# Concurrent Models of Computation

### Edward A. Lee

Robert S. Pepper Distinguished Professor, UC Berkeley
EECS 219D
*Concurrent Models of Computation*
Fall 2011

Week 4: Message Passing Patterns

---

## Message Passing Interface
## MPI

MPI is a collaborative standard developed since the early 1990s with many parallel computer vendors and stakeholders involved.

Realized as a C and Fortran APIs.

First draft of MPI: J. J. Dongarra, R. Hempel, A. J. G. Hey, and D. W. Walker. A proposal for a user-level, message passing interface in a distributed memory environment. Technical Report TM-12231, Oak Ridge National Laboratory, February 1993.

## Anatomy of an MPI Program (in C)

```
/* On each processor, execute the following with different values for rank. */
int main(int argc, char *argv[]) {
    int rank, size, …;

    MPI_Init(&argc, &argv);
    // Find out which process this is (rank)
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    // Find out how many processes there are (size)
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    if (rank == 0) {
        … code for one process …
    } else if (rank == RAMP2) {
        … code for another process …
    }
    MPI_Finalize();
    return 0;
}
```
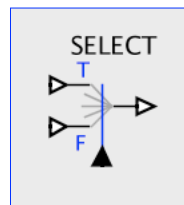
## MPI Implementation of Select Process

```
int control;
while (1) {
    MPI_Recv(&control, 1, MPI_INT, CONTROL_SOURCE, ...);
    if (control) {
        MPI_Recv(&selected, 1, MPI_INT, DATA_SOURCE1, ...);
    } else {
        MPI_Recv(&selected, 1, MPI_INT, DATA_SOURCE2, ...);
    }
    MPI_Send(&selected, 1, MPI_INT, DATA_SINK, ...);
}
```

Rank of the source or destination process

Data type of the handled message

SELECT

●2

## Vague MPI Send Semantics

MPI_Send is a "blocking send," which means that it does not return until the memory storing the value to be sent can be safely overwritten. The MPI standard allows implementations to either copy the data into a "system buffer" for later delivery to the receiver, or to rendezvous with the receiving process and return only after the receiver has begun receiving the data.
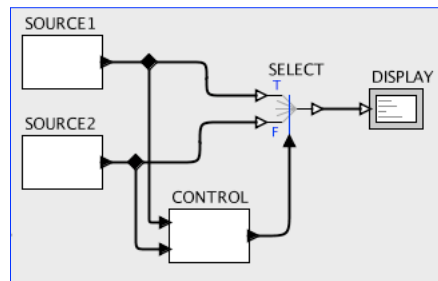
*Discussion: What do you think of this?*

You can force a rendezvous style by using MPI_Ssend instead of MPI_Send

---

## What happens to this program under a rendezvous style of communication?



**CONTROL Process:**

```
MPI_Recv(&data1, 1, MPI_INT, SOURCE1, ...);
MPI_Recv(&data2, 1, MPI_INT, SOURCE2, ...);
while (1) {
  if (someCondition(data1, data2)) {
    MPI_Send(&trueValue, 1, MPI_INT, SELECT, ...);
    MPI_Recv(&data1, 1, MPI_INT, SOURCE1, ...);
  } else {
    MPI_Send(&falseValue, 1, MPI_INT, SELECT, ...);
    MPI_Recv(&data2, 1, MPI_INT, SOURCE2, ...);
  }
}
```
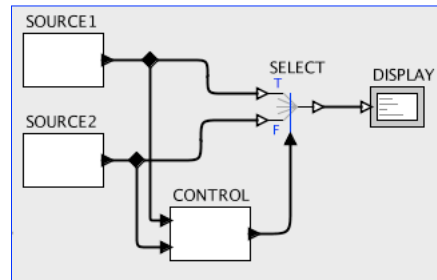
A Design Question:
How to accomplish the fork processes?

Option 1:
   Create a process for each
   fork that copies inputs to
   outputs (in what order?)



Option 2:
   Modify the SOURCE
   processes to do
   successive writes to
   SELECT and CONTROL
   (in what order?).

---

Forcing Buffered Send: MPI_Bsend()

"A buffered send operation that cannot complete
because of a lack of buffer space is erroneous. When
such a situation is detected, an error is signalled that may
cause the program to terminate abnormally. On the other
hand, a standard send operation that cannot complete
because of lack of buffer space will merely block, waiting
for buffer space to become available or for a matching
receive to be posted. This behavior is preferable in many
situations."

Message Passing Interface Forum (2008). MPI: A Message Passing Interface
standard -- Version 2.1, University of Tennessee, Knoxville, Tennessee.

## Irony

"The reluctance of MPI to mandate whether standard sends are buffering or not stems from the desire to achieve portable programs."

Message Passing Interface Forum (2008). MPI: A Message Passing Interface standard -- Version 2.1, University of Tennessee, Knoxville, Tennessee.

Lee 04: 9

## Buffer Size Control in MPI

MPI_Buffer_attach associates a buffer with a process. Any output can use the buffer, and MPI does not limit the buffering to the specified buffers.

The MPI_Send procedure can return an error, so you can write processes that do something when buffers overflow. What should they do?

MPI provides few mechanisms to exercise control over the process scheduling (barrier synchronization seems to be about it).

Lee 04: 10

## MPI_Recv Semantics

MPI_Recv blocks until the message is received.

