

The Parallel Computing Laboratory

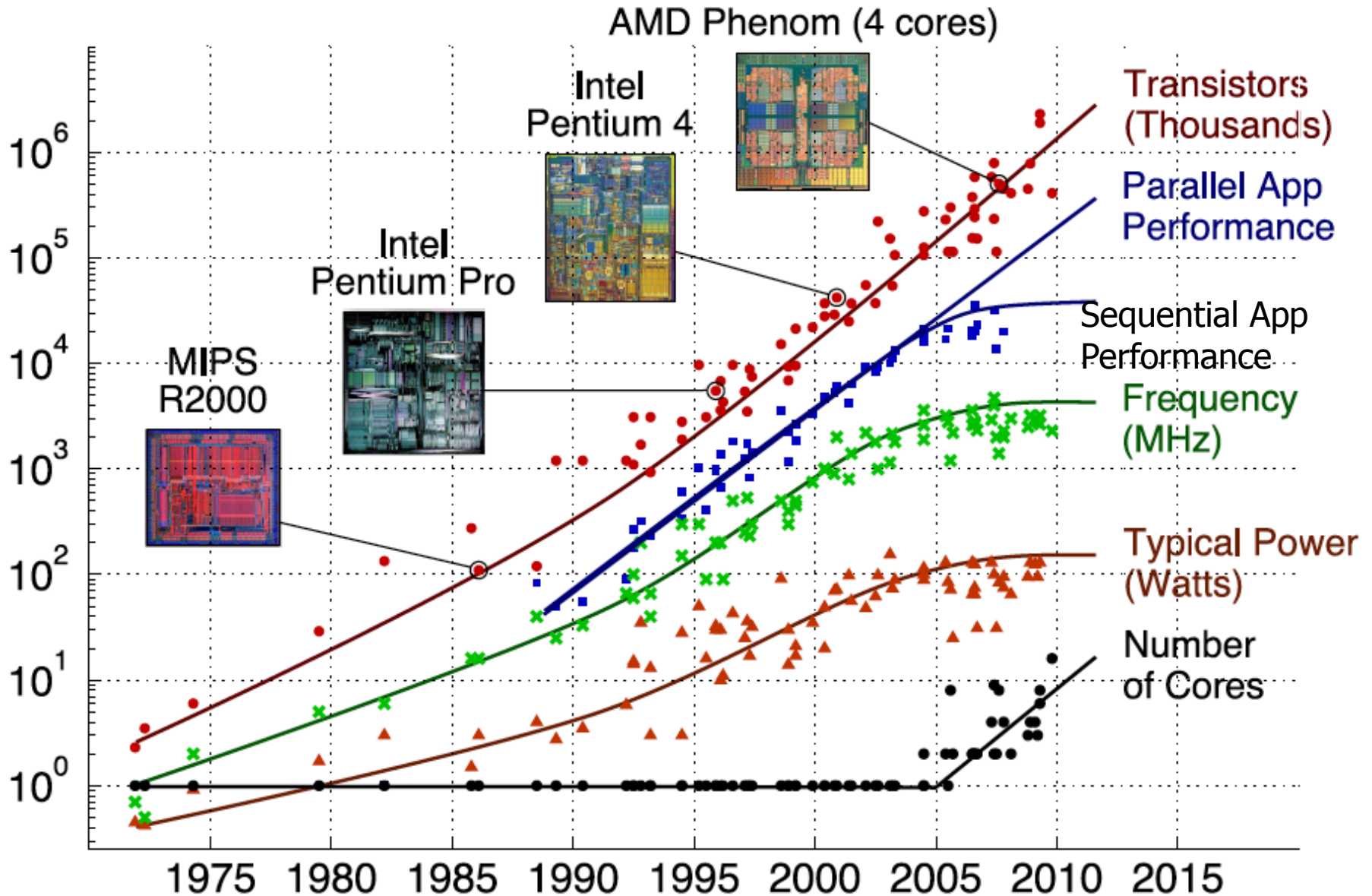
Krste Asanovic, Ras Bodik,
Jim Demmel, Armando Fox, Tony Keaveny,
Kurt Keutzer, John Kubiawicz,
Nelson Morgan, Dave Patterson, Koushik Sen,
David Wessel, and Kathy Yelick
UC Berkeley

Microsoft Faculty Summit

Bellevue, WA

July 17, 2012

Transition to Multicore



- Berkeley researchers from many backgrounds meeting since Feb. 2005 to discuss parallelism
 - Krste Asanovic, Eric Brewer, Ras Bodik, Jim Demmel, Kurt Keutzer, John Kubiawicz, Dave Patterson, Koushik Sen, Kathy Yelick, ...
 - Circuit design, computer architecture, massively parallel computing, computer-aided design, embedded hardware and software, programming languages, compilers, scientific programming, and numerical analysis
 - Tried to learn from successes in high-performance computing (LBNL) and parallel embedded (BWRC)
- Led to “Berkeley View” Tech. Report 12/2006 and new Parallel Computing Laboratory (“Par Lab”)

Goal: To enable most programmers to be productive writing efficient, correct, portable SW for 100+ cores & scale as cores increase every 2 years (!)

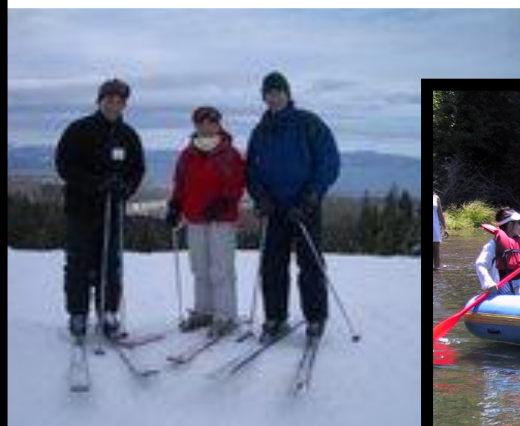
- ❖ Past parallel projects often dominated by hardware architecture:
 - *This is the one true way to build computers, software must adapt to this breakthrough!*
 - E.g., ILLIAC IV, Thinking Machines CM-2, Transputer, Kendall Square KSR-1, Silicon Graphics Origin 2000 ...
- ❖ Or sometimes by programming language:
 - *This is the one true way to write programs, hardware must adapt to this breakthrough!*
 - E.g., Id, Backus Functional Language FP, Occam, Linda, HPF, Chapel, X10, Fortress ...
- ❖ Applications usually an afterthought

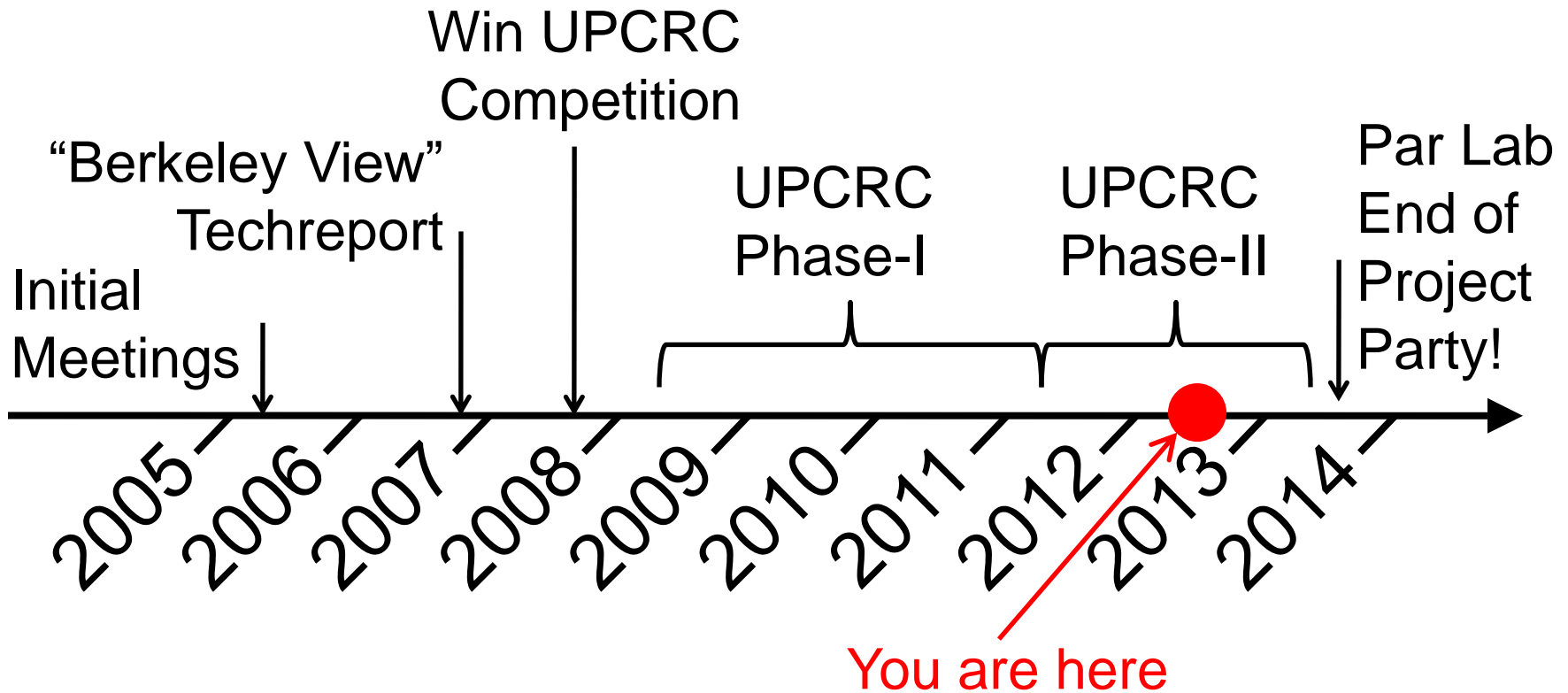
- ❖ ***Let compelling applications drive research agenda***
 - ❖ Software platform: data center + mobile client
 - ❖ Identify common programming patterns
 - ❖ Productivity versus efficiency programmers
 - ❖ Autotuning and software synthesis
 - ❖ Build-in correctness + power/performance diagnostics
 - ❖ OS/Architecture support applications, provide flexible primitives not pre-packaged solutions
 - ❖ FPGA simulation of new parallel architectures: RAMP
 - ❖ Co-located integrated collaborative center
- Above all, no preconceived big idea - see what works driven by application needs.*

Co-located Collaborative Center Approach



- ❖ 60+ students, 8+ faculty in one shared space
- ❖ Faculty in open space, not in offices
- ❖ Off-site retreat every 6 months with ~60 outside visitors (industry sponsors, and other invited experts)





- ❖ Patterns for parallel programming
- ❖ Communication-avoiding algorithms
- ❖ Specializers: Pattern-specific compilers
- ❖ Effective composition of parallel modules



