# Partial Globalization of Partitioned Address Spaces for Zero-copy Communication with Shared Memory

Fangzhou Jiao, Nilesh Mahajan, Jeremiah Willcock,
**Arun Chauhan**, Andrew Lumsdaine
Indiana University

HiPC 2011

# Motivation

- Increasing popularity and availability of many-cores

- Abundance of legacy MPI code

- Simplifying programming model

    - single model, instead of hybrid

- Leveraging shared memory fully for performance

- Proving that shared memory could be used as an optimization for communication

# Partitioned Address Space Programming on Shared Memory

- Avoids having to worry about race conditions

- Encourages programmers to think about locality

- Could make it easier to reason about program correctness

  - if done at the right level of abstraction

Needs special handling to compete in performance with threaded shared memory programs

*Arun Chauhan, Zero-copy communication in partitioned address space programs on shared memory, HiPC 2011*
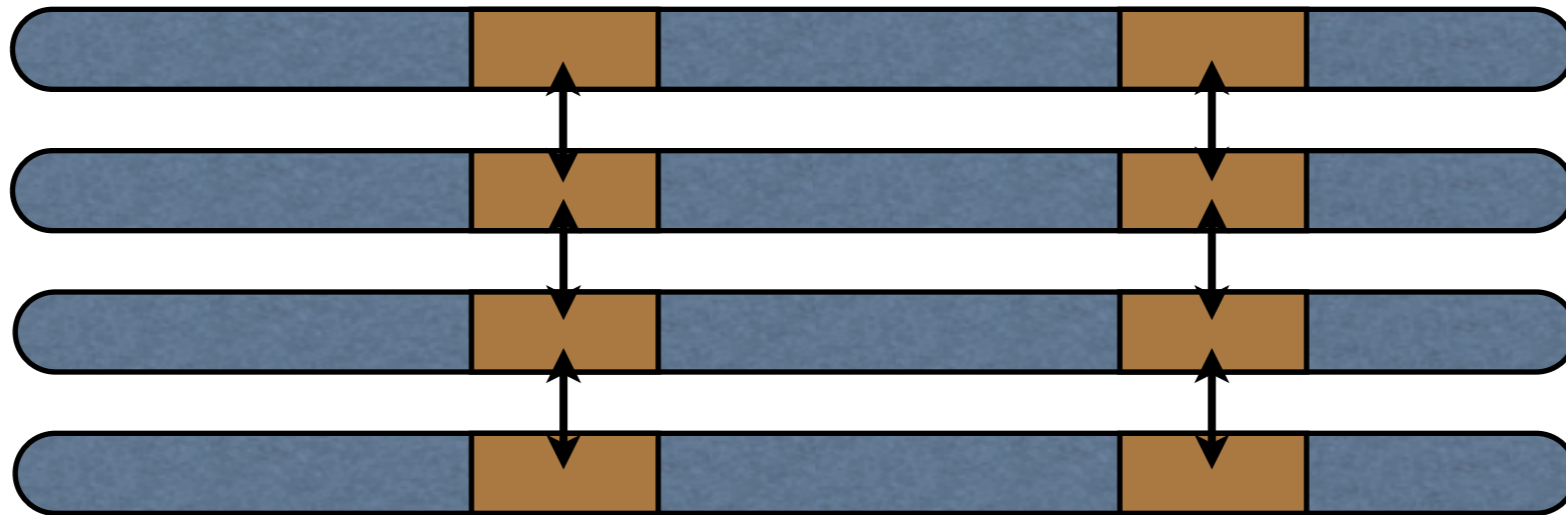
# Declarative Approach

- Originally motivated by Block-synchronous Parallel (BSP) programs, especially for collective communication

  - alternate between computation and communication

  - communication optimization breaks the structure

# Declarative Approach

- Originally motivated by Block-synchronous Parallel (BSP) programs, especially for collective communication

  - alternate between computation and communication

  - communication optimization breaks the structure

# Declarative Approach

- Originally motivated by Block-synchronous Parallel (BSP) programs, especially for collective communication

  - alternate between computation and communication

  - communication optimization breaks the structure

# Declarative Approach

- Originally motivated by Block-synchronous Parallel (BSP) programs, especially for collective communication

  - alternate between computation and communication

  - communication optimization breaks the structure

- Extend to non BSP-style applications

# Kanor for Clusters

$$@communicate \; \{ \; b@recv\_rank <<= a@send\_rank \; \}$$

**Eric Holk, William E. Byrd, Jeremiah Willcock, Torsten Hoefler, Arun Chauhan, and Andrew Lumsdaine.** *Kanor: A Declarative Language for Explicit Communication. In Proceedings of the Thirteenth International Symposium on the Practical Aspects of Declarative Languages (PADL), 2011. Held in conjunction with the ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages (POPL).*

# Kanor for Clusters

$$\textbf{\textit{@communicate}} \ \{ \ b@recv\_rank <<= a@send\_rank \ \}$$

$$e_0@e_1 << op << e_2@e_3 \ \texttt{where} \ e_4$$

**Eric Holk, William E. Byrd, Jeremiah Willcock, Torsten Hoefler, Arun Chauhan, and Andrew Lumsdaine.** *Kanor: A Declarative Language for Explicit Communication. In Proceedings of the Thirteenth International Symposium on the Practical Aspects of Declarative Languages (PADL), 2011. Held in conjunction with the ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages (POPL).*

# Kanor for Clusters

$$@\textbf{\textit{communicate}} \; \{ \, b@recv\_rank <<= a@send\_rank \, \}$$

$$e_0@e_1 << op << e_2@e_3 \; \text{where} \; e_4$$

$$e_0@e_1 <<= e_2@e_3 \; \text{where} \; e_4$$

**Eric Holk, William E. Byrd, Jeremiah Willcock, Torsten Hoefler, Arun Chauhan, and Andrew Lumsdaine.** *Kanor: A Declarative Language for Explicit Communication. In Proceedings of the Thirteenth International Symposium on the Practical Aspects of Declarative Languages (PADL), 2011. Held in conjunction with the ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages (POPL).*

# Kanor for Clusters

$$@\textbf{\textit{communicate}} \; \{ \; b@recv\_rank <<= a@send\_rank \; \}$$

$$e_0@e_1 << op << e_2@e_3 \; \text{where} \; e_4$$

$$e_0@e_1 <<= e_2@e_3 \; \text{where} \; e_4$$

$$\underbrace{\text{A[j]}}_{\substack{storage \\ location}} @ \underbrace{\text{i}}_{\substack{receiver \\ rank}} \quad \underbrace{\text{<<=}}_{\substack{reduction \\ operator}} \quad \underbrace{\text{B[i]}}_{data} @ \underbrace{\text{j}}_{\substack{sender \\ rank}} \quad \text{where} \; \underbrace{\text{i in world}}_{generator}, \underbrace{\text{j in \{0...i\}}}_{generator}, \underbrace{\text{i \% 2 == 0}}_{filter}$$

**Eric Holk, William E. Byrd, Jeremiah Willcock, Torsten Hoefler, Arun Chauhan, and Andrew Lumsdaine.** *Kanor: A Declarative Language for Explicit Communication. In Proceedings of the Thirteenth International Symposium on the Practical Aspects of Declarative Languages (PADL), 2011. Held in conjunction with the ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages (POPL).*
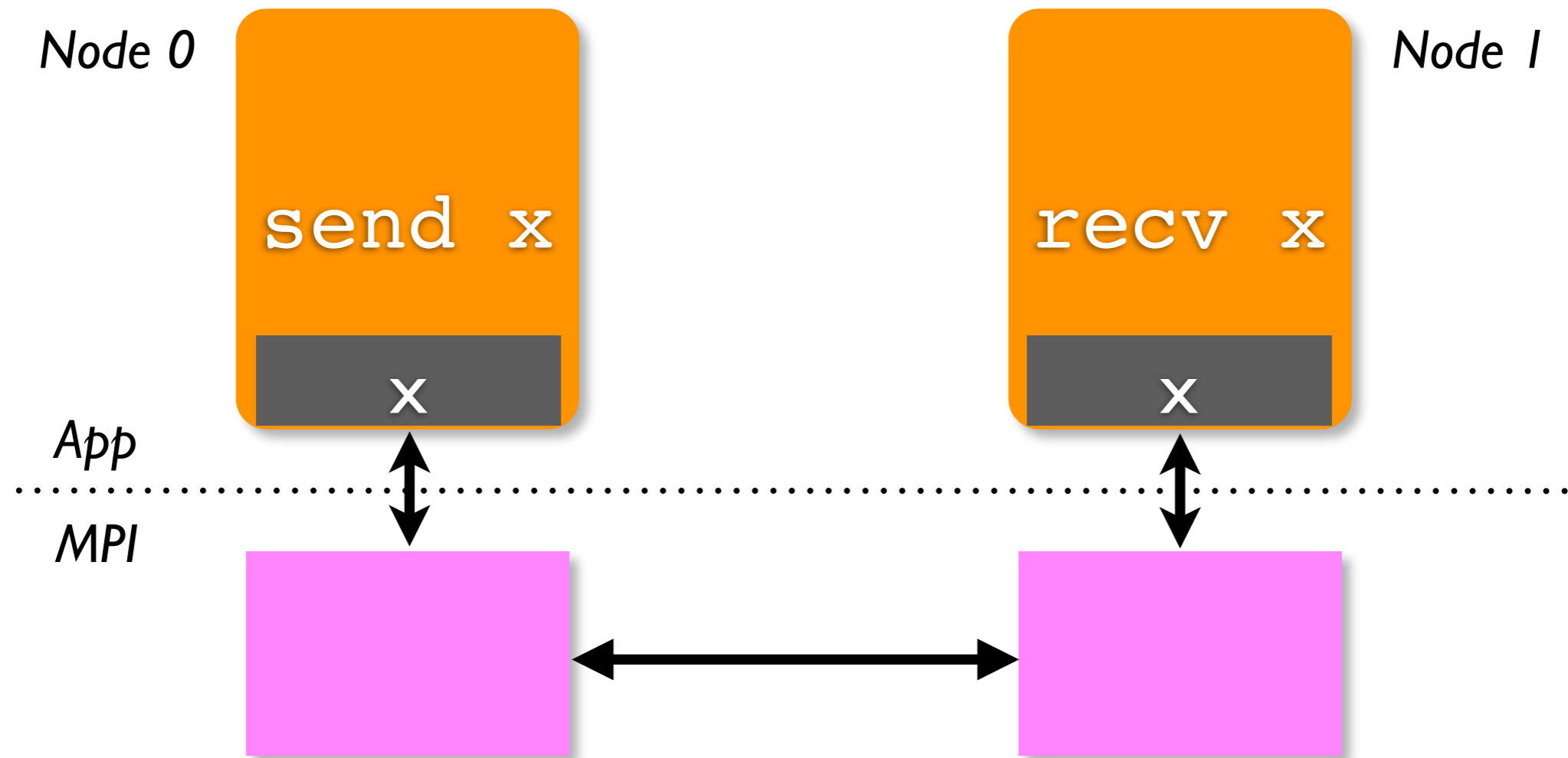
# Kanor for Clusters

$$@communicate \{ b@recv\_rank <<= a@send\_rank \}$$

$$e_0@e_1 << op << e_2@e_3 \text{ where } e_4$$

$$e_0@e_1 <<= e_2@e_3 \text{ where } e_4$$

$$\underbrace{\texttt{A[j]}}_{\substack{\textit{storage} \\ \textit{location}}} \texttt{@} \underbrace{\texttt{i}}_{\substack{\textit{receiver} \\ \textit{rank}}} \quad \underbrace{\texttt{<<=}}_{\substack{\textit{reduction} \\ \textit{operator}}} \quad \underbrace{\texttt{B[i]}}_{\textit{data}} \texttt{@} \underbrace{\texttt{j}}_{\substack{\textit{sender} \\ \textit{rank}}} \texttt{ where } \underbrace{\texttt{i in world}}_{\textit{generator}}, \underbrace{\texttt{j in \{0...i\}}}_{\textit{generator}}, \underbrace{\texttt{i \% 2 == 0}}_{\textit{filter}}$$

