

Structural Simulation Toolkit

Arun Rodrigues

Many
Cores
+
Memory

X

Many
Many
Nodes

X

Many
Many
Many
Threads

Audiences.....

or

X

Application writers
purchasers
designers

X

system procurement
algorithm co-design
architecture research
language research

X

present
future

Complexity.....

Physics Apps
Dynamics Apps

X

Communication Libraries
Run-Times
OS Effects

X

Existing Language
New Language

Constraints.....

Performance

Power

Cooling

Risk

Work/Application Feedback: A static trace or simple statistics will not capture the causal relationships between messages.

Scalability: Many network effects only become apparent at hundreds of thousands of nodes.

Processor/Memory/Network Systems: Local interactions can have global performance implications.

Ability to Model Message Overheads: Overheads in the network (e.g. packetization, protocol overhead) and messaging library (e.g. marshaling, message assembly) can have a major effect on performance.

Ability to Explore Programming Models: Novel hardware will require new programming techniques and capabilities.

Power and Economic Effects: Power and cost are the key limitations on system design. Any system model must be able to provide feedback on the power and cost implications of new architectural

Several Projects

Architecture

CacheLine Gather (ICGL)

on. FU (Wisc./SNL)

Aggregates

Memory Ops

Execution analysis

Memory Footprint

Instruction Usage

Work/MPI

Tradeoffs

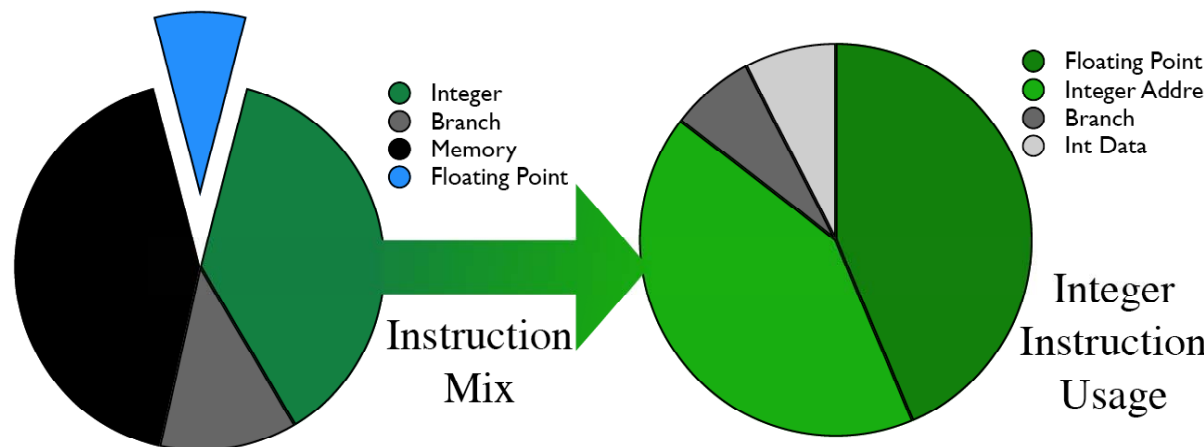
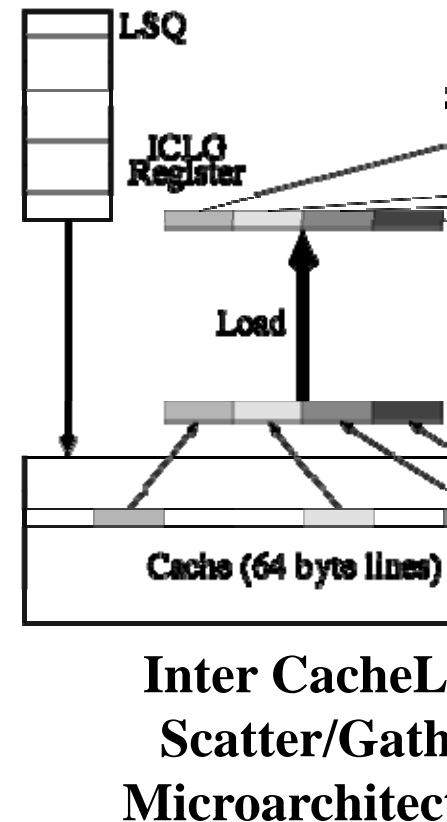
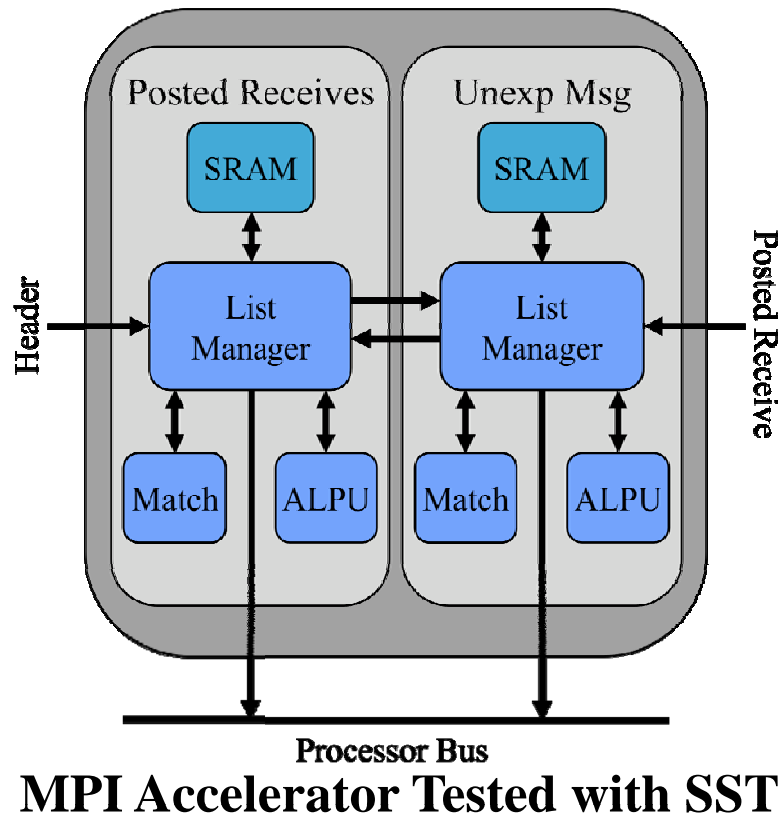
Acceleration

