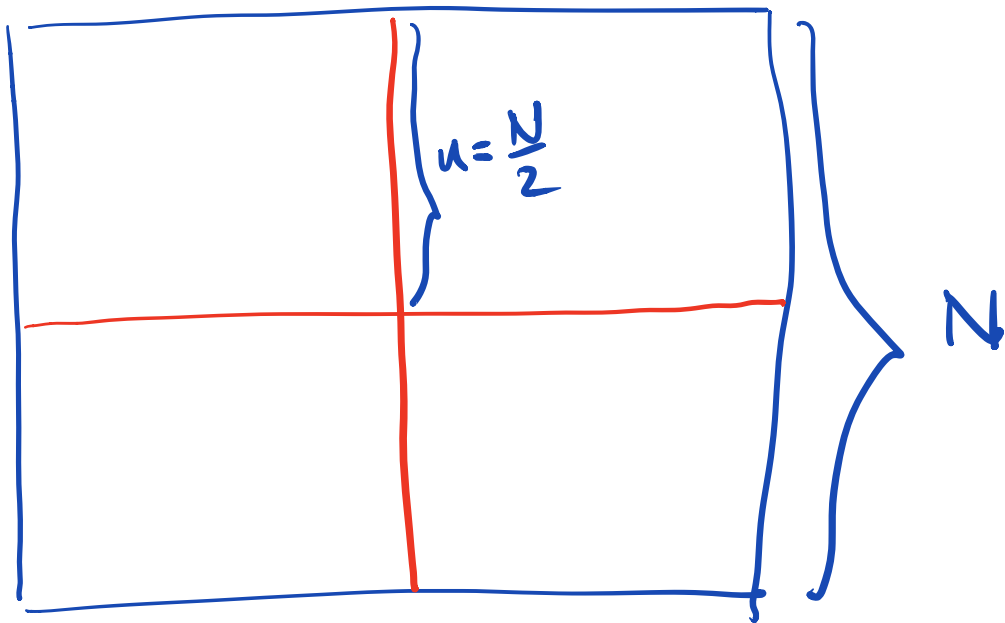


Last time:

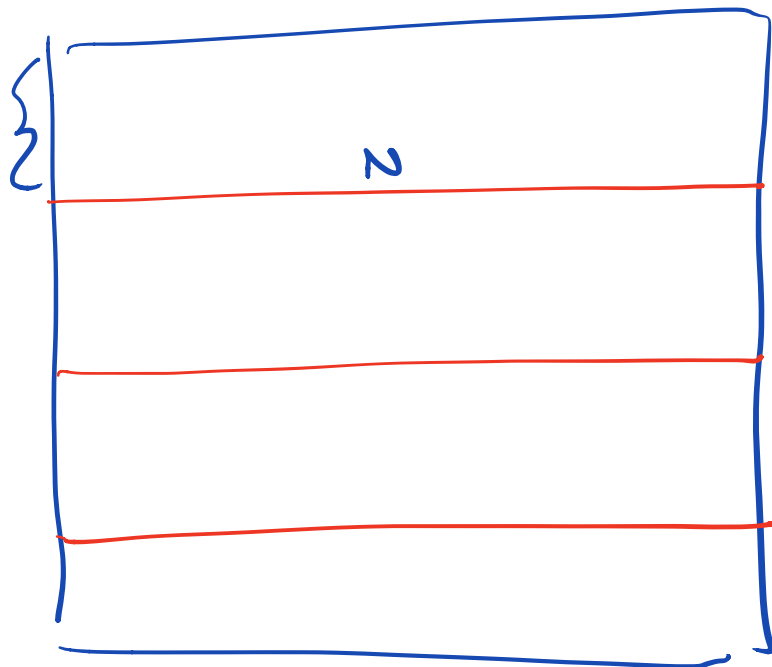
Motivated why we want for finite differences, a decomposition that attempts to minimize surface to volume ratio:

⇒ ideally squares/cubes



⇒ each box has  $\frac{N^2}{4}$  parts and  $2N$  side length.