Communication is point-to-point: Sending and receiving processes refer to each other. According to the MPI standard: "[this] guarantees that message-passing code is deterministic, if processes are single-threaded and the wildcard MPI_ANY_SOURCE is not used in receives."
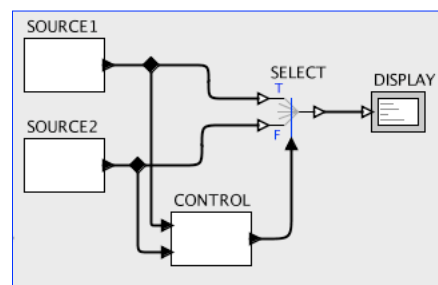
MPI_ANY_SOURCE can be specified in a MPI_Recv()

Messages arrive in the same order sent.

---

## Discussion: Suppose you wanted to implement Parks' algorithm or Geilen and Basten?

How would you do it?
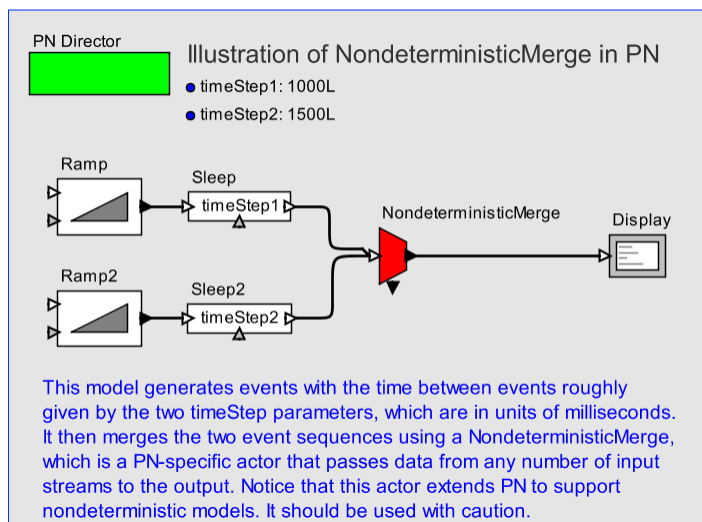
●6

# Threads and Fairness

MPI is used sometimes with threads, where a single process runs in multiple threads. This can

"Fairness MPI makes no guarantee of fairness in the handling of communication. Suppose that a send is posted. Then it is possible that the destination process repeatedly posts a receive that matches this send, yet the message is never received, because it is each time overtaken by another message, sent from another source. Similarly, suppose that a receive was posted by a multi-threaded process. Then it is possible that messages that match this receive are repeatedly received, yet the receive is never satisfied, because it is overtaken by other receives posted at this node (by other executing threads). It is the programmer's responsibility to prevent starvation in such situations."

---

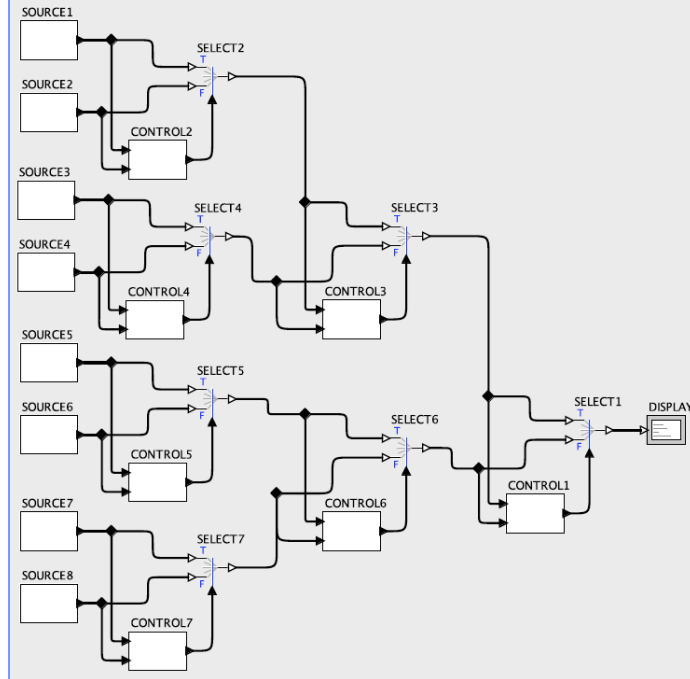# NondeterministicMerge in Ptolemy II is implemented in a multithreaded actor

Two threads perform blocking reads on each of two input channels and write to the same output port.



PN Director

Illustration of NondeterministicMerge in PN
- timeStep1: 1000L
- timeStep2: 1500L

Ramp
Sleep
timeStep1

Ramp2
Sleep2
timeStep2

NondeterministicMerge

Display

This model generates events with the time between events roughly given by the two timeStep parameters, which are in units of milliseconds. It then merges the two event sequences using a NondeterministicMerge, which is a PN-specific actor that passes data from any number of input streams to the output. Notice that this actor extends PN to support nondeterministic models. It should be used with caution.

