

cores: 4

logical processors: 8

private: variable inside and outside parallel environment
are different

=> inside parallel environment, need to specify
value of this variable

- any variable created inside II region
is automatically private

first private: in parallel environment, variable is initialized with value
outside of environment but result will not copy outside
of environment

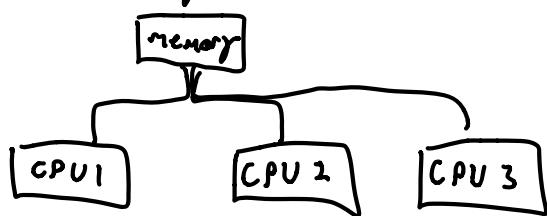
last private: variable must be initialized inside parallel environment
result will be copied to value outside parallel env.

shared: shared outside and inside of parallel environment

Notes on parallel programming and openMP

parallel architectures and programming models

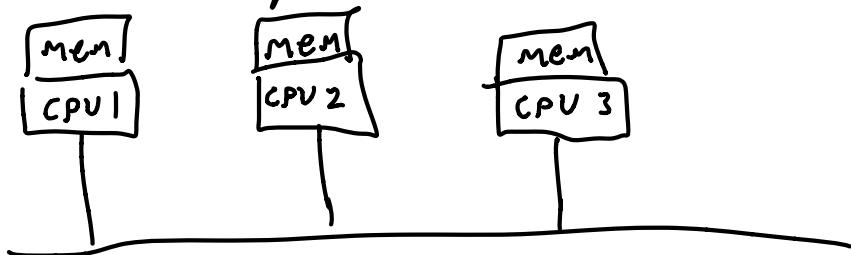
shared memory



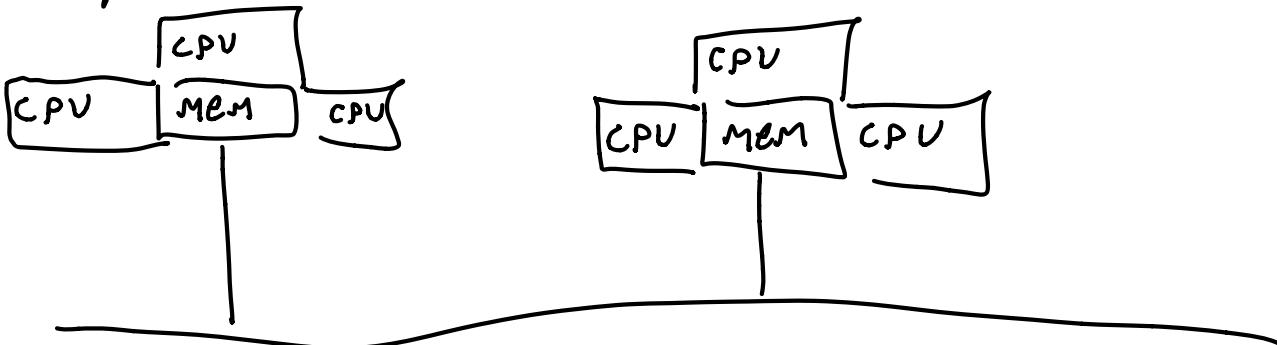
issue: not scalable, limit to # of CPUs

=> this is what openMP works with

distributed memory

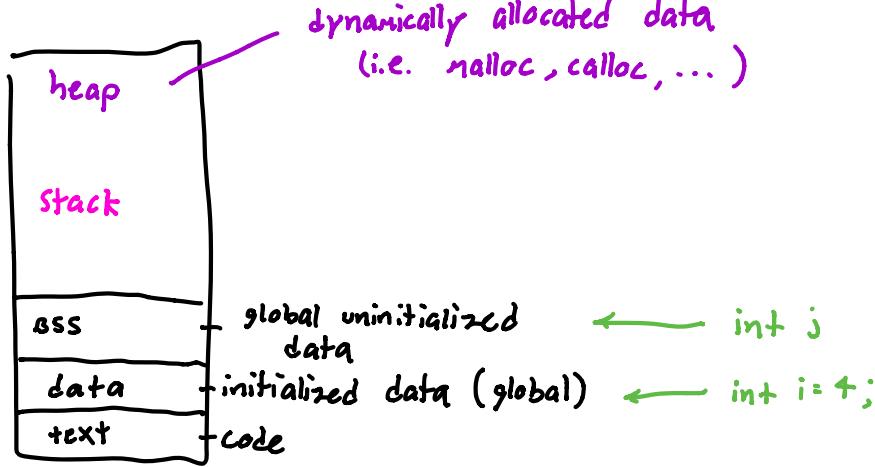


hybrid



Thread into