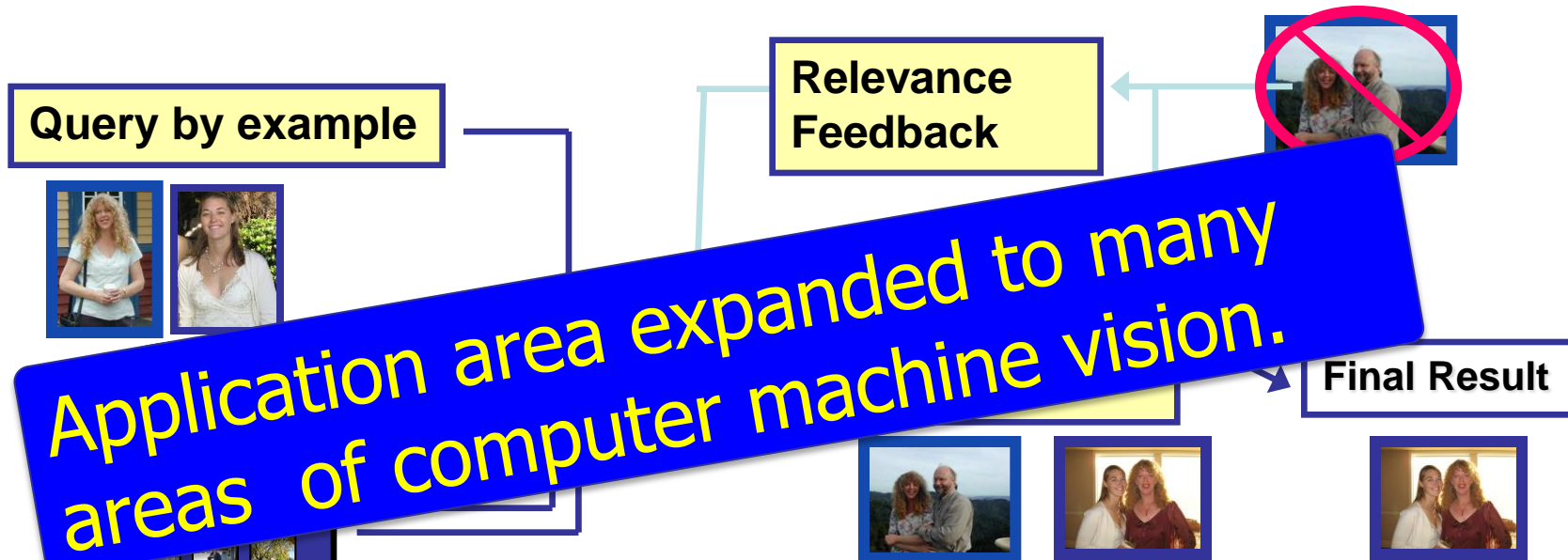
Controversial in 2005,
"Obvious" in 2012

- ❖ Laptop/Handheld ("Mobile Client")
 - Par Lab focuses on mobile clients
- ❖ Data Center or Cloud ("Cloud")
 - RAD Lab/AMPLab focuses on Cloud
- ❖ Both together ("Client+Cloud")
 - ParLab-AMPLab collaborations

Content-Based Image Retrieval

(Kurt Keutzer)

Electrical Engineering and
Computer Sciences



1000's of
images

❖ Built around Key Characteristics of personal databases

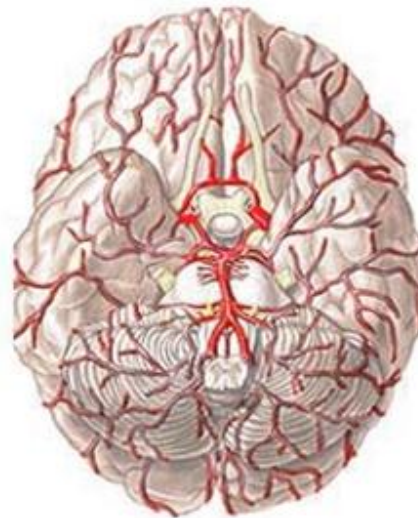
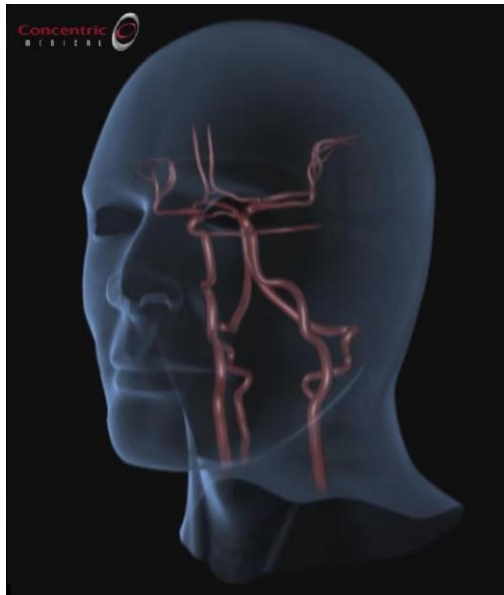
- Very large number of pictures (>5K)
- Non-labeled images
- Many pictures of few people
- Complex pictures including people, events, places, and objects



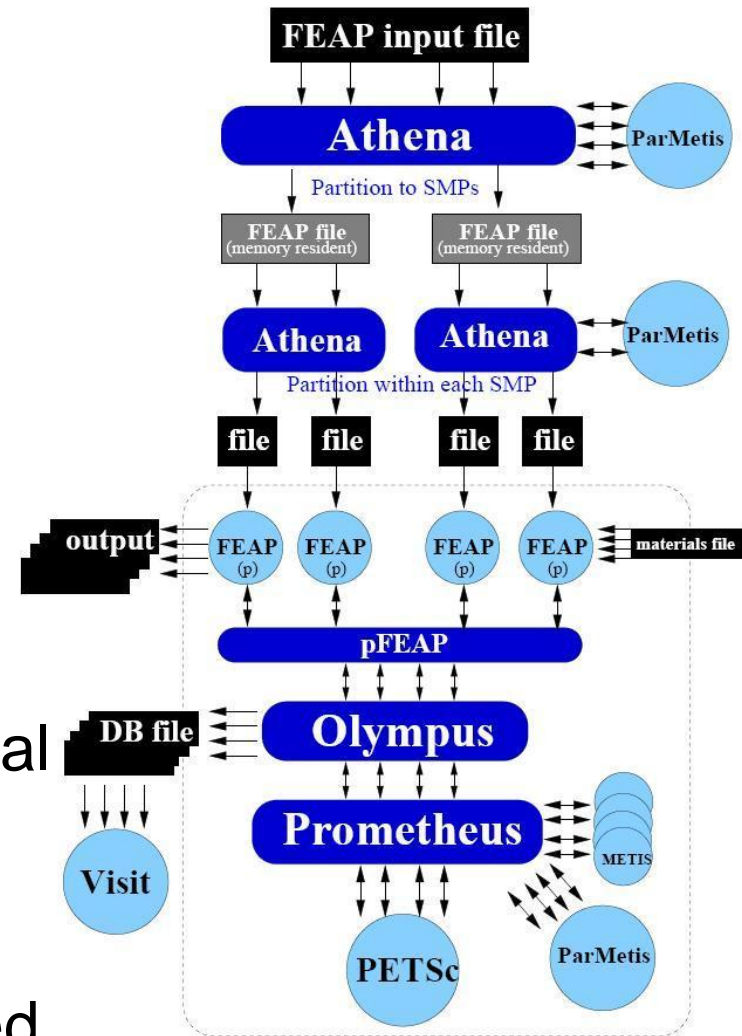
Health Application: Stroke Treatment

(Tony Keaveny, ME@UCB)

Electrical Engineering and
Computer Sciences



Bottom view of brain
© ADAM, Inc.



- Stroke treatment time-critical, need supercomputer performance in hospital
- Goal: 1.5D Fluid-Solid Interaction analysis of Circle of Willis (3D vessel geometry + 1D blood flow).
- Based on existing codes for distributed clusters

Parallel Browser (Ras Bodik)

Old Joe's Showhouse [Add to My Favorite Theaters] Theater Info | Map It

17 Main St., Lilliput

Grapes of Khan, The Rated PG-13, 1 hr 22 min Showtimes: 11:00, 4:15, 5:20, 7:45, 9:55, 10:10	Rainman Forever Rated R, 1 hr 33 min Showtimes: 2:15, 4:45, 7:15, 9:35
Rent and Rentability Rated PG-13, 1 hr 43 min Showtimes: 1:15, 5:15, 9:30	Rent and Rentability Rated PG, 1 hr 43 min Showtimes: 7:00, 9:30, 11:30

Little-End Cinemas [Add to My Favorite Theaters] Theater Info | Map It

47 Main St., Lilliput

Die Hard With More Intensity Rated R, 1 hr 47 min Showtimes: 11:55, 1:15, 2:30, 3:50, 5:05, 6:25, 7:40, 9:05 10:15	Hairy Plumber and the Goomba of Doom Rated PG, 2 hr 10 min Showtimes: 11:30, 1:35, 3:40, 5:45, 7:50, 10:00
Rainman Forever Rated R, 1 hr 33 min Showtimes: 2:15, 4:45, 7:15, 9:35	Rent and Rentability Rated PG, 1 hr 43 min Showtimes: 7:00, 9:30, 11:30

**Readable
Layouts**

★★★★★ Rainman Forever
Action: An autistic man fights crime on the streets of Gotham.
Uwe Boll, Dustin Hoffman, Jim Carrey, Jet Li
"137 interminable minutes. I counted them." (Ebert)
"He's an excellent driver in a terrible movie." (Boston Sun)
"A flop. Definitely a flop." (SF Chronicle)

★★★★★ Die Hard With More Intensity
Drama: A streetwise cop confronts loneliness in Tokyo.
Sofia Coppola, Bill Murray, Bruce Willis, Jet Li
"Extraordinarily powerful ... A masterpiece of cinema." (Ebert)
"Beautiful and haunting." (filmscritic.com)
"Moves slow but packs a punch." (Boston Sun)

★★★★★ Hairy Plumber and the Goomba of Doom
Fantasy: Mario and Luigi attend a school of wizardry.
Steven Spielberg, Daniel Radcliffe, Emma Watson, Jet Li
"Not as good as the others, but still a visual treat." (Ebert)
"The boys are back and better than ever." (filmscritic.com)
"Fans of the series won't be disappointed." (Boston Sun)

★★★★★ Rent and Rentability
Romance: Dissimilar sisters seek husbands in the East Village.
Ang Lee, Emma Thompson, Kate Winslet, Jet Li
"A poor adaptation of the Broadway hit." (Ebert)
"A real tear-jerker. Keep your hanky handy!" (filmscritic.com)

★★★★★ The Little Schemer
Adventure: An elephant journeys to find Lambda the Ultimate.
Friedman & Felleisen, Car, Cdr, Cons, Cond
"Cons is magnificent! ... Add this movie to your list!" (Ebert)

Old Joe's Showhouse	11:55	3:00	7:15
AMC Lilliput	1:25	2:45	4:20
UA Easy Street	12:45	3:00	5:00
Gulliver Theater	1:40	4:25	7:20
Little-End Cinemas	2:15	4:45	7:15

AMC Lilliput	11:45	12:40	2:00	3:30	5:00	6:30	7:30	8:30
UA Easy Street	1:00	2:30	4:30	7:00	8:50			
Landmark Quinbus			4:30	8:00				
Little-End Cinemas	11:55	1:15	2:30	3:50	5:05	6:25	7:40	9:00

AMC Lilliput	12:00	2:25	4:45	7:05
UA Easy Street	12:40	3:55	7:15	
Landmark Quinbus	12:05	2:45		
Gulliver Theater				8:30
Little-End Cinemas	11:30	1:35	3:40	5:45