Programming Models

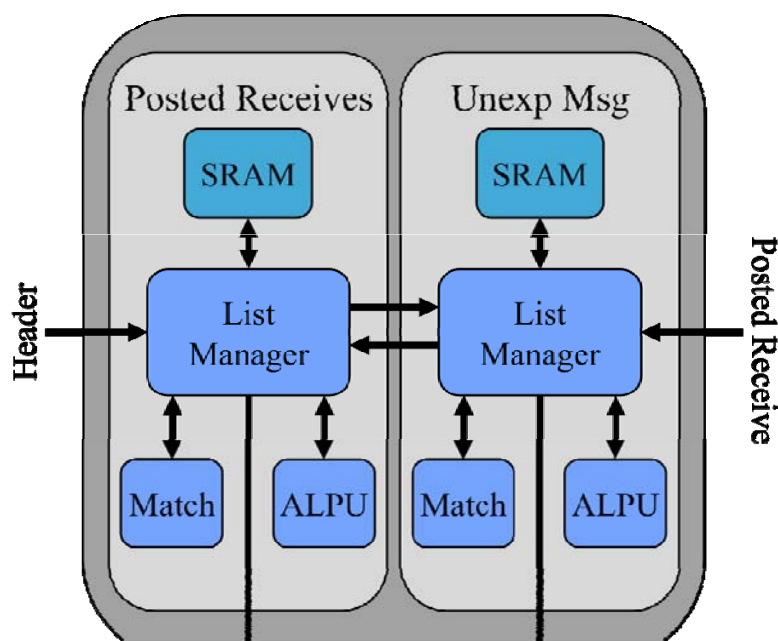
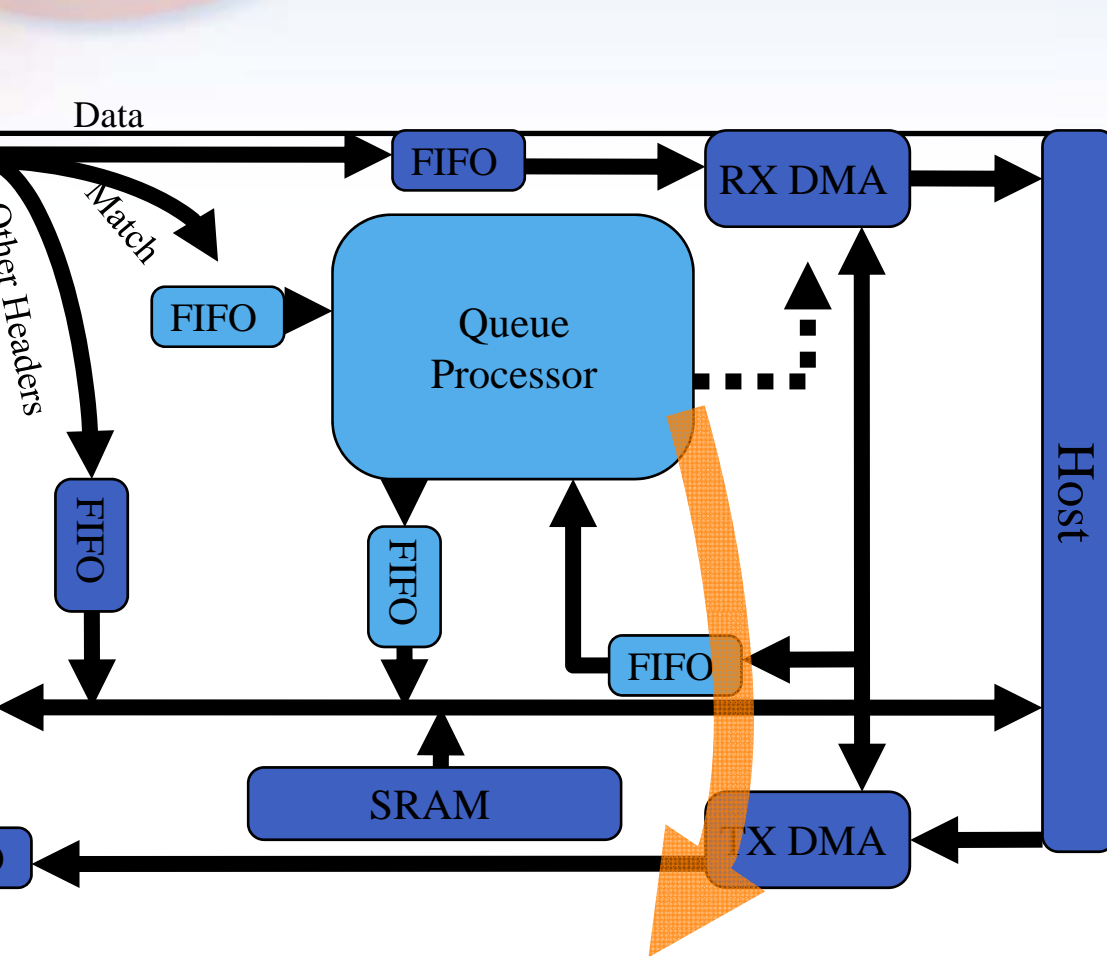
Compiler work (SNL/Rice)

