

## Basics of caches

Tuesday, April 14, 2020 6:28 PM

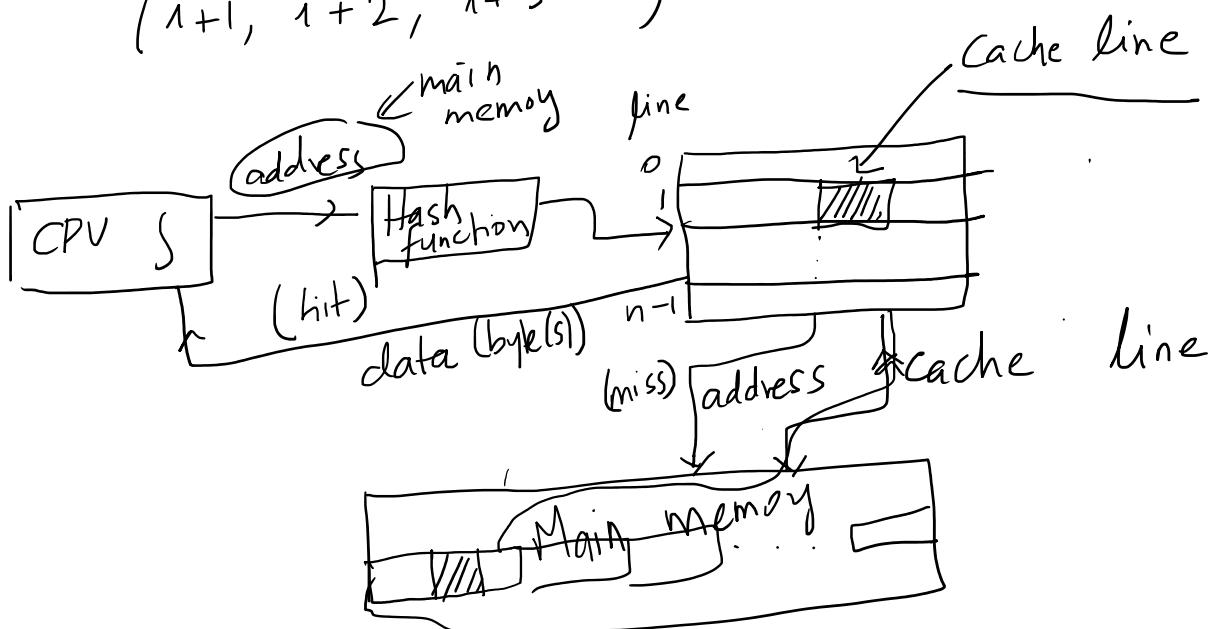
- CPU ↑ fast
  - Main memory ↓ slow
- 
- How do we close gap?
- Cache ↴ <sup>litude</sup> small, fast, on-chip memory
  - Advantage of **locality**\*
- Spatial locality      Temporal locality

Example : SAXPY loop

```
for (i=0; i< N; i++)
    y[i] = a X[i] + y[i];
```

Spatial locality:

( $i+1, i+2, i+3 \dots$ )



Cache line:

Multiple bytes long  
For example, 128 bytes  
32 Floating point items  
4 bytes

\* Minimum unit of data that is moved between main memory and Cache.

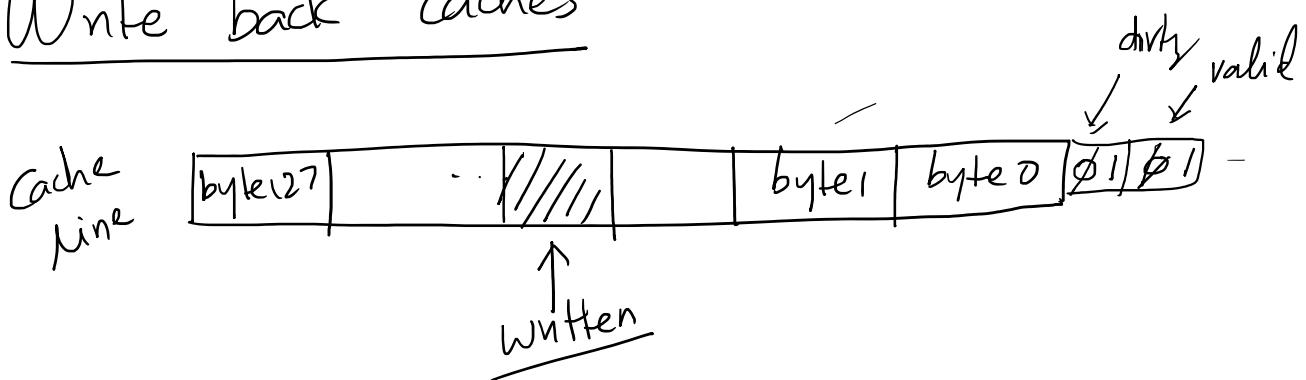
Cache placement problem

- Direct mapped (no flexibility)
- Set associative (flexibility)

Cache replacement problem

- Direct mapped (no flexibility)
- Set associative (Least recently used)

Write back caches



## Write through caches

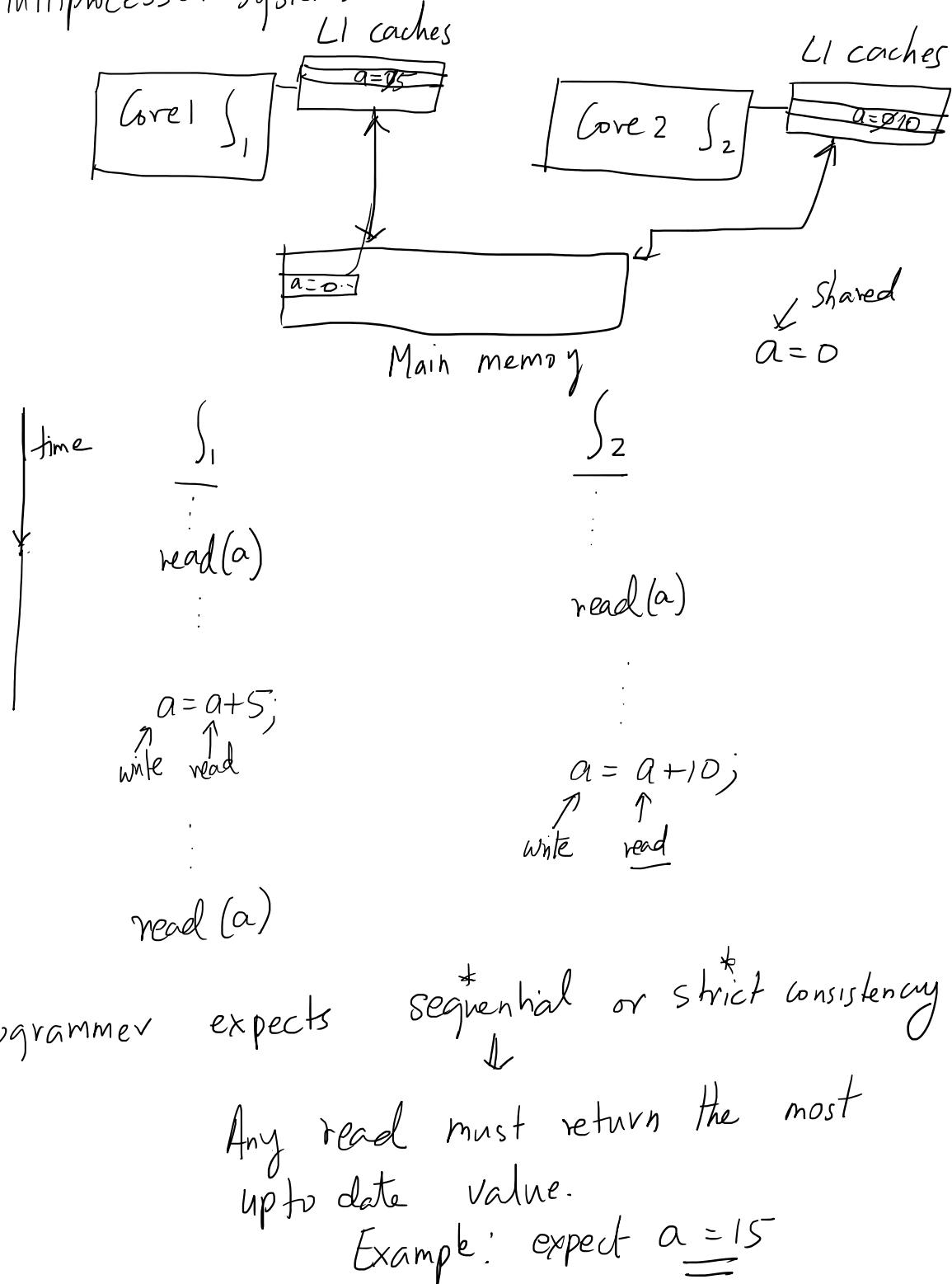
- Writes are made to cache line + propagated to main memory
  - ↓  
Contents of cache always consistent with main memory
  - ↓  
Slow.
- We will assume write back cache.

