# High-Performance Scientific Simulation and Distributed Computing at ORNL:
## **Harness**, **CUMULVS** and the
## Common Component Architecture (**CCA**)

Dr. James Arthur Kohl

Oak Ridge National Laboratory

Wednesday, September 13, 2000

# Scientific Simulation is HOT!

- Flexible, Powerful, Inexpensive (?)
- Remote Collaboration Possible
- Infrastructure Issues:
  - $\Rightarrow$ Interaction / Control of Parallel Simulations
  - $\Rightarrow$ Fault Tolerance ~ Clusters, Networks, etc.
  - $\Rightarrow$ Adaptable to Changing Resources / Technology
- ORNL Projects:
  - $\Rightarrow$ **Harness**, **CUMULVS**, **CCA**

# Distributed Computing History at ORNL

> PVM was developed as a tool to help us explore HDC issues.
> Software has been redesigned from scratch every 2-3 yrs
> to study emerging architecture, network, and user needs.

**1990 Heterogeneity -** architecture, data, power, network...

explore the problems and research issues in heterogeneous computing

**1991 Robust, Portable Programming Environment**

study how to create a robust, portable programming environment for HDC

**1993 Transparent Multiprocessor integration**

efficient multi-protocol handling of distributed and shared memory computers

**1995 Fault Tolerance and Extensibility**

study VM and application fault tolerance. Design first VM plug-in interfaces

**1998 Windows and Unix interoperability**

enable the millions of NT and Win2000 hosts to exploit cluster computing

**Next Step is Harness**

# HARNESS

**Exploring New Capabilities in Heterogeneous Distributed Computing**

Building on our experience and success with PVM we will create a fundamentally new heterogeneous virtual machine environment based on three key research concepts:

- **Parallel Plug-in Environment**
  Extend the concept of a plug-in to the parallel computing world. Dynamic with no restrictions on functions.

- **Distributed Peer-to-Peer Control**
  No single point of failure unlike typical client/server models.

- **Multiple Distributed Virtual Machines Merge/Split**
  Provide a means for short-term sharing of resources and collaboration between teams.

# HARNESS Motivation
**Needs of Simulation Science and Cluster Computing**

• develop applications by plugging together component models.

• customize/tune virtual environment for application's needs and for performance on existing resources.

• support long-running simulations despite maintenance, faults, and migration (dynamically evolving VM).

• adapt virtual machine to faults and dynamic scheduling in large clusters (C-Plant).

• Provide framework for collaborative simulations (in spirit of CUMULVS or a collaborative PSE).

# HARNESS Team

## Collaborative Effort by the Developers of PVM

Al Geist,
Jim Kohl, Stephen Scott, Conrad Albrecht-Buehler, Wael Elwasif
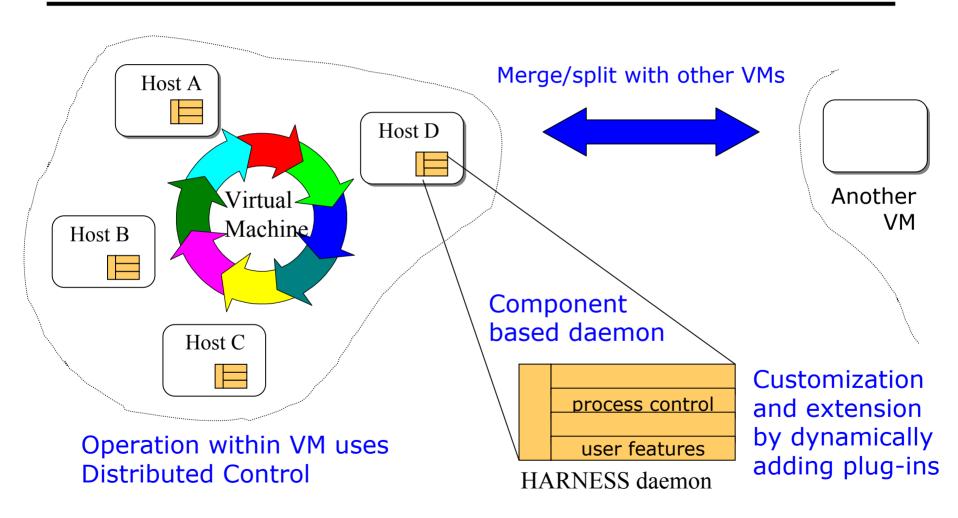Oak Ridge National Laboratory