UNIX (LSU/SNL)(FastOS)

Transactional Memory (ORNL)



Instruction Mix & Usage for



- Problem: Message rate determines effective bandwidth
- SST Analysis: MPI matching limits performance
- Solution: Accelerate list management with hardware
- SST Simulation
 - Hardware
 - Implement NIC/Router based RedStorm SeaStar
 - Implement List Manager, ALPU, Match Unit, integrate with MPI
 - “Translated” from FPGA
 - Validate against FPGA & Router
 - Software
 - Create baseline offload MPI
 - Modify MPI to use accelerated
- Impact: Collaboration w/ Intel

Goals

- Become the standard architectural simulator for the HPC community
- Be able to evaluate future systems on DOE workloads
- Use supercomputers to design supercomputers

Technical Approach

Scale
Accurate to analytic
Transition-based to message-based

of simulated nodes on 100s of
nodes

ted Tech. Models

Consortium

- “Best of Breed” simulation suite
- Combine Lab, academic, & industrial




```

= getNextEvent(queues)) {
    nt->time;
    sClockEvent) {
        mponent->preTic();
        ent->exchange) {
            nds();
        ends();
        ent->checkpoint) {
            oint();
        mponent->handleEvent();
    }
}

```

Parallel Core Pseudocode

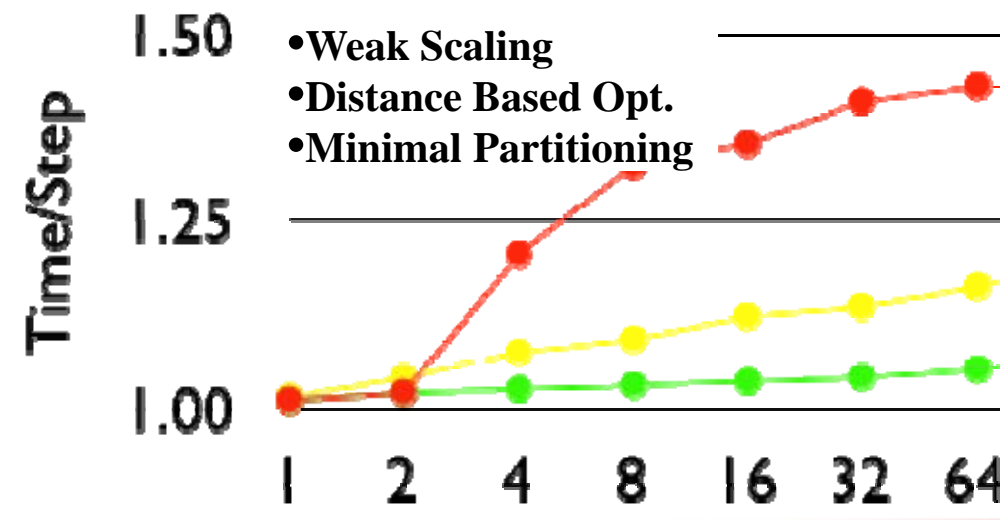
Message Traces, Symbolic Workload Descriptions	Execution Based	Execution w/ FPGA Acceleration
100s-1000s	100-1000s	1-10
10000s-100000s	100s-1000s	1-10
System Scaling Behavior	Cycle-level system performance	Co-design

