# Improvements to the Structural Simulation Toolkit

**Arun Rodrigues, Keren Bergman, David Bunde, Elliot Cooper-Balis, Kurt Ferreira, K. Scott Hemmert, Brian Barrett, Cassandra Versaggi, Robert Hendry, Bruce Jacob, Hyesoon Kim, Vitus Leung, Michael Levenhagen, Mitchelle Rasquinha, Rolf Riesen, Paul Rosenfeld, Maria del Carmen Ruiz Varela, Sudhakar Yalamanchili**

# View of the Simulation Problem

**Scale.....**

| Many Cores + Memory | X | Many Many Nodes | X | Many Many Many Threads |

**Multiple Audiences.....**

| Network Processor System | X | Application writers purchasers designers | X | system procurement algorithm co-design architecture research language research | X | present systems future systems |

**Complexity.....**

| Multi-Physics Apps Informatics Apps | X | Communication Libraries Run-Times OS Effects | X | Existing Languages New Languages |

**Constraints.....**

| Performance Cost | Power Reliability | Cooling Usability | Risk Size |

# Worldwide Impact

"**Total power used by servers [in 2005] represented ... an amount comparable to that for color televisions.** "
-ESTIMATING TOTAL POWER CONSUMPTION BY SERVERS IN THE U.S. AND THE WORLD, Jonathan G. Koomey

| | |
|---|---|
| 3741e9 KW-Hrs | Total US power consumption |
| * 3-4% | used by computers (>2% servers, >1% household computer use) |
| = 112 - 150e9 KW-Hrs | US Computer power consumption |
| * $0.1 $/KW-Hr | Retail cost, US Average 2009 |
| = $11 - $15 | Billion US$ in compute power |
| * 3-5 | in 2005 US was roughly 1/3 of servers, by power. |
| = $33 ( 🇶🇦 ) - $75 ( 🇲🇲 ) | Billion US$ in worldwide computer power |
| * 15-35% | DRAM memory power |
| = $5 ( 🇱🇾 ) - $25 ( 🇱🇻 ) | BIllion in US$ in DRAM power |

Tuesday, March 27, 2012

# Major Simulation Challenges

- **Multiple Objectives**
  - Performance used to be only criteria
  - Now, Energy, cost, power, reliability, etc...

- **Scale & Detail**
  - Many system characteristics require detail to measure
  - Detailed simulation takes too long ($10^4$-$10^5$ slower than realtime)

- **Accuracy**
  - Systems more complex
  - Vendors don't reveal necessary details

# Major Simulation Challenges

- **Multiple Objectives**
  - Performance used to be only criteria
  - Now, Energy, cost, power, reliability, etc...

- **Scale & Detail**
  - Many system characteristics require detail to measure
  - Detailed simulation takes too long ($10^4$-$10^5$ slower than realtime)

- **Accuracy**
  - Systems more complex
  - Vendors don't reveal necessary details

Performance

Energy

# SST Simulation Project Overview

## Goals

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

## Status

- Current Release (2.1) at code.google.com/p/sst-simulator/
- Includes parallel simulation core, configuration, power models, basic network and processor models, and interface to detailed memory model

## Technical Approach

- Parallel
  - Parallel Discrete Event core with conservative optimization over MPI
- Holistic
  - Integrated Tech. Models for power
  - McPAT, Sim-Panalyzer
- Multiscale
  - Detailed and simple models for processor, network, and memory
- Open
  - Open Core, non viral, modular

