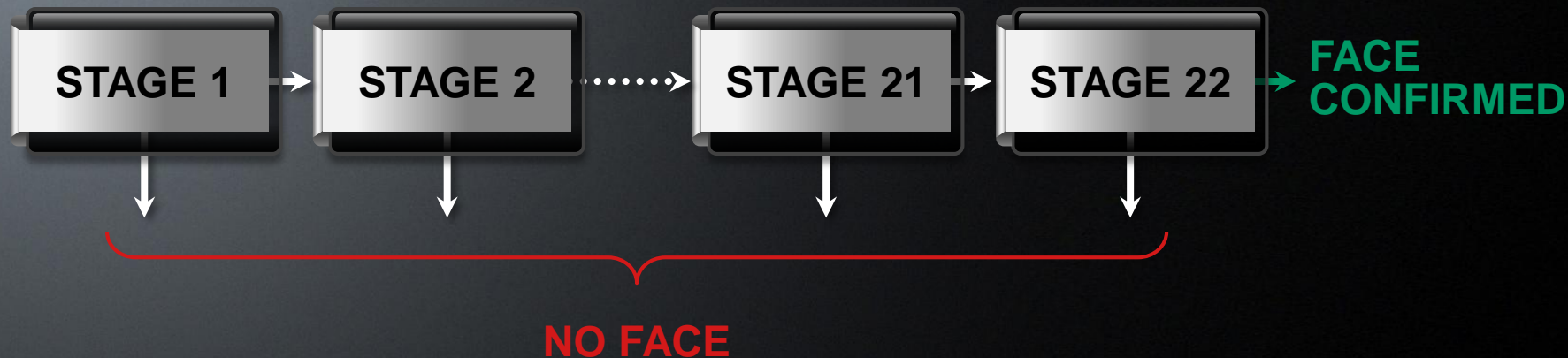
70% scaling in H and V

Total Pixels = 4.07 Million

Search squares = 3.8 Million

AMD

# HAAR CASCADE STAGES

# 22 CASCADE STAGES, EARLY OUT BETWEEN EACH

| STAGE 1 | → | STAGE 2 | ⋯⋯⋯→ | STAGE 21 | → | STAGE 22 | → | **FACE CONFIRMED** |