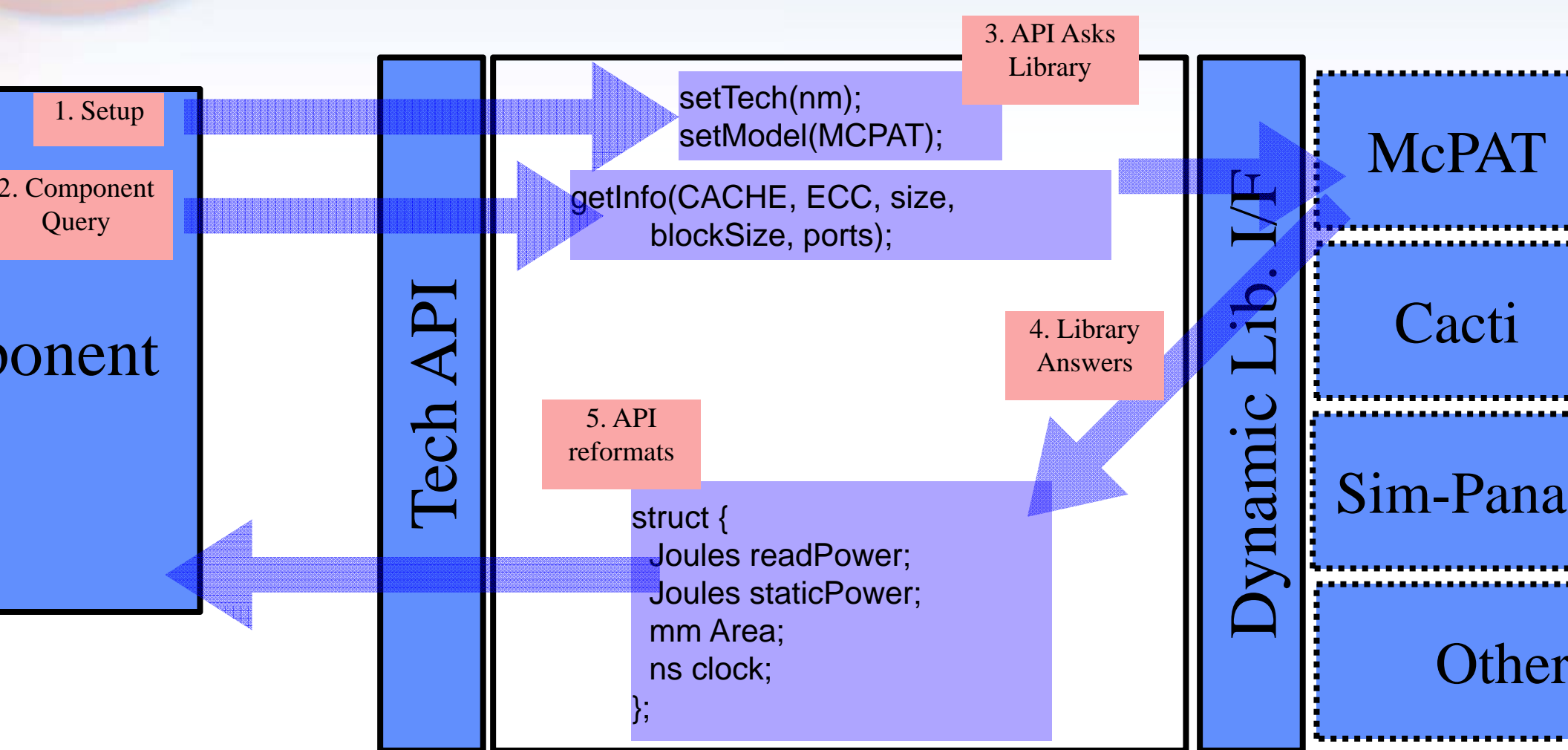
• Requirements

- High speed, Parallel, Scalable
- Multiple clock domains
- Checkpointing

• Implementation

- Conservative distance-based D optimization
- Multi-criteria partitioning
- Built on MPI
- Future: FPGA acceleration