Jack Dongarra,
Graham Fagg, Martin Swany, Nathan Garner
University of Tennessee

Vaidy Sunderam,
Paul Gray, Mauro Migliardi, Gopi Sankar
Emory University

www.csm.ornl.gov/harness
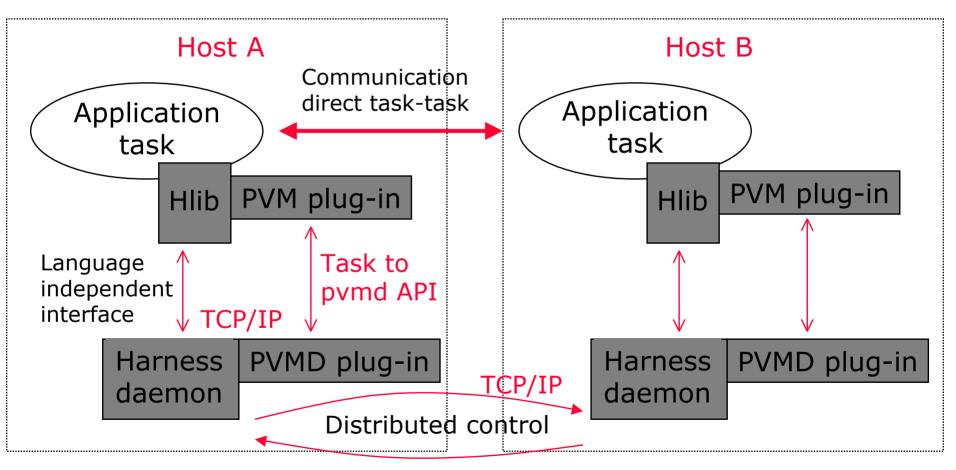
# HARNESS Virtual Machine
## Scalable Distributed Control and CCA-Based Daemon

Host A

Merge/split with other VMs

Host D

Virtual Machine

Another VM

Host B

Host C

Operation within VM uses Distributed Control

Component based daemon

process control

user features

HARNESS daemon

Customization and extension by dynamically adding plug-ins

# HARNESS Architecture

## Hlib and Harness daemon within host

### Example: PVM Application on Harness



Host A

Host B

Application task

Communication direct task-task

Application task

Hlib | PVM plug-in

Hlib | PVM plug-in

Language independent interface

Task to pvmd API

TCP/IP

Harness daemon | PVMD plug-in

Harness daemon | PVMD plug-in

TCP/IP

Distributed control

# HARNESS Development Plan

Plan is to produce the core framework and a couple key plug-ins that will provide a practical environment and illustrate the capabilities of HARNESS:

- **Harness Core**
  Task library and Harness Daemon software
  Provides API to load, unload plug-ins and distributed control.
- **PVM Plug-in**
  Provides PVM API veneer to support exiting PVM applications.
- **Fault tolerant MPI plug-in**
  Provides MPI API for 30 most used functions. Semantics adjusted to allow recovery from corrupted communicator.
- **VIA/FM communication plug-in**
  To illustrate how different low level communication plug-ins can be used within Harness, and to provide high performance.

