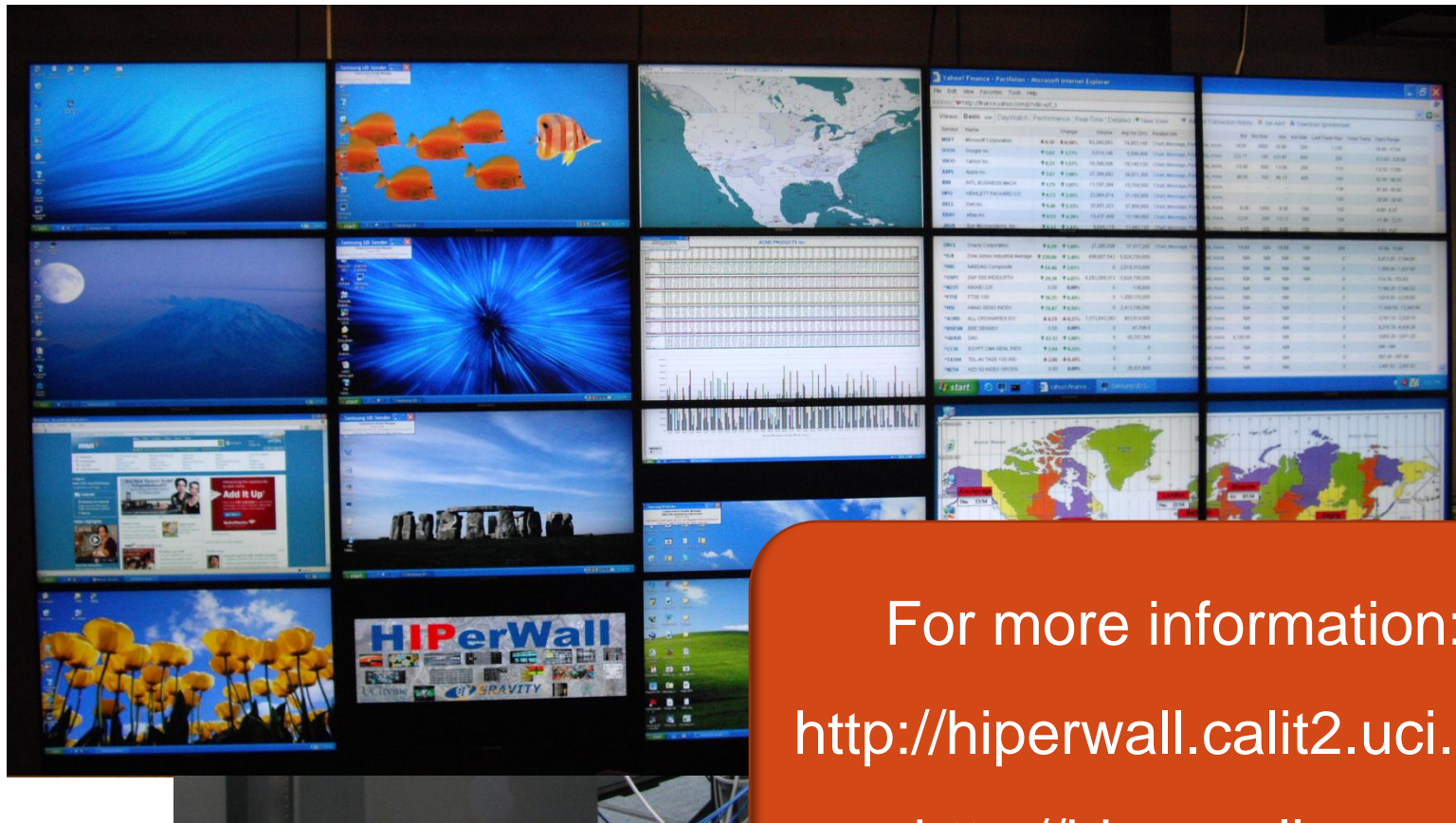# Ubiquitous Parallel Computing:
## Architectures and Programming from GPUs to the Cloud