ate interface to multiple technology libraries

Power/Energy

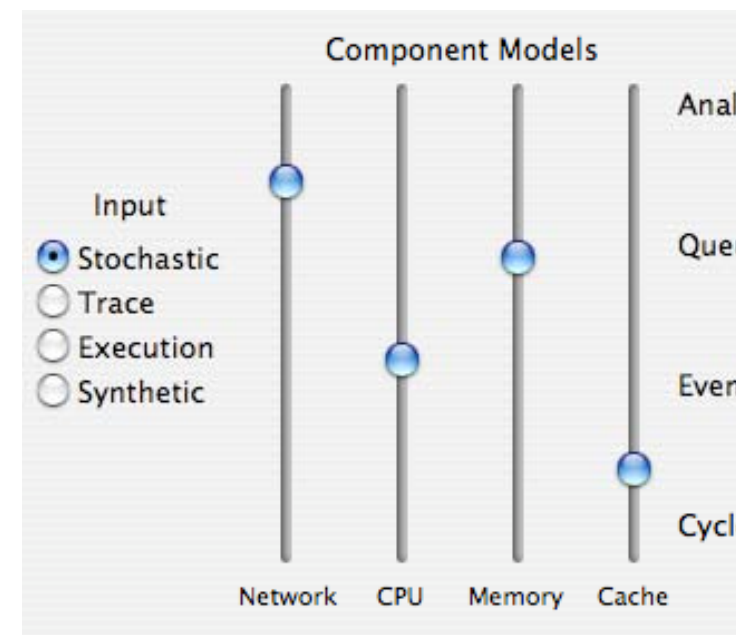
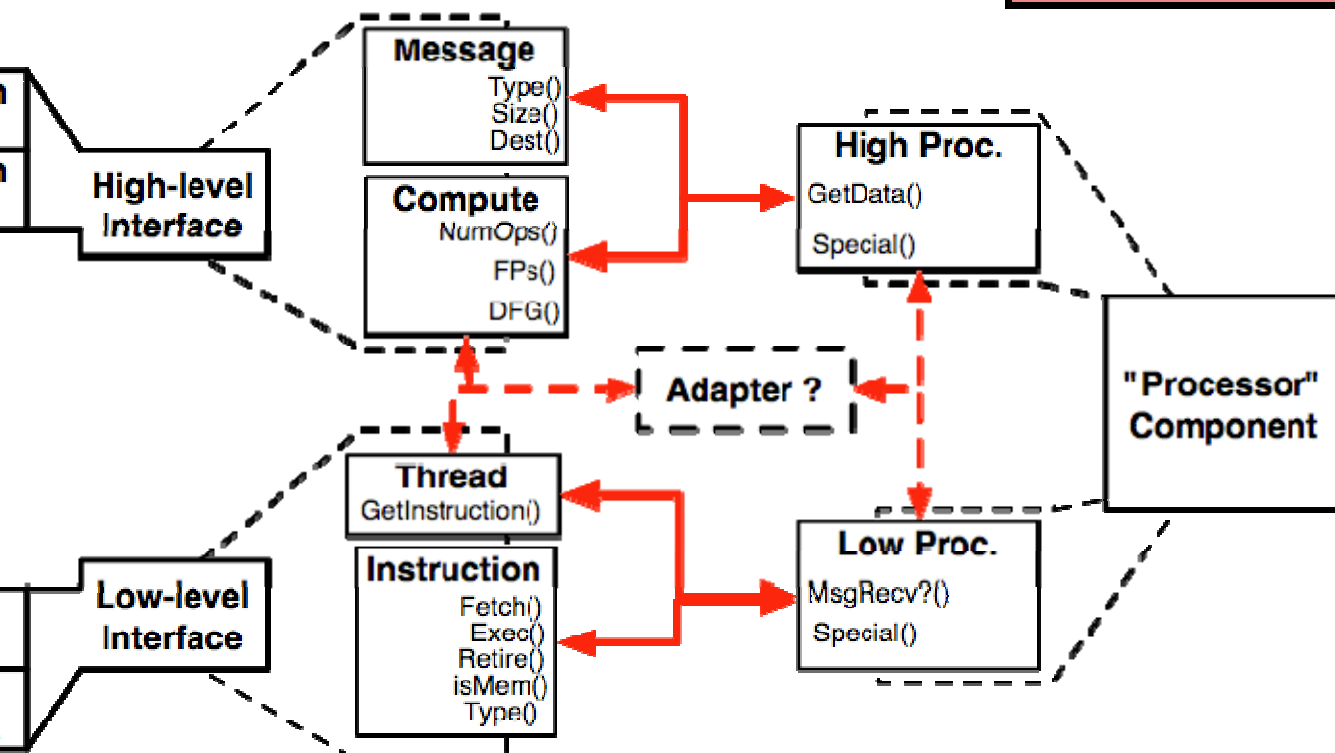
Reliability

