# Performance

**Example 2:** Assume that you want to write a program that should achieve a speedup of 100 on 128 processors.
($i$) What is the maximum sequential fraction of the program when this speedup should be achieved under the assumption of strong scalability?
We start with Amdahl's law and then isolate $f$ as follows:

Thus, only less than 1% of your program can be serial in the strong scaling scenario!
($ii$) What is the maximum sequential fraction of the program when this speedup should be achieved under the assumption of weak scalability whereby the ratio $y$ scales linearly?
We now start with Gustafson's law and then isolate $f$ as follows:

Thus, in this weak scaling scenario a significantly higher fraction can be serial!

$$y \text{ scales linearly}$$
$$=$$
$$y = np$$

$$¿a?$$

i) done

ii)   Gustafson's  Law: Weak scalability,  $n = np * n$

$$S = [1-a] \, np + a \qquad \text{with } Y = np, \ S = 100, \ np = 128$$

Plugging  into  the  Law,

$$100 = (1-a) \, 128 + a$$
$$100 = 128 - 128a + a$$
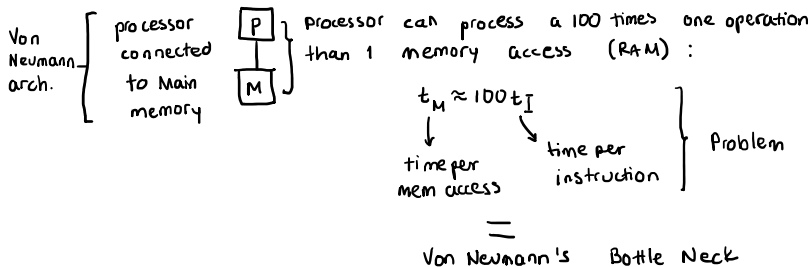$$100 = 128 - 127a$$
$$127a = 128 - 100$$
$$127a = 28$$
$$a = \frac{28}{127} = \boxed{0.22} \rightarrow 22\% \text{ of Time}$$

Speedup → acceleration

## ✱ Memory : Cache

Cache: small memory  since we want quick access
big memory : bigger  distances, electrons
take more  time  to travel = slow

Let's analyze  with 1  processor ( Von  Neumann's Architecture)

Von
Neumann
arch.
$$\begin{bmatrix} \text{processor} \\ \text{connected} \\ \text{to Main} \\ \text{memory} \end{bmatrix}$$
[P]
[M]

Processor can  process  a 100 times  one operation
than 1  memory  access   (RAM) :

$$t_M \approx 100 \, t_I$$

↓ time per mem access
↘ time per instruction

} Problem

$$=$$

Von Neumann's  Bottle Neck
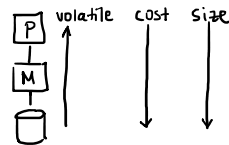
Solution  to  this :  Cache  Memory
↳ another  memory  close to the  processor, commonly  inside its  chip
→ small memory
Fast
=
More cost
than RAM

[P]
[Mᶜ]
[M]

RAM memory  looks for a value  the P asks,
if it's not  on  cache,  it copies it to  cache
so that the processor  access it  from  cache.
the advantage  is ─┐

It is based  in a principle divided in two
a) Temporal:  if x  is  accessed  now  it is  probable
that x is  needed later
b) Spatial:  if x  is accessed  now, it is  probable
that  x's  surroundings  are needed too.
Therefore  you  copy  to cache the  space  around x
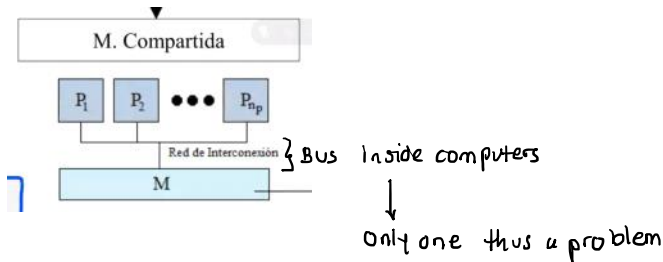(other variables)  in case the processor  needs  it.