Prof. Stephen Jenks

UCI Electrical Engineering
and Computer Science

& Hiperwall, Inc.

stephen@stephenjenks.com

# Hiperwall Overview



For more information:

http://hiperwall.calit2.uci.edu

http://hiperwall.com

9/28/2009

# Outline

- Parallelism Overview
  - Thread and OpenMP Programming
- Parallel Microprocessors (Multi-core)
  - Architecture Bottlenecks and Solutions
- Asymmetric Parallel Computing
  - Cell Processor
  - GPUs
    - CUDA and OpenCL
    - Examples and Performance Benefits
- Cloud Computing with Hadoop

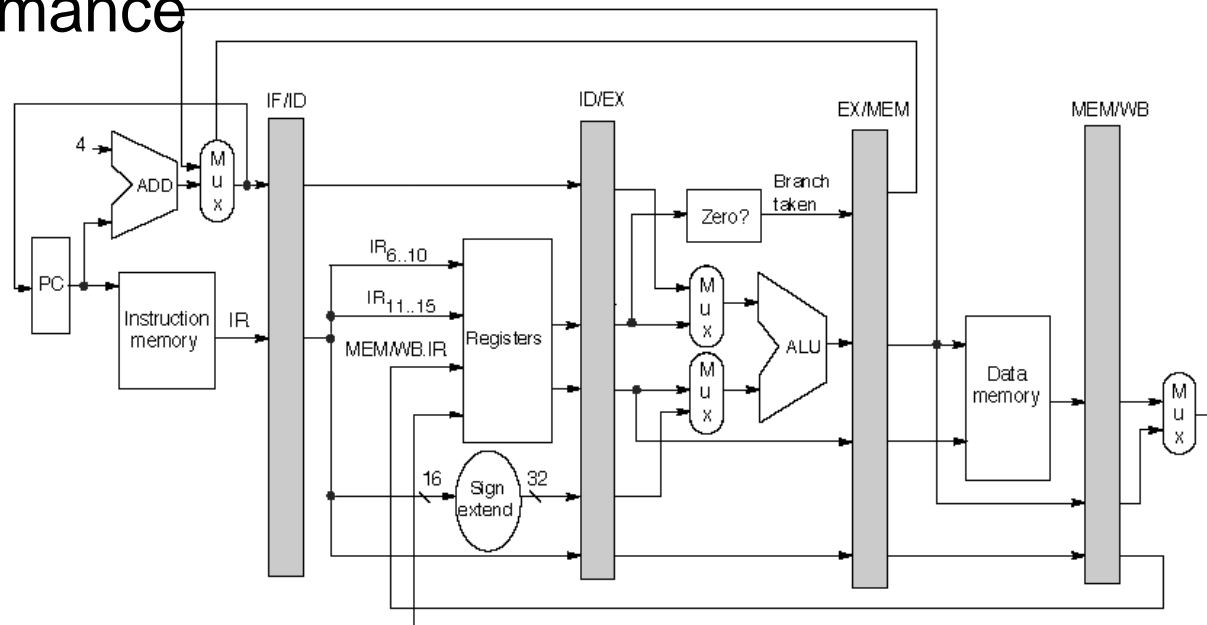9/28/2009

# Overview and Motivation

- Conventional Parallel Computing Was:
  - For scientists with large problem sets
  - Complex and difficult
  - Very expensive computers with limited access
- Parallel Computing is Becoming:
  - Ubiquitous
  - Cheap
  - Essential
  - Different
  - Still complex and difficult
- Where Is Parallel Computing Going?