\* Cache line || \*  
\* Write back || =

## The cache coherence problem

Tuesday, April 14, 2020 5:52 PM

- Multiprocessor systems.



Any read must return the most up-to-date value.

Example: expect  $\underline{\underline{a = 15}}$

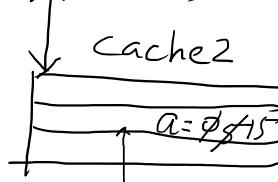
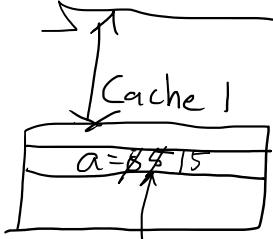
How to achieve strict consistency?

Write  
invalidation  
protocol.

Write Invalidate protocol

bus

S<sub>1</sub>



S<sub>2</sub>



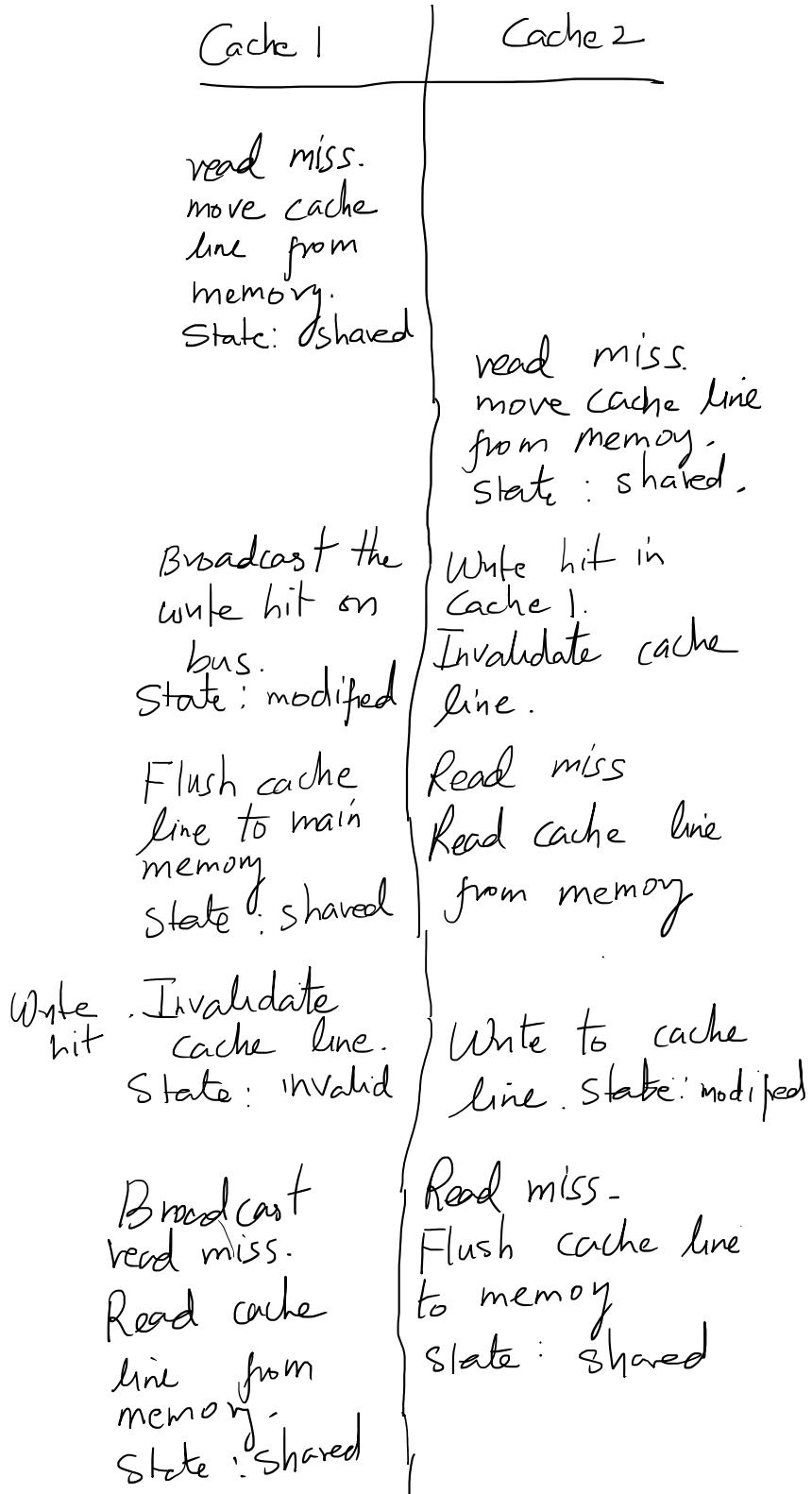
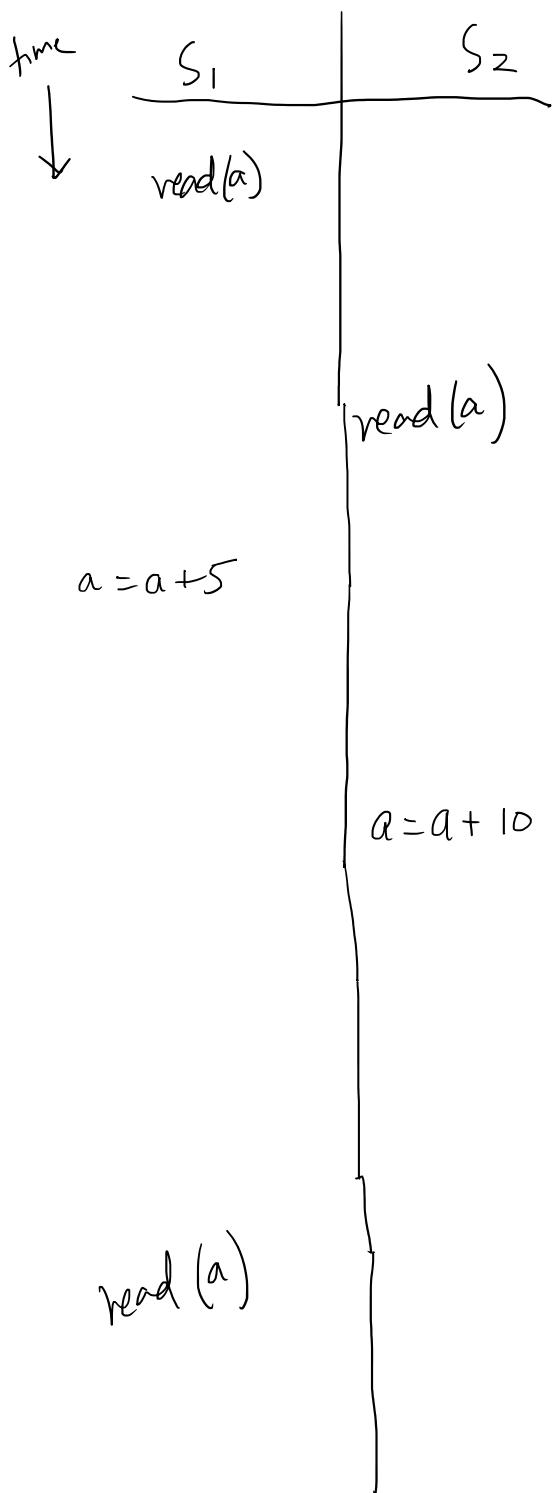
a: shared variable  
 $a = 0;$

