

Today

Distributed Subsystems (Lesson 7)

- ✓ \* Global memory system
- ⇒ \* Distributed Shared Memory

Friday

\* Distributed file Systems

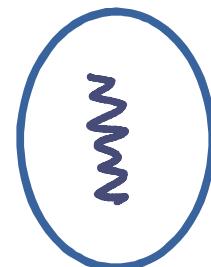
(\* Watch video \*)

# Thought Experiment

## DSM

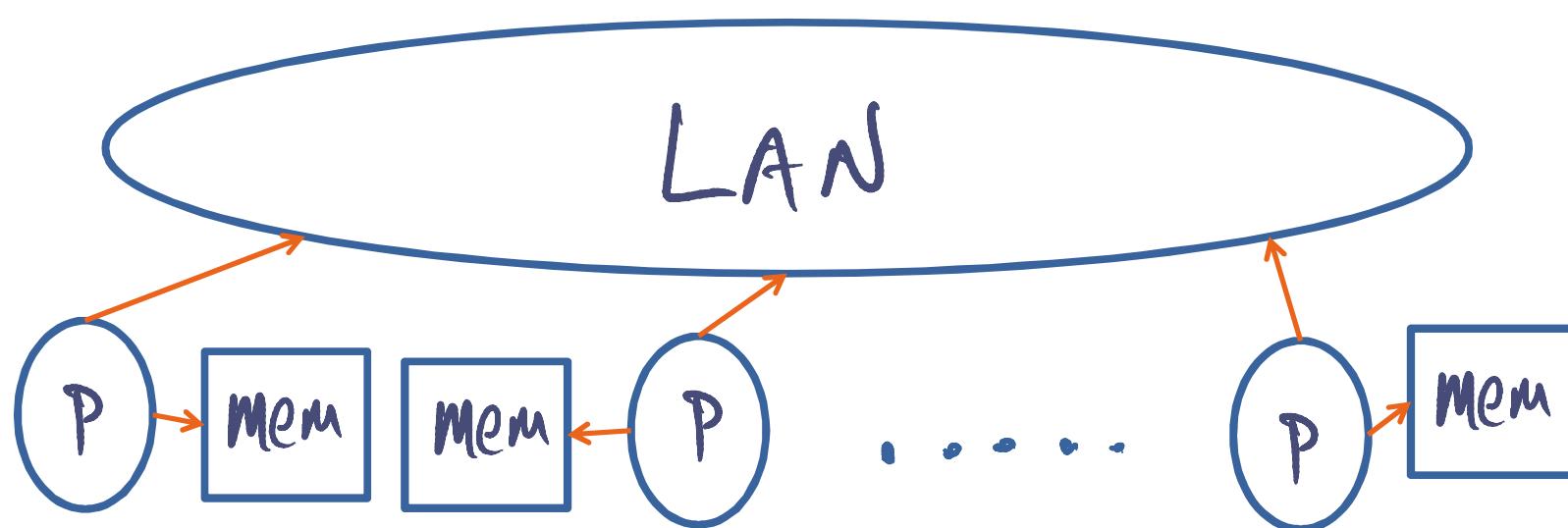
- Can we make the cluster appear like a shared memory machine?

## cluster as a parallel machine



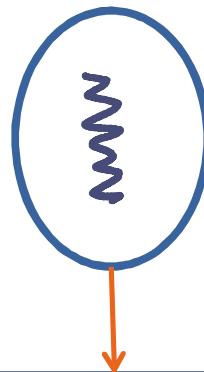
Sequential program

? How to exploit the cluster?