**NO FACE**

**Final HD Calculations**

Search squares = 3.8 million

Average features per square = 124

Calculations per feature = 100

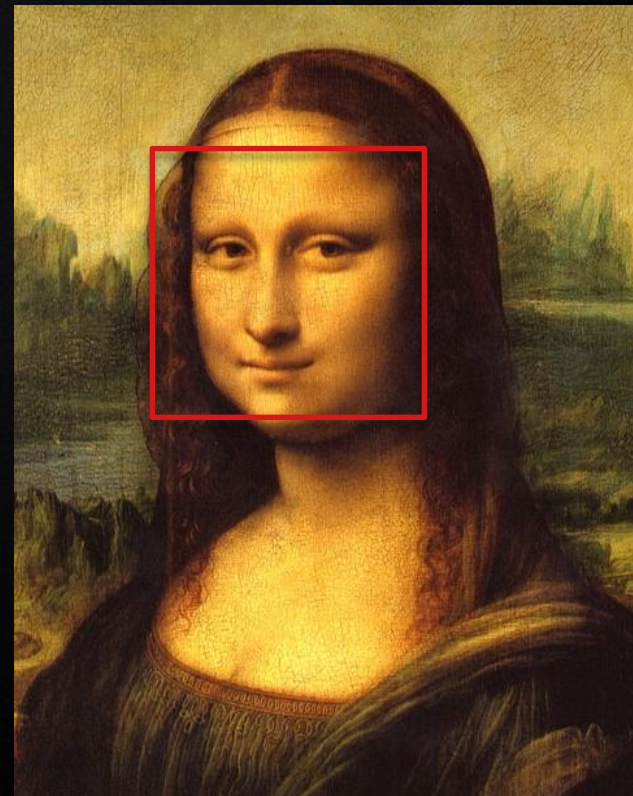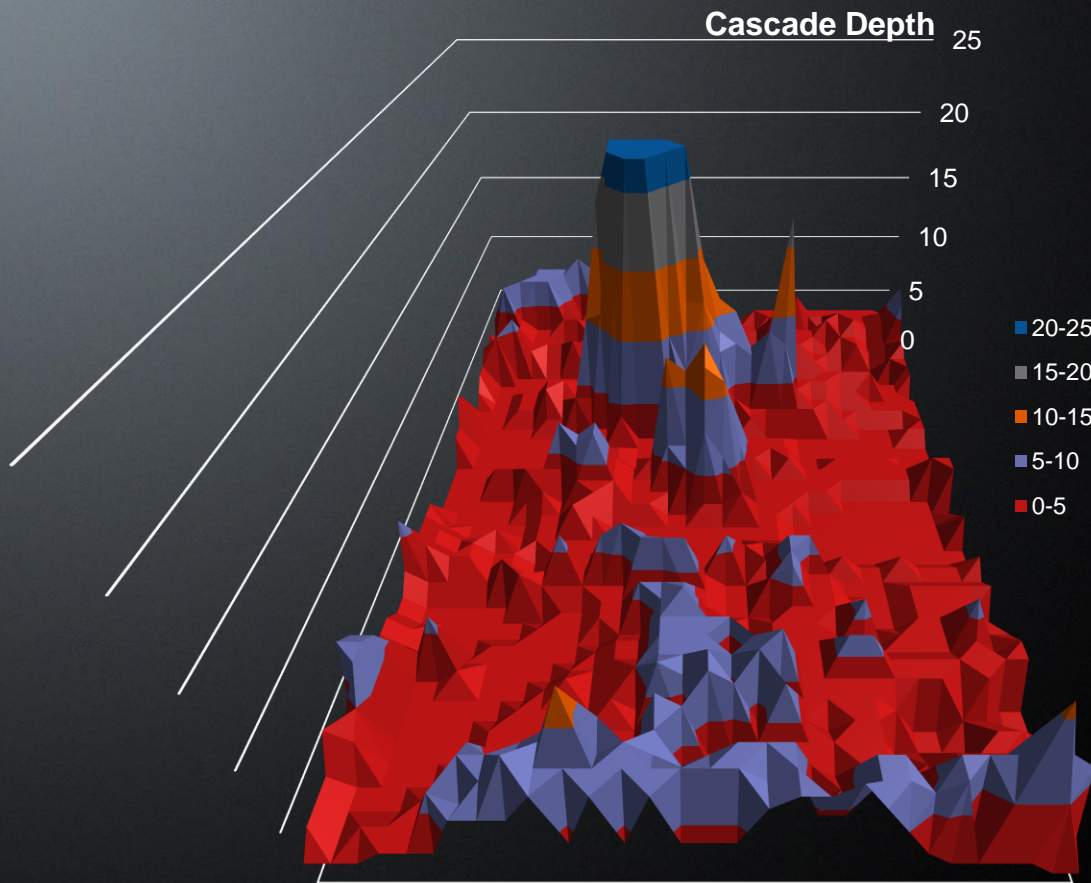Calculations per frame = 47 GCalcs

