

Parallel Programming

Performance Metrics

Performance Speedup: speedup(p_n) = $\frac{\text{performance}(p_n)}{\text{performance}(p_1)}$
 \uparrow , processor

In HPC (high perf. computing) performance = $\frac{\text{work}}{\text{time}}$

$$\text{speeddeep}(pn) = \frac{\text{time}(p1)}{\text{time}(pn)}$$

Speedup based on Throughput:

$$\text{speedup}_t(p_n) = \frac{\text{tpm}(p_n)}{\text{tpm}(p_1)} \quad \begin{matrix} \text{transactions per} \\ \text{minute} \end{matrix}$$

Parallel efficiency: $\text{efficiency}(p_n) = \frac{\text{speedup}(p_n)}{n}$

Amdahl's law

Execution time (assuming parallel regions have full speed):

$$T(p) = (1-f) \cdot T + \frac{f \cdot T}{P}$$

$$25\% = 0,25$$

$$(1-0,25) \cdot \overline{I} + \frac{f \cdot T}{\overline{I}}$$

График
бен. seq.

Оч. - начало,
 максим.

f - fraction of parallel execut.
p - n of parallel tasks/threads/
processes
T - sequential time

We assume that state (resources) stays the same

Maximal speedup

$$SU(p) = \frac{T}{T(p)} \Rightarrow \frac{T}{(1-f) \cdot T + \frac{f \cdot T}{P}} = \frac{1}{1-f + \frac{f}{P}}$$

Superlinear speedup - more than linear speedup, can happen when with num. of processes resources increase

Types of parallelism

- Data parallelism

Спредум дахнве ка векторе, блокове са огни и те ме определене ка веc векторах.

- Functional parallelism

Деини веc ка задачи, задачи могуt блокирато по време реагира

Flynn's Classification

Single Instruction

Single Data
SISD

Multiple Data
SIMD

Multiple Instruction

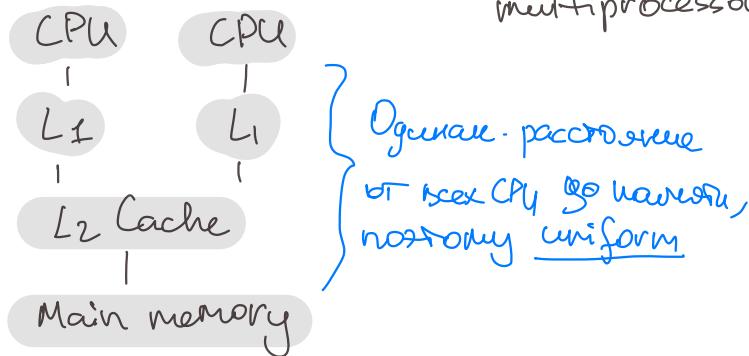
MISD

MIMD

- single instr.
mult. data

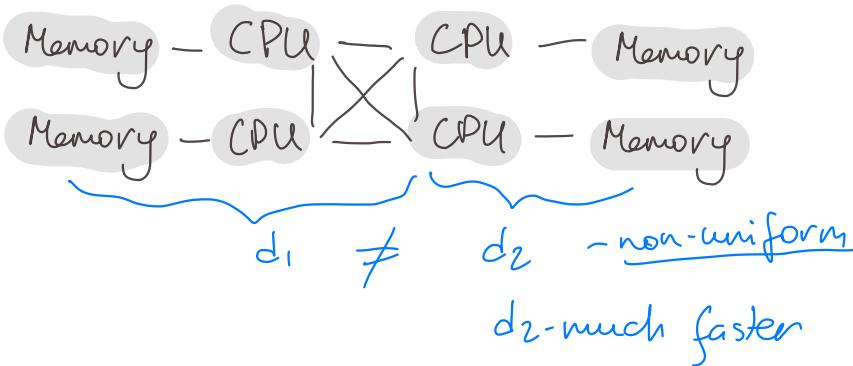
Shared Memory Models

- Uniform Access Memory (UAMA) (aka symmetric multiprocessors - SMP)



- Non-uniform Memory Access (NUMA)

(aka Distributed Shared Memory Systems - DSW-DSM)



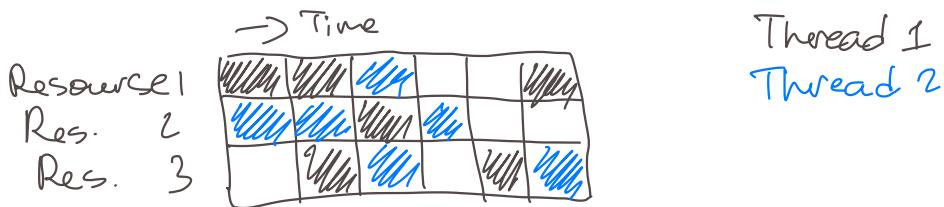
Cache coherent!

Thread is a stream of execution

Software threads are abstracting execution streams

user-level ↗ system-level

Hyper-Threading (SMT)



POSIX Threads

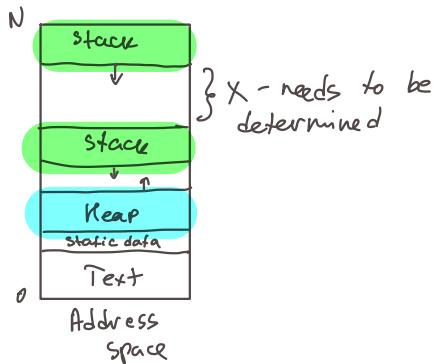
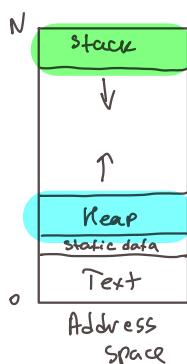
```
int pthread_create(*thread, *attr, *start_routine, *arg)
```

ret. thread id ↑
properties
of threads behaviour
(bound/unbound)
Scheduling policy...)

```
int pthread_join(pthread_t thread, void **retval)
```

id ↑ return value ↑

Hyunso (momoko) set&get stack size & stack addresses,
Stack Structure



Process - a running instance of a program with its own memory address space

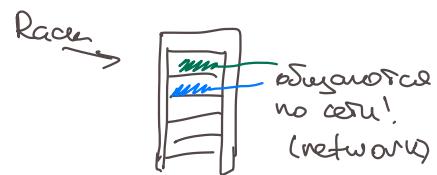
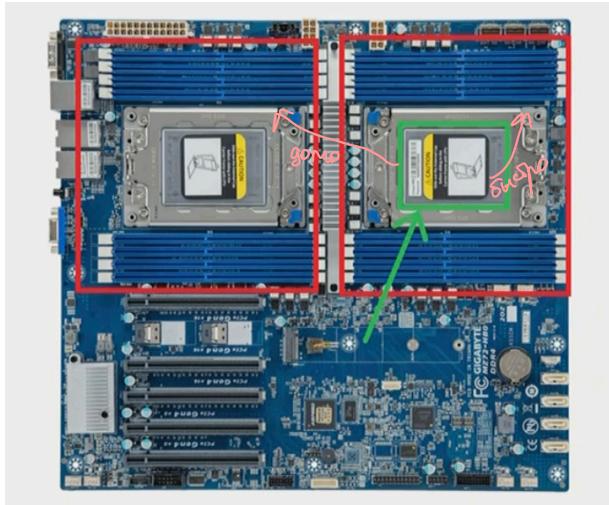
Core - an independent processing unit

Thread - independent stream of execution

Hardware thread - implementation of an execution stream in hardware

Software thread - programming abstraction that represents a stream of execution

Node (multiple per race)



NUMA Socket
(CPU + memory)

