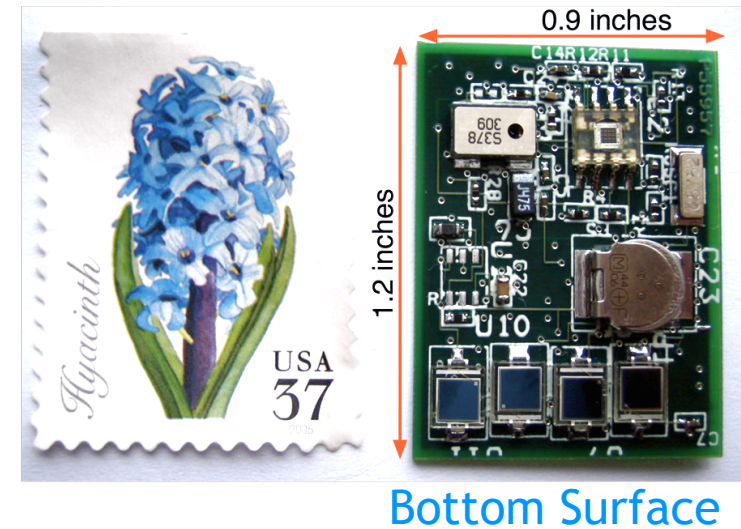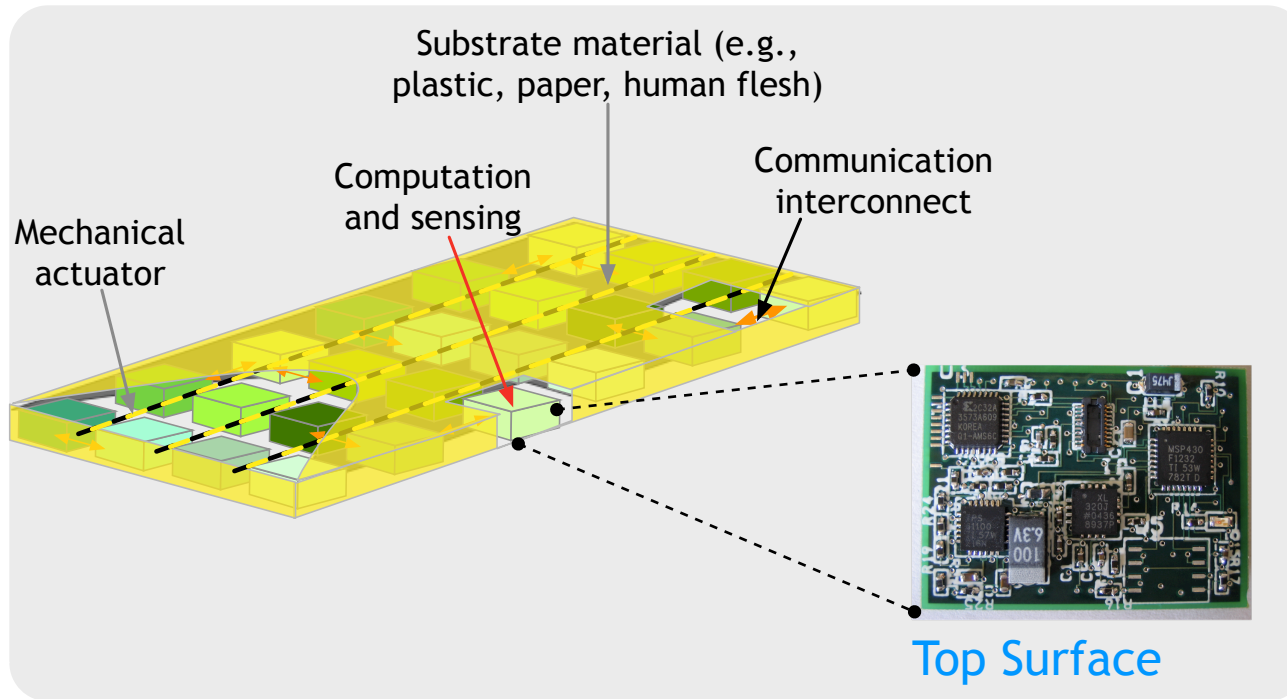# Implementation of a Distributed Full-System Simulation Framework as a Filesystem Server

## Phillip Stanley–Marbell
Carnegie Mellon University

# Motivation



Substrate material (e.g., plastic, paper, human flesh)

Computation and sensing

Communication interconnect

Mechanical actuator

Top Surface

0.9 inches

1.2 inches

Bottom Surface

["An 0.9x1.2 Energy-Harvesting System with Custom Multi-Channel Communication Interface", *IEEE DATE'07*, Nice, France]

- Context

  – Investigating highly-integrated networks of compute/sensing/actuation systems

  – We currently cannot afford ($) to build systems with thousands of nodes

  – Simulation permits the investigation of large scale systems

  – Simulation is not a substitute for actual hardware

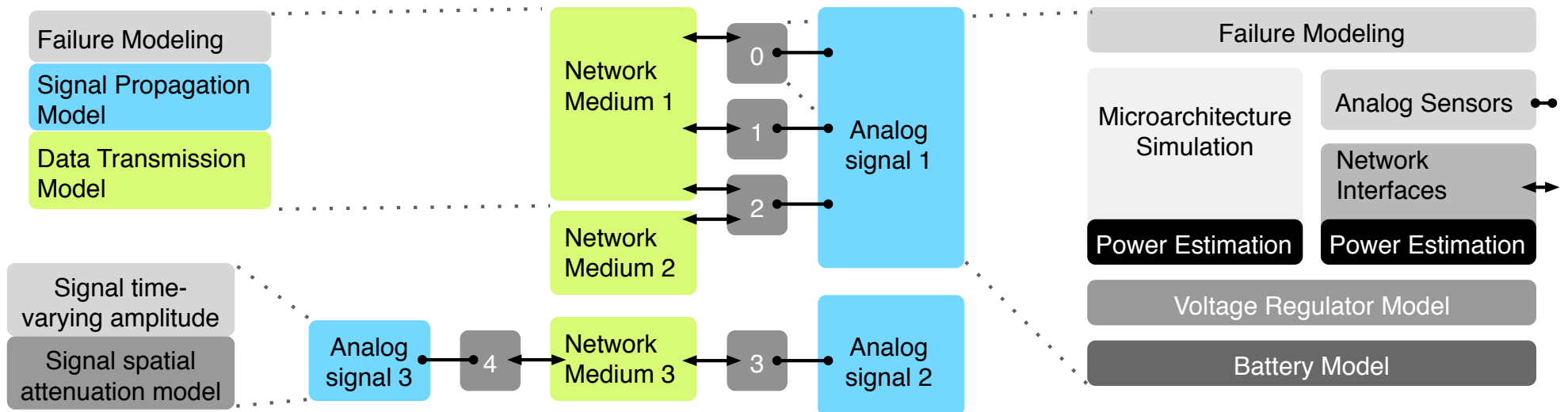- Challenge: simulating large scale (thousand+ node) systems

# Outline

- Motivation

- Simulation Framework Overview

- Distributed Simulation

- Multi-Platform Packaging

- Summary

# Sunflower Full-System Simulator

- ## Simulation Engine Models:

  - Computation — at instruction execution level, for two different ISAs

  - Communication — at the MAC and PHY levels

  - Compute & network power dissipation — includes instruction-level power models

  - Batteries and voltage regulators — models for several batteries and regulators

  - Device and networking faults — bit-level logic upsets and node / network failures