## Scaling Up Designs

Collective operations enable compact representations of certain composite structures (not this one though, at least not in MPI).
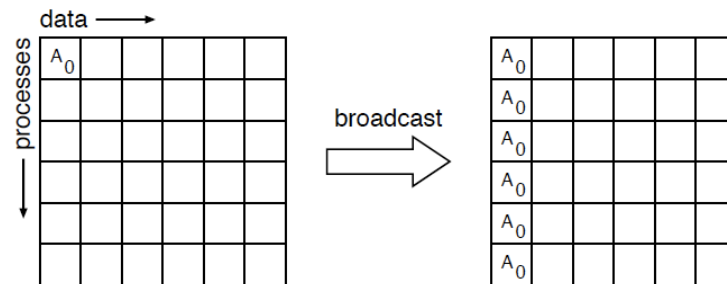
---

## Collective Operations Provided by MPI

- Barrier synchronization in a group
- Broadcast to a group
- Gather from a group (to one member or all members)
- Scatter to a group
- Scatter/Gather all-to-all
- Reduction operations such as sum, max, min, or user-defined functions, where the result is returned to all group members or one member
- Combined reduction and scatter operation
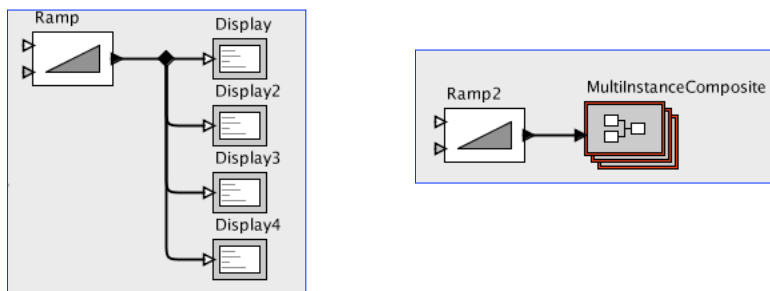- Scan across all members of a group (also called prefix)

●8

# Broadcast

MPI_Bcast is like MPI_Send except that it sends to all members of the group. Data are copied.
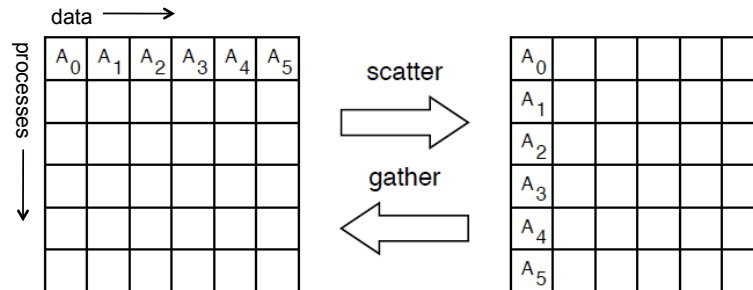
Lee 04: 17

---

# Broadcast in Ptolemy II

Choice of director defines the communication policy.



Lee 04: 18

9

## Gather/Scatter

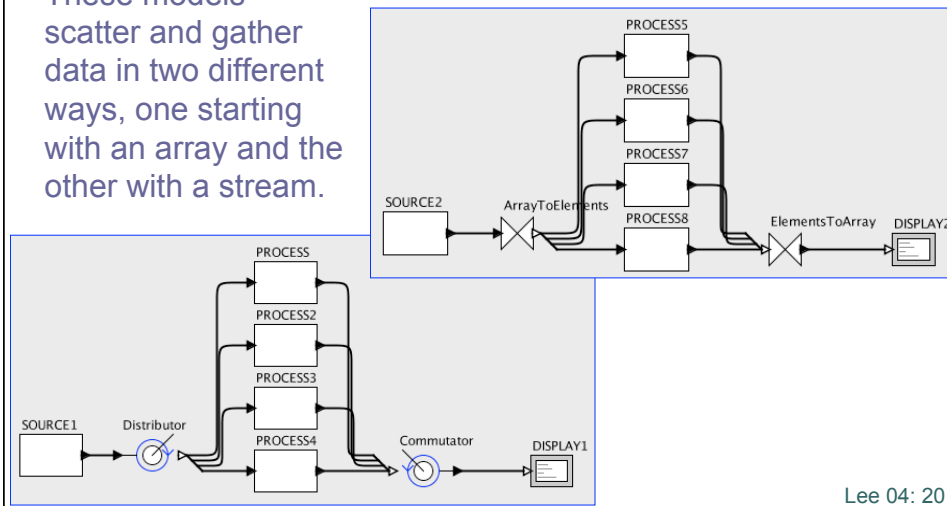E.g., For gather, processes execute

```
MPI_Gather(sendbuf, sendcount, sendtype, recvbuf,
recvcount, recvtype, receivingProcessID, communicator);
```

At the receiving process, this results in recvbuf getting filled with items sent by each of the processes (including the receiving processes).
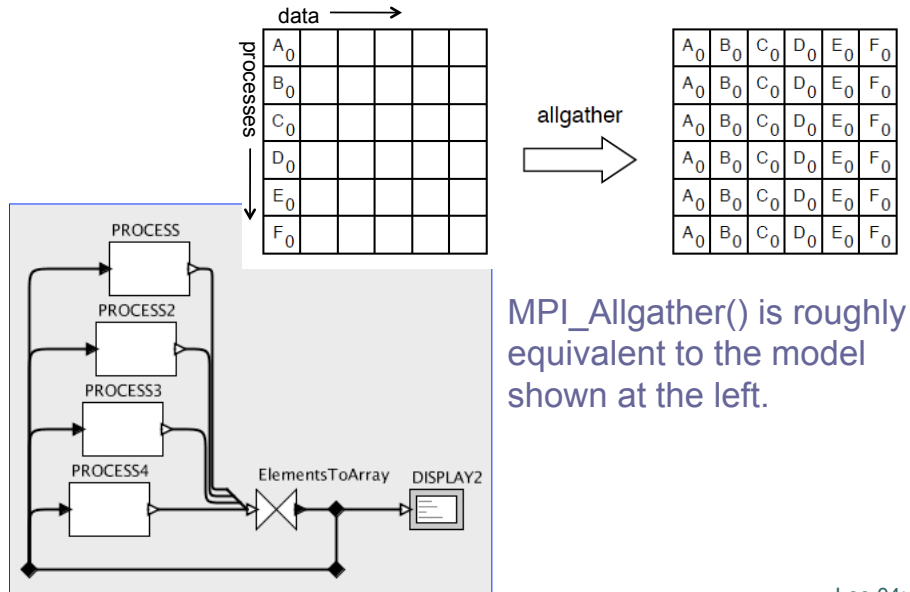
Lee 04: 19