$\frac{N}{4}$



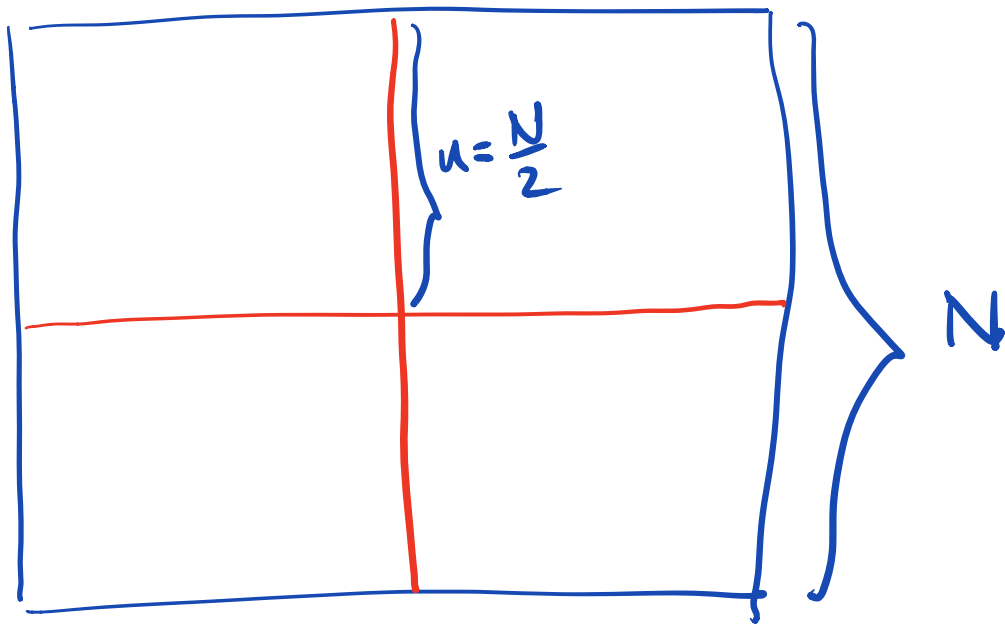
$\frac{N^2}{4}$   
 $\frac{N}{2}$

parts

side length.

Computational time  $\sim N^2$  # local parts

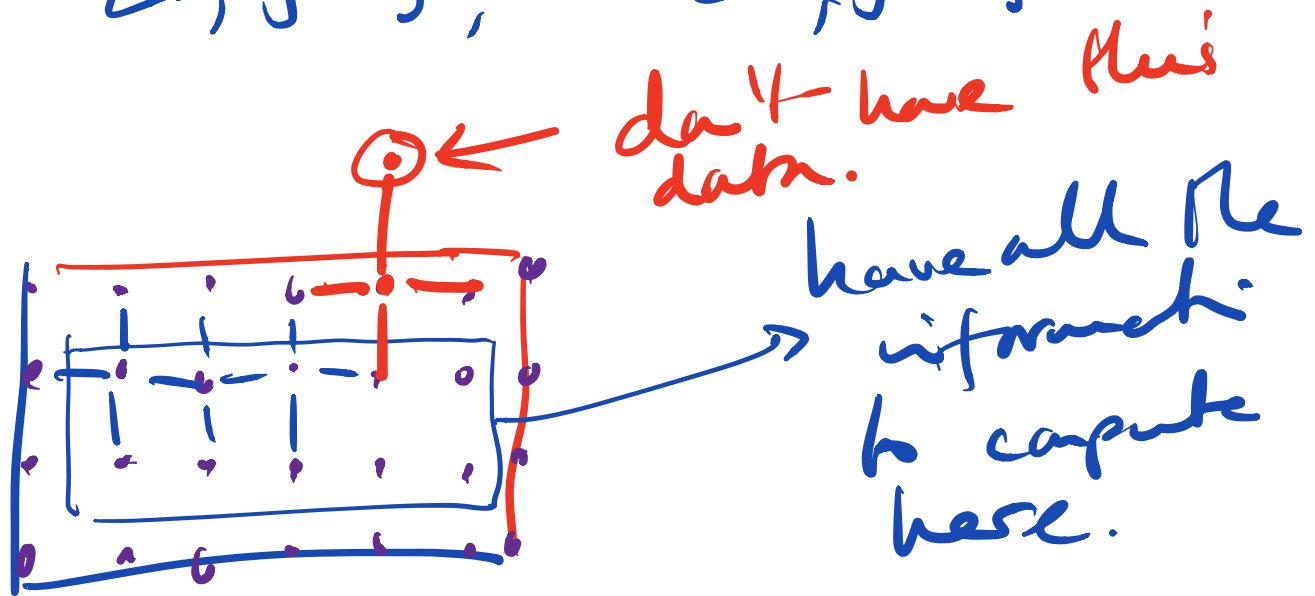
Communication volume  $\sim N$  side length





Accesses

$w[i, j]$ ,  $w[i-1, j]$ ,  $w[i+1, j]$   
 $w[i, j+1]$ ,  $w[i, j-1]$ .