  - Physical phenomena external to hardware — location in three-space, attenuation

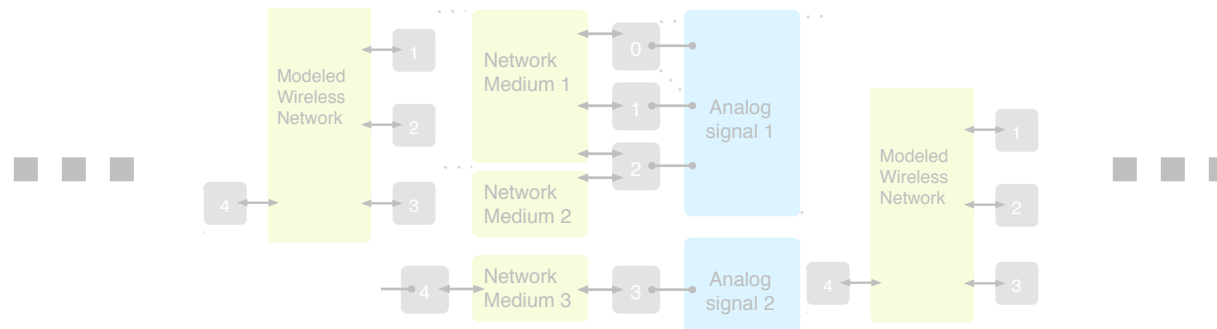- ## Example composition of a simulated system:



| Failure Modeling | | | Network Medium 1 | 0 | Analog signal 1 | Failure Modeling | |
| Signal Propagation Model | | | | 1 | | Microarchitecture Simulation | Analog Sensors |

["Sunflower: Full-System Embedded Microarchitecture Evaluation", *HiPEAC '07*, Ghent, Belgium]

# Sunflower Full-System Simulator

- ## Simulation engine implemented in C
  - – Originally implemented as a standalone application on Unix

- ## Challenges
  - – Simulating large systems (thousands of nodes) can be tasking, even on a high-end workstation
  - – Especially so, when modeling all the aforementioned components



  - – Can split simulation across multiple workstations
  - – Nodes simulated on separate hosts may communicate : must handle that
  - – Must keep passage of time synchronous across portions of partitioned simulation

# Outline

# Distributed Simulation

- ## Challenges
  - An efficient / easy way to connect state of portions of sim across hosts
  - A means of ensuring the state of simulation across hosts is consistent
  - An interface to keep track of hosts, global simulation state, etc.
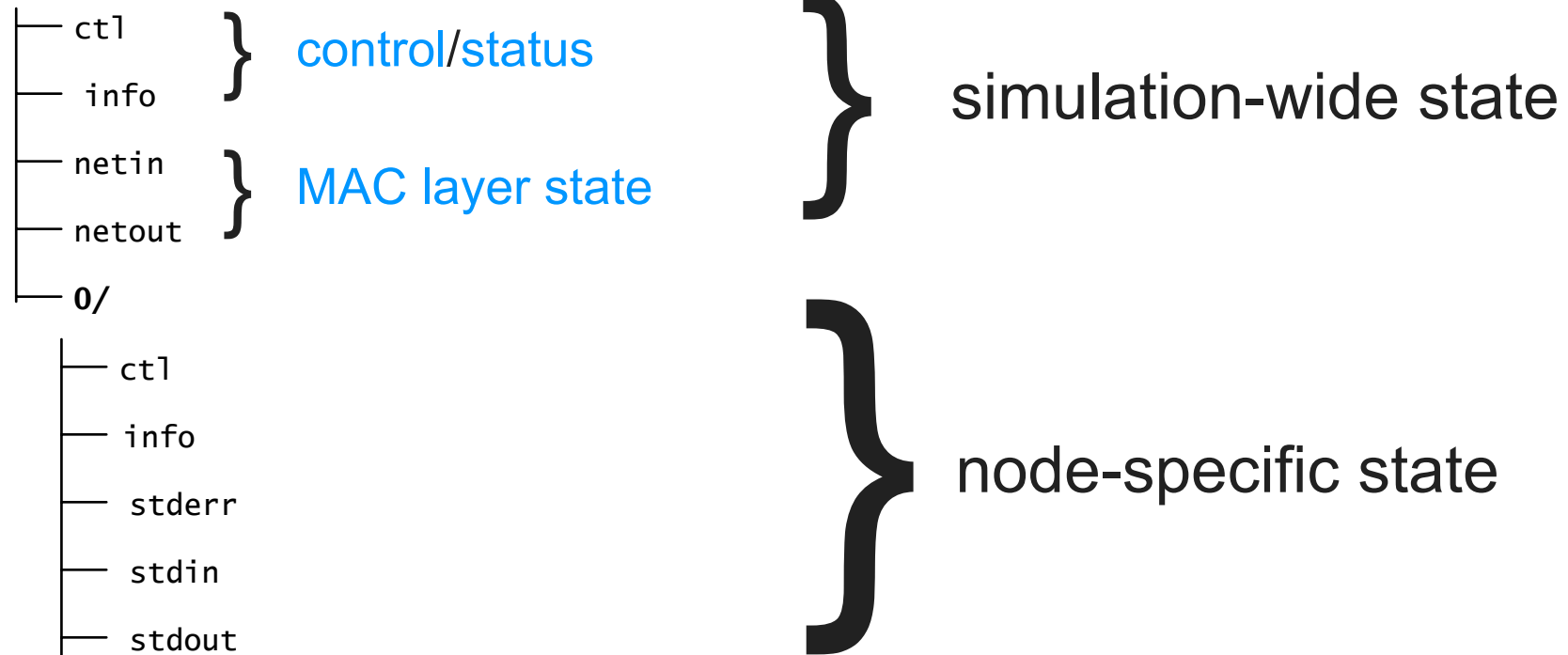
- ## Approach
  - *There are many ways you could do this...*
  - Chose to take advantage of ease of connecting systems with Inferno

- ## Implementation
  - Simulation engine compiled as a library, linked against emu (chose not to use libstyx)
  - All relevant state on each host exposed as a filesystem, via a device driver interface
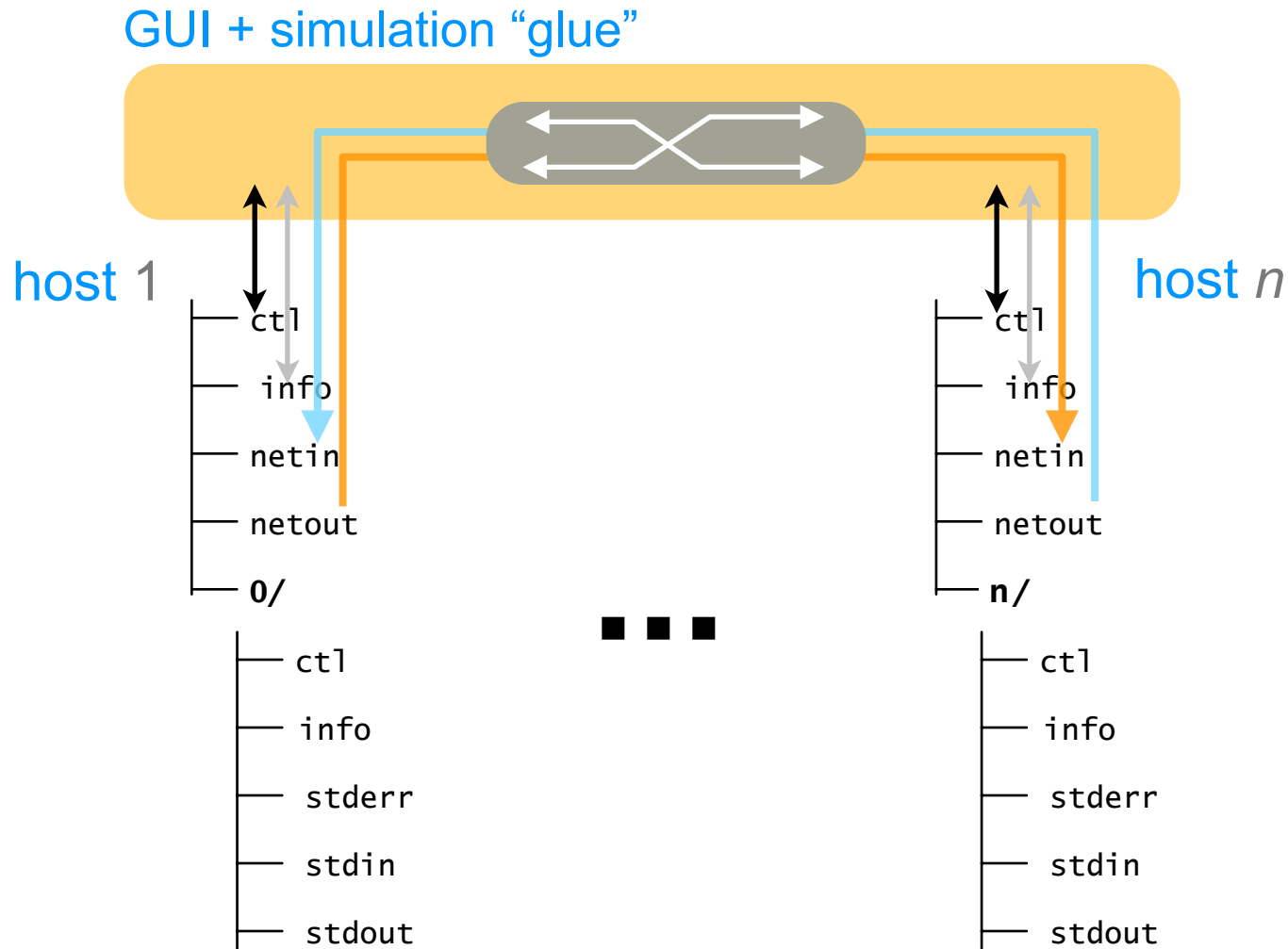  - "Glue" application connects together name spaces, keeps engine state consistent

