# Distributed Shared Memory Model
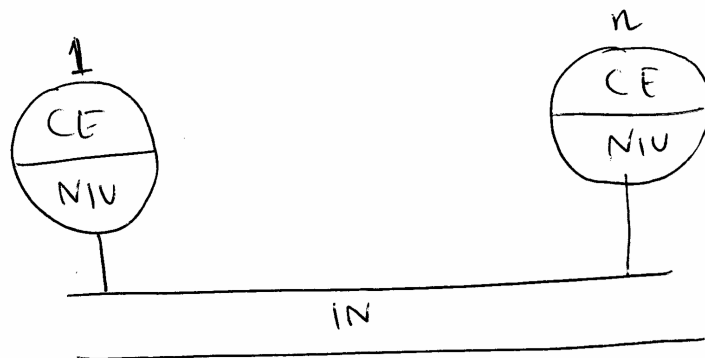
NIU: Network Int. Unit
CE: Computing Element $\{Cache, PE, MM, D\}$

Local components.

- All nodes use a single global address space.

- If a processor issues a <u>read</u> request (load data from memory) it will specify the address from where data is to be retrieved.

address
— a node #
— address of location within the node.

Thus, if the node # is that of the processor issuing the read command, the data is retrieved from the memory attached to it. Otherwise, the node # is used to send a message via IN to the appropriate node's memory from where the data is to be retrieved & delivered to the requesting node.

We summarize the procedure below:

## Procedure A :

Command : Load from specified address into a register

1) Translate address { Addr ← CE# + $\genfrac{}{}{0pt}{}{\text{memory}}{\text{addr in}}$ }
                                                                              CE

2) Is CE local ?

3) If Yes, load from local CE. Else, send request via n/w to remote CE

4) Retrieve data from remote CE's mem.

5) Transport data via n/w to requesting CE

6) Load data in specified reg.

---

## Procedure B:

Command: Store register in address.

(1) Translate address

(2) Is CE local ?

(3) Yes ⇒ done ; Else, send contents of reg via n/w to remote CE.

(4) Store contents of reg in the specified location in the mem of remote CE.

(5) Send message to CE which issued store command that the task is completed.

__formats:__

1) __Load Req__:

(Source CE addr, Dest. CE addr, address in mem from where data is to be loaded)

2) __Retrieval data__ (sent over IN to the requesting CE)

(Source CE addr, Dest. CE addr., data)

3) __Store instruction__ (store req.)

(Source CE addr, Dest. CE addr, Addr. in mem (where data is to be stored), Data)

(4) Ack packet ( back to originator )

( Source CF , Dest. CF , Store
       addr        addr.      successful )

- Error detecting bits may be added to the above packets.

---

## Timing : ( for data retrieval from a remote memory)

- fixed time $T$ needed by the system program @ the host node to issue a command over the n/w. This will include time to decode (to which remote node the req. is to be sent) & formatting a packet.

- Time taken by the load req. packet to travel via IN. This depends on bandwidth of n/w & packet size, i.e., $\left(\frac{n}{B}\right)$ secs, where $n$: pkt size in bytes & $B$ is the bandwidth (bytes/sec)

- Time taken to retrieve word from remote memory, $W$.

- Fixed time $T$ needed by the destination CE system program to access the n/w.

- Time taken by dest. CE to transport the reply packet ( which contains the data retrieved ) over the n/w, i.e., $(m/B)$ sec where $m$ : size of the packet.

Thus, total time $= 2T + q + \dfrac{(n+m)}{B}$

Note: For store operation time taken is similar, however $m$ & $n$ may be different.

Note: Above model does not claim to be an exact model. Can communication & computation be time-overlapped?

Another issue is how frequently service requests to a remote processor are issued by a CE.

# A full-fledged Example

System: NUMA parallel computer

#q nodes: 256 CEs

Each CE has 16MB memory

In a set of programs, 10% of instructions are loads and 15% are stores.

memory access time for
local load/store : 5 clock cycles

Overhead to initiate transmission of a request to a remote CE : 20 clock cycles

Bandwidth of IN : 100 MB/sec

Assume 32 bit words & a clock cycle time of 5 nsec

Now, if 400,000 instructions are to be executed, compute:

1. Load/store time if all accesses are to local CEs

2. Repeat (i) if 25% of accesses are to a remote CE.

Solution:

① No. of load/store instructions

$$= 400,000 \times \frac{1}{4} = 100,000$$

Time to execute load/store locally

$$= 100,000 \times 5 \times 5 = 2500 \,\mu secs \cdots \text{(1)}$$

---

② · #q load instructions : 40,000.  (10%.)

# q local loads : $40,000 \times \frac{3}{4} = 30,000$

Time taken for local loads $\left.\begin{array}{c}\\ \end{array}\right\} = 30,000 \times 25$
$$= 750 \,\mu secs.$$

#q remote loads : 10,000

recv. packet format : $\left(\begin{array}{l} Src : 8 \,bits \\ dest : 8 bits \\ addr : 24 bits \end{array}\right\} 40 bits.$

∴ request packet length = 5 bytes.

Response packet length $= \underline{6\,bytes} \left\{\begin{array}{c} 32\,bits \\ per \\ word \end{array}\right\}$

$\left( Src : 8, \ Dest : 8 ; \ Data : 32\,bits \right)$

Time taken for remote <u>load</u> of one
word =

$$20 \times 5 \times 10^{-9} \quad (\text{overhead, fixed}) \leftarrow$$

$$+ \frac{5}{100} \times 10^{-6} \quad (\text{transmit a req. pack})$$

$$+ 5 \times 5 \times 10^{-9} \quad (\text{data retrieval})$$

$$+ 20 \times 5 \times 10^{-9} \quad (\text{overhead, fixed}) \leftarrow$$

$$+ \frac{6}{100} \times 10^{-6}$$

$$= \underline{335 \text{ nsecs}}.$$

# of requests to remote CE = 10,000

Total time: 3350 μsecs

Time for local loads = 750 μsecs (see page 7)

Total time taken for loads : 3350 + 750

$$= \underline{4100 \text{ μsecs}}$$

\# q Store instructions:

$$400,000 \times 0.15 = 60,000$$

\# q local stores : $60,000 \times 0.75$

$$= 45,000$$

\# q remote Stores = $15,000$

Time for local stores : $45,000 \times 25 = \underline{1125 \; \mu secs}$.

Time taken for $\underline{1}$ remote store :

$$20 \times 5 \times 10^{-9}$$

$$+ \frac{9}{100} \times 10^{-6} \quad ( \text{remote store packet length : 9 bytes})$$

$$+ 25 \times 10^{-9}$$

$$+ 20 \times 5 \times 10^{-9}$$

$$+ \frac{3}{100} \times 10^{-6} \quad (\text{ack packet length 3 bytes})$$

$$\overline{\underline{345 \; nsecs}}$$

Time to store 15,000 words

$$= 345 \times 15,000 = \underline{5175 \; \mu secs}$$

Time for local stores = 1125 $\mu$secs (see page 9)

$\Rightarrow$ Total time for stores =
$$\begin{array}{r} 5175 + \\ 1125 \\ \hline 6300 \; \mu secs \end{array}$$

Total time for loads + stores :
$$\begin{array}{r} 4100 + \\ 6300 \\ \hline 10,400 \; \mu secs \;\; (2) \end{array}$$

$$\frac{\text{Total time for load \& store (local + remote)}}{\text{Total time for load \& store if entirely local}} = \frac{(1)}{(2)} =$$

$$= \frac{10,400}{2500} = \underline{4.16} \; .$$

<u>Inference</u> : ?