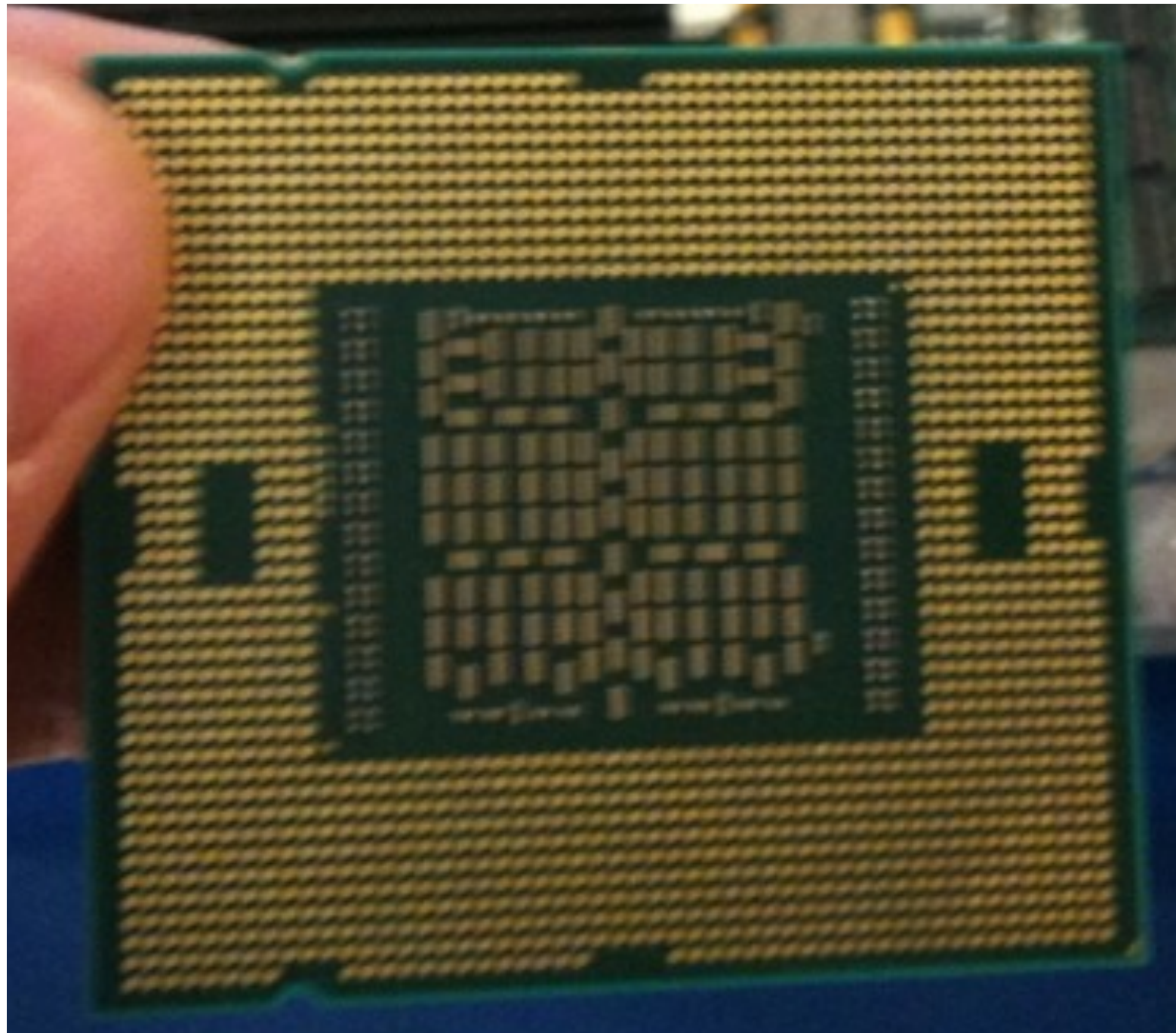


X10 on the Single-chip Cloud Computer (SCC)

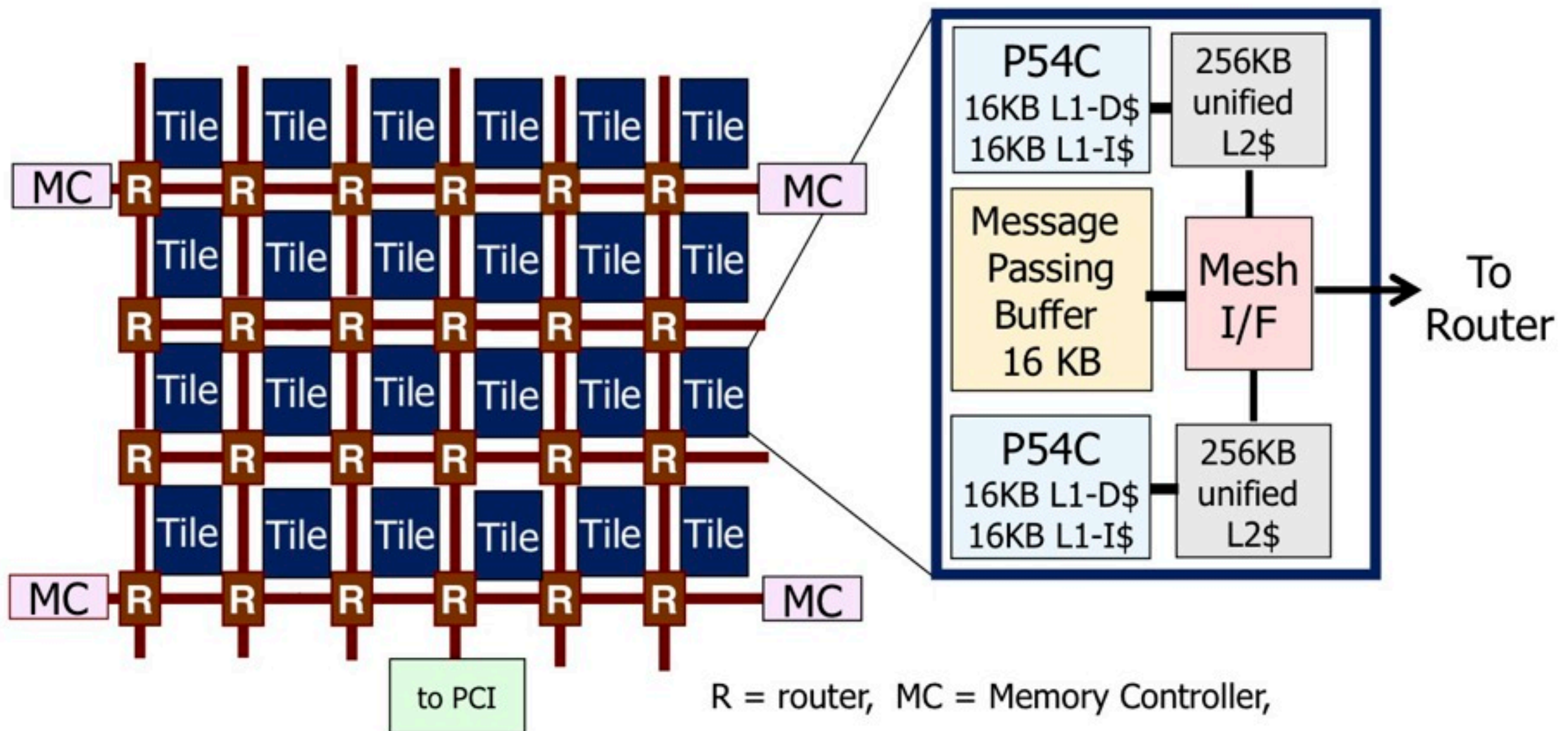
Keith Chapman, Ahmed Hussein, Antony Hosking