If hard drive speed was faster, we wouldn't need RAM Memory:
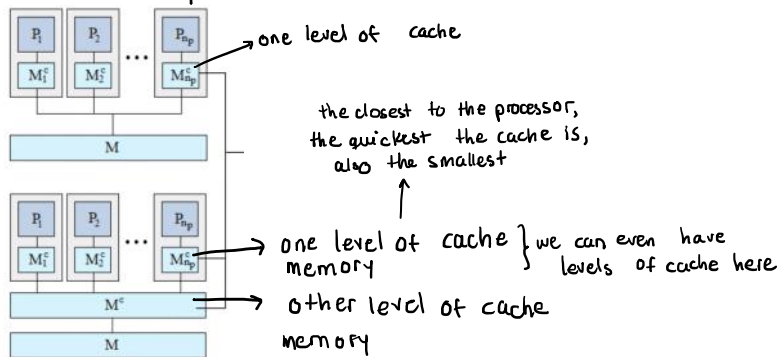


volatile cost size

→ RAM Memory is the cache of the hard drive.
→ The source program is distributed into pieces between RAM and cache.
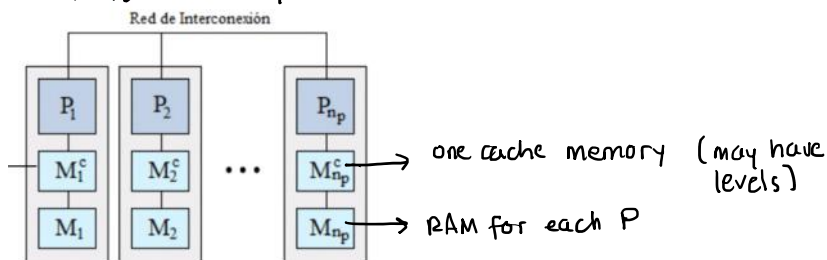
→ For the case of ∞ processors:



Bus Inside computers
↓
only one thus a problem

Solution: cache memory, in two ways
Shared Memory



→ one level of cache
the closest to the processor, the quickest the cache is, also the smallest

→ one level of cache memory } we can even have levels of cache here
→ other level of cache memory

Distributed memory



→ one cache memory (may have levels)
→ RAM for each P
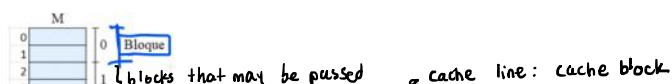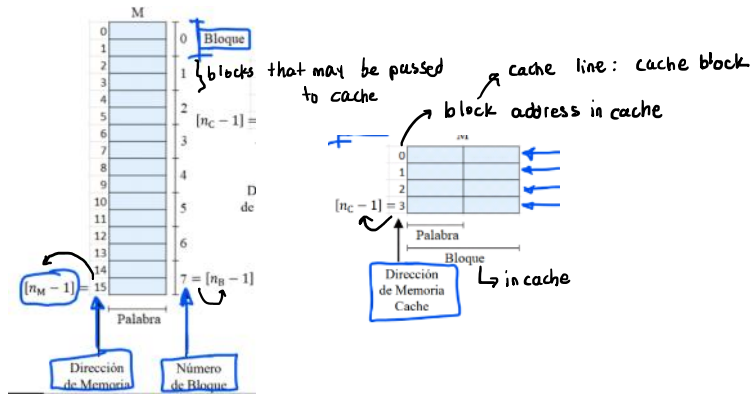
How is memory organized?
For 1 processor:

→ Memory usually is organized in words, which is a block with a number that is the address (locality number
   ↳ To respect the spatial principle: we copy pieces of ram to the cache, and this means we need to divide RAM into blocks
      ↓
      sets of words } also receive an address



blocks that may be passed ≃ cache line: cache block

M

0 Bloque

blocks that may be passed to cache

cache line: cache block

block address in cache

$[n_C - 1] =$

D de

$[n_M - 1] = 15$

Palabra

$7 = [n_B - 1]$

Dirección de Memoria

Número de Bloque

Dirección de Memoria Cache

Palabra

Bloque

in cache

Now the cache memory has also lines, but each line is a block : they are usually called cache lines

↳ in this case has two words space

↳ each RAM block goes to one block / cache line.

If x data value is required by the processor, RAM copies that x word's block to cache.