Main memory

Snooping based protocol.

Cache activity of all controllers is globally visible

→ read hits/misses  
write hits/misses



## Key idea:

- Only one cache maintains a cache line in the modified state
  - All other caches invalidate that cache line.

## Overhead:

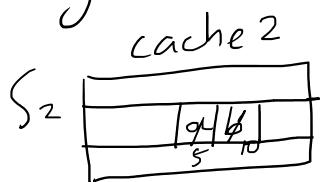
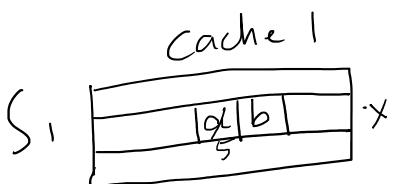
$\sim 64$  cores

- Since adding more cache controllers to the bus significantly increased traffic.

False sharing  $\rightarrow$  impacts performance of parallel code.

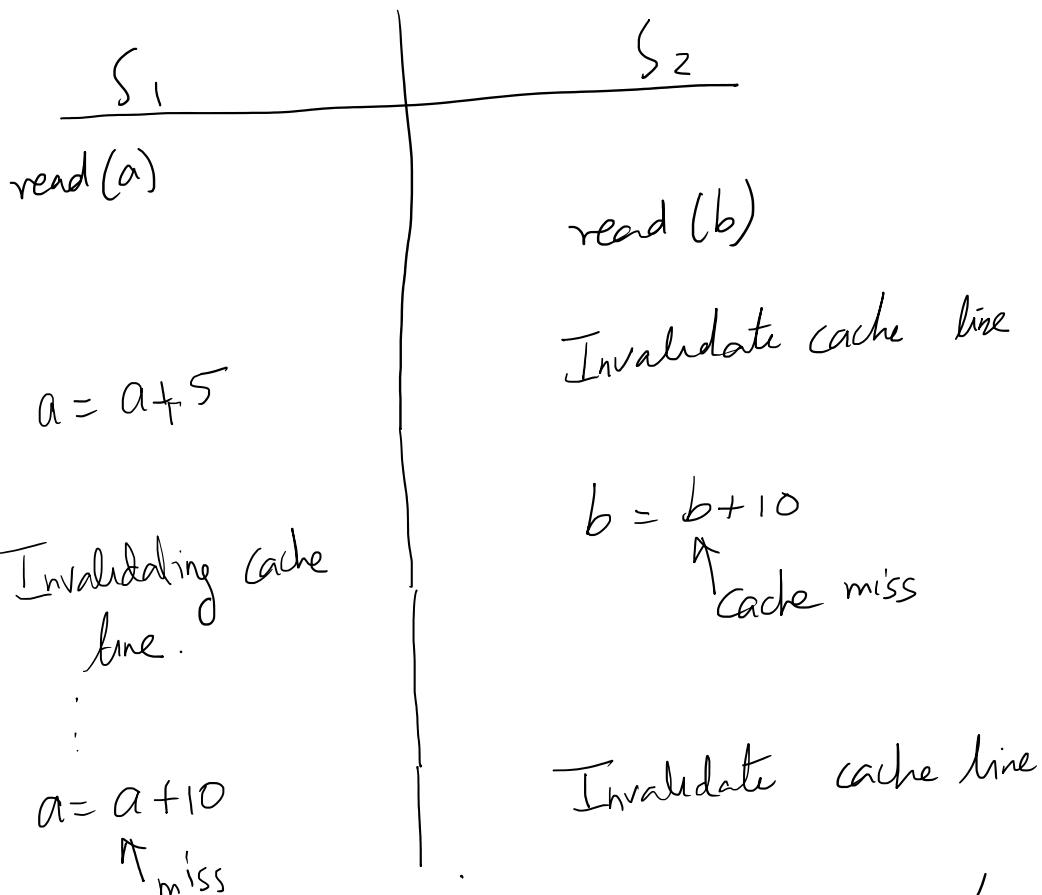
## False sharing

- We are using write invalidate protocol to ensure strict consistency.



$$a = b = 0$$

$s_1$  and  $s_2$  share a cache line  
 $a$  is used by  $s_1$        $b$  is used by  $s_2$        $\Rightarrow$  Threads don't share variables ; but the cache line .

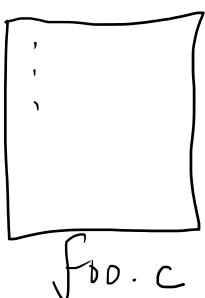


$s_1$  and  $s_2$  keep invalidating each other's cache lines if they continue to write to variables  $a, b$ .