when you run a code, it has access to some memory that is divided into regions



```
function main
{
```

```
    int x;
```

```
=
```

```
=  
g(x);
```

```
=
```

```
}
```

```
g(int x) {
```

```
=
```

```
=  
int y
```

```
=
```

```
}
```

program counter stores the line we are currently computing... but when g(x) is called, counter must go down to bottom line and the location of counter before it started searching for g is saved in STACK

STACK also stores local variables
i.e. variables created in functions

Multiple threads and multitasking

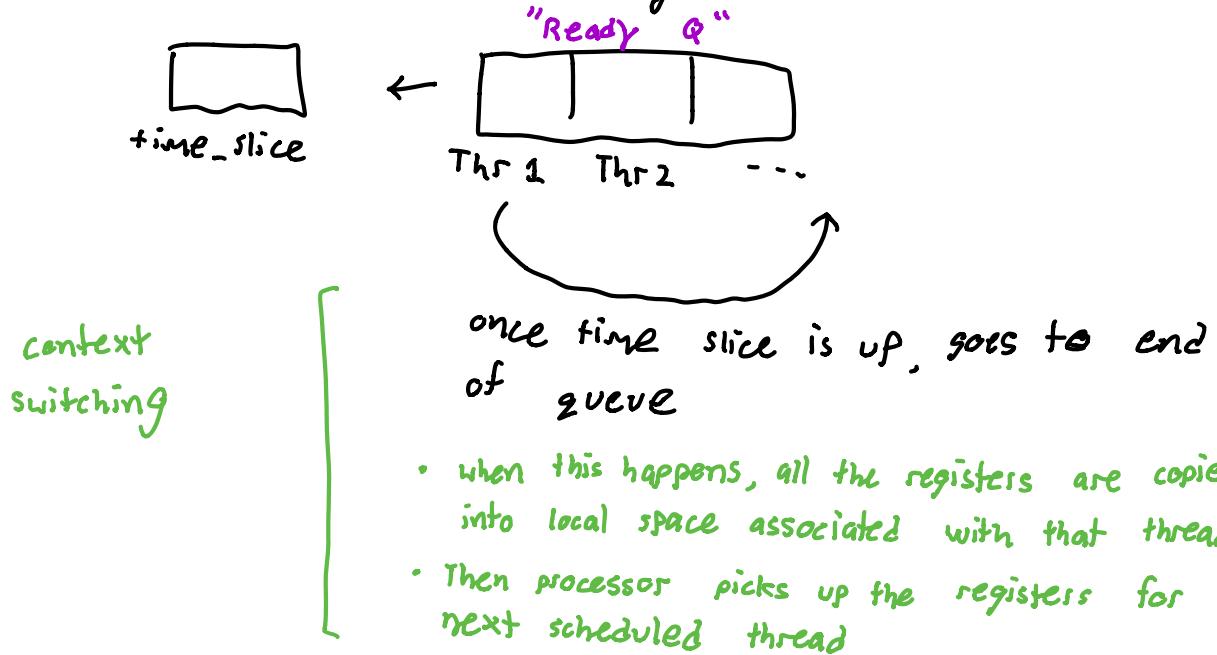
each thread computes independent data

⇒ each thread has its own stack

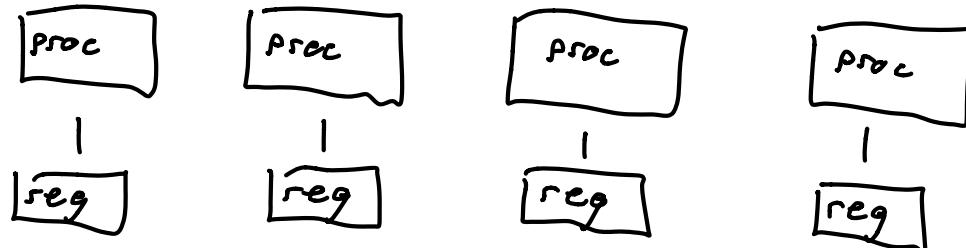
context switching

each processor has fixed time-slice

each thread can only be continued for a max time of time-slice
then next thread in line gets computed



multicore processing



each processor has its own register set

Basic thread functions

each thread has an ID (starting with zero)