Physical CPU

Threads creation

```
1 ...
2 // Start numThreads threads
3 std::thread threads[numThreads];
4 for (int i = 0 ; i < numThreads ; ++i){
5     threads[i] = std::thread(oops_kernel, &i);
6 }
7 for (int i = 0 ; i < numThreads ; ++i){
8     threads[i].join();
9 }
10 ...
```

```
1 void oops_kernel (int *args) {
2     int argument = *args;
3     ...
4 }
```

From user's moment ~~basepointer~~
parameter, new thread ~~basepointer~~, u
is created out of scope

Myself generate race

```
1 ...
2 // Start numThreads threads
3 int ids[numThreads];
4 std::thread threads[numThreads];
5 for (int i = 0 ; i < numThreads ; ++i){
6     ids[i] = i;
7     threads[i] = std::thread(a_kernel, &ids[i]);
8 }
9 for (int i = 0 ; i < numThreads ; ++i){
10    threads[i].join();
11 }
12 ...
```

```
1 void a_kernel (int* args) {
2     int argument = *args;
3     ...
4 }
```

Lock Implementations

Kriterien:

1. Correctness : Guarantees mutual exclusion
2. Fairness: Every process eventually gets the mutex

Implementations:

• Spin Lock

Spin on a flag, if lock is taken
Typically via atomic operations

Advantage: fast response time

Disadvantage:

- blocks the hardware threads,
- uses resources, costs energy
- generally bad for concurrent execution (race condition, deadlocks)

(Race condition, deadlock => buggy
non-deterministic behavior)

• Yielding locks

If lock is taken - do nothing, give back hardware thread. Usually implemented via runtime system.

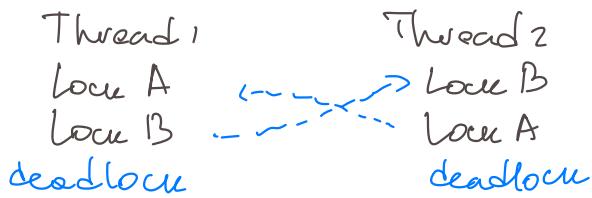
Advantage: low resource usage High performance computing

Disadvantage: slow response time
• generally bad for HPC

Lock Granularity

- + No race condition, a smaller lock (granularity)
- + Maximizes parallelism
- requires multiple locks management

Danger: Deadlocks



Deadlock avoidance Strategies

- only hold one lock at a time - bad!
- have a lock manager, ask permission to lock
- order the locks - charana A, norm sanga B

Threading Performance

Overheads:

- Thread creation & destruction can be expensive

That's why we have thread pools - create once in the beg.

locks:

→ lock
→ O - low overhead
→ ↑

→ lock
→ O - large overhead
→ ↑↑↑

more threads accessing
more often

Choose wisely between locking all data under
one lock VS lock granularity

Thread pinning:

Укога ходим "занурено" тез ка огаси ере,
множу overhead ка cache move узен

Open MP

MP - multiprocessing

export OMP_NUM_THREADS=2
gcc -O3 -fopenmp openmp.c

#include <omp.h>

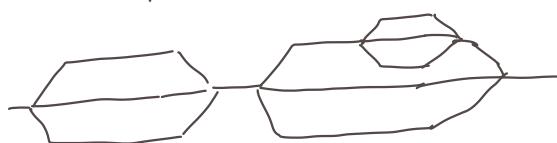
...
#pragma omp parallel
{
...
}

}

bce parsed
fork-join
noz kospom

+ ucnovbyz yet
pthreads

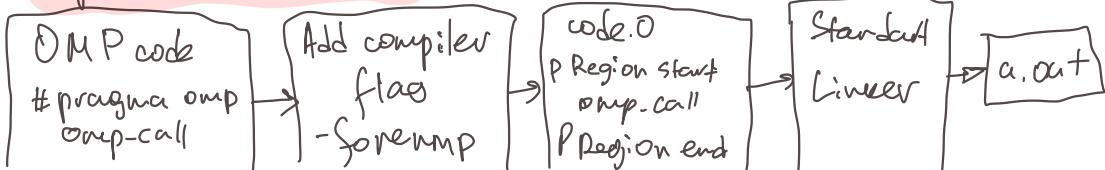
Nested parallelism:



No ne bce nengopen
ew nozgerumusand

na canom gene noz kospom ono cozaet tyi
max-n Tregos, res ucn. 3. Nozau ne yganeer, g
uz OC neperenobz yet.

Open MP Toolchain:



Hluge rečenice kcn. univerza, karere.c omp...
↗
namespace omp

```
main() {
    int a[1000], b[1000];

    #pragma omp parallel sections
    {

        #pragma omp section
        for (int i=0; i<1000; i++)
            a[i] = 100;
        #pragma omp section
        for (int i=0; i<1000; i++)
            b[i] = 200;
    }
}
```

} nesynchronizované
} napájení,
} konečný okraj
} stejné

```
main () {
    int a[100];
    #pragma omp parallel
    {
        #pragma omp for
        for (int i= 1; i<n;i++)
            a[i] = i;
    }
}
```

here guaranteed the sync

} should be very simple
+ independent
iteration,
i.e. order is not
guaranteed

(unless) #pragma omp for nowait

Manned nacipants poydzenie o użyciu:

#pragma omp for schedule(**static**)

reguparne
c O

- **static**

Fixed sized chunks, distributed in round-robin fashion
(default = $\frac{\text{num_iter}}{\text{num_thr}}$) lower overhead
no race or contention!

- **dynamic**

Fixed size chunks (default = 1), distr. one by one at runtime as chunks finished (ognie zakonczeni, nazyw. rezusive)
schedule(dynamic, 4)

- **guided**

Start with large chunk size, then exponent. decrease

- **runtime**

Controlled at runtime using control variable

- **auto**

Compiler/runtime can choose

Shared data

No gener. Bce reodawcze.

Wszystkie zmienne uwalniane:

int t;



Manned nacipants
shared(t)

uni
default(private)
shared/none

#pragma omp parallel private(t)

Но! Эта соргает только неправильные, не конечные!

Т.е.

int t=3;

... private(t)

{

t = t+1;

}

private(t) - соргает неправильное счищением
t при замкн. потока, но
правильное. мусором

содержит мусор! + соргает,
исчезаю, если
раз на один фраг

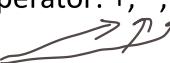
firstprivate(t) - соргает конец t до первого соргаса
или всех фрагов

lastprivate(t) - делает конец неправильные соргаса
t, но $\Rightarrow t=t+1$ заменяется
результатом последнего фрага

reduction(op, t) - инициализ. t первым,
в конце каждого фрага
берет значение t и
применяет оператор op.

o Reductions Операции:

- Aggregate results of region to master thread
- reduction(operator: list)
 - operator: +, *, -, &, ^, |, &&, ||

иерархия 0  иерархия 1

```

int a=0;
#pragma omp parallel for reduction(+: a)
for (int j=0; j<4; j++)
{
    a = j; ← Но у каждого Треда
}
printf("Final Value of a=%d\n", a);

```

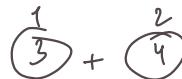
уничтожен предыдущий кусок,
потому что +
это конец
Синхронизация
одинаковое значение a
каждого Треда

OMP_NUM_THREADS=4

Final value of a:

Было неправильное 4, а
Треда 2, то есть сумма
была 2x увеличена!