# Harness Core Implementations
## Two Different Schemes are Being Explored

---

- ## C Scheme (ORNL)
  - ⇒ Based on dynamically linked / shared libraries
  - ⇒ Advantage ~ DLL / lib can be written in the language of the app
  - ⇒ Potential for higher performance plug-ins (compiled binary)

- ## Java (2) Scheme (Emory)
  - ⇒ Based on JVM dynamic class loader
  - ⇒ Advantage ~ fast prototyping
  - ⇒ Leverages wealth of existing Java specs, JINI, JavaBeans, RMI, Java Spaces, etc.
  - ⇒ Good integration with emerging Java SC apps and interfaces
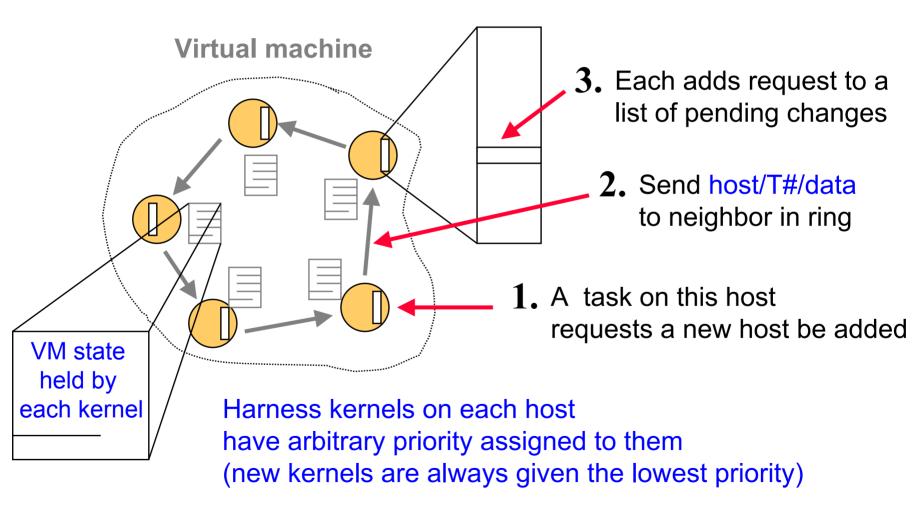  - ⇒ Beta Release!

# Symmetric Peer-to-Peer Distributed Control
## Requirements

- No single point of failure for Harness.
  - $\Rightarrow$ It survives as long as one member still lives.

- All members know the state of the virtual machine
  - $\Rightarrow$ Knowledge is kept consistent w.r.t. the order of changes of state (Important parallel programming requirement!)

- No member is more important than any other
  - $\Rightarrow$ At any instant
  - $\Rightarrow$ There's no "control token" being passed around

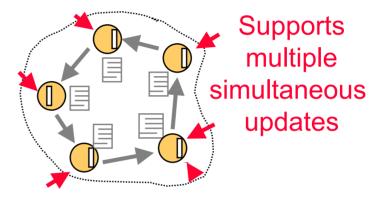# Distributed Control
## Harness Overlapping Two Phase Arbitration

**Virtual machine**

**3.** Each adds request to a list of pending changes

**2.** Send host/T#/data to neighbor in ring

**1.** A task on this host requests a new host be added

VM state held by each kernel

Harness kernels on each host have arbitrary priority assigned to them (new kernels are always given the lowest priority)

# Harness Distributed Control
## Control is Asynchronous and Parallel

Supports
fast host
adding

add
host

Supports
multiple
simultaneous
updates

Fast
host delete
or recovery
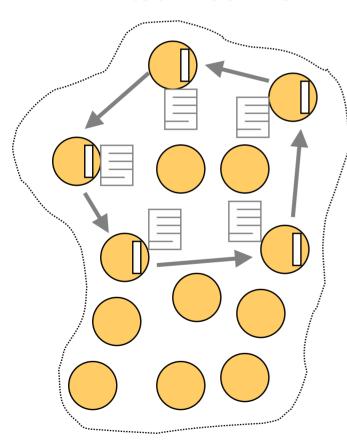from fault

Parallel recovery
from multiple
host failures

# HARNESS Scalability
## Variable Distributed Control Loop Size

**Virtual machine**



**Size of the Control Loop
1 <= S <= (size of VM)**

For small VM and ultimate fault tolerance S = (size of VM)

For large VM a random selection of a few  hosts (f.e. S = 10) gives a balance of multi-point  failure and performance.