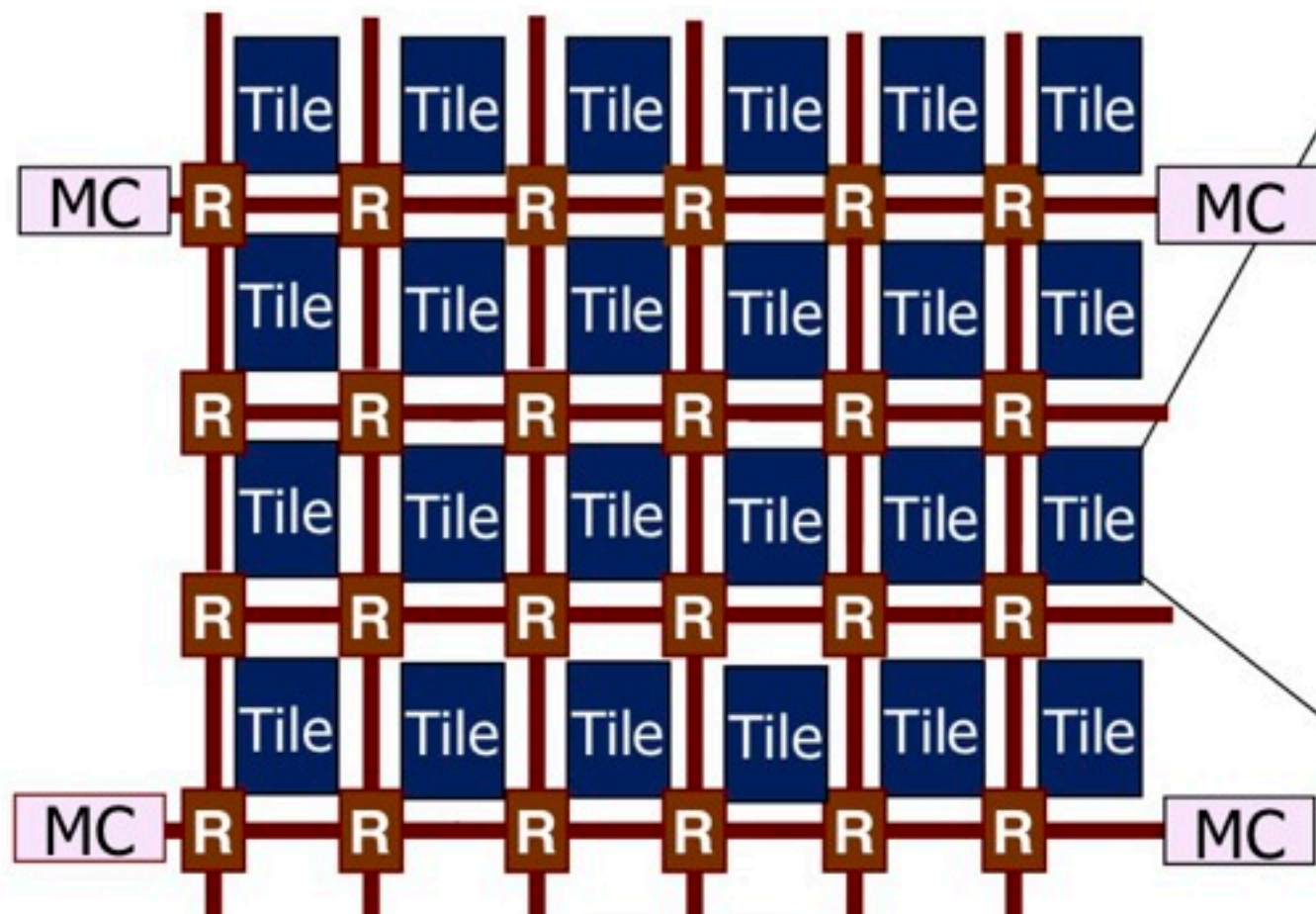
What is the SCC



What is the SCC



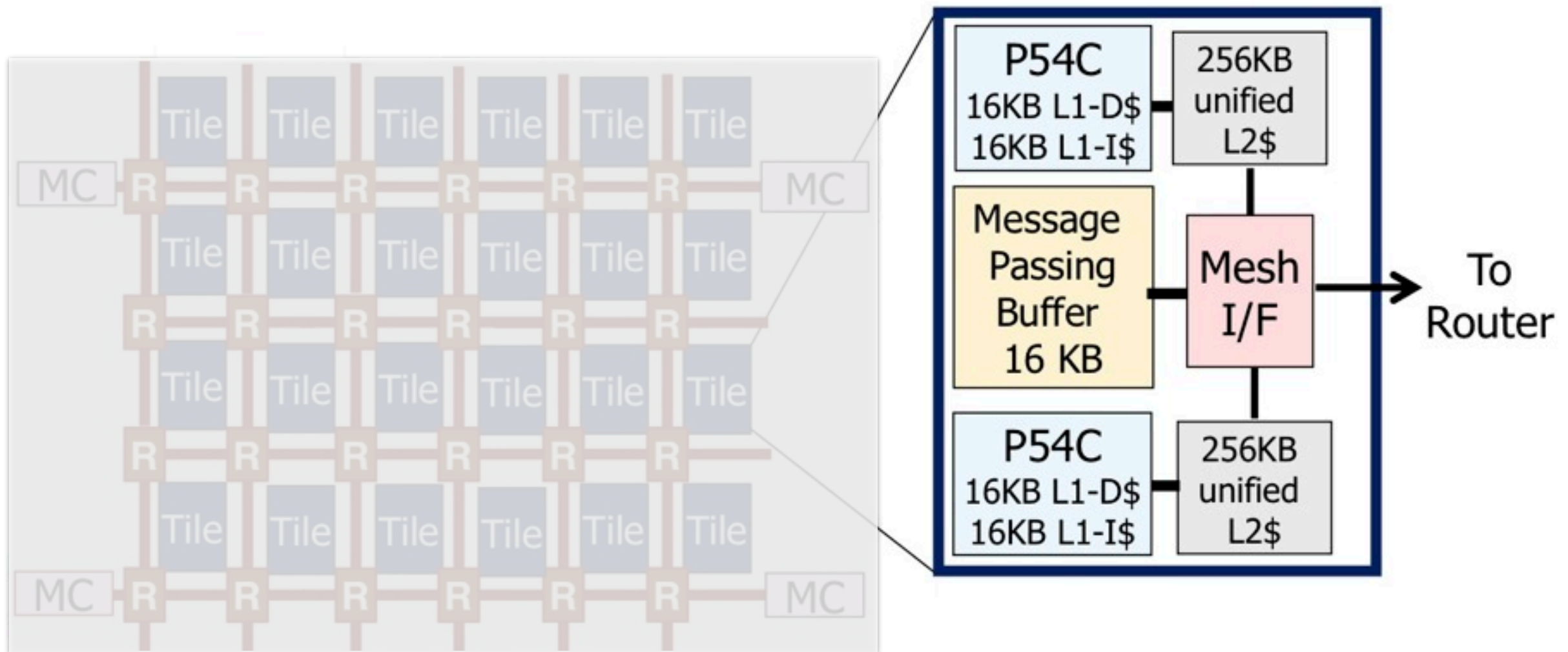
What is the SCC



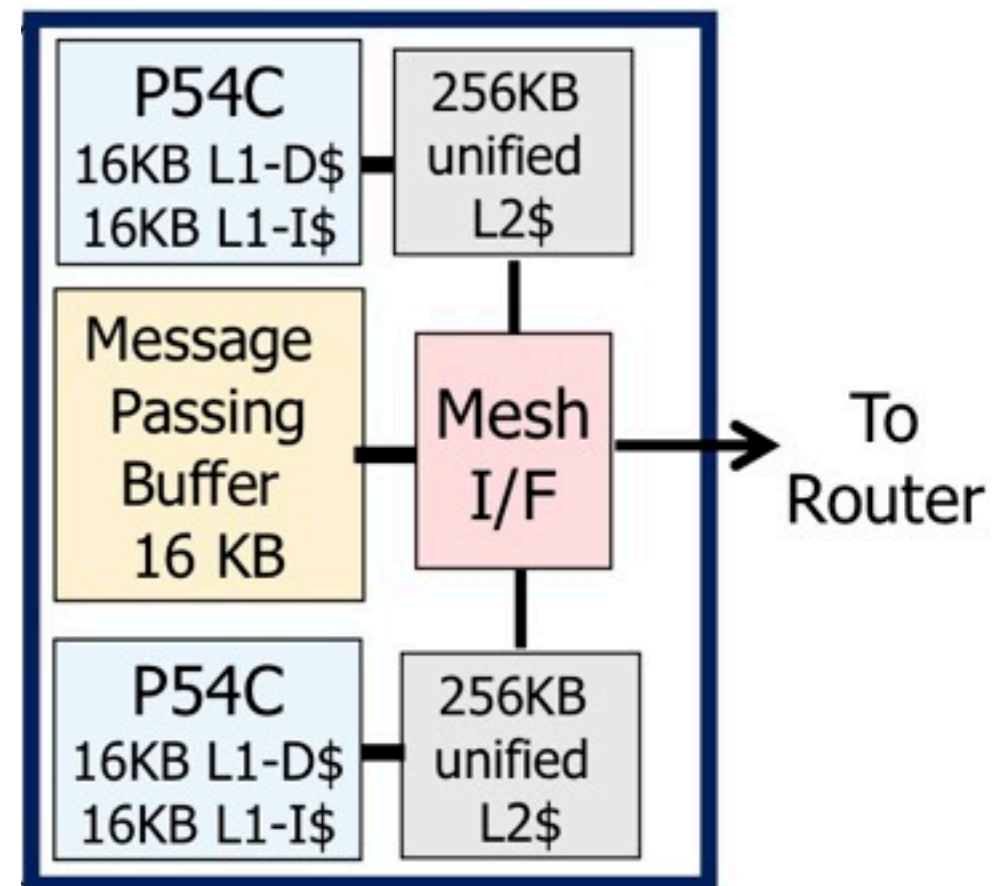
What is the SCC



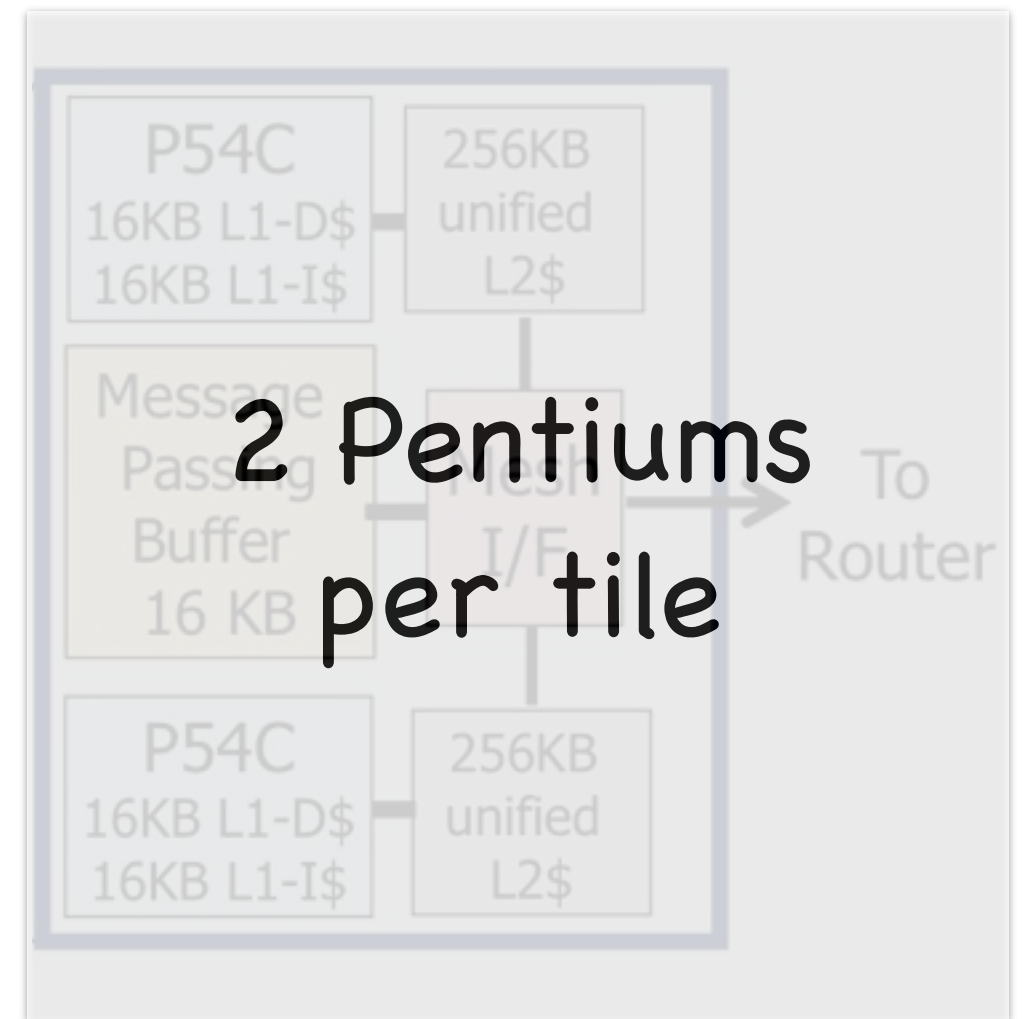
What is the SCC



What is the SCC

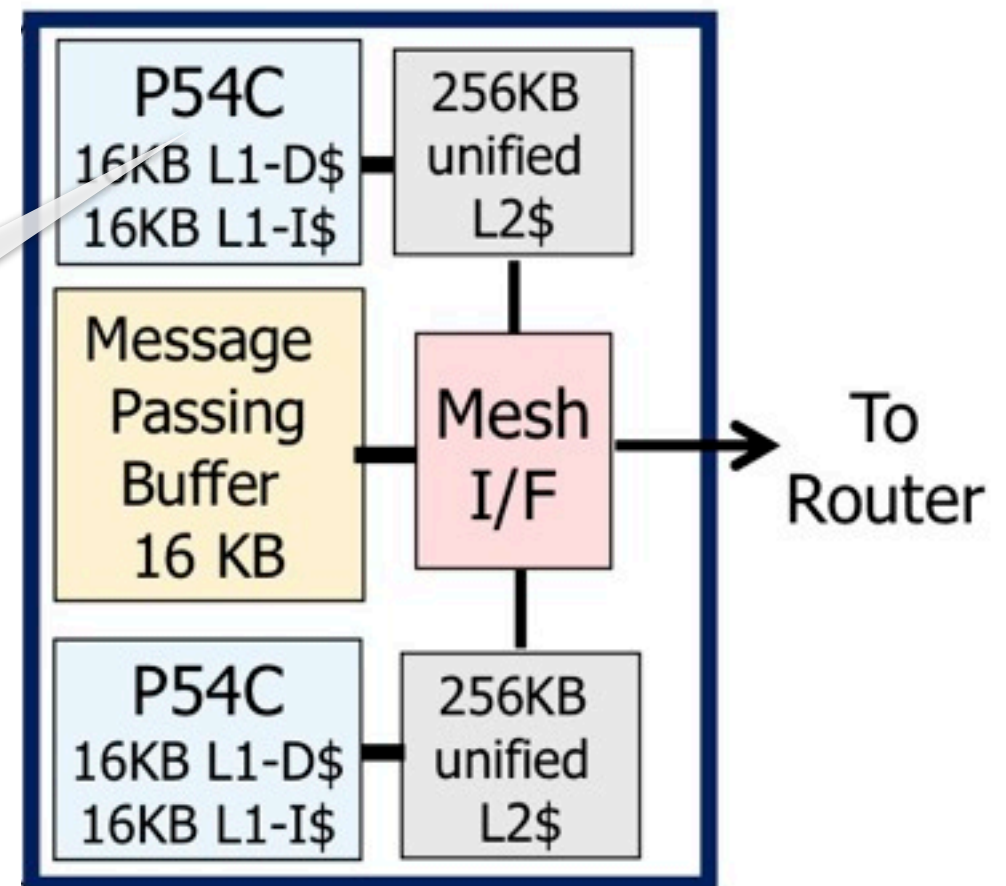


What is the SCC

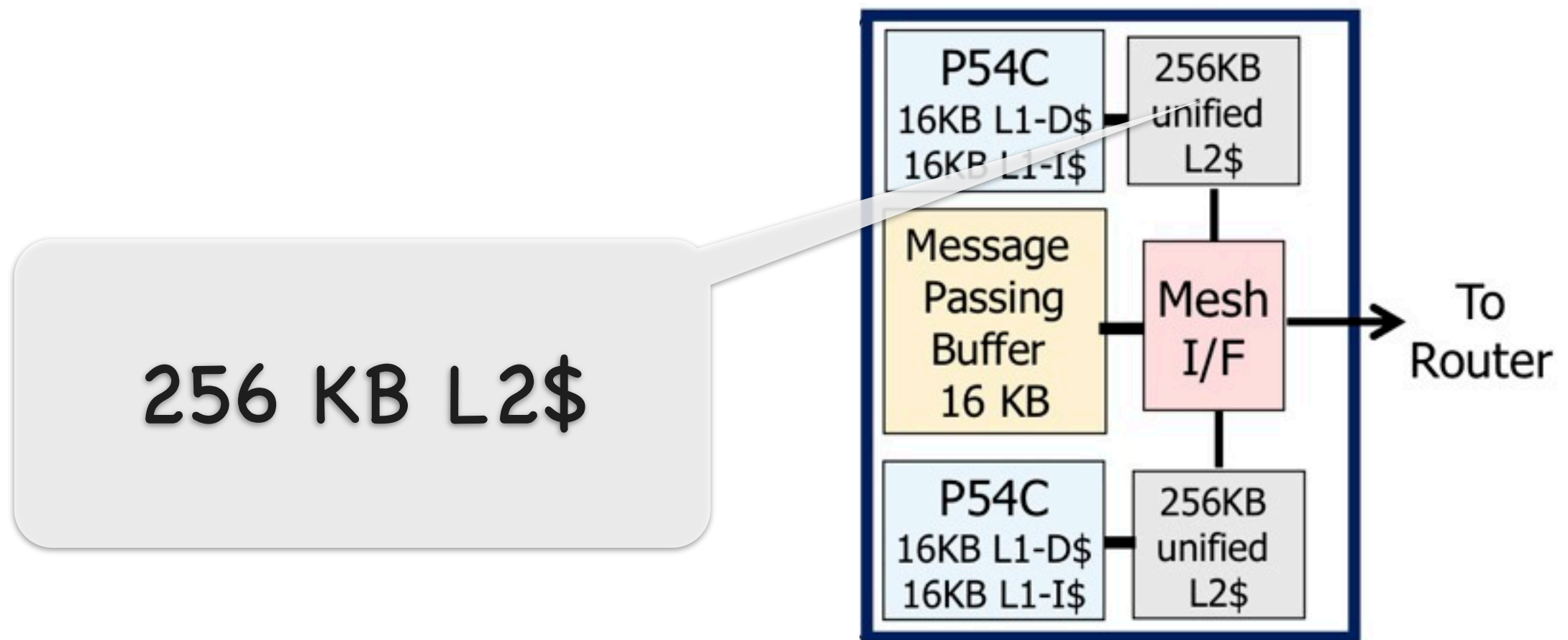


What is the SCC

16 KB L1 D\$ + I\$

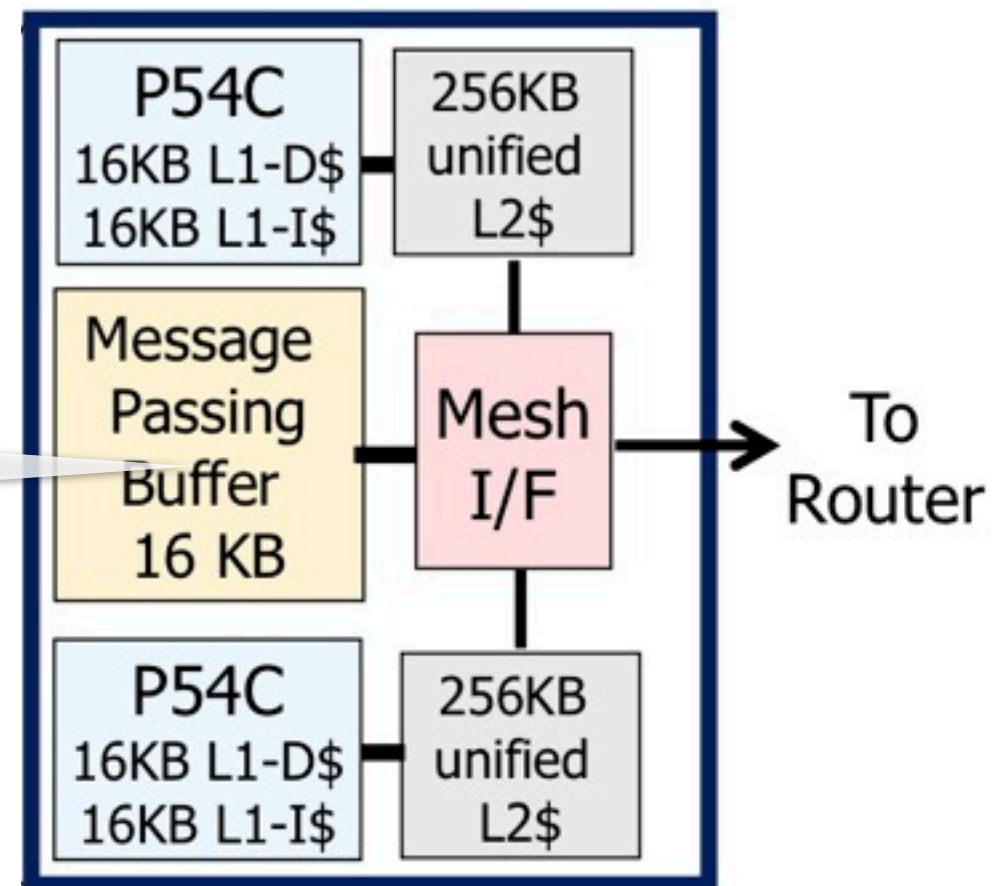


What is the SCC

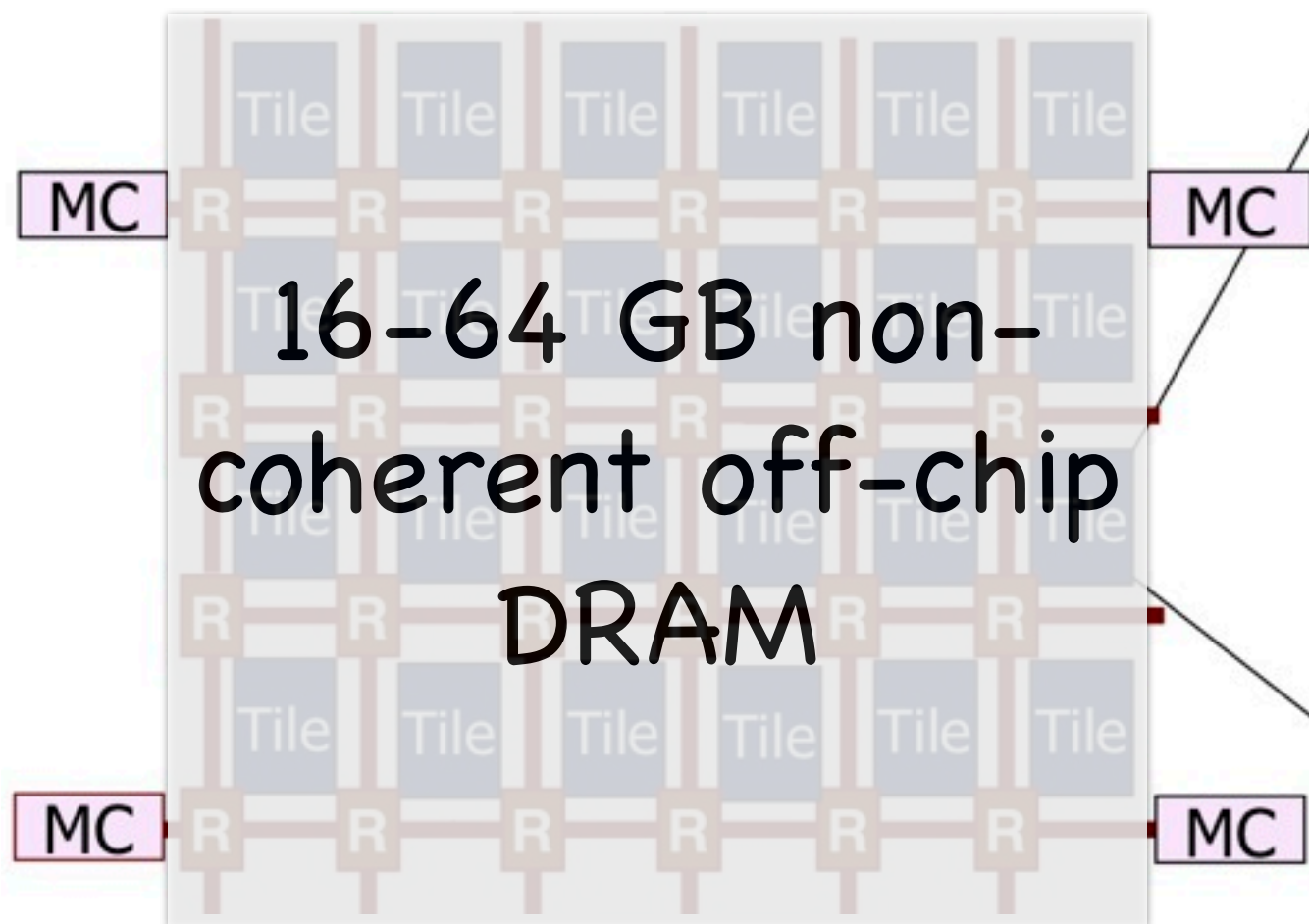


What is the SCC

16 KB message
passing buffer

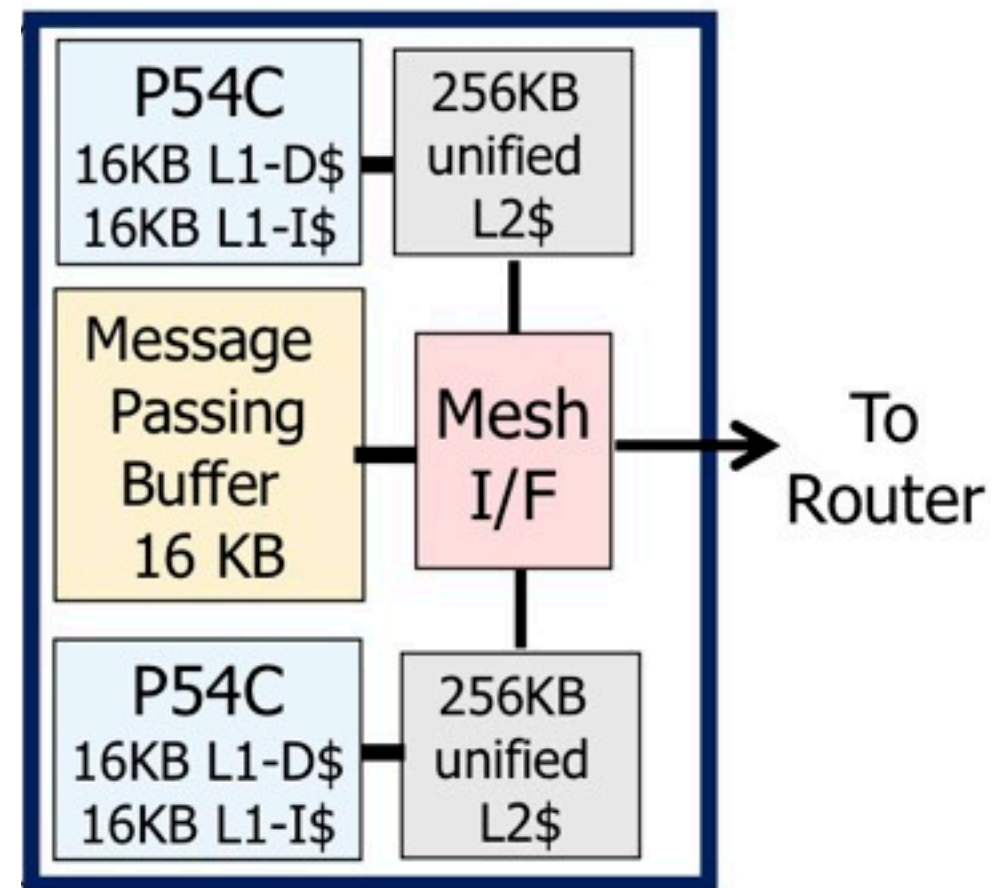


What is the SCC



Message Passing Buffer

- Shared on-chip SRAM
- Mapped by all cores
- Cached in L1
- Tagged for fast invalidate
- 8 KB per core
- 384 KB total accessible



RCCE native messaging

- Fast, thin layer over MPB
- One-sided *put/get* API