Area/Timing estimation

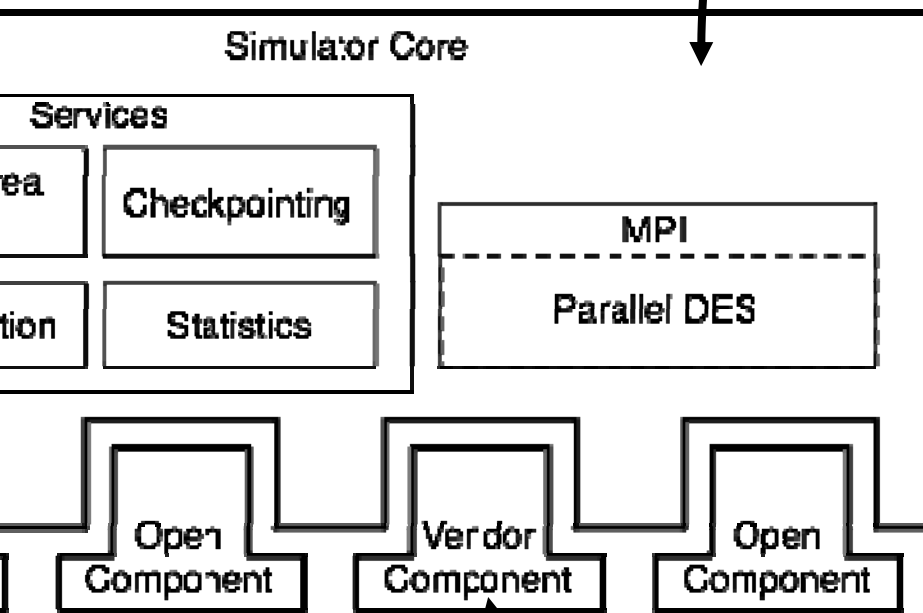
Make it easier for components to model technology

interface component reuse at
 different scales
 & Low-level interfaces (more?)
 shows multiple input types
 shows multiple input sources
 traces, stochastic, state-
 machines, execution...
 adapter objects to translate?

	High-Level	Low
Detail	Message	Instr
Fundamental Objects	Message, Compute block, Process	Instr Th
Static Generation	MPI Traces, MA Traces	Instr Tr
Dynamic Generation	State Machine	Exe



Multiscale Parameters



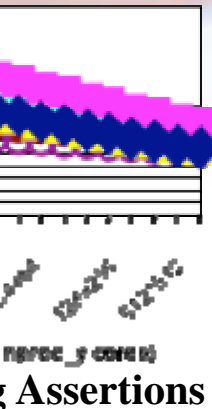
• Simulator Core will provide

- Power, Area, Cost modeling
- Checkpointing
- Configuration
- Statistics gathering
- Parallel Component-Based Discrete Event Simulation
 - MPI hidden from user
 - Multiple clocks

• Components

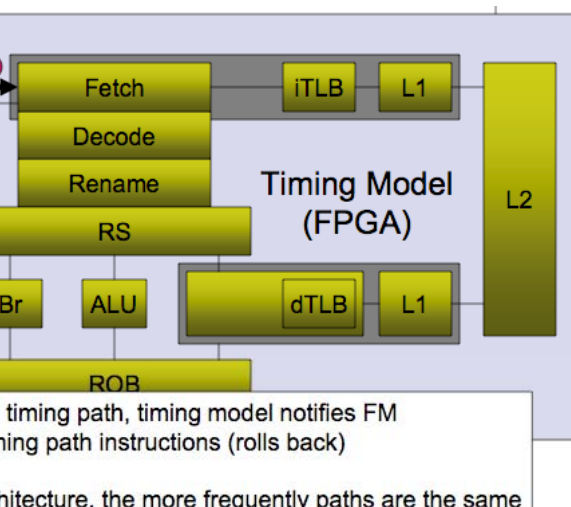
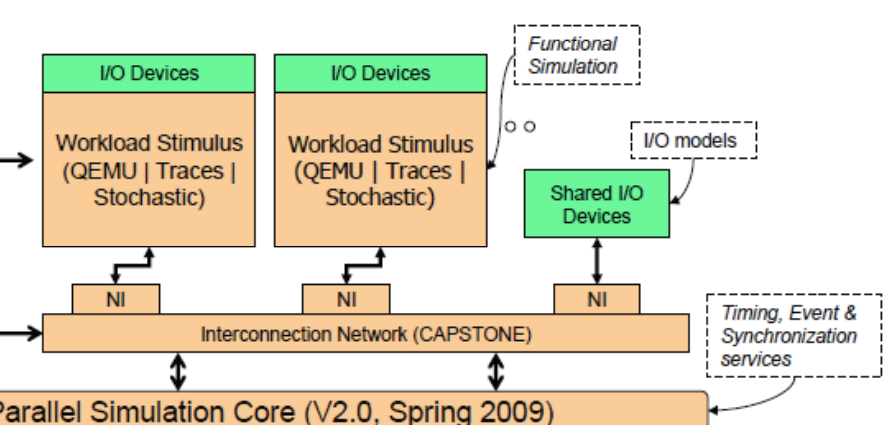
- Ships with basic set of open components
- Industry can plug in their own models