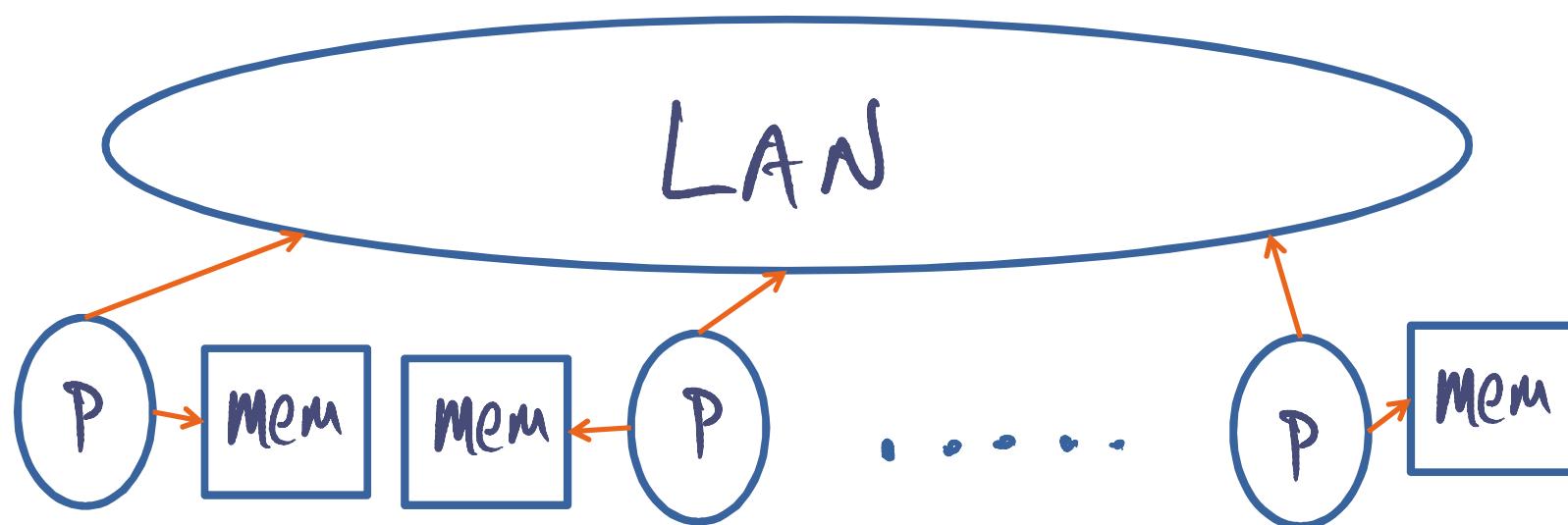


# cluster as a parallel machine

Sequential Program



Automatic  
parallelization



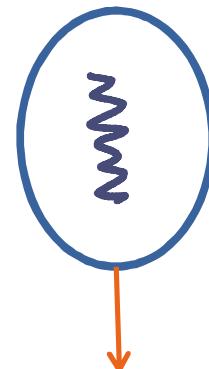
# cluster as a parallel machine

HPF  
- directives for

data/computation  
distribution

- data parallel  
programs

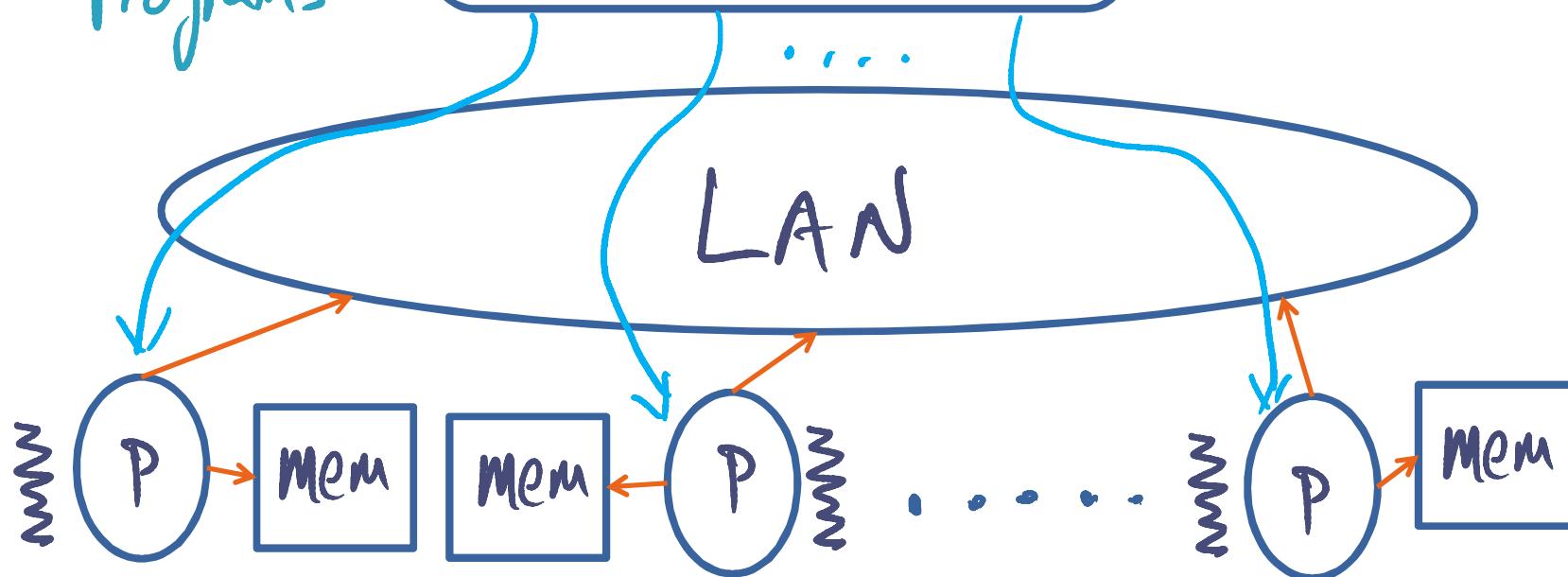
Sequential program



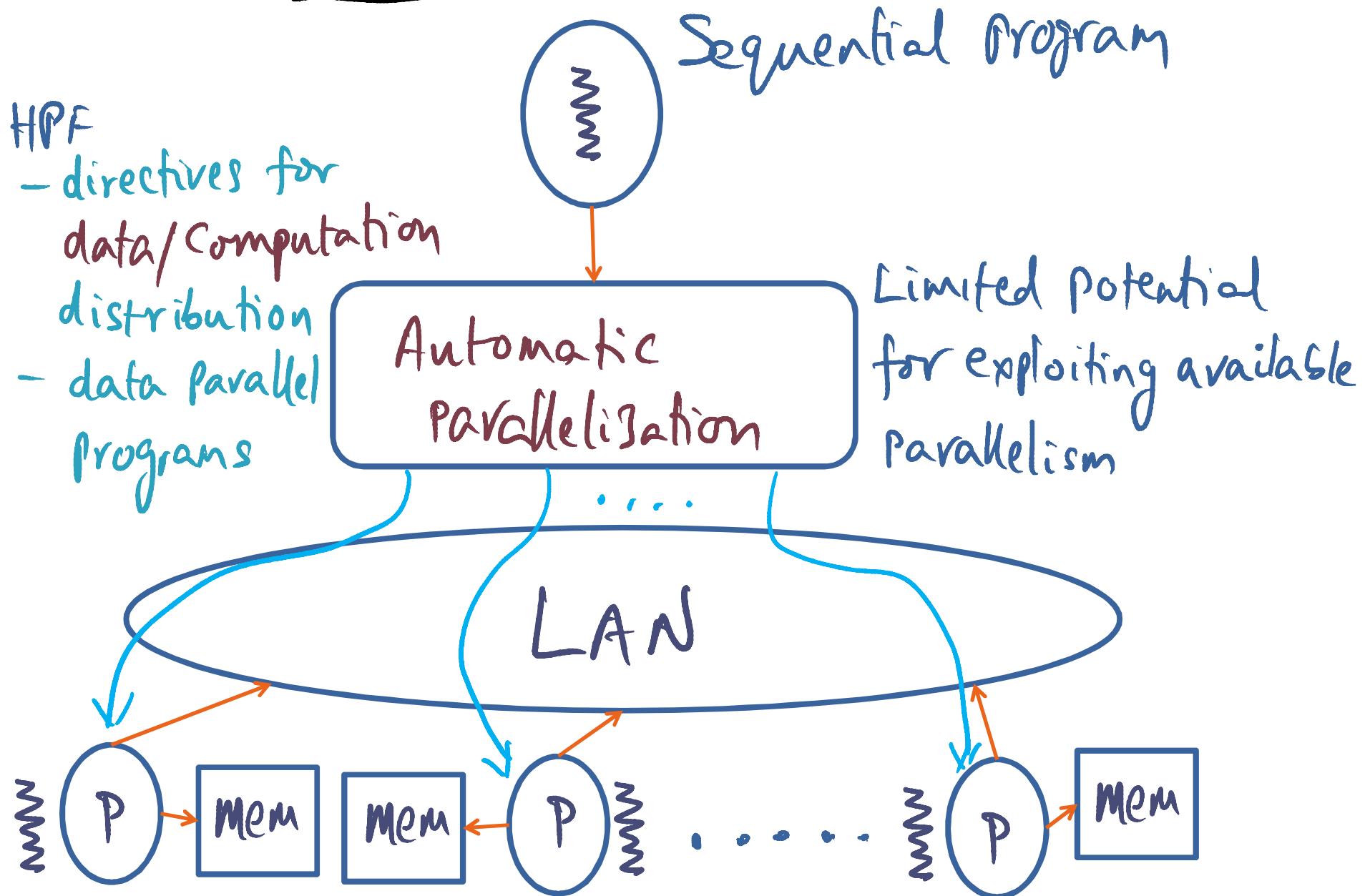
Automatic  
parallelization

...