- RCCE flags allocatable in each MPB for coordination (one T&S register per core)
- Two-sided *send/rcv* built on one-sided API using RCCE Flags

XIO (C++) on the SCC

- Relatively straightforward
 - cross-compile using SCC tool chain
 - MPI or sockets over MPB out of the box
- What about RCCE?

XIORT API

- Used by XI0 runtime to communicate between *places*
- Send serialized object graphs
- Registered receiver callbacks for different types of message
- Probe occasionally for incoming messages
- Don't block waiting for messages to arrive

XIORT + MPI + SCC

- XIORT with MPI binding
- SCC MPI library with alternative channels
 - mpi-mpb: direct MPB
 - mpi-sock: sockets via mesh TCP/IP driver
 - mpi-shm: shared memory

mpi-mpb MPB management

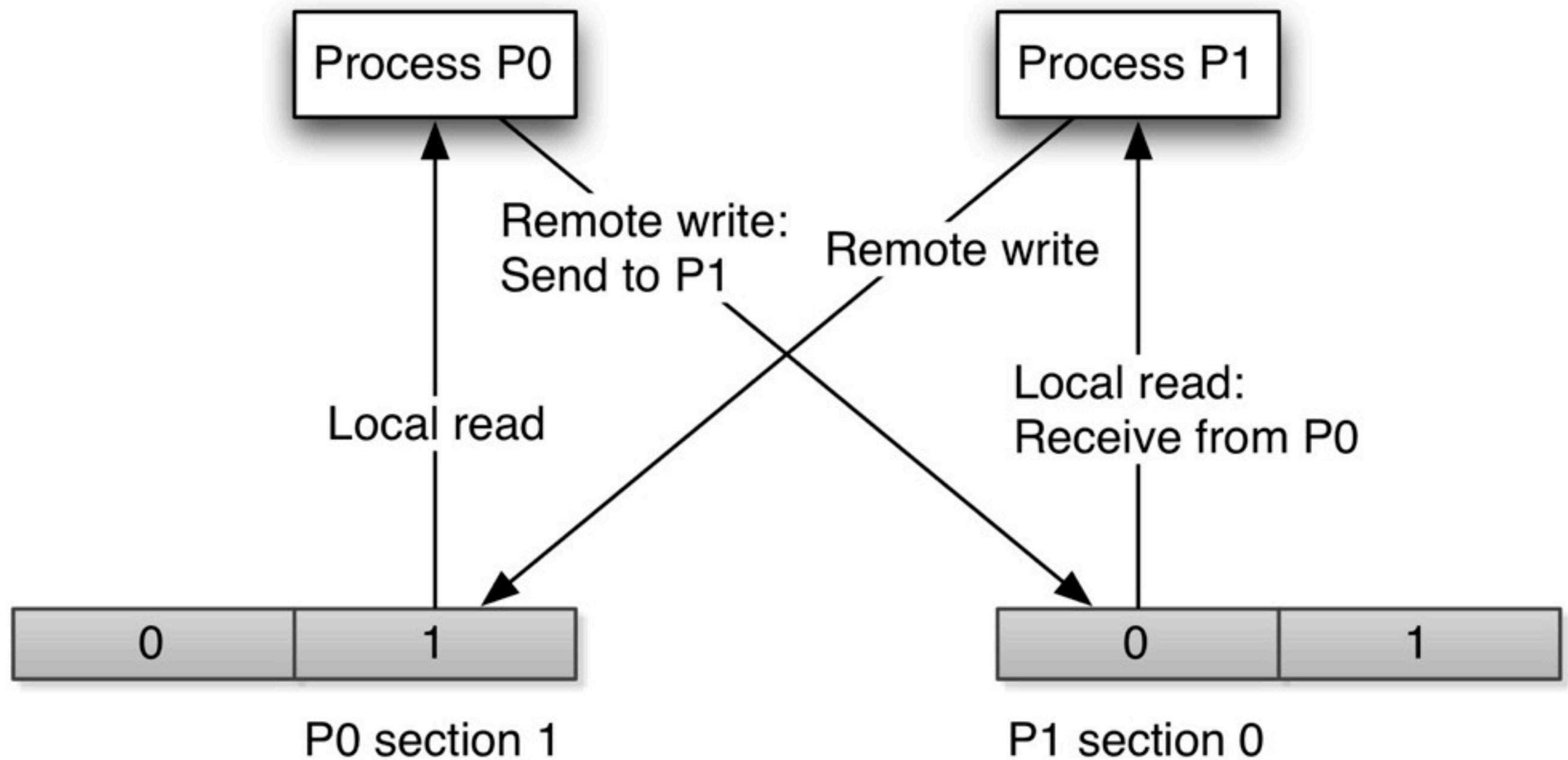
48 MPB areas of 8 Kbytes per core



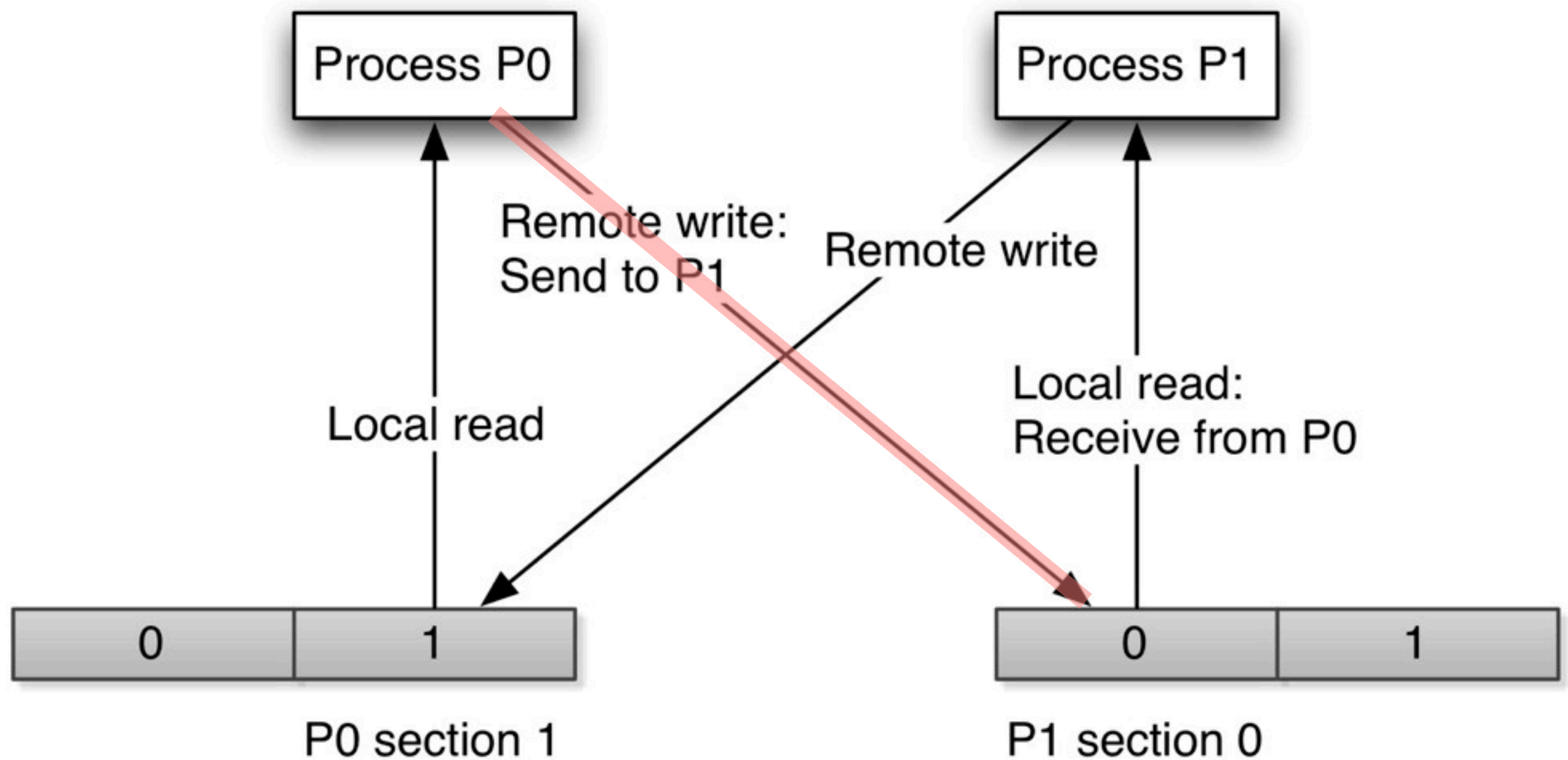
48 MPB write sections per core of 160 bytes