In the case of ∞ processors, i.e., Shared Memory



$M_1^c$          $M_{n_p}^c$

$[n_C - 1] =$

Palabra

Bloque

Cache Memories (n for n processors)

Dirección de Memoria Cache

Same idea: if x is required, we copy Block 4 into Cache (x and y)

the bigger the block, the more probable it is to find something in the cache

M

0 Bloque

$[n_M - 1] = 15$

Palabra

$7 = [n_B - 1]$

Dirección de Memoria

Número de Bloque

In distributed Memory:

$M_1^c$

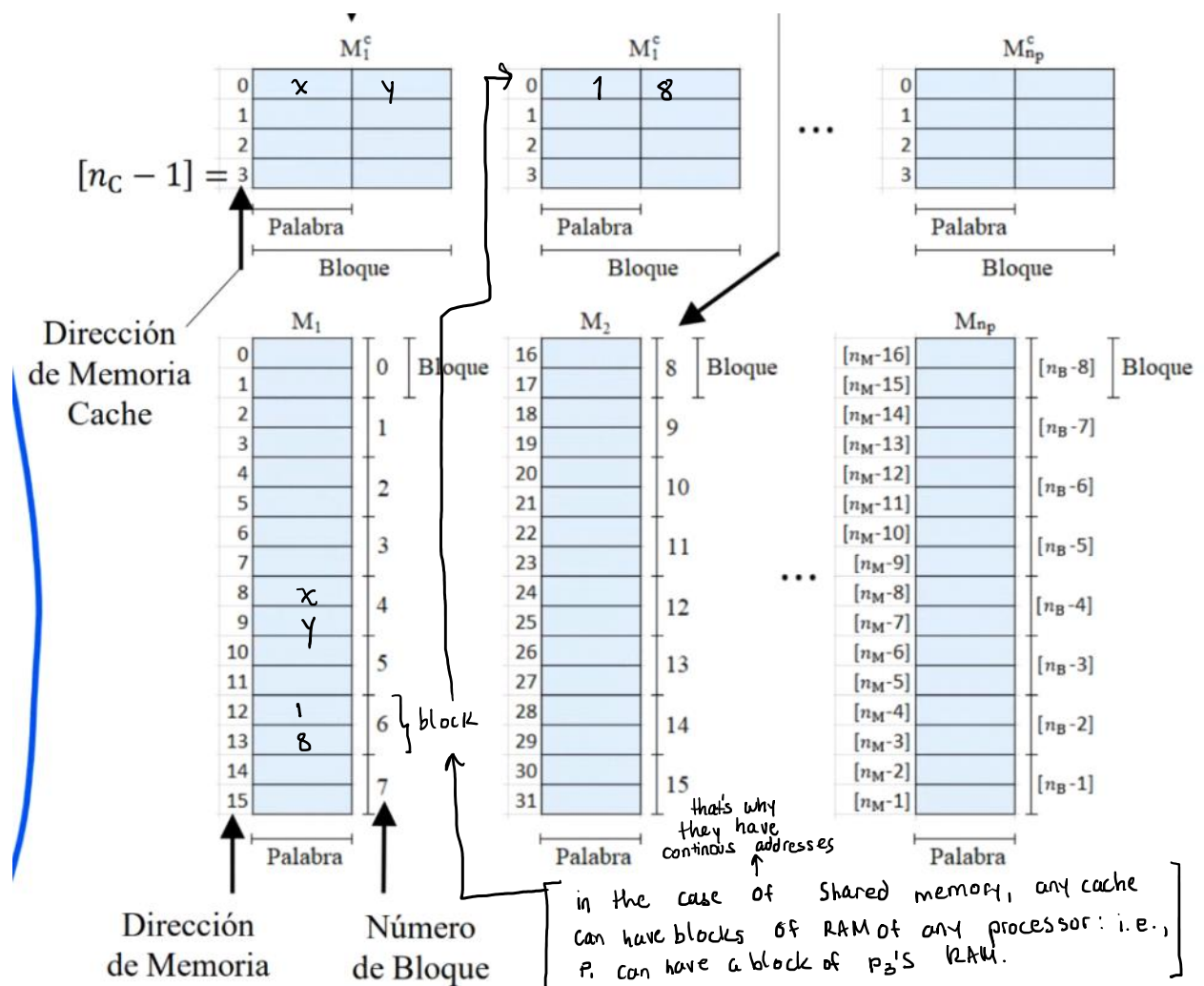| | |
|---|---|
| 0 | x ∥ y |
| 1 | |
| 2 | |
| $[n_C - 1] = $ 3 | |

Palabra
Bloque

$M_1^c$

| | |
|---|---|
| 0 | 1 ∥ 8 |
| 1 | |
| 2 | |
| 3 | |

Palabra
Bloque

· · ·

$M_{n_p}^c$

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |

Palabra
Bloque

Dirección
de Memoria
Cache

$M_1$

| | | |
|---|---|---|
| 0 | | 0 ] Bloque |
| 1 | | |
| 2 | | 1 |
| 3 | | |
| 4 | | 2 |
| 5 | | |
| 6 | | 3 |
| 7 | | |
| 8 | x | 4 |
| 9 | y | |
| 10 | | 5 |
| 11 | | |
| 12 | 1 | 6 } block |
| 13 | 8 | |
| 14 | | 7 |
| 15 | | |