# Computing is Changing

- Parallel programming is essential
  - Clock speed not increasing much
  - Performance gains require parallelism
- Parallelism is changing
  - Special purpose parallel engines
  - CPU and parallel engine work together
  - Different code on CPU & parallel engine → Asymmetric Computing

9/28/2009

# Conventional Processor Architecture

- Hasn't Changed Much For 40 Years
  - Pipelining and Superscalar Since the 1960s
  - But has become integrated → Microprocessors
- High Clock Speed
  - Great performance
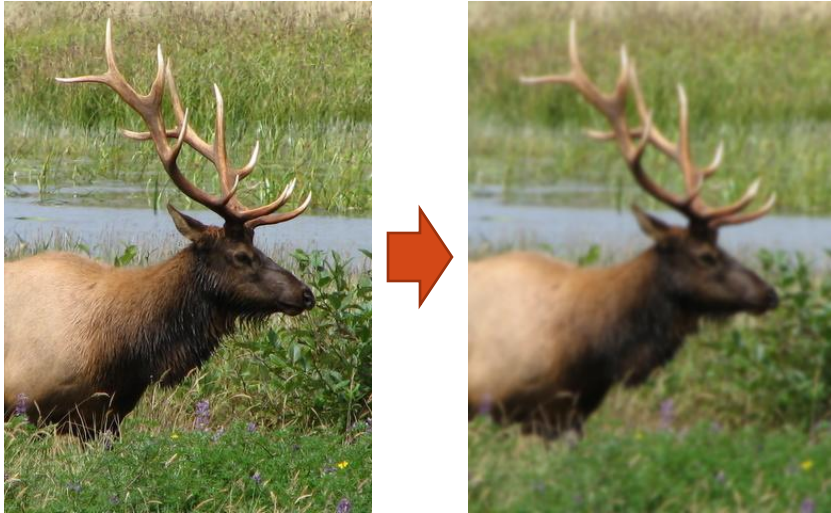  - High Power
  - Cooling Issues
  - Various Solutions

From Hennessy & Patterson, 2nd Ed. 9/28/2009
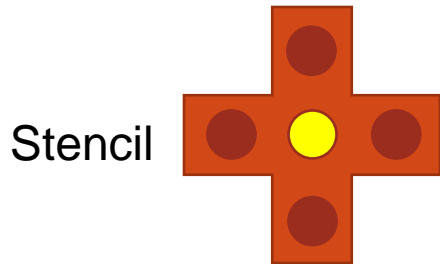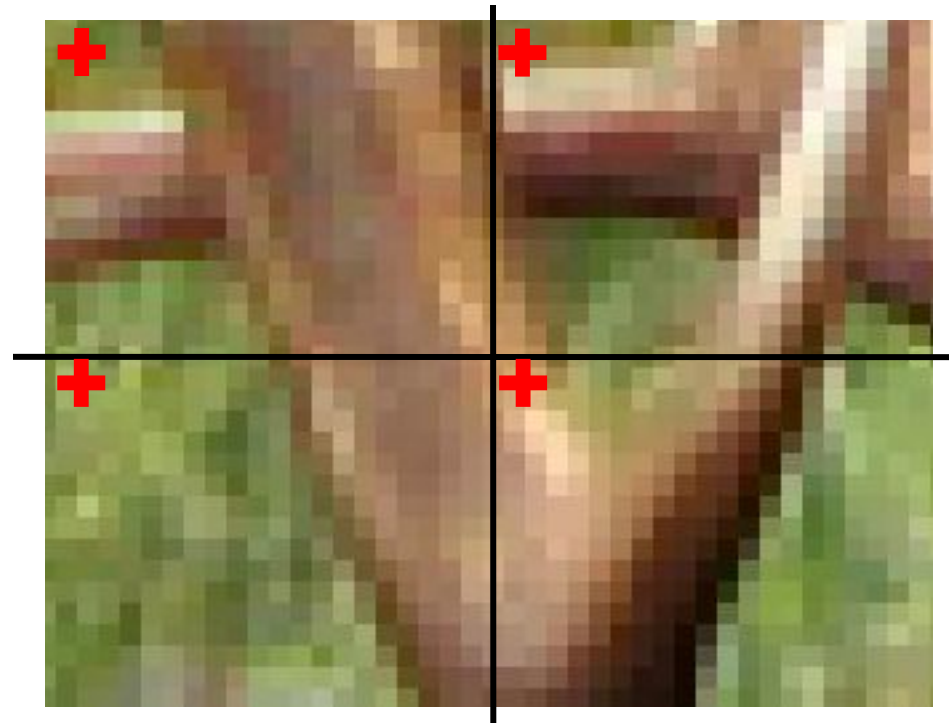
# Parallel Computing Problem Overview