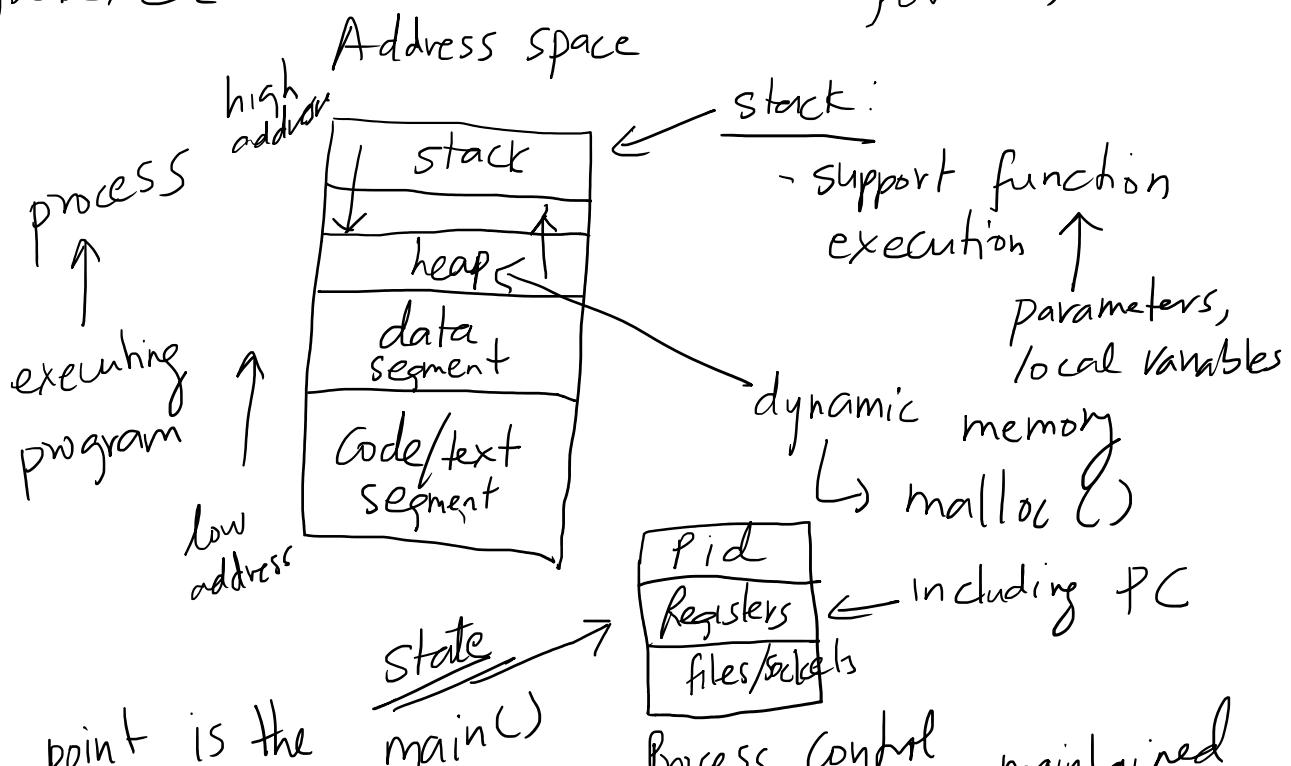
Don't strictly share a variable(s). They share a cache line.

Performance : it will be worse than if we had no cache at all.

Examine ~~your~~<sup>threads'</sup> data access patterns to  
make sure that eliminates / minimize false  
sharing.

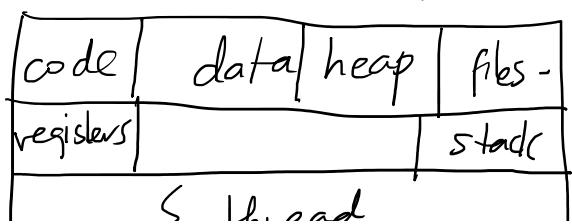
PthreadsSingle-threaded process

\$ ./foo.exe ↴

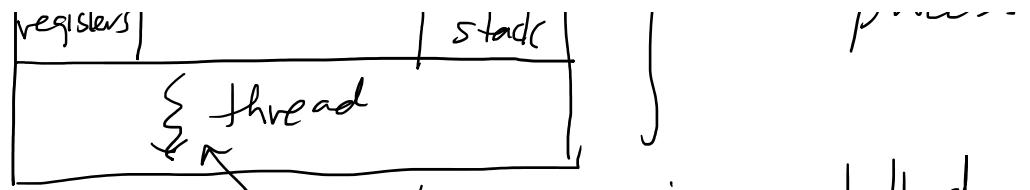


- Entry point is the function main()
- PC value is loaded....

Address space

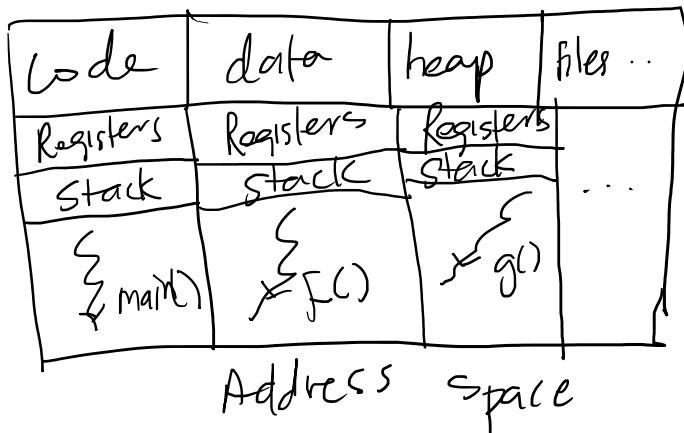


} Single threaded process



Execution path is controlled  
by program counter (PC)  
↑  
address of instruction to  
be execution.