LAN

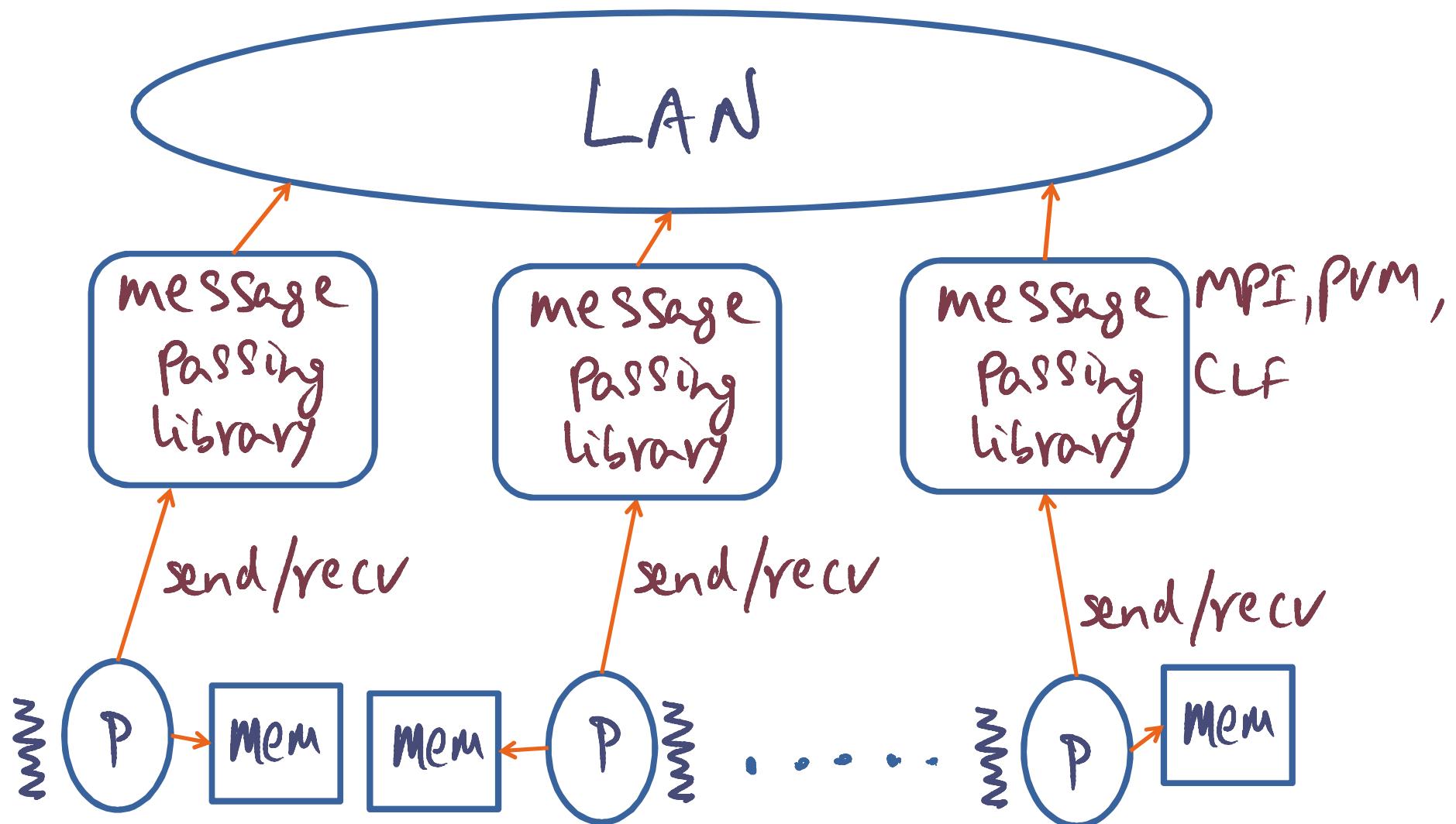


# cluster as a parallel machine



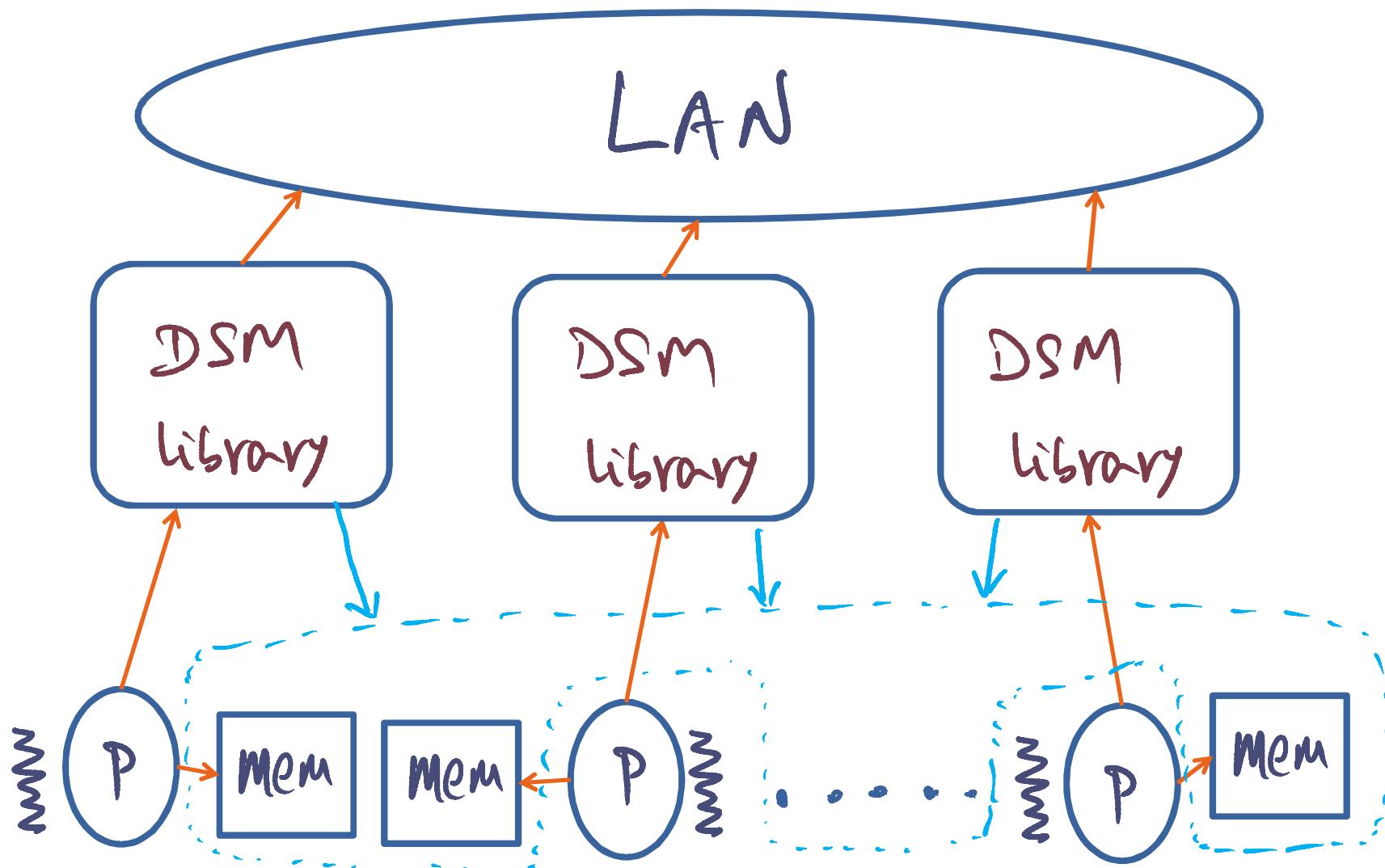
# cluster as a parallel machine

## Parallel program: message passing

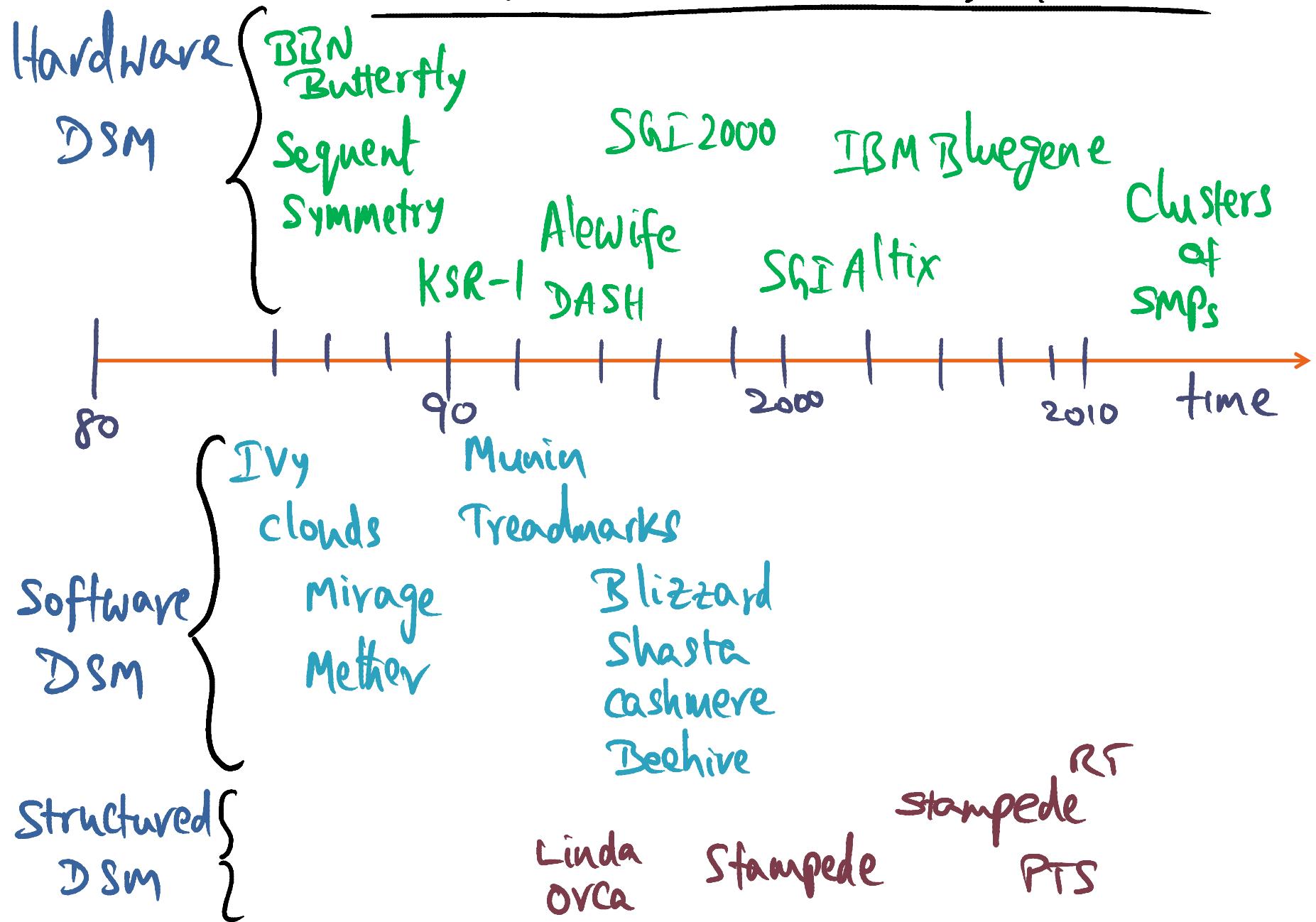


# cluster as a parallel machine

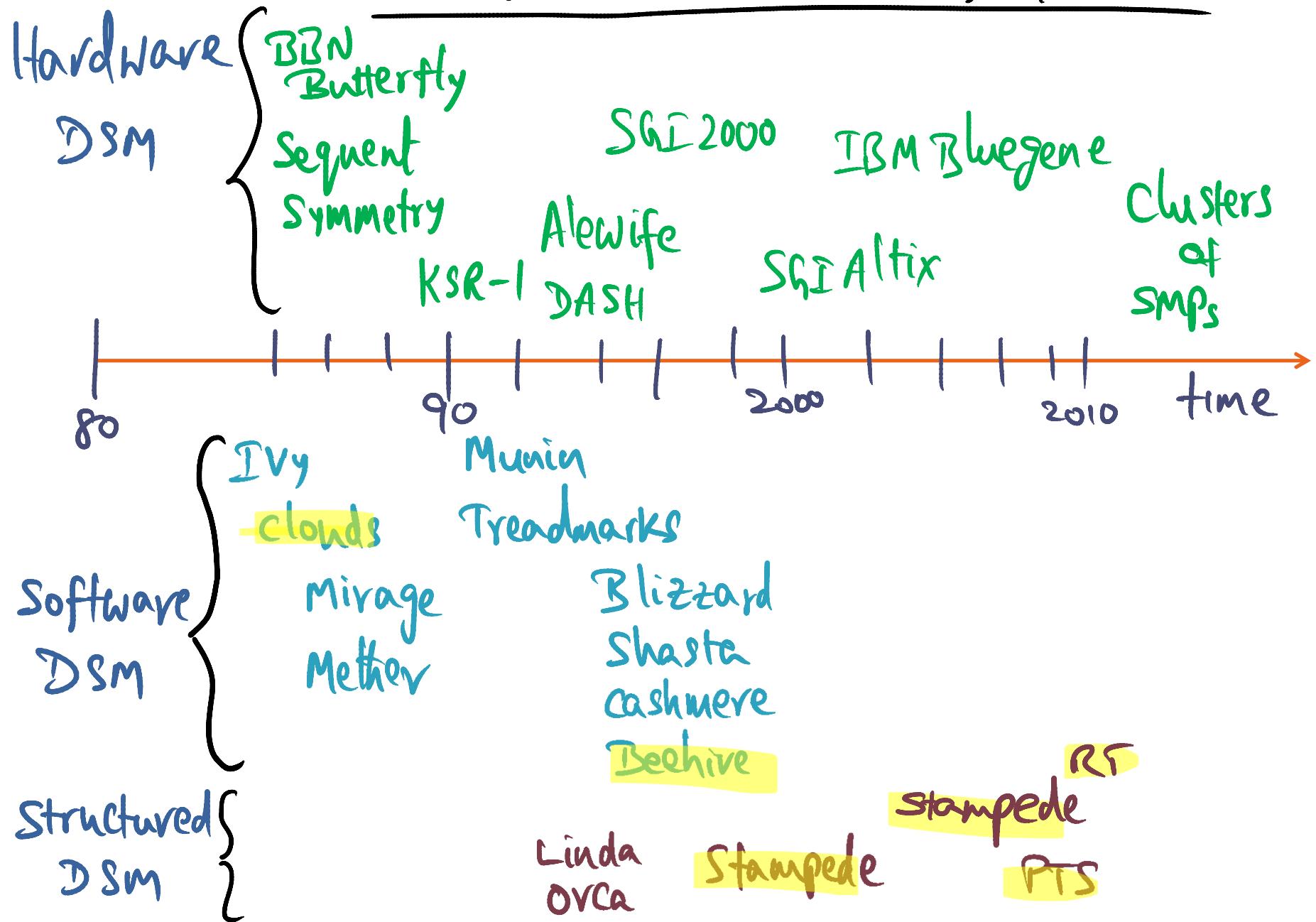
Parallel program: Distributed shared memory