---

## Ptolemy II Mechanisms a bit like Gather/Scatter

These models scatter and gather data in two different ways, one starting with an array and the other with a stream.
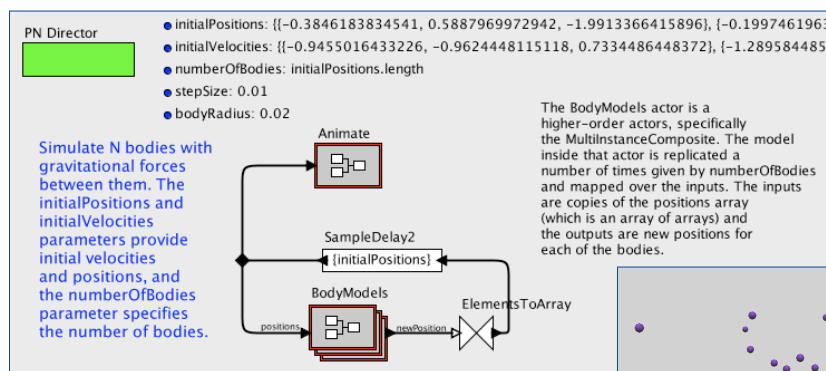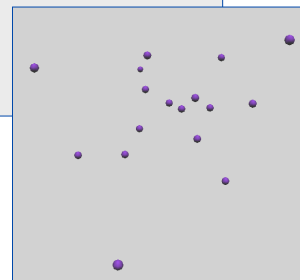


Lee 04: 20

●10

## Gather to all

data →

processes ↓

| $A_0$ | | | | | |
| $B_0$ | | | | | |
| $C_0$ | | | | | |
| $D_0$ | | | | | |
| $E_0$ | | | | | |
| $F_0$ | | | | | |

allgather →

| $A_0$ | $B_0$ | $C_0$ | $D_0$ | $E_0$ | $F_0$ |
|---|---|---|---|---|---|
| $A_0$ | $B_0$ | $C_0$ | $D_0$ | $E_0$ | $F_0$ |
| $A_0$ | $B_0$ | $C_0$ | $D_0$ | $E_0$ | $F_0$ |
| $A_0$ | $B_0$ | $C_0$ | $D_0$ | $E_0$ | $F_0$ |
| $A_0$ | $B_0$ | $C_0$ | $D_0$ | $E_0$ | $F_0$ |
| $A_0$ | $B_0$ | $C_0$ | $D_0$ | $E_0$ | $F_0$ |

PROCESS
PROCESS2
PROCESS3
PROCESS4
ElementsToArray
DISPLAY2

MPI_Allgather() is roughly equivalent to the model shown at the left.

Lee 04: 21

---

## Application of Gather to All:
## Gravitation Simulation

PN Director

- initialPositions: {{−0.3846183834541, 0.5887969972942, −1.9913366415896}, {−0.1997461963
- initialVelocities: {{−0.9455016433226, −0.9624448115118, 0.7334486448372}, {−1.289584485
- numberOfBodies: initialPositions.length
- stepSize: 0.01
- bodyRadius: 0.02

Simulate N bodies with gravitational forces between them. The initialPositions and initialVelocities parameters provide initial velocities and positions, and the numberOfBodies parameter specifies the number of bodies.

The BodyModels actor is a higher-order actors, specifically the MultiInstanceComposite. The model inside that actor is replicated a number of times given by numberOfBodies and mapped over the inputs. The inputs are copies of the positions array (which is an array of arrays) and the outputs are new positions for each of the bodies.

Animate

SampleDelay2
{initialPositions}

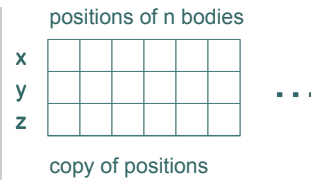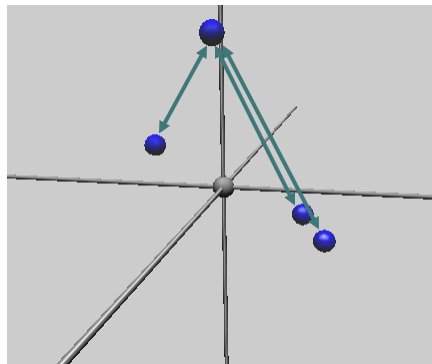BodyModels
positions → newPosition

ElementsToArray

Here, BodyModels receives an array of positions and broadcasts it to one process per body. Each process computes the position of the body at the next time step and the ElementsToArray gathers these into an array.

