

Optimizing Collective Communication on Multicores

Rajesh Nishtala¹ Katherine Yelick¹

¹University of California, Berkeley

(2009)

Algorithms for Scalable Synchronization on Shared-Memory Multiprocessors

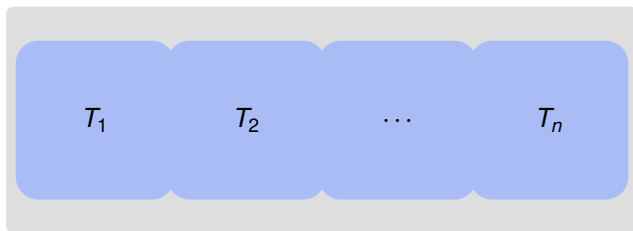
John M. Mellor-Crummey, Michael L.Scott
(1991)

PGAS Languages

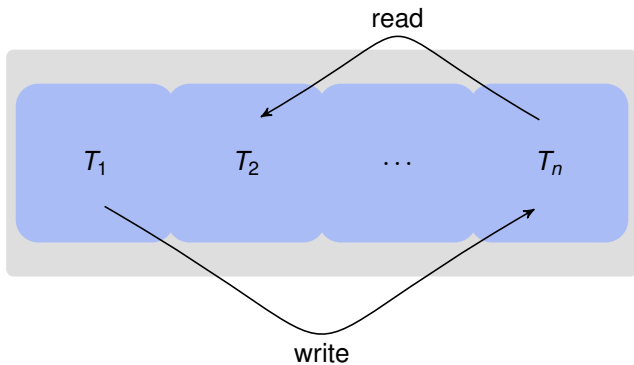
- ▶ Focus on **P**artitioned **G**lobal **A**ddress **S**pace languages

Partitioned Addressspace

one address space



One Sided Communication

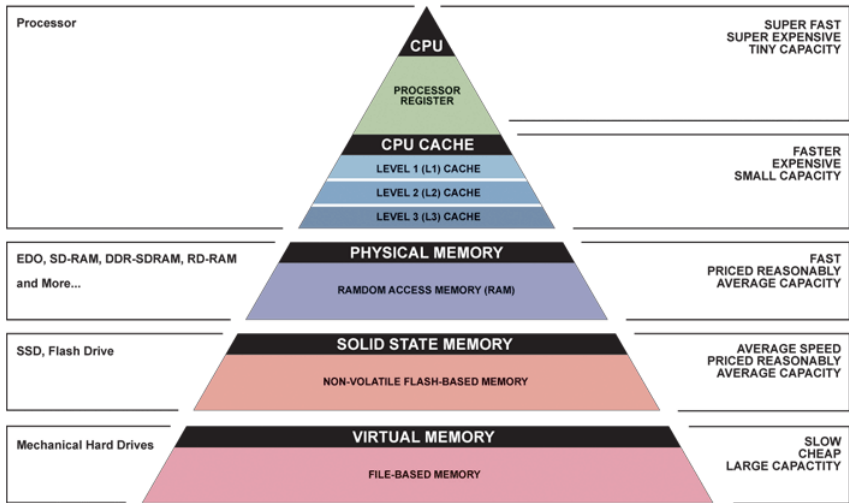


PGAS Languages

- ▶ UPC, Unified Parallel C
- ▶ CAF, Co-array Fortran
- ▶ Titanium, a Java dialect

Context

- ▶ The gap between processors and memory systems is still enormous

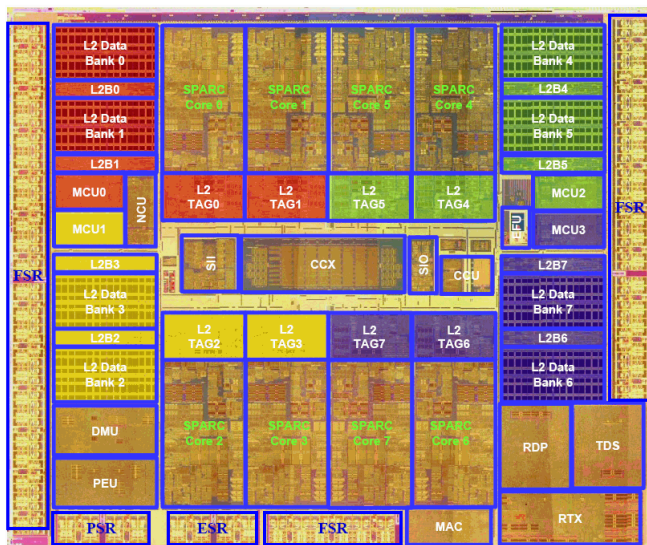


- ▶ Today: processors don't get faster, but we see more and more processors on a single chip