# Simulation Engine Interface

```
engine.attachname/
├── ctl          } control/status    }  simulation-wide state
├── info
├── netin        } MAC layer state
├── netout
└── 0/                                }  node-specific state
    ├── ctl
    ├── info
    ├── stderr
    ├── stdin
    └── stdout
```

- Interface to simulation engine is a device driver, `devsf`, #j
  - Dynamic filesystem interface w/ one line directory per simulated node (processor+batt, etc)
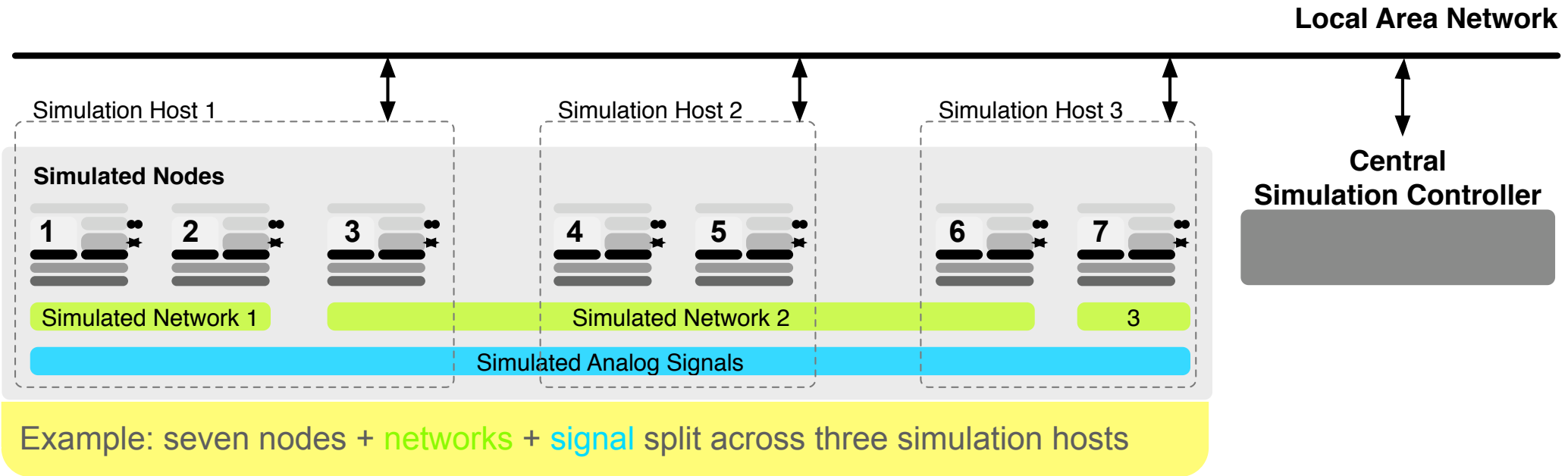  - Exposes simulation engine state (e.g., simulated MAC-layer frames via `netin`, `netout`)

# Simulation Glue Logic

GUI + simulation "glue"



- "Glue" application (implemented in Limbo)
  - Engine filesystems from different simulation hosts mounted in name space of glue app
  - Interconnects simulated networks on different hosts
  - Implements facilities for synchronizing virtual time across portions of simulated system

# Distributed Simulation

**Simulation Host 1**

**Simulated Nodes**

**1**   **2**   **3**

Simulated Network 1

**Simulation Host 2**

**4**   **5**

Simulated Network 2

**Simulation Host 3**

**6**   **7**

3

Simulated Analog Signals

**Central Simulation Controller**

Example: seven nodes + networks + signal split across three simulation hosts

- ## Synchronization issues
  - – Simulation rates across hosts may vary, but global timebase must remain synchronized
  - – Well known issues, already investigated in the area of parallel discrete-event simulation

- ## Synchronization facilities
  - – Time synchronization
  - – Simulation rate synchronization