Context:

Shared memory model (OpenMP)

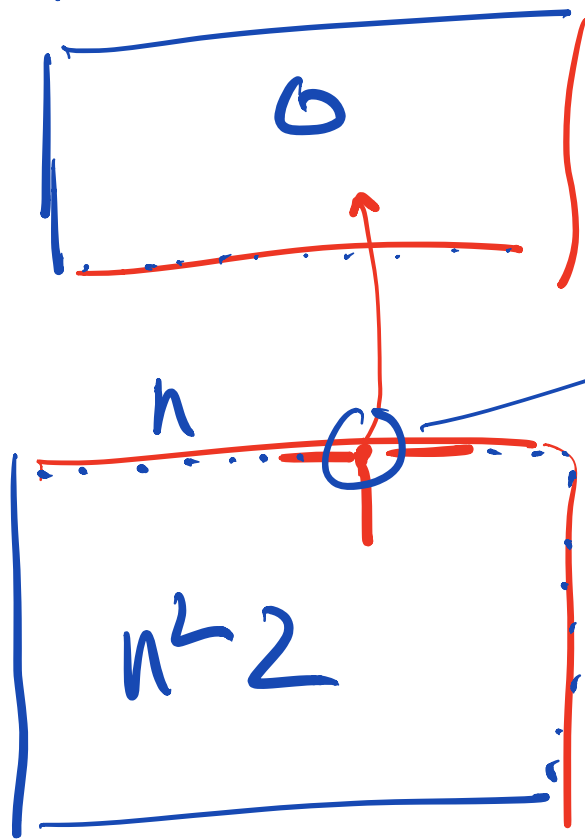
→ Not a problem.

→ just read this bit of the array.

→ still need to synchronize

in case someone else is updating that part.

MPI :



update reads  
from process 0.

message exchange.

Aside: there are some  
interfaces that offer  
distributed one-sided  
comm  $\Leftrightarrow$  shared memory  
 $\rightarrow$  synchronization  
we'll basically ignore these.

Message exchange synchronic  
"pairwise".

→ messaging is two sided.

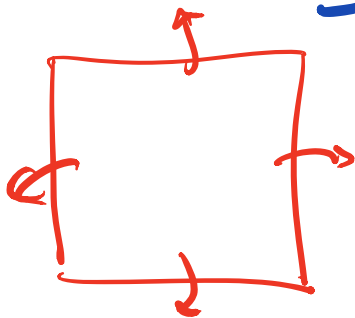
Process 0 says "I am sending  
this value" MPI\_Send

Process 2 says "I am  
receiving" MPI\_Recv

"Thanks, I've got it".

---

How much data?  
many messages?



$4n$  messages.  
each of size  
1 double.

Individual messages:

$$4n (t_{\alpha} + \beta)$$

message latency inverse bandwidth

All at once.

$$t_{\alpha} + \beta 4n$$

$$4nt_{\alpha} > t_{\alpha}$$

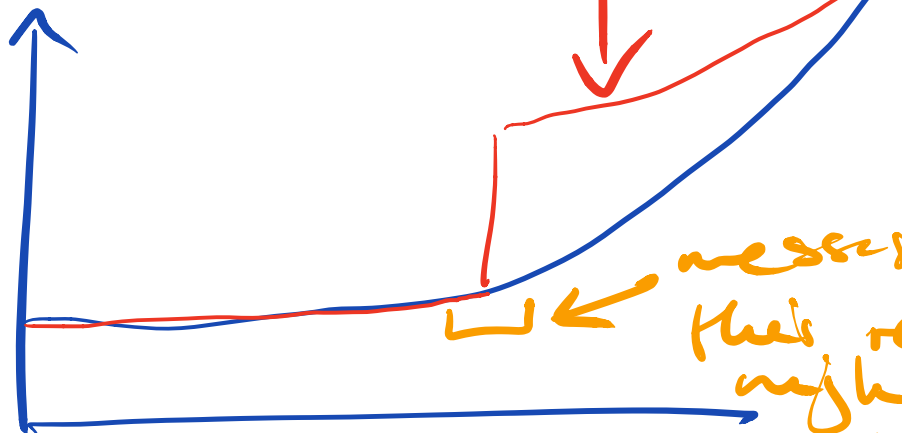
Calculate "batch" messages

up.