## Consortium

- "Best of Breed" simulation suite
- Combine Lab, academic, & industry

# 1. Parallel Implementation

- Implemented over MPI
- Configuration, partitioning, initialization handled by core
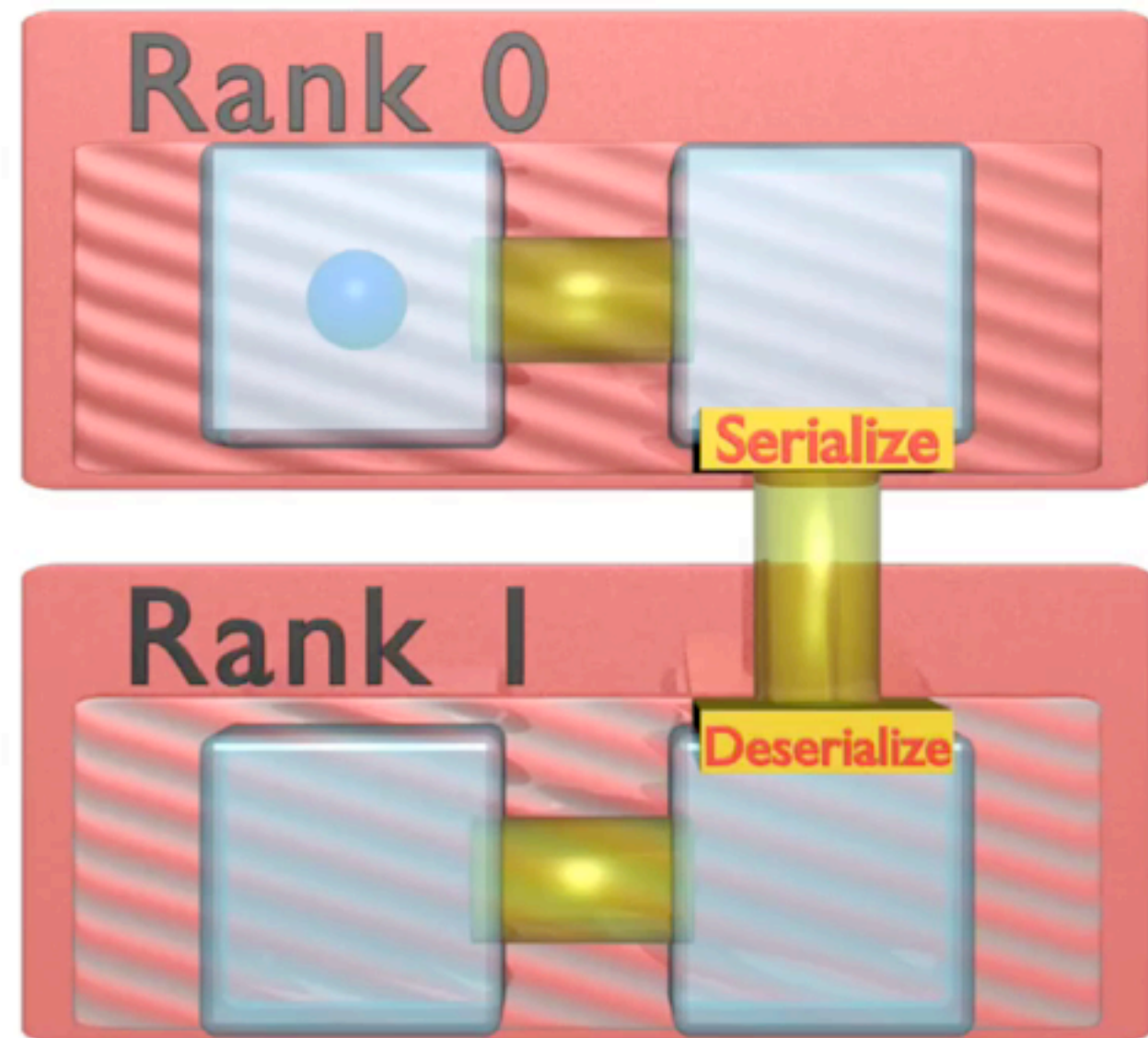- Conservative, distance-based optimization

# 1. Parallel Implementation

- Implemented over MPI
- Configuration, partitioning, initialization handled by core
- Conservative, distance-based optimization
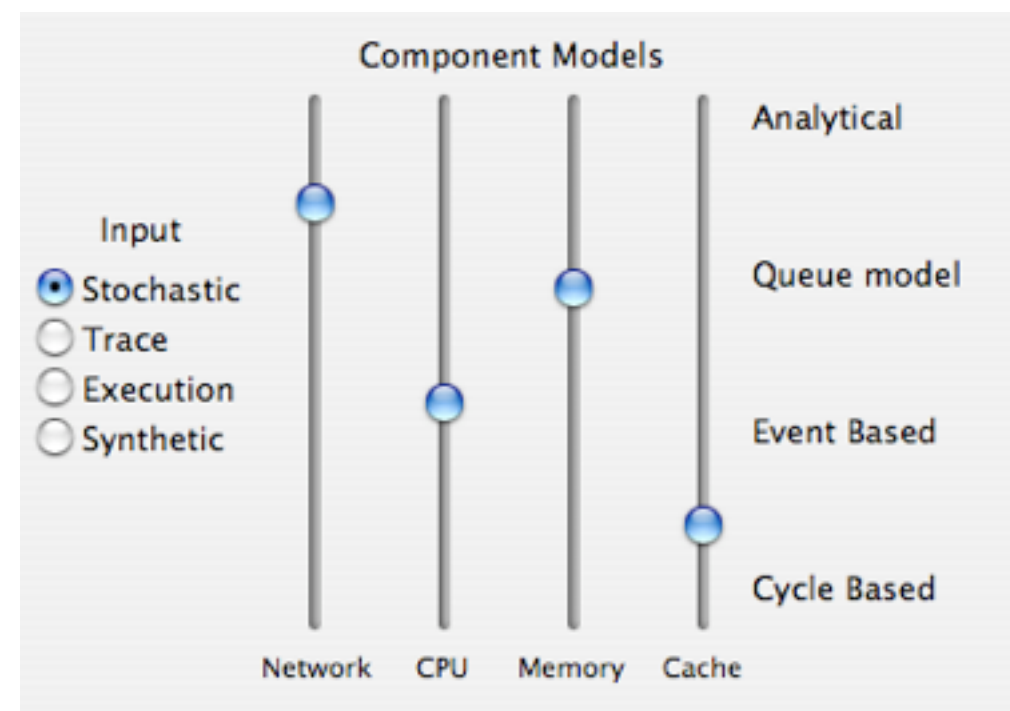
# Message Handling