# GUI + Glue Logic Application

Pull-down menu with shortcuts for common commands

Button shortcuts for "glue" commands

Message output window

clicking on a node makes it the current

node # @ host #

Command input

Summary of statistics for current node

Warning messages

Connected local and remote simulation engines

Error messages

# Outline

- Motivation

- Simulation Framework Overview

- Distributed Simulation

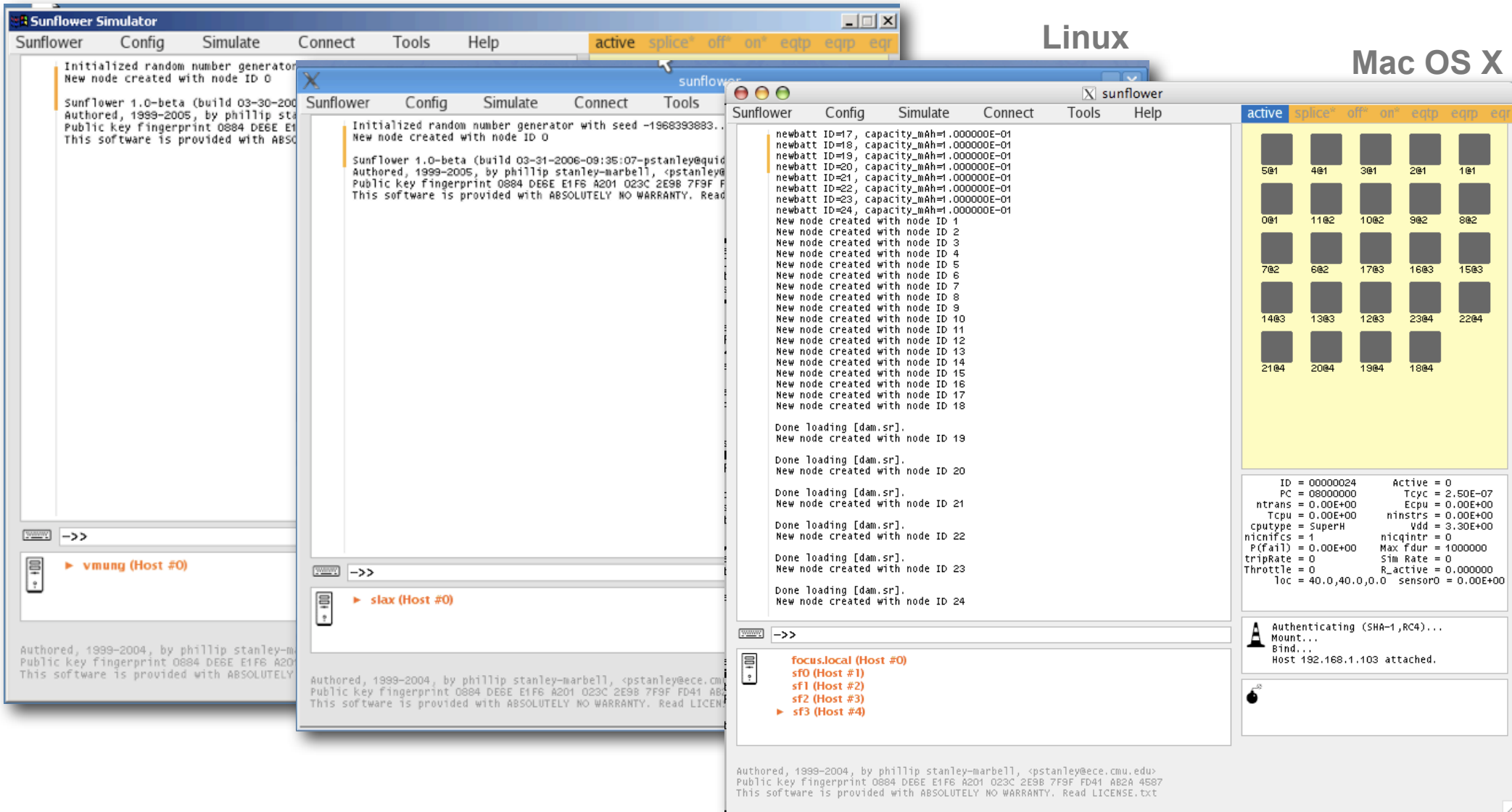- **Multi-Platform Packaging**

- **Summary**

# Multi-Platform Packaging

- Goal — Implementation transparency
  - A single executable binary, indistinguishable from a native app. on host platform
  - Should not require (explicit or automatic) installation of Inferno

- Implementation
  - All the necessary dis executables, fonts, etc. compiled into in-memory root filesystem
  - Simulation engine library & driver interface (#j)
  - Server mode or client mode, w/ GUI or server app. instead of `emuinit.dis` → `sh.dis`
  - Reads / writes from host filesystem like any other host application (via #U)
  - On most platforms, binary is ~2—5 MB; self-contained executable

- Alternative — don't roll filesystem into emu binary image
  - Distribute emu + Inferno filesystem tree (acme-sac does this)
  - I currently think a single-executable approach is cleaner for our purposes

↵

# Multi-Platform Packaging



- Distributed as a single executable, no installation reqd.
  - 2.5MB binary on MacOS, 5.2MB binary on Linux, 2.7MB .exe on Windows ↵

# Outline

- Motivation

- Simulation Framework Overview

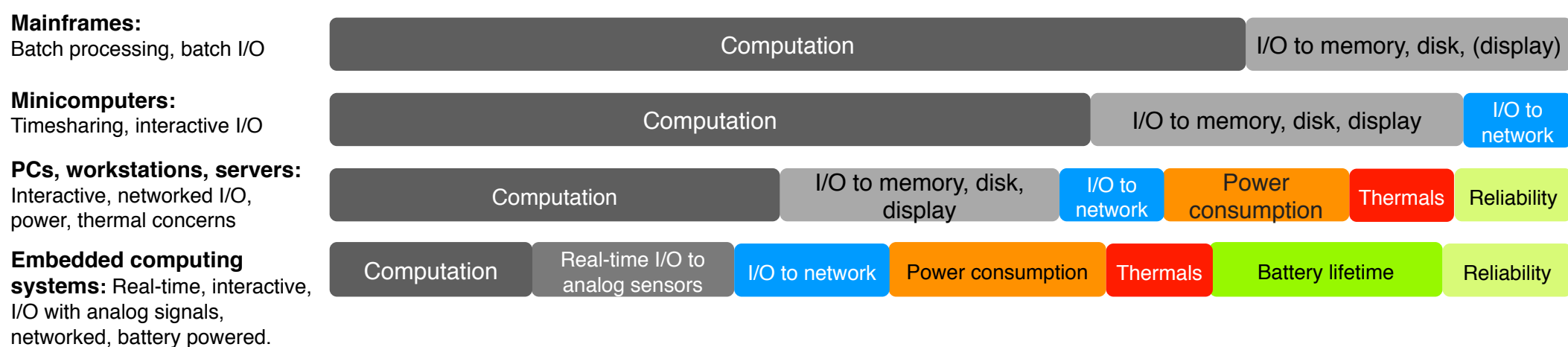- Distributed Simulation

- Multi-Platform Packaging

- **Summary**

# Summary

- ## Sunflower
  - A full-system simulator for networks of embedded systems

- ## Problem
  - Simulation of large networks (thousand+ nodes) is compute- and memory intensive
  - Simulation can be split at the level of individual nodes (with some added work)

- ## Distributed simulation
  - Simulation engine compiled as a Inferno emulator library
  - Device driver (`#j`) interface to simulation state
  - GUI+glue application interconnects simulated state across multiple simulation hosts
  - Integration into Inferno enables easy multi-platform packaging/GUI
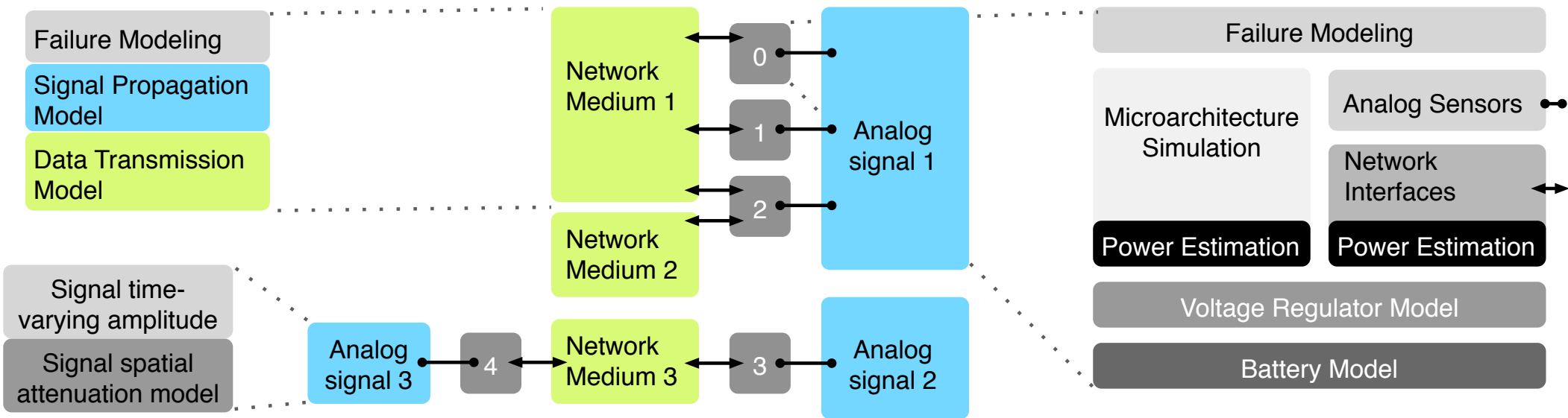
- ## Sources, binaries, documentation
  - `http://www.ece.cmu.edu/~pstanley/sunflower`