If model is

time



Bytes

messages - this region might benefit from splitting.

Conclude:

Communicate to exchange  
all messages

Compute.

Communicate again.

Grid data structure.

Solver stage: "trickery".

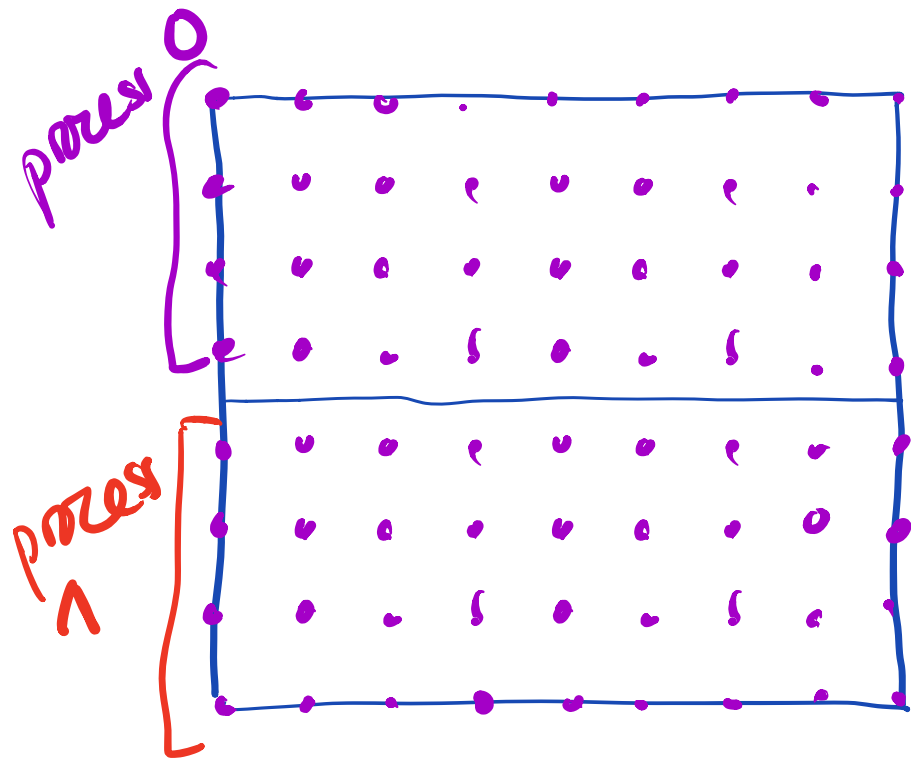
"Assembly" stage:  $\underbrace{Ax}$   
apply FD  
stencil.

Time step update.

$$u_{n+1} = u_n + \Delta t \underbrace{A u_n}_{\substack{\text{if we have} \\ A u_n \text{ as a vector}}}$$



This is Patrick.



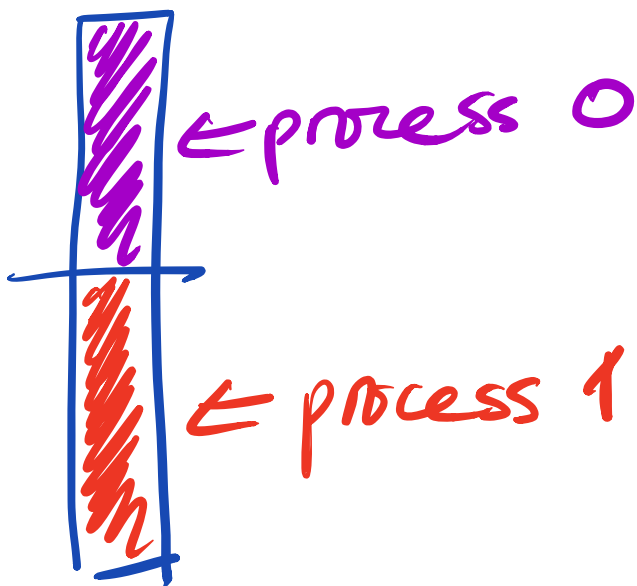
two-step update  
doesn't  
communicate  
over  $A_{nn}$   
is calculated.

$$u_{n+1} = \underset{\substack{\uparrow \\ \text{Patrick} \\ \text{stencil}}}{\cdot} u_n$$

Update can work on

"global" vector.