For  S = 1, distributed control becomes simple client/server model.

# Fault Tolerant MPI

**Motivation**

---

**Two major drawbacks to MPI** are:
- lack of interoperability - being addressed by IMPI and MPI-Connect
- lack of fault tolerance - any failure is catastrophic

Being used
on Blue at LLNL

As application and machine sizes grow the
**MTBF is less than the application run time**.

MPI standard is based on a static model any decrease in tasks
leads to **corrupted communicator** (MPI_COMM_WORLD).

Develop MPI plugin that **takes advantage of Harness
robustness**  to allow a range of recovery alternatives
to an MPI application. Not just another MPI implementation.

FT-MPI follows the syntax of MPI standard

# HARNESS Research Status

Java-based HARNESS prototype created at Emory
- used to test parallel plug-in concepts
- integration with JINI underway

IceT package developed by Paul Gray (Iowa St.)
- demonstrates merging and splitting of virtual machines
- dynamically switching communication (MPI to CCTL)

C-based HARNESS kernel and distributed control
- feasibility demonstrated at ORNL
- production release in progress

C-based FT-MPI plug-in prototype developed at UTK
- built on top of PVM 3.4 API

# Harness Enables New Kinds of Applications

## Thinking "Outside the Box" -- It's not just for Scientific Computing

Harness is still just a research project but its potential is great



Applications follow user roaming wearable computers



On-the-fly simulation tuning plug-in different methods if simulation evolves to need them



Teams of  tasks patrol and monitor local network for performance or security
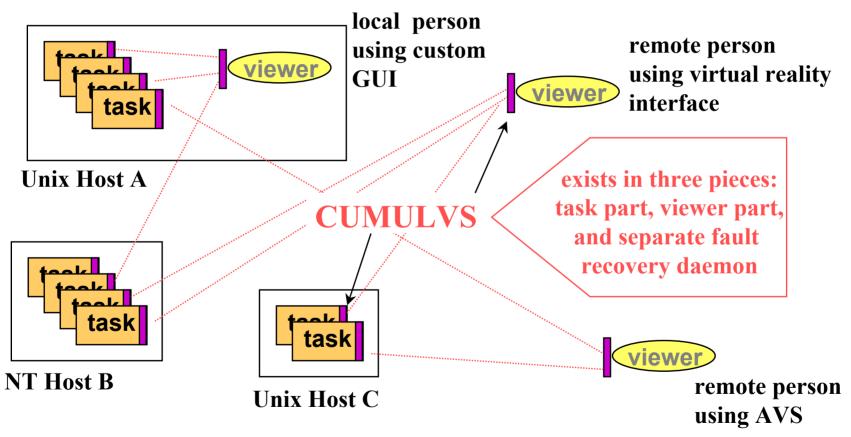
(<u>C</u>ollaborative,  <u>U</u>ser <u>M</u>igration,  <u>U</u>ser <u>L</u>ibrary  for  <u>V</u>isualization  and  <u>S</u>teering)

- Collaborative Infrastructure for Interacting with Scientific Simulations:

  $\Rightarrow$ Run-Time Visualization by Multiple Viewers

  $\rightarrow$ Dynamic Attachment

  $\Rightarrow$ Coordinated Computational Steering

  $\rightarrow$ Model  &  Algorithm

  $\Rightarrow$ Heterogeneous Fault Tolerance

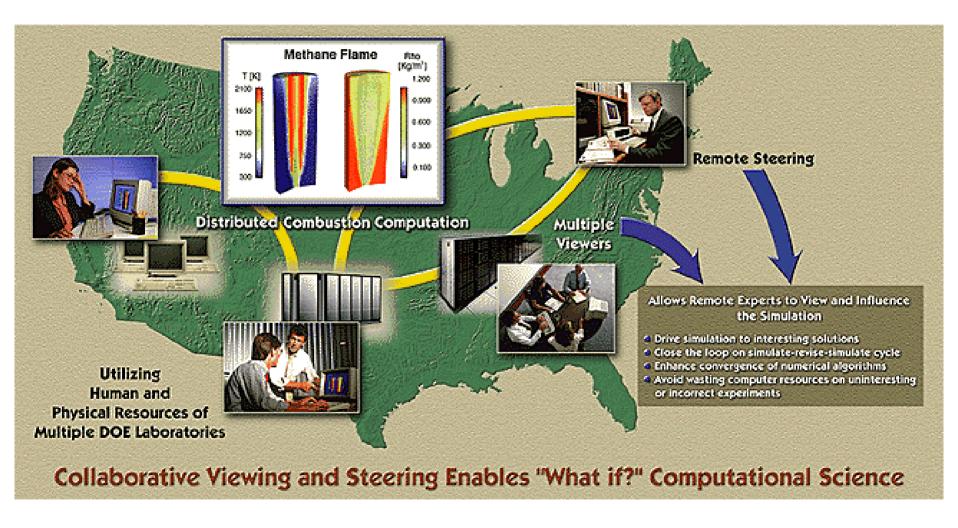  $\rightarrow$ Automatic Fault Recovery and Task Migration

  $\Rightarrow$ Coupled Models…

# CUMULVS

## coordinates the consistent collection and dissemination of information to/from parallel tasks to multiple viewers

local person using custom GUI

remote person using virtual reality interface

**task**
**task**
**task**
**task**

**viewer**

**viewer**

**Unix Host A**

**CUMULVS**

exists in three pieces: task part, viewer part, and separate fault recovery daemon

**task**
**task**
**task**
**task**

**task**
**task**

**viewer**

**NT Host B**

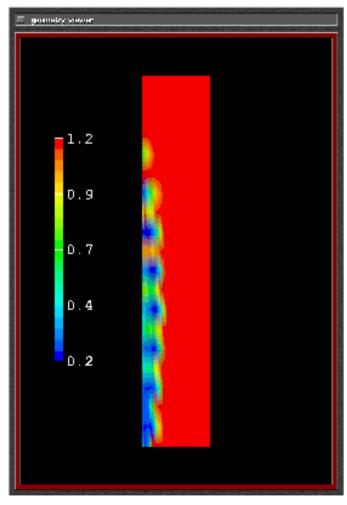**Unix Host C**

remote person using AVS

distributed parallel application or simulation
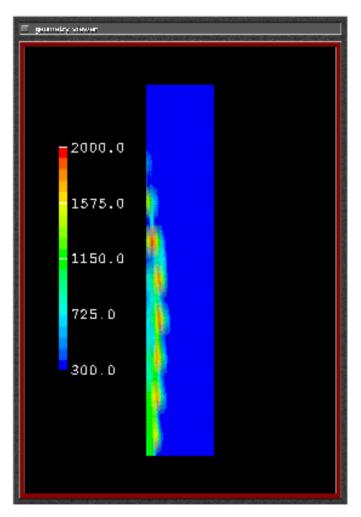supports most target platforms (PVM/MPI, Unix/NT, etc.)