## Multi-threaded process



```
main() {
}
→ create thread( , f, )
→ create thread( , g, )
```

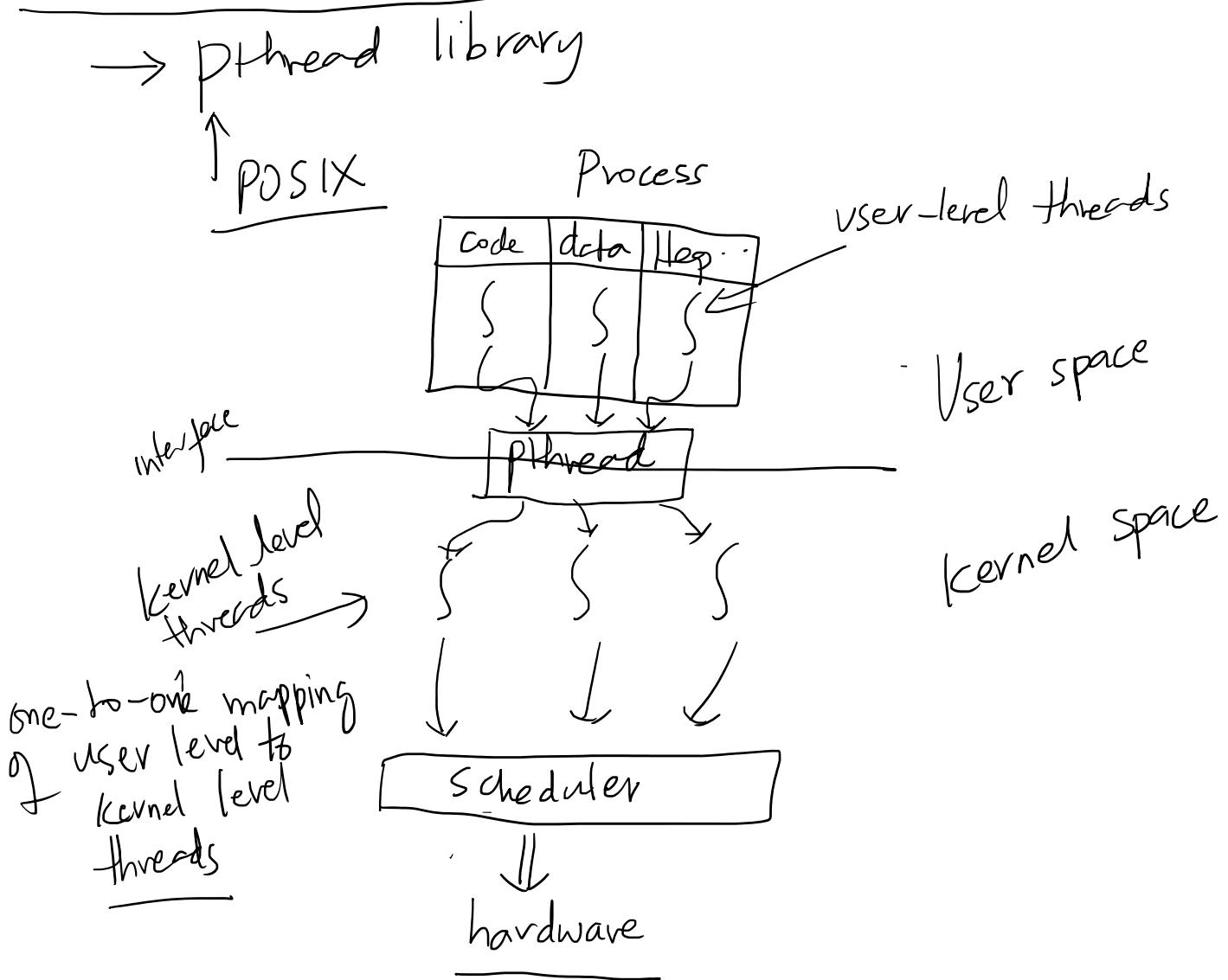
## Threads share:

- Code segment
- Data segment
- Heap segment

## Separate:

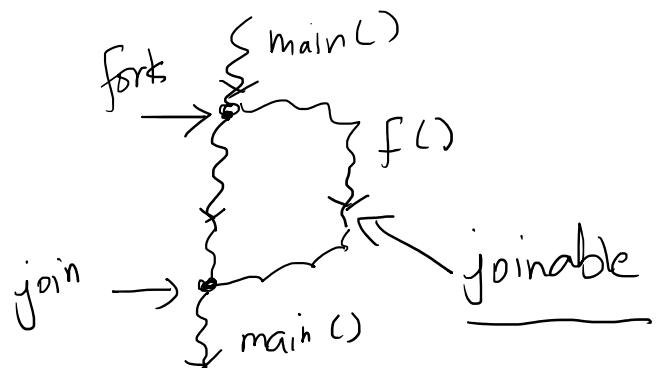
- Registers
- Stack

# How do we create threads?



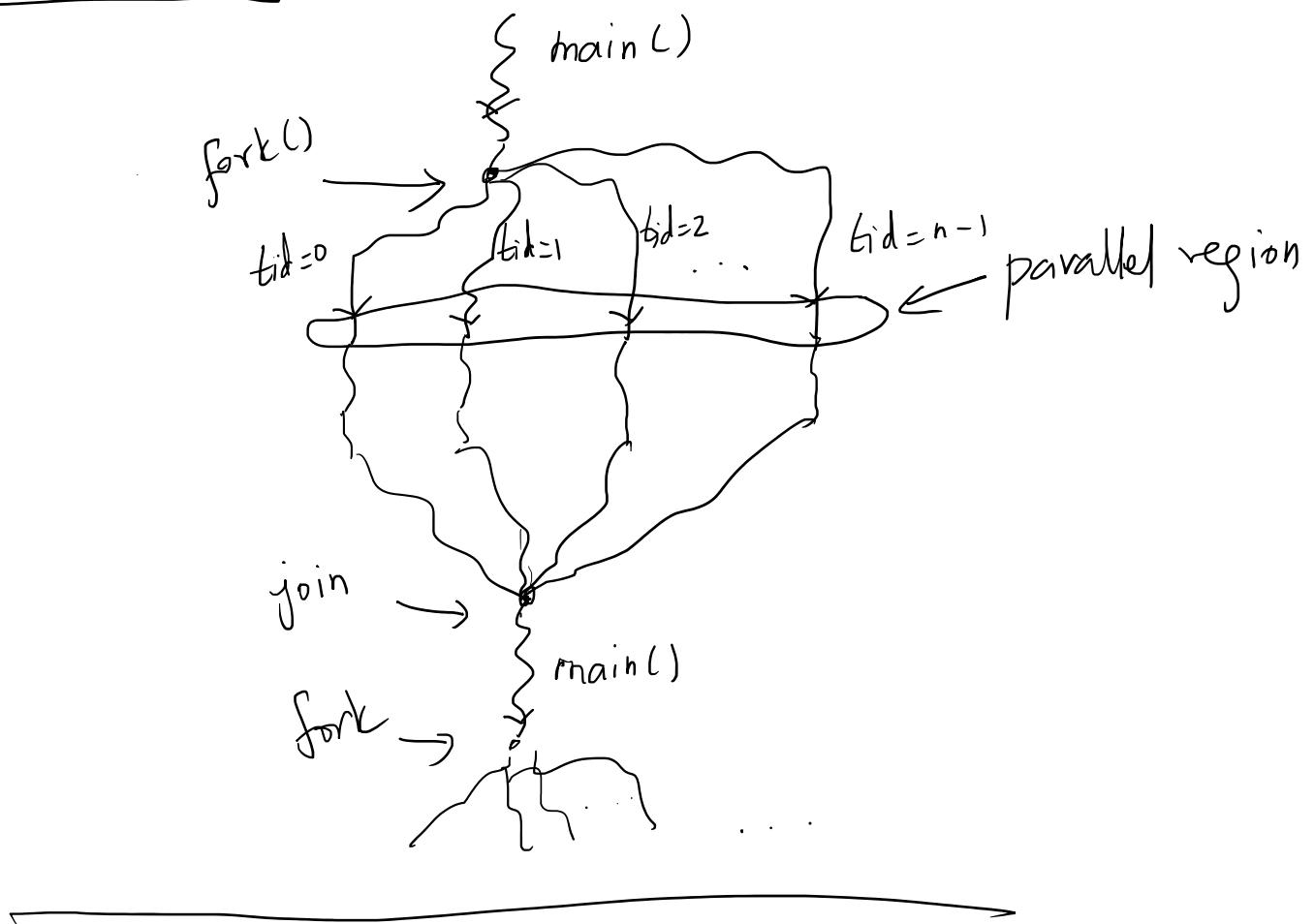
`pthread-create (tid, , f, )`

`pthread-join (tid, status)`



Every `pthread_create()` must have  
a corresponding `pthread_join()`.