# Backup Slides

Backup Slides

# Hardware Emulation Environment — Sunflower Simulator

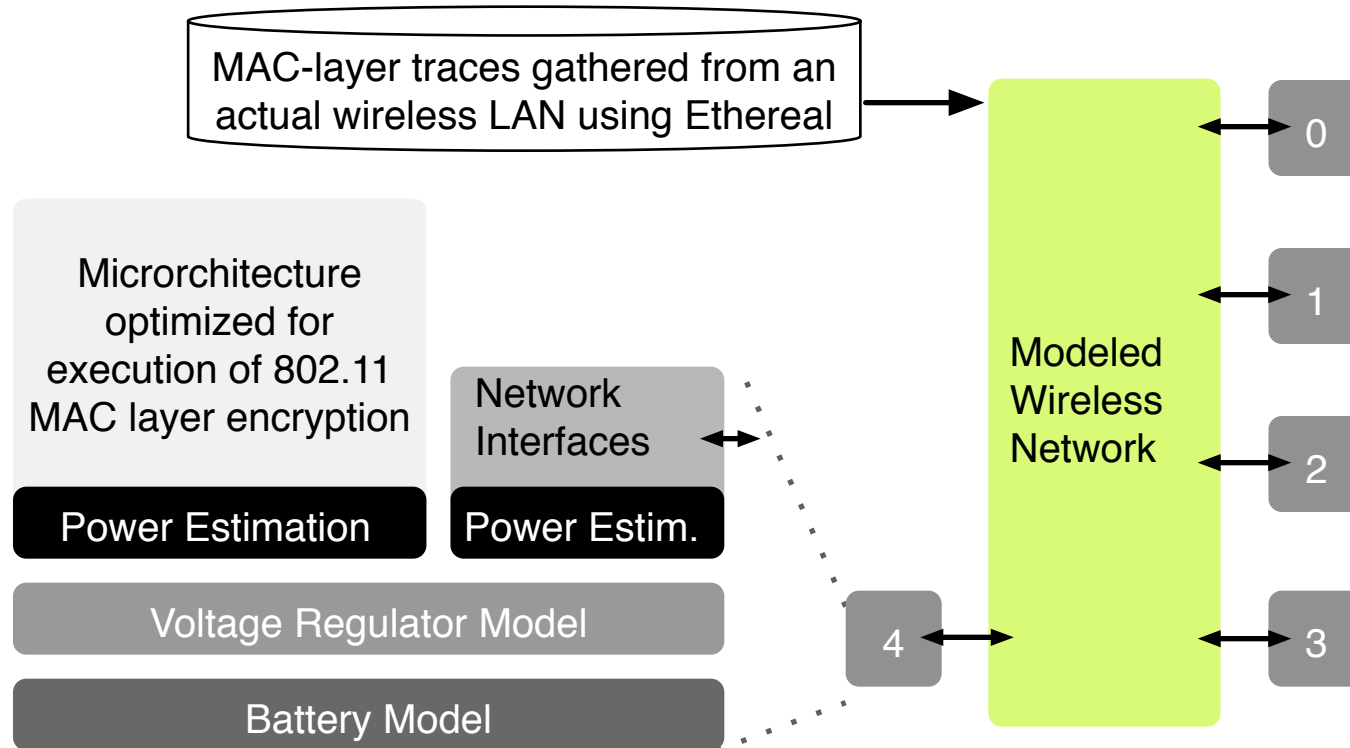Evaluation metrics relevant to microarchitecture's performance

**Mainframes:**
Batch processing, batch I/O

| Computation | I/O to memory, disk, (display) |
|---|---|

**Minicomputers:**
Timesharing, interactive I/O

| Computation | I/O to memory, disk, display | I/O to network |
|---|---|---|

**PCs, workstations, servers:**
Interactive, networked I/O, power, thermal concerns

| Computation | I/O to memory, disk, display | I/O to network | Power consumption | Thermals | Reliability |
|---|---|---|---|---|---|

**Embedded computing systems:** Real-time, interactive, I/O with analog signals, networked, battery powered.

| Computation | Real-time I/O to analog sensors | I/O to network | Power consumption | Thermals | Battery lifetime | Reliability |
|---|---|---|---|---|---|---|

Backup Slides

# Hardware Emulation Environment — Sunflower Simulator

Backup Slides

# Using Real Network Traces

```
┌─────────────────────┐     ┌─────────────────────┐     ┌─────────────────────┐
│   MAC layer trace   │     │  tracetool: conversion of │ │                     │
│ collection (e.g., with │ ──▶ │   MAC layer trace in  │ ──▶ │     Simulation      │
│     Ethereal)       │     │    libpcap format   │     │                     │
└─────────────────────┘     └─────────────────────┘     └─────────────────────┘
```

Backup Slides

# Example: Using Real Network Traces

MAC-layer traces gathered from an actual wireless LAN using Ethereal

Microrchitecture optimized for execution of 802.11 MAC layer encryption

**Power Estimation**

Network Interfaces

**Power Estim.**

Voltage Regulator Model
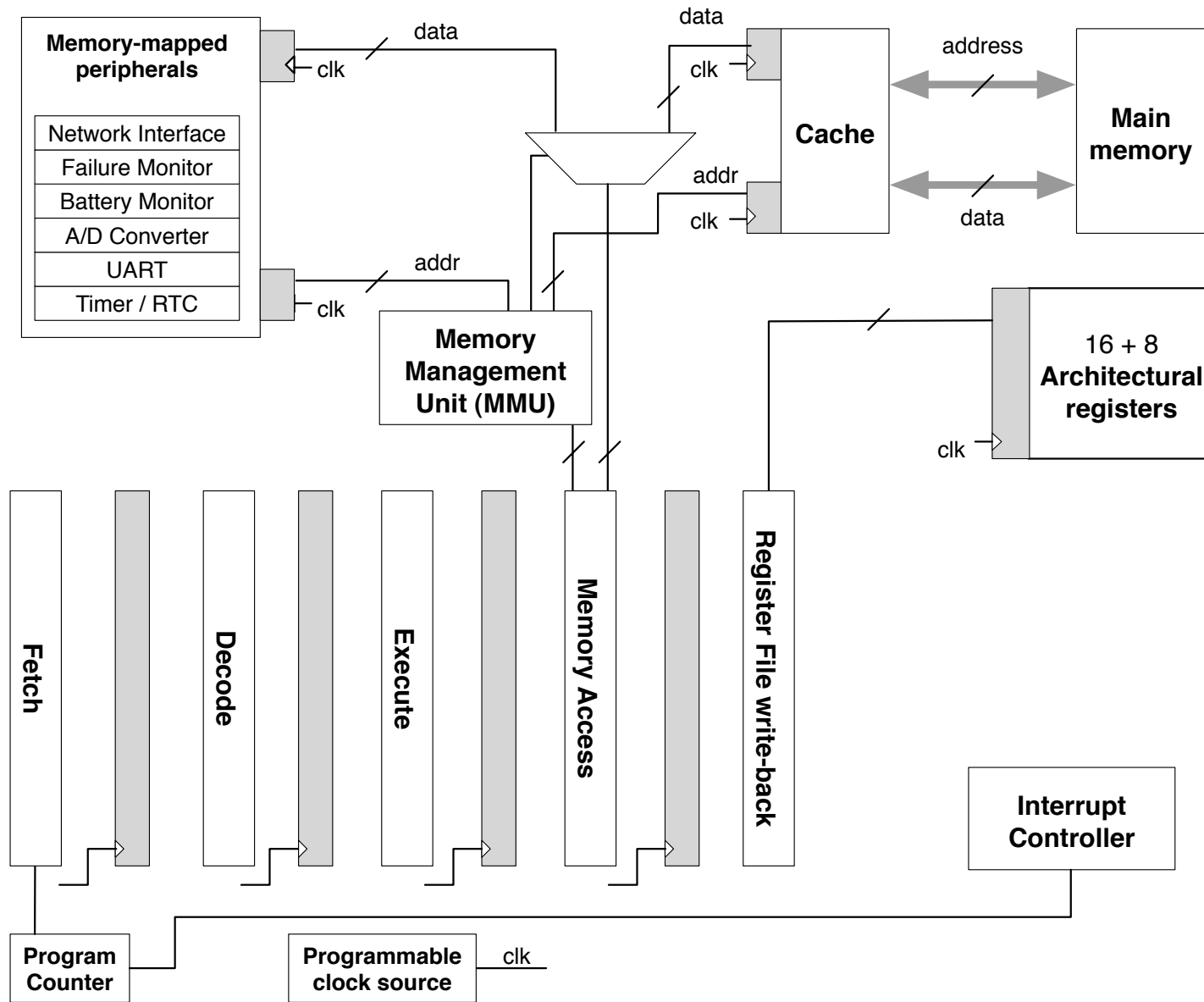
Battery Model

Modeled Wireless Network

0

1

2

3

4

# Modeled 16-bit Microarchitecture (TI MSP430)



= Structures modeled at bit-level, enabling signal transition activity and SEU modeling