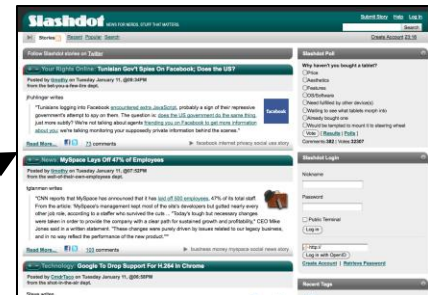
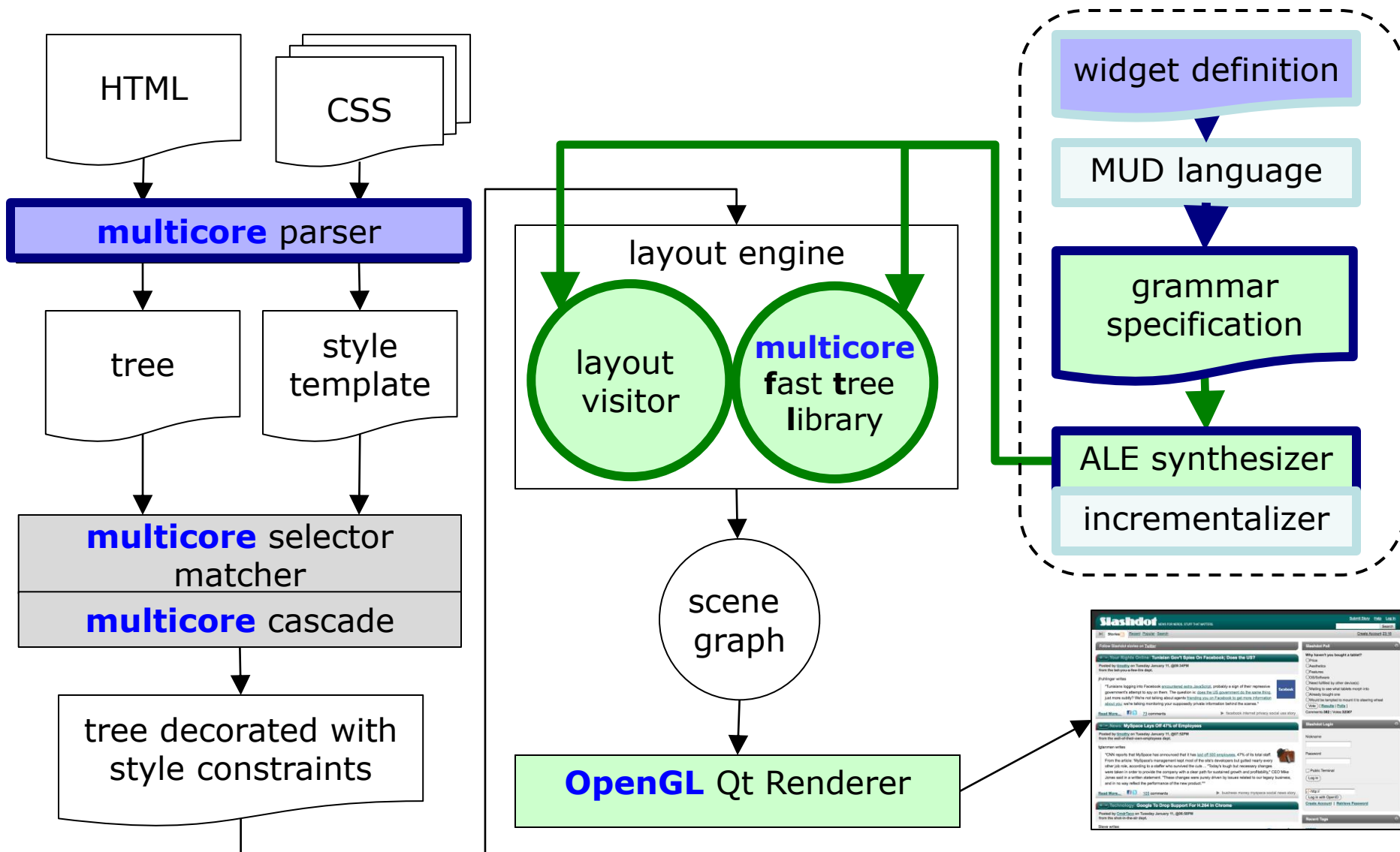
Old Joe's Showhouse	1:15	5:15
Gulliver Theater	12:35	2:20
Little-End Cinemas	4:55	8:20

PLT Arthouse	7:30
--------------	------



Augmented Reality

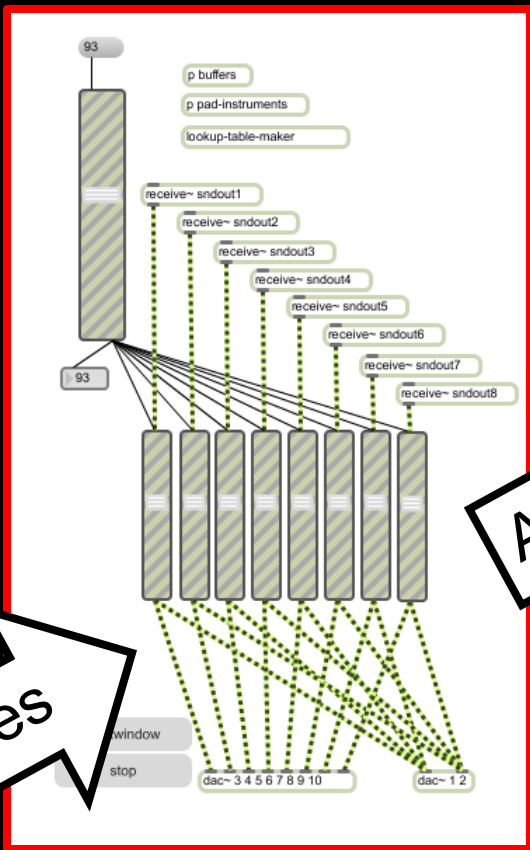
- ❖ Original goal: Desktop-quality browsing on handhelds (Enabled by 4G networks, better output devices)
- ❖ Now: Better development environment for new mobile-client applications, merging characteristics of browsers and frameworks (Silverlight, Qt, Android)



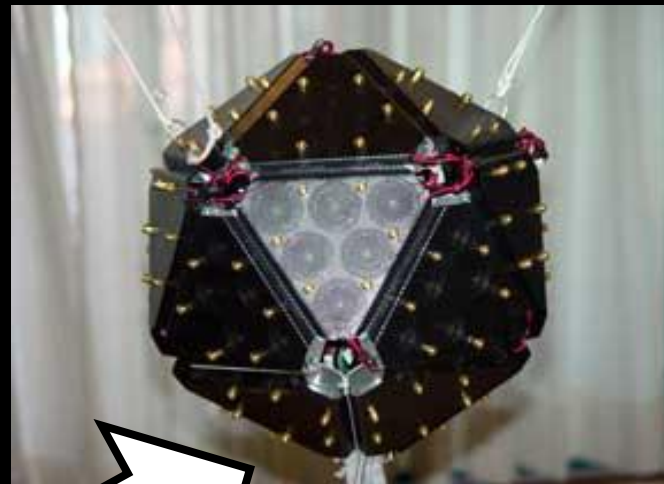
New user interfaces
with pressure-sensitive
multi-touch gestural
interfaces



Gestures

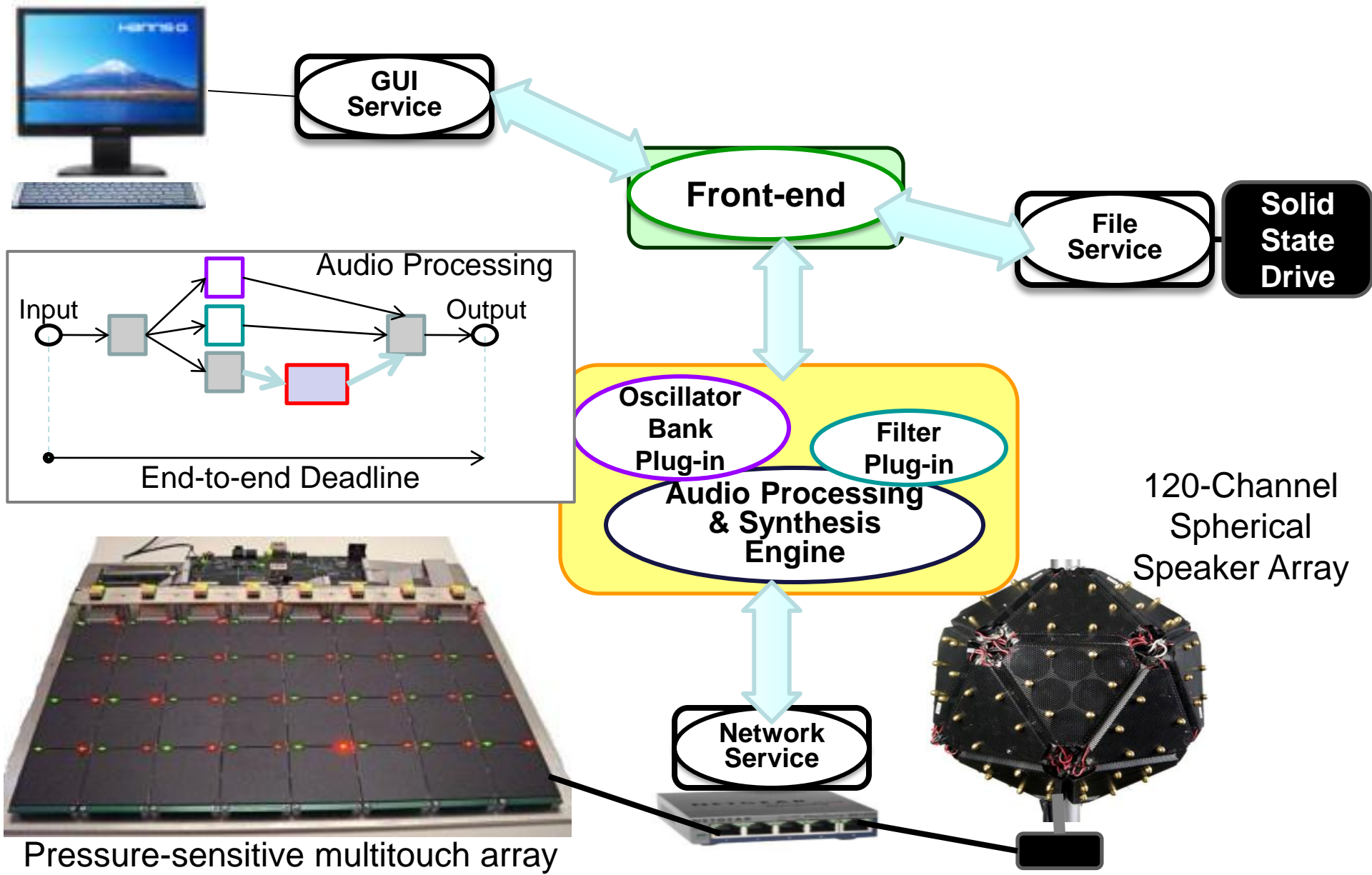


Audio



120-channel
speaker array

Programmable virtual instrument
and audio processing



Pressure-sensitive multitouch array

120-Channel
Spherical
Speaker Array

Speech: Meeting Diarist

(Nelson Morgan, Gerald Friedland, ICSI/UCB)



- Laptops/ Handhelds at meeting coordinate to create speaker identified, partially transcribed text diary of meeting

Winner ACM Multimedia Grand Challenge 2009
- find best punchlines in Seinfeld episodes

Speedup progress in Par Lab:

2006 0.3x realtime, original code

2008 1.5x realtime, optimized serial code

2010 14.3x realtime, multicore CPU+GPU

2011 250x realtime, pure GPU, from Python
code, changed the field!