General case :

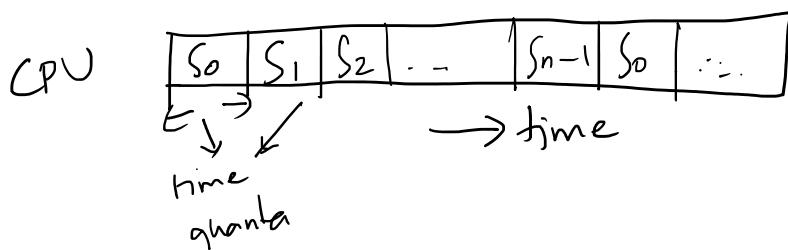


## Concurrent execution

$s_0, s_1, s_2, \dots, s_{n-1}$

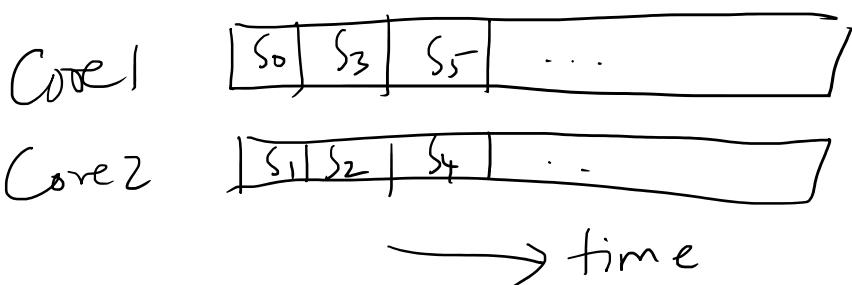


Time share the single CPU between threads.



## Multi-core

→ Space / time sharing



pthread code ...

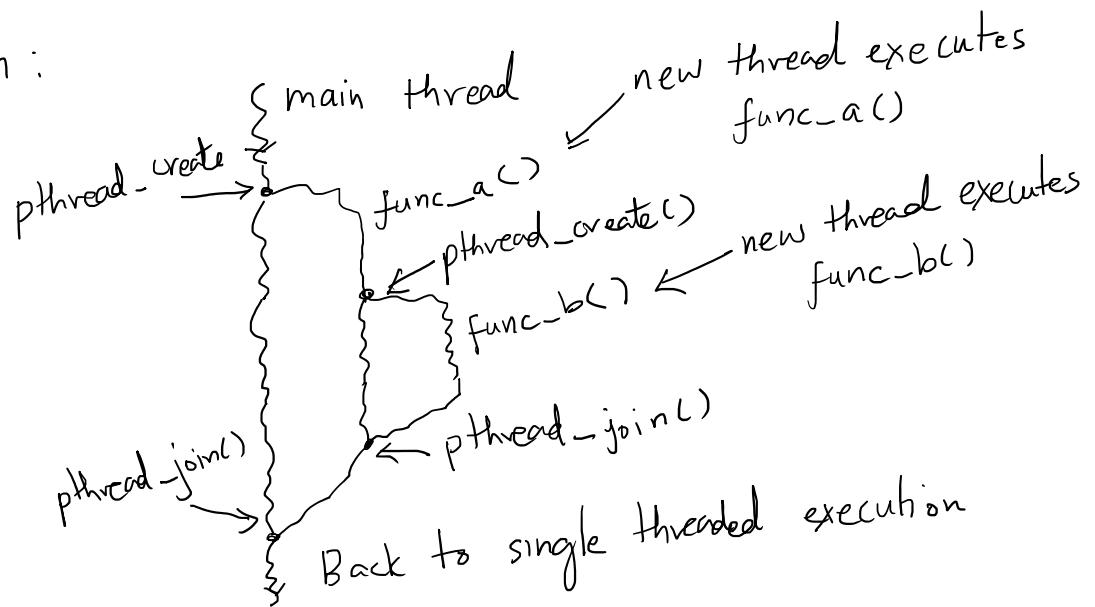
## Thread creation and management

Thursday, April 16, 2020 9:22 AM

Show the use of `pthread_create()` and `pthread_join()`.

### Simple\_thread.c

Visualization:

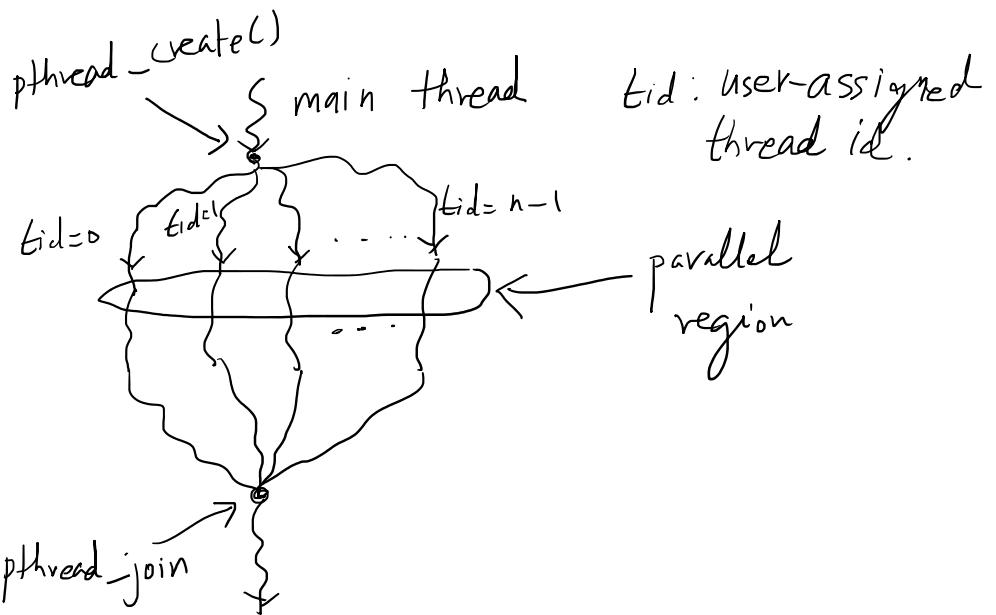


Key things to note:

- No parameters are passed to `func-a()`
- Integer parameter is passed to `func-b()`
  - ↑  
type cast to `(void *)` before  
passing
  - type cast `(void *)` back to `(int)`  
within the function

## multiple\_threads.c

Visualization:



Key points:

- Every `pthread_create()` must have a corresponding `pthread_join()` if the threads are meant to be "joinable".
- Multiple arguments are passed to threads in the form of a structure
  - ↑
    - pack arguments into fields of a structure;
    - pass pointer to structure to the function.
- Every thread executes the same function.  
Only thing different between the threads are the arguments passed to them.
  - ↳ Example of SIMD code