`int tid = omp_get_thread_num();` current thread
 $\in \{0, 1, 2, \dots, 7\}$

`int numt = omp_get_num_threads();` total # of threads
 $= 8$

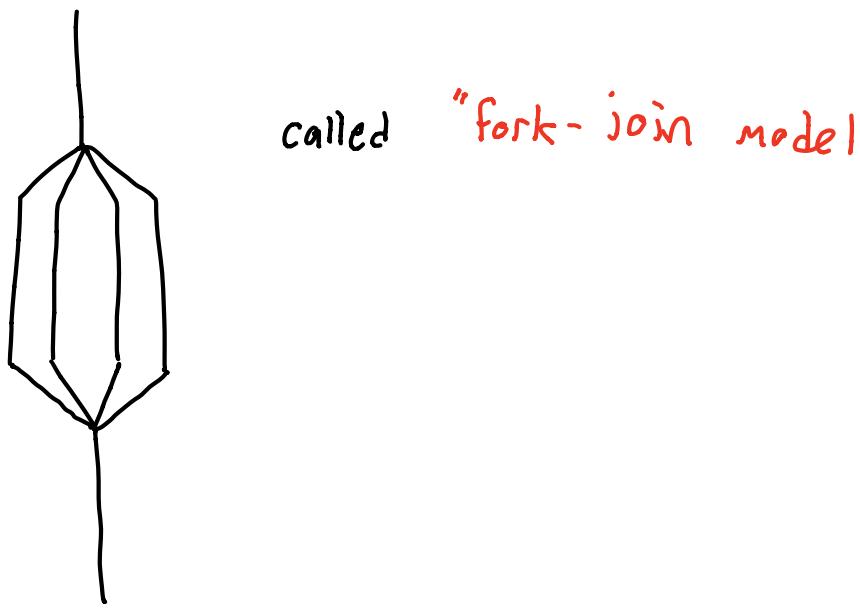
atomic operation: appear as if nothing else is happening within that operation

* if I have a print statement in each of the threads, the prints could happen at the same time
(i.e. print Hello world to world)

but there are ways to make this atomic => locks, mutual exclusion

Forking

going from a master thread to launching additional threads



sequential code and consistency model

compilers sometimes move around lines of code when computing for optimization purposes

X = a
for (i=0; i < n; i++)
 ~
 ~

sometimes will move this if not in for loop

sequential consistency model



processor 1

line 1
line 2
line 3

processor 2

line 1
line 2
line 3

processor 3

line 1
line 2
line 3

operations of each individual processes appear in the order specified by its program

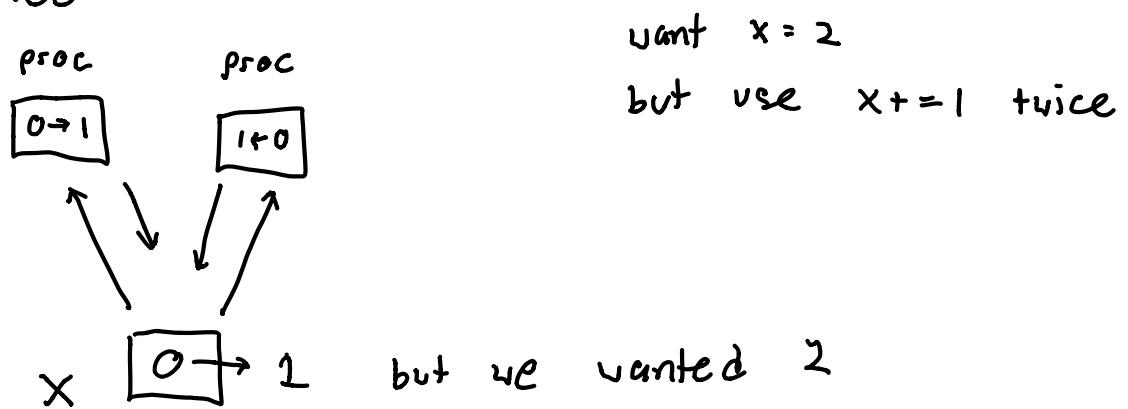
i.e.