$\Rightarrow$  all dots are unique.



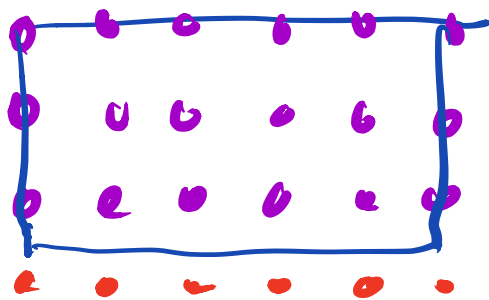
"non-overlapping"

$$\tau_n = -\frac{1}{A} u_n$$

stencil deriv.

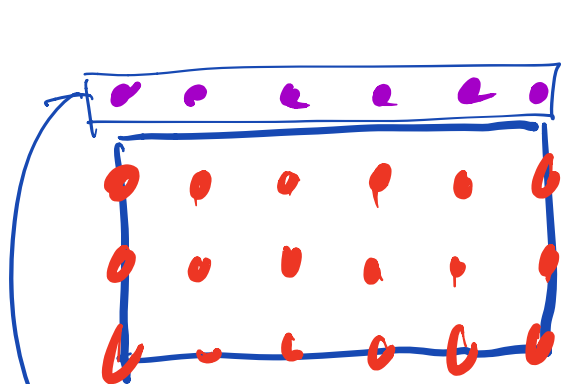
Needs "local" vector.

⇒ includes shared dots  
 ↳ compute stencil updates.



process 0

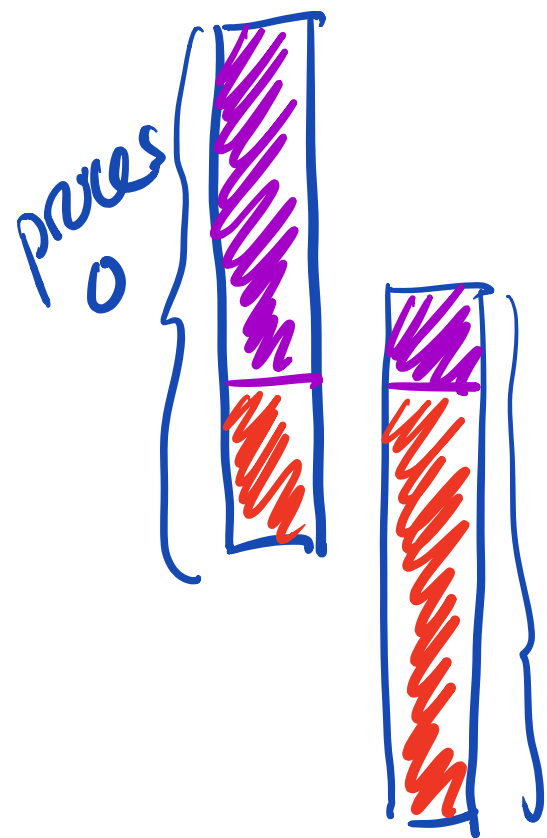
← values from proc 1



← values from proc 0

process 1

~ "ghost pairs"  
 ~ "halo regions"



overlapping  
"local" vectors.

process 1.

enough information to  
compute  $Ax$  without  
further communication.

Goal: decompose dots  
into "owned" (global)  
partitions.

- determine global parts  
so that overlapped (local)

vectors can be produced.

def Au (u<sub>g</sub>):

insert ghost  
pts.

compute | u<sub>L</sub> = grid.global\_to\_local(u<sub>g</sub>)

compute | loop over local grid pts:

$$r_L[i, j] = A(\text{stencil}, u_L)$$

← produces correct entries for all owned points.