Image Relaxation (Blur)


Stencil

newimage[i][j] = (image[i][j] +
image[i][j-1] + image[i][j+1] +
image[i+1][j] + image[i-1][j]) / 5

9/28/2009

# Shared Memory Multiprocessors

CPUs see common address space

| CPU | CPU | CPU | CPU |
|-----|-----|-----|-----|
| Cache | Cache | Cache | Cache |

**Shared Memory**



- Each CPU computes results for its partition
- Memory is shared so dependences satisfied

9/28/2009

# Shared Memory Programming

- Threads
  - POSIX Threads
  - Windows Threads

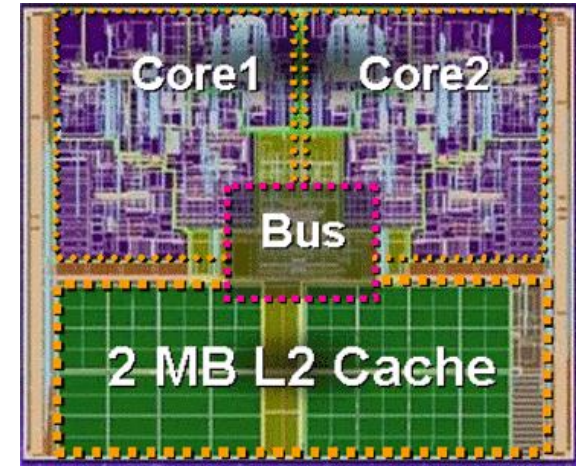Not Integrated with Compiler or Language
→No idea if code is in thread or not
→Poor optimizations

- OpenMP
  - User-inserted directives to compiler
  - Loop parallelism
  - Parallel regions
  - Visual Studio
  - GCC 4.2

```
#pragma omp parallel for
for (i=0; i<n; i++)
  a[i] = b[i] * c[i];
```
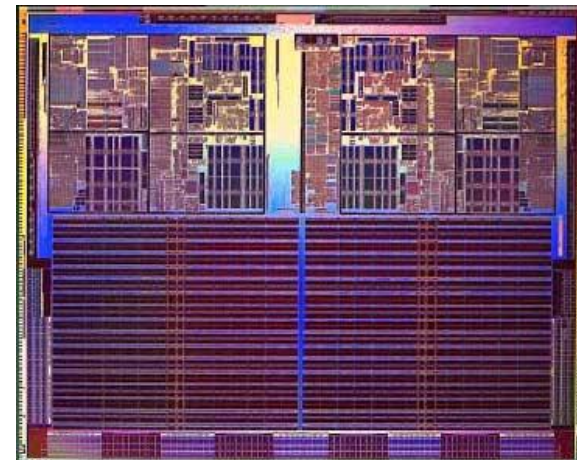
9/28/2009

# Multicore Processors

- Several CPU Cores Per Chip
- Shared Memory
- Shared Caches (sometimes)
- Lower Clock Speed
  - Lower Power & Heat
  - But Good Performance
- Program with Threads
- Single Threaded Code
  - Not Faster (except on Core i7)
  - Majority of Code Today