Interactive GUI



Types of Programming (or “types of programmer”)

Example Languages

Example Activities

**Domain-Level
(No formal CS)**

Max/MSP, SQL,
CSS/Flash/Silverlight,
Matlab, Excel

Builds app with DSL
and/or by optimizing
app for

Productivity-Level

Python/Ruby/Perl
Haskell
Scala

programming
frameworks, writes
application
frameworks (or apps)

**Efficiency-Level
(MS in CS)**

Java/C#
C/C++/FORTRAN
assembler

Uses hardware/OS
primitives, builds
programming
frameworks (or apps)

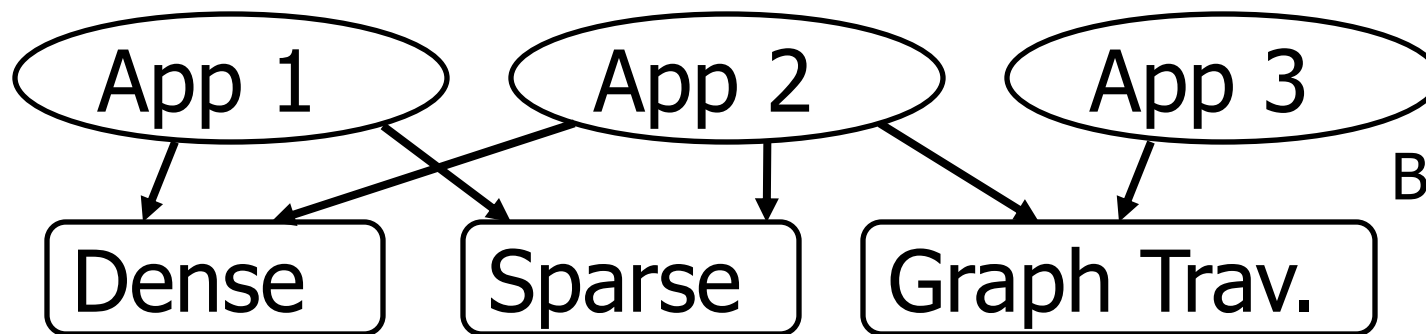
Hardware/OS

Provides hardware
primitives and OS services

Where & how to make parallelism visible?

How to make parallelism visible?

- ❖ In a new general-purpose parallel language?
 - An oxymoron?
 - Won't get adopted
 - Most big applications written in >1 language
- ❖ **Par Lab is betting on Computational and Structural Patterns at all levels of programming (Domain thru Efficiency)**
 - Patterns provide a good vocabulary for domain experts
 - Also comprehensible to efficiency-level experts or hardware architects
 - *Lingua franca* between the different levels in Par Lab

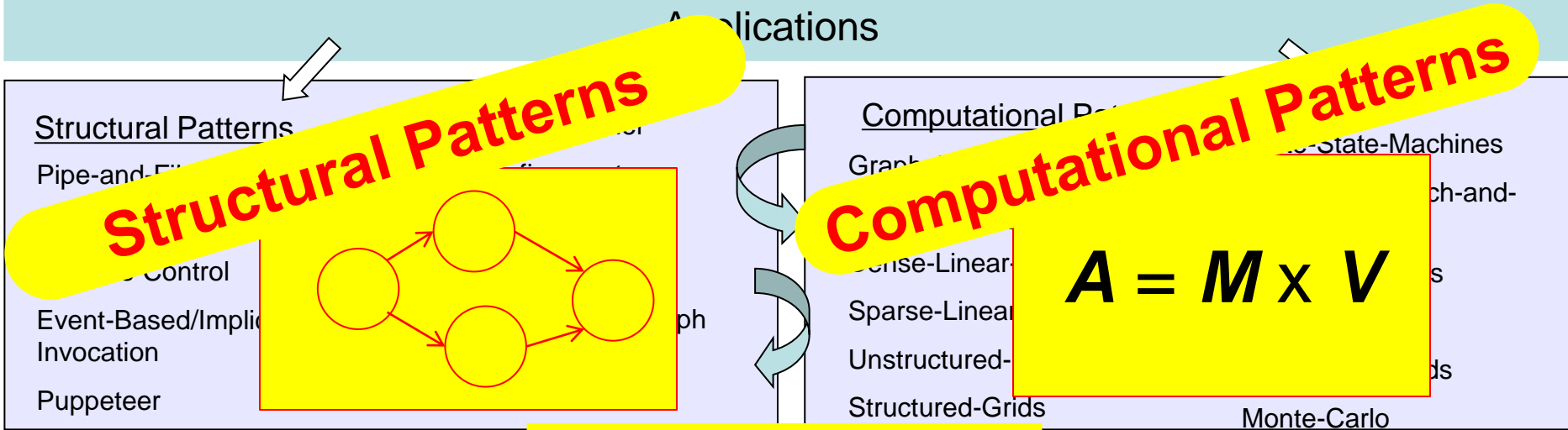


Berkeley View
Motifs
("Dwarfs")

How do compelling apps relate to 13 motifs?

	Embed	SPEC	DB	Games	ML	CAD	HPC	Health	Image	Speech	Music	Browser	
1 Finite State Mach.	Red	Red	Red	Yellow	Yellow	Yellow	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Red
2 Circuits	Red	Light Blue	Light Green	Light Blue	Light Green	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Red
3 Graph Algorithms	Red	Yellow	Light Blue	Yellow	Red	Red	Light Blue	Red	Light Blue	Red	Light Green	Light Green	Light Green
4 Structured Grid	Red	Red	Light Blue	Yellow	Light Blue	Light Blue	Red	Light Blue	Red	Light Blue	Light Blue	Light Blue	Light Blue
5 Dense Matrix	Red	Red	Yellow	Red	Red	Red	Light Blue	Light Blue	Red	Red	Red	Light Blue	Light Blue
6 Sparse Matrix	Yellow	Yellow	Light Blue	Red	Red	Red	Light Blue	Red	Light Blue	Light Blue	Red	Light Blue	Light Blue
7 Spectral (FFT)	Yellow	Light Blue	Light Blue	Yellow	Yellow	Yellow	Red	Light Blue	Light Green	Red	Red	Red	Red
8 Dynamic Prog	Yellow	Light Blue	Red	Light Blue	Red	Red	Light Blue	Light Blue	Light Blue	Yellow	Light Blue	Light Blue	Red
9 Particle Methods	Light Blue	Yellow	Light Blue	Yellow	Light Blue	Light Blue	Red	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue
10 Backtrack/ B&B	Light Blue	Light Blue	Yellow	Light Blue	Red	Red	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Yellow	Light Blue
11 Graphical Models	Light Blue	Light Blue	Yellow	Light Blue	Red	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Red	Light Blue	Light Blue
12 Unstructured Grid	Light Blue	Light Blue	Light Blue	Yellow	Yellow	Yellow	Red	Red	Light Blue	Light Blue	Light Blue	Red	Light Blue
13 Monte Carlo	Light Blue	Light Green	Light Blue	Red	Yellow	Yellow	Red	Red	Yellow	Light Blue	Light Blue	Yellow	Light Blue

"Our" Pattern Language (OPL-2010) (Kurt Keutzer, Tim Mattson)



Concurrent Algorithm Strategies

- Task-Parallelism
- Divide and Conquer

Refine Towards Implementation

Discrete-Event

- Geometric-Decomposition
- Speculation

Implementation Strategy Patterns

- SPMD
- Data-Par/index-space
- Fork/Join
- Actors
- Task-Graph

Program structure

Implementation Strategy Patterns

- Queue
- Shared-map
- Partitioned Graph
- Distributed-Array
- Shared-Data

Data structure

Parallel Execution Patterns

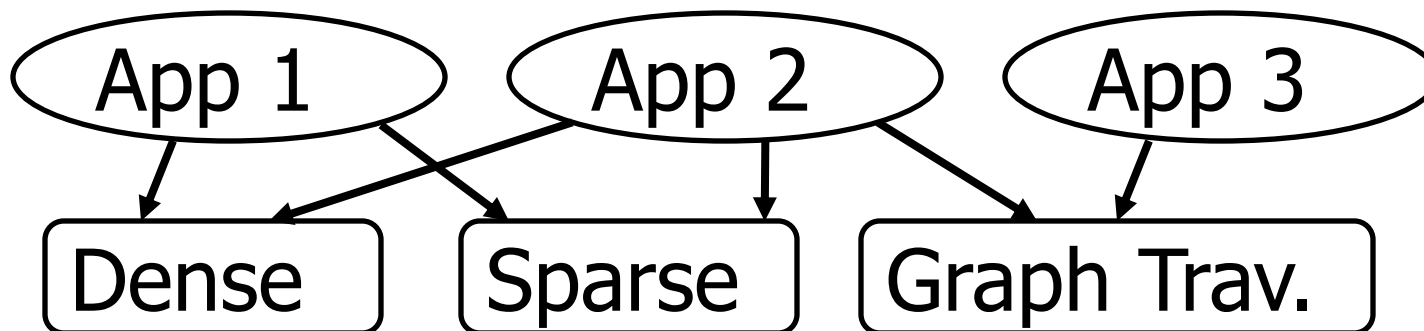
- MIMD
- SIMD
- Thread-Pool
- Task-Graph
- Transactions

Concurrency Foundation constructs (not expressed as patterns)

Thread creation/destruction
Process creation/destruction

Message-Passing
Collective-Comm.

Point-To-Point-Sync. (mutual exclusion)
collective sync. (barrier)



Only a few types of hardware platform

Multicore

GPU

“Cloud”