$\downarrow$

## Source-level compiler (using ROSE)

$\downarrow$

*standard C++ code*

**Eric Holk, William E. Byrd, Jeremiah Willcock, Torsten Hoefler, Arun Chauhan, and Andrew Lumsdaine.** *Kanor: A Declarative Language for Explicit Communication. In Proceedings of the Thirteenth International Symposium on the Practical Aspects of Declarative Languages (PADL), 2011. Held in conjunction with the ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages (POPL).*

# Design Principles

- Users must think in parallel (creativity)

  - but not be encumbered with optimizations that can be automated, or proving synchronization correctness

- Compiler focuses on what it can do (mechanics)

  - not creative tasks, such as determining data distributions, or creating new parallel algorithms

- Incremental deployment

  - not a new programming language

  - more of a *coordination language* (DSL)
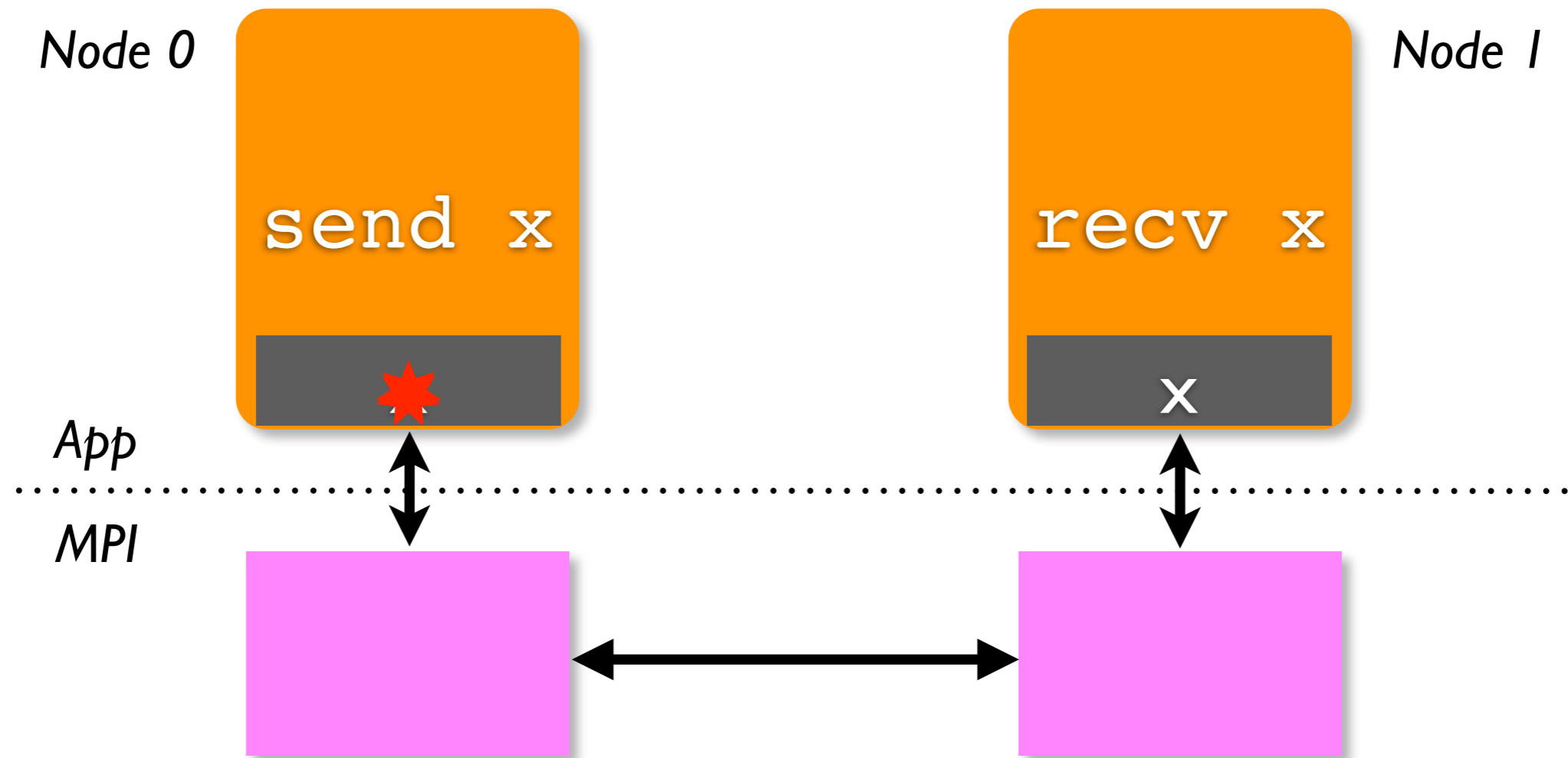
- Formal semantics