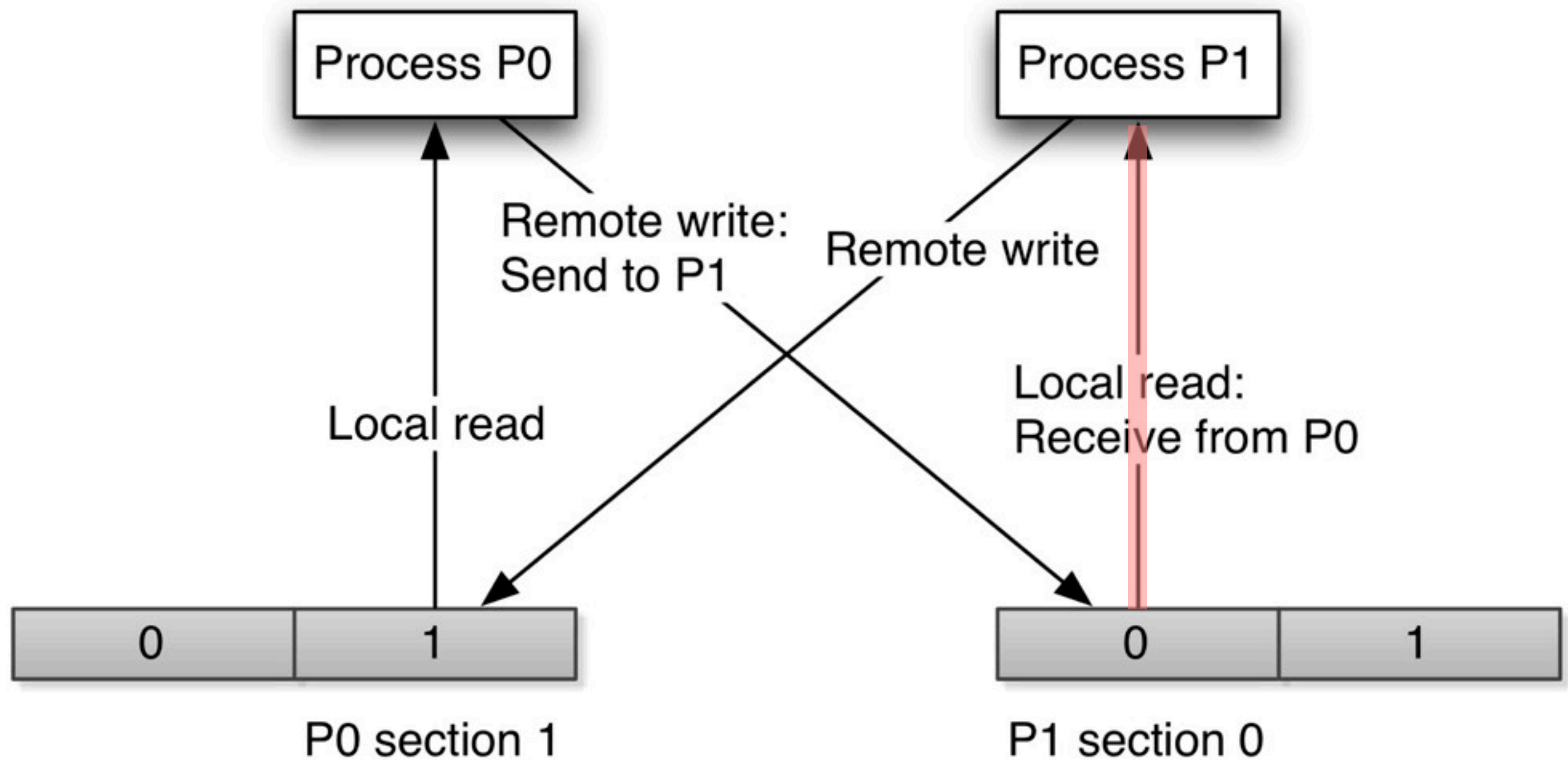
mpi-mpb send/receive



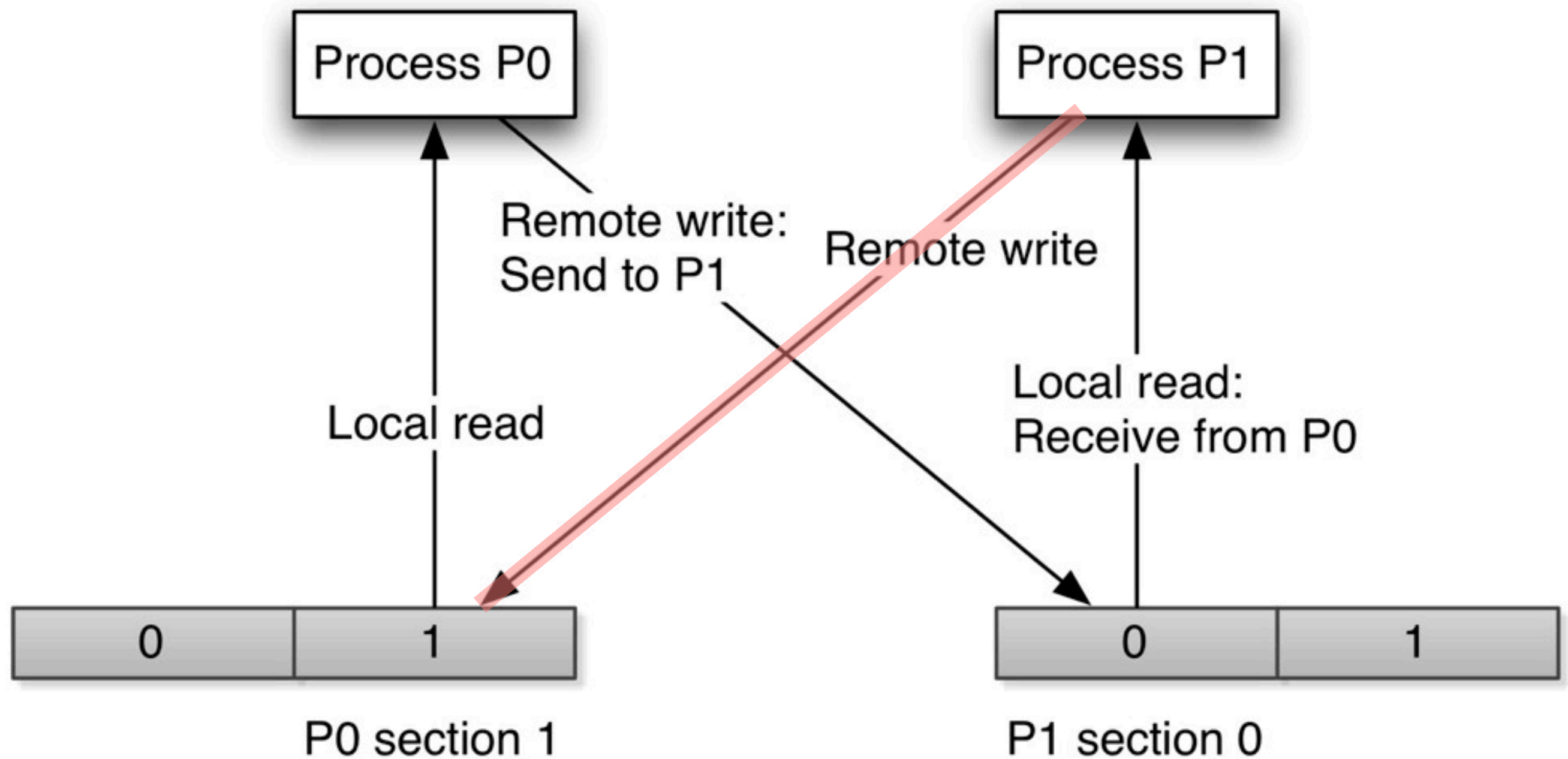
mpi-mpb send/receive



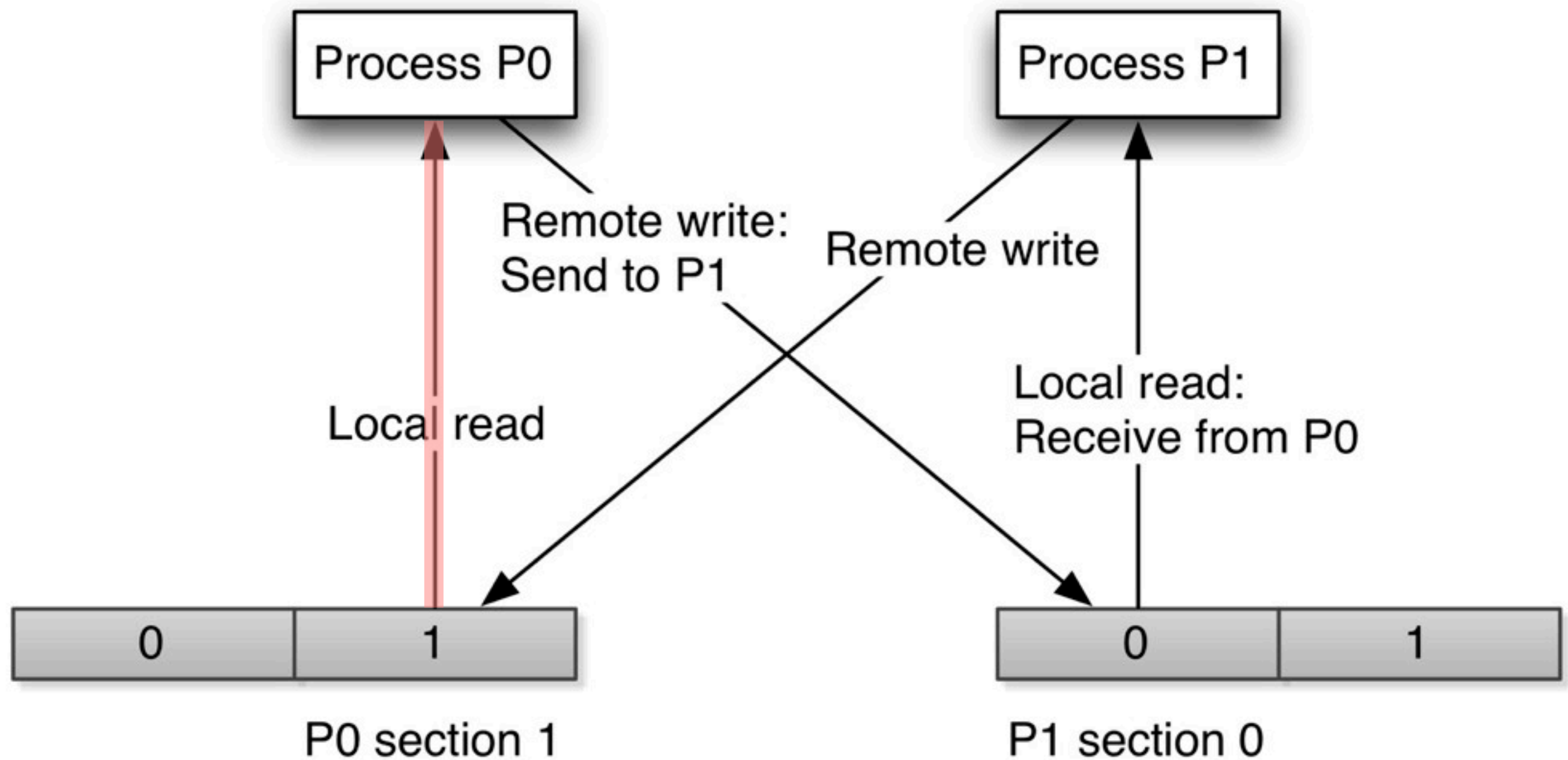
mpi-mpb send/receive



mpi-mpb send/receive



mpi-mpb send/receive



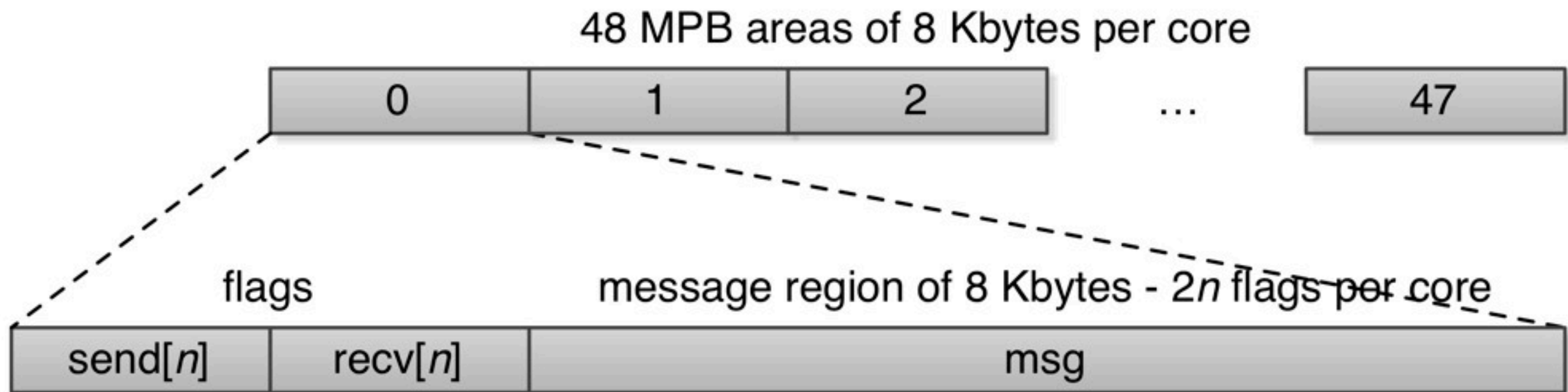
XIORT+“RCCE”+SCC

- RCCE send/recv are blocking
- RCCE non-blocking recv_test requires knowing sender to test
- RCCE does not support XIORT tags