Intel Core Duo



AMD Athlon 64 X2

9/28/2009

# Conventional Processors are Dinosaurs

- So much circuitry dedicated to keeping ALUs fed:
  - Cache
  - Out-of-order execution/reorder buffer
  - Branch prediction
  - Large Register Sets
  - Simultaneous Multithreading
- ALU (Arithmetic Logic Unit) tiny by comparison
- Huge power for little performance gain

With thanks to Stanford's Pat Hanrahan for the analogy

9/28/2009

# AMD Phenom X4 Floorplan

9/28/2009

Source: AMD.com

# Parallel Microprocessor Problems



- Memory interface too slow for 1 core/thread
- Now multiple threads access memory simultaneously, overwhelming memory interface
- Parallel programs can run as slowly as sequential ones!

9/28/2009

# SPPM: Producer/Consumer Parallelism Using The Cache

| Data in Memory | | Data in Memory | |
|---|---|---|---|

Memory Bottleneck

Communications Through Cache

| Thread 1 Half the Work | Thread 2 Half the Work | Producer Thread | Consumer Thread |
|---|---|---|---|

9/28/2009

# Multicore Architecture Improvements

- Cores in Common Multicore Chips Not Well Connected
  - Communicate Through Cache & Memory
  - Synchronization Slow (OS-based) or Uses Spin Waits
- Register-Based Synchronization
  - Shared Registers Between Cores
  - Halt Processor While Waiting To Save Power
- Preemptive Communications (Prepushing)
  - Reduces Latency over Demand-Based Fetches
  - Cuts Cache Coherence Traffic/Activity
- Software Controlled Eviction
  - Manages Shared Caches Using Explicit Operations
  - Move Data to Shared Cache Before Needed From Private Cache
- Synchronization Engine
  - Hardware-based multicore synchronization operations

9/28/2009

# Single Nehalem (Intel I7) Core



New SSE4.2 Instructions

Improved Lock Support

Superior Caching Hierarchy

Execution Units

L1 Data Cache

L2 Cache & Interrupt Servicing

Memory Ordering & Execution

Paging

Deeper Buffers

Instruction Reordering, Scheduling & Retirement

Instruction Decode & Microcode

Branch Prediction

Instruction Fetch & L1 Cache

Improved Loop Streaming

Simultaneous Multi-Threading

Better Branch Prediction

Memory Controller

Misc IO

Core

Core

Queue

Core

Core

Misc IO

QPI 0

Shared L3 Cache

QPI 1

*Nehalem - Everything You Need to Know about Intel's New Architecture*
http://www.anandtech.com/cpuchipsets/intel/showdoc.aspx?i=3382

9/28/2009

# Asymmetric Parallel Accelerators

- Current Cores are Powerful and General
  - Some Applications Only Need Certain Operations
  - Perhaps a Simpler Processor Could Be Faster
- Pair General CPUs with Specialized Accelerators
  - Graphics Processing Unit
  - Field Programmable Gate Array (FPGA)
  - Single Instruction, Multiple Data (SIMD) Processor

Different code runs on CPU cores and accelerators: Asymmetric Computing

| Athlon 64 CPU | ATI GPU |
|---|---|
| XBAR | |
| Hyper-Transport | Memory Controller |

Possible Hybrid AMD Multi-Core Design

9/28/2009

# Cell Broadband Engine

- PowerPC Processing Element with Simultaneous Multithreading at 3.2 GHz

- 8 Synergistic Processing Elements at 3.2 GHz
  - Optimized for SIMD/Vector processing (100 GFLOPS Total)
  - 256KB Local Storage - no cache

- 4x16-byte-wide rings @ 96 bytes per clock cycle

From *IBM Cell Broadband Engine Programmer Handbook*, 10 May 2006

9/28/2009

# NVIDIA GPU Floorplan

10 multiprocessors
24 threads each
240 simultaneous threads!