  - provable correctness

# Compiling to MPI

`@communicate {x@1 <<= x@0}`

# Compiling to MPI

`@communicate {x@1 <<= x@0}`



Node 0

Node 1

`send x`

`recv x`

x

App

MPI

# Compiling to MPI

`@communicate {x@1 <<= x@0}`

# Compiling to MPI

`@communicate {x@1 <<= x@0}`

# Compiling to MPI

`@communicate {x@1 <<= x@0}`

# Compiling to MPI

`@communicate {x@1 <<= x@0}`



3 copies

# MPI Optimized for Shared Memory

`@communicate {x@1 <<= x@0}`



Node 0

Node 1

send x

recv x

x

x

App

MPI

# MPI Optimized for Shared Memory

`@communicate {x@1 <<= x@0}`

*Node 0*

`send x`

x

*Node 1*

`recv x`

x

*App*

*MPI*

# MPI Optimized for Shared Memory

`@communicate {x@1 <<= x@0}`

Node 0

`send x`

Node 1

`recv x`

App

MPI

# MPI Optimized for Shared Memory

`@communicate {x@1 <<= x@0}`

# MPI Optimized for Shared Memory

`@communicate {x@1 <<= x@0}`

# MPI Optimized for Shared Memory

`@communicate {x@1 <<= x@0}`



2 copies

# Optimizing for Shared Memory

`@communicate {x@1 <<= x@1}`

# Optimizing for Shared Memory

`@communicate {x@1 <<= x@1}`

Node 0

`send x`

Node 1

`recv x`

x

App

MPI

# Optimizing for Shared Memory

`@communicate {x@1 <<= x@1}`

# Optimizing for Shared Memory
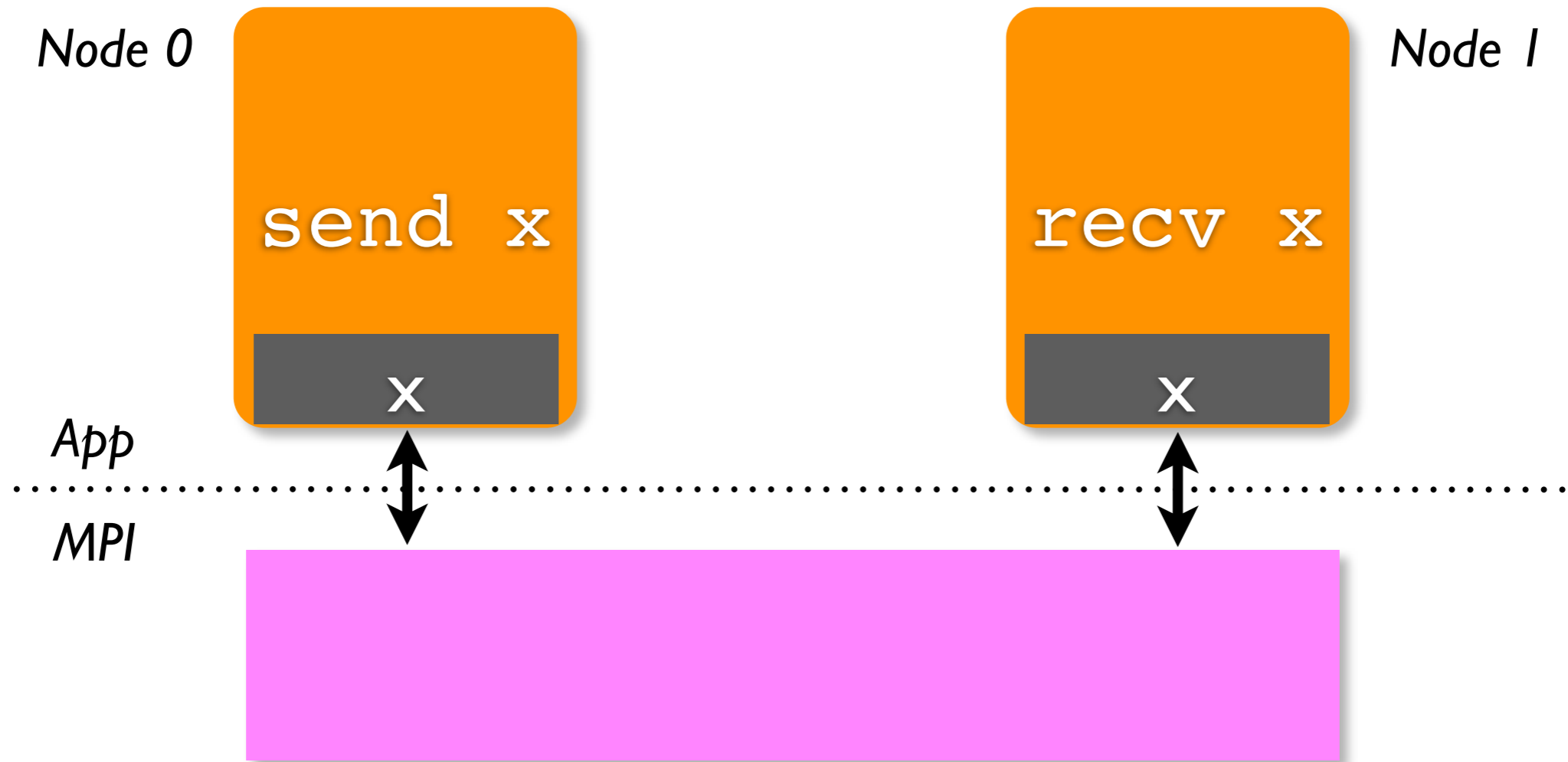
`@communicate {x@1 <<= x@1}`



I copy

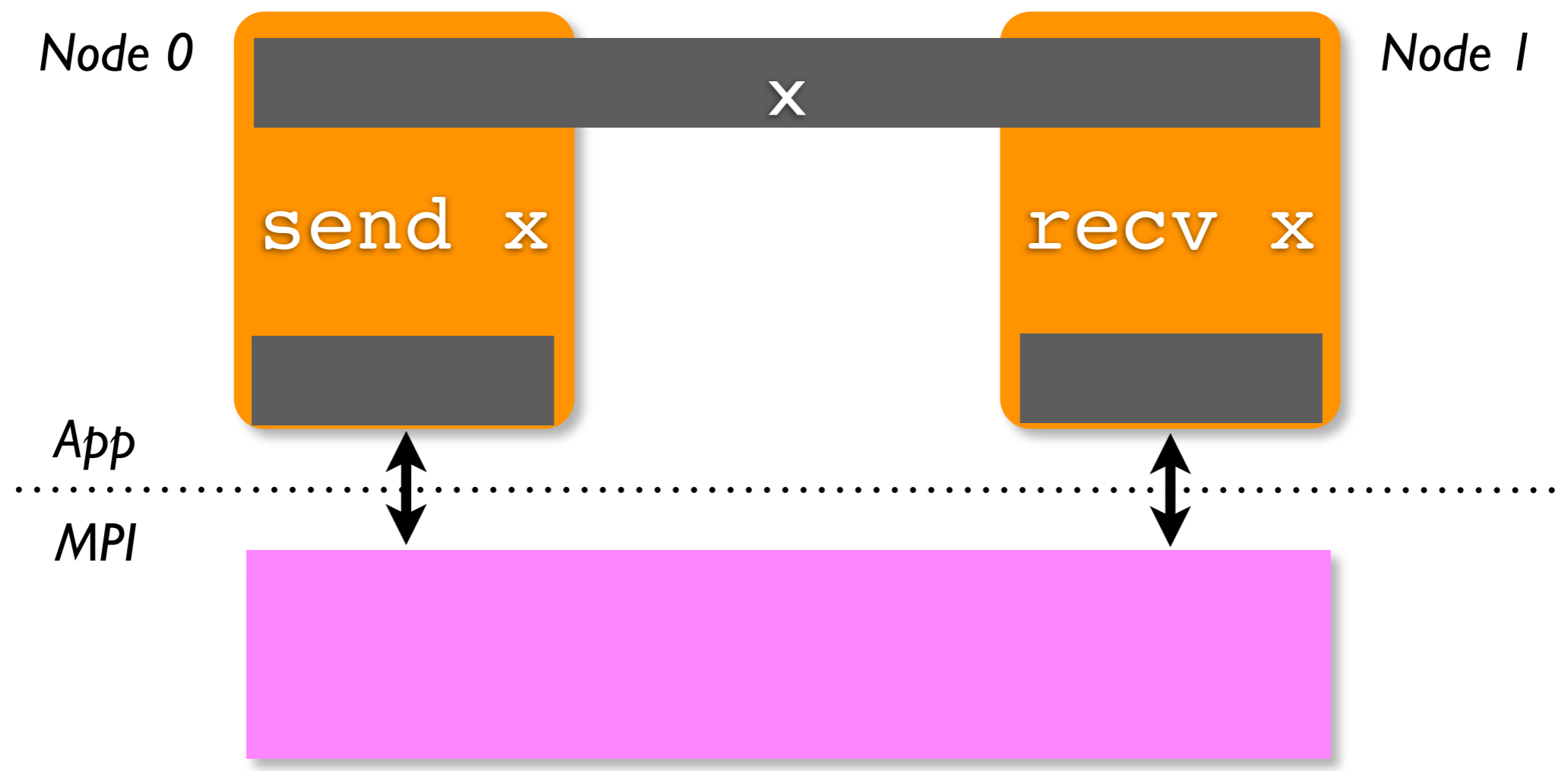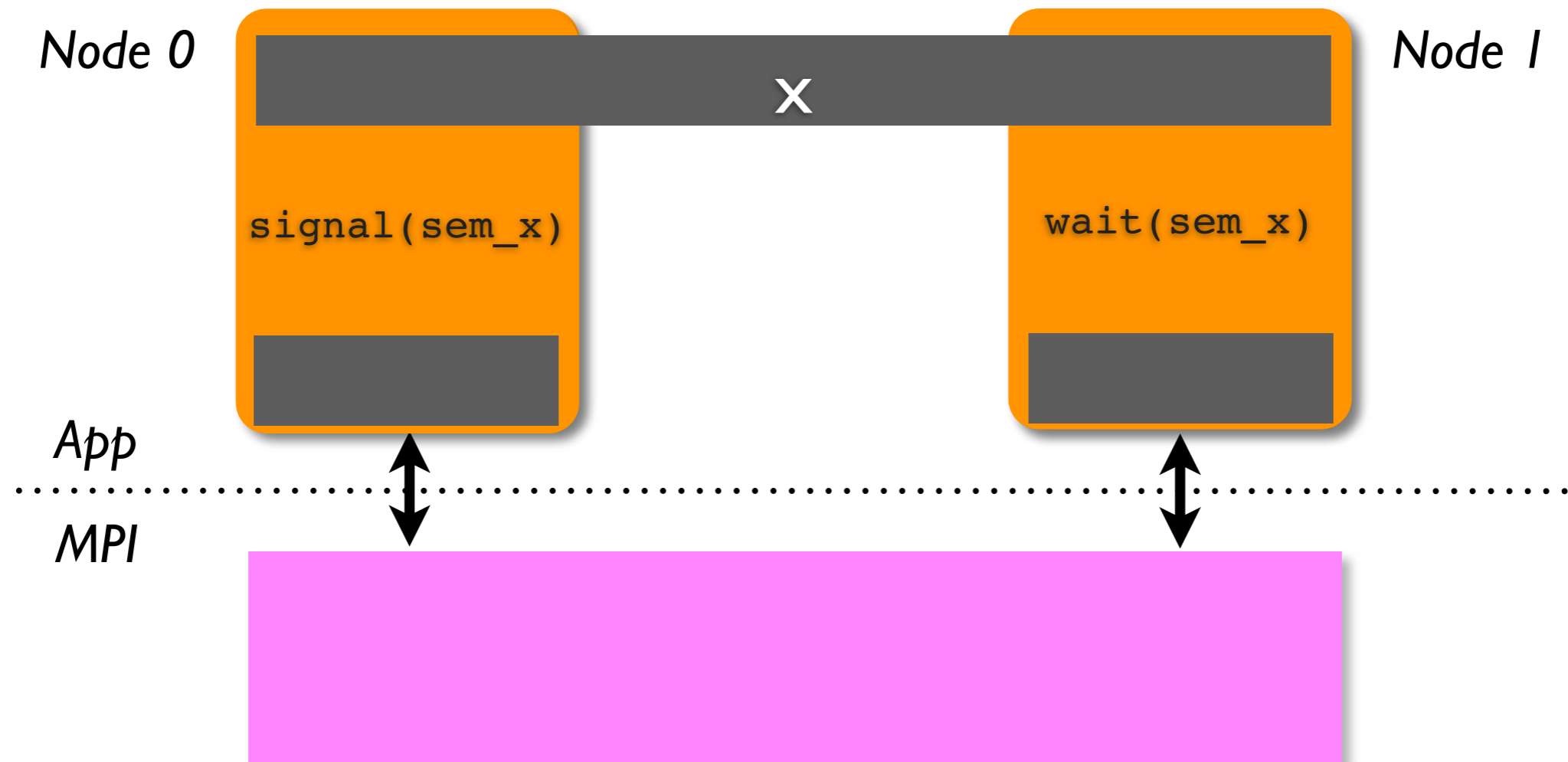*(requires rendezvous or compiler intervention)*