Lee 04: 22

## How the Gravitation Simulation is a Gather-to-all Pattern
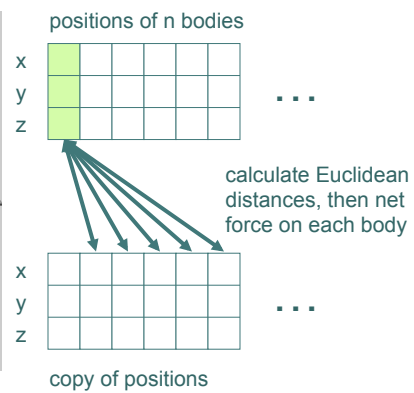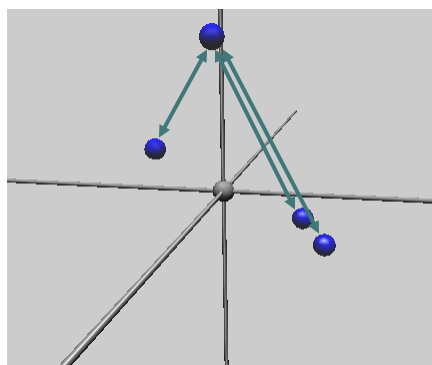
3-D gravitational simulation of *n* bodies



positions of n bodies

x
y
z

· · ·

copy of positions

Thanks to Rodric Rabbah, IBM Watson Center, for suggesting this example.

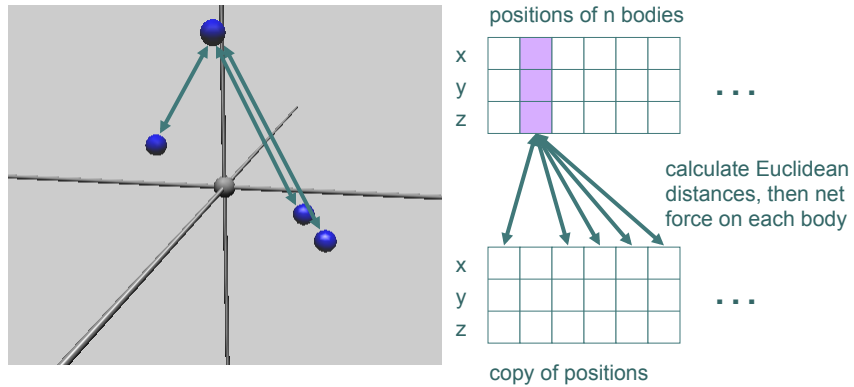## How the Gravitation Simulation is a Gather-to-all Pattern

3-D gravitational simulation of *n* bodies



positions of n bodies

x
y
z

· · ·

calculate Euclidean distances, then net force on each body

x
y
z

· · ·

copy of positions

# How the Gravitation Simulation is a Gather-to-all Pattern

3-D gravitational simulation of *n* bodies



positions of n bodies

calculate Euclidean distances, then net force on each body

copy of positions
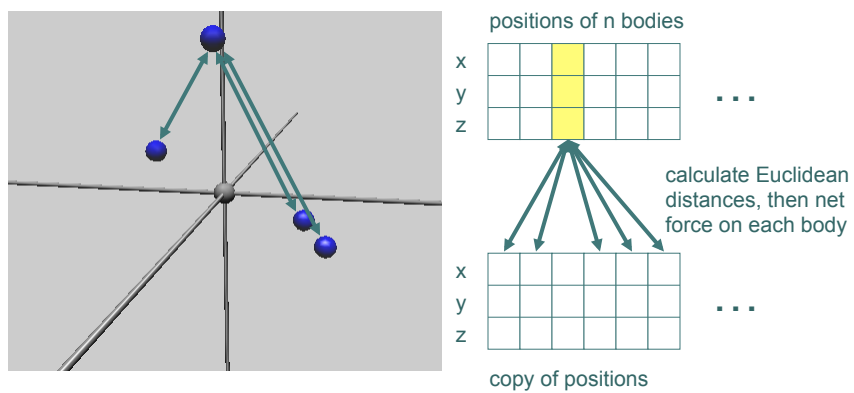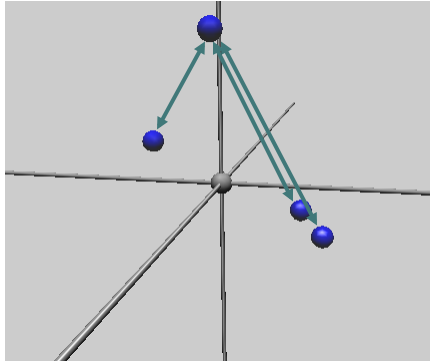
Lee 04: 25


# How the Gravitation Simulation is a Gather-to-all Pattern

3-D gravitational simulation of *n* bodies



positions of n bodies

calculate Euclidean distances, then net force on each body

copy of positions

Lee 04: 26

●13

## How the Gravitation Simulation is a Gather-to-all Pattern

3-D gravitational simulation of $n$ bodies



$$F(t) = ma(t)$$
$$a(t) = F(t)/m$$
$$v(t) = \int_0^t a(\tau)d\tau + v(0)$$
$$p(t) = \int_0^t v(\tau)d\tau + p(0)$$

**A simple (naïve) approximation:**

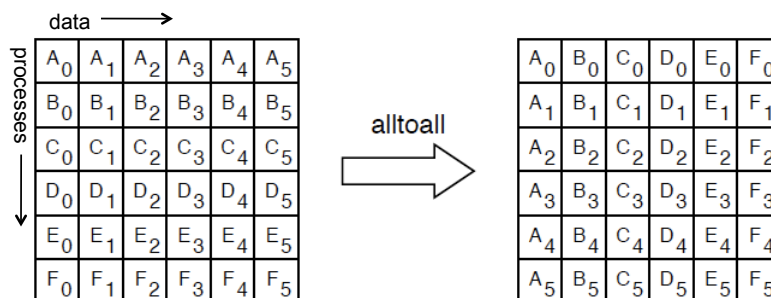$$v(t + \Delta) = v(t) + \Delta a(t)$$
$$p(t + \Delta) = p(t) + \Delta v(t)$$

**Each process computes this approximation.**

Lee 04: 27

---

## All to all Gather/Scatter



Exercise: Realize this pattern in Ptolemy II.