# History of shared memory systems

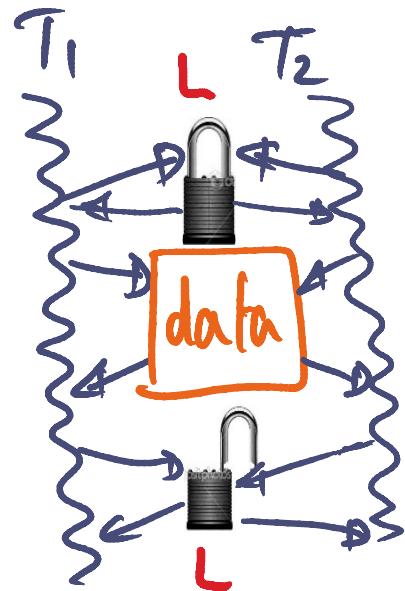


# History of shared memory systems

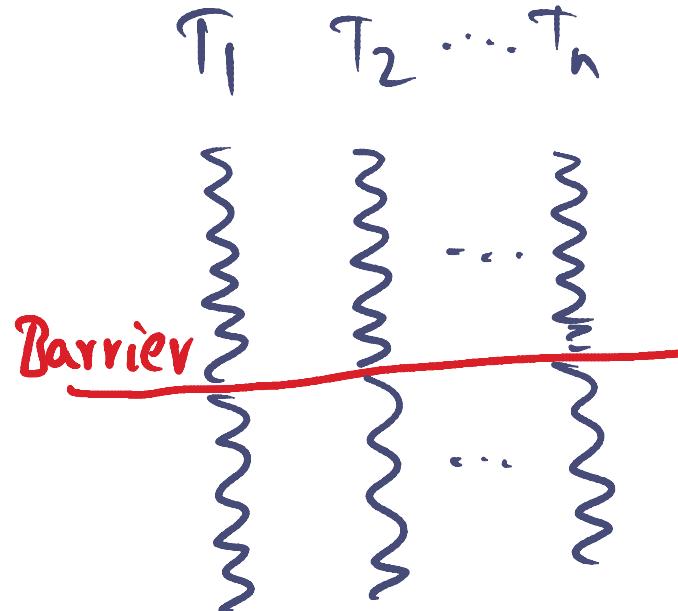


# Shared memory Programming

## LOCK



## Barrier



Two types of memory accesses

- normal r/w to shared data
- r/w to synchronization variables

Memory Consistency



Cache Coherence



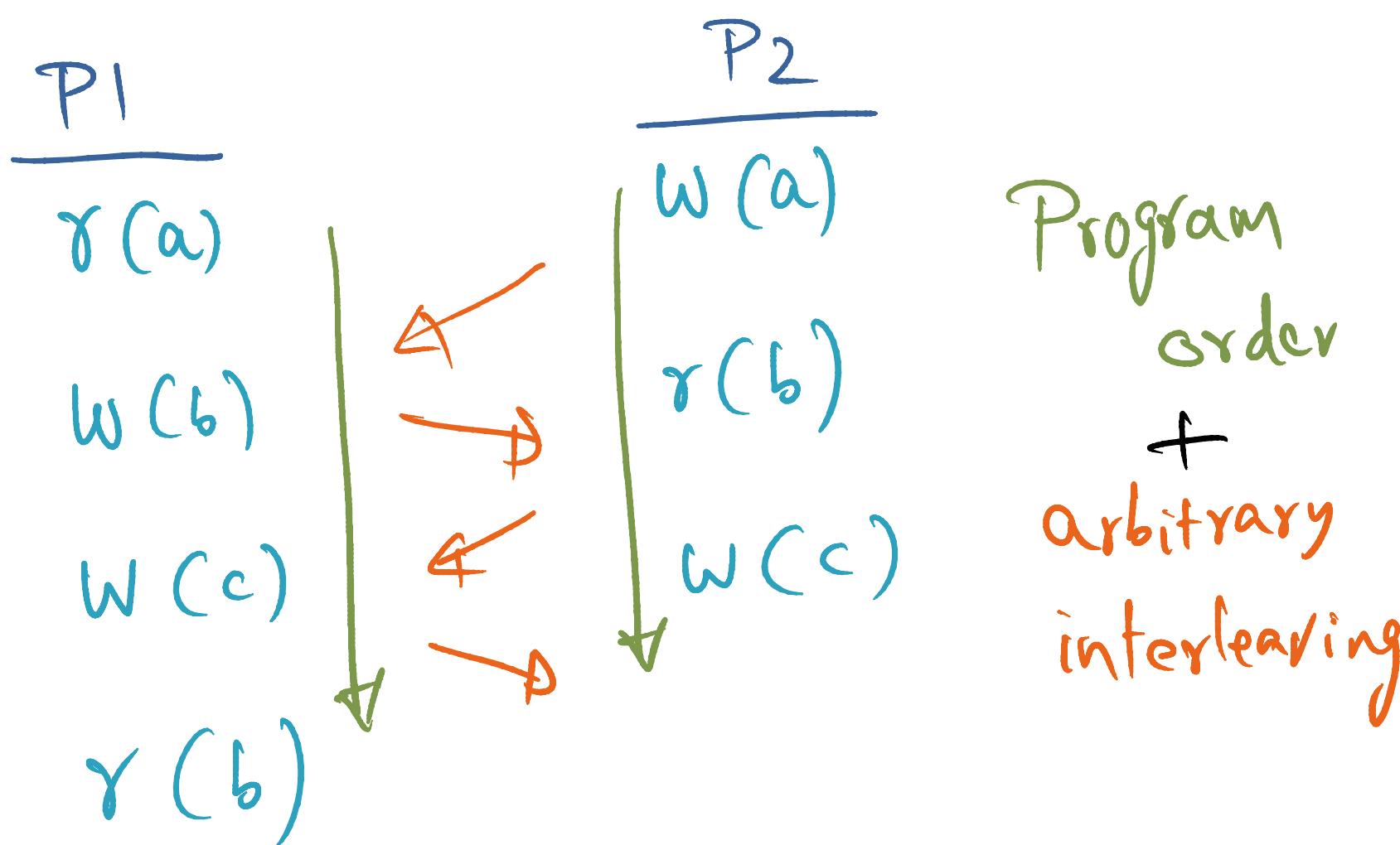
What is the Model  
Presented to the

Programmer?



How is the System  
implementing the  
Model in the presence  
of Private Caches?

# Sequential consistency (SC)



# Sequential consistency (SC)

P1

$r(a)$

$w(b)$

$w(c)$

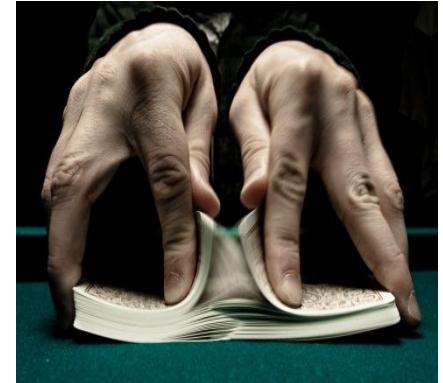
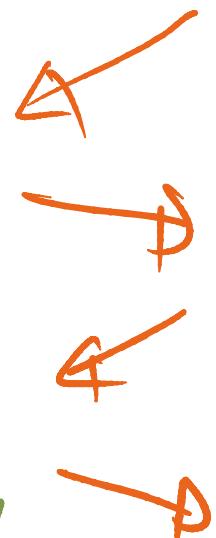
$r(b)$