line 1
line 2
line 1
line 1
line 2
line 3
line 3
:

too stringent and aren't used in practice

Race Condition

correctness of the program depends on how the threads are scheduled



private versus shared variables

private: each processor has its own copy and each copy is only accessible to a specific processor

shared: same memory location for all threads

critical section: section of the code that can only be executed one thread at a time

Linear Algebra

BLAS (Basic Linear Algebra Subroutines)

BLAS - 1 vector , vector - vector

BLAS - 2 matrix - vector

BLAS - 3 matrix - matrix

Scheduling

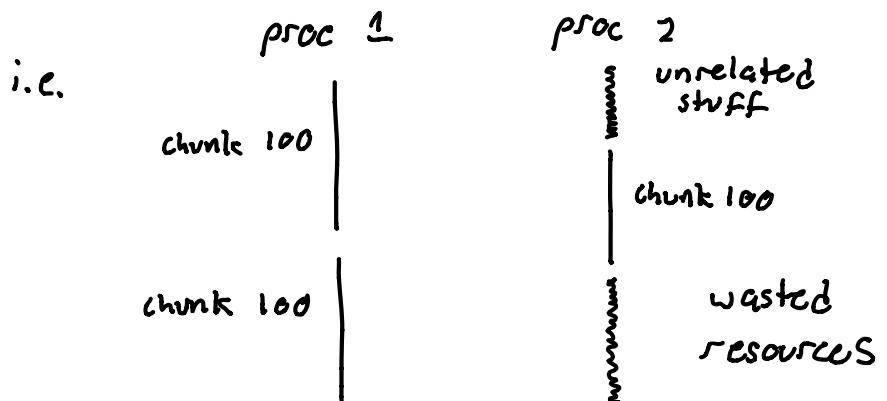
schedule (dynamic
static , chunk_size)

static : all threads scheduled during compile time
(before run-time)

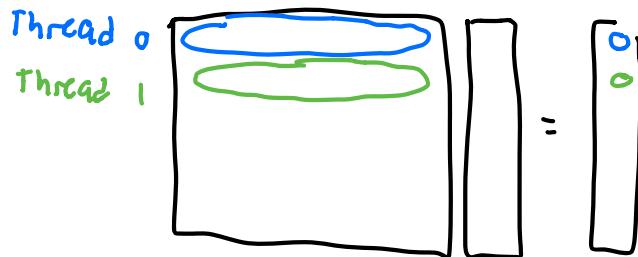
dynamic : some threads scheduled during run time

chunk_size : number of threads in a row that will all
be assigned to the same processor

- * dynamic scheduling better because you can't typically know when each processor will finish its work nor how long each chunk will take
- * don't want small chunk_size so we can take advantage of cache but don't want it too big because then I'll have processors that are free for a long time waiting for the last chunk to finish



matrix -vector operations



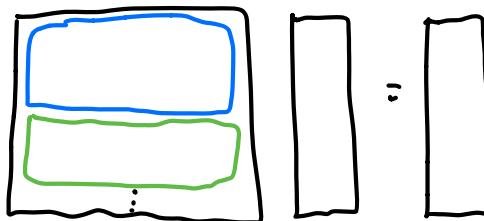
optimal approach for C
because of row-major order,
doesn't need to lock