Kohl/2000-19

# Collaborative Combustion Simulation
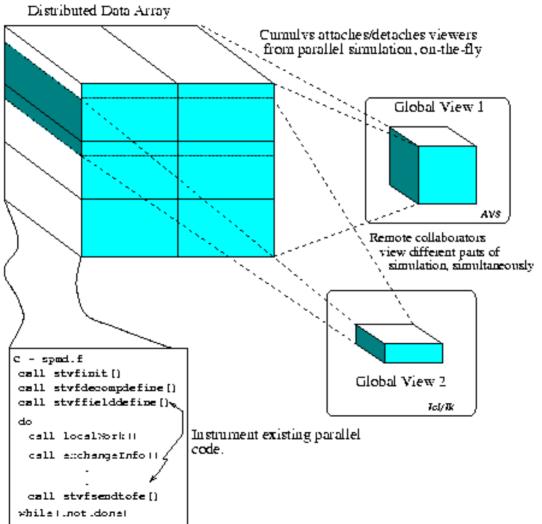
# Multiple Simultaneous Views



Density



Temperature

# Multiple Distinct Views



Distributed Data Array

Cumulvs attaches/detaches viewers from parallel simulation, on-the-fly

Global View 1

AVS

Remote collaborators view different parts of simulation, simultaneously

Global View 2

Tcl/Tk

```
c - spmd.f
 call stvfinit[]
 call stvfdecompdefine[]
 call stvffielddefine[]
 do
   call localWork[]
   call exchangeInfo[]
     .
     .
   call stvfsendtofe[]
 while[.not.done]
```
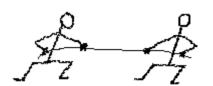
Instrument existing parallel code.

# CUMULVS Particle Handling

- **Particle Data Fundamentally Different**
  $\Rightarrow$ Nested Data Fields, Explicit Coordinates

- **Particle-Based Decomposition API**
  $\Rightarrow$ User-Defined, Vectored Accessor Routines

- **Viewing Particle Data**
  $\Rightarrow$ AVS Module Extensions
  $\Rightarrow$ Tcl/Tk Slicer Particle Mode

# Coordinated Steering

- Multiple, Remote Collaborators

- Simultaneously Steer Different Parameters

    $\Rightarrow$ Physical Parameters of Simulation

    $\Rightarrow$ Algorithmic Parameters ~ e.g. Convergence Rate

- "What If ?" Analyses

    $\Rightarrow$ Explore Non-Physical Effects

- Efficient Experimentation Cycle

    $\Rightarrow$ Keep Simulation On Track

    $\Rightarrow$ Crop Off  Experiments Gone Awry…

# Heterogeneous Checkpointing

- ## Application Defines Its Own State
  $\Rightarrow$ Provide CUMULVS with Semantic Information

  $\Rightarrow$ CUMULVS Handles Checkpoint Collection

  $\Rightarrow$ Automatic Fault Recovery System

- ## Efficient & Flexible Checkpoints
  $\Rightarrow$ Not Full Core Image

  $\Rightarrow$ Semantics $\rightarrow$ Heterogeneous Restart & Migration

  $\Rightarrow$ On-The-Fly Reconfiguration Also Possible…

# Run-Time System Architecture

- One Checkpointing Daemon (CPD) Per Host
  - $\Rightarrow$ Ckpt Collector / Provider
  - $\Rightarrow$ Run-Time Monitor
  - $\Rightarrow$ Console for Restart / Migrate
- CPDs Comprise Fault-Tolerant Application…
  - $\Rightarrow$ Handle Failure of Host / CPD
  - $\Rightarrow$ Coordinate Redundancy
  - $\Rightarrow$ Ring Topology



Virtual Machine

Physical hosts

App App App

App App

CPD CPD

Migrate

App CPD

App

New CPD

Replacement host added on failure

Spare Host

# Manual Software Instrumentation

- SPDT 98 Case Study ~ SW Instrumentation Cost

| Instrumentation: | Seismic: | Wing Flow: |
|---|---|---|
| Original Lines of Code | 20,632 | 2,250 |
| Vis / Steer System Init | 3 | 3 |
| Vis / Steer Variable Decls | 48 | 73 |
| CP Restart Initialization | 21 | 12 |
| CP Rollback Handling | 41 | 34 |
| Total Instrumentation | 204 ~ 1.0 % | 188 ~ 7.7 % |