- **SST core transparently handles message delivery**

- **Detects if destination is local or remote**

- **Local messages delivered to local queues**

- **Remote messages stored for later serialization and remote delivery**
  - **Boost Serialization Library used for message serialization**
  - **MPI used for transfer**

- **Ranks synchronize based on partitioning**

# 2. Multi-Scale

- **Goal: Enable tradeoffs between accuracy, flexibility, and simulation speed**
  - **No single "right" way to simulate**
  - **Support multiple audiences**
- **High- & Low-level interfaces**
  - **Allows multiple input types**
  - **Allows multiple input sources**
    - **Traces, stochastic, state-machines, execution...**

| | High-Level | Low-Level |
|---|---|---|
| **Detail** | Message | Instruction |
| **Fundamental Objects** | Message, Compute block, Process | Instruction, Thread |
| **Static Generation** | MPI Traces, MA Traces | Instruction Trace |
| **Dynamic Generation** | State Machine | Execution |



**Multiscale Parameters**

# Diversity!

**Complexity/Detail** (vertical axis)

**Execution Time** (horizontal axis)

**Zesto**

**MacSim**

**M5 O3**

**DRAMSim**

**genericProc**

**RS Router**

**Stochastic**

**simpleRouter**

**Macro**

**scheduleSim**

**ResiliencySim**

**Commonalities**
- **Discrete Event or time stepped**
- **Amenable to event counting for power modeling**
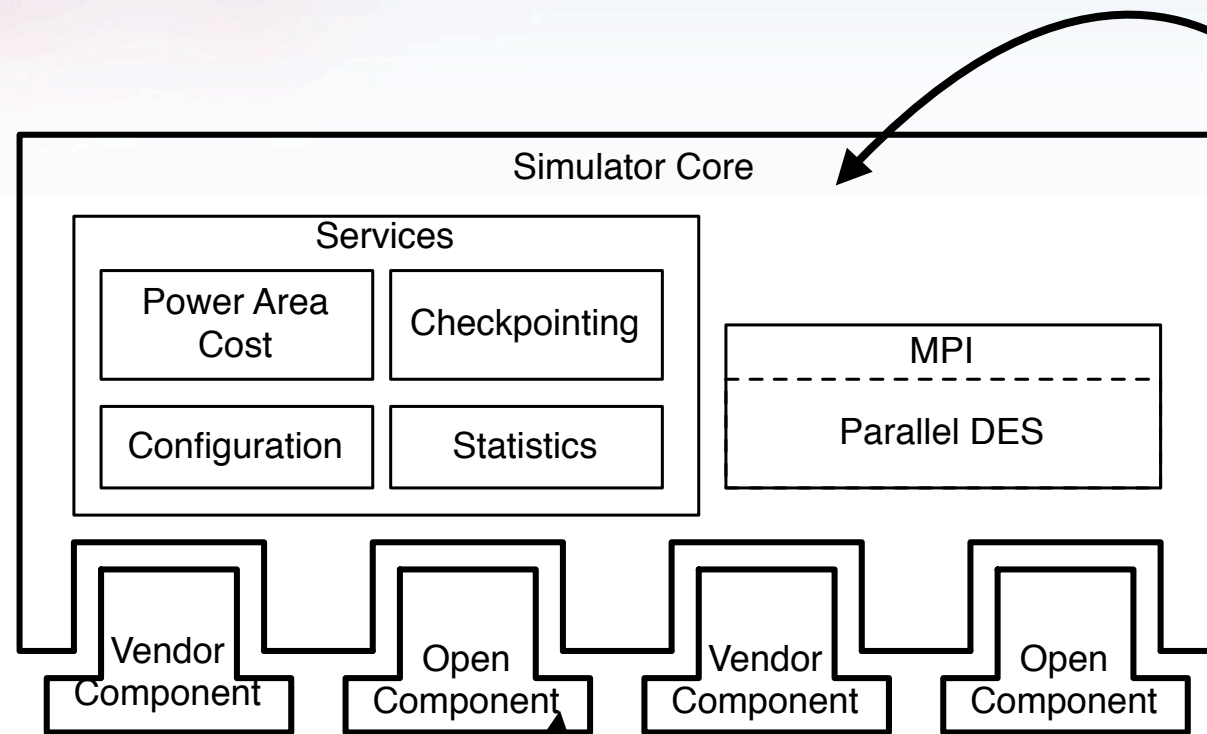
# 3. Holistic Simulation