Processor	GHz	Cores (Threads)	Sockets
Intel Clovertown	2.66	8 (8)	2
AMD Barcelona	2.3	32 (32)	8
Sun Niagara 2	1.4	32 (256)	4

Table: Experimental Platforms

Sun Niagara 2

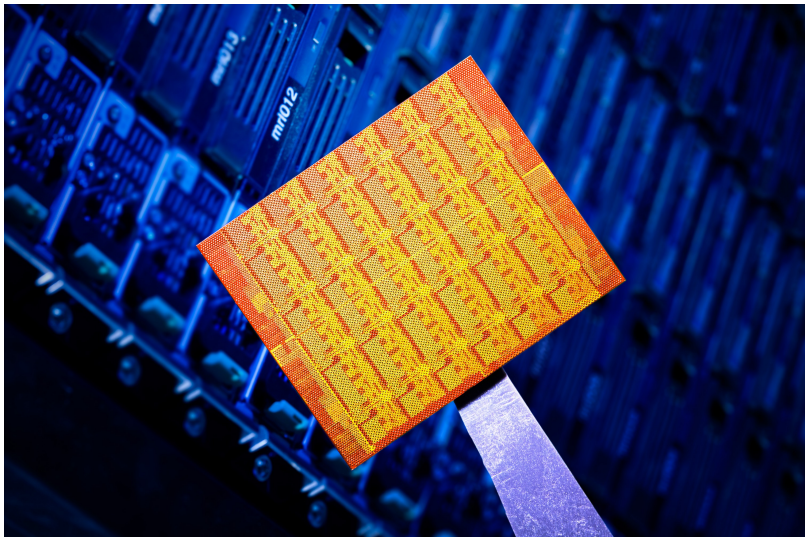


<http://www.rz.rwth-aachen.de/aw/cms/rz/Themen/hochleistungsrechnen/>

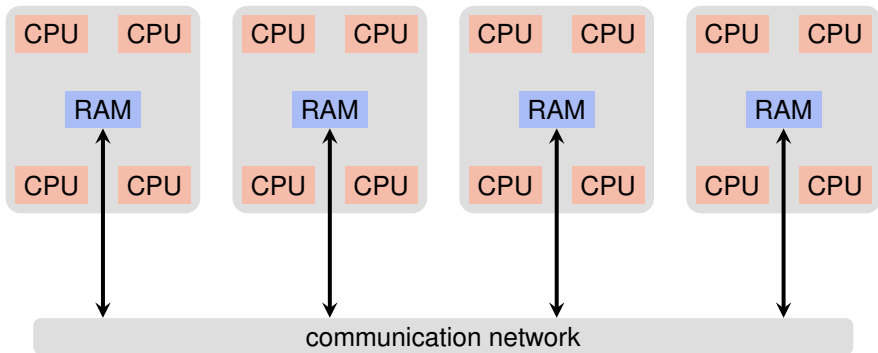
[rechnersysteme/beschreibung_der_hpc_systeme/ultrasparc_t2/rba/ultrasparc_t2_architectural_details/?lang=de](http://www.rz.rwth-aachen.de/aw/cms/rz/Themen/hochleistungsrechnen/rechnersysteme/beschreibung_der_hpc_systeme/ultrasparc_t2/rba/ultrasparc_t2_architectural_details/?lang=de)

- ▶ The number of processors on a chip grows at an exponential pace

Intel Single-Chip Cloud Computer (48 Cores)



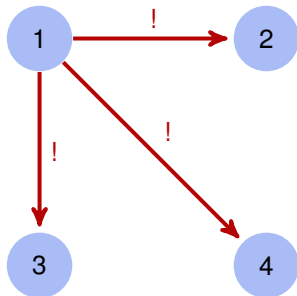
- ▶ Communication in its most general form is the movement of data within cores, between cores or within memory systems