High-level pattern constrains space of reasonable low-level mappings

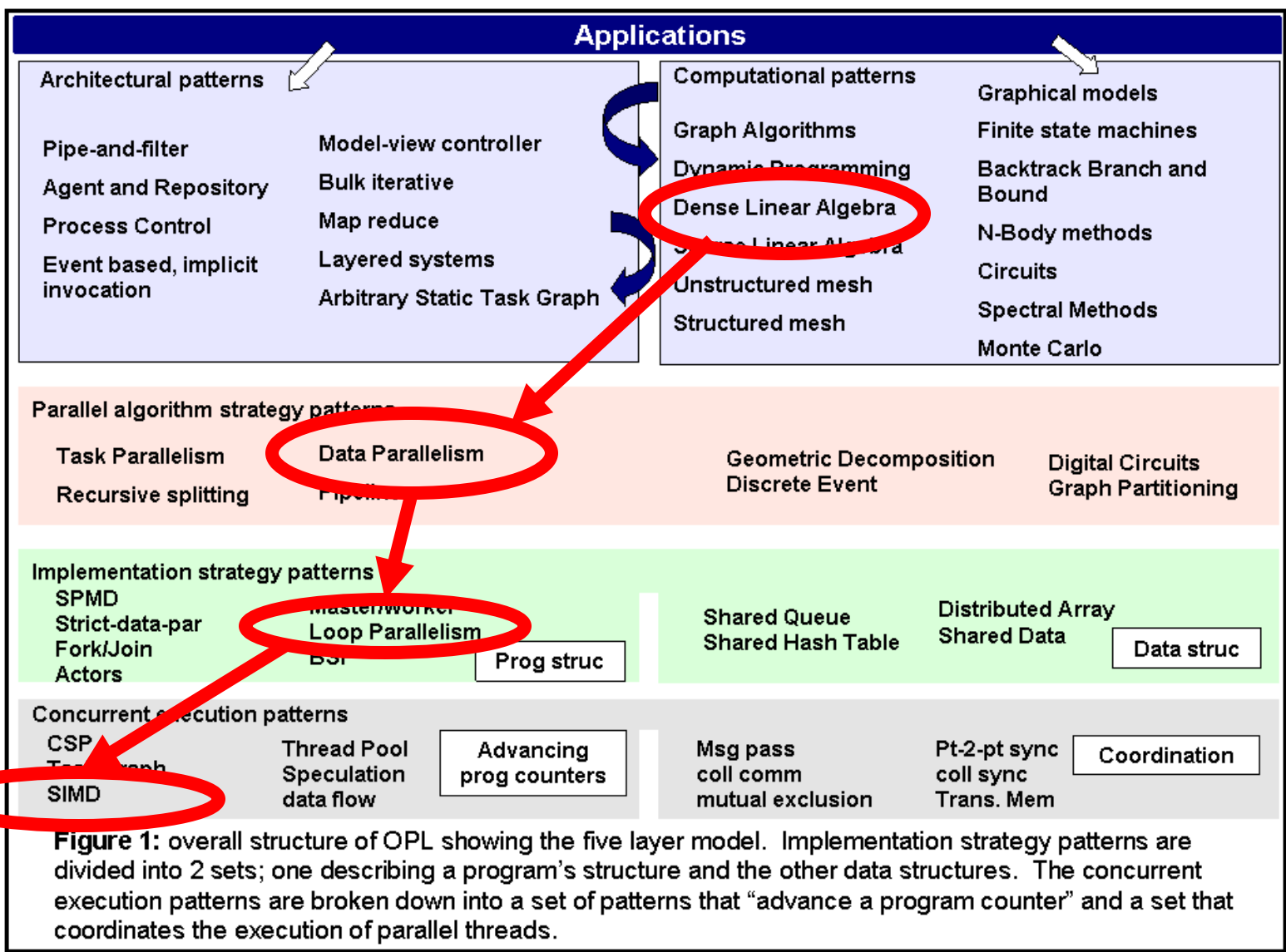
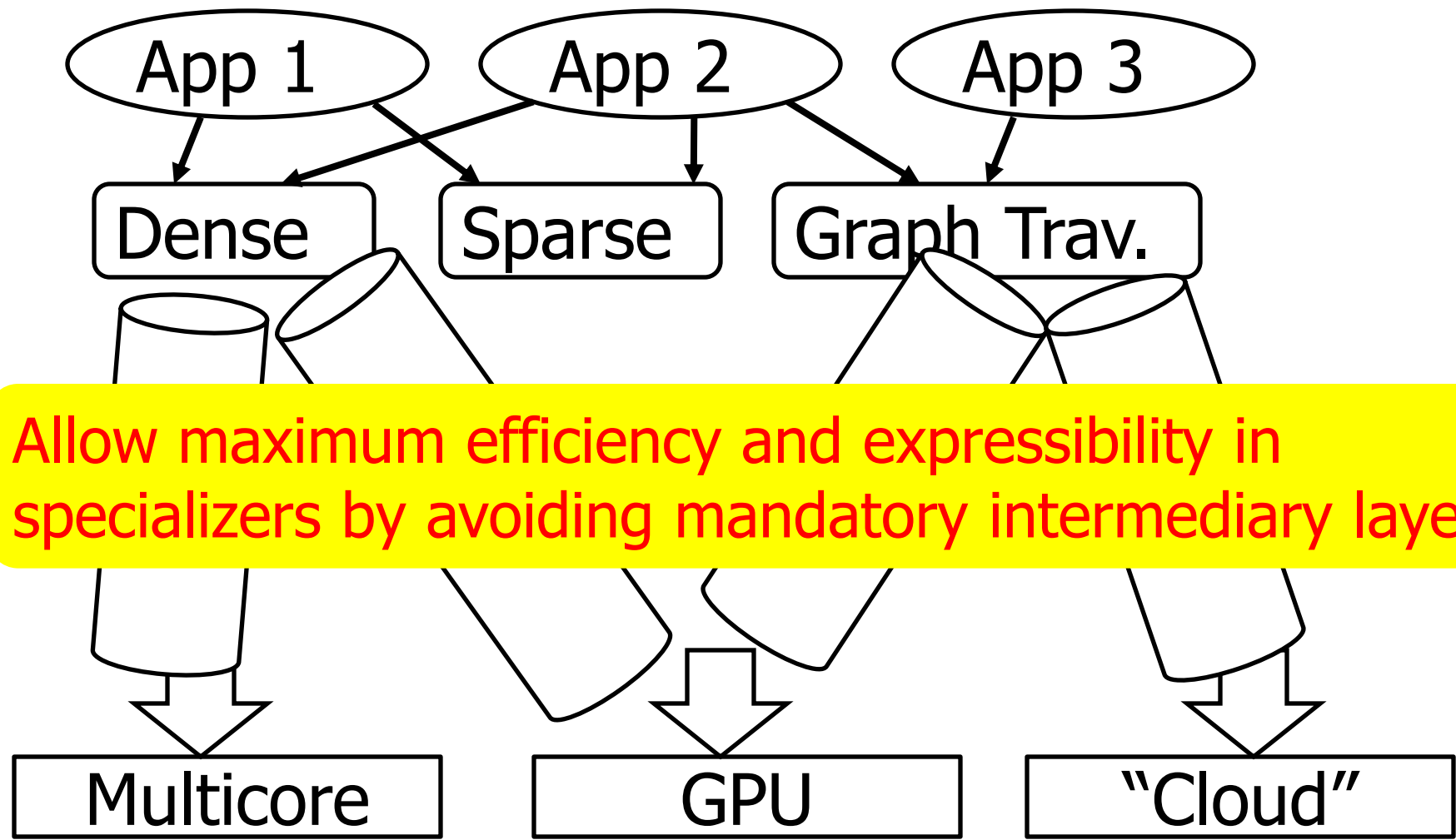


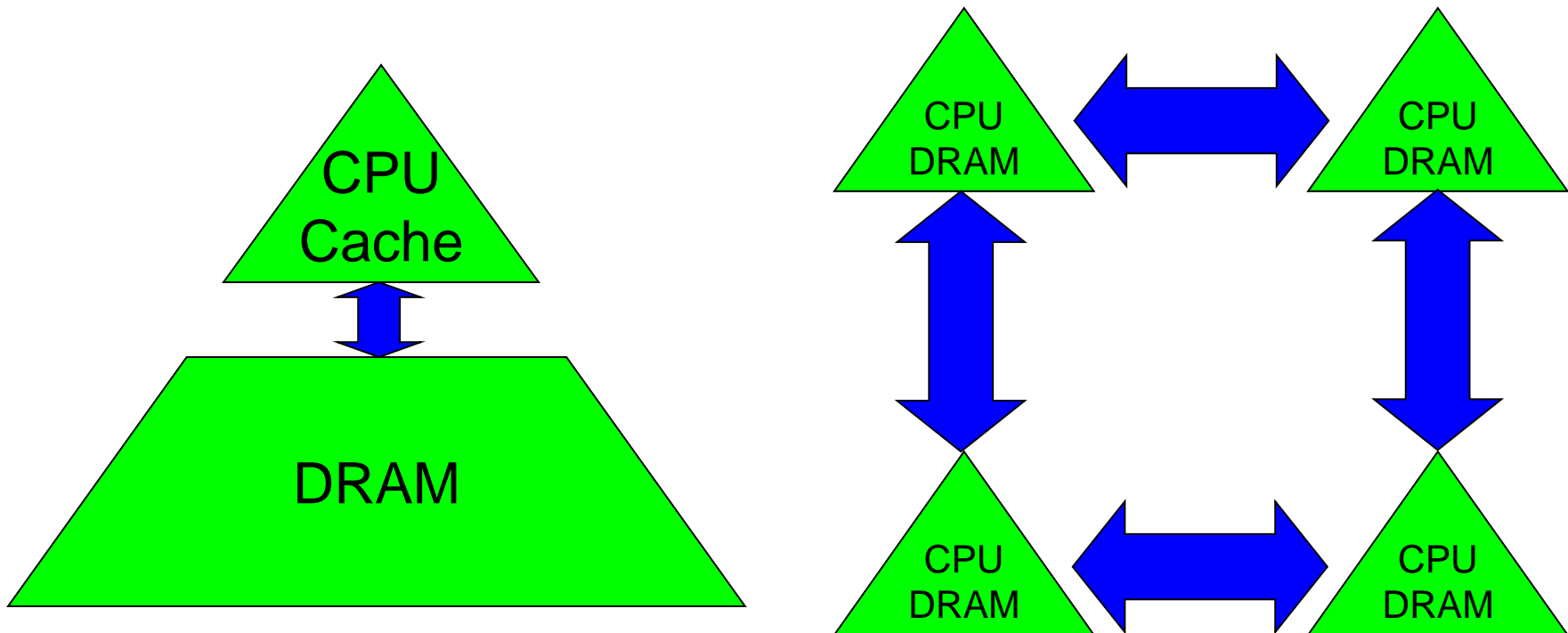
Figure 1: overall structure of OPL showing the five layer model. Implementation strategy patterns are divided into 2 sets; one describing a program's structure and the other data structures. The concurrent execution patterns are broken down into a set of patterns that "advance a program counter" and a set that coordinates the execution of parallel threads.

Specializers: Pattern-specific and platform-specific compilers

aka. "Stovepipes"



1. Arithmetic (FLOPS)
2. Communication: moving data between
 - levels of a memory hierarchy (sequential case)
 - processors over a network (parallel case).



- ❖ Cost of communication \gg cost of arithmetic
 - True for cost = time, or cost = energy per operation
 - Cost gap growing over time
- ❖ Goals
 - Identify lower bounds on *communication* required by widely used algorithms
 - Many widely used libraries (eg Sca/LAPACK) communicate asymptotically more than necessary
 - Design new algorithms that attain lower bounds
 - Possible for dense and sparse linear algebra, n-body, ...
 - Big speedups and energy savings possible

- ❖ Matrix multiplication
 - Up to **12x** on IBM BG/P for $n=8K$ on 64K cores; **95% less communication**
- ❖ QR decomposition (used in least squares, data mining, ...)
 - Up to **8x** on 8-core dual-socket Intel Clovertown, for $10M \times 10$
 - Up to **6.7x** on 16-proc. Pentium III cluster, for $100K \times 200$
 - Up to **13x** on Tesla C2050 / Fermi, for $110k \times 100$
 - Up to **4x** on Grid of 4 cities (Dongarra, Langou et al)
 - “infinite speedup” for out-of-core on PowerPC laptop
 - LAPACK thrashed virtual memory, didn’t finish
- ❖ Eigenvalues of band symmetric matrices
 - Up to **17x** on Intel Gainestown, 8 core, vs MKL 10.0 (up to **1.9x** sequential)
- ❖ Iterative sparse linear equations solvers (GMRES)
 - Up to **4.3x** on Intel Clovertown, 8 core
- ❖ N-body (direct particle interactions with cutoff distance)
 - Up to **10x** on Cray XT-4 (Hopper), 24K particles on 6K procs.

- ❖ SIAM Linear Algebra Prize 2012, for best paper in previous 3 years, deriving lower bounds
- ❖ SPAA'11 Best Paper Award, for Strassen lower bounds
- ❖ EuroPar'11 Distinguished Paper Award, for asymptotically faster “2.5D” matmul and LU
- ❖ Citation in 2012 DOE Budget Request ...