→ that's all we need for

tristep.

stencils → not for FD. <sup>used for</sup> FEM.

$$r_g = \text{grid.local_to_global}(r_L)$$

For our purposes

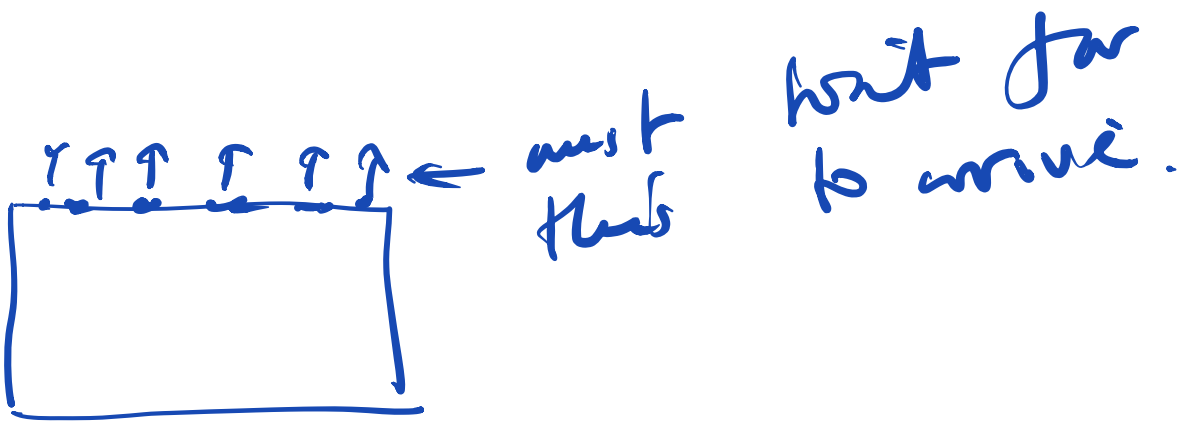
$$r_g = r_L$$

return r<sub>g</sub>.

This gives an algorithm that  
takes  $T_C$  time to communicate  
&  $T_W$  time to compute

$$\Rightarrow T_{\text{total}} = T_C + T_W.$$

Why? work must wait  
for communication to be  
finished before it starts.

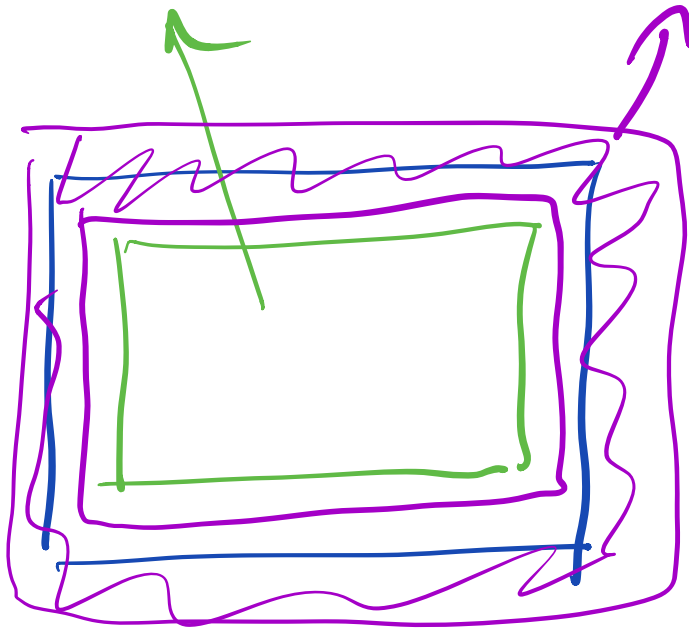


---

Can arrange for  
 $T_{\text{total}} = \max(T_C, T_W)$

Has:

Split local domain into  
"interior" & "boundary"



interior doesn't need result of  
communication.

$u_i = \text{gnd.g2l\_begin}(u_g)$  ← visits local points (tail)

$r_i = \text{compute\_interior}(u_i)$

$u_i = \text{gnd.g2l\_end}(u_g)$  ← finishes communication

$r_{i+} = \text{compute\_boundary}(u_i)$

→ Asynchronous communication  
in background.

⇒ MPI offers this:

getting it to work in practice  
is black magic

→ "asynchronous progress"  
needs a background thread  
inside the library  
that is scheduled onto a core  
by the OS.

Hamilton doesn't offer this.

Some supercomputing platforms  
add a separate chip for  
this purpose.

---

• Design:

Separate comms &  
compute.

• Allows work on global  
vectors (pointers)

↔ local vectors (overlapped).

Next time:

• See how to realize this with MPI.

• Scaling analysis after Fischer (2015) for Jacobi iteration.  
→ scaling limits for PDE-based simulation.

• Communication latency  
→ consequences for multigrid.

→ Read introduction & upto end of section 11. B. 1

~ "Jacobi iteration". (3½ pages total).