Source: Dr. Sumit Gupta - NVIDIA    9/28/2009

# Graphics Processing Unit (GPU)

- GPUs Do Pixel Pushing and Matrix Math

From
*NVIDIA CUDA*
*Compute Unified*
*Device Architecture*
*Programming Guide*
11/29/2007

# CUDA & OpenCL Programming Model

- Data Parallel
  - But not Loop Parallel
  - Very Lightweight Threads
  - Write Code from Thread's Point of View
- No Shared Memory
  - Host Copies Data To and From Device
  - Different Memories
- Hundreds of Parallel Threads (Sort-of SIMD)

Block of Threads

9/28/2009

# OpenCL Programming Details

- Supports GPUs (NVIDIA and ATI) and CPUs
- Built into Apple's Snow Leopard Mac OS X 10.6
- On-the-fly Compilation
- Supports floats (doubles optional and not supported on all hardware)
- Pattern:

**Overhead**

1. Acquire "device" and get capabilities
2. Initialize (compile) OpenCL "kernel"
3. Move data to "device" memory from "host" memory
4. Set kernel parameters and execution size, start kernel
5. When done, copy results back from device memory
6. Repeat prior 3 steps, as needed

9/28/2009

# Simple OpenCL Kernel Code

```
__kernel void openclmin(__global float *a,
        __global float *b, __global float *c)
{

    int gid = get_global_id(0);

    c[gid] = fmin(a[gid], b[gid]);
}
```

*Get Position In Global Index Space*

100 Iterations in OpenCL:
0.031 secs/iteration (or nearly
6 times faster than sequential!)

Problem size: 24 million element arrays
(or 3*24*4 = 288 Mbytes)



| | Sequential | OpenMP | OpenCL-1(1) | OpenCL-2(1) |
|---|---|---|---|---|
| CPU | 0.174 | 0.092 | 0.532 | 0.577 |
| GPU | | | 1.983 | 1.768 |

9/28/2009

# GPU Performance/Architecture Issues (CUDA and OpenCL)

- Avoid branching

```
__kernel void openclmin(__global float *a,
      __global float *b, __global float *c)
{

    int gid = get_global_id(0);
    float local_a = a[gid];
    float local_b = b[gid];
    if (local_a < local_b)
        c[gid] = local_a;
    else
        c[gid] = local_b;
}
```

Branch instead of fmin()

Runs 32% slower on GPU than before

- Memory is fast, but long latency
  - Cache data in shared space
  - Avoid bank conflicts
- Only new GPUs support Doubles
- Memory limited (256MB to 4GB vs 16 or 32 GB)

**GPU Computing is Great, but be aware of the limitations.**

9/28/2009

# GPU Usage on Hiperwall

- QuickTime renders movies as YUV frames
  - 16 bits per pixel rather than 32 for RGBA
  - Apple OpenGL supports YUV textures natively
- CPU takes 5-8ms to convert 720p frame to RGBA
- Solution GPU computing with Cg
  - Modern NVIDIA card – 100 times faster
  - Older ATI chip (shared memory) – 15 times faster
  - My laptop (Intel GPU w/ shared) – 3 times slower!

# GPU Programming Choices

- CUDA (NVIDIA)
  - Mature. Probably best tool support
  - Supported on Windows, Linux, Mac
  - Only on NVIDIA hardware
- OpenCL (Open standard coalition)
  - Stable, but less tool support (for now)
  - Built into Mac OS X 10.6, supported on Linux & Windows
  - ATI and NVIDIA hardware
- DirectCompute (Microsoft)
  - New, but drivers available
  - Built into DirectX 11 (no Mac or Linux support)
  - ATI and NVIDIA hardware

9/28/2009

# Intel Larrabee

| Simple x86 Core | Simple x86 Core | Simple x86 Core | • • • | Simple x86 Core | Simple x86 Core |
|---|---|---|---|---|---|