- **Design space includes much more than simple performance**

- **Create common interface to multiple technology libraries**
  - **Power/Energy**
  - **Area/Timing estimation**

- **Make it easier for components to model technology parameters**

# 4. Open Simulator Framework



- **Simulator Core will provide...**
  - –**Power, Area, Cost modeling**
  - –**Checkpointing**
  - –**Configuration**
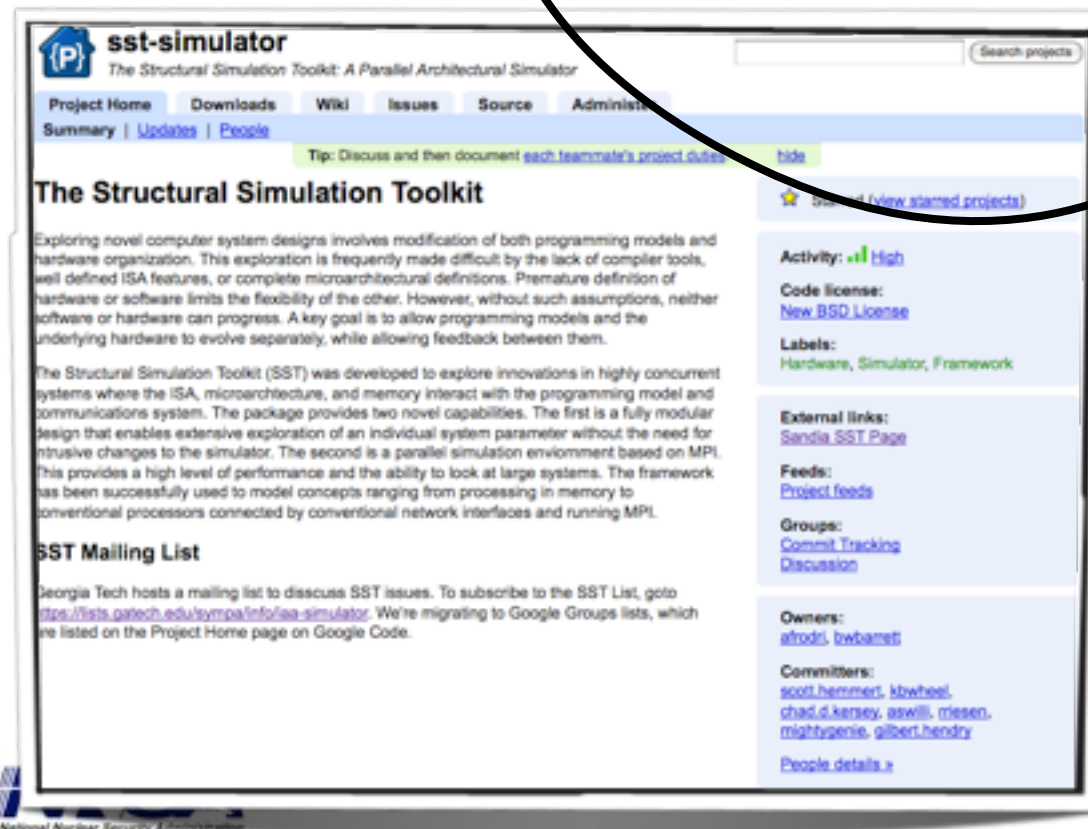  - –**Parallel Component-Based Discrete Event Simulation**
- **Components**
  - –**Ships with basic set of open components**
  - –**Industry can plug in their own models**
    - •**Under no obligation to share**
- **Open Source (BSD-like) license**
- **SVN hosted on Google Code**