P2

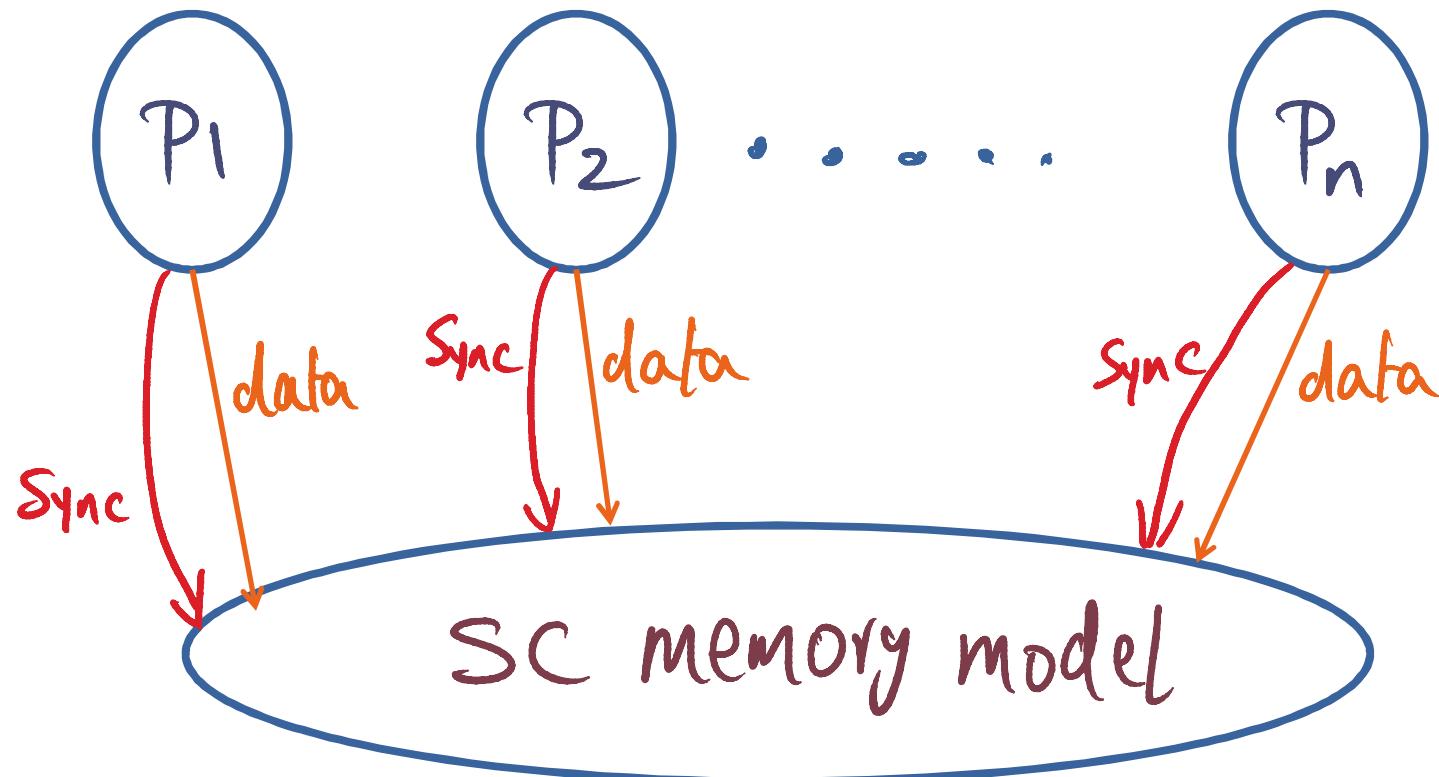
$w(a)$

$r(b)$

$w(c)$



Program  
order  
+  
arbitrary  
interleaving



SC does not distinguish between  
data r/w and synchronization r/w

upshot?

Coherence action on every r/w access

## Typical Parallel Program

P1

lock(L); //acq

read(a);

write(b);

!  
unlock(L); //rel

P2

lock(L); //acq

write(a);

read(b);

!  
unlock(L); //rel

## Typical Parallel Program

P1

Lock(L); //acq

read(a);

write(b);

!  
Unlock(L); //rel

P2

⊕ Lock(L); //acq

write(a);

read(b);

!  
Unlock(L); //rel

P2 does not access  
data P1 releases L

## Typical Parallel Program

P1

lock(L); //acq

read(a);

write(b);

!

unlock(L); //rel

P2

lock(L); //acq

write(a);

read(b);

!

unlock(L); //rel

P2 does not access

data P1 releases L

⇒ no need for coherence action for a and b  
until release

## Release Consistency

P1:

$a_1: \text{acq}(L)$

data  
accesses

$r_1: \text{rel}(L)$

P2:

$a_2: \text{acq}(L)$

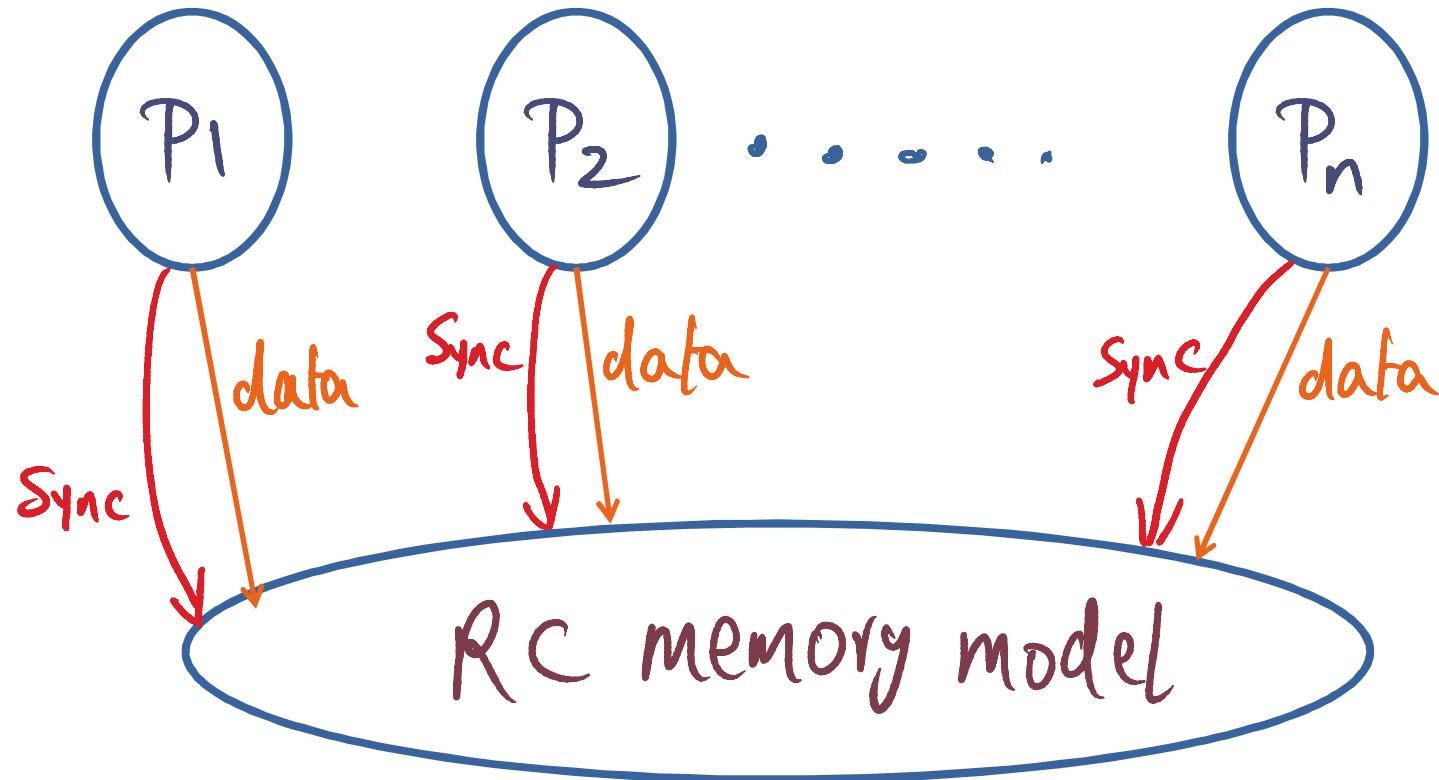
data  
accesses

$r_2: \text{rel}(L)$

If  $P1: r_1 \xrightarrow{hb} P2: a_2$

- All Coherence actions Prior to  $P1: r_1$

Should be complete before  $P2: a_2$



RC distinguishes between  
data  $\wedge$   $/w$  and Synchronization  $\wedge$   $/w$

upshot?