**Interprocessor Ring Network**

| Coherent L2 Cache | Coherent L2 Cache | Coherent L2 Cache | | Coherent L2 Cache | Coherent L2 Cache |
|---|---|---|---|---|---|
| Coherent L2 Cache | Coherent L2 Cache | Coherent L2 Cache | | Coherent L2 Cache | Coherent L2 Cache |

Memory & I/O Interfaces

| Simple x86 Core | Simple x86 Core | Simple x86 Core | • • • | Simple x86 Core | Simple x86 Core |
|---|---|---|---|---|---|

## Many simple, fast, low power, in-order x86 cores

*Larrabee: A Many-Core x86 Architecture for Visual Computing*
ACM Transactions on Graphics, Vol. 27, No. 3, Article 18, Publication date: August 2008.
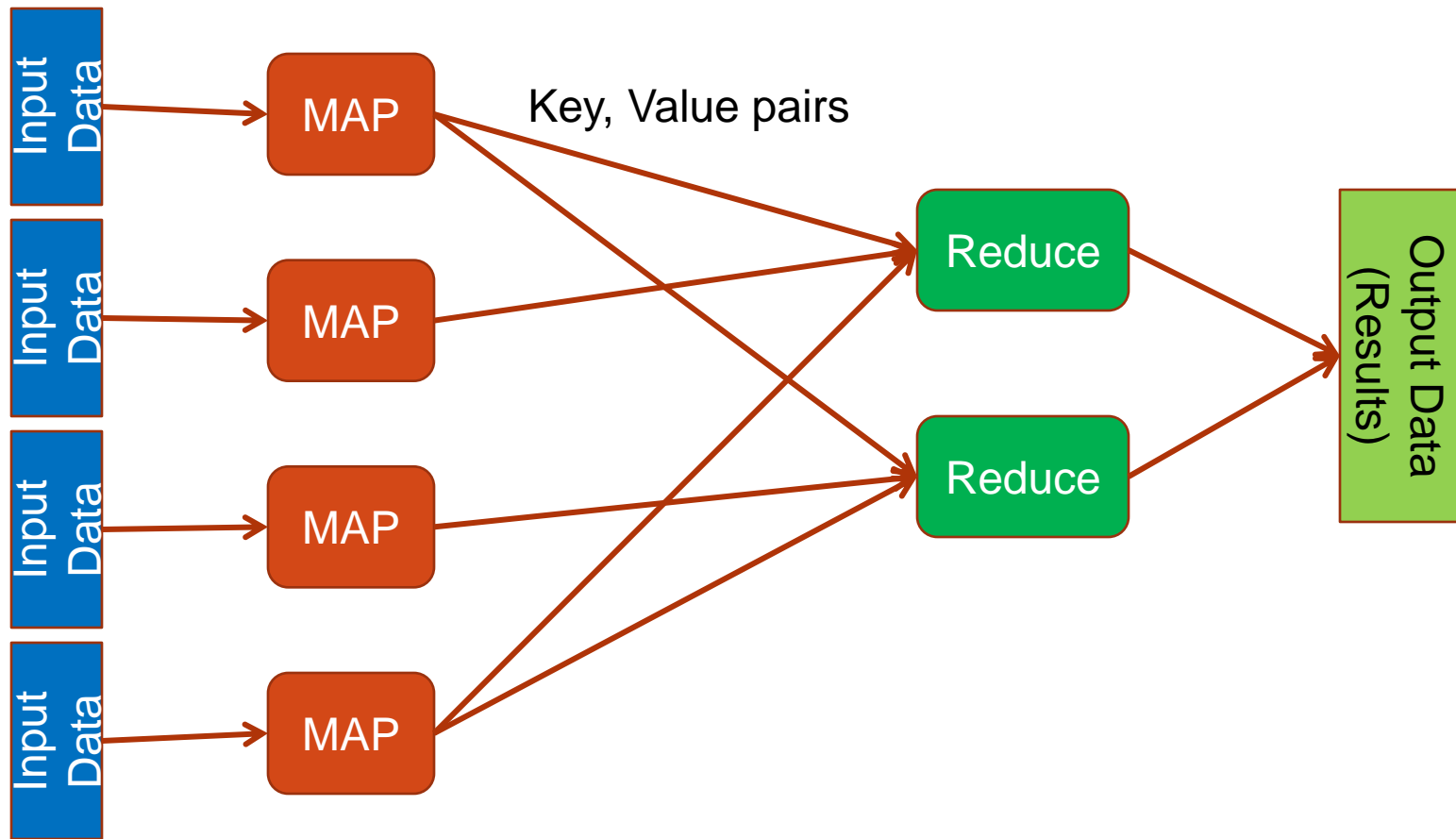
9/28/2009