# Improvements
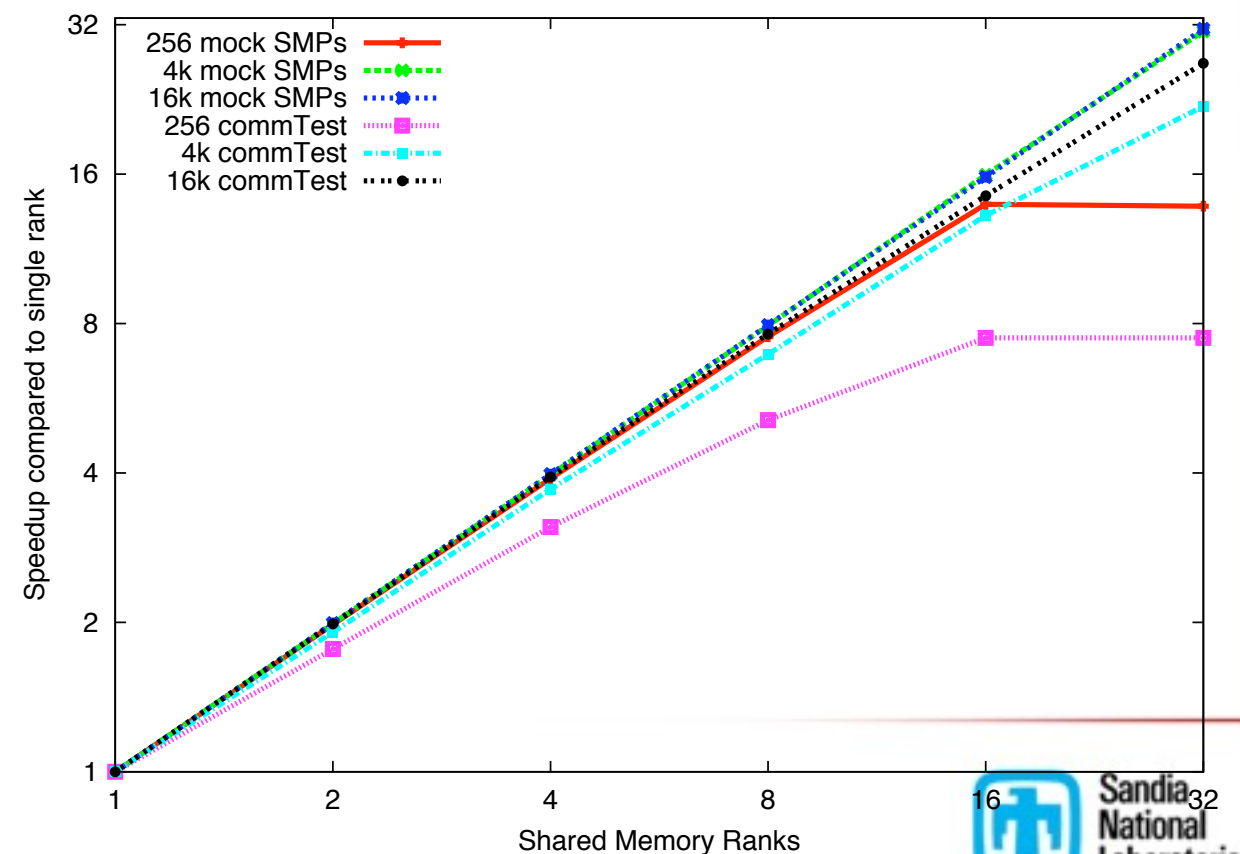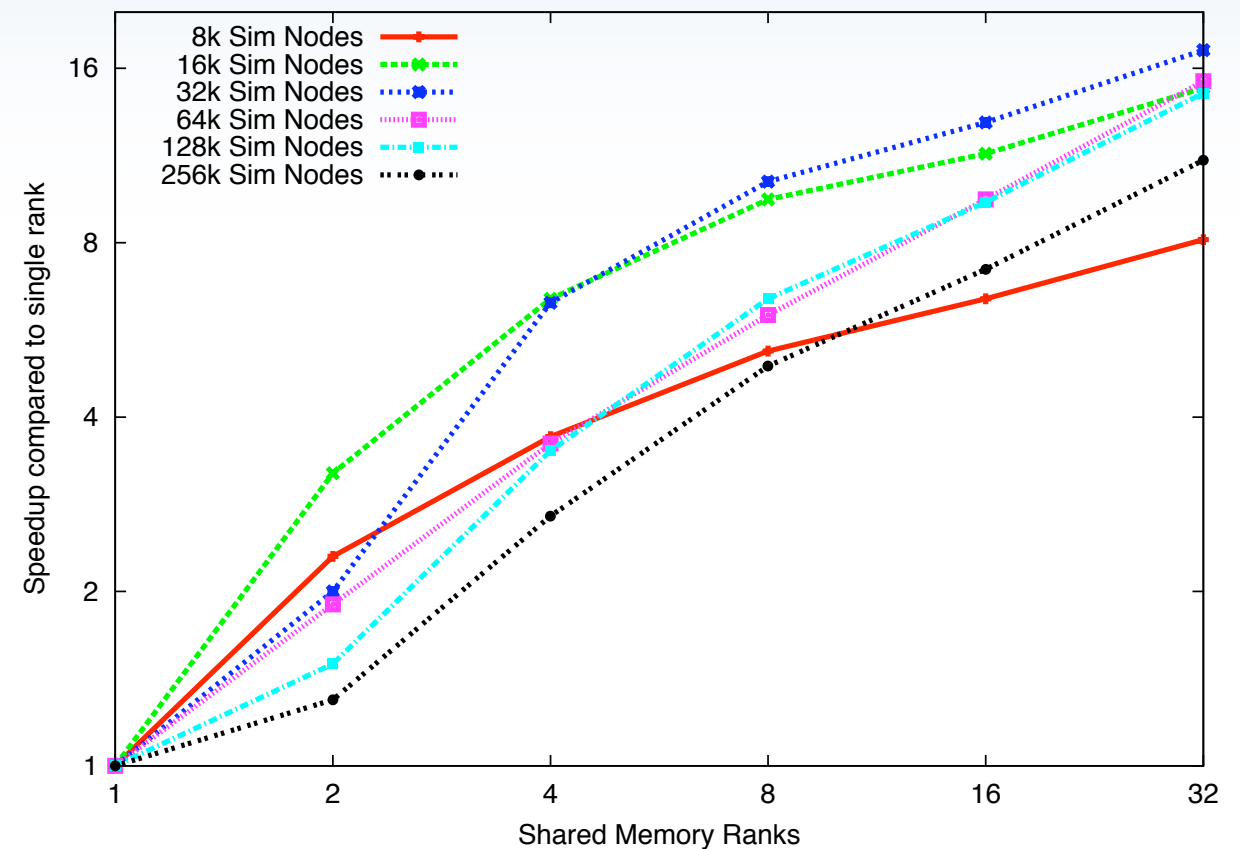
# Scalability

- **Test setup: OpenMPI 1.4.3 Shared memory**

- **Portals NIC Simulation**
  - **State machine application, detailed NIC & Router**
  - **Best scaling w/ 32K simulated nodes (17.2x speed up on 32 host nodes)**
  - **May be limited by memory bandwidth**

- **Synthetic components**
  - **Poor scaling with 256 simulated SMP nodes**
  - **85%+ scaling efficiency with 4K nodes**

# IRIS & PhoenixSim

- **IRIS**
  - **Simulates pipelined, cycle-accurate router**
  - **Models variety of Network-on-Chip (NoC) and inter-node interconnects**
  - **Parameterized buffer sizes, virtual channels, routing & arbitration**
  - **Modular architecture**

- **Phoenix Sim**
  - **Physically detailed photonic network simulator**
  - **Component: waveguides, modulators, detectors, filters, switches, lateral couplers, and lasers**