**Calculation Rate**
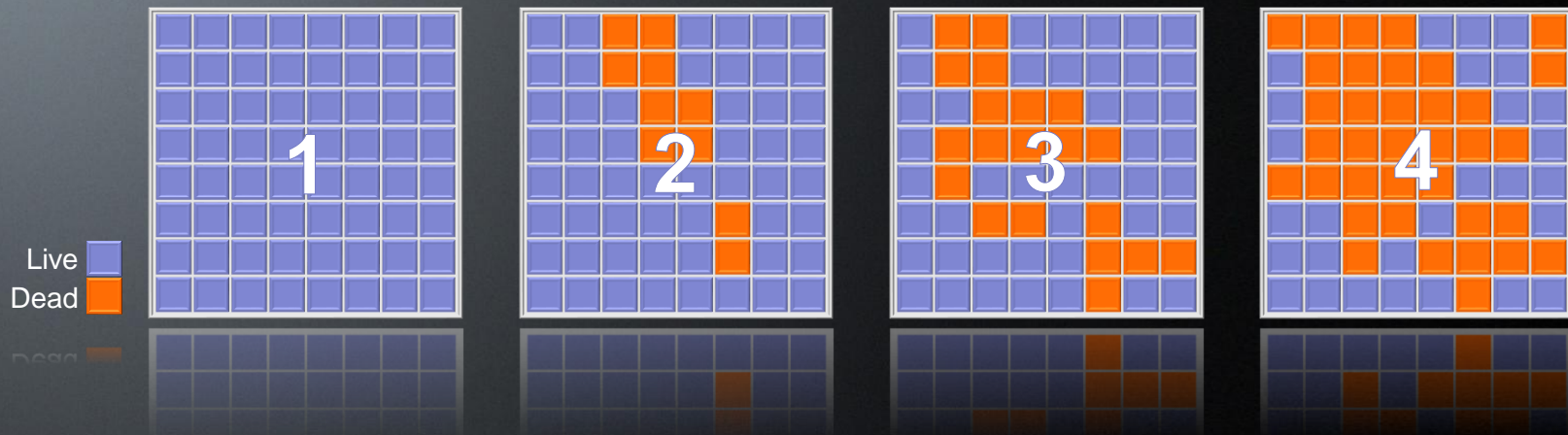
30 frames/sec = 1.4TCalcs/second

60 frames/sec = 2.8TCalcs/second

*…and this only gets front-facing faces*

AMD

# CASCADE DEPTH ANALYSIS



Cascade Depth

25
20
15
10
5
0

- 20-25
- 15-20
- 10-15
- 5-10
- 0-5

AMD

# UNBALANCING DUE TO EXITS IN EARLIER CASCADE STAGES
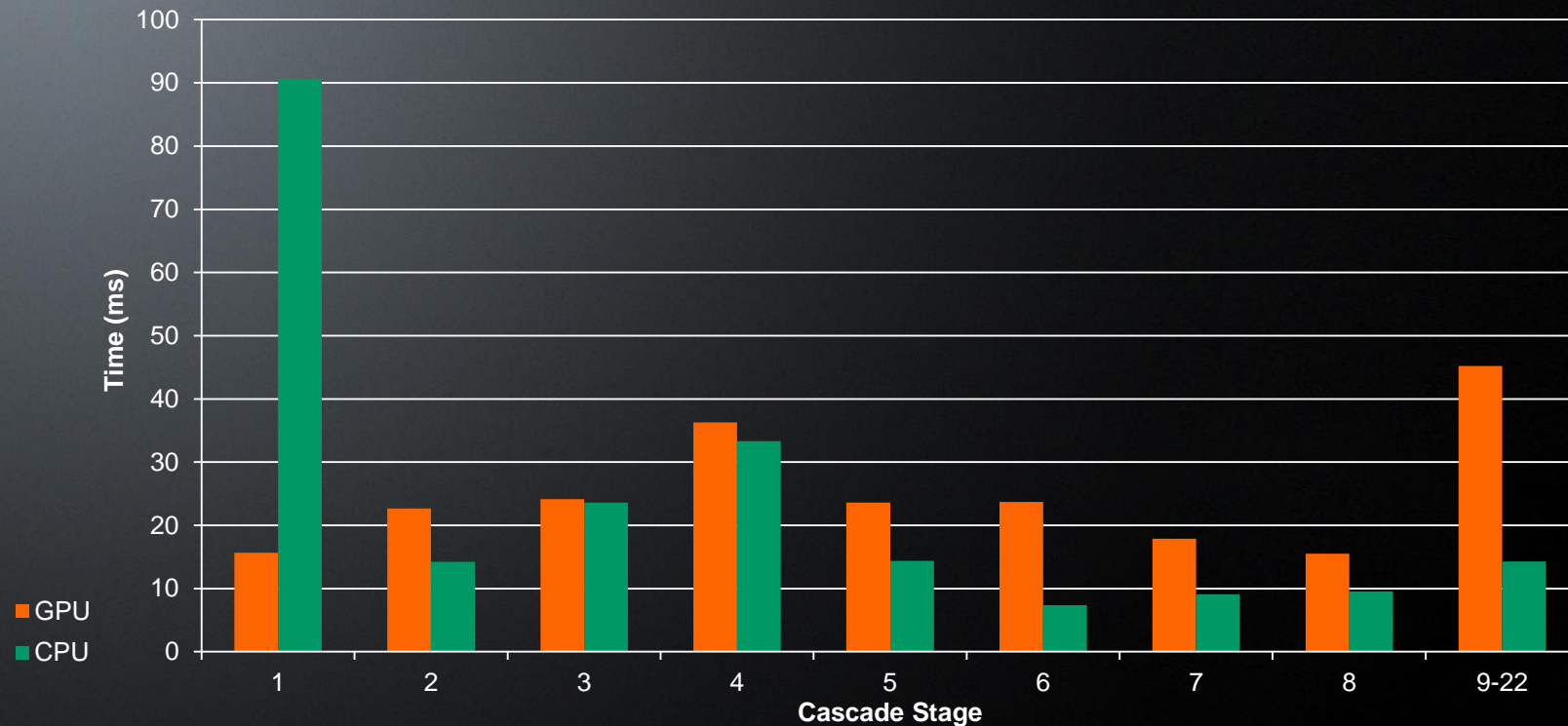
Live
Dead