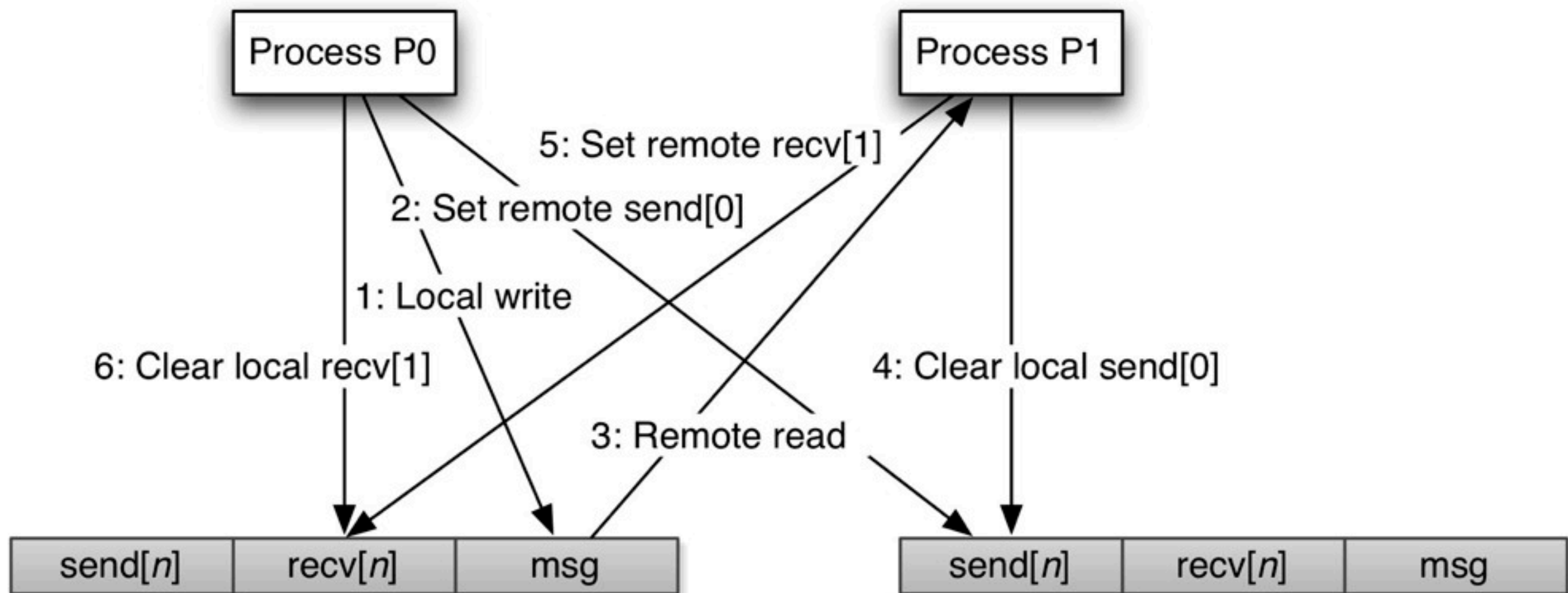
XIORT+“RCCE”+SCC

- Re-implement as RCCE-XIO
 - one-sided non-blocking send
 - non-blocking probe

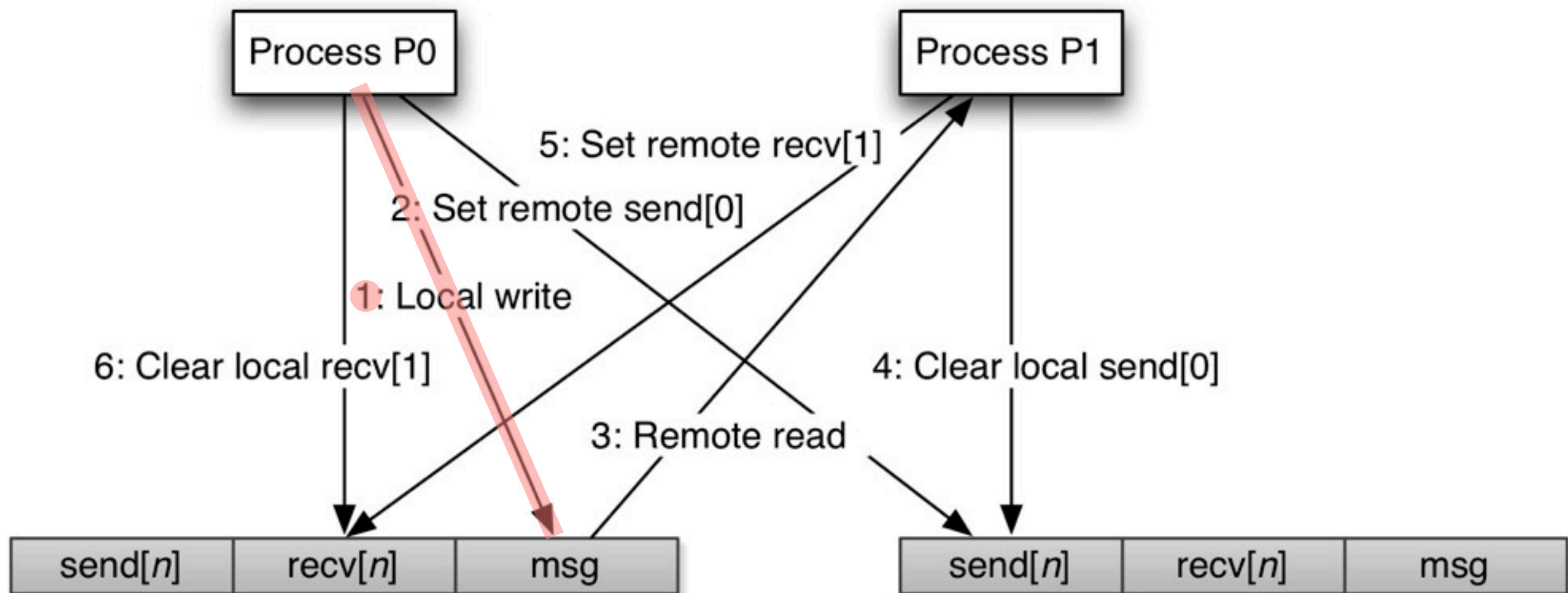
RCCE-XIO MPB Management



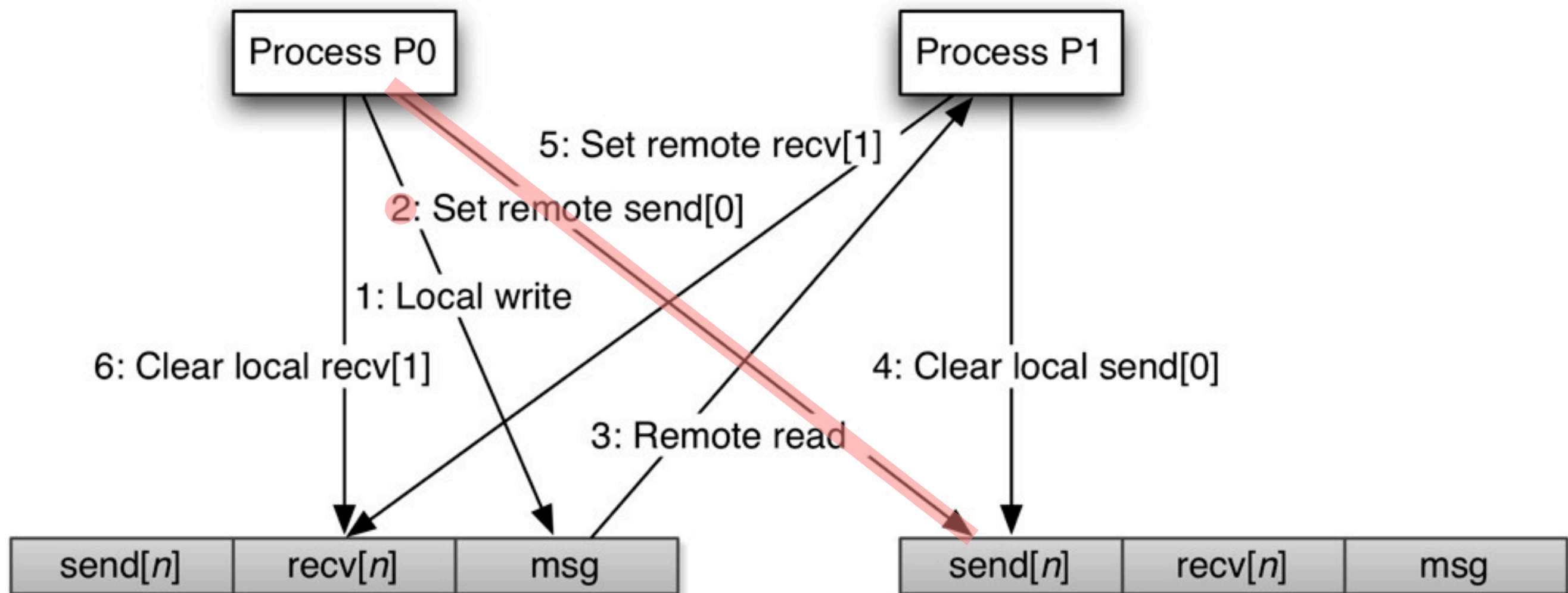
RCCE-XIO send/ receive



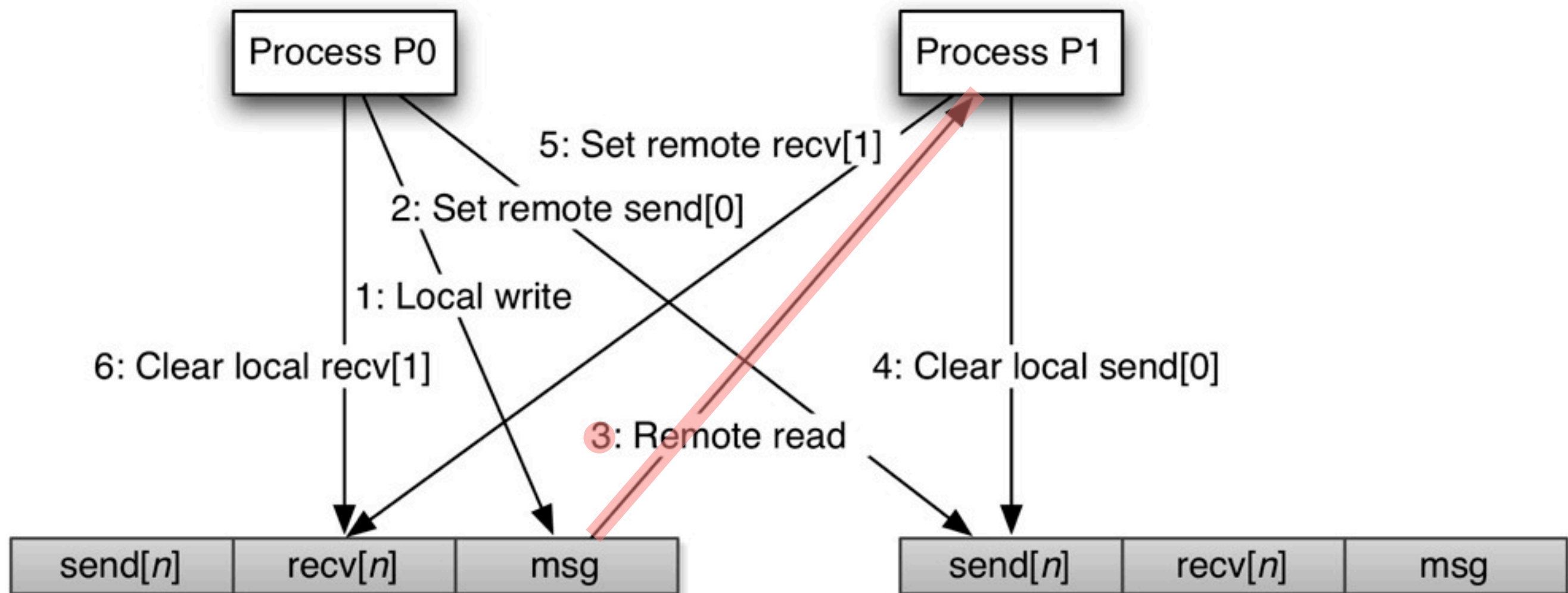
RCCE-XIO send/ receive



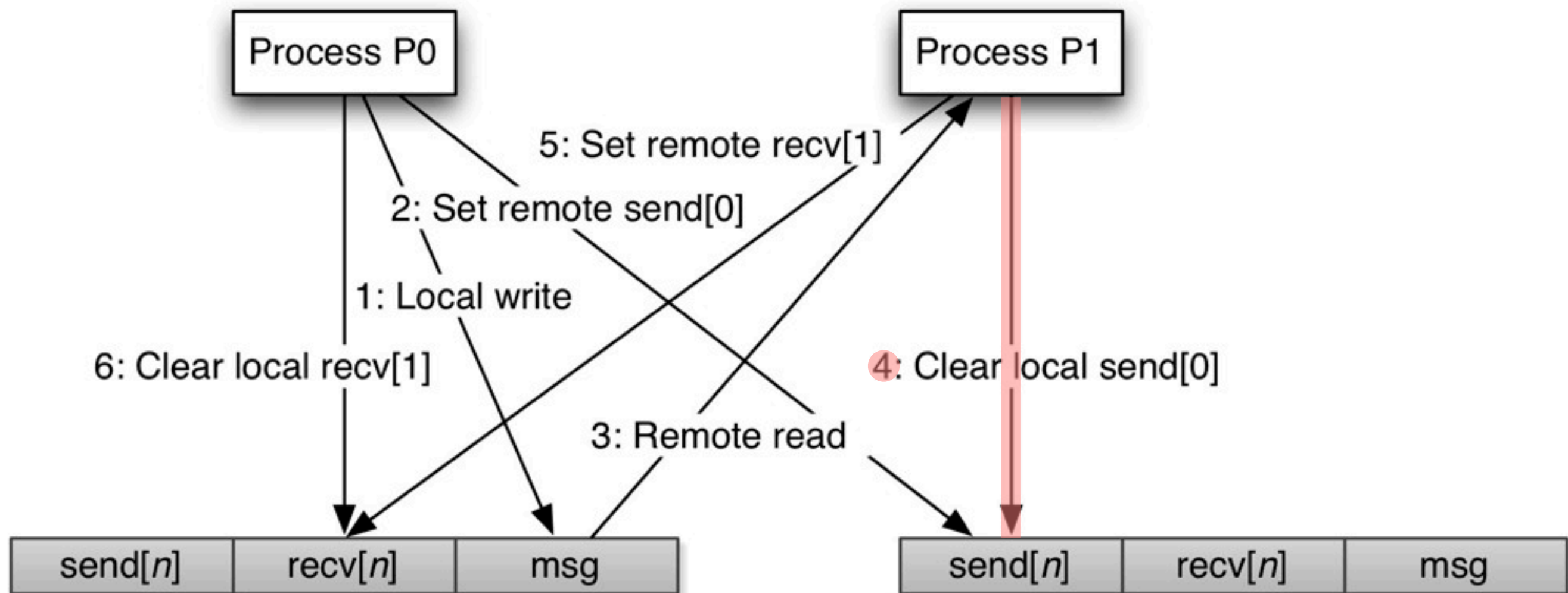
RCCE-XIO send/ receive



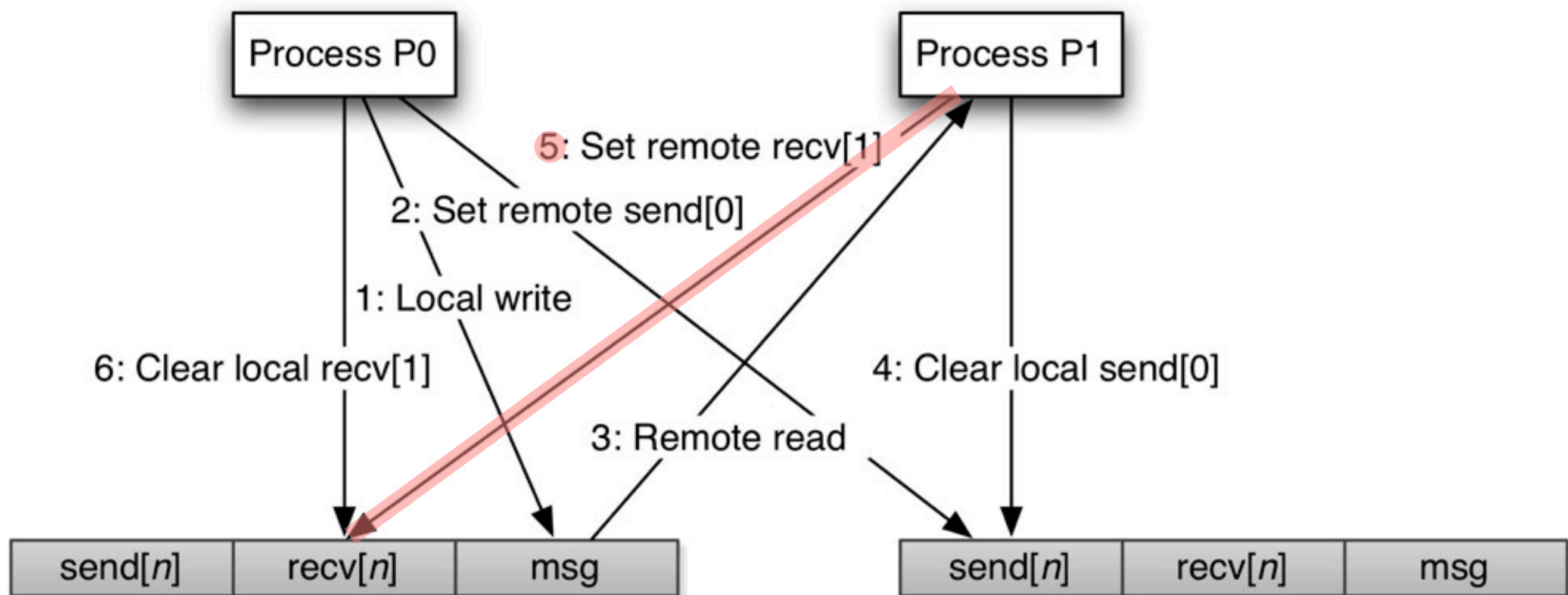
RCCE-XIO send/ receive



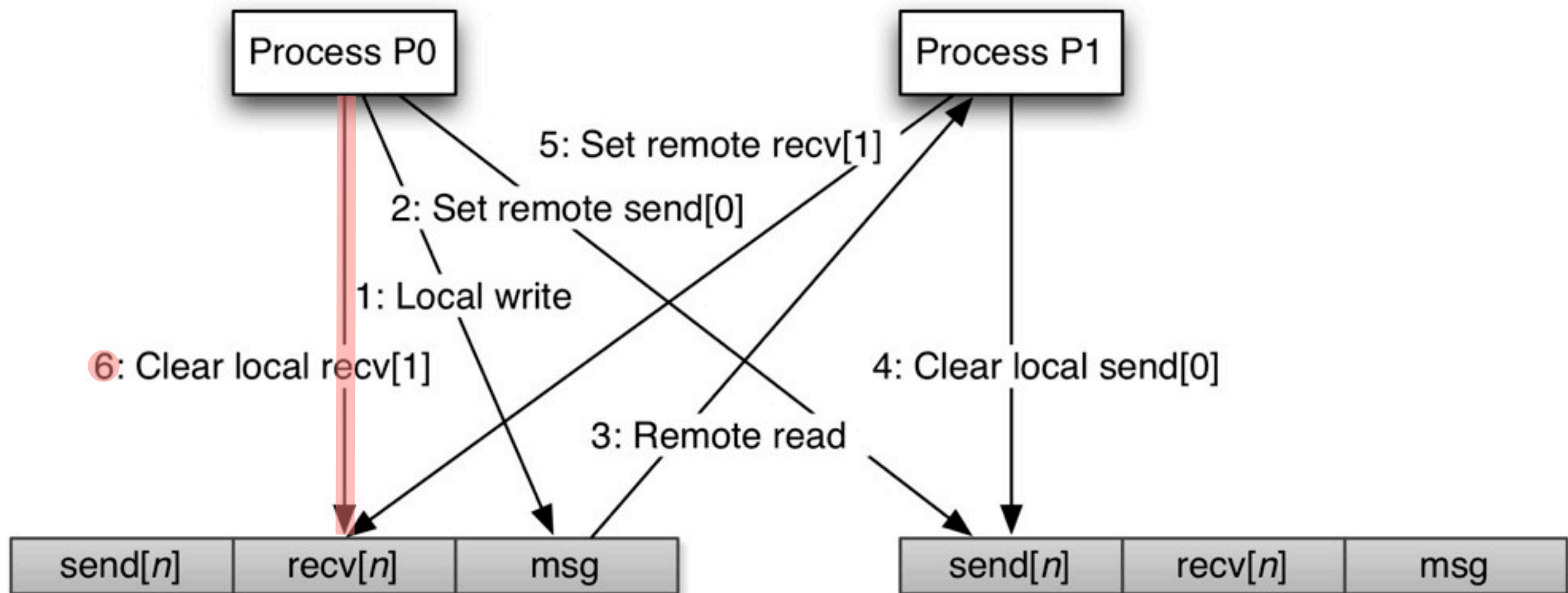
RCCE-XIO send/ receive



RCCE-XIO send/ receive