* Divide the work given to different threads given the output

Also in C, data is stored in Row-major order

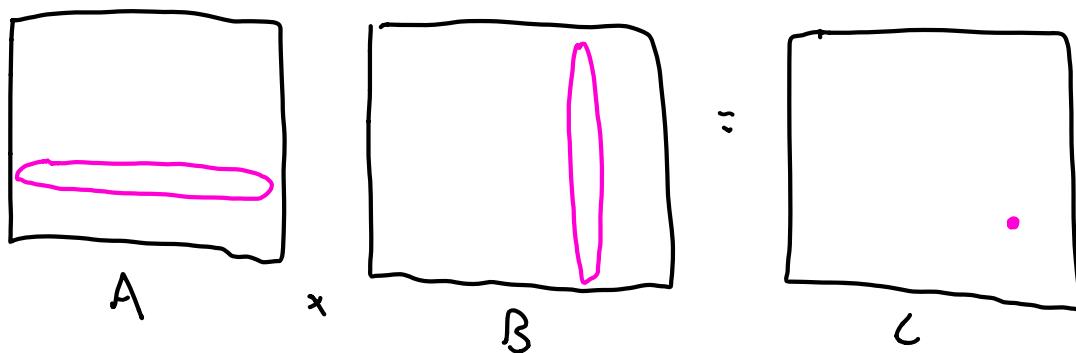
Best to store results in local variables and update it at the end because of "false sharing" where the caches for the two threads are so close, the caches both store data for location so computer thinks it has to reload data

Better solution:



assign multiple contiguous rows to same thread

matrix-matrix multiply



all $\in \mathbb{R}^{n \times n}$

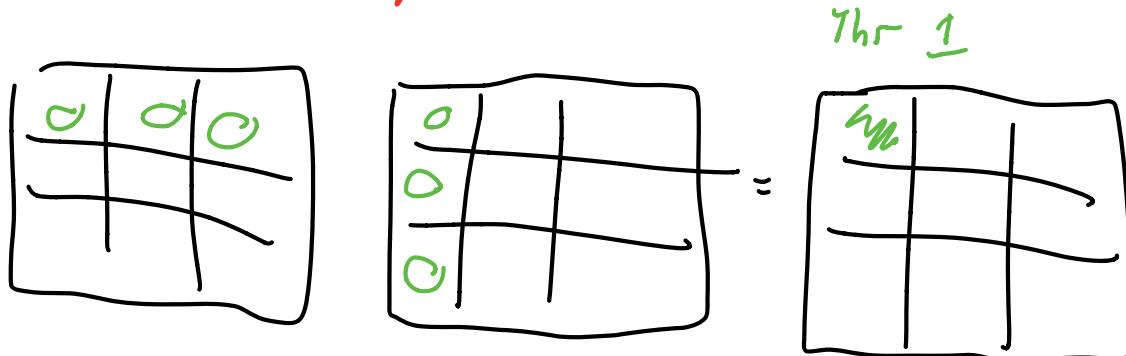
```
#pragma omp parallel for collapse (2)
```

```
for ( ;  
      for (j  
            for (k
```

```
c[i,j] += A[i,k]*B[k,j];
```

will choose i,j pair for each thread

But to take full advantage of cache



Make into submatrices that are $b \times b$
s.t. 3 $b \times b$ matrices can fit into the
cache at the same time

OpenMP tasks

Piece of code that doesn't have to be executed right now,
just whenever you have time

Tasks were introduced in OpenMP 3.0
but VS uses OpenMP 2.0

```
#pragma omp parallel
```

```
{
```

```
    int i;  
    i = omp_get_thread_num();
```

pragma omp task

{

int k = omp_get_thread_num();
printf ("%d, %d", i, k);

initial thread
that executes first
few lines of parallel
region and also puts
the task into task queue



thread that ends
up executing this
task

LU factorization

$$A = LU$$

$$\begin{bmatrix} \text{wavy lines} \end{bmatrix} = \begin{bmatrix} \text{triangle} \end{bmatrix} \cdot \begin{bmatrix} \text{triangle} \end{bmatrix}$$

$$Ax = b$$

$$LUx = b$$

$$Lx' = b$$
$$Ux = x'$$
$$\begin{bmatrix} \text{triangle} \end{bmatrix} \begin{bmatrix} x' \end{bmatrix} = \begin{bmatrix} b \end{bmatrix}$$

Triangular Solve Vector (TRSV) }

BLAS

$$\begin{bmatrix} \text{triangle} \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} x' \end{bmatrix}$$

TRSV

But how do we get L, U? Let's do example

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 1 & 4 \\ 5 & 3 & 1 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad b = \begin{bmatrix} 14 \\ 17 \\ 14 \end{bmatrix}$$

$$A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 2 & 3 \\ 3 & 1 & 4 \\ 5 & 3 & 1 \end{bmatrix}$$