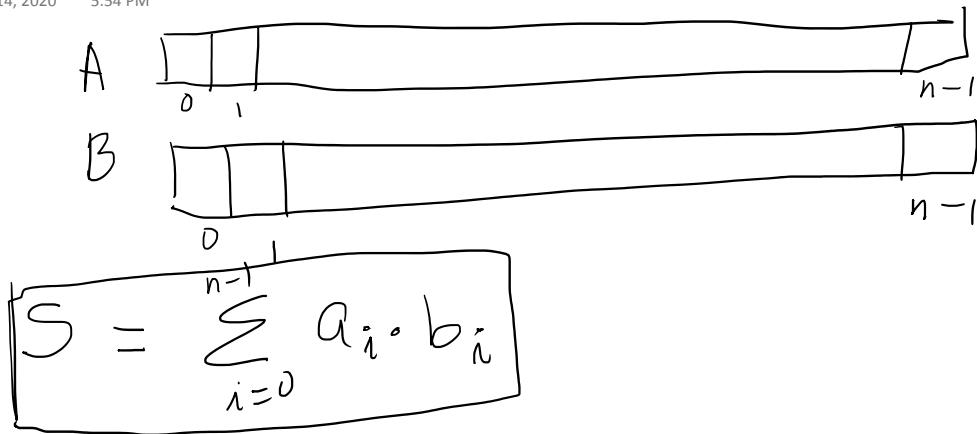
# Writing parallel programs using pthreads

Tuesday, April 14, 2020 5:53 PM

- Vector dot product
- Vector matrix multiplication

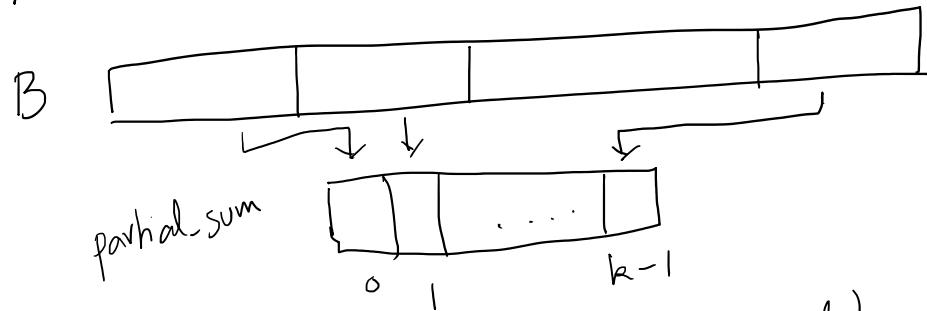
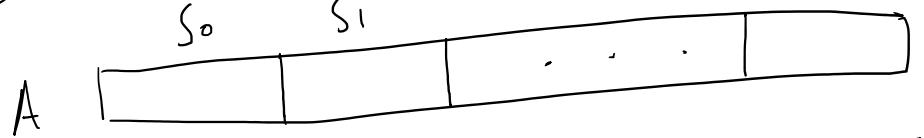
## Vector dot product

Tuesday, April 14, 2020 5:54 PM



### Version 1 : "Chunking method"

- Create  $k$  threads ; each thread calculates partial sum on a chunk.



- Each thread gets its thread id (tid)

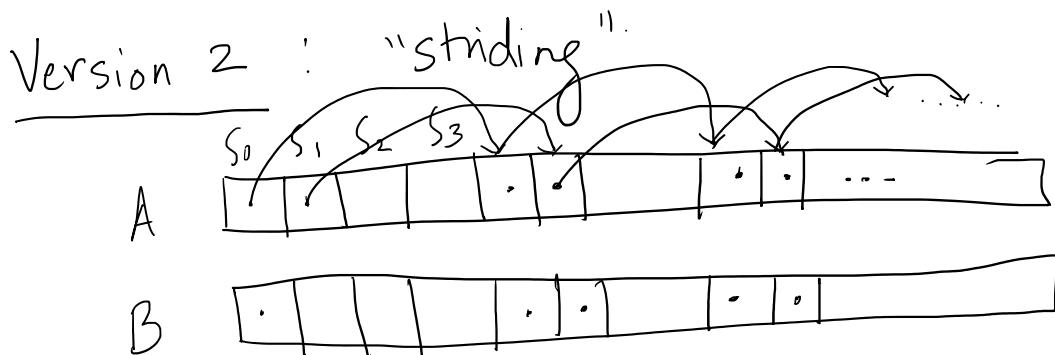
$$\text{Chunk size} = \left\lfloor \frac{n}{k} \right\rfloor \xrightarrow{\text{floor operator}}$$

- Offset within arrays A and B for each thread =  $\text{tid} \times \text{chunk\_size}$

- Threads 0 through  $k-2$  compute partial sum over  $(\cancel{\text{tid}} \text{ offset})$  to  $(\cancel{\text{tid}} \text{ offset} + \text{chunk\_size})$

$$\text{offset} = \text{tid} \times \text{chunk size}$$

- Last thread  $tid = k-1$  may need to process more elements. That is,  
~~( $\leftarrow$  offset to  $n-1$ )~~
- Main thread generates the final sum by accumulating the partial sums.



Each thread strides over elements of A and B, calculating partial sum.

Example:

```

tid ← thread ID ; partial_sum = 0
k ← stride length
while (tid < n) {
    partial_sum += A[tid] × B[tid];
    tid += k;
}

```

Main thread accumulates the partial sums.

## Matrix-vector multiplication

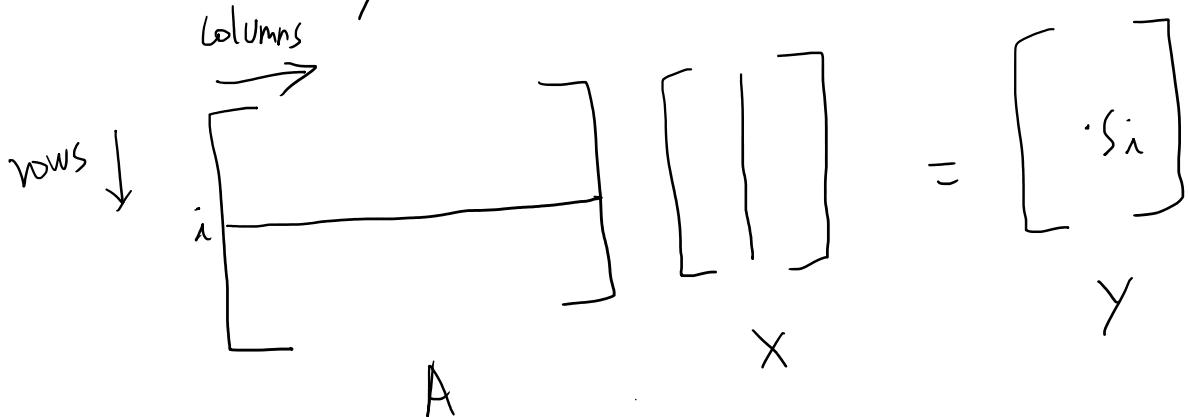
Tuesday, April 14, 2020 5:53 PM

$$Ax = y$$

A is a  $n \times n$  matrix (to keep discussion simple)