# Checkpointing Efficiency

- SPDT 98 Case Study ~ Execution Overhead

Seconds per Iteration

| Experiment: | SGI: | Cluster: | Hetero: |
|---|---|---|---|
| Seismic - No Checkpointing | 2.83 | 6.23 | 9.46 |
| Seismic - Checkpoint for Restart | 2.99 | 6.50 | 10.76 |
| Seismic - Checkpoint for Rollback | 3.03 | 6.66 | 10.90 |
| Wing - No Checkpointing | 0.69 | 1.58 | 6.14 |
| Wing - Checkpoint for Restart | 0.77 | 1.71 | 7.10 |
| Wing - Checkpoint for Rollback | 0.79 | 1.71 | 7.30 |

**(Checkpointing Every 20 Iters.)**

**Seismic Overhead: 4-14% Restart, +1-3% Rollback.**
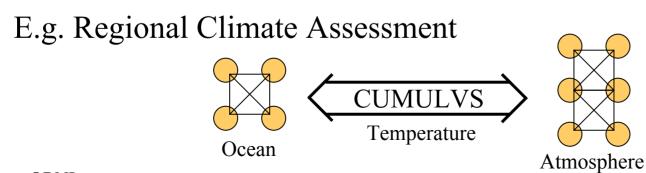**Wing Overhead: 8-15% Restart, +0-2.5% Rollback.**

# Seismic Example ~ 3D (AVS)

# Air Flow Over Wing Example ~ 3D (AVS)

# Coupling Data Fields in Simulation Models Using CUMULVS

- ## Natural Extension to Viewer Scenario

    $\Rightarrow$ Promote "Many-to-1" $\rightarrow$ "Many-to-Many"



- ## Translate Disparate Data Decompositions

    $\Rightarrow$ Complements PAWS Coupling Work

    $\Rightarrow$ Builds on CCA (<u>C</u>ommon <u>C</u>omponent <u>A</u>rchitecture) Forum

E.g. Regional Climate Assessment



Ocean

CUMULVS

Temperature

Atmosphere

# Future CUMULVS Plans

- Application Interface:

  $\Rightarrow$ Assist Manual Instrumentation of Applications

  $\rightarrow$ GUI, Pre-Compiler…

- Checkpointing Efficiency:

  $\Rightarrow$ Tasks Write Data in Parallel / Parallel File System

  $\Rightarrow$ Variable Redundancy Levels, Improve Scalability

- Portability:

  $\Rightarrow$ Other Messaging Substrates

  $\rightarrow$ Reduced Functionality for MPI / MPI-2…

http://www.csm.ornl.gov/cs/cumulvs.html

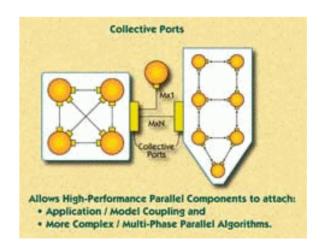# Common Component Architecture (CCA)

- Component Architecture for Scientific Simulation

  ⇒ Special Emphasis on High-Performance / Parallelism

- Reusable "Components" Connect Via "Ports"
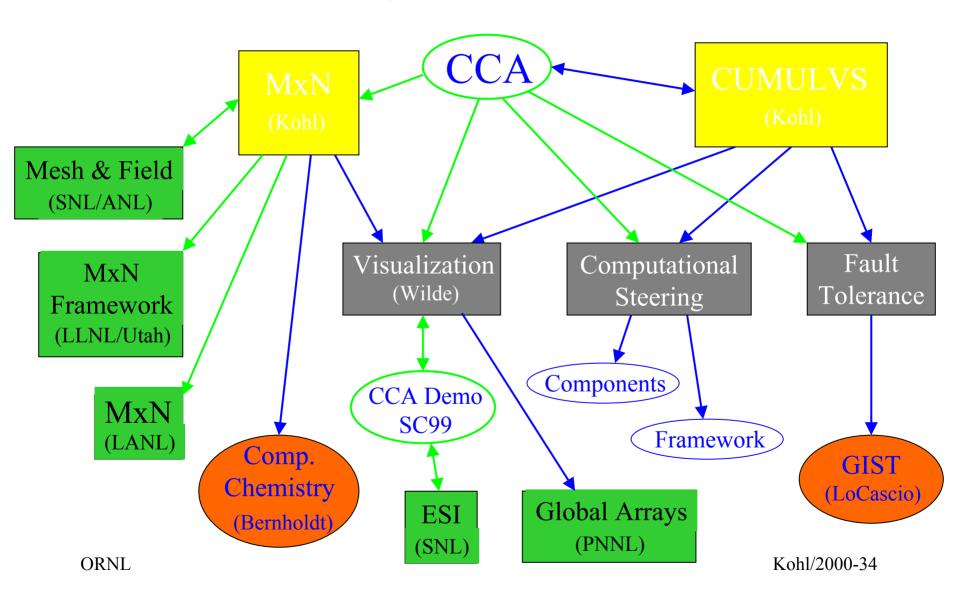
  ⇒ Forum Creating Specification and Reference Framework