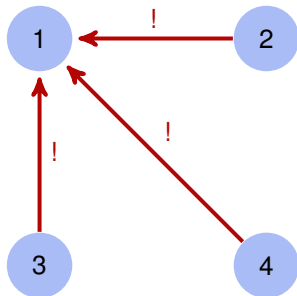
Collective Communication

- ▶ Communication-intensive problems often involve global communication

Broadcast



Gather



- ▶ These operations are thought of as **collective communication operations**

Example: Sum of Vector Elements

1 2 3 4 5 6 7 8 9 10

Example: Sum of Vector Elements

- ▶ Create workers

1 2 3 4 5 6 7 8 9 10

W_1

W_2

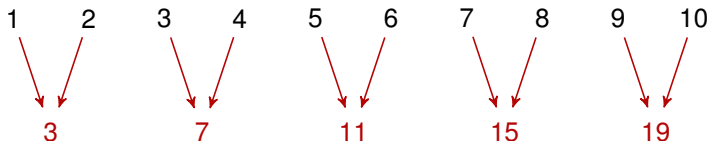
W_3

W_4

W_5

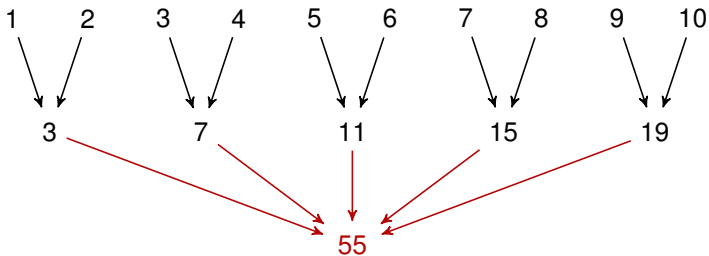
Example: Sum of Vector Elements

- ▶ Every worker **sums up** it's part of the vector



Example: Sum of Vector Elements

- ▶ The main thread gathers the partial results and **sums them up**



Example: Sum of Vector Elements

Pseudocode (main thread):

```
double [] vector = read_vector();  
Thread [] workers = spawn_workers();  
  
start_workers(workers);  
double result = calculate_result(workers);
```

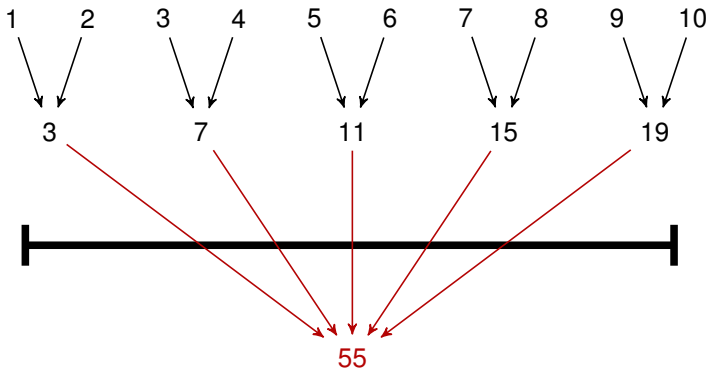

Example: Sum of Vector Elements

Pseudocode (main thread):

```
double [] vector = read_vector();  
Thread [] workers = spwan_workers();  
  
start_workers(workers);  
wait_until_everything_finished(workers);  
double result = calculate_result(workers);
```

Barrier

- ▶ Synchronization method for a group of threads
- ▶ A thread can only continue its execution after every thread has called the barrier



Collective Communication Operation

“... group of threads works together to perform a global communication operation ...”

Reduce

- ▶ Divide a problem into smaller subproblems
- ▶ Every thread contributes it's part to the solution
- ▶ Example: Calculate the smallest entry of a vector

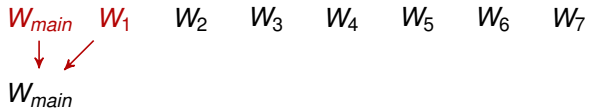
Flat vs. Tree

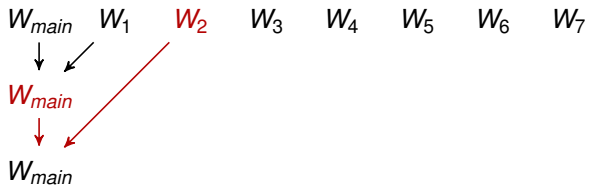
- ▶ For communication among threads, different topologies can be used