T₁ T₂



Мы уже видели так:

```

int a=0;
#pragma omp parallel for reduction(+: a)
for (int j=0; j<100; j++)
{
    a += j;
}
printf("Final Value of a=%d\n", a);

```

CV кумулятивно
не может параллельно!
guard-lock {^{error?}_{unlock}}
модификатор local() {^{error?}_{doesn't unlock!}}

```

int data=0;
#pragma omp parallel for
for (int i=0; i<2; i++) {
    data++;
}

```

race condition

предыдущий
keep secret,
модификатор 1,
модификатор 2

- `int data = 0`

установить `data = 0` где нек
треугольник

```
#pragma omp parallel for firstprivate(data)
lastprivate(data)
for (int i=0; i<2; i++) {
    data++;
}
```

занести в `data`
новые значения
из каждого потока

`output(data)` → `data = 1`

Но это не единственный способ
редукции (+, data)

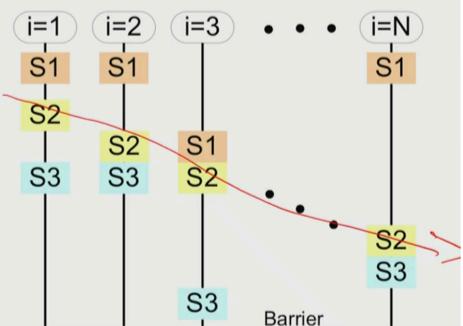
- `#pragma omp critical`
 - { ... }
- норма на использование
- ```
#pragma omp parallel - организовать
{
 #pragma omp for - непараллельный
 {
 } - для датас
}
```
- `omp barrier` - explicit "join", можно использовать с `nowait`
  - `omp master` - creates regions, forced - `masked (threads_ids)` - то же,
    - Only master thread executes code (within parallel region) то же
    - Others skip and continue иначе
  - `omp single`
    - Only one ~~one~~ thread executes code один, который не делит
    - Barrier at end конец

- **omp critical [Name]** - *некто ровно один раз одновременно*
  - can only be executed by one thread at a time
  - block executed sequentially
- **atomic**
  - memory updated atomically - *один раз, не делится на части*
  - only simple updates
    - increment/decrement, assignment
- **omp\_lock\_to object**
  - same as pthread mutex
  - init/destroy, set/unset, test
- **omp ordered** - *один раз в блоке*
  - Define block within parallel region *блоки между {}*
  - Block executed in same order as if sequential
- **section** - *блоки, которые выполняются параллельно, в пределах параллельной области*

|                                                       |                                  |
|-------------------------------------------------------|----------------------------------|
| <code>omp_init_lock(&amp;lockvar)</code>              | initialize a lock                |
| <code>omp_destroy_lock(&amp;lockvar)</code>           | destroy a lock                   |
| <code>omp_set_lock(&amp;lockvar)</code>               | set lock                         |
| <code>omp_unset_lock(&amp;lockvar)</code>             | free lock                        |
| <code>logicalvar = omp_test_lock(&amp;lockvar)</code> | check lock and possibly set lock |

```
#pragma omp for ordered
for (...)

 S1
 #pragma omp ordered
 { S2 }
 S3
}
```



all implicit barriers:

- single
- for
- parallel
- Sections
- workshare

# Data Races

- Two unsynchronized access, and one of them is write
- Hard to detect
  - Debugging changes timing
  - Use tools (but they make code slow)
- Dependencies – guarantee program correctness
  - Control dependence
    - Need to synchronize control flow
  - Data dependence
    - Need to synchronize data
  - Types:
    - Flow dependence: write  $\rightarrow$  read
    - Anti dependence: read  $\rightarrow$  write
    - Output dependence: write  $\rightarrow$  write

Проблема конкурентности  
race condition  
независимость от  
исполнения  
 $b = da$

## True Dependence Flow Dependence

$$O(S1) \cap I(S2) \neq \emptyset$$

output      input  
 $(S1 \delta^1 S2)$

The output of S1 overlaps with the input of S2

$$\begin{array}{l} S1: x = \dots \\ \dots \\ S2: \dots = x \end{array}$$

## Anti Dependence

$$I(S1) \cap O(S2) \neq \emptyset$$

$(S1 \delta^{-1} S2)$

The input of S1 overlaps with the output of S2

$$\begin{array}{l} S1: \dots = x \\ \dots \\ S2: x = \dots \end{array}$$

## Output Dependence

$$O(S1) \cap O(S2) \neq \emptyset$$

$(S1 \delta^o S2)$

The output of S1 and S2 write to same location

$$\begin{array}{l} x = \dots \\ x = \dots \end{array}$$

## Loop Independent Dependencies

All dependencies are with iterations  
No dependencies across iterations

### Example:

```
for (i = 0; i < 4; i++)
 S1: b[i] = 8;
 S2: a[i] = b[i] + 10;
```

Dependencies across iterations  
Computation in one iteration depends on data written in another

### Example:

```
for (i = 0; i < 4; i++)
 S1: b[i] = 8;
 S2: a[i] = b[i-1] + 10;
```

## Loops transformations

- Transformations:

- Loop interchange (nested loops)

- Interchange is safe if:

- dependencies are within same iteration ( $i=i$ )
  - or depends on previous steps ( $i+1=i$ )

- Move the chunkier loop outside & parallelize outer loop

- Loop distribution/fission

- If loop has multiple statements, one has inter-iteration dependencies, statements have no dependence cycle

- Split statements into two separate loops

- Increases granularity, useful if it makes serial loops parallelizable

- Loop fusion

- #### ○ Combine loops

- Can introduce loop carried dependencies & reduce parallelism

- Increases granularity, reduce thread overhead (?)

- Loop alignment

- Shift loop carried dependencies -> independence

- Peel off first and last iteration

```
do i=2,n
S1: a(i)= b(i)+2
S2: c(i)= a(i-1)
enddo
```

```

do i=1,n
S1: if (i>1) a(i)= b(i)+2
S2: if (i<n) c(i+1)= a(i) * 2
enddo

```

- Can eliminate loop carried dependencies if:

- Loop has no dependence cycle

- Distance of dependence is constant throughout iterations

# Memory Models

Memory model - describes the interactions of threads through memory and their shared use of data

- Memory/Cache Coherency: reasoning about updates to one memory location
- Memory Consistency: reasoning about updates to several memory locations

Sequential Consistency - A multiprocessor system is sequentially consistent if the result of any execution is the same as if the operations of all processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by the program

Processor consistency — Ensures writes from a single processor are seen in order by all processors. - ~~same hardware thread consistency~~  
~~No global order~~  
Requires cache coherency. Writes issued by a single processor are seen by all processors in the order they were given.

Weak consistency — Guarantees memory consistency only at synchronization points.

Release consistency — Ensures memory consistency at acquire and release.  
(Same as weak but with more relaxed)

## Memory synchronization points (or flush points)

- Entry and exit of parallel regions
- Barriers (implicit and explicit)
- Entry and exit of critical regions
- Use of OpenMP runtime locks
- Every task scheduling point (see Segment 19)

However, note the following are NOT synchronization points:

- Entry and exit of work sharing regions
- Entry and exit of masked regions

→ open spec?  
genus

## Sequential Consistency:

### 1. Program order

Within a single thread, memory operations (reads and writes) must appear to be executed in the order specified by the program.

T.e. read  $x$   
write  $t$

↓ should not be shifted!

А тае монга ингээсээ нь зүйл  
ондмын асуулт хамаалагдах  
и нь зүйл зажигжьец (бэсэн,  
харчилж)

write atomicity



### 2. Atomicity

Memory operations must be atomic.

### 3. Visibility

A write by one processor must be visible to all processors.

А чийн орчин төгрөгийн тохиолдсан, то бие даасанын төгрөгийн  
голчилж иржигжээ илрүүлж  
тийн нисээж, а нь чадигүй

### 4. Total order

All memory operations must be globally ordered.

T.e. ийнхийн салбарын нийгэмлэлийн  
и төслийн сэргүүлжсэн төгрөгийн  
нисээж, а нь чадигүй

## Weak Consistency:

Two categories of memory operations:

- Synchronization Operations

task scheduler, <sup>entries and exits</sup> parallel regions, barriers, etc. Cannot be reordered.

All preceding instructions must be completed before the execution of a synchronization operation.