Performance Release			
Component	September	Dec 09	Future
CORE	DES, Checkpointing, Configuration	Improved Power Models, Optimized DES	Manifold Compatibility
Processor	SimpleScalar-based	Stochastic NMSU Models	Detailed
Network	Off-load NIC, Simple Router, Macroscale Models	Improved NIC, multiple topologies	Improved Router
Memory	DRAMSim II Integration	Improved	Improved

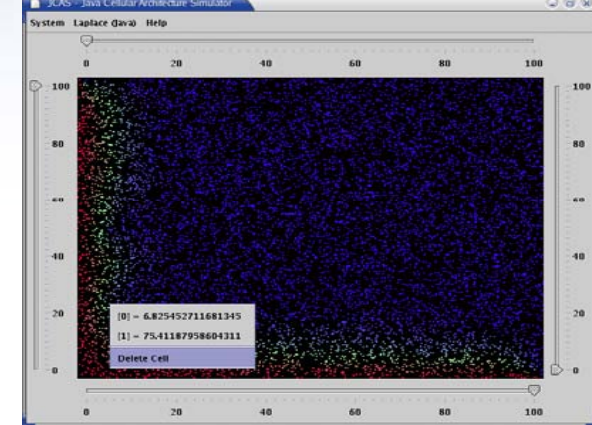


g Assertions

Overview

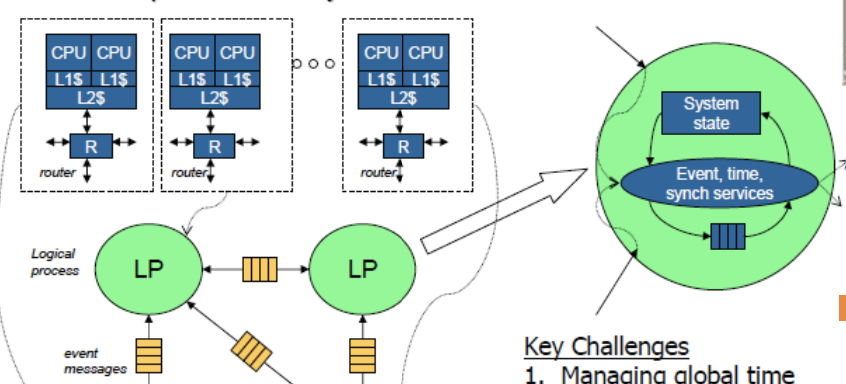


FAST



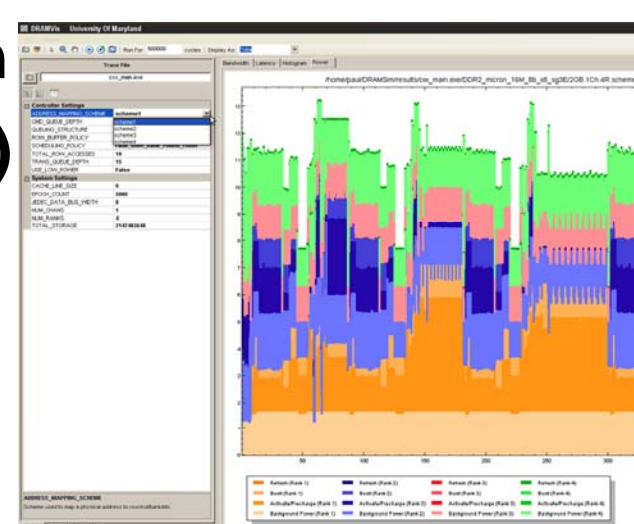
JCAS Vizualizer

Example Modeled System



Key Challenges
1. Managing global time

- IAA Simulation effort is a community effort
- Seeking more partners...
- Current consortium
 - Sandia (Structural Simulation Tool)
 - ORNL (Scalable application model)
 - U. Maryland (DRAMSim II)
 - U.Texas-Austin (FAST)
 - Georgia Tech
 - JCAS (ORNL)
 - Seshat (SNL)