Lee 04: 28

## Reduction Operations

Reduce operations gather data from multiple processes and reduce them using an associative operation (like sum, maximum, …). The operation need not be commutative. The order of reduction is by process ID (called "rank" in MPI).

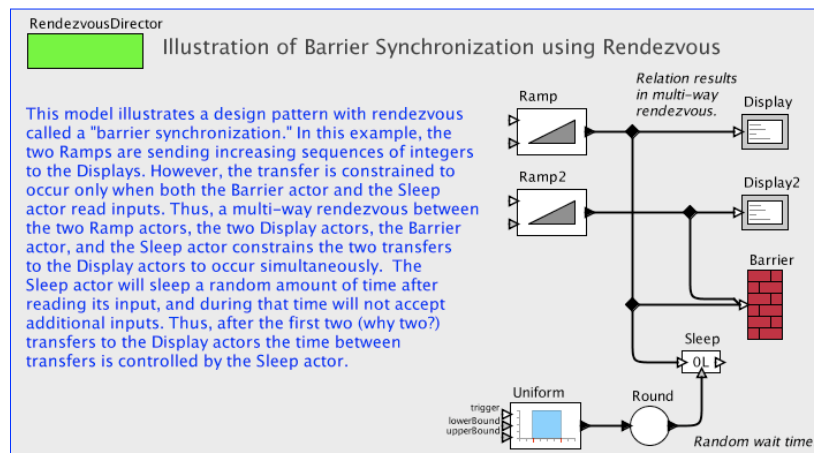Result may be returned to one process or to all.

E.g.,

```
    MPI Reduce(sendbuf, recvbuf, count, type, operation
receivingProcessID, communicator);
```
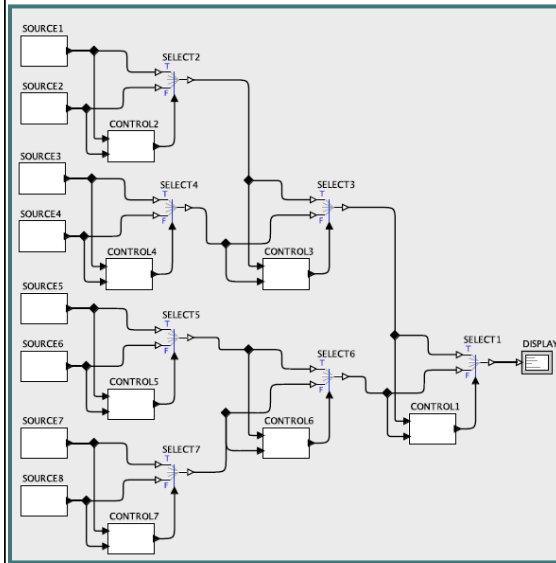
## A Rather Different MPI Pattern:
## Barrier Synchronization

MPI_Barrier() blocks until all members of a group have called it. Ptolemy II equivalent uses the Rendezvous director:



RendezvousDirector

Illustration of Barrier Synchronization using Rendezvous

This model illustrates a design pattern with rendezvous called a "barrier synchronization." In this example, the two Ramps are sending increasing sequences of integers to the Displays. However, the transfer is constrained to occur only when both the Barrier actor and the Sleep actor read inputs. Thus, a multi-way rendezvous between the two Ramp actors, the two Display actors, the Barrier actor, and the Sleep actor constrains the two transfers to the Display actors to occur simultaneously. The Sleep actor will sleep a random amount of time after reading its input, and during that time will not accept additional inputs. Thus, after the first two (why two?) transfers to the Display actors the time between transfers is controlled by the Sleep actor.

Relation results in multi-way rendezvous.

Ramp — Display
Ramp2 — Display2
Barrier
Sleep
Uniform — Round — Random wait time.

●15

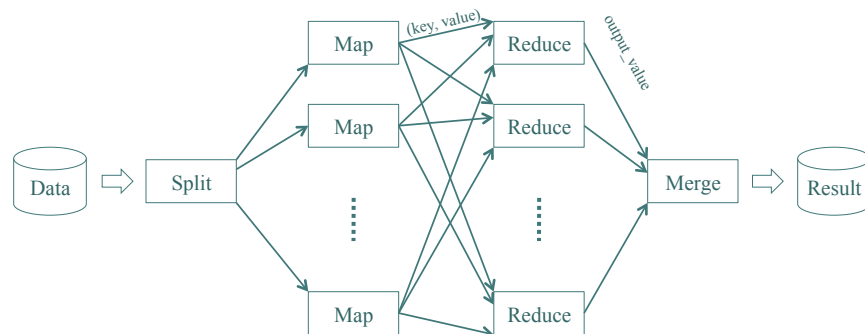## Not provided Directly by MPI: Sorting Trees



Consider collecting time-stamped trades from commodities markets around the world and merging them into a single time-stamped stream. The CONTROL actors could compare time stamps, with logic like this:

```
data1 = topPort.get();
data2 = bottomPort.get();
while (true) {
  if (data1.time < data2.time)) {
    output.send(true);
    data1 = topPort.get();
} else {
    output.send(false);
    data2 = bottomPort.get();}
}
```