1    2    3    4

- When running on the GPU, we run each search rectangle on a separate work item

- Early out algorithms, like HAAR, exhibit divergence between work items

  – Some work items exit early

  – Their neighbors continue

  – SIMD packing suffers as a result

AMD

# PROCESSING TIME/STAGE



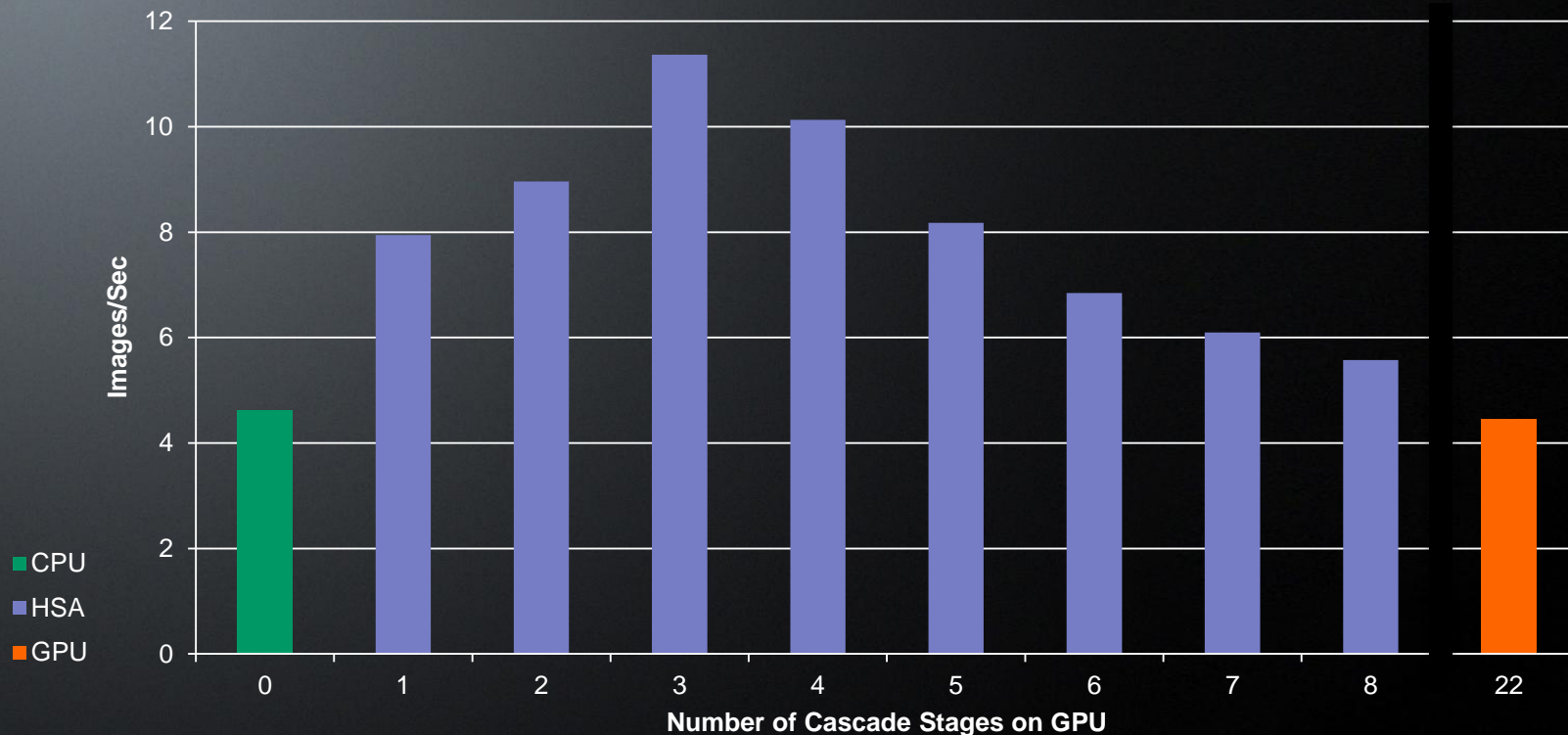"Trinity" A10-4600M (6CU@497Mhz, 4 cores@2700Mhz)