DRAMSim II



Seeks to be...

multiscale

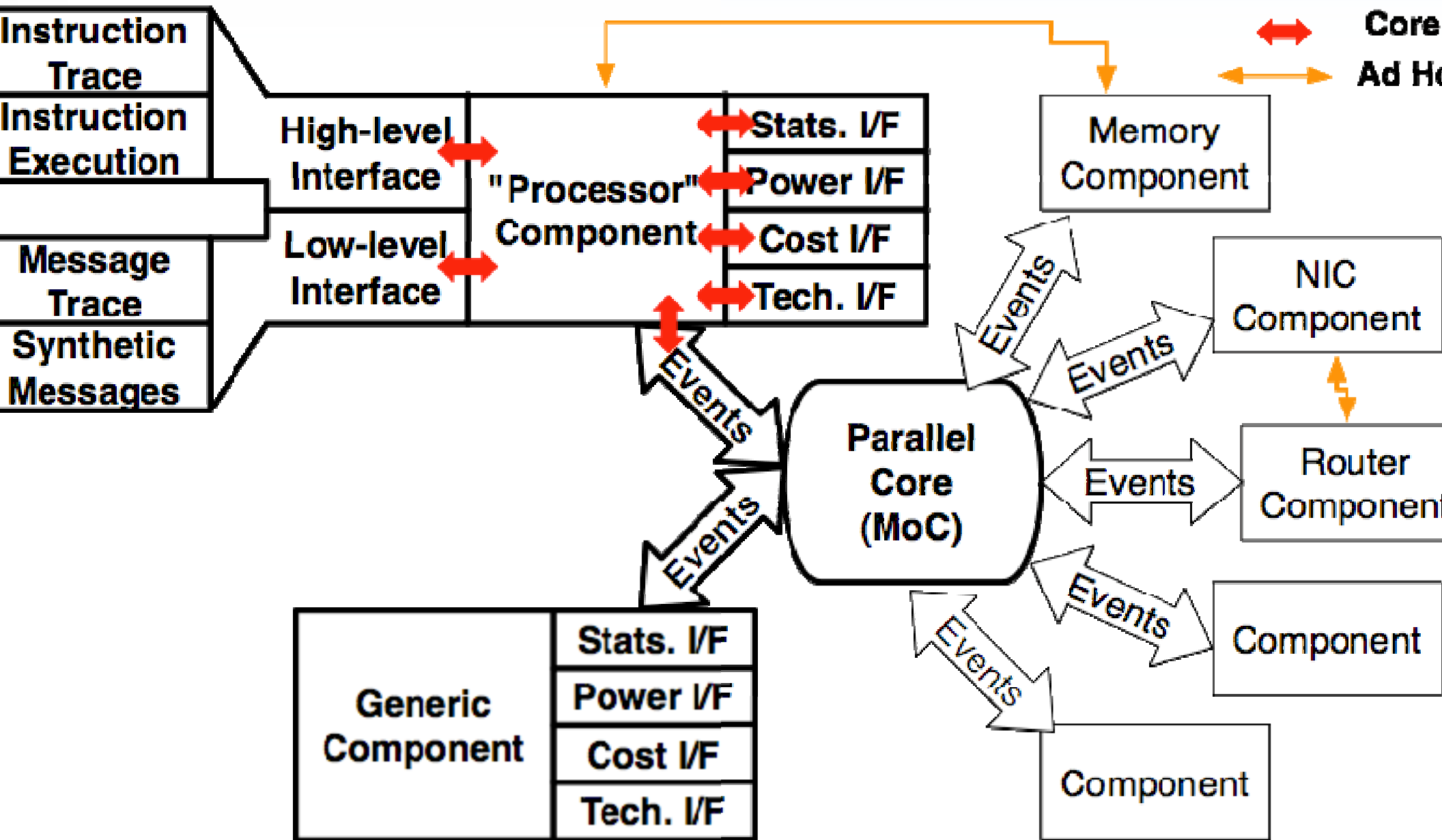
parallel

realistic

**Media is actively seeking partners, requirements, and
input on the SST**

How can the SST be useful for YOU?

Bonus



Separate Software/Front-End from Hardware/Timing/Back-End
Standard interfaces for power, area, cost?