Lee 04: 31

---

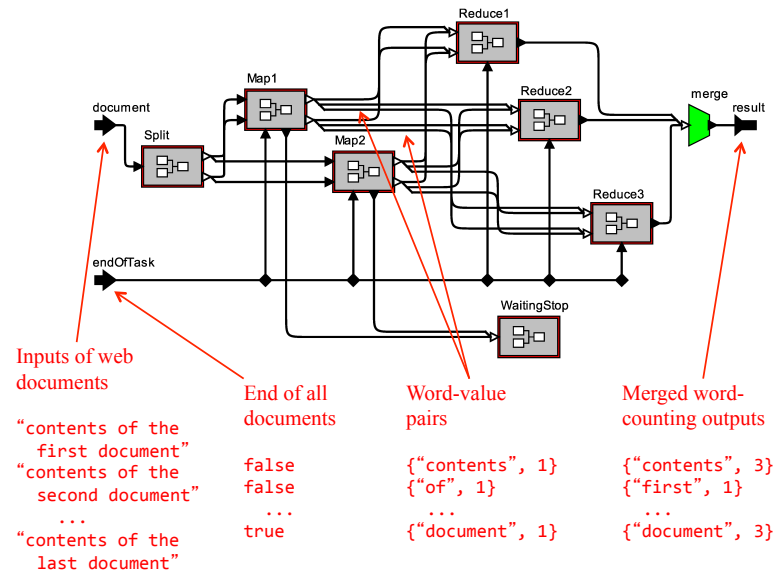## Not provided directly by MPI: Map/Reduce

Dean, J. and S. Ghemawat (2004). {MapReduce}: Simplified Data Processing on Large Clusters. Symposium on Operating System Design and Implementation (OSDI).



This pattern is intended to exploit parallel computing by distributing computations that fit the structure. The canonical example constructs an index of words found in a set of documents.
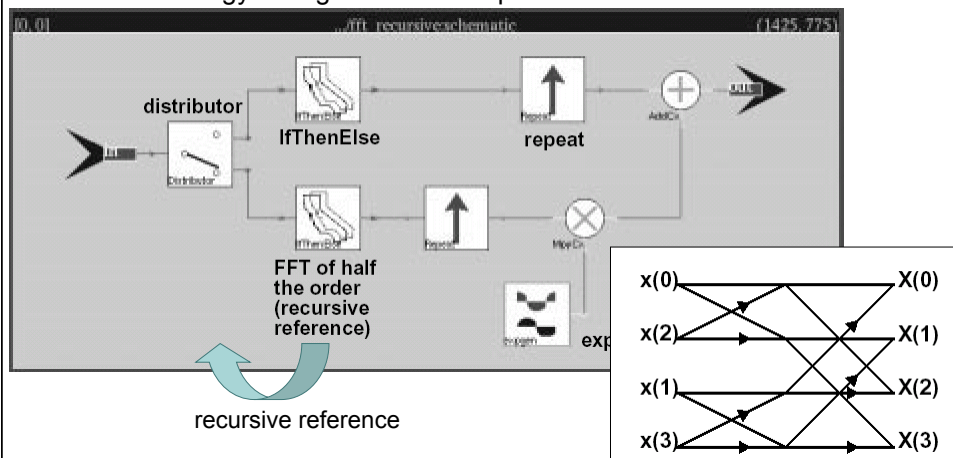
Lee 04: 32

●16

# A MapReduce Model in Ptolemy II

Reduce1

Map1

document

Split

Map2

Reduce2

merge

result

Reduce3

endOfTask

WaitingStop

Inputs of web
documents

"contents of the
  first document"
"contents of the
  second document"
   ...
"contents of the
  last document"

End of all
documents

false
false
...
true

Word-value
pairs

{"contents", 1}
{"of", 1}
   ...
{"document", 1}

Merged word-
counting outputs

{"contents", 3}
{"first", 1}
   ...
{"document", 3}

Lee 04: 33

# Not provided by MPI:
# Recursion

FFT implementation in Ptolemy Classic (1995) used a partial
evaluation strategy on higher-order components.



distributor

IfThenElse

repeat

FFT of half
the order
(recursive
reference)

exp

recursive reference

x(0)  X(0)
x(2)  X(1)
x(1)  X(2)
x(3)  X(3)

Lee 04: 34

●17

## Not provided by MPI: Dynamically Instantiated Processes

```
Process SIFT in QI => QO;
   Vars PRIME;
   repeat
      GET(QI) → PRIME; PUT(PRIME,QO)
      doco FILTER(PRIME,QI)→QI; CONTINUE closeco
   forever
Endprocess;
```
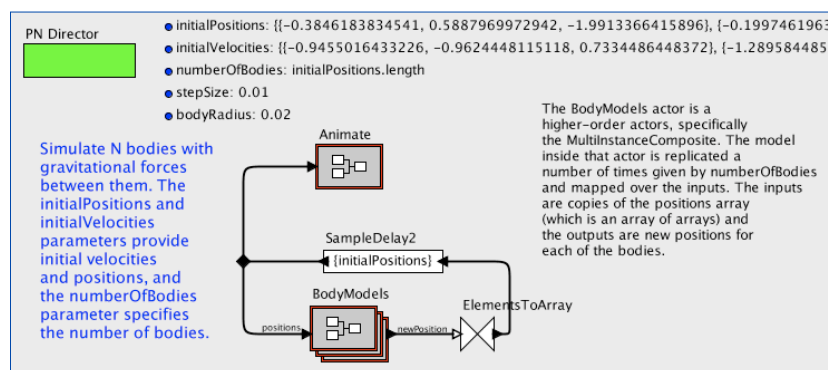
Recall Kahn & MacQueen (1977). Above, a new instance of FILTER is spliced into the pipeline ahead of this process each time a new input arrives.