Accesses to synchronisation instructions are sequentially consistent

Requires explicit sync. operations to ensure that memory operations are consistent and visible.

- Data Operations (memory instructions)

Just reads and writes. Can be reordered

## Release Consistency:

Same as weak but more relaxed

Instead of synchronization barrier, provides 2 barriers - Acquire and Release operations.

↑  
lock              ↑  
                    unlock

Acquire doesn't have to wait for the preceding instructions!

Release doesn't have to be completed before the following instructions

(But if there are instructions that need to wait until Release completes, we need to Acquire before them)

Why is the program order requirement hard?

- Network can reorder
- Use of write-buffers
- Compiler optimizations

Why is the in-sequence/atomicity requirement hard?

- A write has to take place instantly with respect to all processors
- Can lead to a system-wide serialization of writes

We need to relax consistency requirements

- Necessary to achieve well performing systems
  - Increase the level of concurrency
- Compensate with synchronization primitives
  - Need them anyway for proper synchronization of contents
  - Can implicitly execute consistency operations

Create the “illusion” of sequentially consistency

### Manual Synchronisation

Thread 1

Thread 2

→  
a = foo();  
→ flag = 1;

while(flag);  
b = a;

2 separate  
variables, compiler doesn't know  
they belong

What can go wrong?

- Compiler can reorder memory accesses
- Compiler can keep variables in registers

#pragma omp flush [list of variables]

### Thread 1

```
a = foo();
#pragma omp flush(a,flag)
flag = 1;
#pragma omp flush(flag)
```

modern  
reorder  
mechanism  
↓

### Thread 2

```
while (flag)
{
 #pragma omp flush(flag)
} // ensures there is no
 ↓ reordering
#pragma omp flush(a,flag)
b = a;
```

## Explicit Tasking

Tied Tasks - once started, will remain on the same thread.

Default behaviour, easy to reason about

Untied Tasks - can move to a different thread.  
 Execution can be interrupted and task moved. More flexibility  
 for runtime system and better resource utilization (che)  
 (but you can lose your ca

#pragma omp parallel {  
 #pragma omp single {  
 for (elem = first; elem; elem = elem->next)  
 #pragma omp task  
 process(elem); } }  
 } ← barrier  
 } ← default-tied tasks

organise tasks TONED  
over noTOK

no boundaries  
use noTOK!  
Onoverrunsafe

```

1 #include <iostream>
2 #include <omp.h>
3
4 int main(int argc, char *argv[])
5 {
6 #pragma omp parallel ← организуй преду
7 {
8 #pragma omp task ← организуй задачу, замените Spec
9 std::cout << "Hello World from task"
10 << std::endl;
11 }
12 return 0;
13 }

```

```

1 int main(int argc, char *argv[])
2 {
3 #pragma omp parallel
4 {
5 #pragma omp task
6 {
7 #pragma omp critical
8 std::cout << "Hello World from task, "
9 << "executed by thread: "
10 << omp_get_thread_num()
11 << std::endl;
12 }
13 }
14 return 0;
15 }

```

```

1 int main(int argc, char *argv[])
2 {
3 #pragma omp parallel
4 {
5 #pragma omp single ← один.Spec.войти задачу
6 {
7 for (int t = 0; t < omp_get_num_threads(); t++)
8 {
9 #pragma omp task
10 {
11 #pragma omp critical ← это одна организуй
12 std::cout << "Hello World from task, "
13 << "executed by thread: "
14 << omp_get_thread_num()
15 << std::endl;
16 }
17 }
18 }
19 }
20 return 0;
21 }

```

`#pragma omp task [clause [list]]`

Select clauses:

`if (scalar-expression)`

← если False, то это не  
вызывает task, и блок  
#pragma не наработан  
исключительно  
(например if n>10,  
но это вызывает разбивку  
на задачи)

`untied` — чтобы не параллеллизировать

Default (shared/move), private, firstprivate, shared  
default is firstprivate! ← не параллелизовать

`priority (value)`

execute this task before another task

Must NOT be used to rely on task ordering

Example: Tree Traversal

```
struct node {
 struct node *left;
 struct node *right;
};
```

```
void traverse(struct node *p) {
 if (p->left)
 #pragma omp task // p is firstprivate by default
 traverse(p->left);
 if (p->right)
 #pragma omp task // p is firstprivate by default
 traverse(p->right);
 process(p);
}

#pragma omp parallel
#pragma omp single
traverse(root);
```

#pragma omp taskwait - чтобы параллелизация не  
параллелизовала

`#pragma omp taskwait` - merges nonconcurrent tasks  
generally (i.e. tasks, created up  
since the beginning of  
the current parallel scope)

`#pragma omp taskyield` - specifies that the  
current task can be  
suspended  
This is an explicit  
task scheduling point

Implicit task scheduling points:

- task creation
- end of task
- taskwait
- barrier sync

Tasks dependencies:

`#pragma omp task shared(x) depend(out:x)`  
task1();

`#pragma omp task shared(x) depend(in:x)`  
task2();

For one can use names,  
as various no longer  
nonconcurrent tasks

Considerations:

- task granularity
- NUMA optimization

```

1 int fib(int n) {
2 int i, j;
3
4 if (n < 2) return n;
5
6 #pragma omp task shared(i) firstprivate(n)
7 i = fib(n - 1);
8
9 #pragma omp task shared(j) firstprivate(n)
10 j = fib(n - 2);
11
12 #pragma omp taskwait
13 return i + j;
14 }

```

```

1 int main(int argc, char** argv) {
2 int n = 30;
3
4 if(argc > 1)
5 n= atoi(argv[1]);
6
7 omp_set_num_threads(4);
8
9 #pragma omp parallel shared(n)
10 {
11 #pragma omp single
12 printf("fib(%d) = %d\n", n, fib(n));
13 }
14 }

```

Orphaned task uses `firstprivate` by default

Non-orphaned task inherit variables as shared

doesn't have parent, first task that was created  
task invoked within a task

```

1 #define T 30 // THRESHOLD
2
3 int fib(int n)
4 {
5 int i, j;
6
7 if (n < 2)
8 return n;
9
10 #pragma omp task shared(i) firstprivate(n)
11 i = fib(n - 1);
12
13 #pragma omp task shared(j) firstprivate(n) final(n < T)
14 j = fib(n - 2);
15
16 #pragma omp taskwait
17 return i + j;
18 }

```

запущен, но больше  
не организован  
↓ child-task,  
a hygen  
запускает  
пескоструй  
запускает  
такие

не корректные tasks! Тому же parallel  
Omp\_set\_nested(1); — если хотим организоватьested nested  
parallelism — parallel region inside  
parallel region

```
#pragma omp parallel num_threads(2)
{
 #pragma omp parallel num_threads(2) ← changed no
 { ... } ←rega
} ←внутри
каждогоrega
```

```
int a=1;
void parallel_function()
{
 int b=2, c=3;
 #pragma omp parallel shared(b)
 #pragma omp parallel private(b)
 { →private - orgame в ауге ковуо
 int d=4;
 #pragma omp task ← no умнр. firstprivate quo next repeat,
 { →объекту. внутрь parallel ауга
 int e=5;
 a // Shared - обьекты гарюо ✓
 // a= 1 ✓
 b // Firstprivate - оргаме в ауге ковуо ✓
 // b= ? ✓
 c // Shared - обьекты go сюда ✓
 // c= 3 ✓
 d // Firstprivate - обьекты → ауге
 // d= 4 ✓
 e // Private - обьекты внутрь task ✓
 // e= 5 ✓
 }
 }
}
```

*Annotations:*

- a // Shared - обьекты гарюо**
- b // Firstprivate - оргаме в ауге ковуо**
- c // Shared - обьекты go сюда**
- d // Firstprivate - обьекты → ауге**
- e // Private - обьекты внутрь task**