# Optimizing for Shared Memory

## @communicate {x@0 <<= x@1}
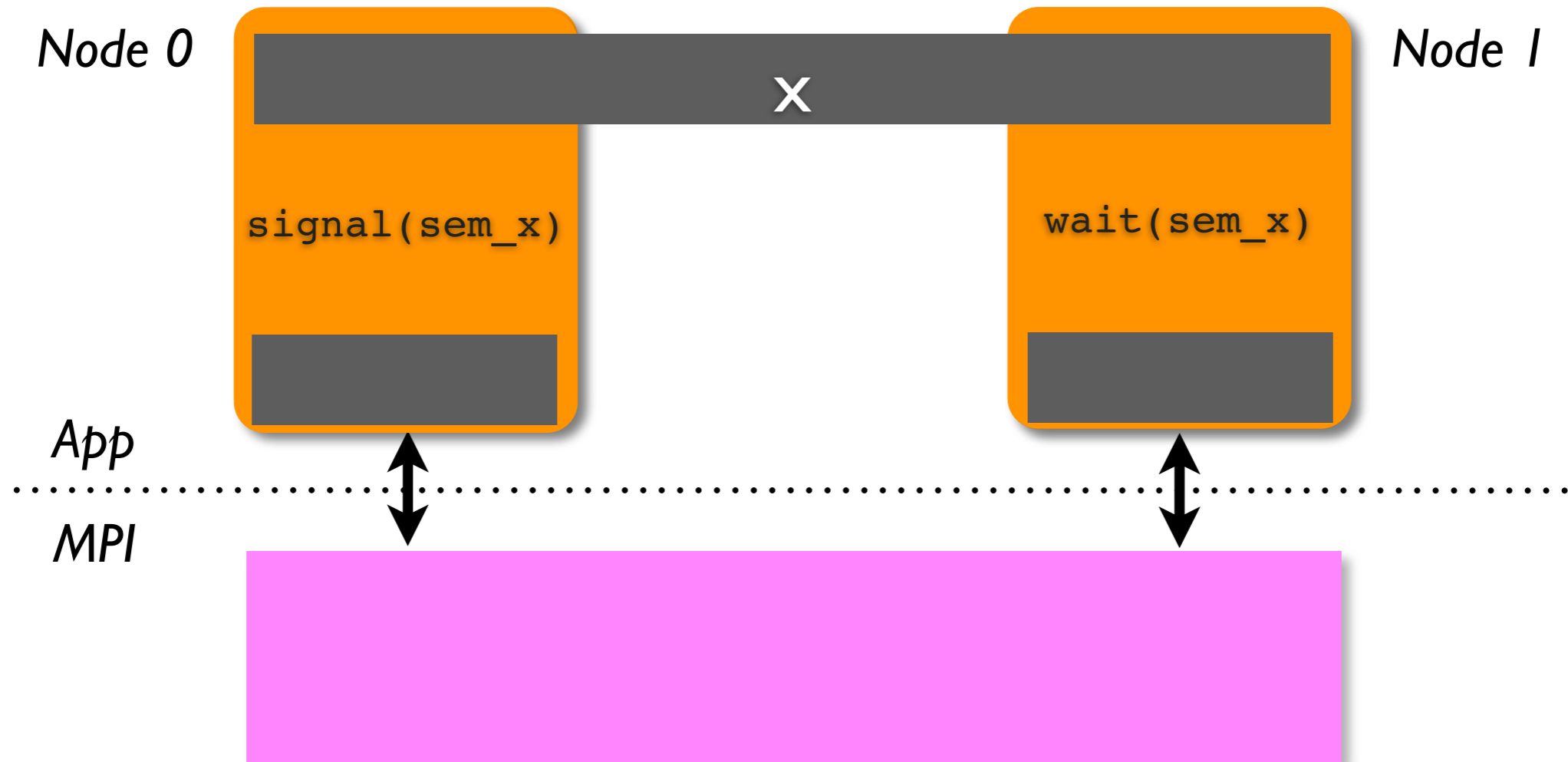
# Optimizing for Shared Memory

`@communicate {x@0 <<= x@1}`

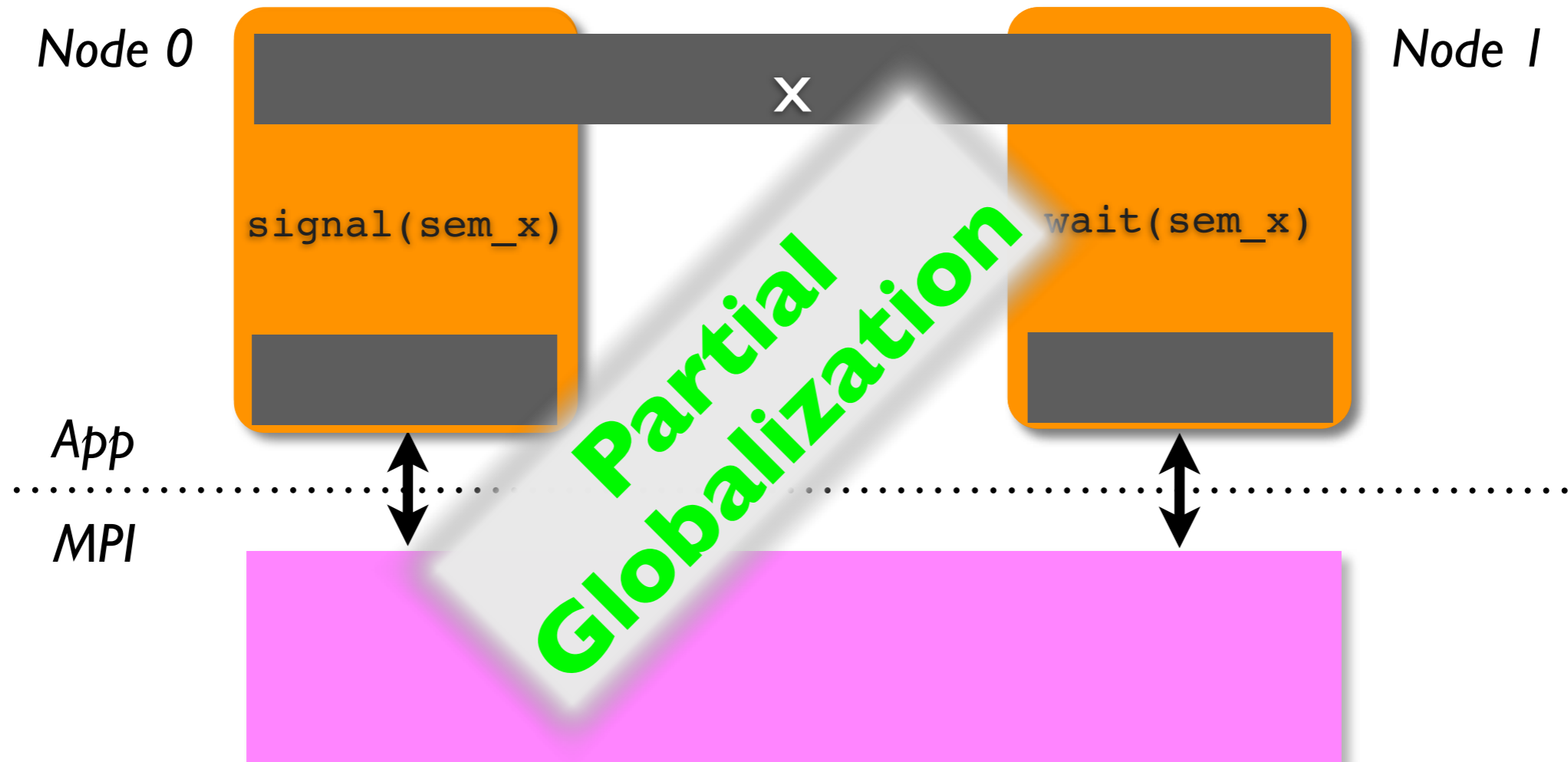# Optimizing for Shared Memory

## @communicate {x@0 <<= x@1}

# Optimizing for Shared Memory

`@communicate {x@0 <<= x@1}`



Node 0

Node 1

x

`signal(sem_x)`

`wait(sem_x)`

App

MPI

## 0 copy
*(requires compiler intervention)*

# Optimizing for Shared Memory

`@communicate {x@0 <<= x@1}`



*Node 0*

`signal(sem_x)`

*Node 1*

`wait(sem_x)`

**Partial Globalization**

*App*

*MPI*

0 copy
*(requires compiler intervention)*

# Steps for Optimizing Communication with Shared Memory

- Identify globalization candidates

- Ensure correctness

  - insert appropriate synchronization

- Minimize contention

  - minimize synchronization points

  - minimize synchronization overheads

    - *using a run-time trick*

# Globalization Candidates

- Contiguous chunks of memory

    - excluding strided array sections, for example

    - contiguous array sections OK (but not implemented)

- Large buffers

    - communication inside loops

- Small local reuse

# Ensuring Correctness

```
@communicate {x@i <<= x@i+1,
              where i in Kanor::WORLD}
…
consume(A);      // consume communicated data
…
overwrite(A);    // reuse A for local data
…
consume(A);      // consume local data
```

# Ensuring Correctness

```
@communicate {x@i <<= x@i+1,
               where i in Kanor::WORLD}
…
consume(A);        // consume communicated data
…
overwrite(A);      // reuse A for local data
…
consume(A);        // consume local data
```
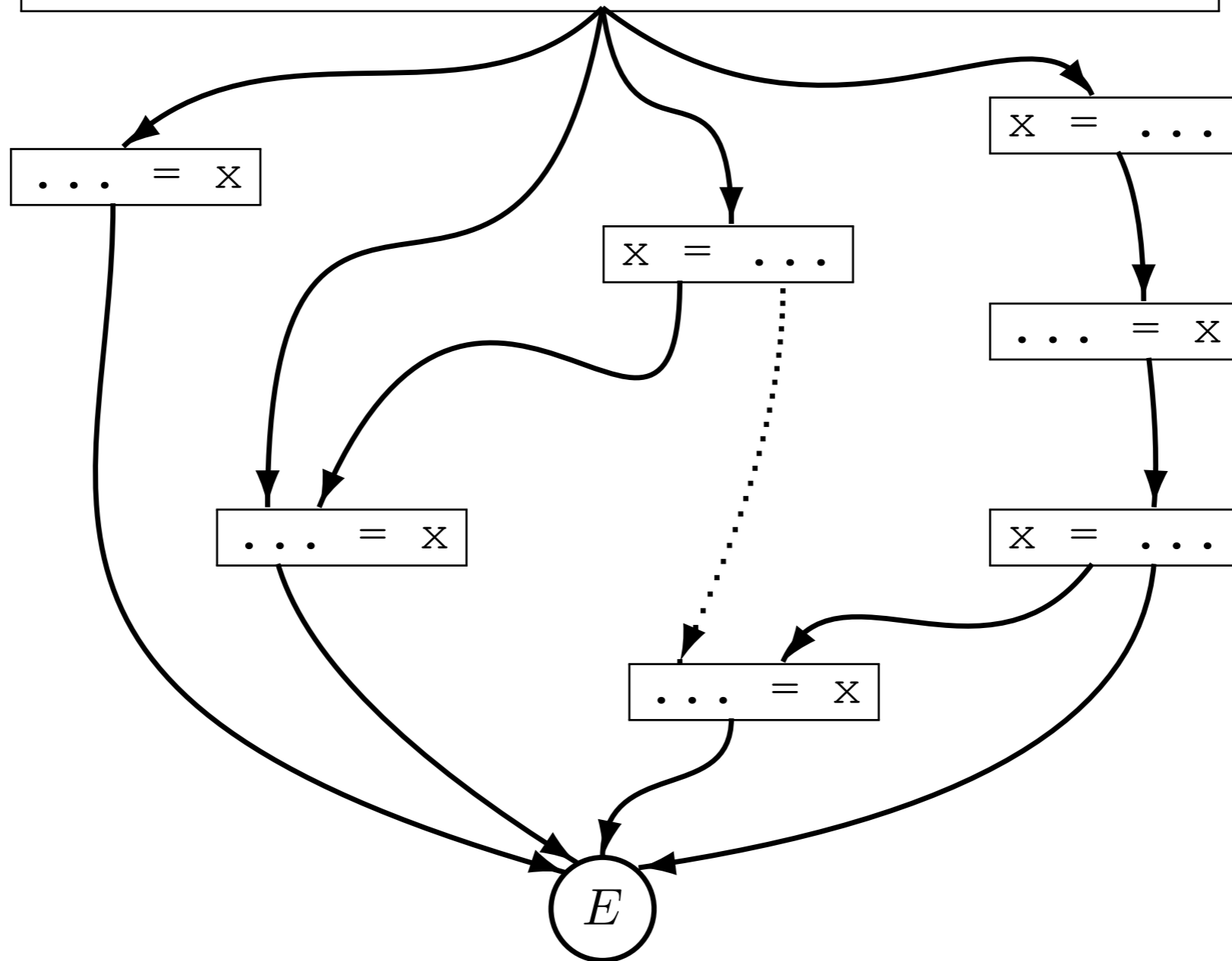
# Correctness Issues



@communicate{x@i <<= x@0}, where i > 0

# Observations

**Definition:** *Locking Set: The set of CFG nodes that lie on a path from a node containing local write into a globalized variable to a node containing read of that value*

# Observations

**Definition:** *Locking Set: The set of CFG nodes that lie on a path from a node containing local write into a globalized variable to a node containing read of that value*
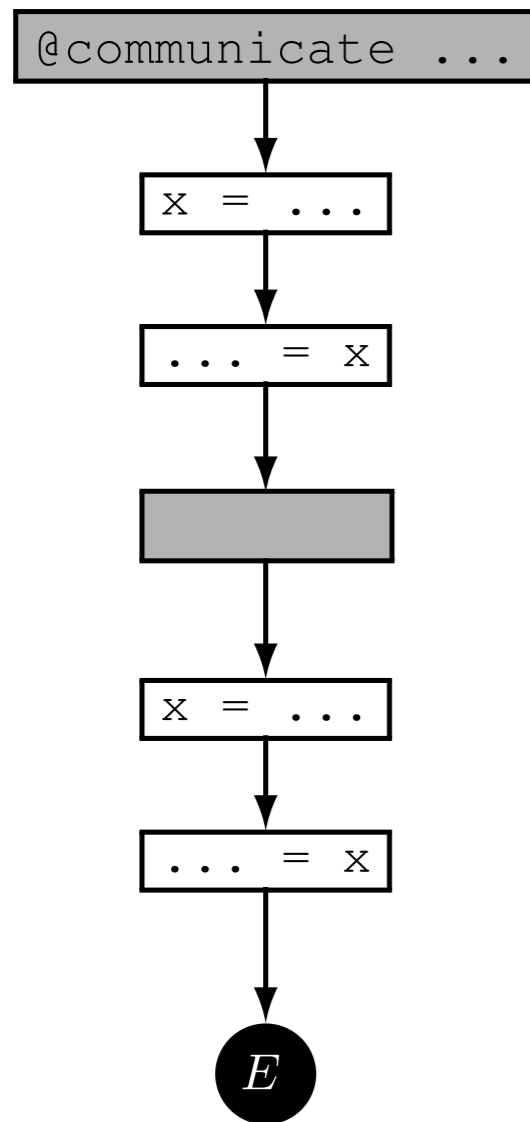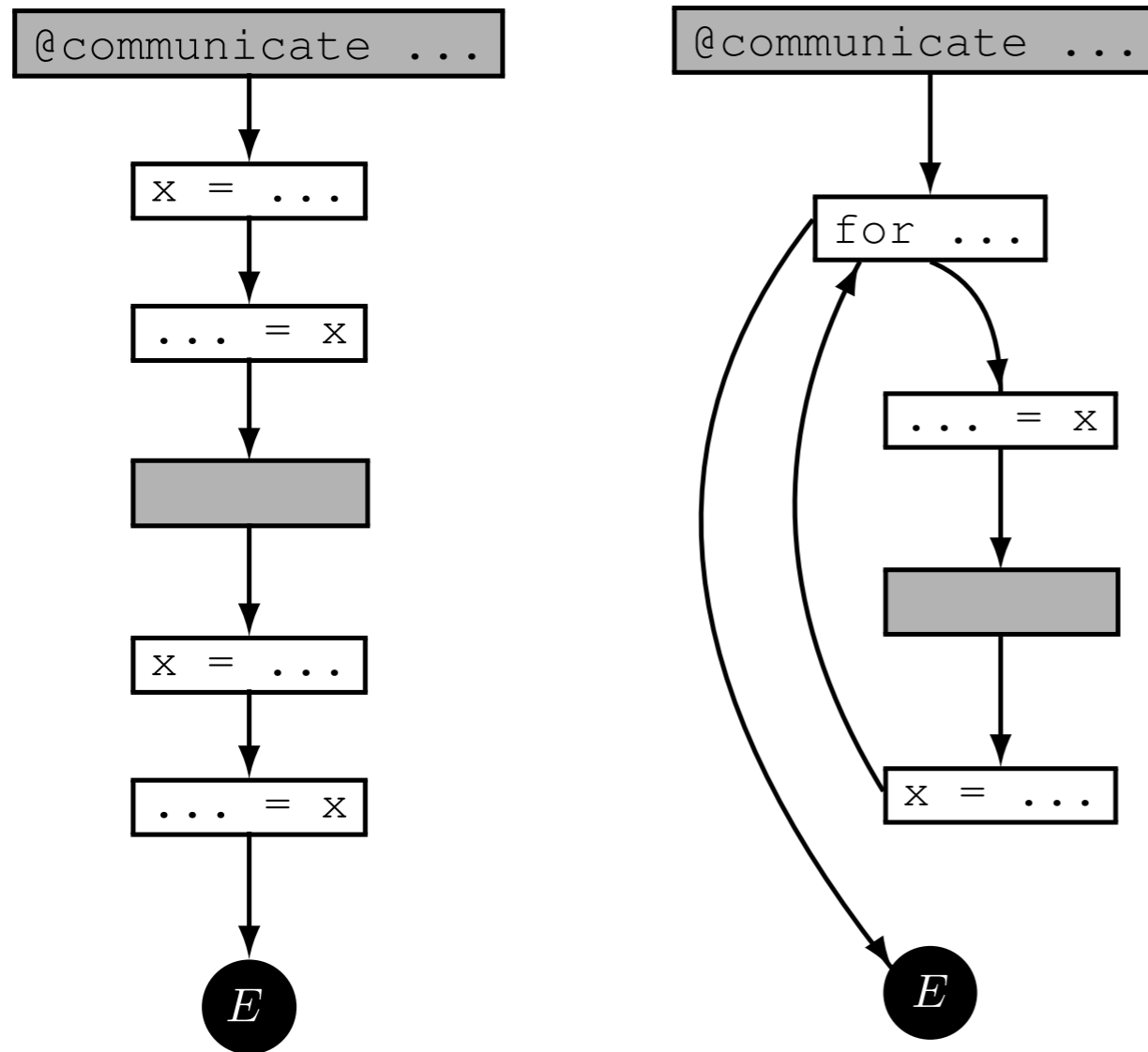
**Theorem:** *If the locking set belongs to a critical section then the partitioned address space semantics are maintained*
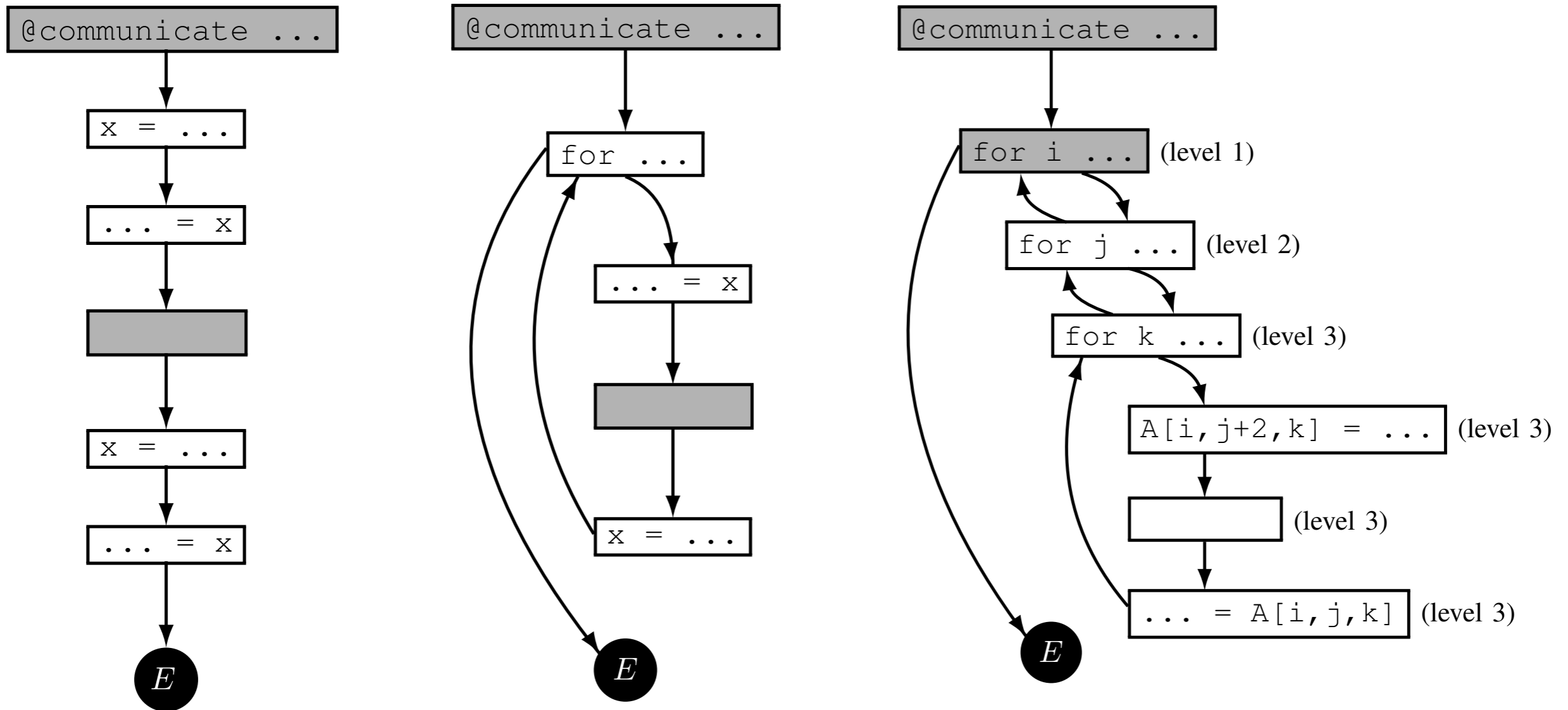
# Correctness: Examples



```
@communicate ...
```

```
x = ...
```

```
... = x
```

```
x = ...
```

```
... = x
```