Kahn, G. and D. B. MacQueen (1977). Coroutines and Networks of Parallel Processes. Information Processing, North-Holland Publishing Co.
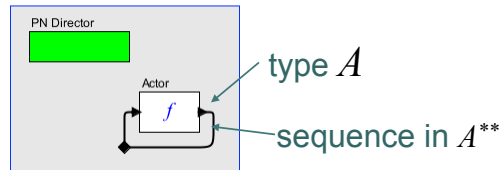
Lee 04: 35

---

## Patterns as Higher-Order Components



BodyModels here is an instance of MultiInstanceComposite, an actor in Ptolemy II that has two parameters: one specifying the number of instances, and one specifying the model to instantiate. This is a "higher-order-component" because it operates on components, not just data.

Lee 04: 36

18

## Reexamining Kahn MacQueen Blocking Reads
or "do we need MPI_Probe()?"

PN Director

Actor $f$

type $A$

sequence in $A^{**}$

Recall: Semantics of a PN Model is the Least Fixed Point of a Monotonic Function:

○ Chain: $C = \{ f(\bot), f(f(\bot)), \ldots , f^n(\bot), \ldots \}$

○ Continuity: $f(\vee\, C) = \vee\, \hat{f}(C)$

Limits

## Kahn-MacQueen Blocking Reads vs. Kahn Continuity

Following Kahn-MacQueen [1977], actors are threads that implement *blocking reads*, which means that when they attempt to read from an empty input, the thread stalls.

● This restricts expressiveness more than continuity

## PN Implementation in Ptolemy II

Body of a process:

```
while (!stopRequested()) {
  …
  if (inputPort.hasToken(channelNo)) {
    …
    Token input = inputPort.get(channelNo);
    …
  }
}
```

When using the PN Director, hasToken() always returns true. Why?

---

## Blocking reads realize
## *sequential* Functions [Vuillemin]

Let $f: A^n \rightarrow A^m$ be an $n$ input, $m$ output function.

Then $f$ is *sequential* if it is continuous and for any $a \in A^n$ there exists an $i \in \{1, \dots n\}$, such that for all $b \in A^n$ where $a \leq b,$
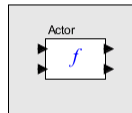
$$a\mid_{\{i\}} = b\mid_{\{i\}} \Rightarrow f(a) = f(b)$$

Intuitively: At all times during an execution, there is an input channel that blocks further output. This is the Kahn-MacQueen blocking read!

## Continuous Function that is not Sequential

Two input identity function is not sequential:



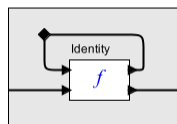Let $f: A^2 \rightarrow A^2$ such that for all $a \in A^2$, $f(a) = a$.
Then $f$ is not sequential.

## Cannot Implement the Two-Input Identity with Blocking Reads
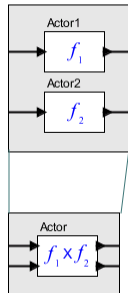
Consider the following connection:



This has a well-defined behavior, but an implementation of the two-input identity with blocking reads will fail to find that behavior.
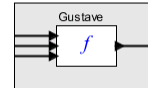
## Sequential Functions do not Compose

If $f_1 : A \to B$ and $f_2 : C \to D$ are sequential then $f_1 \times f_2$ may or may not be sequential. Simple example: suppose $f_1$ and $f_2$ are identity functions in the following:

## Gustave Function
## Non Sequential but Continuous



Let $A = T^{**}$ where $T = \{t, f\}$ .
Let $f : A^3 \to N^{**}$ such that for all $a \in A^3$ ,

$$f(a) = \begin{cases} (1) & \text{if } ((t),(f),\perp) \sqsubseteq a \\ (2) & \text{if } (\perp,(t),(f)) \sqsubseteq a \\ (3) & \text{if } ((f),\perp,(t)) \sqsubseteq a \end{cases}$$

This function is continuous but not sequential.

## Linear Functions [Erhard]

Function $f: A \to B$ on CPOs is *linear* if for all joinable sets $C \subseteq A$, $\hat{f}(C)$ is joinable and

$$\vee \hat{f}(C) = f(\vee C)$$

Intuition: If two possible inputs can be extended to a common input, then the two corresponding outputs can be extended to the common output.

Fact: Sequential functions are linear.
Fact: Linear functions are continuous (trivial)

## Stable Functions [Berry]

Function $f: A \to B$ on complete semilattices (CPOs where every subset has a greatest lower bound) is *stable* if it is continuous and for all joinable sets $C \subseteq A$, $\hat{f}(C)$ is joinable and

$$\wedge \hat{f}(C) = f(\wedge C) \qquad \longleftarrow \quad \text{NOTE: meet! not join!}$$

Intuition: If two possible inputs do not contain contradictory information, then neither will the two corresponding outputs.

Fact: Sequential functions are stable.

## Summary

- MPI is an underspecified standard (buffering issues)
- MPI programs are not modular
- Collective operations in MPI are useful
- There are useful collective operations not specified in MPI
- Collective operations can be viewed as higher-order components.
- Constraint to blocking reads makes process networks non-compositional.
- Constraint to blocking reads precludes implementing certain continuous functions (but are any of those useful?)