Backup Slides
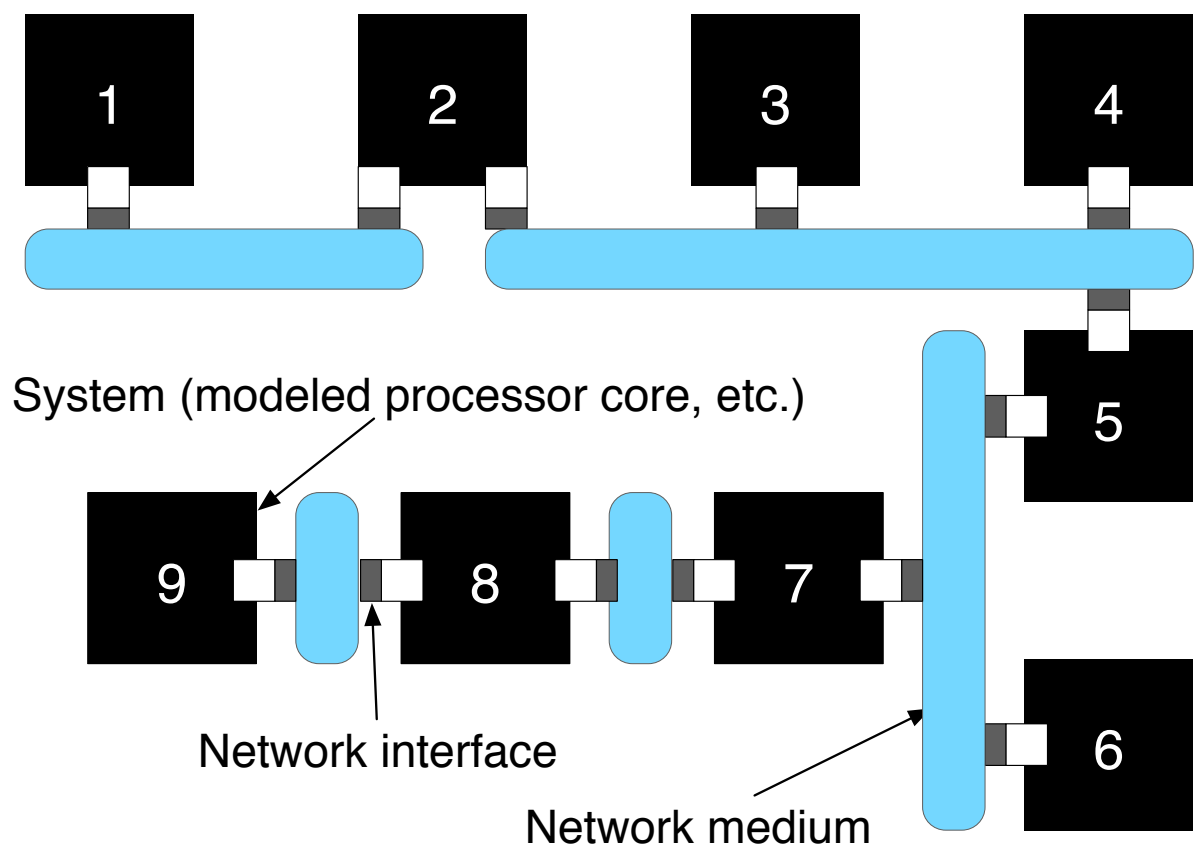
# Modeled 32-bit Microarchitecture (Hitachi SH)

**Memory-mapped peripherals**

| Network Interface |
| Failure Monitor |
| Battery Monitor |
| A/D Converter |
| UART |
| Timer / RTC |

clk

data

data

clk

**Cache**

address

**Main memory**

addr

clk

data

addr

clk

**Memory Management Unit (MMU)**

**16 + 8 Architectural registers**

clk

Fetch

Decode

Execute

Memory Access

Register File write-back

**Interrupt Controller**

**Program Counter**

**Programmable clock source**

clk

= Structures modeled at bit-level, enabling monitoring of signal transition activity and SEUs during simulation

Backup Slides

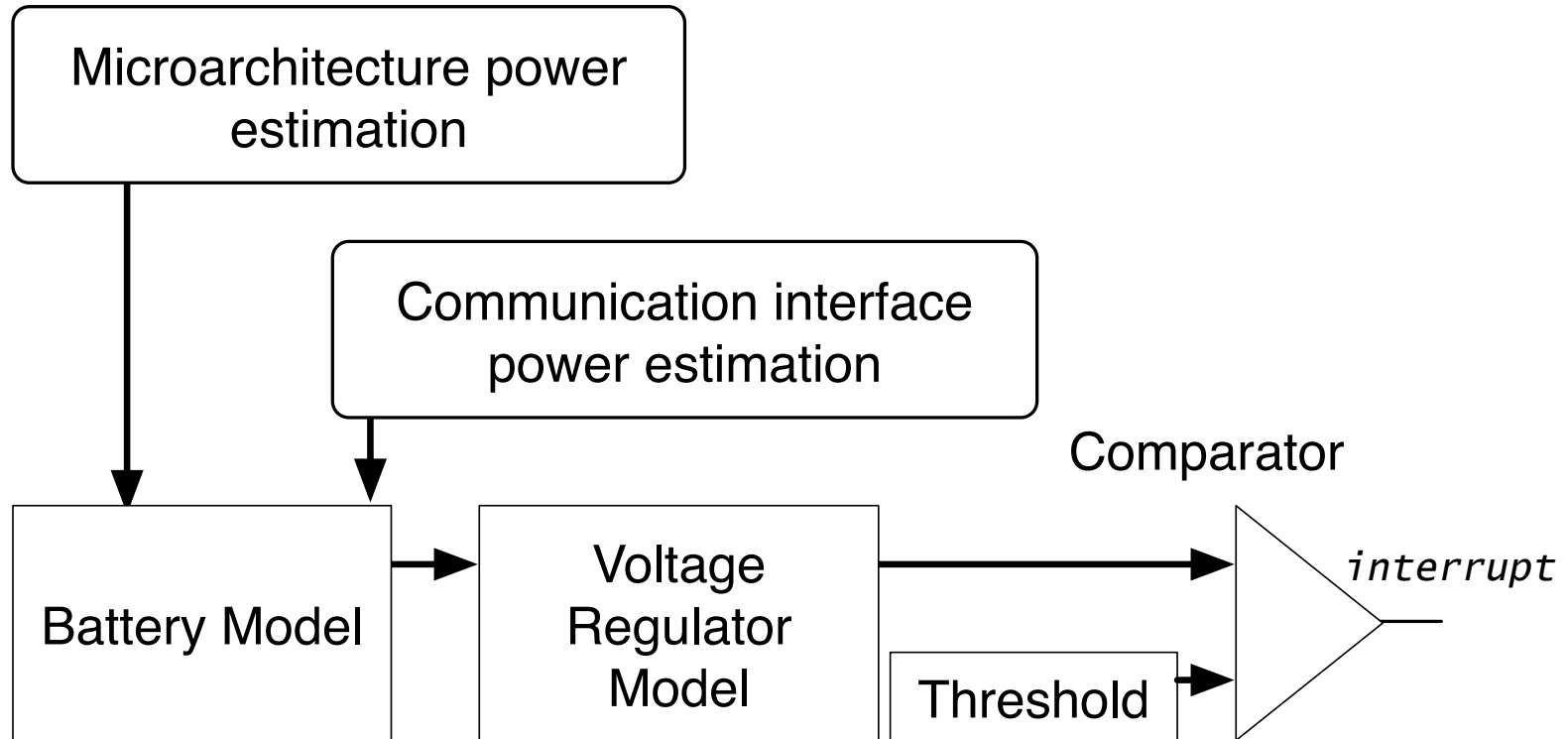# Communication Interface Modeling

Network interface

Communication medium

MAC-layer collision retry algorithm

PHY-layer signal propagation model

Network Interface Failure Model

Transmit power consumption

Receive power consumption

Idle power consumption

**TX FIFO**

size configurable in simulation

**RX FIFO**

size configurable in simulation

**Memory-mapped interface registers:**

| TX Data Register |
| TX Status Register |
| ... |
| |
| Collision Count Register |

# Creating Arbitrary Topologies



System (modeled processor core, etc.)