$E$

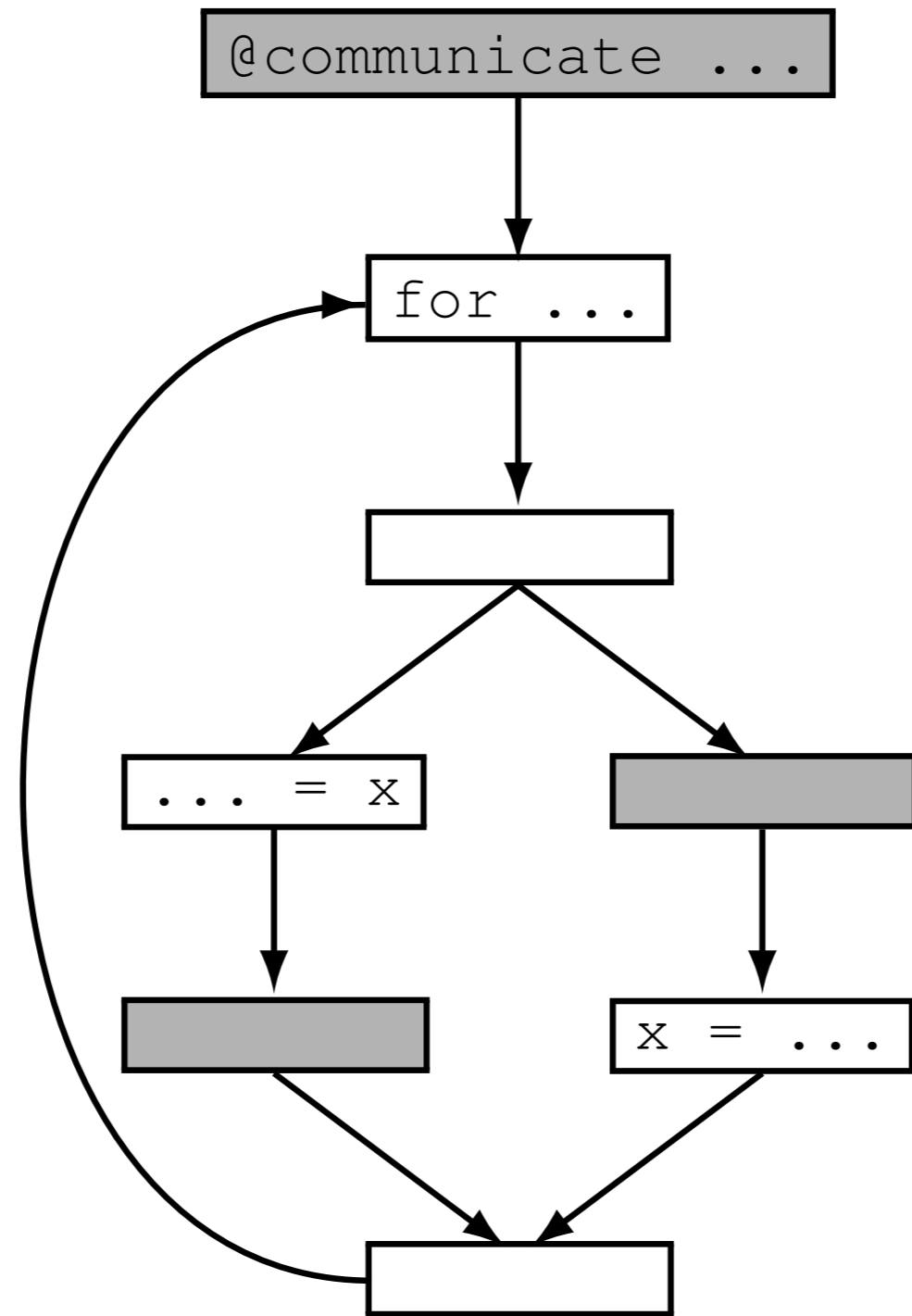# Correctness: Examples

# Correctness: Examples

# Overall Algorithm

- Identify globalization candidates

- For each globalized variable

    - compute the *locking set*

    - divide the locking set into *connected components*, $C_i$

        - CFG edge into $C_i$ $\Rightarrow$ insert `lock_acquire`

        - CFG edge out of $C_i$ $\Rightarrow$ insert `lock_release`

# Example of sub-optimal Behavior



```
@communicate ...
```

```
for ...
```

```
... = x
```

```
x = ...
```

# Copy-on-conflict

```
1  void acquire_or_copy (Buffer& a, Lock& lock)
2  {
3     if (Localized[a]) return NULL;
4     Condition cond;
5     enum {COPY_THRD, LOCK_THRD} notifier;
6     a_cpy = new Buffer;
7
8     Thread l_thrd =
9        spawn(acquire_lock, lock, cond, &notifier);
10    Thread c_thrd =
11       spawn(buf_copy, a, a_cpy, cond, &notifier);
12    wait(cond);
13
14    if (notifier == LOCK_THRD) {
15       c_thrd.kill();
16       free(a_cpy);
17    } else {
18       l_thrd.kill();
19       if (lock.held()) lock.release();
20       delete a;
21       a = a_cpy;
22       Localized[a] = true;
23    }
24  }
```

# Copy-on-conflict

```
1  void acquire_or_copy (Buffer& a, Lock& lock)
2  {
3    if (Localized[a]) return NULL;
4    Condition cond;
5    enum {COPY_THRD, LOCK_THRD} notifier;
6    a_cpy = new Buffer;
7
8    Thread l_thrd =
9      spawn(acquire_lock, lock, cond, &notifier);
10   Thread c_thrd =
11     spawn(buf_copy, a, a_cpy, cond, &notifier);
12   wait(cond);
13
14   if (notifier == LOCK_THRD) {
15     c_thrd.kill();
16     free(a_cpy);
17   } else {
18     l_thrd.kill();
19     if (lock.held()) lock.release();
20     delete a;
21     a = a_cpy;
22     Localized[a] = true;
23   }
24 }
```
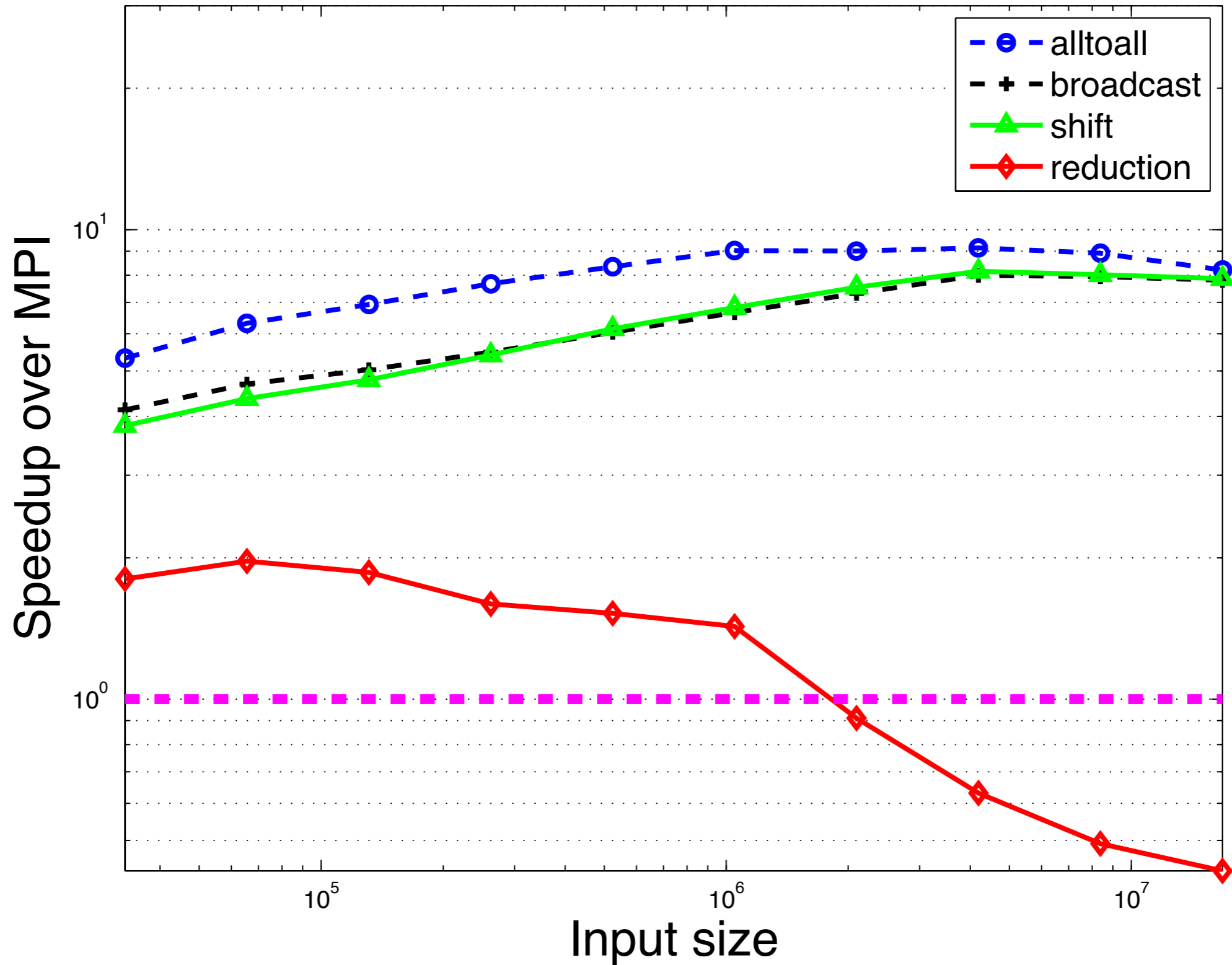
# Experimental Evaluation