President Obama cites Communication-Avoiding Algorithms in the FY 2012 Department of Energy Budget Request to Congress:

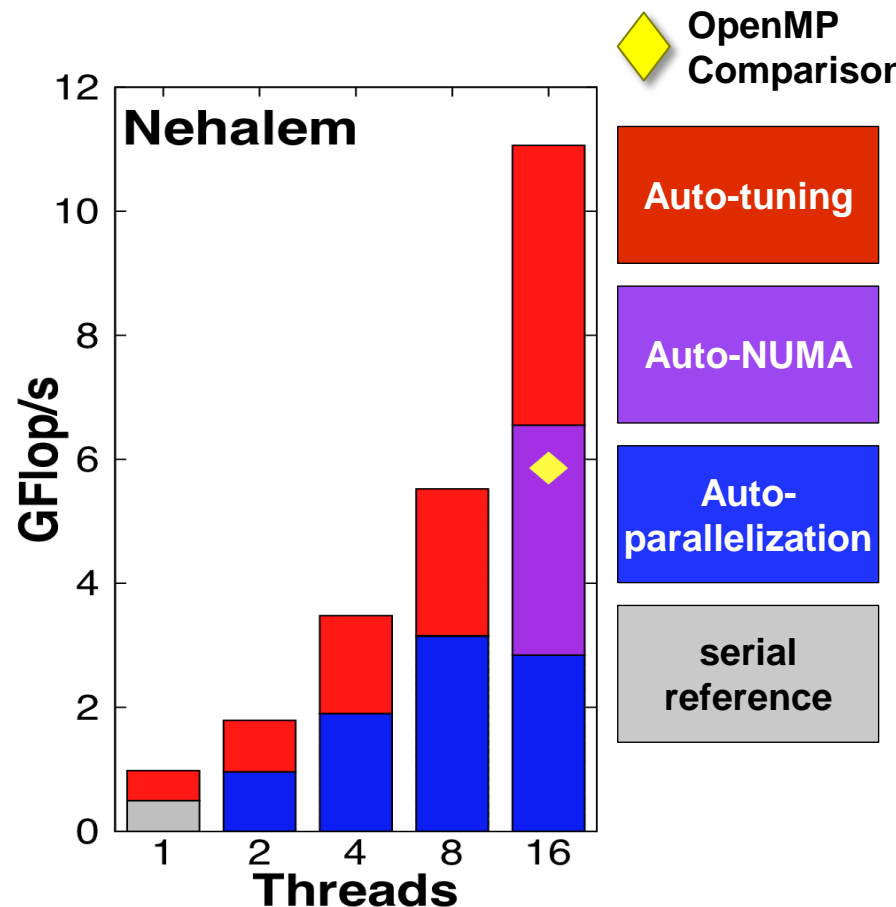
“New Algorithm Improves Performance and Accuracy on Extreme-Scale Computing Systems. **On modern computer architectures, communication between processors takes longer than the performance of a floating point arithmetic operation by a given processor.** ASCR researchers have developed a new method, derived from commonly used linear algebra methods, to **minimize communications between processors and the memory hierarchy, by reformulating the communication patterns specified within the algorithm.** This method has been implemented in the TRILINOS framework, a highly-regarded suite of software, which provides functionality for researchers around the world to solve large scale, complex multi-physics problems.”

FY 2010 Congressional Budget, Volume 4, FY2010 Accomplishments, Advanced Scientific Computing Research (ASCR), pages 65-67.

CA-GMRES (Hoemmen, Mohiyuddin, Yelick, Demmel)
“Tall-Skinny” QR (Grigori, Hoemmen, Langou, Demmel)

- ❖ New algorithm for Breadth-First Search
- ❖ Highest single-node performance in November 2011, Graph500, using Intel Xeon E7-8870 (Mirasol)
- ❖ #15: BlueGene 2048 cores 6.93 GTEPS
- ❖ #16: Jaguar 1024 cores 6.26 GTEPS
- ❖ **#17: Mirasol 40 cores 5.12 GTEPS**
- ❖ #18: Blacklight 512 cores 4.45 GTEPS
- ❖ #19: Todi 176 TESLA GPUs 3.05 GTEPS
- ❖ #20: Convey 4 FPGAs 1.76 GTEPS

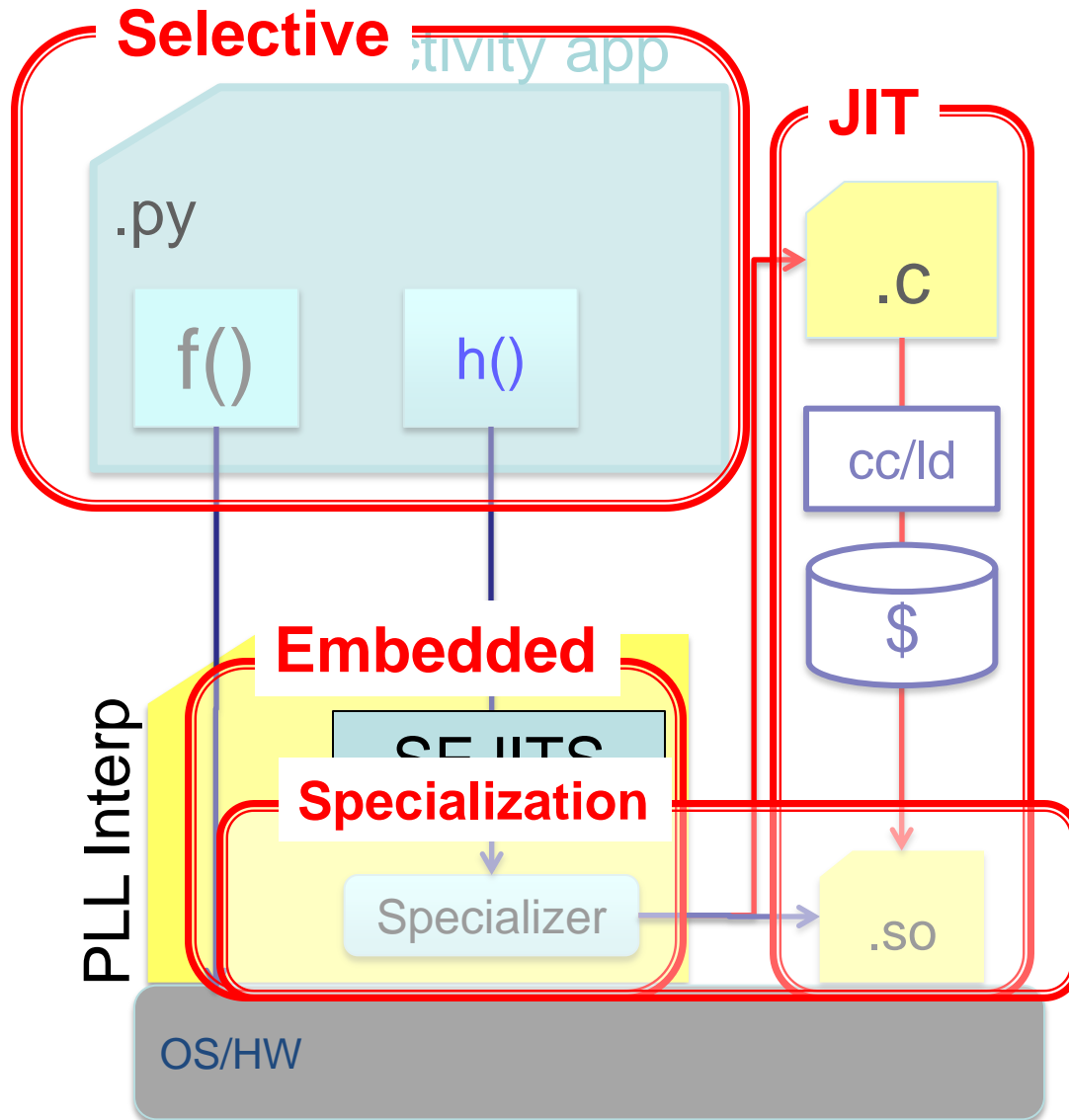
- Problem: generating optimized code is like searching for needle in haystack; use computers rather than humans
- **Auto-tuners** approach: program *generates* optimized code and data structures for a “motif” (~kernel) mapped to some instance of a family of architectures (e.g., x86 multicore)
- Use empirical measurement to select best performing.
- ParLab autotuners for stencils (e.g., images), sparse matrices, particle/mesh, collectives (e.g., “reduce”), ...



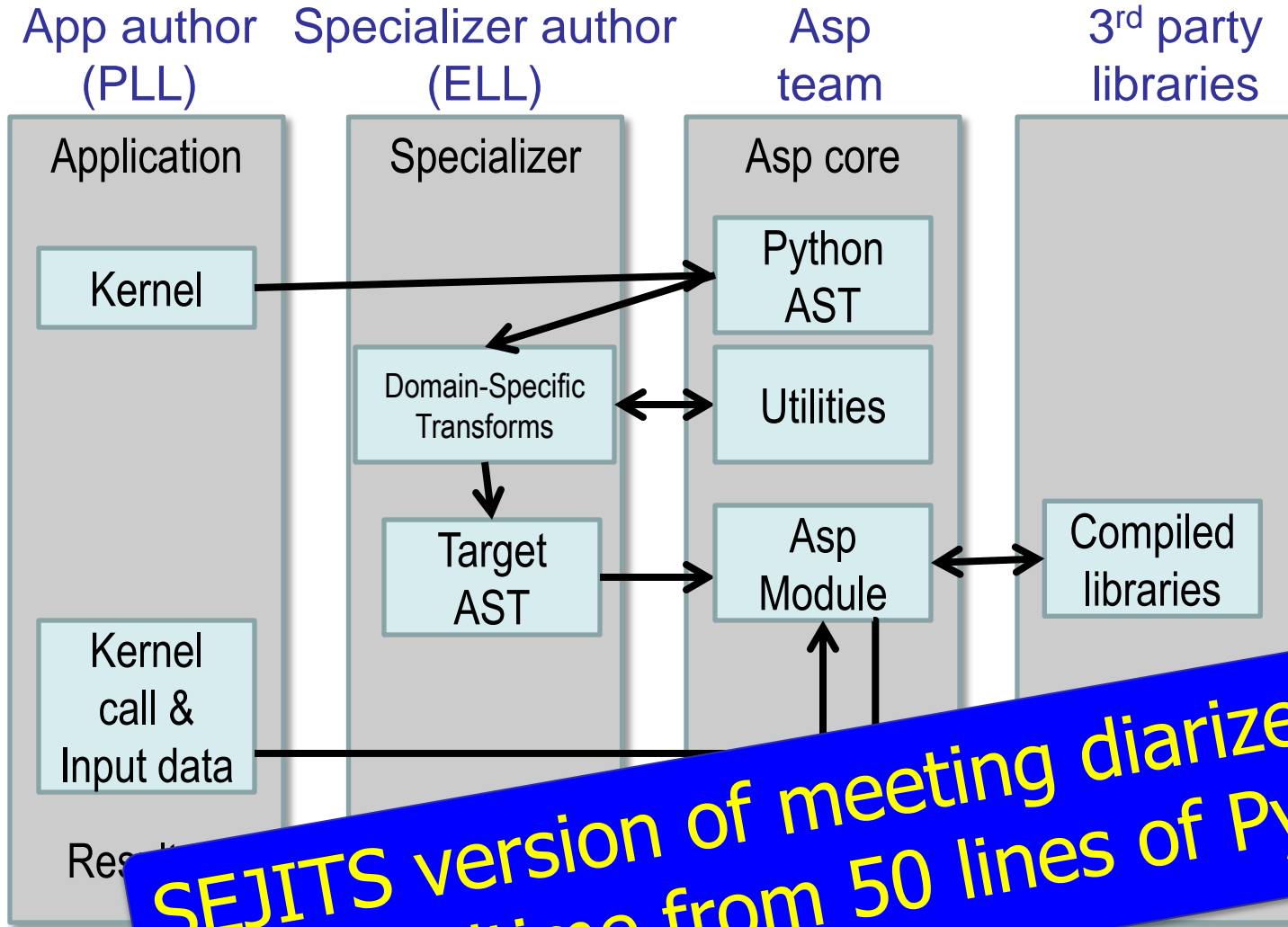
- ❖ SEJITS bridges productivity and efficiency layers through specializers embedded in modern high-level productivity language (Python, Ruby, ...)
 - Embedded “specializers” use language facilities to map high-level pattern to efficient low-level code (at run time, install time, or development time)
 - Specializers can incorporate/package autotuners

Two ParLab SEJITS projects:

- ❖ **Copperhead**: Data-parallel subset of Python, development continuing at NVIDIA
- ❖ **Asp**: “Asp is SEJITS in Python” general specializer framework
 - Provide functionality common across different specializers

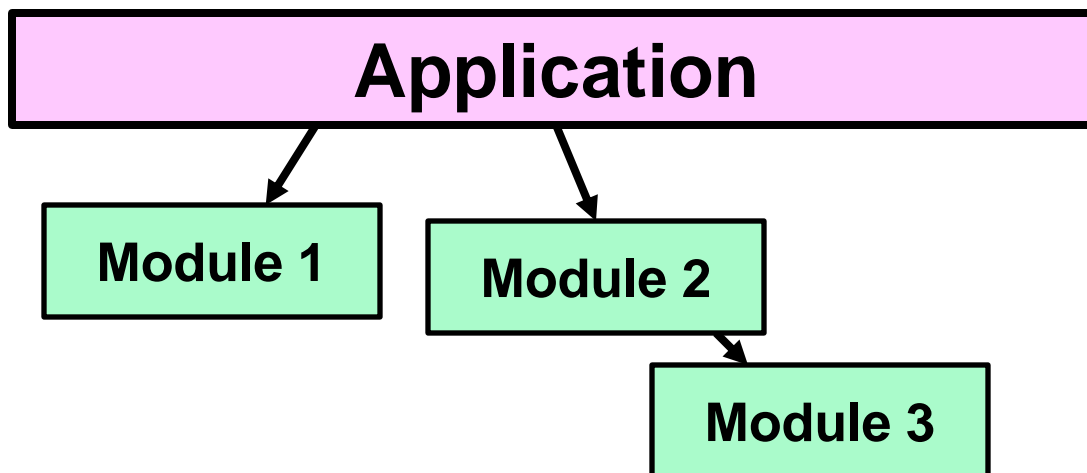


Asp: Who Does What?