## Important Runtime Routines

Determine the number of threads for parallel regions

- `omp_set_num_threads(count)`

Query the maximum number of threads for team creation

- `numthreads = omp_get_max_threads()`

Query number of threads in the current team

- `numthreads = omp_get_num_threads()`

Query own thread number (0..n-1)

- `iam = omp_get_thread_num()`

Query the number of processors

- `numprocs = omp_get_num_procs()`

internal control variables

## Relevant Environment Variables (ICVs):

### `OMP_NUM_THREADS=4`

- Number of threads in a team of a parallel region

### `OMP_SCHEDULE="dynamic" OMP_SCHEDULE="GUIDED,4"`

- Selects scheduling strategy to be applied at runtime
- Schedule clause in the code takes precedence

### `OMP_DYNAMIC=TRUE`

- Allow runtime system to determine the number of threads

### `OMP_NESTED=TRUE`

- Allow nesting of parallel regions
- If supported by the runtime

Runtime routines to query and set ICVs

## Typical performance considerations

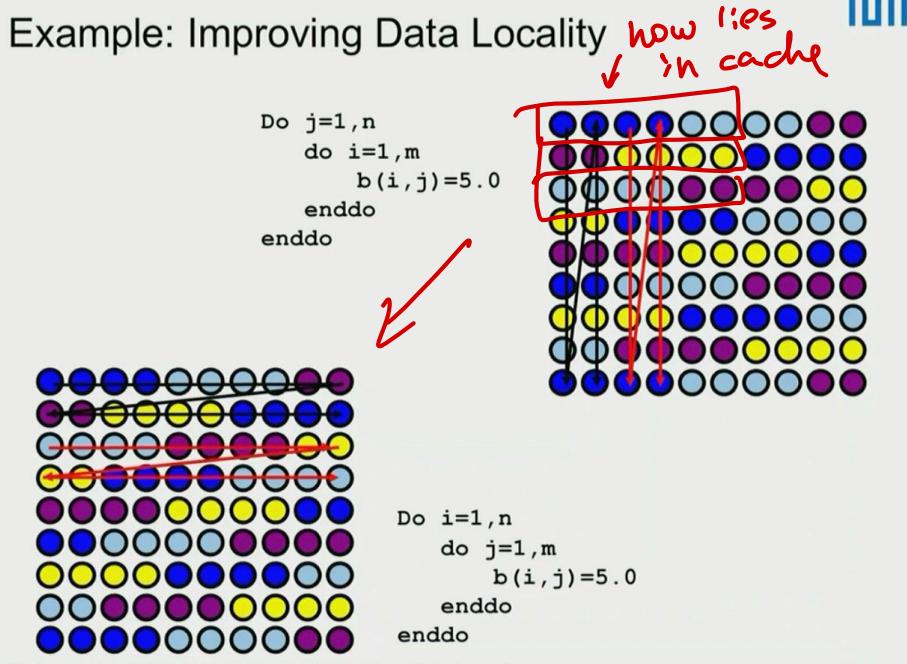
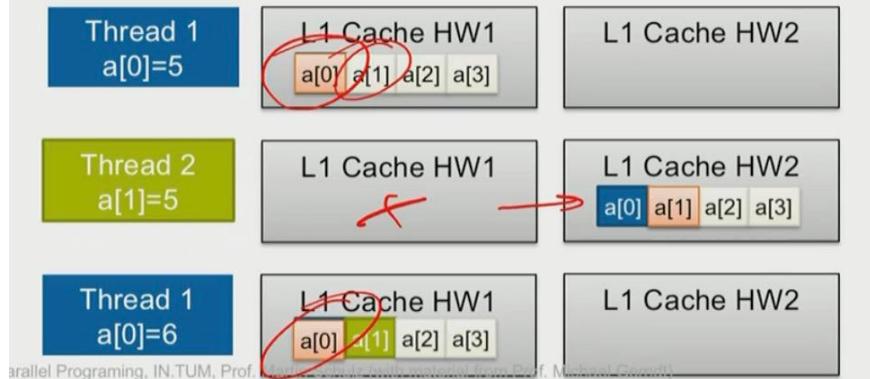
- Threading overhead
  - Expensive to starting/ending a thread - better do once (as omp does)
  - Scheduling is also not that cheap - better not to have too little work per thread
- Synchronisation overhead
  - Lock too often
  - Lock too much (lots of waiting)
- Cache behaviour and locality
  - Have to keep cache to reuse
  - Locality - better take from nearest cache
- Thread mapping
  - How far to put thread from memory

Compare-and-swap ← another solution  
atomic operation → in loop

### ABA Problem

- A tries to remove an element
  - Has to wait
- B removes the element
- Frees it
- Reallocates it (new element, same address)
- B inserts element again
- A compares against the old/new element

## Problem of false sharing

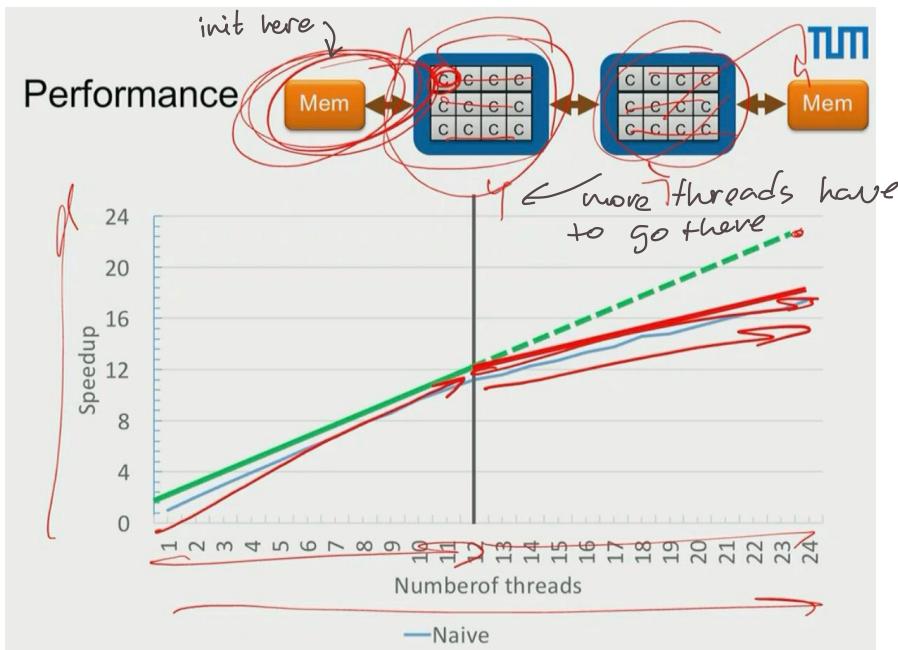


# Placing Memory

First touch:

- Allocation creates virtual space
- First access allocates physical memory located close to accessing thread
- No later migration

That's why this happens:



## Control by ICV **OMP\_PROC\_BIND**

- **true**: threads are bound to cores/hw-threads
- **false**: threads are not bound and can be migrated
- **master**: new threads are always co-located with master
- **close**: new threads are located "close" to their master threads
- **spread**: new threads are spread out as much as possible

Description of available locations, also known as places: **OMP\_PLACES** ICV

- Describes ordered list and hierarchy of all available hardware threads
- Example: {0:4}, {4:4}, {8:4}, {12:4}
- Master thread executes on first place

## SIMD

To support SIMD, processor needs special vector registers and instruction set support

SSE - Streaming SIMD extensions - 128 bit registers available on all x64-86 intel processors

AVX - SSE + SSE = 256 bit registers

`<immintrin.h>` - intrinsics define datatypes and operations with them  
↑  
AVX gcc