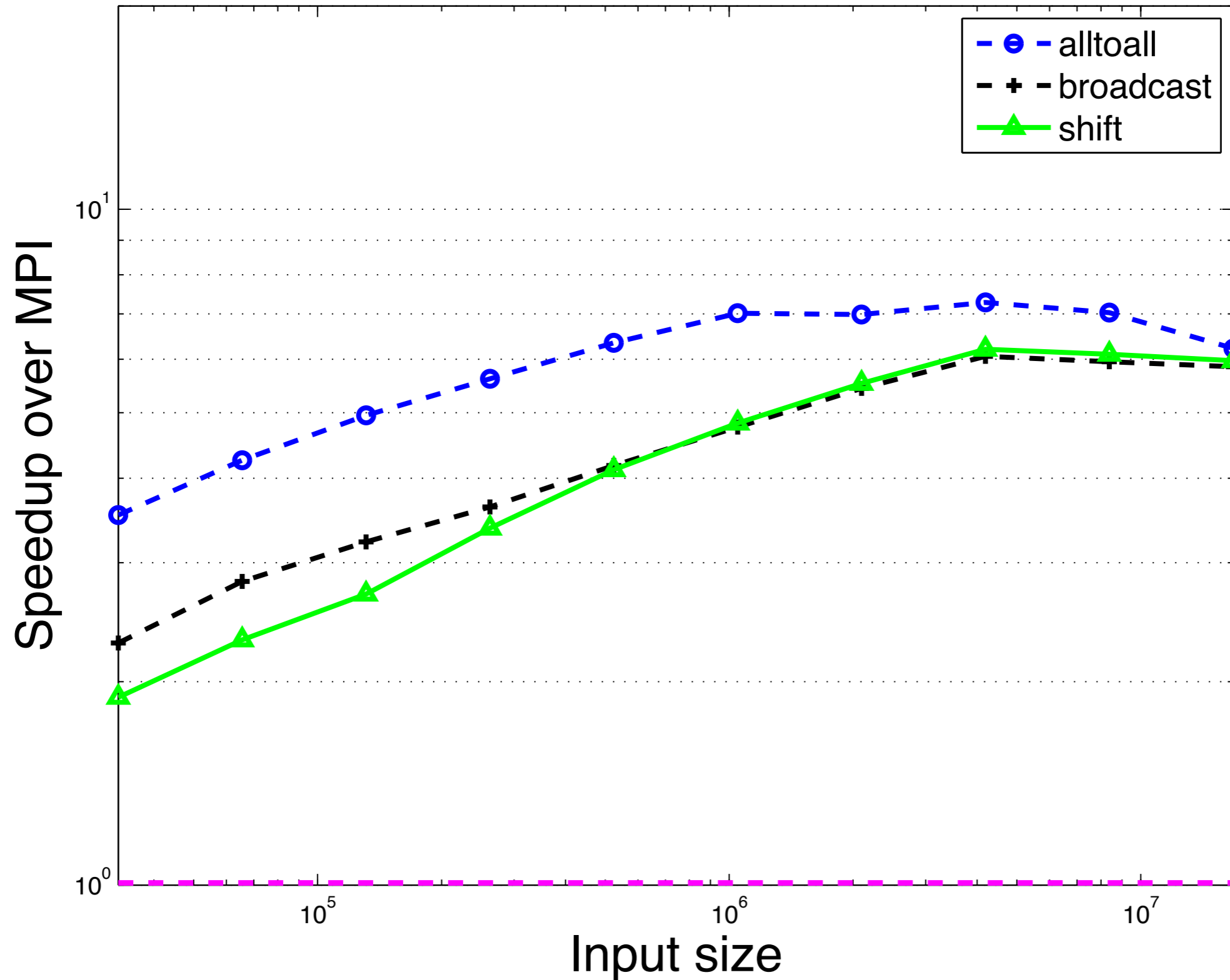| Op | Kanor | MPI | Shared Memory |
|---|---|---|---|
| all | `A[j]@i <<= A[i]@j`<br>`  where i,j in WORLD` | `MPI_Alltoall (...)` | `barrier();` |
| b'cast | `A@i <<= A@0`<br>`  where i in WORLD` | `MPI_Bcast(A, ..., ..., 0, ...);` | `barrier();` |
| shift | `A@i <<= A@i+1` | `if (Rank == (numprocs - 1)) dest = 0;`<br>`else dest = Rank + 1;`<br>`MPI_Send(A, array_size, ...);`<br>`MPI_Recv(A, array_size, ...);` | `barrier();` |
| reduce | `A@0 <<op<< A@i`<br>`  where i in WORLD` | `MPI_Reduce (...)`<br>`// or specialized code for`<br>`// tree-reduction of ``op''` | `// loop for tree-reduction`<br>`for (i ...) {`<br>`    A[i] = op(..);`<br>`}` |

- 8-core AMD Opteron, Gentoo Linux, OpenMPI 1.4.3



Case 1
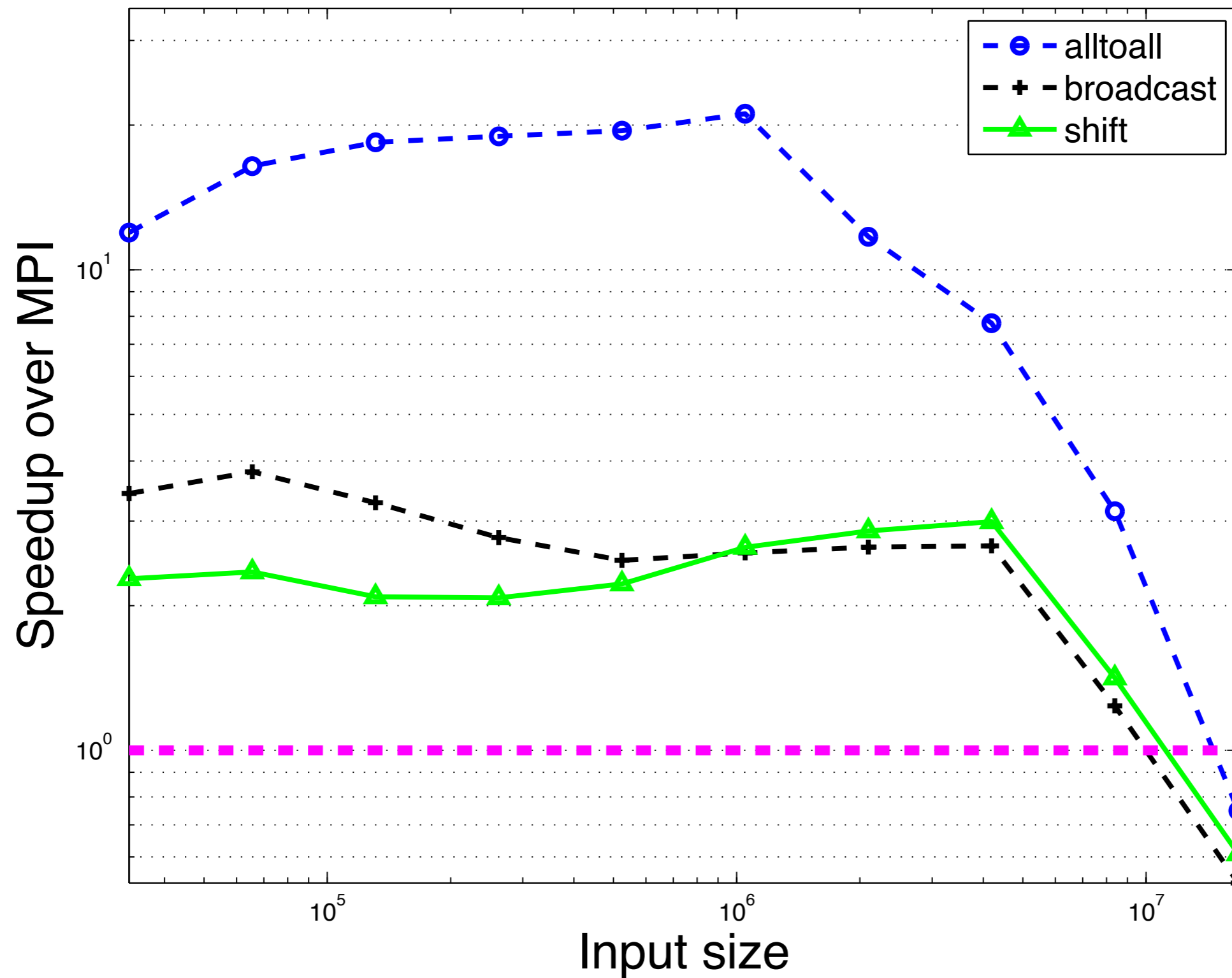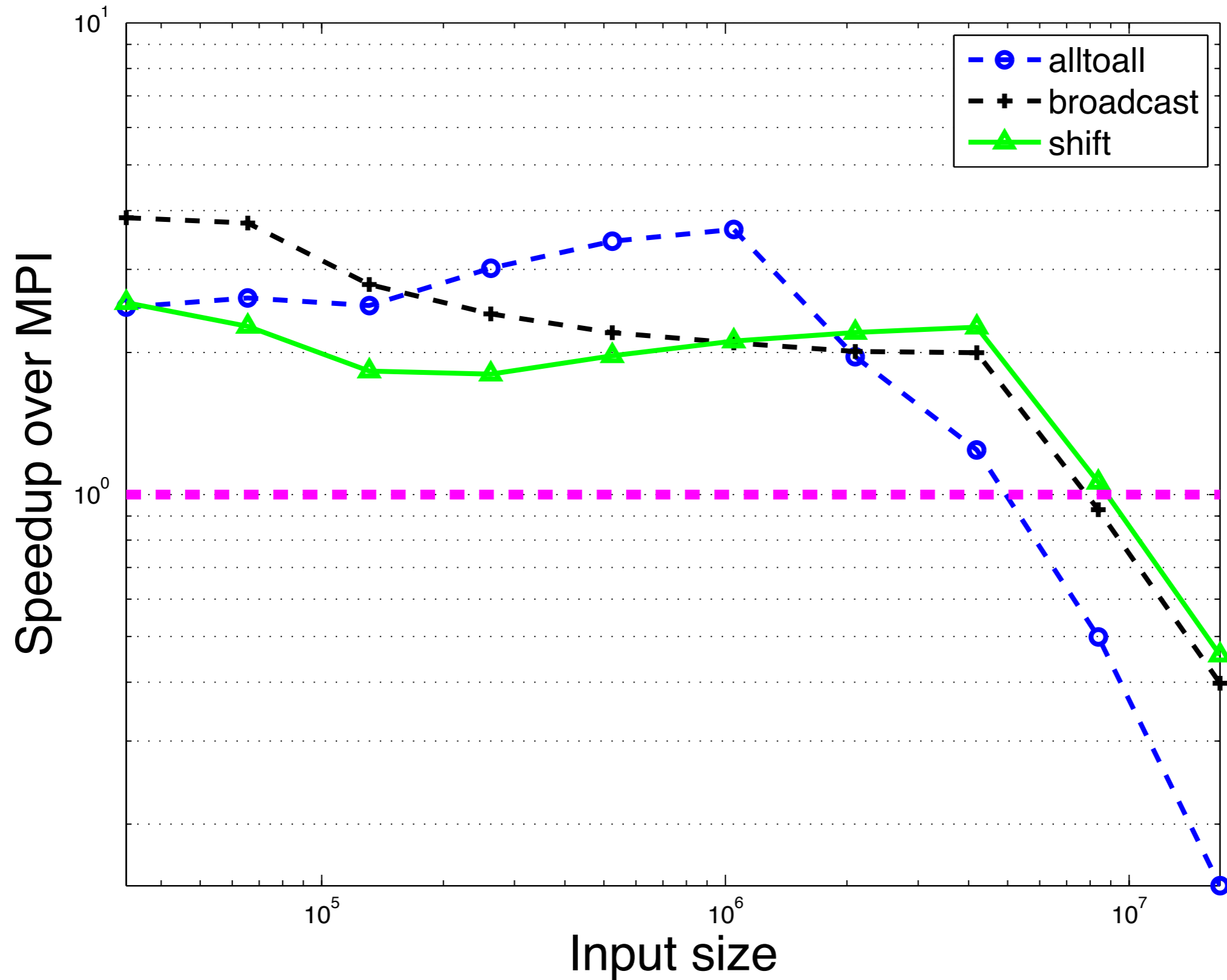
Case 2a