Coherence action only when lock released

## Example

Programmer's intent

P1

modify struct(A)



P2

wait for mod  
use struct (A)

## Example Programmer's intent

P1

modify struct(A)



P2

wait for mod  
use struct(A)

P2

LOCK(L);  
while(flag == 0)

P2 here  $\rightarrow$  Wait(C, L);  
unlock(L);  
use(A);

## Example Programmer's intent

P1

modify struct(A)



P1

modify(A);

P2

wait for mod  
use struct(A)

P2

LOCK(L);

while(flag == 0)

P2 here  $\rightarrow$  Wait(C, L);

unlock(L);

use(A);

## Example

### Programmer's intent

P1

modify struct(A)



P1

modify(A);

LOCK(L);

flag = 1;

signal(c);

unlock(L);

P2

wait for mod  
use struct(A)

P2

LOCK(L);

while(flag == 0)

P2 here  $\rightarrow$  Wait(C, L);

unlock(L);

use(A);

## Example

### Programmer's intent

P1

modify struct(A)

—X—

P1

modify(A);

LOCK(L);

flag = 1;

signal(C);

unlock(L);

Coherence

actions complete

P2

wait for mod  
use struct(A)

P2

LOCK(L);

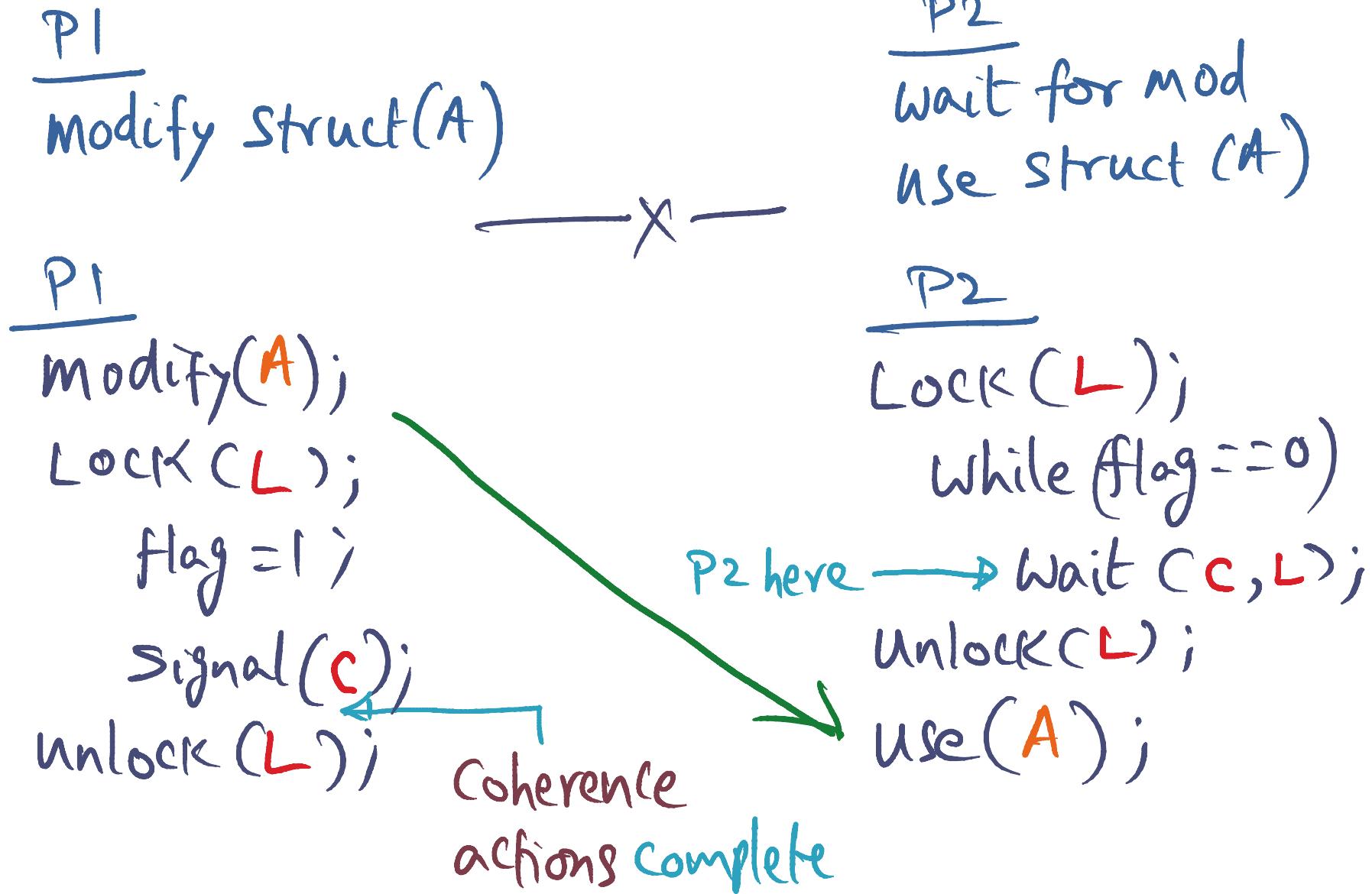
while(flag == 0)

P2 here  $\rightarrow$  Wait(C, L);

unlock(L);

use(A);

## Example Programmer's intent



## Advantage of RC over SC

— no waiting for coherence actions on  
every memory access

⇒ Overlap Computation with

Communication

⇒ better performance for RC over SC

## Lazy RC

P1:

$a_1: acq(L)$

data  
accesses

$r_1: rel(L)$

P2:

$a_2: acq(L)$

data  
accesses

$r_2: rel(L)$

## Lazy RC

P1:

$a_1: \text{acq}(L)$

data  
accesses

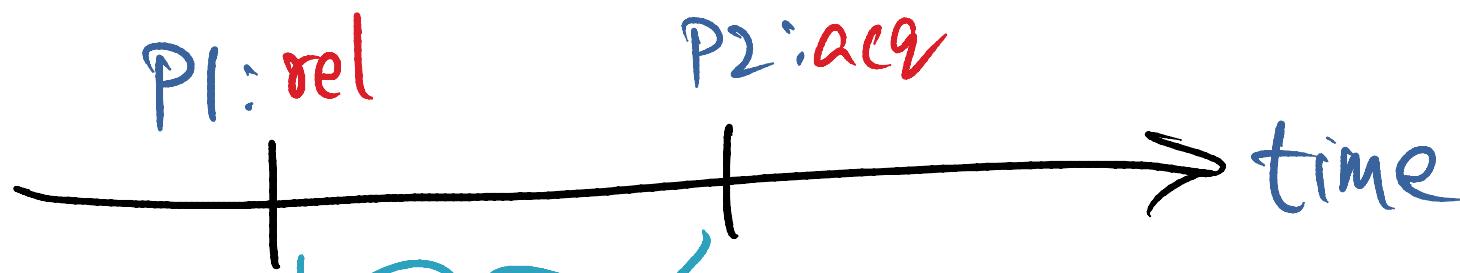
$r_1: \text{rel}(L)$

P2:

$a_2: \text{acq}(L)$

data  
accesses

$r_2: \text{rel}(L)$



opportunity for procrastination!

## Lazy RC

P1:

$a1: \text{acq}(L)$

data  
accesses

$r1: \text{rel}(L)$

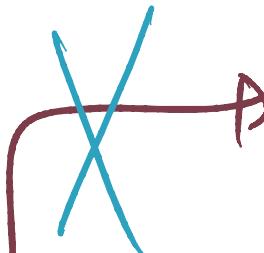
P2:

$a2: \text{acq}(L)$

Coherence  
action  
at  
acquire

data  
accesses

$r2: \text{rel}(L)$



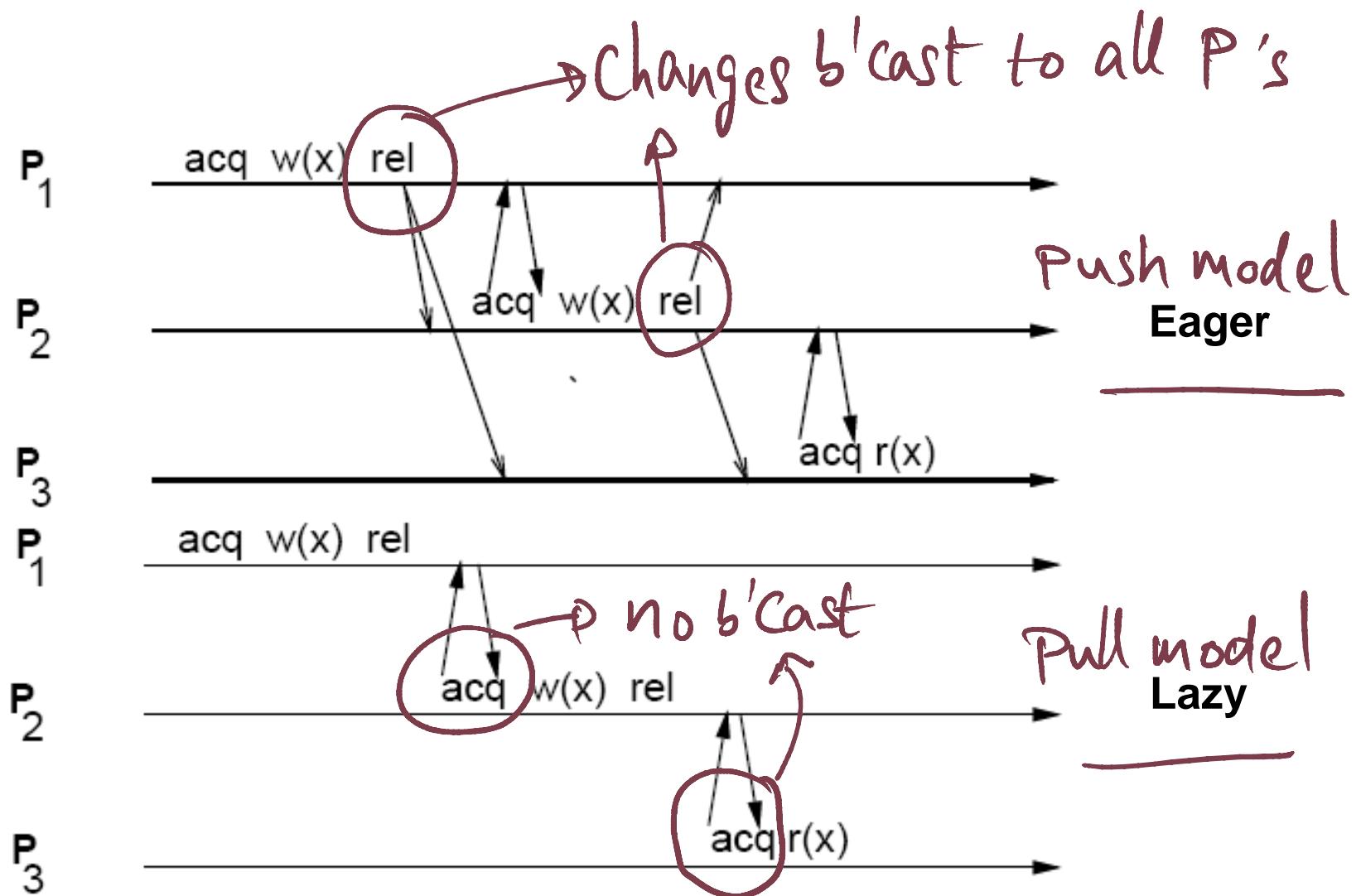
P1:  $\text{rel}$

P2:  $\text{acq}$

time

opportunity for procrastination!