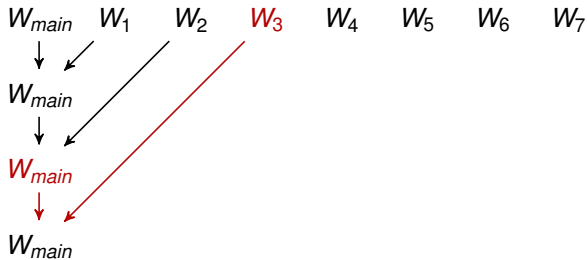
Flat Topology

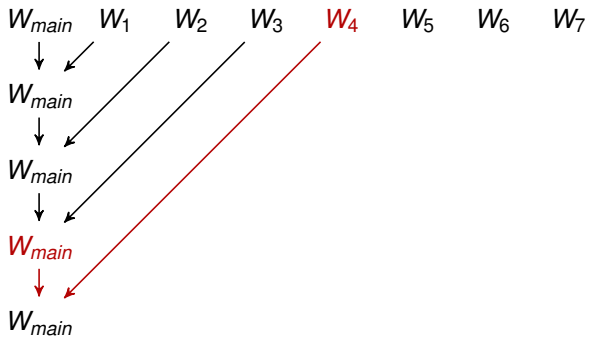
- ▶ Example: we have a reduce operation
- ▶ in the end the **main thread** W_{main} has to wait for every worker thread W_1, \dots, W_7