**SEJITS version of meeting diarizer,
250x realtime from 50 lines of Python**

- ❖ All applications built as a hierarchy of modules, not just one kernel

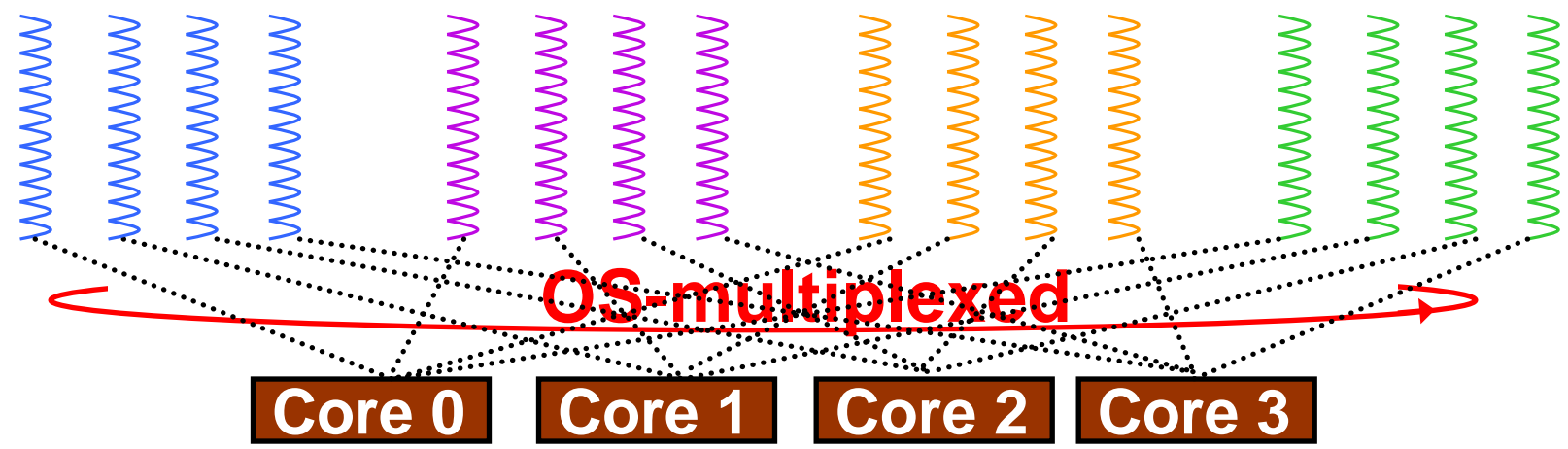
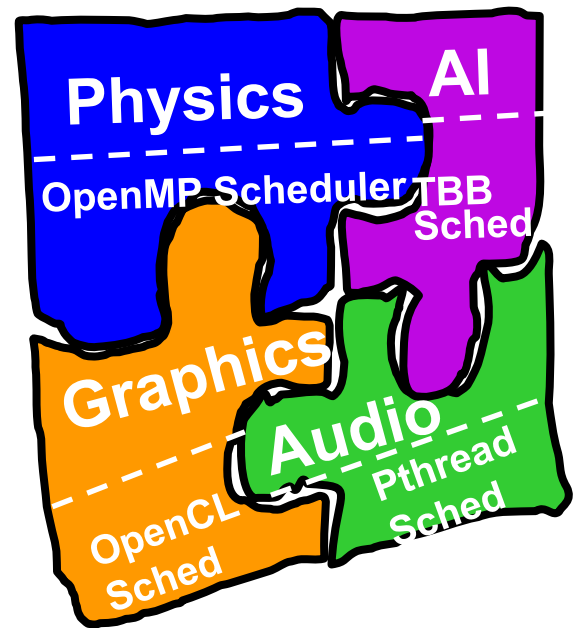


Structural patterns describe the common forms of composing sub-computations:

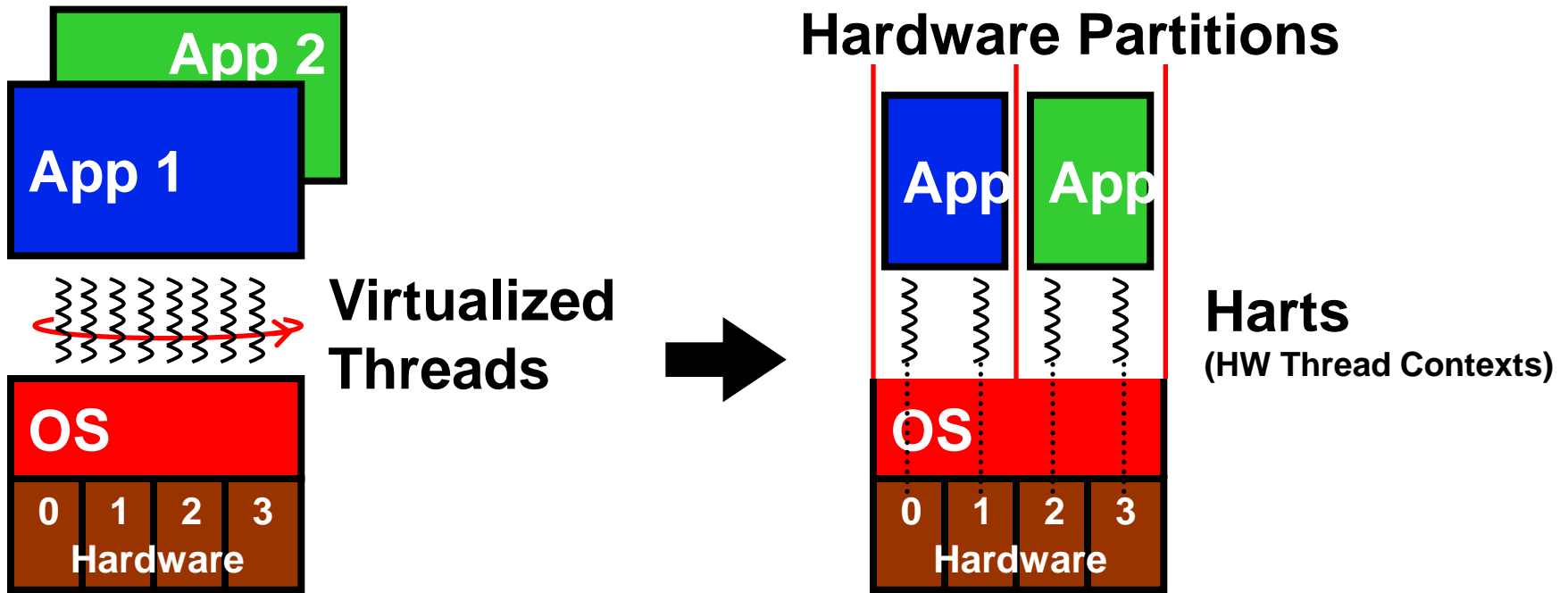
E.g., task graph, pipelines, agent&repository

- ❖ **Data format/layout:** Must translate between data formats or layouts expected by different components
- ❖ **Synchronization:** Must correctly synchronize data passing between or shared by multiple components
- ❖ **Resource management:** Must share hardware resources to execute components in parallel

**Gaming
App
Example**



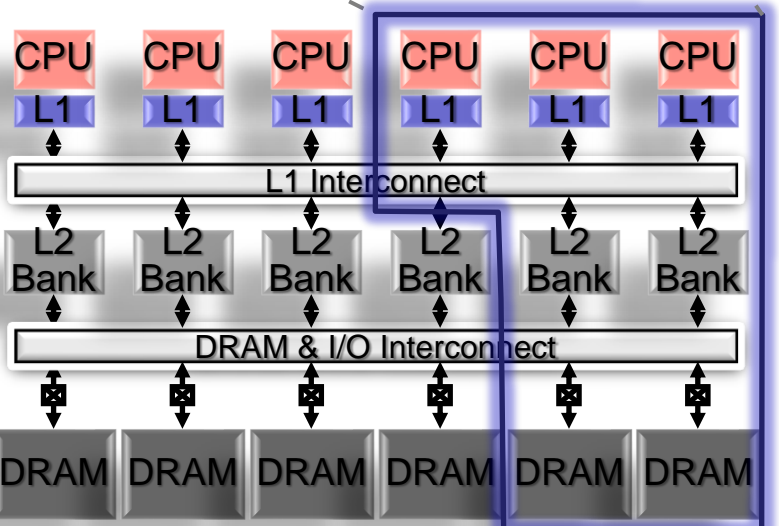
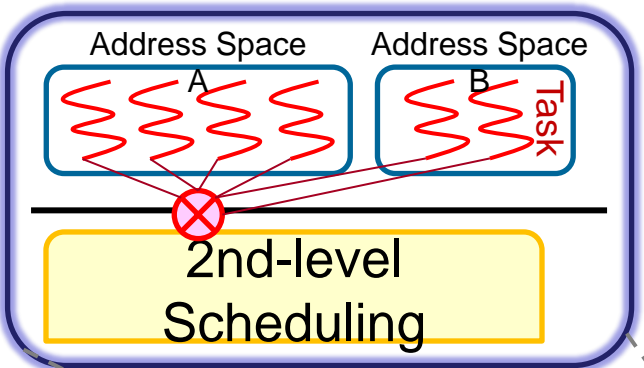
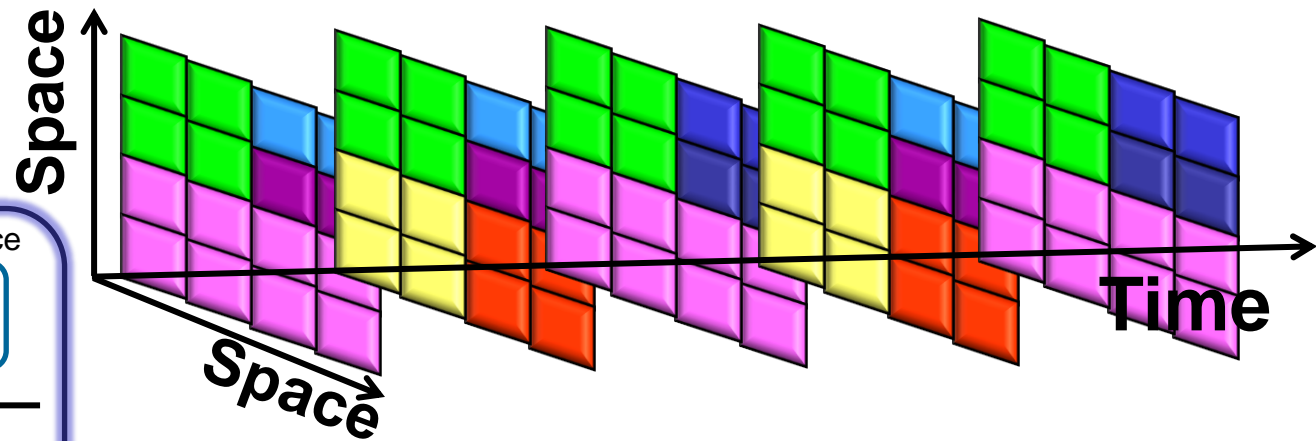
Libraries compete unproductively for resources!