# Eager Vs. Lazy RC



What are the Pros of Lazy over Eager?

What are the Cons of Lazy over Eager?

What are the Pros of Lazy over Eager?

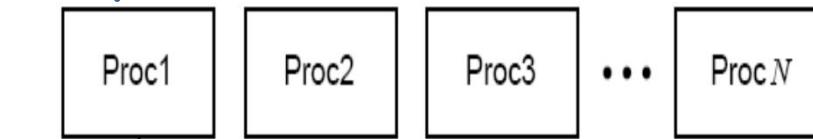
Less messages

What are the Cons of Lazy over Eager?

More latency at acquire

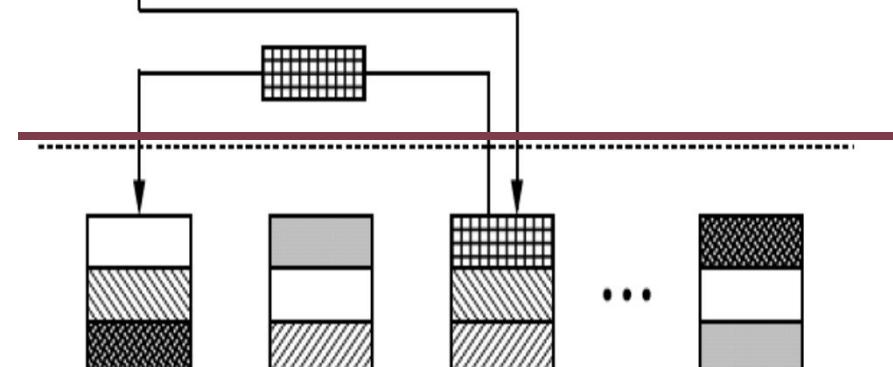
# Software DSM

Application



Global Virtual  
memory  
abstraction

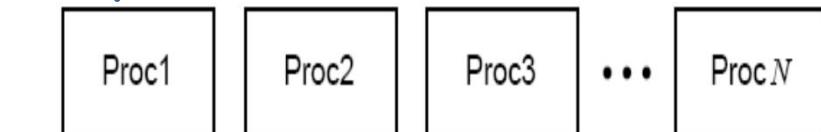
Address space partitioned  
Address equivalence  
Distributed ownership



Local physical Memories

# Software DSM

Application

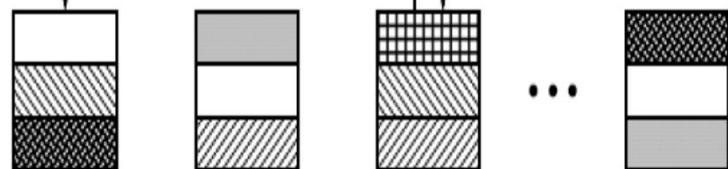


Global Virtual  
memory  
abstraction

Address space partitioned  
Address equivalence  
Distributed ownership

Contact owner of page  
to get current copy

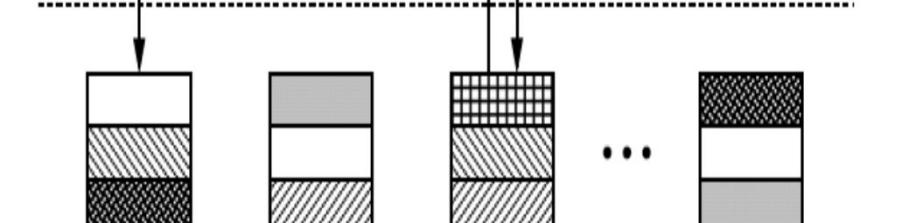
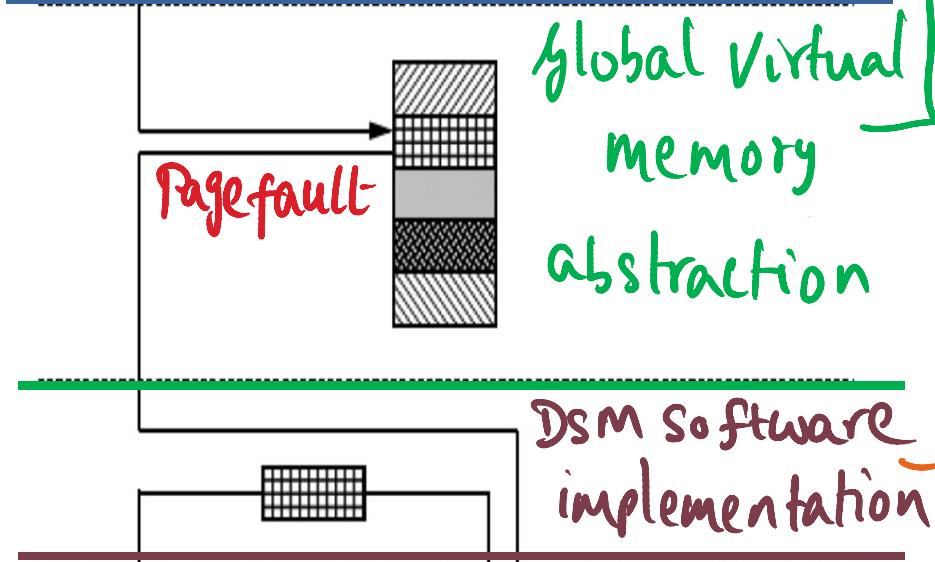
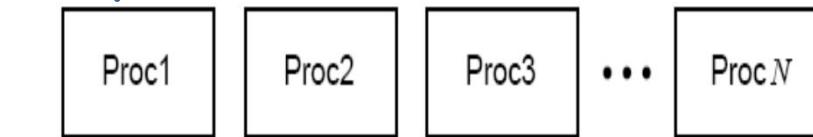
DSM Software  
implementation



Local physical Memories

# Software DSM

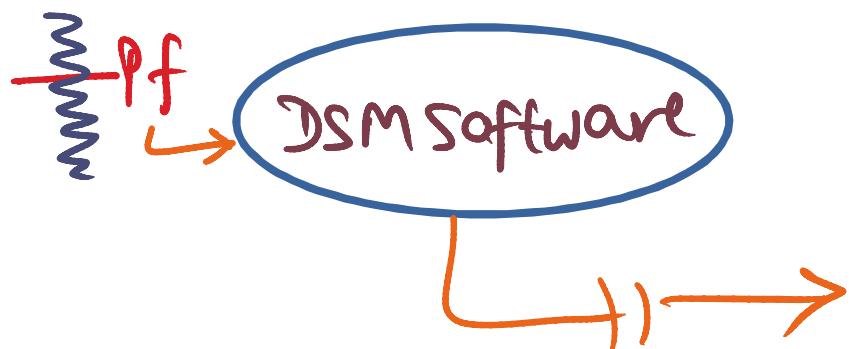
Application



Local physical Memories

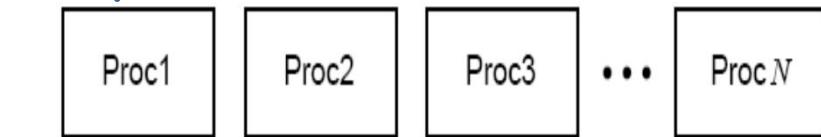
Address space partitioned  
Address equivalence  
Distributed ownership