RCCE-XIO send/ receive



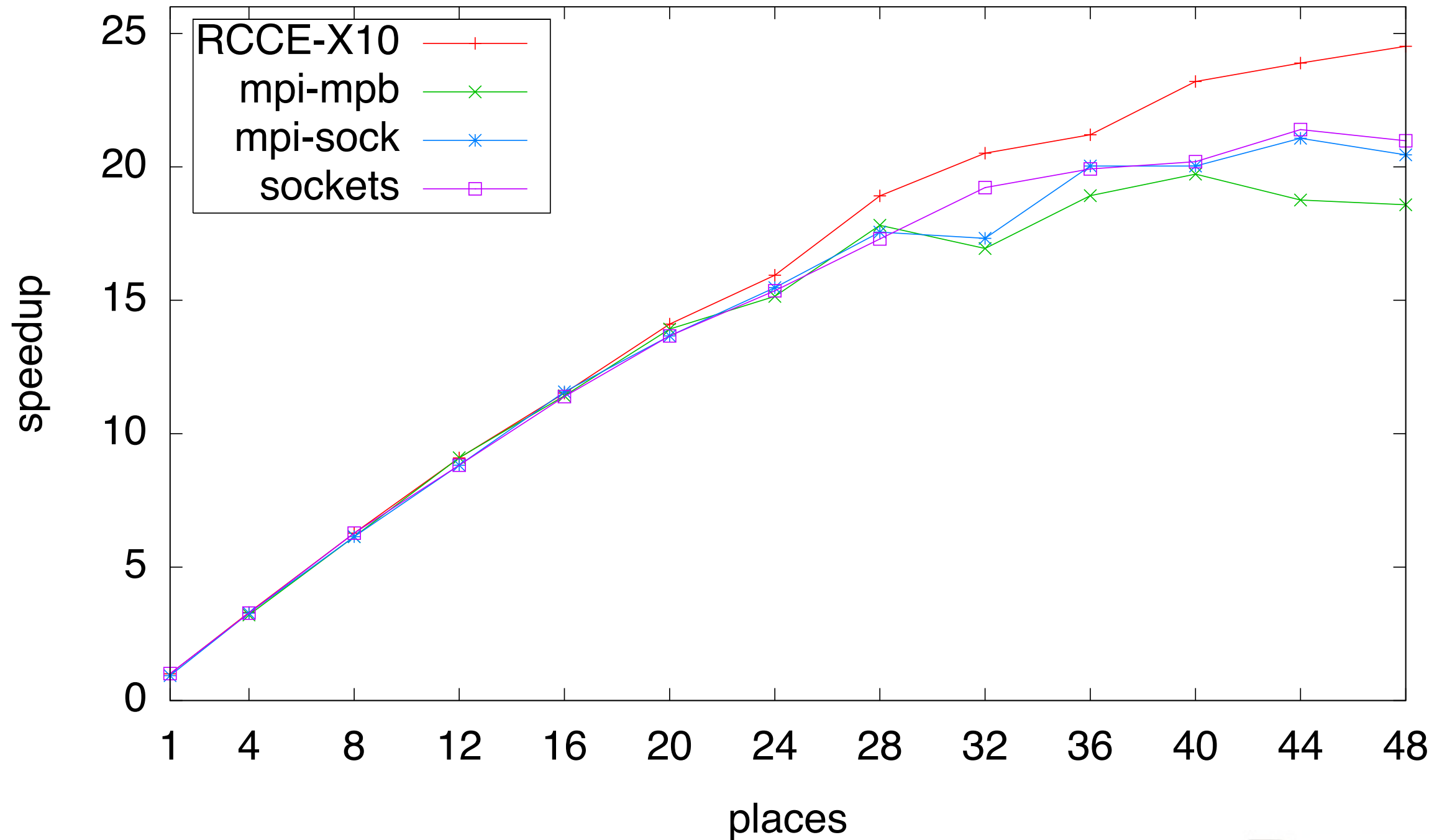
Experiments

- One place per SCC core
- One worker thread per place
- Results in paper without GC
- 322 MB private memory per core
- 2 benchmarks: ANUChem-HF and BC
- X10 post-2.1.1 trunk dated January 13th

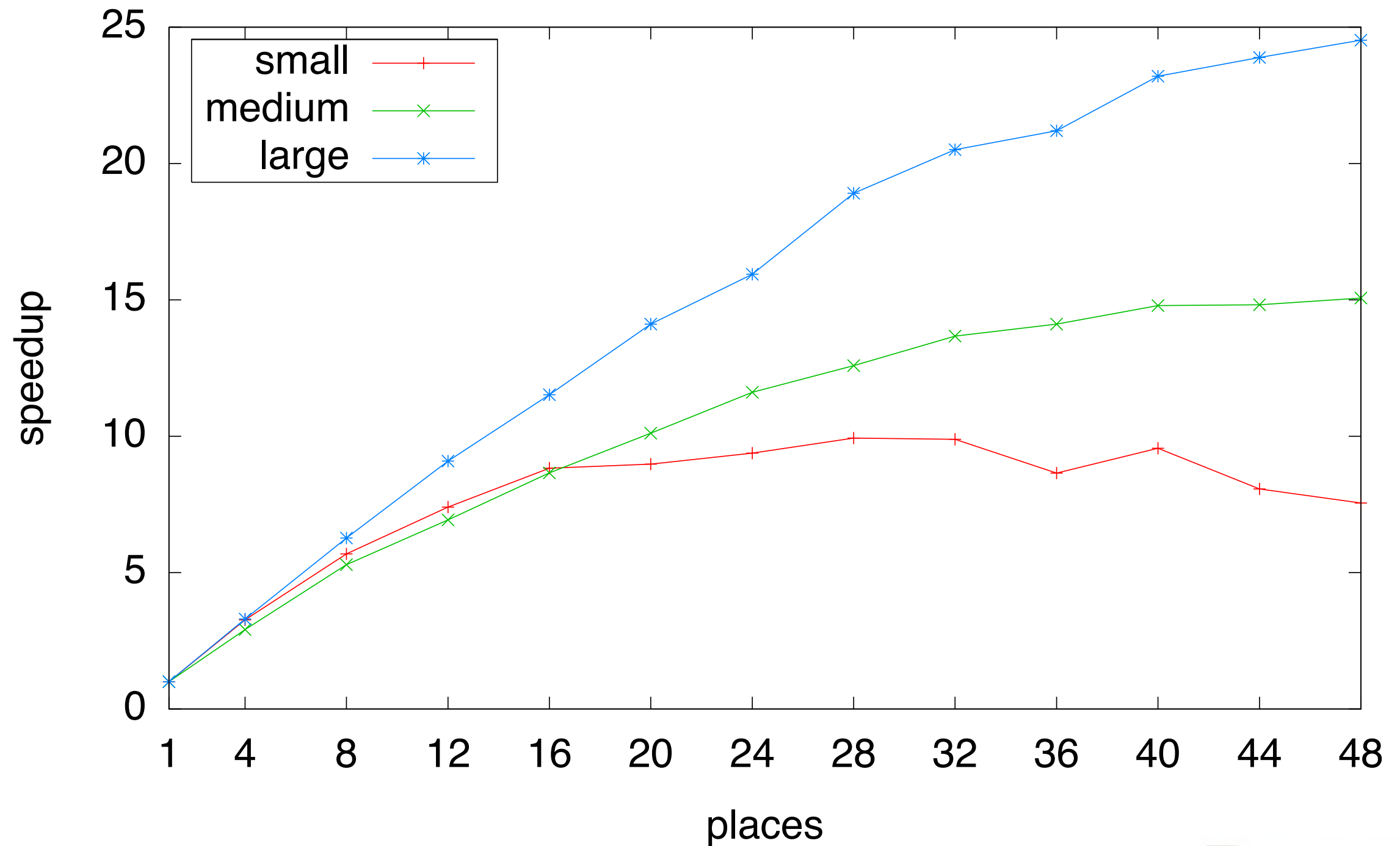
Benchmarks

- ANUChem Hartree-Fock (HF) quantum chemistry benchmark
 - Benzene STO-3G workload (in paper)
 - Benzene 3-21G with GC
- Betweenness Centrality (BC) social network actor centrality benchmark
 - Dijkstra's shortest path algorithm
 - small, medium, large workloads