# GeM5

- **M5: Modular platform for computer system architecture research, encompassing system-level architecture as well as processor microarchitecture.**

- **Provides detailed, full-system CPU models for x86, ARM, SPARC, Alpha**

- **Integrated at SST Component, allows interaction with SST models, and parallel execution**

- **Currently tested up to 256 nodes.**

# SST/GeM5 Integration

- **Goals: Run GeM5 in parallel, connect with other SST components**
- **High parallel efficiency**
- **Changes**
  - **Replaced Python-based configuration w/ XML or C++-based system**
  - **Encapsulated GeM5 as an SST Component, GeM5 event Q driven by SST clock**
  - **Created translator SimObjects to connect to SST links**
  - **Changes made to GeM5 loader to avoid use of async (untimed out-of-band) messages**
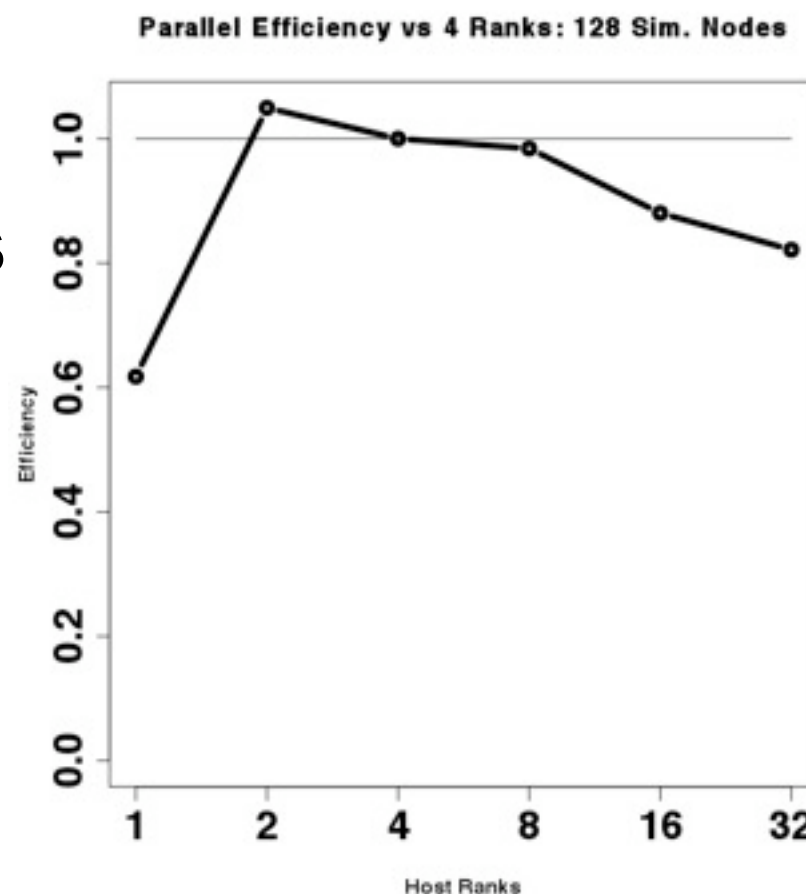


gem5 SST Component

CPU → L1, L1 → BUS → L2 → MemBus+ (Memory BUS) → PhysMem; DRAMSim ↔ Memory BUS; Memory BUS ↔ IO Bridge ↔ IO Bus; Syscall Handler; Translator; SS Router ↔ Portals NIC