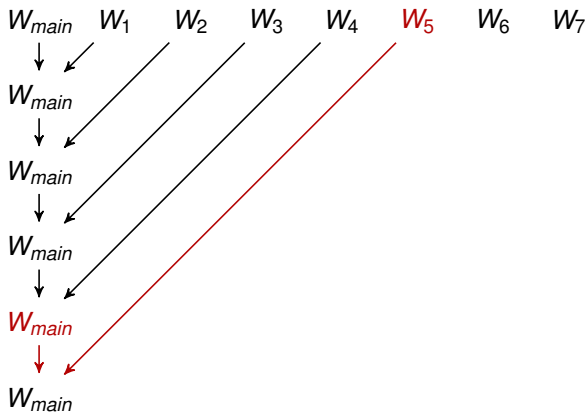
W_{main} W_1 W_2 W_3 W_4 W_5 W_6 W_7

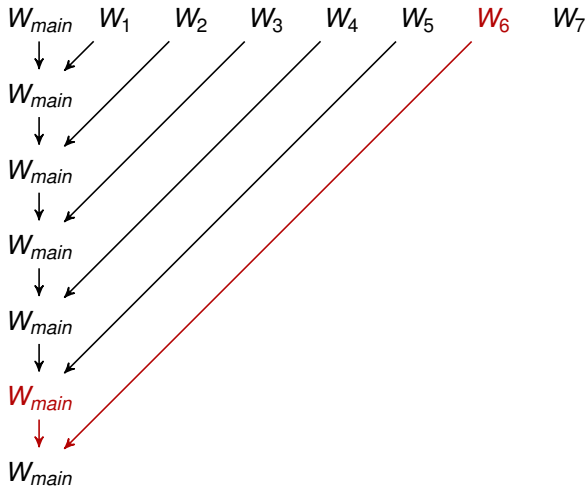


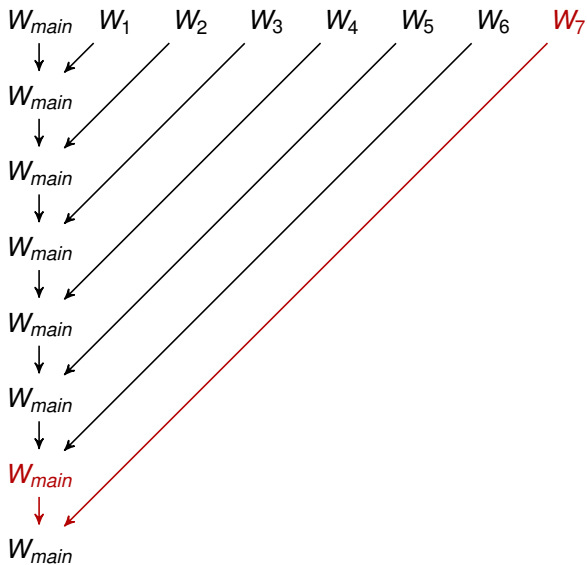








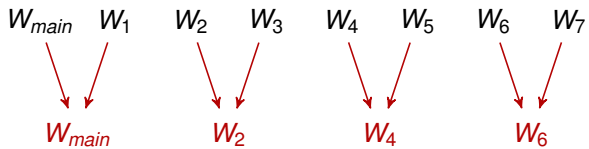


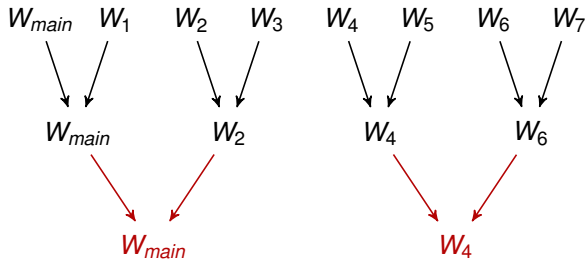


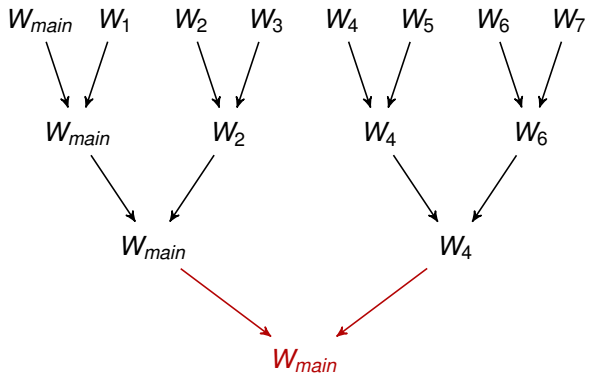
Tree Topology

- ▶ Example: we have a reduce operation
- ▶ in the end the **main thread** W_{main} has to wait for every worker thread W_1, \dots, W_7