X is  $n \times 1$  vector

y is  $n \times 1$  result vector

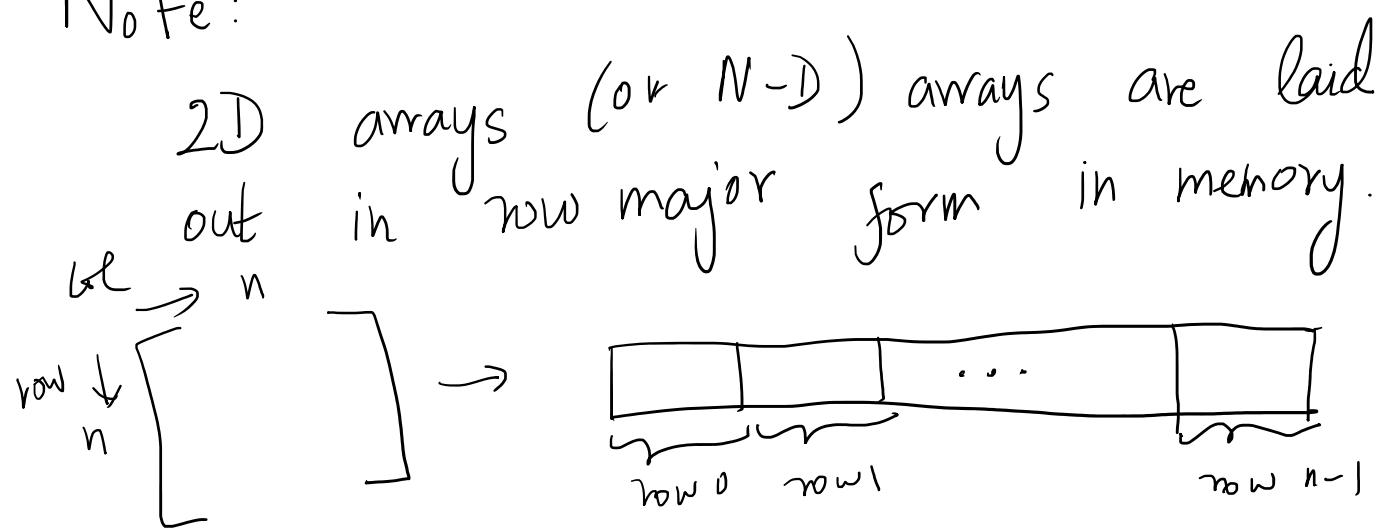


- Problem decomposition
  - Each element of y can be calculated in parallel.

$y[i] = \text{dot product of } i^{\text{th}} \text{ row of } A$   
and X vector

- Can use concept of chunking or striding to calculate each output value.

Note:



So  $A[i][j] \rightarrow A[i \times n + j]$

The diagram shows the mapping of a 2D index  $(i, j)$  to a 1D index  $i \times n + j$ . On the left,  $A[i][j]$  is shown with arrows pointing from  $i$  to the label "row" and from  $j$  to the label "col". On the right,  $A[i \times n + j]$  is shown with an arrow pointing from the expression to the label "width of matrix".

## Exceptions

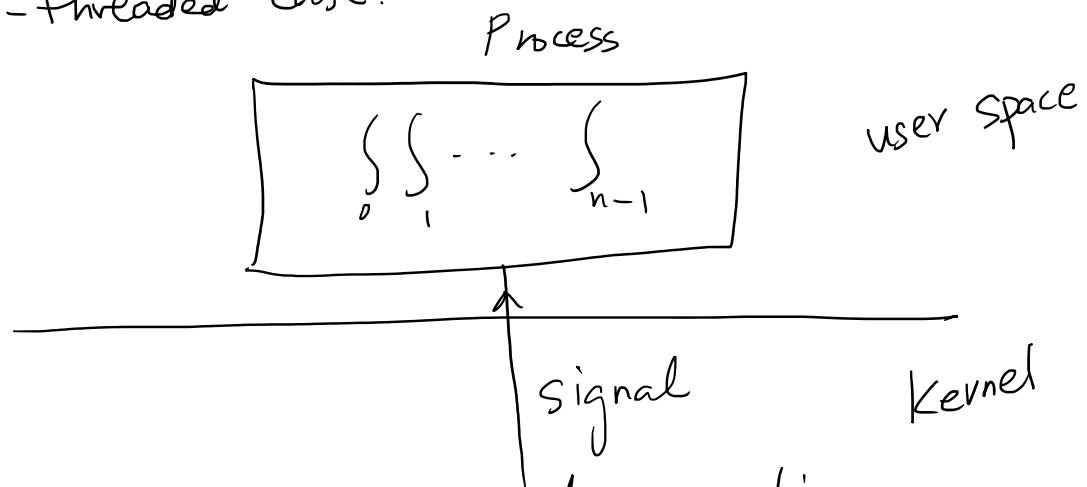
Wednesday, April 15, 2020 10:01 AM

A process can cause various types of exceptions over course of its execution.

Example :

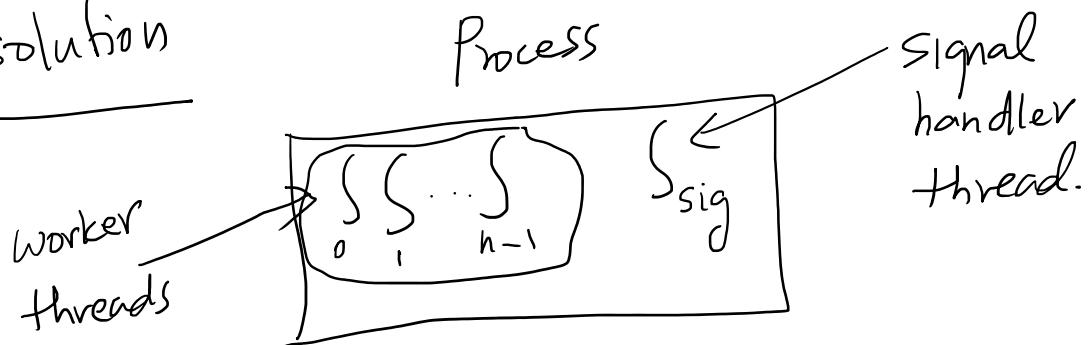
1.  $a = b + c \leftarrow$  value of  $a$  could be too large to store in defined data type
2.  $a = b / c \leftarrow$  what happens if  $c=0$ ?
  - (1) and (2) generate SIGFPE exceptions
- Floating point exception .
3.  $\text{int } A[100];$ 
  - $A[100] = \dots$  out of bounds
  - $\uparrow$  generates segmentation fault
4. Other signals :
  - . Control + C to terminate process
  - $\uparrow$  SIGINT signal

In multi-threaded case:



Any thread may generate exception.  
Kernel posts signal to all threads.

Typical solution



- Designate and set up one thread to be the "signal handler" thread.
- Other threads ignore signals posted to them.

Note: the process cannot ignore certain signals

e.g.: SIGKILL



Used by kernel to terminate  
a "misbehaving" process.

See signal-handler-thread.c for code.