BC large workload

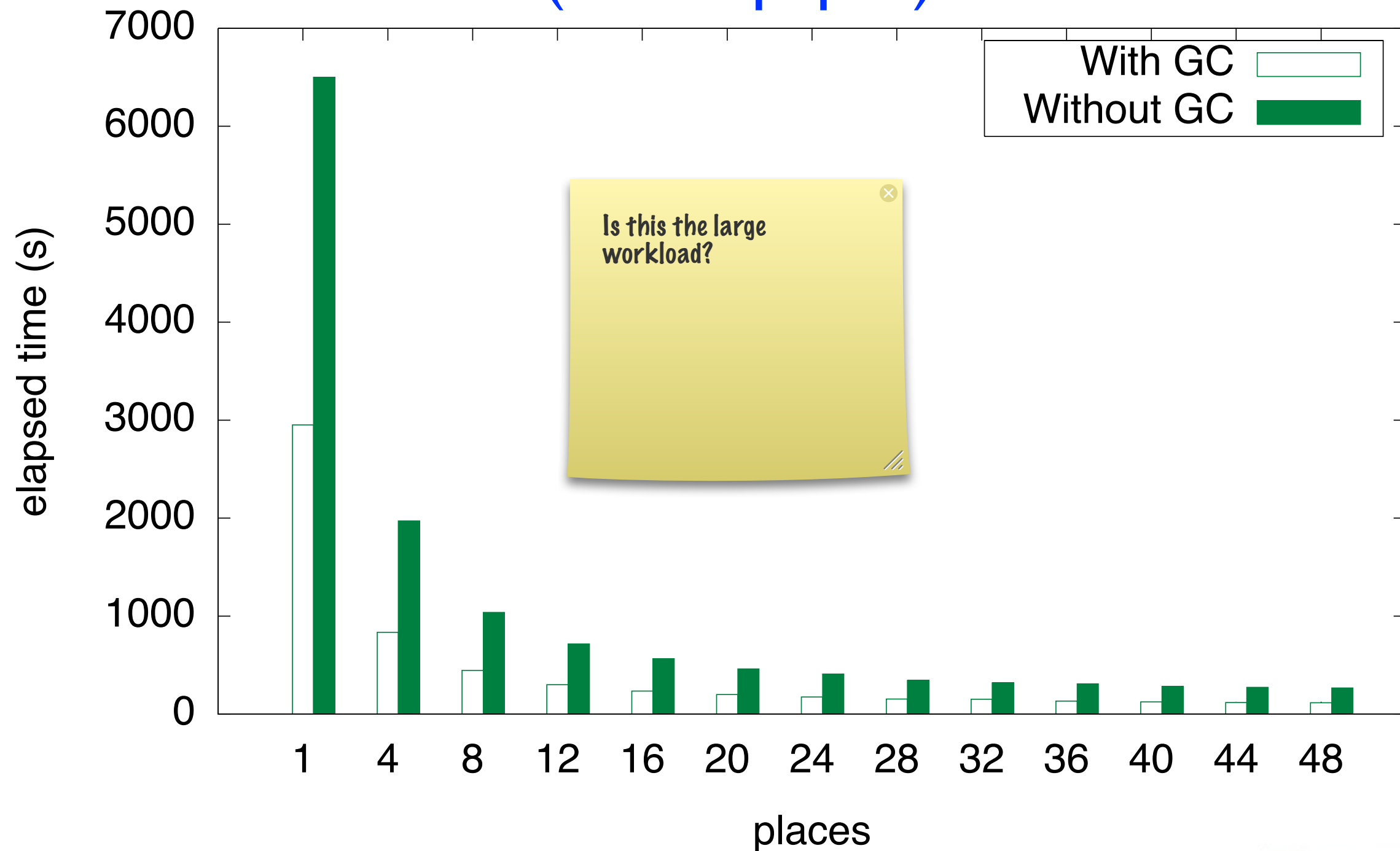


BC varying workload

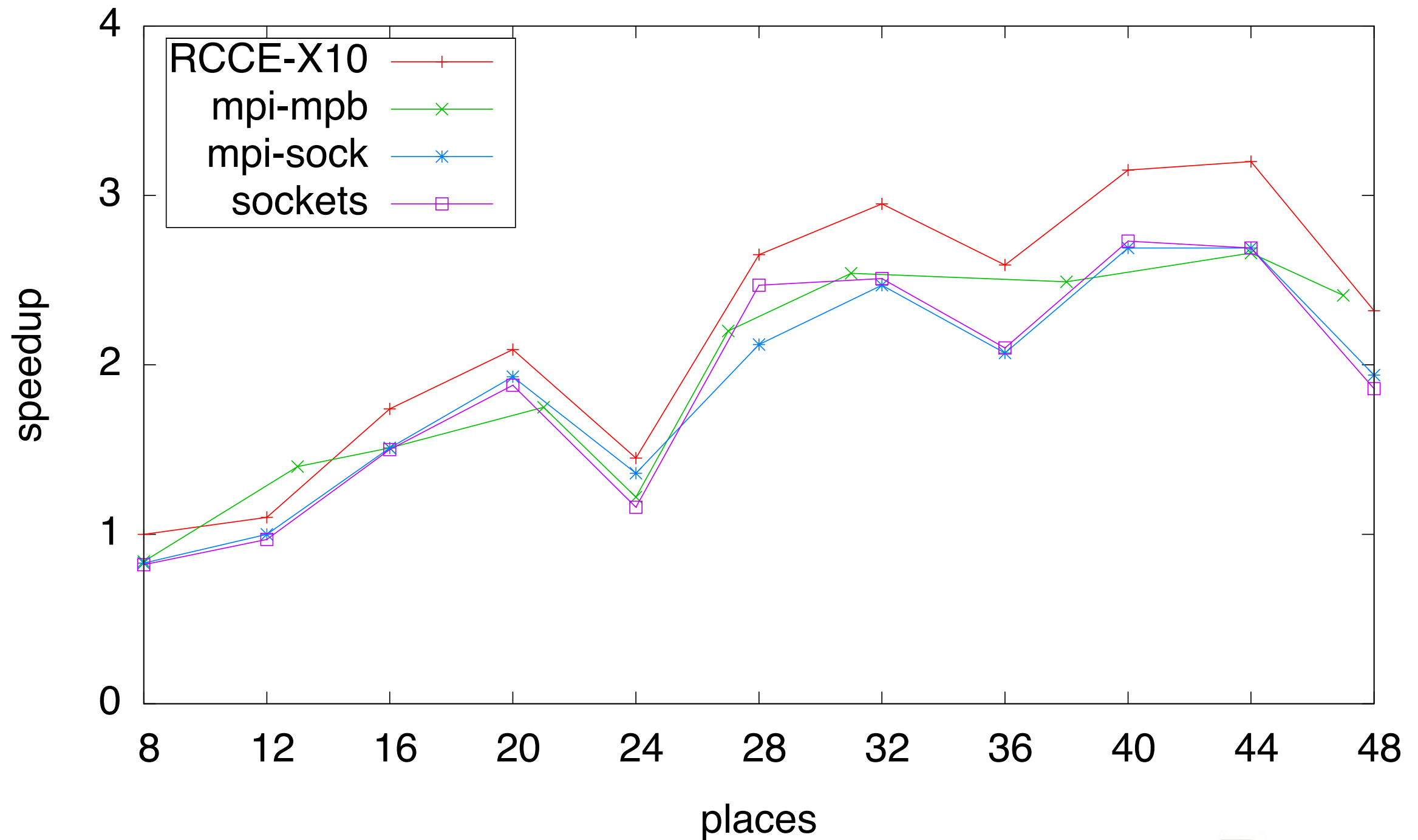


BC with RCCE-X10

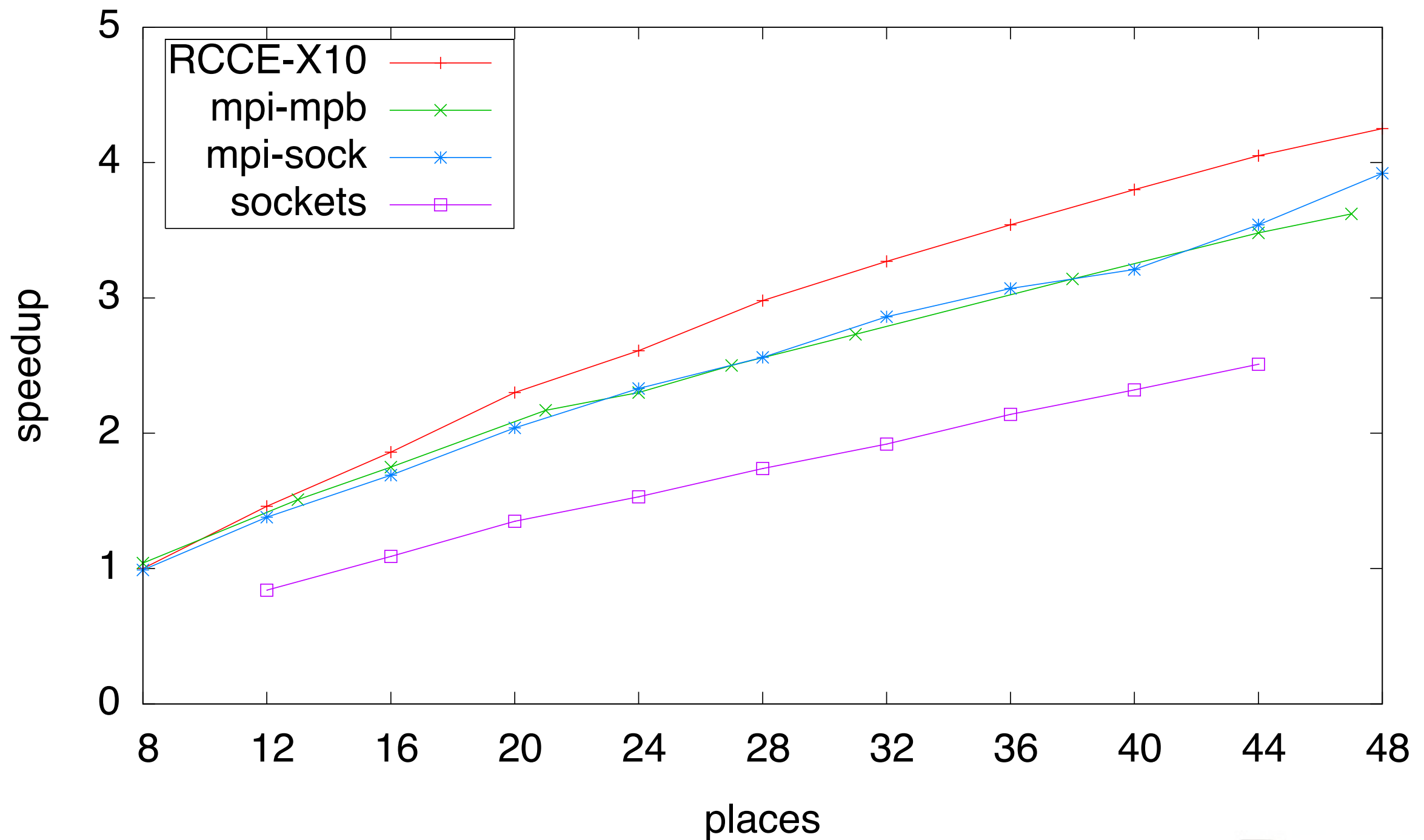
(not in paper)