## AVX datatypes:

- -m256 -  $8 \times 32$  bit float
- -m256d -  $4 \times 64$  bit double
- -m256i -  $8 \times 32$  bit int or  $4 \times 64$  bit int
  
- -mm256 - loadu\_ps
  - unaligned
  - ↑
  - ↖
  - packed single
  - i.e. one packed
  - on mem
  - scope, we ka
  - scope
  
- -mm-add\_ps - add packed floating point numbers
  
- -mm-movelh\_ps - byte operations
  
- -sd - scalar operations on the least significant data element (0-31 bits for floats)

Things to consider:

- Latency - number of clock cycles
- Throughput - rate of executions of instructions

# MPI: Message Passing Interface

aiming 10 microsec latency max



network transportation

Types of operations:

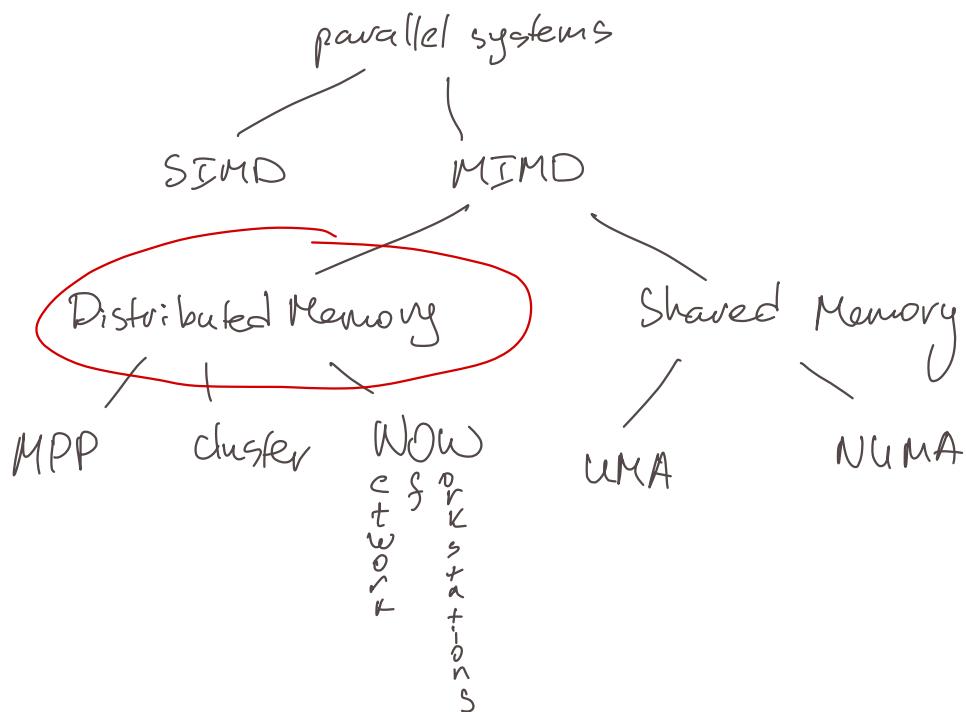
Point to Point

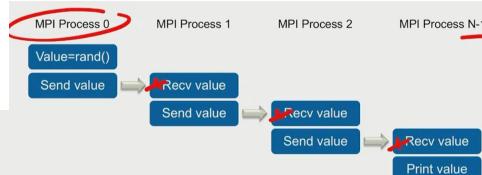


Collectives



One Sided





```

#include <stdio.h>
#include <stdlib.h>
#include <mpi.h>

int main (int argc, char** argv)
{
 double value;
 int size, rank;
 MPI_Status s;
 MPI_Init (&argc, &argv); // init the job. Replaces main()
 MPI_Comm_size (MPI_COMM_WORLD, &size); // strips MPI args
 MPI_Comm_rank (MPI_COMM_WORLD, &rank); // any number from argc, argv
 value=MPI_Wtime(); // from
 printf("MPI Process %d of %d (value=%f)\n", rank, size, value);
 if (rank>0)
 MPI_Recv(&value, 1, MPI_DOUBLE, rank-1, 0, MPI_COMM_WORLD, &s);
 if (rank<size-1)
 MPI_Send(&value, 1, MPI_DOUBLE, rank+1, 0, MPI_COMM_WORLD); // context
 if (rank==size-1)
 printf("Value from MPI Process 0: %f\n",value);
 MPI_Finalize (); // shut down MPI - resources released
}

```

How many arguments remain after MPI...?

return error code

what to send new meaning type where to send to

not try same approach

any number from argc, argv

from

Should be called by all MPI processes exactly once!

Parallel Programming. IN.TUM. Prof. Martin Schulz (with material from Prof. Michael Gerndt)

Memory management requires MPI...

Communicators:

**MPI\_COMM\_WORLD**

- Group of MPI processes
- Work with communication context - always!

Communications across communicators are not possible!

- Different context :C

But! one MPI process can be in multiple communicators

Default:

**MPI\_COMM\_WORLD**: all initial MPI processes

**MPI\_COMM\_SELF**: only own MPI process

## MPI Processes vs OS Processes vs Rank

↗

MPI concept!

Communication endpoint that can send and receive messages. Members of MPI groups and Communicators.

But:

Usually implemented as processes (OS)

Rank is identifier of an MPI process relative to an MPI communicator

MPI\_Send - initiate message send and blocks until send buffer can be reused

↑ f.e. nongraceful nonreturn  
possible deadlocks!

### MPI\_Send

Send one message to another MPI process on a given communicator

```
int MPI_Send (void *buf, int count, MPI_Datatype dtype,
 int dest, int tag, MPI_Comm comm)
```

|          |                                                               |
|----------|---------------------------------------------------------------|
| IN buf   | Address of the send buffer                                    |
| IN count | Number of data elements of type <code>dtype</code> to be sent |
| IN dtype | Data type                                                     |
| IN dest  | Receiver (rank of target MPI process in <code>comm</code> )   |
| IN tag   | Message tag                                                   |
| IN comm  | Communicator                                                  |

Initiates message send and blocks until send buffer can be reused.

- Note: possibilities of deadlocks (!)

## MPI\_Recv

Receives a message from another MPI process on a given communicator

```
int MPI_Recv (void *buf, int count, MPI_Datatype dtype,
int source, int tag, MPI_Comm comm, MPI_Status *status)
```

*source tag have to match with sender*

|            |                                                                   |
|------------|-------------------------------------------------------------------|
| OUT buf    | Address of the receive buffer                                     |
| IN count   | Number of data elements of type <code>dtype</code> to be received |
| IN dtype   | Data type                                                         |
| IN source  | Sender (rank of source MPI process in <code>comm</code> )         |
| IN tag     | Message tag                                                       |
| IN comm    | Communicator                                                      |
| OUT status | Status information                                                |

Blocking receive of a message

- Blocks execution until a matching message has been received
- Note: possibilities of deadlocks (!)

Sender and tag can be specified as wild cards

- `MPI_ANY_SOURCE` and `MPI_ANY_TAG`
- Can be used instead of `source` and `tag` arguments
- There is no wild card for the communicator/context

We have a guaranteed messages order!  
(but this can be overwritten)

Messages are supposed to have matched types  
Casting types can be dangerous!