- *Merged* resource and computation abstraction.

- More accurate resource abstraction.
- Let apps provide own computation abstractions

- ❖ Lithe is an ABI to allow application components to co-operatively share hardware threads.
- ❖ Each component is free to map computational to hardware threads in any way they see fit
 - No mandatory thread or task abstractions
- ❖ Components request but cannot demand harts, and must yield harts when blocked or finished with task



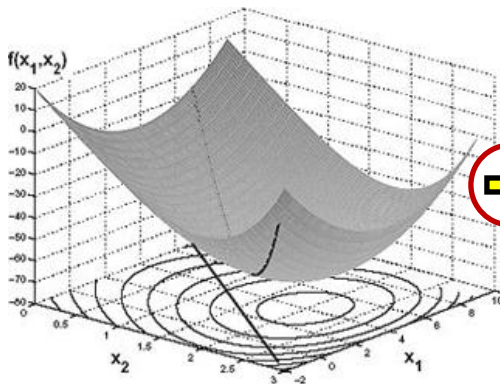
1st level: OS determines coarse-grain allocation of resources to jobs over space and time

2nd level: Application schedules component tasks onto available "harts" (hardware thread contexts) using Lithe

Resource Management using Convex Optimization (Sarah Bird, Burton Smith)

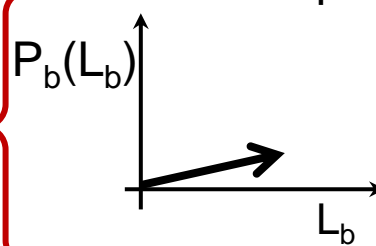
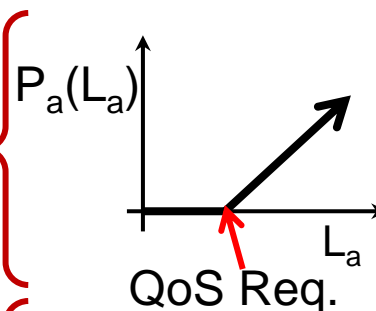
- ❖ Each process receives a **vector of basic resources** dedicated to it
 - e.g., fractions of cores, cache slices, memory pages, bandwidth
- ❖ Allocate minimum for QoS requirements
- ❖ Allocate remaining to meet some system-level objective
 - e.g., best performance, lowest energy, best user experience

Continuously
Minimize
(subject to restrictions
on the total amount of
resources)



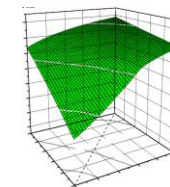
Convex Surface

Penalty Function
Reflects the app's
importance

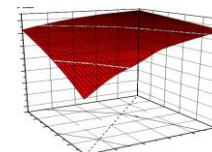


⋮
Performance Metric (L), e.g., latency

Resource Utility Function
Performance as function of
resources

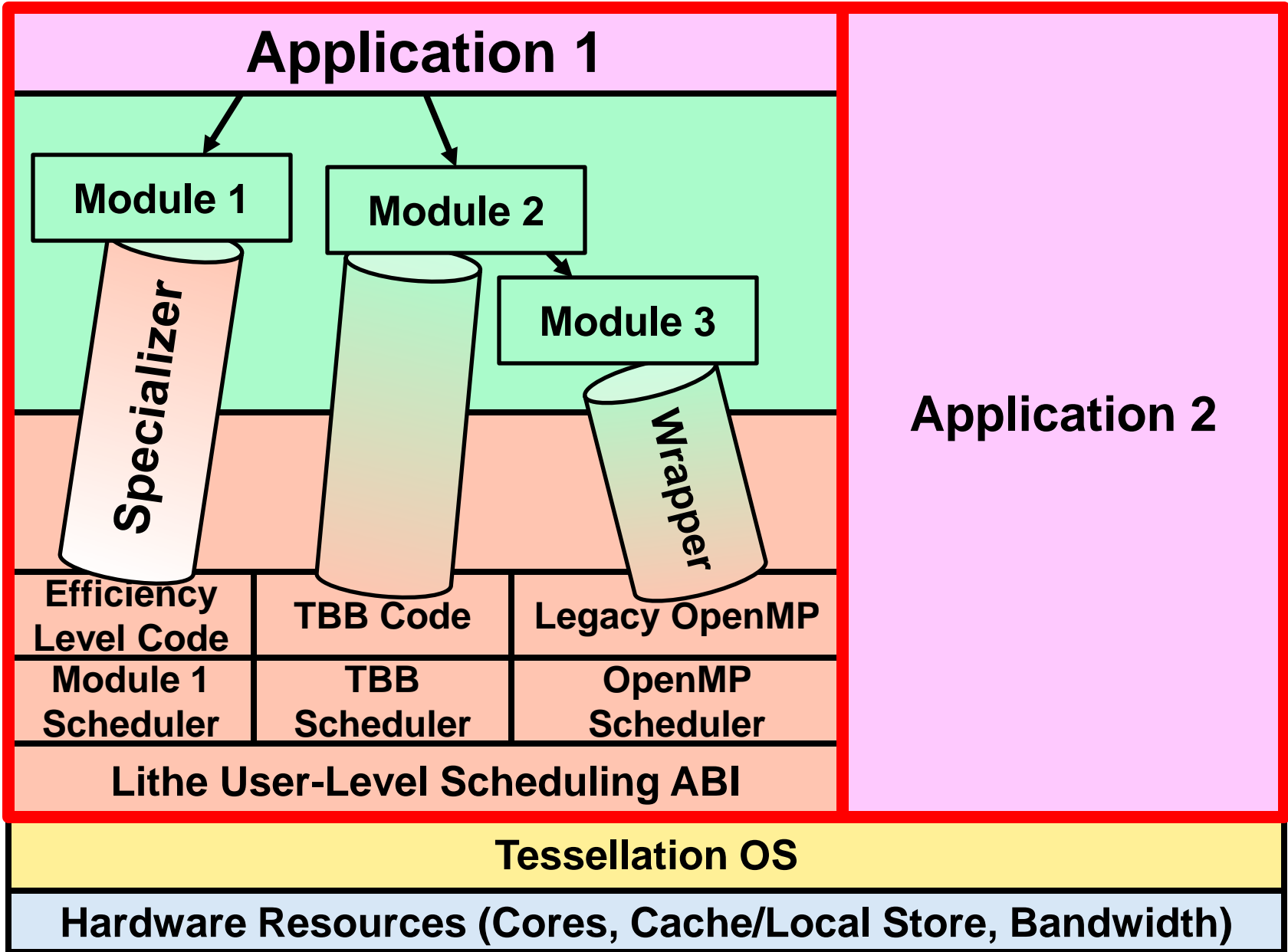


$$L_a = RU_a(r_{(0,a)}, r_{(1,a)}, \dots, r_{(n-1,a)})$$



$$L_b = RU_b(r_{(0,b)}, r_{(1,b)}, \dots, r_{(n-1,b)})$$

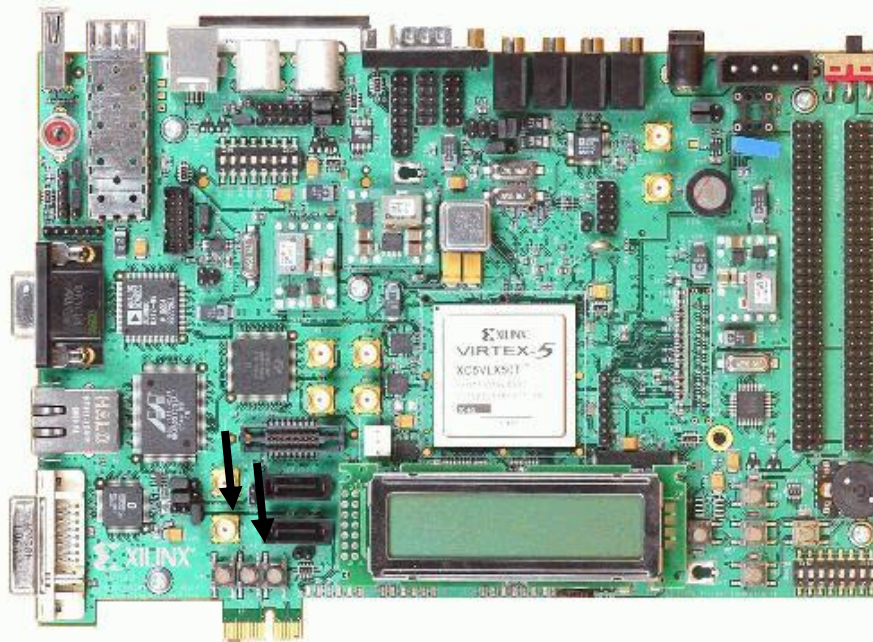
- ❖ Organize software around parallel patterns
 - Maximize reuse since patterns common across domains
- ❖ Each pattern implemented with efficient algorithms packaged as SEJITS specializers using autotuners
- ❖ Programmer composes functionality at high-level using productivity language
- ❖ System composes resource usage at low-level using 2-level scheduling
 - Tessellation OS at coarse-grain
 - Lithe user-level scheduler ABI at fine-grain



- ❖ What about GPUs versus CPUs?
- ❖ These architectures are closely related, and converging.
- ❖ Both have multiple multithreaded cores each with many SIMD lanes
 - original vision was “manycore” – more accurate to say future is “manylane”.
- ❖ Most of our techniques can be applied to both

- ❖ Focus on supporting application and OS needs:
 - Hardware partitioning support
 - Performance counters
 - High-performance FPGA-based simulators
- ❖ New architecture ideas:
 - New data-parallel execution engines
 - Hardware+software managed memory hierarchy
 - Specialized accelerators (e.g., graph machines)
- ❖ Extensive development of VLSI flow to allow real layout of various data-parallel accelerators
 - Accurate cycle time, area, energy

- ❖ Rapid accurate simulation of manycore architectural ideas using FPGAs
- ❖ Initial version models 64 cores of SPARC v8 with shared memory system on \$750 board
- ❖ Hardware FPU, MMU, boots our OS and Par Lab stack!



	Cost	Performance (MIPS)	Time per 64 core simulation
Software Simulator	\$2,000	0.1 - 1	250 hours
RAMP Gold	\$2,000 + \$750	50 - 100	1 hour

- ❖ Research supported by Microsoft (Award #024263) and Intel (Award #024894) funding and by matching funding by U.C. Discovery (Award #DIG07-10227).
- ❖ Additional support comes from Par Lab affiliates National Instruments, NEC, Nokia, NVIDIA, Samsung, and Oracle/Sun.

Questions?