Network interface

Network medium

# Rich Set of Stochastic Distribution Generators

Random Variables from some common distributions

| 64-bit pseudo-random number generator | → Uniform RV on $[0, 2^{64} - 1]$ → | Inverse Transform Method / Accept / Reject Method |

| Gaussian | χ2, χ , β, F | Log Normal |
| Pareto | Beta Prime | Student t, z |
| Weibull | Erlang | Maxwell |
| Exponential | Fermi-Dirac | Logistic |
| Cauchy | Fisher-z | Log Series |
| Gumbel | Extremal Value | Rayleigh |
| Gamma | Negative Binomial | Gibrat |
| Laplace | Pearson Type III | |

- SEU Modeling
- Node Failure Modeling
- Communication Failure Modeling
- Microarchitectural parameters

- All built upon a 64-bit pseudo-random number generator with very large period [Nishimura]

Backup Slides

# Battery Subsystem Model

Backup Slides

# Distributed Simulation





**Local Area Network**

**Simulation Host 1**  **Simulation Host 2**  **Simulation Host 3**

**Central Simulation Controller**

**Simulated Nodes**

Simulated Network 1          Simulated Network 2          3

Simulated Analog Signals

# Example Memory Map



(a) Memory map of simulated processing elements.

(b) Example flow of instruction execution in the presence of interrupts from modeled peripherals

# Sunflower Simulator / HW EMulator

# Modeling Computation

- ## Modeled ISA
  - Hitachi SH ISA, support for new ISAs being added

  - Applications (e.g., SPEC CPU 2000) compiled with GCC 3.x toolchain, Hitachi/Renesas HEW compiler, Microsoft VC for Hitachi SH, ...

  - Should be able to compile apps in any source language that GCC or any of the compilers support (C, C++, FORTRAN, Java, Chill, Ada)

- ## Microarchitecture
  - Detailed simulation of SH3 family pipeline

- ## Processor peripherals and interrupt sources
  - SH3's UART, Timer Unit

  - Added new peripherals: Network Interface, Sensors, Random Number Source, Logging, Simulator Control Interface

Backup Slides

# Application's View

## Memory Map

| | |
|---|---|
| Memmory Mapped Registers | 0xFFFFFFF0 |
| | 0xFFFF0000 |
| | 0x80FFFFF |
| Stack | |
| Heap | |
| Application | |
| | 0x8003000 |
| Monitor | |
| | 0x8001000 |
| Interrupt Vector Base | |
| | 0x8000600 |
| | 0x8000000 |

- Peripherals and interrupt sources
  - Interaction with modeled peripherals and simulation control and is through memory-mapped registers

  - Interrupt sources: network, device failures, battery status, timer

- Typical application style:
  - `main()`+interrupt handler

  - E.g., received network frames handled in network intr handler, periodically scheduled tasks triggered by timer interrupts

  - Rudimentary device drivers and interrupt handling code is often reused

Computation   Communication

Power Estimation   Failures

Battery Modeling

Physical Phenomena

Backup Slides

# Modeling Physical Phenomena



- Modeling both computation and the physical phenomena that drive computation is important

- If we're modeling a sensor network, would like to model the physics of signal propagation

  - Attenuation of signals in space

  - Interference between signals

  - Location and motion of signal sources in space
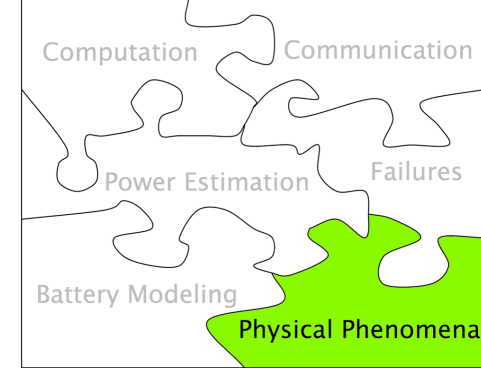
Backup Slides

# Modeling Physical Phenomena

- Example
  - 2 light sources with Gaussian spread of intensity from peak
  - The light sources move according to trajectories specified by a LUT

Example user–specified signal source attenuation

Example user–specified signal source trajectories



Signal strength

Radial distance from signal location

x

y

Source 1

Source 2

# Modeling Physical Phenomena

- ## Signal attenuation

  - Attenuation with radial distance, $x$, specified by providing coefficients for

  $$S(Ax^m + Bx^n + Cx^o + Dx^p + EK^{(Fxq + Gxr + Hxs + Ixt)})$$

  - Example, for Gaussian spread light source with peak intensity S, E=1, K= e, F = -0.5, q = 2, all other coefficients are 0

- ## Signal trajectory

  - Trajectory and speed specified as a list of way-points and sampling rate

  - Notion of time in physical models is synchronized with ISA simulation, communication modeling, battery models, etc.

Backup Slides

# Modeling Communication

- ## Model
  - Each processor can have multiple **network interfaces** (NICs) instantiated
  - Within a simulation, multiple **network segments** are instantiated
  - NICs are connected to network segments to create arbitrary topologies
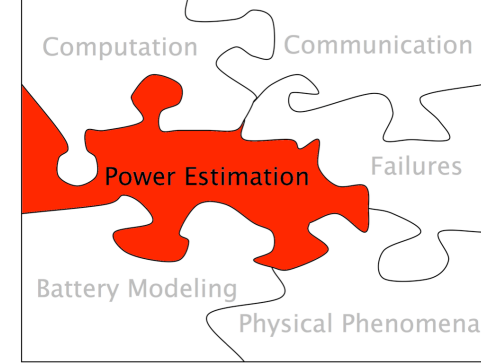


- ## Application's view
  - Code running over simulator sees NIC as a memory-mapped peripheral
  - Interrupts generated for TX/RX and various errors
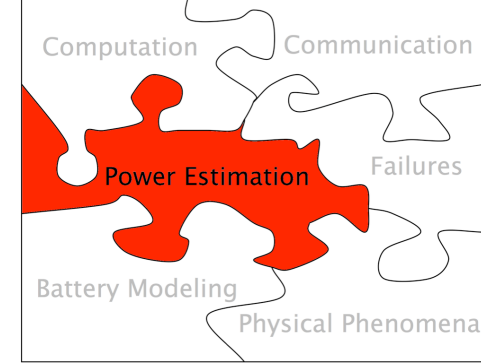
Backup Slides

# Modeling Communication

- **Network Segment** (Defines properties of physical link)
  - Bit rate

  - Frame size

  - Number of simultaneous transmissions that can be accommodated

  - Rate and distribution of intermittent failures

  - Signal propagation properties, same model as described for phy. sources

  - Minimum SNR before introducing bit errors

- **Network Interface**
  - Transmit, receive and idle power consumption
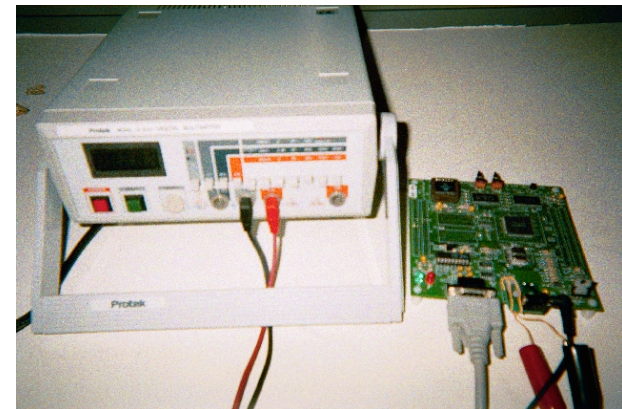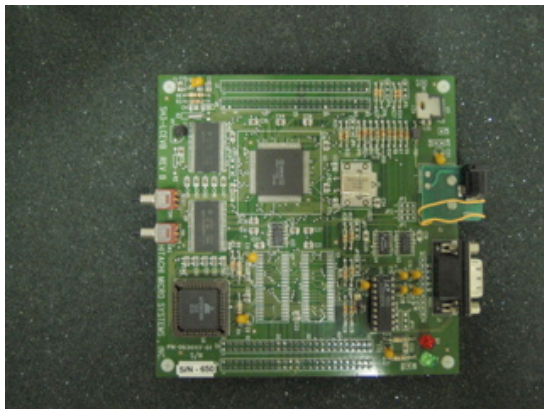  - Number of TX and RX FIFO entries