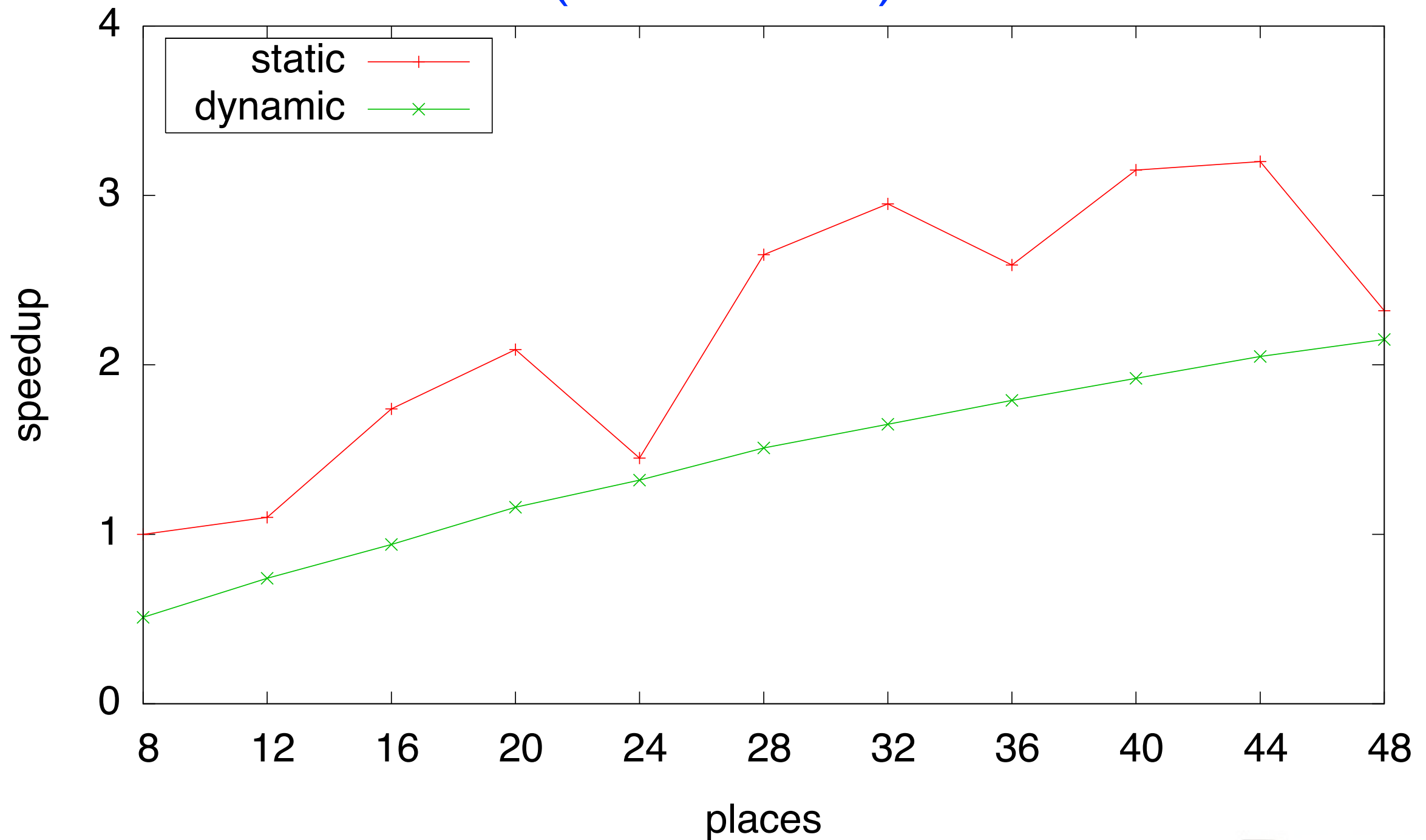
HF static balancing



HF dynamic balancing

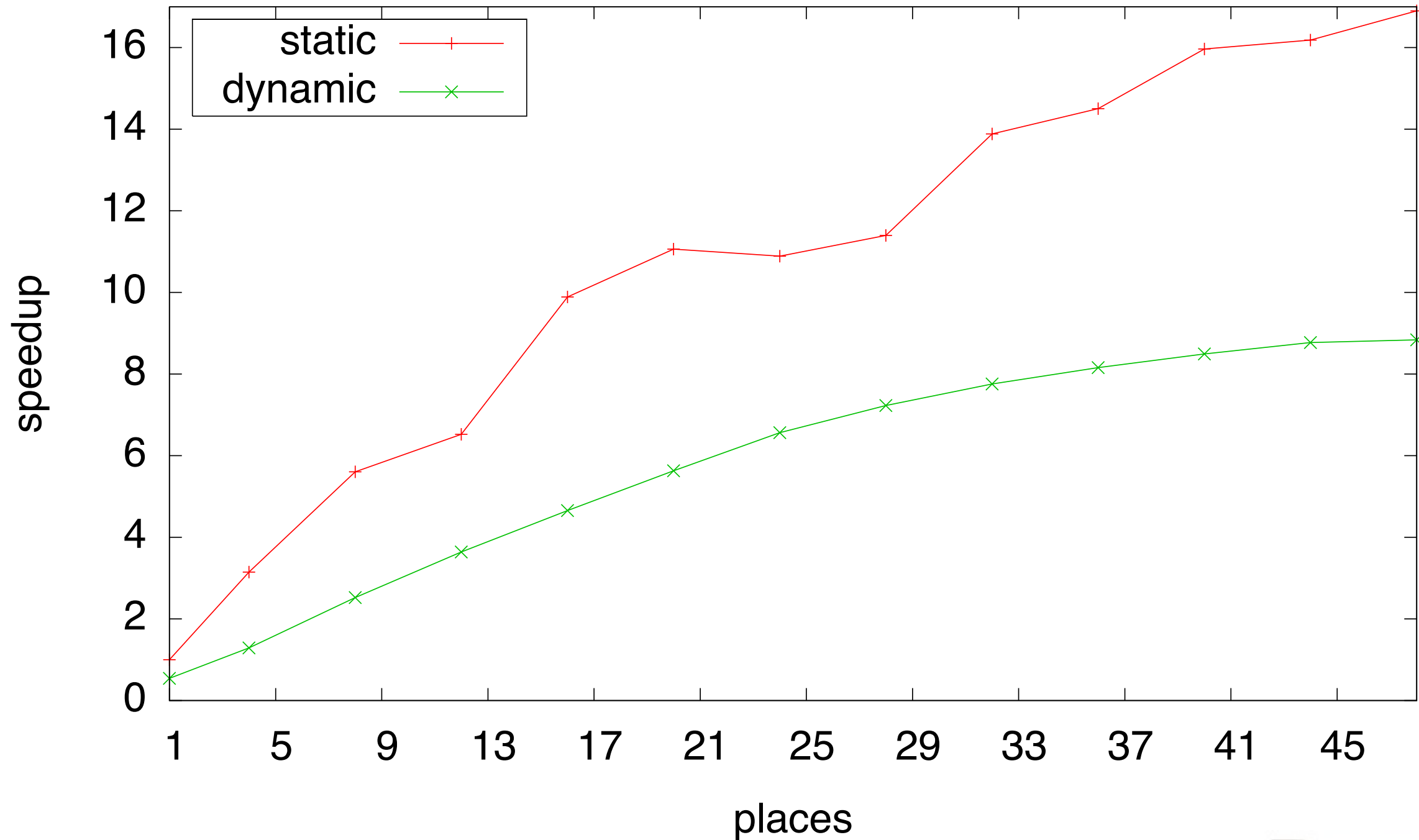


HF dynamic vs static (RCCE-X10)



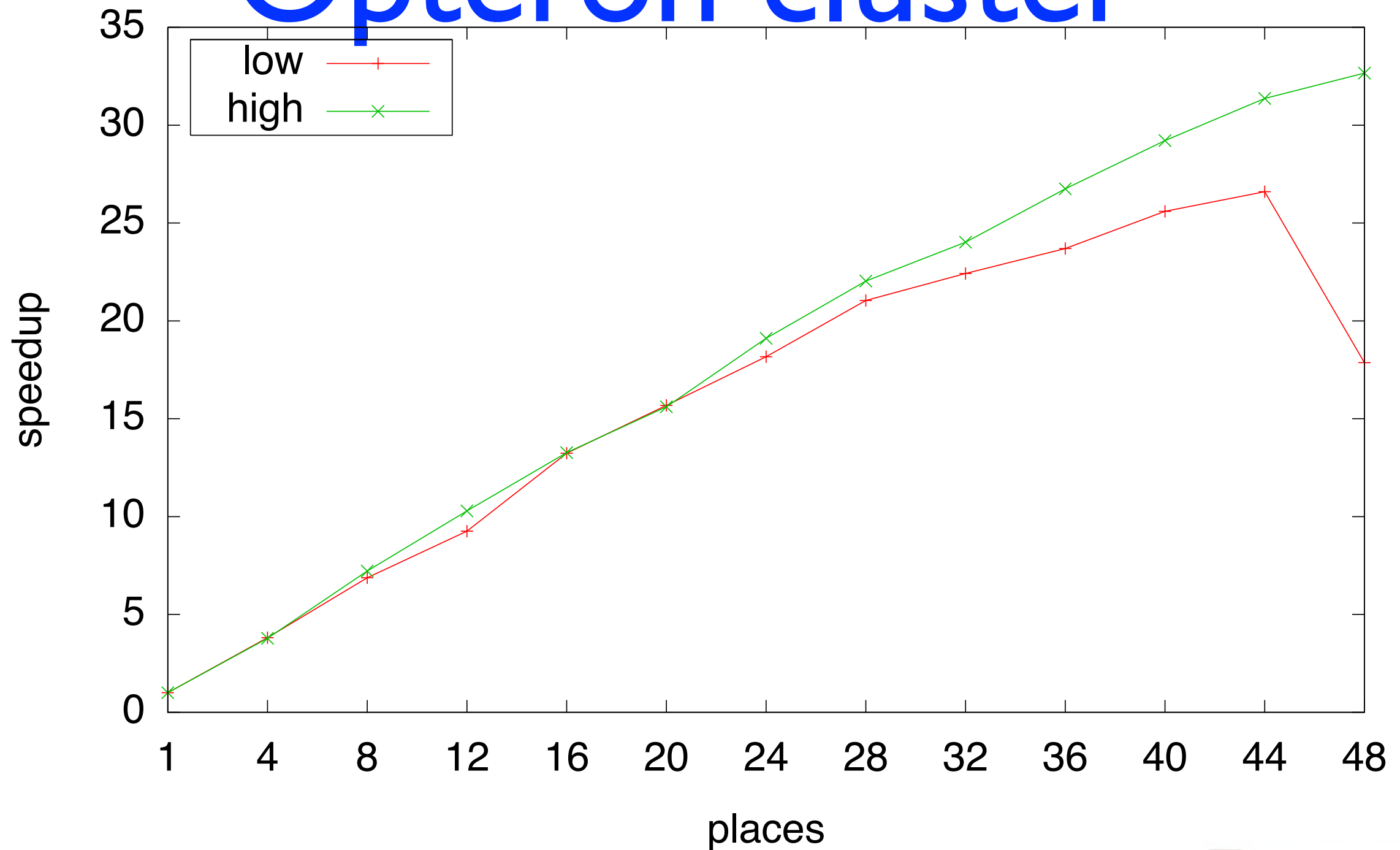
HF with GC

(not in paper)



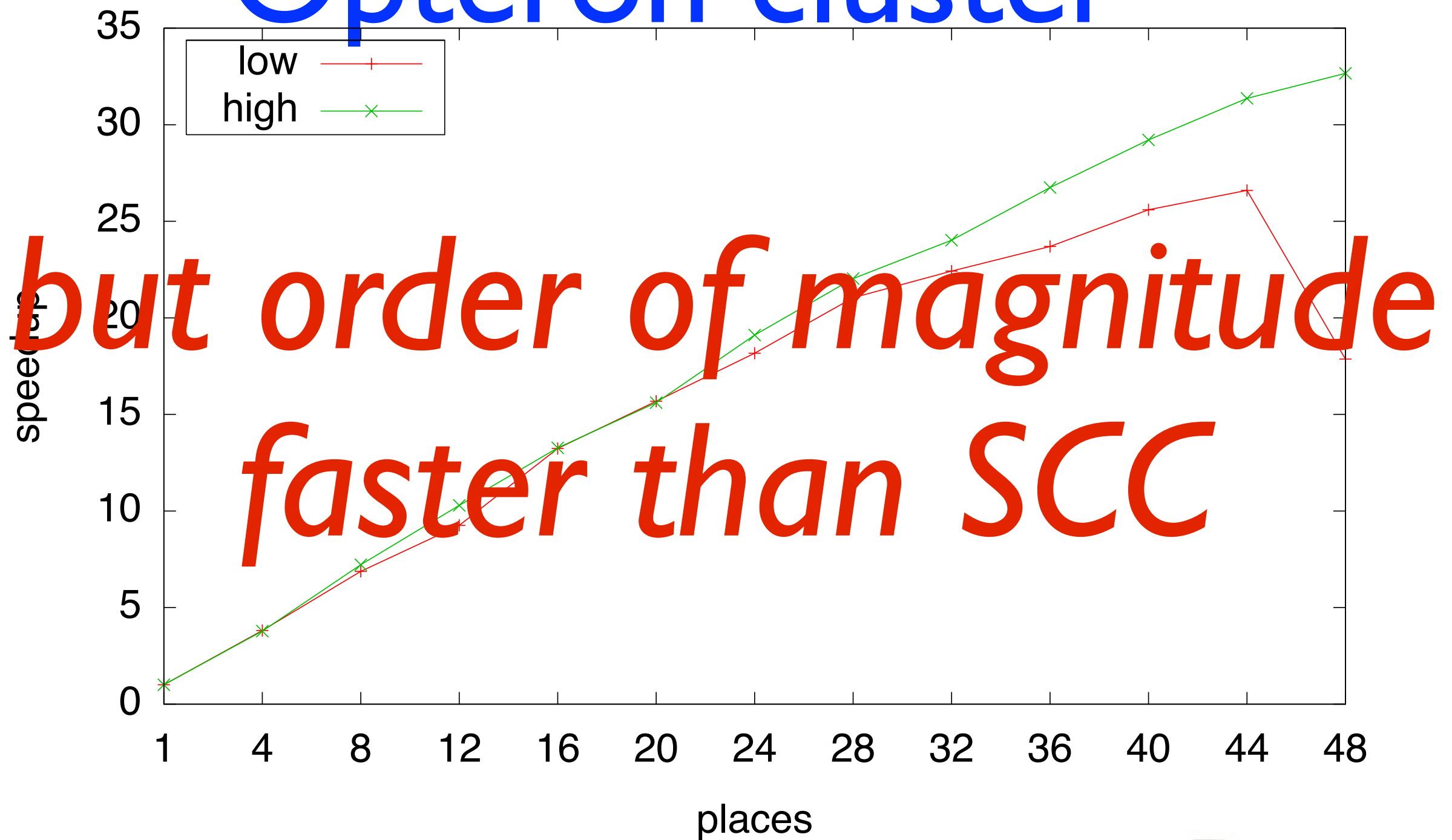
BC with MPI on 3x8x2

Opteron cluster



BC with MPI on 3x8x2

Opteron cluster



Conclusions

- SCC memory constraints are limiting:
 - tradeoff work vs memory to get scaling
 - running with GC allows larger workloads (and better performance!)
- But, X10 + SCC is a nice match
- X10 benchmarks (can) show good scaling
 - order of magnitude slower than cluster