Parallel Efficiency vs 4 Ranks: 128 Sim. Nodes

Efficiency vs Host Ranks

# MacSim GPU Simulator

- **Trace driven simulator**
  - Internal RISC-like uops, with threading information
  - X86 Captured w/ PIN tool
  - CUDA Captured through Ocelot emulation tool
- **Models for**
  - OoO & GPU-like cores
  - Cache memory (inc. texture & Constant caches)
  - Timing & Power
- **Allows ties in to MacPAT, thermal modeling, RAMP reliability model**

# Additional Components

- **BOBSim**
  - **Models Buffer-on-Board Memory systems**
  - **Hardware verified using RTL from Micron**
  - **Includes detailed power model**

- **Reliability Simulation**
  - **extend SSTfor simulating resilience mechanisms of a large-scale system (e.g. Checkpointing, use of NVRAM, burst buffers)**
  - **Models: Router, Storage, Node**
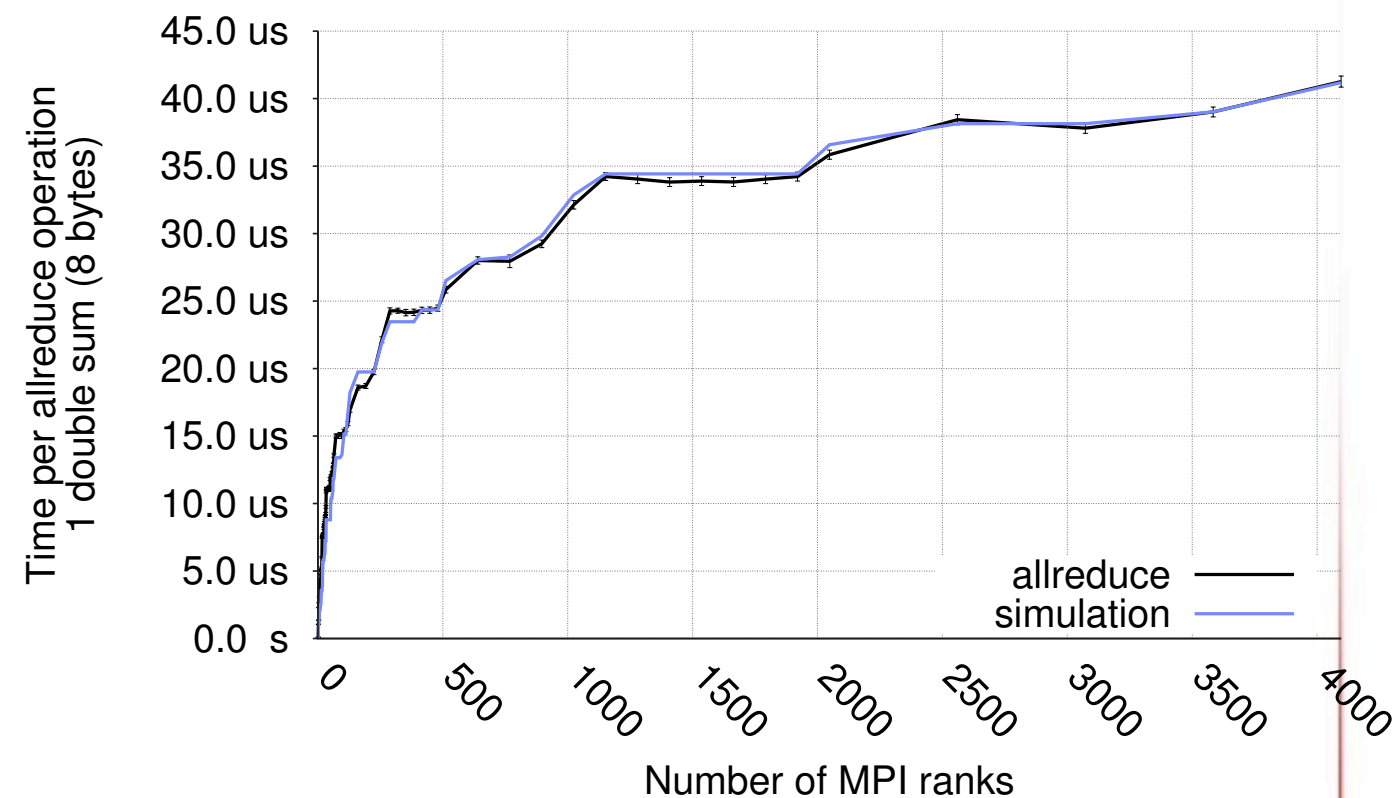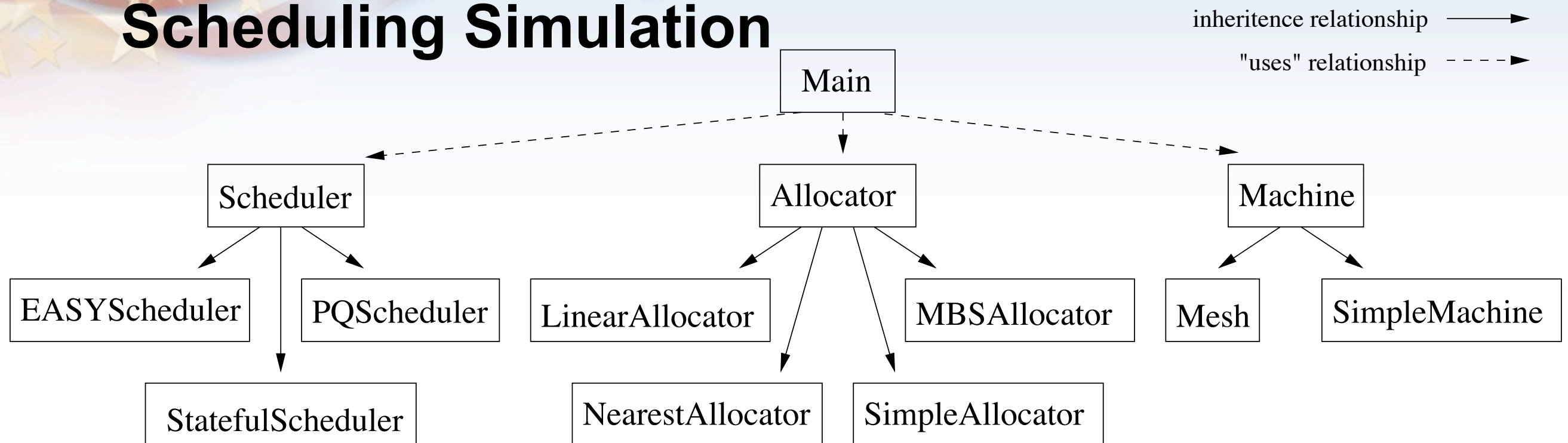  - **Validated against large systems**