$R_2 \rightarrow R_2 - 3R_1$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 0 & -5 & -5 \\ 5 & 3 & 1 \end{bmatrix} \quad R_3 \rightarrow R_3 - 5 \cdot R_1$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 0 & -5 & -5 \\ 0 & -7 & -14 \end{bmatrix} \quad R_3 \rightarrow R_3 - \frac{7}{5} R_2$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 0 & \frac{7}{5} & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 0 & -5 & -5 \\ 0 & 0 & -7 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 0 & \frac{7}{5} & 1 \end{bmatrix} \begin{bmatrix} x_1' \\ x_2' \\ x_3' \end{bmatrix} = \begin{bmatrix} 14 \\ 17 \\ 14 \end{bmatrix} \quad 0 - 3S + x_3' = 14$$

$$3S + x_3' = 14$$

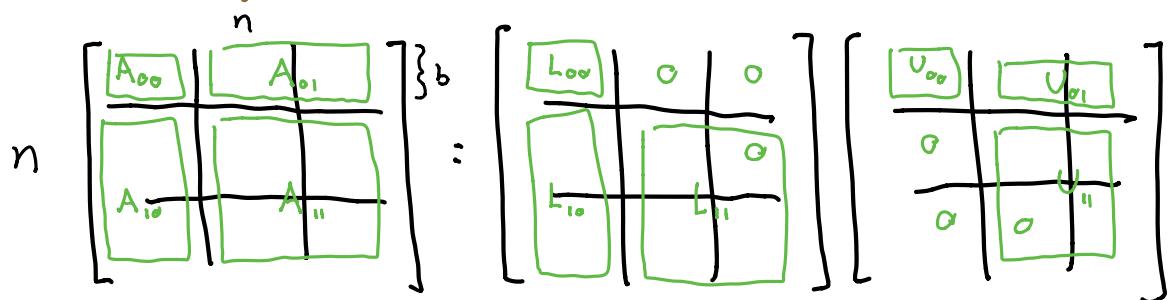
$$x_3' = -21$$

$$\begin{bmatrix} x_1' \\ x_2' \\ x_3' \end{bmatrix} = \begin{bmatrix} 14 \\ -2S \\ -21 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 0 & -5 & -5 \\ 0 & 0 & -7 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 14 \\ -2S \\ -21 \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

How do we implement this in parallel? Blocking!



$$A_{00} = L_{00} U_{00}$$

$$A_{01} = L_{00} U_{01}$$

$$A_{10} = L_{01} U_{00}$$

$$A_{11} = L_{10} U_{01} + L_{11} U_{11}$$

1) compute L_{00}, U_{00} by factorizing $L_{00} U_{00}$ (easy shown above)

2) compute U_{01} using $A_{01} = L_{00} U_{01}$

$$\begin{bmatrix} \text{wavy} \\ \text{wavy} \end{bmatrix} = \begin{bmatrix} \Delta \\ \Delta \end{bmatrix} \begin{bmatrix} \text{wavy} \\ \text{wavy} \end{bmatrix}$$

if these were square, this is simply triangular solve with a matrix TRSM

But they look more like this

$$b \begin{bmatrix} n-b & b \\ \square & \square & \square \end{bmatrix} = b \begin{bmatrix} \Delta \\ \Delta \end{bmatrix} \cdot b \begin{bmatrix} n-b \\ \square & \square & \square \end{bmatrix}$$

But we can break into $b \times b$ blocks then use TRSM methods

3) compute L_{10} using $A_{10} = L_{10} U_{00}$ (same as step 2)

4) $L_{11} U_{11} = A_{11} - L_{10} U_{01} = A_{11}'$

in this step I just need to solve A_{11}'

5) Recursively solve $A_{ii}' = L_{ii} U_{ii}$ until A_{ii}' is of size $b \times b$

how to parallelize: - steps 2, 3 can happen at the same time and even within those individual steps, the $b \times b$ blocks can be computed using different threads

- step 4 is most computationally heavy but we know how to parallelize

`#pragma omp atomic`

- acts like critical section but only works for simple operations (increment, addition, ...)
- essentially puts a lock on the memory address of whatever is being updated