Contact owner of page  
to get current copy



# Software DSM

Application

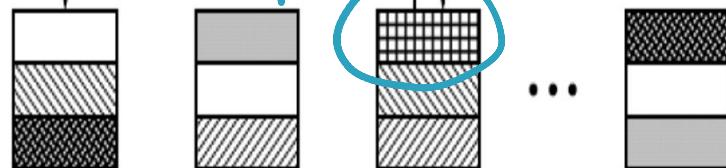


Global Virtual  
memory  
abstraction

Pagefault

DSM Software  
implementation

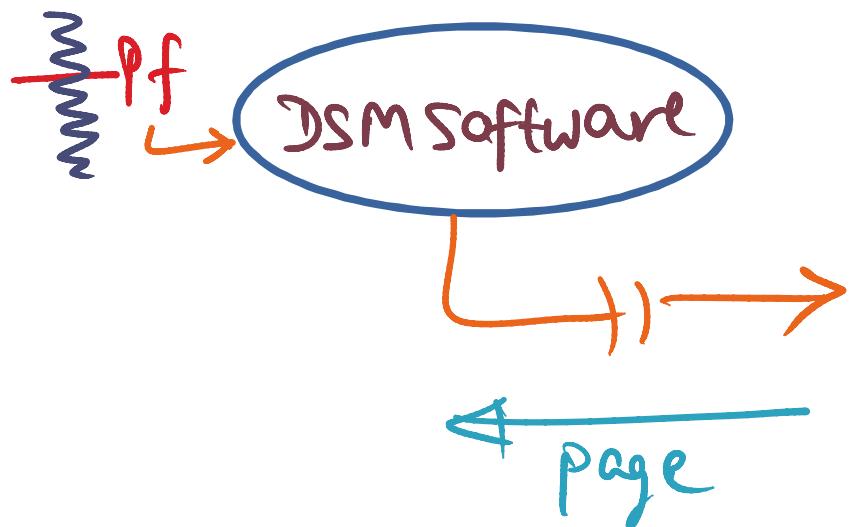
data xfer



Local physical Memories

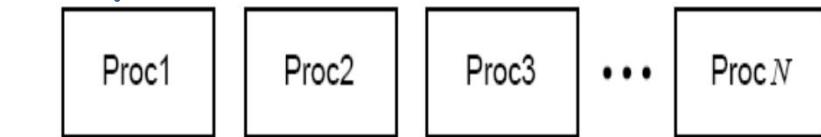
Address space partitioned  
Address equivalence  
Distributed ownership

Contact owner of page  
to get current copy



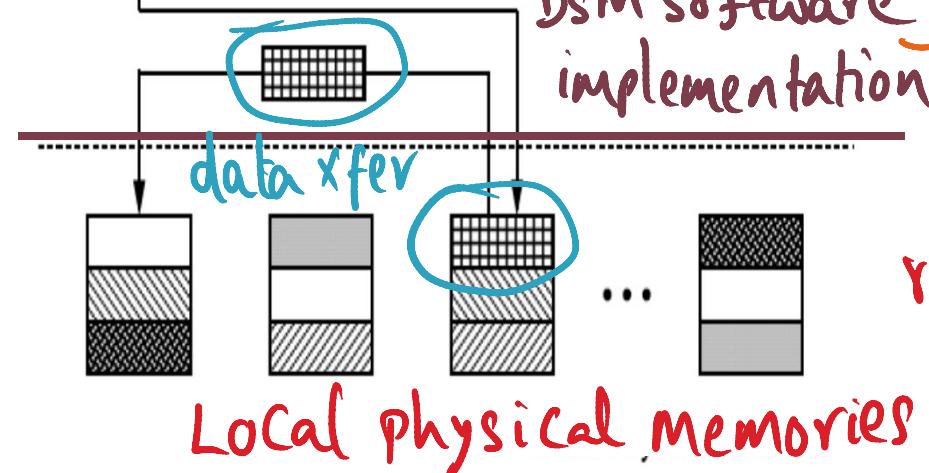
# Software DSM

Application



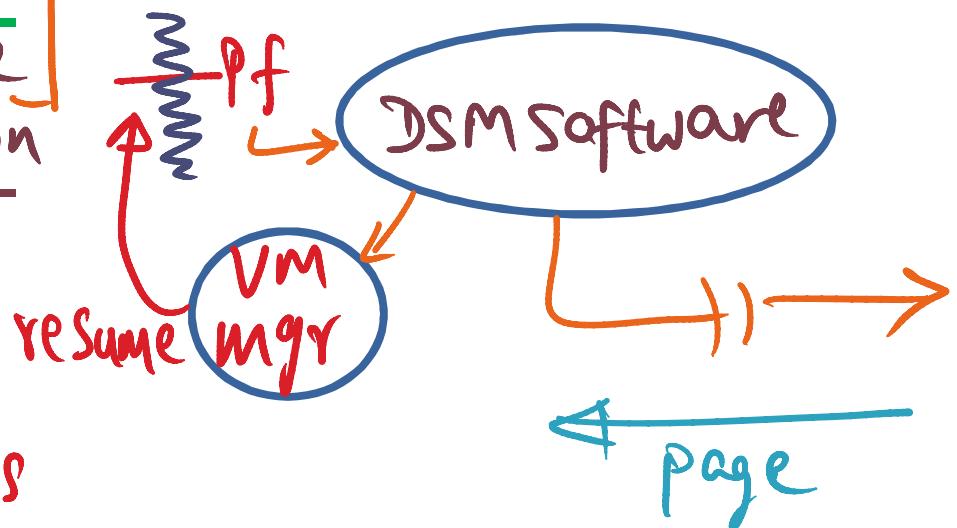
Global Virtual  
memory  
abstraction

Pagefault



Address space partitioned  
Address equivalence  
Distributed ownership

Contact owner of page  
to get current copy



## LRC with Multi-Writer Coherence Protocol

P1

LOCK(L);

$x \leftarrow \begin{cases} x, y, z \text{ pages} \\ \text{modified in C.S} \end{cases} \Rightarrow x_d, y_d, z_d \} \text{diffs}$

Unlock(L);

## LRC with Multi-Writer Coherence Protocol

P1

LOCK(L);

$x \leftarrow \begin{cases} x, y, z \text{ pages} \\ y \leftarrow \text{modified in C.S} \end{cases} \Rightarrow x_d, y_d, z_d \} \text{ diff's}$

Unlock(L);

P2

invalidate  $\rightarrow$  LOCK(L);  
x, y, z at  
lock acquisition  
(CLRC)

## LRC with Multi-Writer Coherence Protocol

P1

LOCK(L);

$x \leftarrow \begin{cases} x, y, z \text{ pages} \\ y \leftarrow \text{modified in C.S} \end{cases} \Rightarrow x_d, y_d, z_d \} \text{ diff's}$

Unlock(L);

P2

invalidate  $\rightarrow$  LOCK(L);

$x, y, z$  at  
lock acquisition  
(CLRC)

$\leftarrow x \leftarrow$   
fetch at  
point of  
access

# LRC with Multi-Writer Coherence Protocol

P1

LOCK(L);

$x \leftarrow \begin{cases} x, y, z \text{ pages} \\ y \leftarrow \text{modified in C.S.} \\ z \leftarrow \end{cases} \Rightarrow \{x_d, y_d, z_d\} \text{ diff's}$

Unlock(L);

P3

LOCK(L);

$x \leftarrow \Rightarrow x'_d$

unlock(L);

invalidate  $\rightarrow$  LOCK(L);  
 $x, y, z$  at  
lock acquisition  
(CLRC)

P2

$\leftarrow x \leftarrow$   
fetch at  
point of  
access

# LRC with Multi-Writer Coherence Protocol

P1

LOCK(L);

$x \leftarrow \begin{cases} x, y, z \text{ pages} \\ y \leftarrow \text{modified in C.S.} \\ z \leftarrow \end{cases} \Rightarrow \{x_d, y_d, z_d\} \text{ diff's}$

Unlock(L);

P3

LOCK(L);

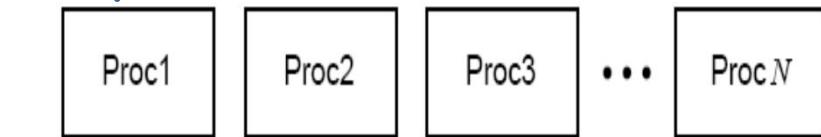
$x \leftarrow \Rightarrow x'_d$

unlock(L);

invalidation  $\rightarrow$  LOCK(L);  
 $x, y, z$  at  
lock acquisition  
(CLRC)  
 $Q \leftarrow$   $x \leftarrow$   $x'_d$   $\leftarrow$   $x$   $\leftarrow$   $Q$   $\leftarrow$   $x$   
fetch at  
point of  
access

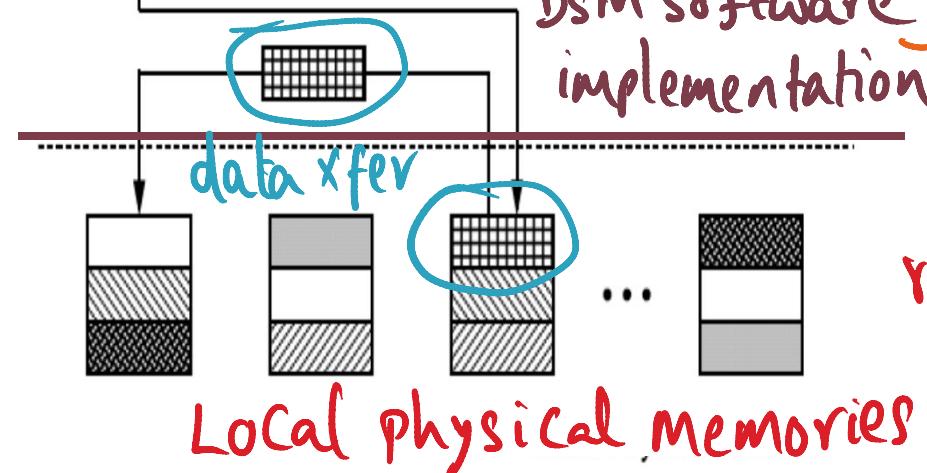
# Software DSM

Application



Global Virtual  
memory  
abstraction

Pagefault



Address space partitioned  
Address equivalence  
Distributed ownership

Contact owner of page  
to get current copy