# Cloud Computing

- What is cloud computing?
  - Latest buzzword in computing
  - Replaces Grid, Utility Computing, … as latest craze
- Multiple definitions
  - Web-based applications (Google Docs)
  - On-demand Computing Resources
    - Virtual Machines in a server farm (Amazon, Google, IBM)
    - Parallel/Distributed Computing Paradigm to use them

9/28/2009

# Cloud-Based Map-Reduce

- Special purpose computation with LOTS of data
  - Used by Google and many others
  - Based on Lisp's Map and Reduce functions
- Examples: Distributed Grep, Count of URL Access Frequency, Reverse Web-Link Graphs, Term-Vector per Host, Inverted Index, Distributed Sort
  - Most produce small results from large input
  - Simple computation per element, but lots of them
- Open source implementation: Hadoop
  - Yahoo and Apache
  - Java-based, includes distributed file system

# Cloud (Hadoop) Application Topology



Input Data → MAP

Input Data → MAP

Input Data → MAP

Input Data → MAP

Key, Value pairs

Reduce

Reduce

Output Data (Results)

9/28/2009

# Hadoop Word Count Example

Input Data → MAP → Key, Value pairs

Input Data → MAP

Input Data → MAP → Reduce

Input Data → MAP → Reduce → Output Data (Results)

Book text

For each word, emit <word, 1> pair

Partitioner sends data to right Reduce based on key (hash or alpha)

Count how many of each word received (add up the "values" for each key)

Table of words and their counts

9/28/2009

# Summary

- Parallel Computing will soon be required for good performance
  - Parallel programming is neither free nor easy
  - New tools make it better than before
  - Architecture influences performance
  - Parallelism at various layers: chip to cloud
    - Thread parallelism
    - Data-parallelism (top-down view): OpenMP
    - Data-parallelism (bottom-up view): GPU programming
    - Distributed memory parallelism: Hadoop

# Resources

- HIPerWall: http://hiperwall.calit2.uci.edu/
- Hiperwall, Inc.: http://hiperwall.com/
- Vadlamani, S. & Jenks, S. "Architectural Considerations for Efficient Software Execution on Parallel Microprocessors," *21st IEEE International Parallel & Distributed Processing Symposium,* **2007**
- Fide, S. & Jenks, S. "Architecture Optimizations for Synchronization and Communication on Chip Multiprocessors," *Workshop on Multithreaded Architectures and Applications (MTAAP08) Held in Conjunction With International Parallel and Distributed Processing Symposium (IPDPS 2008),* **2008**
- Fide, S. & Jenks, S. "Proactive Use of Shared L3 Caches to Enhance Cache Communications in Multi-Core Processors," *IEEE Computer Architecture Letters,* **2008**
- CUDA: http://www.nvidia.com/object/cuda_home.html#
- OpenCL: http://www.khronos.org/opencl/
- OpenMP: http://openmp.org/wp/
- Hadoop: http://hadoop.apache.org/
- OpenCL Tutorials (David Gohara): http://www.macresearch.org/opencl
- Di Blas, A. & Kaldewey, T. Data Monster: Why graphics processors will transform database processing. *IEEE Spectrum,* **2009***, 46*, 46-51

9/28/2009