*Arun Ch*

# Case 1

# Case 2a

# Case 2b

# Case 3

# Concluding Remarks

- Parallel programming with partitioned address spaces has advantages

- Appropriate abstraction makes parallel programming more accessible to intermediate-level programmers

  - Kanor demonstrates the effectiveness of this approach

- Advantages of shared memory can be obtained through compiler optimizations

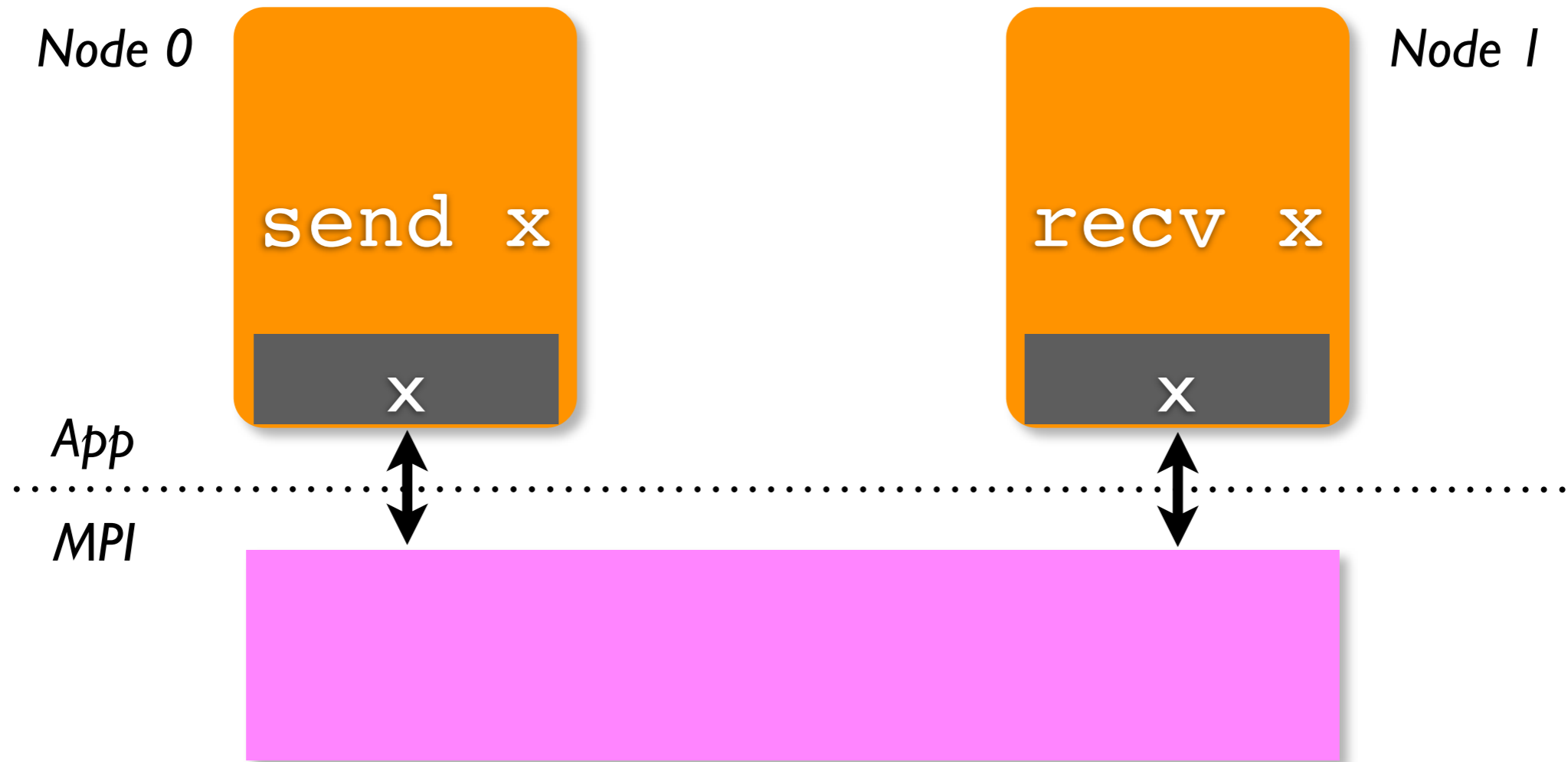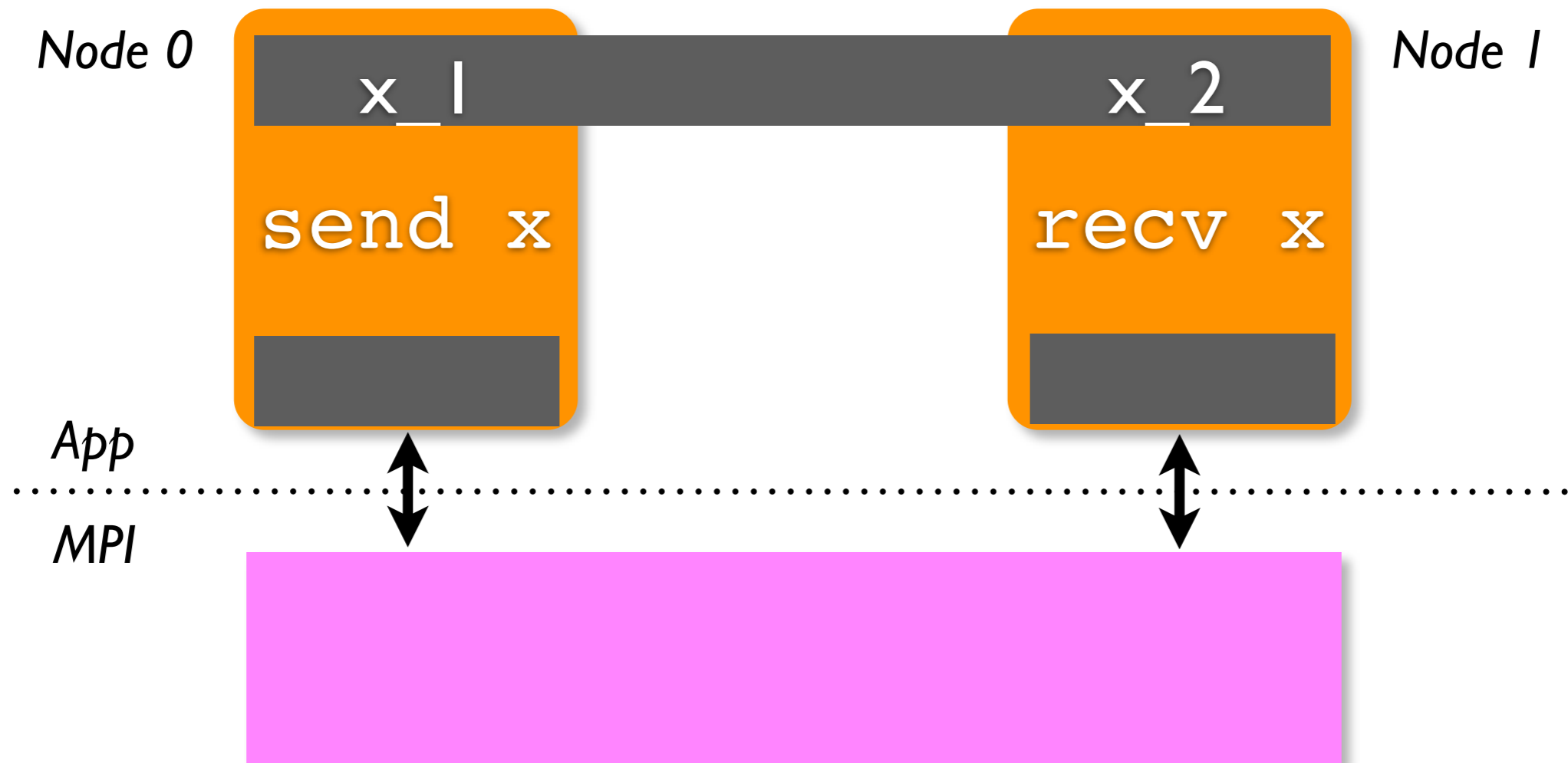  - our compiler algorithms and experimental evaluation substantiate this claim