AMD A10 4600M APU with Radeon™ HD Graphics; CPU: 4 cores @ 2.3 MHz (turbo 3.2 GHz); GPU: AMD Radeon HD 7660G,
6 compute units, 685MHz; 4GB RAM; Windows 7 (64-bit); OpenCL™ 1.1 (873.1)

# PERFORMANCE CPU-VS-GPU

## "Trinity" A10-4600M (6CU@497Mhz, 4 cores@2700Mhz)



Bar chart. Y-axis: Images/Sec (0 to 12). X-axis: Number of Cascade Stages on GPU.

Legend:
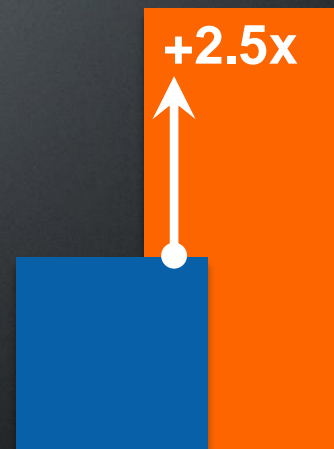- CPU (green)
- HSA (purple)
- GPU (orange)

AMD A10 4600M APU with Radeon™ HD Graphics; CPU: 4 cores @ 2.3 MHz (turbo 3.2 GHz); GPU: AMD Radeon HD 7660G, 6 compute units, 685MHz; 4GB RAM; Windows 7 (64-bit); OpenCL™ 1.1 (873.1)
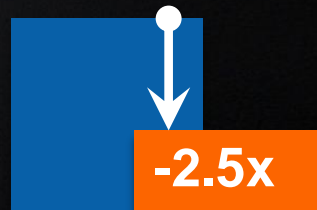
AMD