W_{main} W_1 W_2 W_3 W_4 W_5 W_6 W_7







Analysis

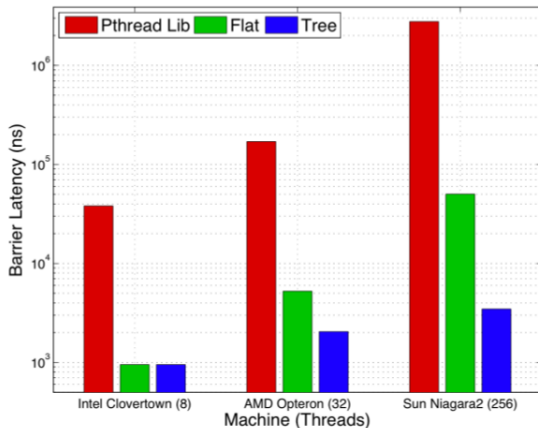


Figure: Barrier Performance

Barrier Implementation

```
#define N 4

pthread_t threads[N];
volatile int ready[N];
volatile int go[N];

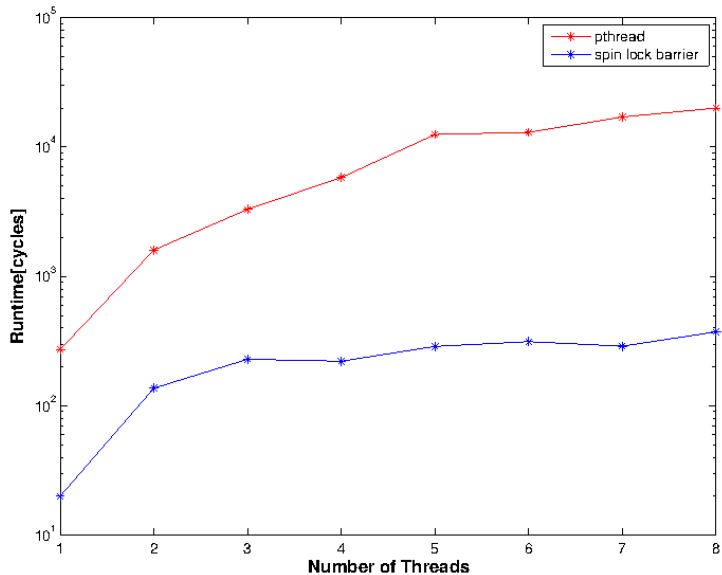
void barrier(int id) {

    if (id == 0) {
        //wait for each thread
        for (int i = 1; i < N; i++)
            while (ready[i] == 0);

        //reset the ready flags
        for (int i = 0; i < N; i++)
            ready[i] = 0;

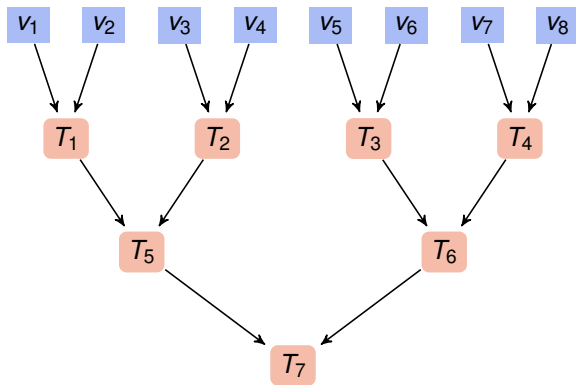
        //signal each thread
        for (int i = 0; i < N; i++)
            go[i] = 1;
    }
    else {
        ready[id] = 1;
        //wait until thread is signalled
        while (go[id] == 0);
        go[id] = 0;
    }
}
```

Experiment: Barrier Implementation



- ▶ **Strict synchronization**: Data movement can only start after all threads have entered the collective and must be completed before the first thread exits the collective