# Scheduling Simulation

```
                          ┌──────┐
                          │ Main │
                          └──────┘
         ┌───────────────────┼───────────────────┐
         ↓                    ↓                    ↓
   ┌───────────┐        ┌───────────┐        ┌───────────┐
   │ Scheduler │        │ Allocator │        │ Machine   │
   └───────────┘        └───────────┘        └───────────┘
     ↓      ↓            ↓    ↓    ↓            ↓      ↓
┌──────────────┐ ┌──────────────┐ ┌──────────────┐ ┌──────────────┐ ┌──────┐ ┌───────────────┐
│EASYScheduler │ │ PQScheduler  │ │LinearAllocator│ │ MBSAllocator │ │ Mesh │ │ SimpleMachine │
└──────────────┘ └──────────────┘ └──────────────┘ └──────────────┘ └──────┘ └───────────────┘
         ↓                          ↓          ↓
┌──────────────────┐      ┌──────────────────┐ ┌──────────────────┐
│ StatefulScheduler│      │ NearestAllocator │ │ SimpleAllocator  │
└──────────────────┘      └──────────────────┘ └──────────────────┘
```

- **new high-level system simulation capabilities**
- **using job traces from large systems**
- **simulates the decisions being made**
  - **scheduling**
  - **allocation**
- **working on task mapping and tighter integration with the lower-level simulations that SST is already able to do**

# Summary

- **The SST is parallel, scalable, modular, and open**
- **Recent Improvements**
  - **Scalability**
  - **Processor, network, memory, IO system models**
- **Future**
  - **Improving component interoperability**
  - **Validation**
    - **Emphasis on mixed-resolution simulation**
  - **Scaling**
    - **Emphasis on distributed memory systems**
  - **New models**
    - **Stacked memory**
    - **Exascale Supercomputers**

# Questions?

## [code.google.com/p/sst-simulator/](code.google.com/p/sst-simulator/)

# Component Library

- **Parallel Core v2**
  - **Parallel DES layered on MPI**
  - **Partitioning**
  - **Configuration & Checkpointing**
  - **Power modeling**
- **Technology Models**
  - **McPAT, Sim-Panalyzer, IntSim, Orion and custom power/energy models**
  - **HotSpot Thermal model**
  - **Working on reliability models**
- **Components**
  - **Processor: Macro Applications, Macro Network, NMSU, genericProc, state-machine, Zesto, GeM5, GPGPU**
  - **Network: Red Storm, simpleRouter, GeM5**
  - **Memory: DRAMSim II, Adv. Memory, Flash, SSD, DiskSim**



**SST Simulator Core**