Message buffers must be at least as long as the message

## Most common MPI types

|                                                           |                                       |
|-----------------------------------------------------------|---------------------------------------|
| MPI_CHAR                                                  | char (treated as printable character) |
| MPI_SHORT                                                 | signed short int                      |
| MPI_INT                                                   | signed int                            |
| MPI_LONG                                                  | signed long int                       |
| MPI_LONG_LONG                                             | signed long long int                  |
| MPI_UNSIGNED_SHORT                                        | unsigned short int                    |
| MPI_UNSIGNED                                              | unsigned int                          |
| MPI_UNSIGNED_LONG                                         | unsigned long int                     |
| MPI_UNSIGNED_LONG_LONG                                    | unsigned long long int                |
| MPI_FLOAT                                                 | float                                 |
| MPI_DOUBLE                                                | double                                |
| MPI_C_COMPLEX                                             | float _Complex                        |
| MPI_C_DOUBLE_COMPLEX                                      | double _Complex                       |
| MPI_BYTE                                                  |                                       |
| MPI_PACKED                                                |                                       |
| Custom datatypes can be defined through special routines. |                                       |

In C, **MPI\_Status** is a structure that includes the following fields:

- MPI\_SOURCE: sender of the message
- MPI\_TAG: message tag
- MPI\_ERROR: error code

Important for wild card communication

- Identifies the actual channel the message was transmitted on

The actual length of the received message can be determined via

```
int MPI_Get_count(const MPI_Status *status,
 MPI_Datatype dtype, int *count)
```

|           |                                                  |
|-----------|--------------------------------------------------|
| IN status | MPI_Status object returned from the receive call |
| IN dtype  | Data type used in receive                        |
| OUT count | Number of data elements actually received        |

```
double MPI_Wtime(void)
```

Return wall clock time since a reference time stamp in the past

- Unit: seconds
- Reference time stamp is fixed for one execution
- Time can be global, but doesn't have to be
  - In fact, in most cases it is not
  - Users have to synchronize, if needed
- Note: return value is not an error code

## Collective Operations

Properties

- Must be executed by all processes of the communicator
- – Must be executed in the same sequence
- Arguments across all MPI processes must be consistent (or equal)

MPI provides three classes of collective operations

- Synchronization
  - Barrier
- Communication, e.g.,
  - Broadcast
  - Gather – collect from all nodes
  - Scatter – send chunks to all nodes (distribute work)
- Reduction, e.g.,
  - Global value returned to one or all processes.
  - Combination with subsequent scatter.
  - Parallel prefix operations

## MPI\_Barrier

```
int MPI_Barrier (MPI_Comm comm)
```

IN comm: Communicator

This operation synchronizes all MPI processes, i.e., no MPI process can return from the call until all MPI processes have called it.

```
#include "mpi.h"

main (int argc,char *argv[]){
...
MPI_Comm_size(MPI_COMM_WORLD, &np);
MPI_Comm_rank(MPI_COMM_WORLD, &myid);

MPI_Barrier(MPI_COMM_WORLD);
...
```

## Broadcast Communication: MPI\_Bcast

```
int MPI_Bcast (void *buf, int count, MPI_Datatype dtype,
 int root, MPI_Comm comm)
```

IN buf: Address of send/receive buffer

IN count: Number of elements

IN dtype: Data type

IN root: Rank of send/broadcasting MPI process - our rank

IN comm: Communicator (send to everyone but root)

The contents of the send buffer is copied to all other MPI processes

Collective vs P2P operations

- No tag
- Number of elements sent must be equal to number of elements received

The routines do not necessarily synchronize processes wrt. P2P

Broadcast doesn't guarantee everyone received the message!

## Collective vs P2P operations

- No tag !!
- Number of elements sent must be equal to number of elements received

→ y collectives

```
switch (rank) {
 case 0:
 MPI_Bcast (buf1, ct, tp, 0, comm);
 MPI_Send (buf2, ct, tp, 1, tag, comm);
 break;
 case 1:
 MPI_Recv (buf2, ct, tp, MPI_ANY_SOURCE, tag, ...);
 MPI_Bcast (buf1, ct, tp, 0, comm);
 MPI_Recv (buf2, ct, tp, MPI_ANY_SOURCE, tag, ...);
 break;
 case 2:
 MPI_Send (buf2, ct, tp, 1, tag, comm);
 MPI_Bcast (buf1, ct, tp, 0, comm);
 break;
}
```

Communicators

Why want anything but MPI\_COMM\_WORLD

- Creation of subgroups for collective communication
- Isolation between communication operations

Duplicating  
Communicators

Isolation important for libraries

- Different libraries should use different communicators
- Avoids "cross-talk" and side effects

Typical strategy for libraries

- Libraries get a communicator passed at initialization (often WORLD)
- The library then clones the communicator and stores the handle

```
int MPI_Comm_dup(MPI_Comm comm, MPI_Comm *newcomm)
 IN comm communicator (handle)
 OUT newcomm copy of comm (handle)
```

↑ my hero even Berger!

Collective operation

Communicators should be freed, when no longer in use

```
int MPI_Comm_free(MPI_Comm *comm)
```

← so application  
doesn't starve

## Creating Subcommunicators

0 1 2  
3 4 5  
6 7 .  
8 9 .  
10 11 .  
12 13 14  
15 16

```
int MPI_Comm_split(MPI_Comm comm, int color, int key, MPI_Comm *newcomm)
```

highest/  
lowest  
values

IN comm      Parent communicator  
IN color     Subset color  
IN key        Key to determine rank order in new communicator  
OUT newcomm   New communicator

is using to say me  
newrank 1-  
new rank order  
is saying 1)

Creates new communicator(s)

- Collective operation over parent communicator
- All MPI processes that pass the same **color**, will be in same new communicator
- **Key** argument determines the rank order in the new communicator

MPI processes can opt out

- Pass **MPI\_UNDEFINED** as color
- Will return **MPI\_COMM\_NULL** as new communicator

Collective operation

Tome ayoko  
is always do merge!

every newrank order is to one - place bigger  
value is always to say me!

## Example: Row and Column Communicators

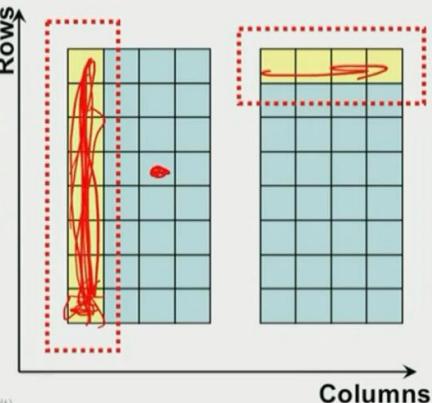
```
int rank, size, col, row;
MPI_Comm row_comm, col_comm;

MPI_Comm_size (MPI_COMM_WORLD, &size);
MPI_Comm_rank (MPI_COMM_WORLD, &rank);

Row = rank % N;
Col = rank / N;

MPI_Comm_split(MPI_COMM_WORLD,
 row, col, &row_comm);
MPI_Comm_split(MPI_COMM_WORLD,
 col, row, &col_comm);
...
MPI_Bcast(buf, 1, MPI_INT, 0, col_comm);

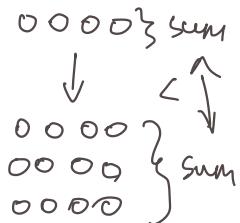
(works also for Send & Recv)
```



# Weak Scaling

Keeping the problem size per HW thread/core/node constant

- Larger machine -> larger problem
  - Expanding the problem domain
  - Increasing refinement
- Traditionally most common way to deal with scaling in HPC



## Assumptions

- Machines are never big enough
- Fully loading a node/core/HW thread is the right thing to do
  - Exploiting the resources on one unit and keeping that constant ↗

## Advantages

- Execution properties (like cache behavior) often stay fixed
- Easier to scale, as overheads stay roughly constant as well

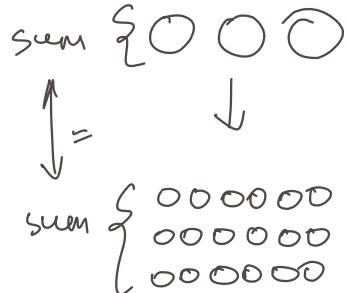
## Challenges

- Keeping load constant can be tricky for complex applications (2D/3D problems)
- Multiple ways to repartition a workload