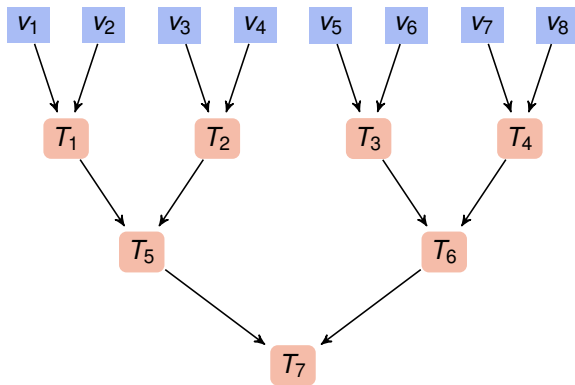
Strict Synchronization

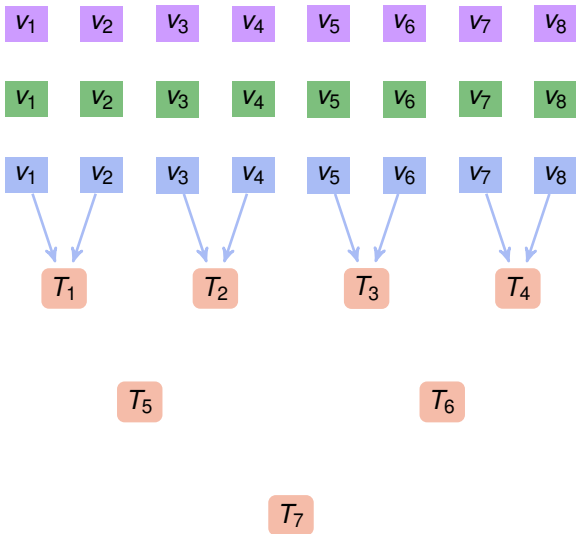


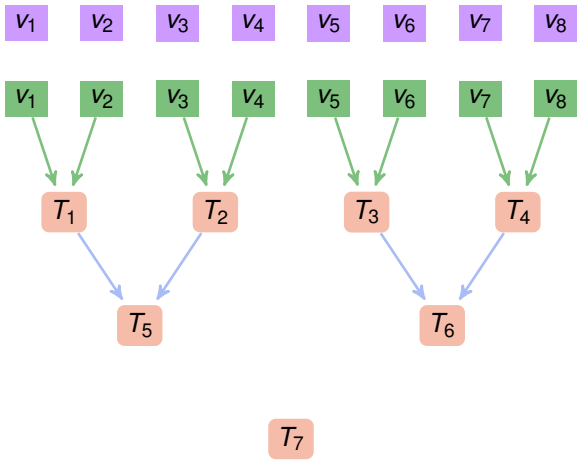
Loosening Synchronization Requirements

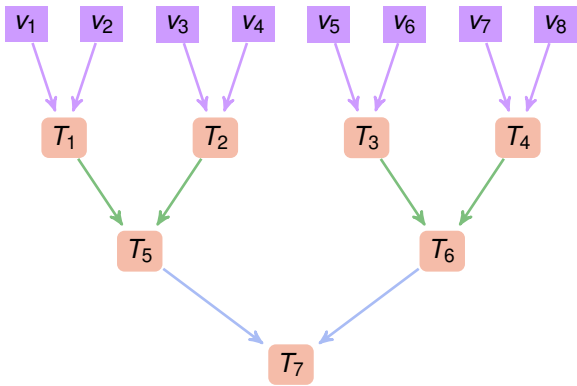
- ▶ **Loose synchronization:** Data movement can begin as soon as any thread has entered the collective and continue until the last thread leaves the collective

Loose Synchronization









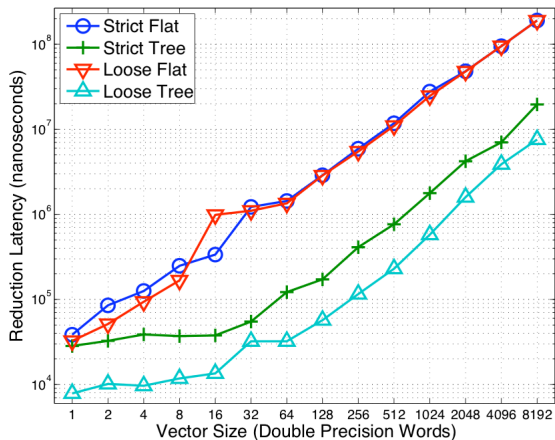


Figure 5: Optimal Algorithm Selection on Niagara2
(32 cores, 256 threads)

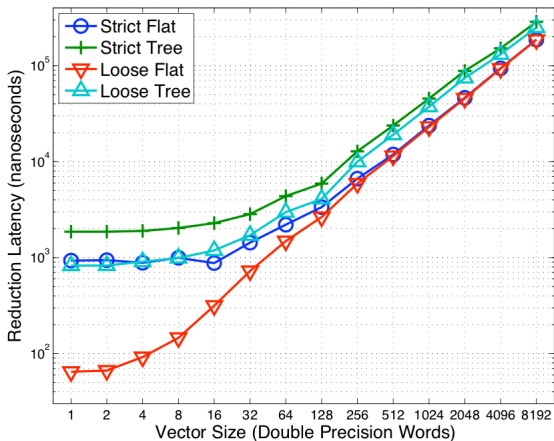


Figure 6: Optimal Algorithm Selection on Clovertown
(8 cores, 8 threads)

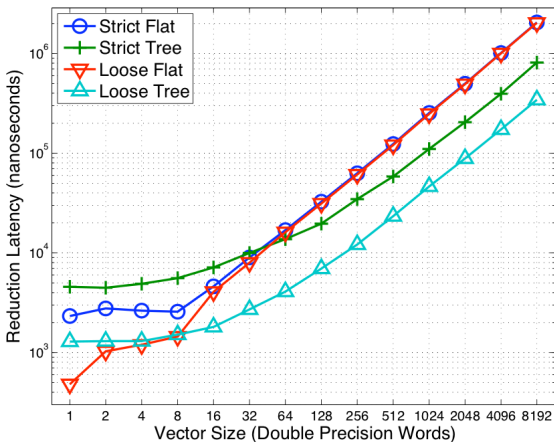


Figure 4: Optimal Algorithm Selection on Barcelona
(32 cores, 32 threads)

Summary

- ▶ Best strategy depends on the hardware and on the problem
- ▶ Using a library that can automatically adapt to a given situation can bring a great performance improvement, since hand tuning takes far too long

Words on the Paper

- ▶ Very high level
- ▶ Description of the problem without concrete solution
- ▶ No implementation
- ▶ Plots aren't always clear and precise