# HAAR SOLUTION – RUN DIFFERENT CASCADES ON GPU AND CPU

By seamlessly sharing data between CPU and GPU,
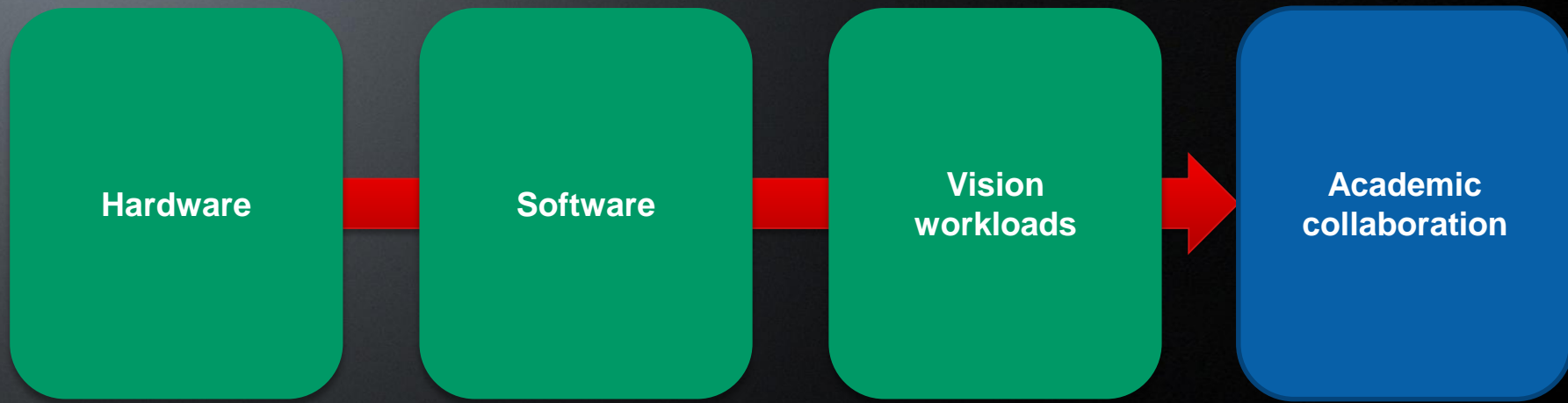HSA allows the right processor to handle its appropriate workload

**+2.5x**

**-2.5x**

**INCREASED
PERFORMANCE**

**DECREASED ENERGY
PER FRAME**

AMD

Hardware → Software → Vision workloads → Academic collaboration

University Collaboration Projects
SLAM: Simultaneous Localisation and Mapping

Collaborators:
A. Davison, P. Kelly, R. S. Moreno et al.
Imperial College London

**AMD**

# *SLAM*

- SLAM: Simultaneous Localisation and Mapping

  - Given an new scene, camera tries to generate a map of the environment and track its own position from it

- Previous SLAM systems used low-level map elements like points or fiducial markers.

- We now use higher-level Objects

- Objects provides more meaningful scene understanding

  - e.g. With objects, a table detection can 'infer' the floor it lies on.

  - Natural occlusion handling

  - Physical predictions: Constrain objects to lie on the floor plane → improves detection quality

- Reduced complexity on map optimisation (#Objects <<< # Points)

- More realistic Augmented Reality applications

AMD

# *SLAM | OBJECT LEVEL*

- We construct a database of scanned objects

- At runtime objects are detected on the GPU

  - Generating a map of objects locations/orientations

  - The map is immediately used to track the pose of the camera in real-time

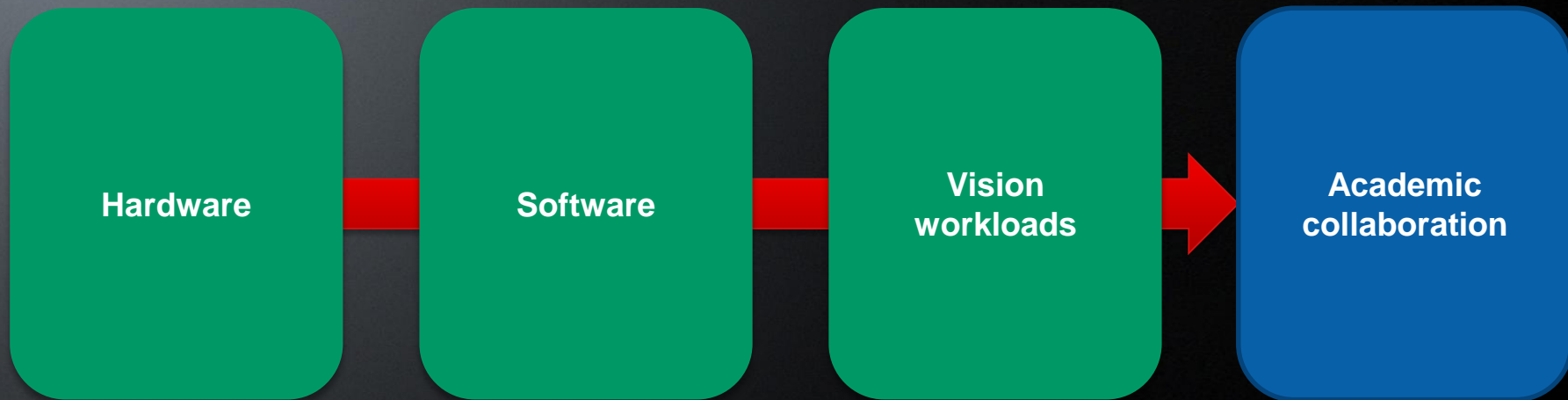- A batch optimisation routine executes on CPU to eliminate drift in camera motion