Backup Slides

# Power Estimation



- 3 different methods for power estimation

- Requires most effort, slower, not always practicable:
  - Signal transition activity reported per cycle for pipeline latches, buses, register file, cache ports, usage of FUs
  - Requires capacitance values to obtain power estimates

- Better tradeoff:
  - Characterized instruction-level power model
  - Requires empirical measurements on actual hardware

- Simplest:
  - Specify average active and sleep power consumption, e.g., from a device manufacturer's data sheet
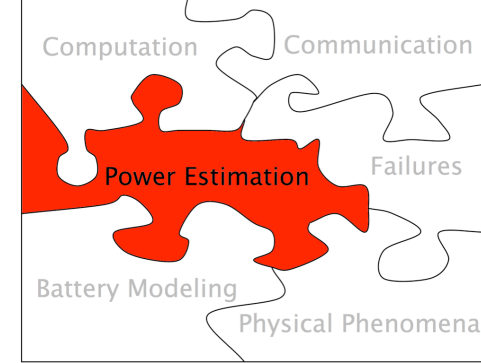
Backup Slides

# Power Estimation

- Simulator includes an instruction-level power model for the Hitachi SH7708

- Empirical measurements using SH7708 eval. board
  - Measured average current draw for a fixed set of operand values, for each instruction in the ISA (~160 instructions / different addressing modes)

  - No inter-instruction effects

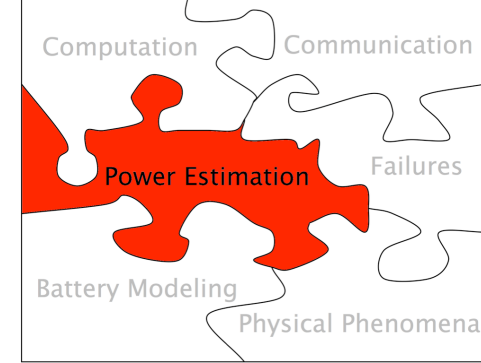  - Non-idempotent instructions like TRAP can't be subjected to this method
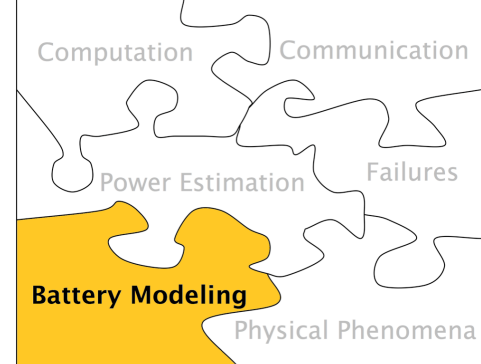




Backup Slides

# Power Estimation

- ## Signal transition activity estimation
  - Signal transitions as encoded/decoded instructions move through pipeline

  - Register file, program counter

  - Data and address buses

  - Cache read/write ports

  - No modeling of internals of functional units

- ## When is transition activity estimation useful ?
  - Comparative studies of signal transition activity in the modeled structures

  - If you have capacitance values; this may not be available until after floor-planning and layout
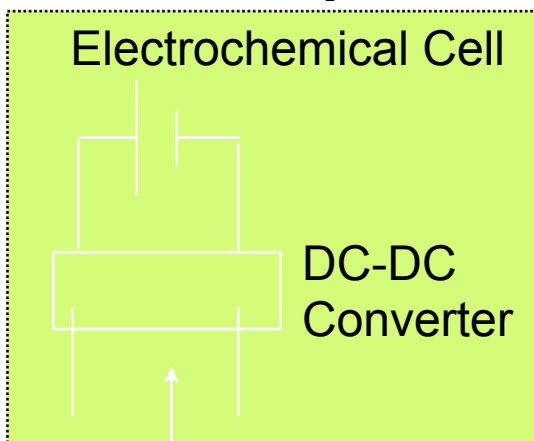
Backup Slides

# Power Estimation

- Other details

- Sleep mode
  - Reduced power consumption when application executes a sleep instruction (until next interrupt)

- Voltage and frequency scaling during simulation
  - Under simulated application's control via simulator control register

  - Interactively via command prompt

  - Independent voltage and frequency scaling or scale VDD and freq. in tandem for a specified Vt and technology-dependent α
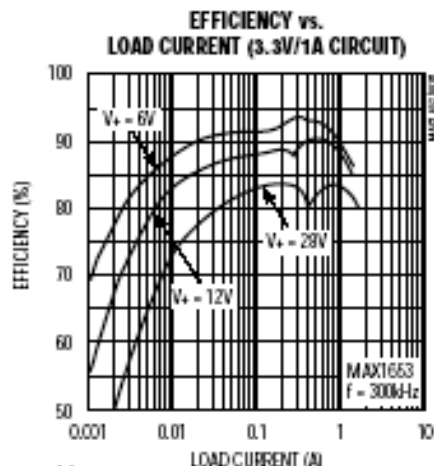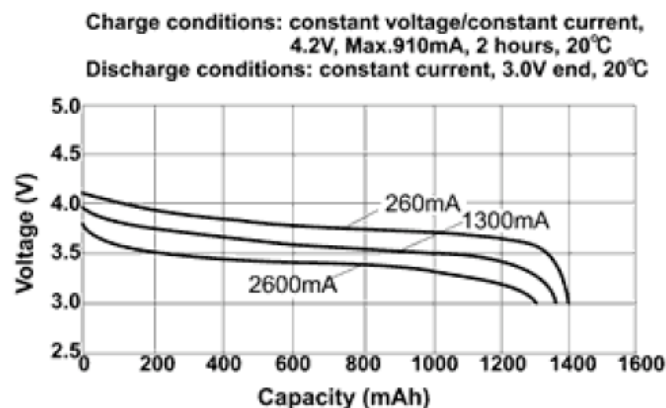
Backup Slides

# Battery Subsystem

**"Battery"**

Electrochemical Cell

DC-DC Converter

Network interface power consumption + Processor power consumption

DC-DC Converter LUT from data sheet



Efficiency vs. current for Maxim MAX1652
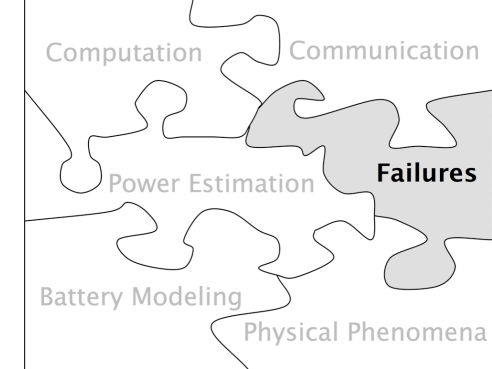
Electrochemical cell discharge profile



- One or more devices can be attached to a battery system

- Sampled current draw periodically supplied to battery subsystem

- Current passed to a model for a DC-DC converter, then to model for electrochemical cell

- Models based on discrete-time battery model from [L. Benini et al., TVLSI '01]
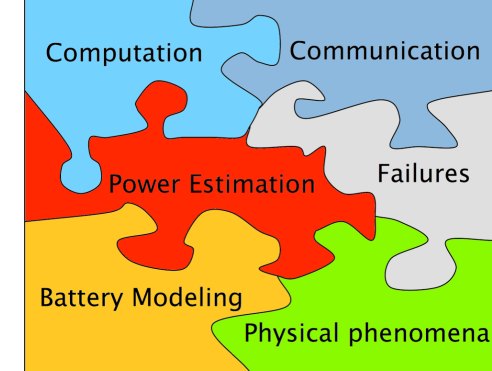
Backup Slides

# Modeling Failures

- ## Computation Failures

  - Event upsets within microarchitecture or treat each processor as a unit that may fail *en-masse* with some probability

  - Intermittent failures with a specified rate and distribution for which processor is temporarily inactive

  - Correlated failures between processing elements and network segments

- ## Communication failures

  - Also permit configuration of intermittent failures of network segments

  - If a network is associated with a physical signal model, induce bit errors when SNR drops below a specified threshold

Backup Slides

# Implementation



- Simulator core written in C, interactive command parser and built-in assembler specified in Yacc

- Runs on *BSD, Linux, MacOS, Windows

- Console application or with GUI

- GUI and simulation parallelization implemented with the Inferno OS and Limbo programming language

Backup Slides