# End

# Optimizing for Shared Memory

`@communicate {x@0 <<= x@1}`

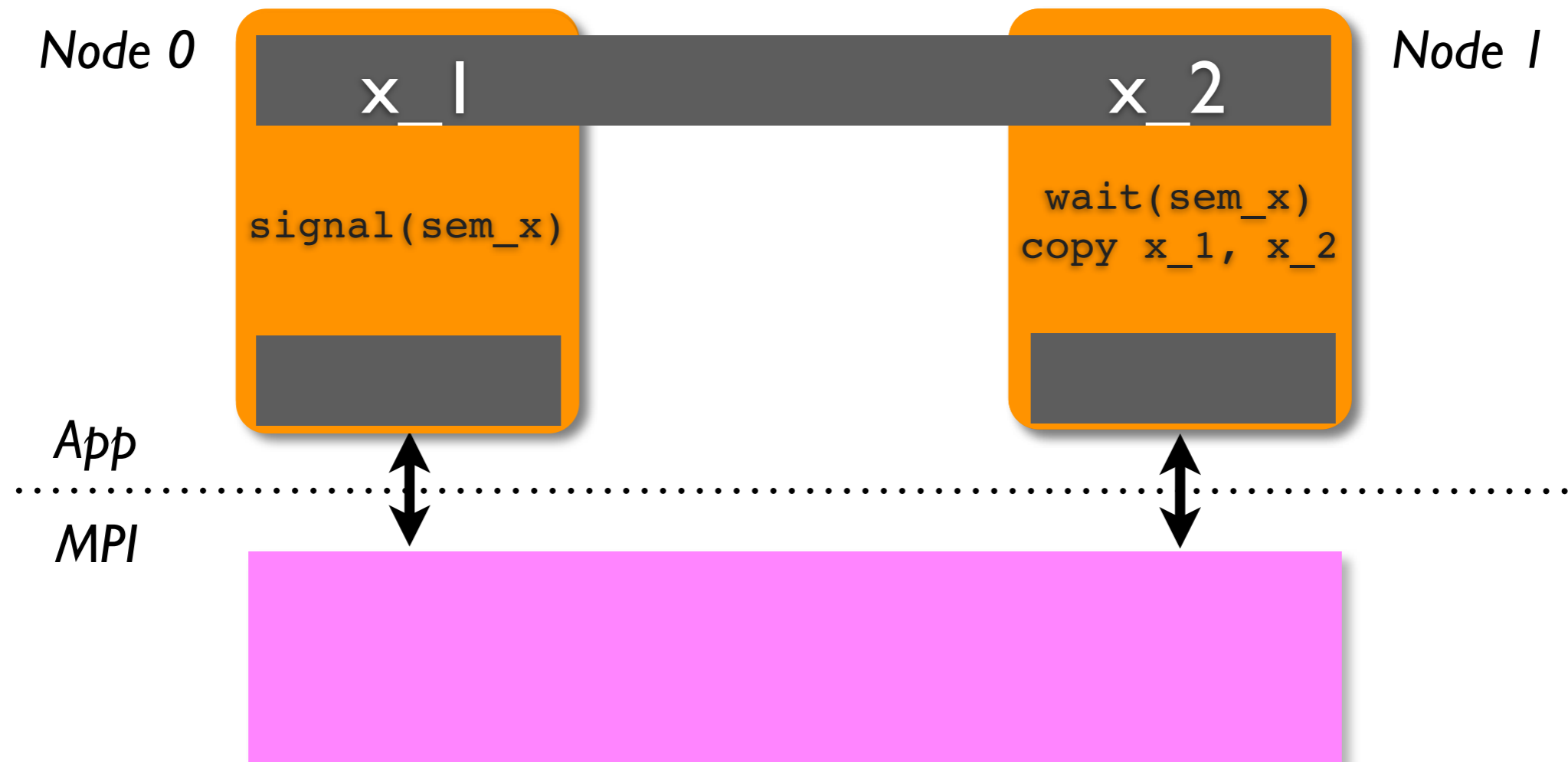# Optimizing for Shared Memory

`@communicate {x@0 <<= x@1}`
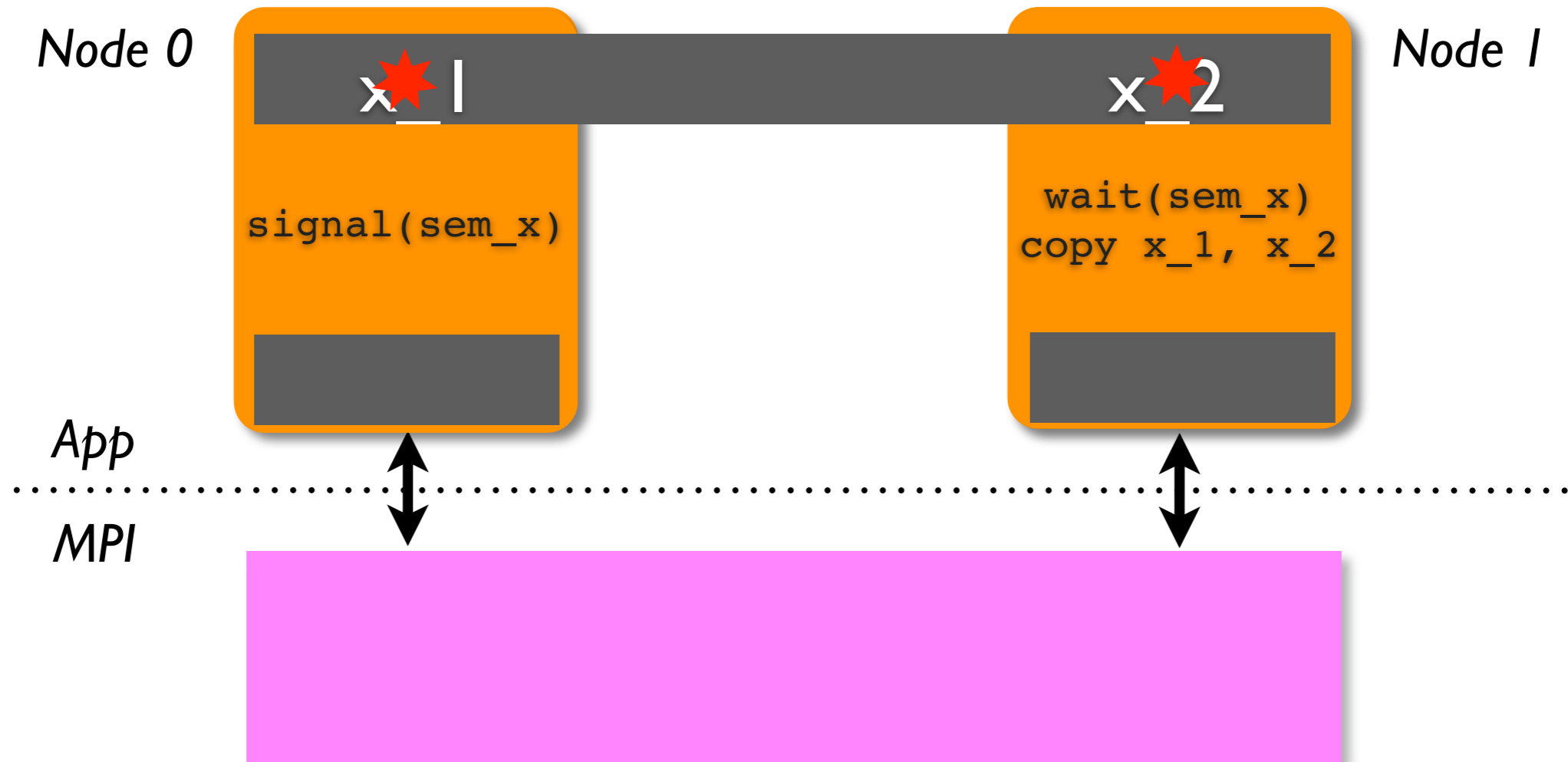
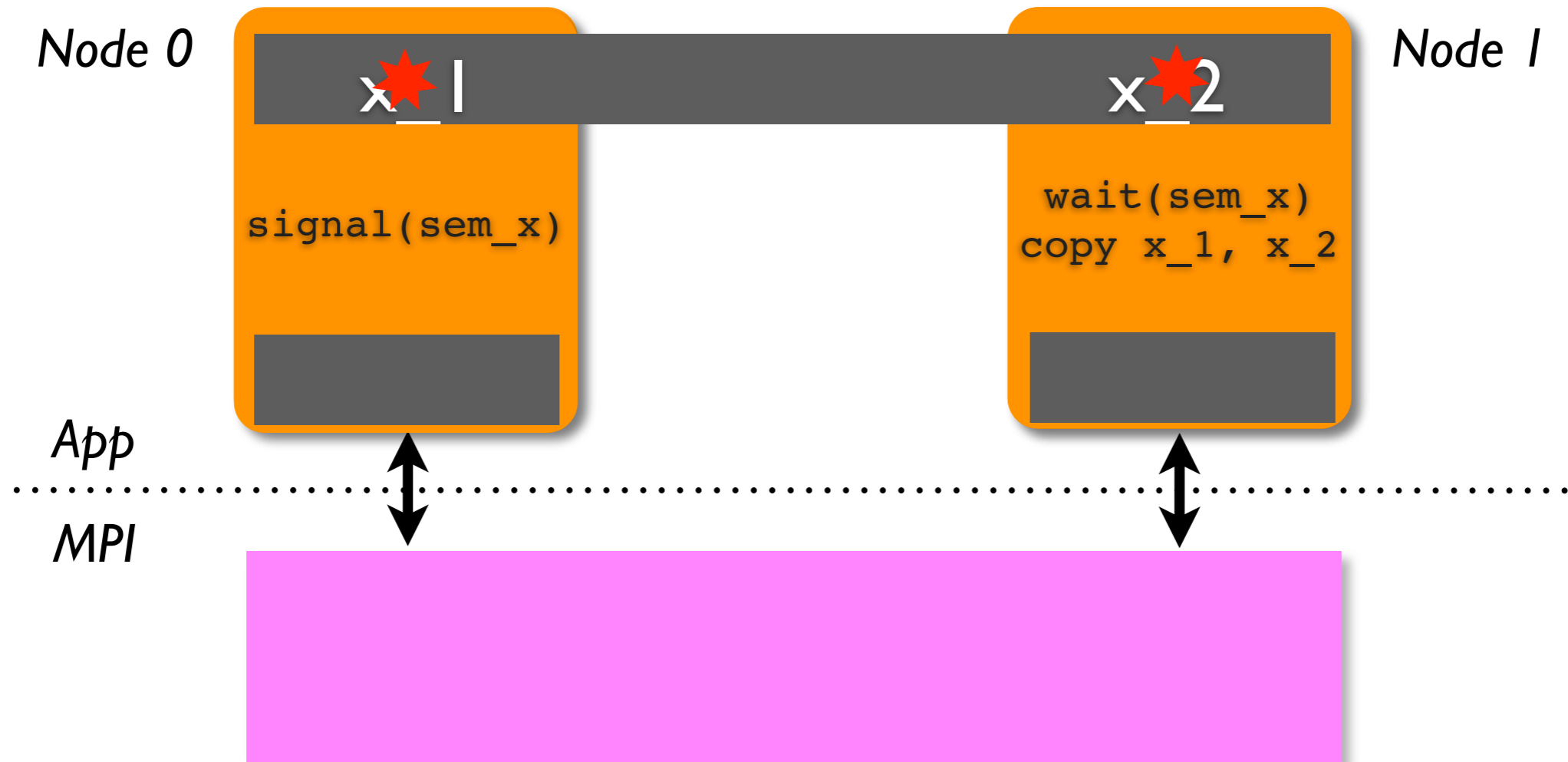# Optimizing for Shared Memory

## @communicate {x@0 <<= x@1}

# Optimizing for Shared Memory

## @communicate {x@0 <<= x@1}

# Optimizing for Shared Memory

## @communicate {x@0 <<= x@1}

# Optimizing for Shared Memory

## `@communicate {x@0 <<= x@1}`

*Node 0*

`x_1`

`signal(sem_x)`

*Node 1*

`x_2`

`wait(sem_x)`
`copy x_1, x_2`

*App*

*MPI*

## 1 copy
*(requires compiler intervention)*

# Computing all Paths from s to t

1 **Algorithm:** PATHS

2 **Input**: Directed graph $G(V, E)$
   Start node $s$
   End node $t$

3 **Output**: Set $P$ of nodes that lie on any path from $s$ to $t$

---

4 $P \leftarrow \phi$

5 **for** *each node $n$ in $G$* **do**

6    $n$.color $\leftarrow$ "white"

7 $Q \leftarrow [s]$

8 **while** *not $Q$.empty* **do**

9    $q \leftarrow Q$.extract

10    **for** *each edge $(q, v) \in E$* **do**

11       **if** *$v$.color $\neq$ "red"* **then**

12          $v$.color $\leftarrow$ "red"

13          $Q$.add($v$)

14 $Q \leftarrow [t]$

15 **while** *not $Q$.empty* **do**

16    $q \leftarrow Q$.extract

17    **for** *each edge $(v, q) \in E$* **do**

18       **if** *$v$.color $\neq$ "black"* **then**

19          **if** *$v$.color = "red"* **then**

20             $P \leftarrow P \cup \{v\}$

21          $v$.color $\leftarrow$ "black"

22          $Q$.add($v$)

23 **return** $P$

# Computing Locking Sets

1    **Algorithm:** COMPUTE-LOCKING-SET

2    **Input**: CFG $G(V, E)$ of code region over which variable x is globalized, with level-annotated nodes; dependence levels, $l_x$, for dependencies involving x; dep. distances, $d_x$, for dependencies involving x;

3    **Output**: Locking set $L$

---

4    $L = \phi$

5    **for** *each node pair* $(w, r)$ *with an entry in* $l_x$ **do**

6      **if** $d_x(w, r) = 0$ **then**

7        **if** $l_x(w, r) = 0$ **then**

8          $L \leftarrow L \cup \text{PATHS}(G, w, r)$

9        **else**

10          $G'(V', E') \leftarrow G$ without any looping back-edges at level $l_x(w, r)$ and lower

11          $L \leftarrow L \cup \text{PATHS}(G', w, r)$

12      **else if** $d_x(w, r) = 1$ **then**

13        $h \leftarrow$ head node of loop at level $l_x(w, r)$

14        $G'(V', E') \leftarrow G$ restricted to levels $l_x(w, r)$ and higher

15        $L \leftarrow L \cup \text{PATHS}(G', w, h) \cup \text{PATHS}(G', h, r)$

16      **else**

17        $G'(V', E') \leftarrow G$ restricted to levels $l_x(w, r)$ and higher

18        $L \leftarrow L \cup \text{PATHS}(G', w, r)$

19    **return** $L$