Palabra

$M_2$

| | | |
|---|---|---|
| 16 | | 8 ] Bloque |
| 17 | | |
| 18 | | 9 |
| 19 | | |
| 20 | | 10 |
| 21 | | |
| 22 | | 11 |
| 23 | | |
| 24 | | 12 |
| 25 | | |
| 26 | | 13 |
| 27 | | |
| 28 | | 14 |
| 29 | | |
| 30 | | 15 |
| 31 | | |

Palabra

· · ·

$M_{n_p}$

| | | |
|---|---|---|
| $[n_M-16]$ | | $[n_B-8]$ ] Bloque |
| $[n_M-15]$ | | |
| $[n_M-14]$ | | $[n_B-7]$ |
| $[n_M-13]$ | | |
| $[n_M-12]$ | | $[n_B-6]$ |
| $[n_M-11]$ | | |
| $[n_M-10]$ | | $[n_B-5]$ |
| $[n_M-9]$ | | |
| $[n_M-8]$ | | $[n_B-4]$ |
| $[n_M-7]$ | | |
| $[n_M-6]$ | | $[n_B-3]$ |
| $[n_M-5]$ | | |
| $[n_M-4]$ | | $[n_B-2]$ |
| $[n_M-3]$ | | |
| $[n_M-2]$ | | $[n_B-1]$ |
| $[n_M-1]$ | | |

Palabra

Dirección
de Memoria

Número
de Bloque

that's why
they have
continous addresses

in the case of shared memory, any cache
can have blocks of RAM of any processor: i.e.,
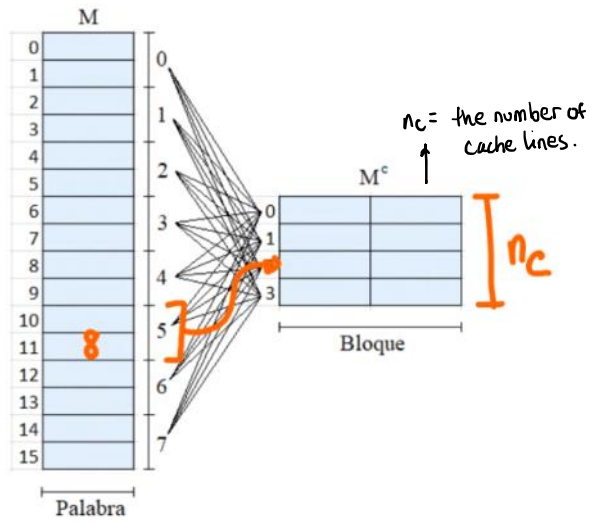$P_1$ can have a block of $P_3$'s RAM.

→ We have one RAM and one cache per processor.
→ If $P_1$ wants address 8 data, say x. It is in block 4, so it is copied to cache (the block)
   into its cache, the processor's cache($P_1$'s cache, in this case)

→ What happens if the cache gets full? We delete the cache block that has the most time
   inside the cache. In this way, we are allowing each <u>cache block</u>, to be there.
   
   <u>cache line</u>

→ The idea of a distributed memory system the idea is to make memory look like one:
   that's why they have continous memory addresses, and that's why we can copy to any
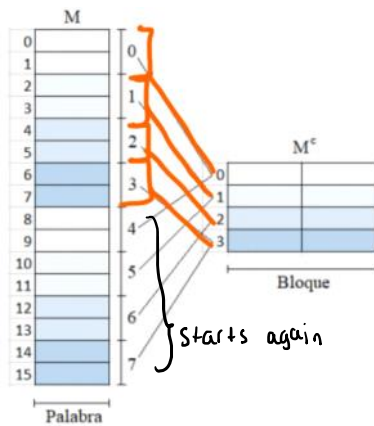   processor's cache.

<u>Memory mappings</u>

Mapping: defines in which cache address will a Ram address be copied in

Above, we used the Completely Associative Mapping, where $n = n_c$
The other extreme is: Associative of 1 via

M / Palabra

Mᶜ / Bloque

$n_c$ = the number of cache lines.

In the direct case and its intermediate case : the first RAM block always goes to the first cache line, and so on. When the cache is full the following RAM address start again in cache line Ø.



M / Palabra

Mᶜ / Bloque

Starts again

As we go further, the electronics and cost increase towards Completely Associative