---

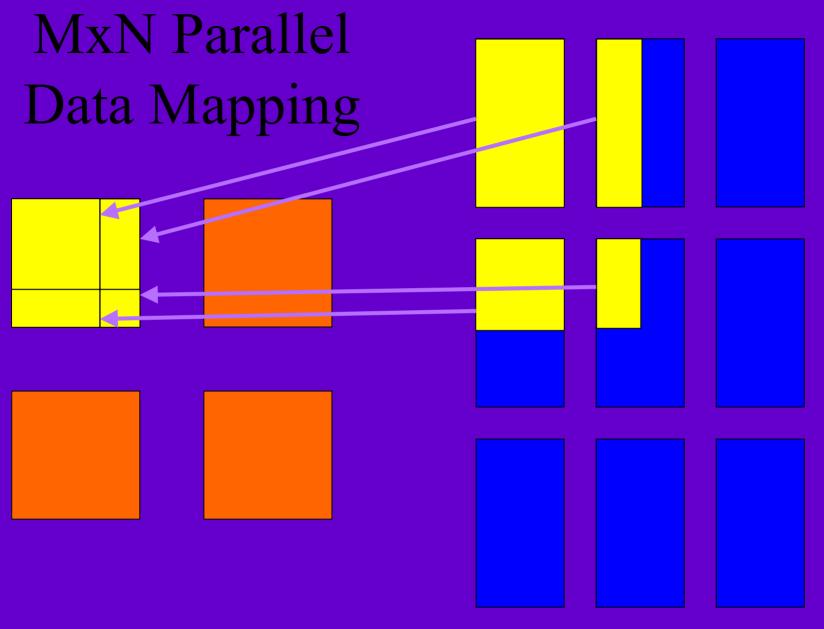## Direct Connect Ports
(Local Components Share Memory)

## Collective / MxN
(Parallel Data Exchange)



Direct Connect Ports
Same Address Space
Allows for direct, low overhead interaction among components in same address space.



Collective Ports
Mx1
MxN
Collective Ports
Allows High-Performance Parallel Components to attach:
- Application / Model Coupling and
- More Complex / Multi-Phase Parallel Algorithms.

# Common Component Architecture (CCA) / ACTS Toolkit
## Oak Ridge National Laboratory

# MxN Parallel
# Data Mapping

# Collective MxN Example

Ocean

Atmosphere

getDataField( space, time )

Collective Port

# High-Performance Parallel Collective Port

Ocean

Atmosphere

getDataField( space, time )

getDataField( space, time )

getDataField( space, time )

getDataField(  space, time )

Kohl/2000-37

# Summary of ORNL Scientific Simulation

- Harness ~ Next Generation HDC Environment
  $\Rightarrow$ Pluggable Virtual Machine, Distributed Control

- CUMULVS ~ Interacting with Simulations On-The-Fly
  $\Rightarrow$ Visualization, Steering, Fault Tolerance

- Common Component Architecture (CCA)
  $\Rightarrow$ Harness Pluggability Builds on CCA Foundation
  $\Rightarrow$ CUMULVS Technology Used for MxN / Coupling

CUMULVS

Harness

CCA