AMD

# *Real-time Depth Maps from Video*

- A pose of a monocular camera is first tracked with PTAM on the CPU
  - [Klein et al. ISMAR 08]
- Given a set of frames with known camera poses we build a depth map on the GPU
  - Camera must move to provide disparity information to figure depth
  - Generated depth maps used as a first step towards a full 3D object reconstruction pipeline
- Uses photo-consistency across frames
  - Assume lighting remains constant with small camera motion
- Look for a similar looking pixel of Frame A in Frames B, C, D, E ...
- True depth of pixel should generated similar intensities over frames
- Regularise solution to return smooth depths over large pixel areas
  - NewCombe et al. ICCV 2011

AMD

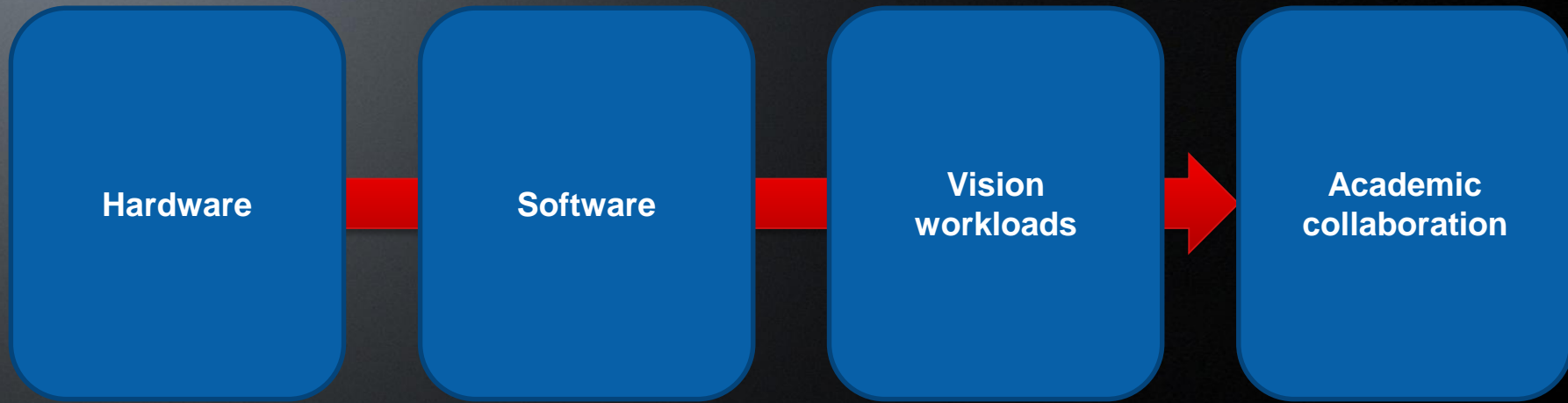# CBIR | Image Re-ranking and Rank Aggregation



Contextual rank aggregation algorithm

"Efficient Image Re-ranking and Rank Aggregation Computation on GPUs", 10th IEEE International Symposium on Parallel and Distributed Processing with Applications, July, 2012
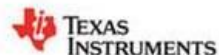
AMD Opteron 6168 1.9GHz, 12 Cores; ATI FirePro V7800; Ubuntu 10.04; AMD-APP-SDK-v2.4

# SUMMARY

Hardware → Software → Vision workloads → Academic collaboration

AMD

# 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, the AMD arrow logo, the HSA logo, and combinations thereof are trademarks of Advanced Micro Devices, Inc. OpenCL™ is a trademark of Apple Corp. which is licensed to the Khronos Organization. All other names used in this presentation are for informational purposes only and may be trademarks of their respective owners.

© 2012 Advanced Micro Devices, Inc.