# Strong Scaling

Keeping the total problem constant

- Larger machine -> (hopefully) faster execution
- Need to adjust problem distribution
- Traditionally not the most common type of scaling
  - But becoming rapidly more and more relevant



Assumptions:

- The machine is big enough for the problem
- Goal: reducing time to solution and increasing throughput

Needed for time critical tasks

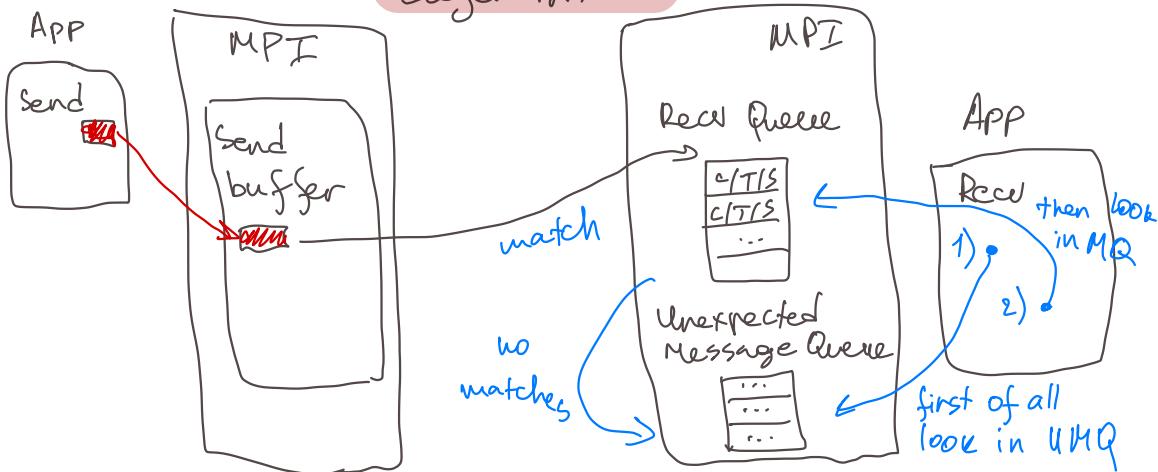
- Emergency response to natural disasters
- Real-time simulations
- Large-scale ensemble calculations

Challenges

- Harder to scale, as overheads grow with smaller per HW thread workloads
- Changing executing characteristics (cache miss rates, etc.)

## How Messages are Transferred (typically)

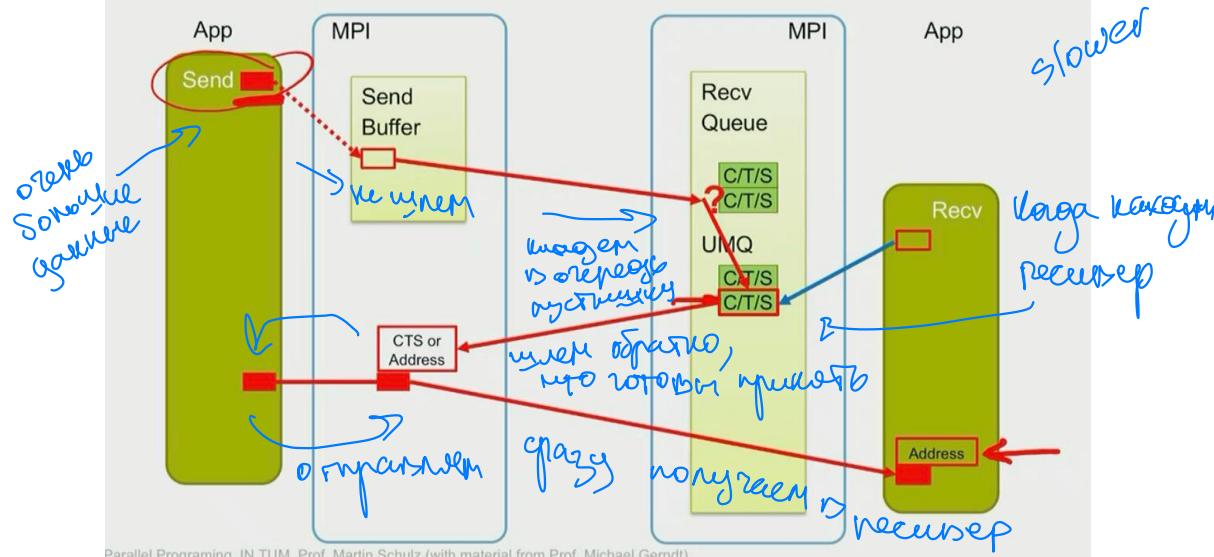
### Eager Protocol



# Rendevouz Protocol

Sender

Receiver



Parallel Programming, IN.TUM, Prof. Martin Schulz (with material from Prof. Michael Gerndt)

## Eager vs. Rendezvous Protocol

### Eager protocol

- Send messages directly
  - Advantage: avoids extra handshake
  - Disadvantage: adds extra copies
- Suitable for short messsages with low copy overhead

### Rendezvous protocol

- Send header to retrieve address, then send message
  - Advantage: Can deposit message directly
  - Disadvantage: extra handshake
- Suitable for long messsages with high copy overhead

### Cross-over point

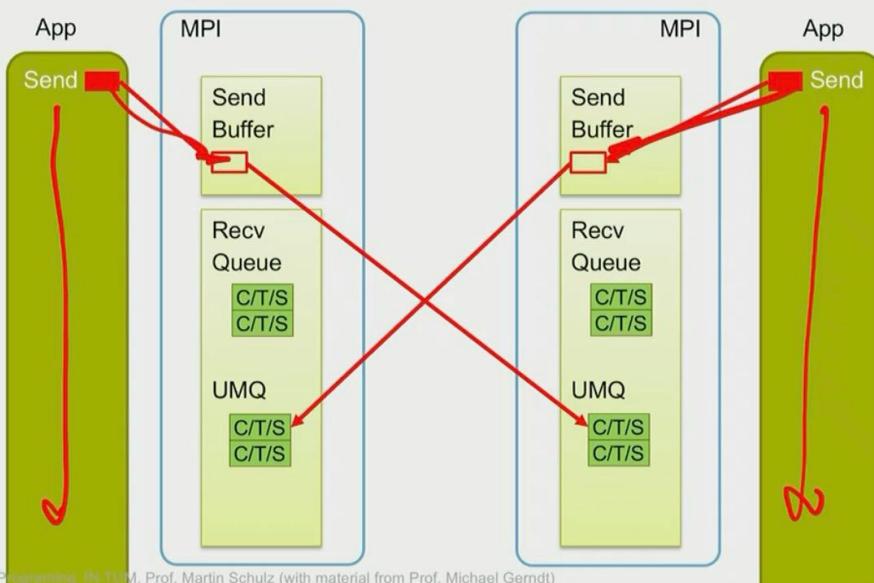
- MPI implementations typically switch protocol based on message size
- Typically defined as eager limit
- Depends on system setup
- Can often be modified

Однако иск. сразу оба, ожидая двух других соединений, грузовой трафик ман.

Последний этап иск. последовательно при выполнении разрешения соединений, возможно, это произойдет

## Deadlock with Rendezvous Protocol

Sender



Receiver

Send

Parallel Programming, Dr. Hartmann, Prof. Martin Schulz (with material from Prof. Michael Gerndt)

7

Transfer data elements between two "adjacent"

- Same or similar data exchanged between two MPI processes
- Example: halo exchange



Bidirectional  
exchange

MPI\_Send( to: 0 )      MPI\_Send( to: 0 )  
MPI\_Recv( from: 1 )      MPI\_Recv( from: 0 )

Problem: deadlock possible

- MPI\_Send may need resources on receiver side

Changing the order of send and receive

long before receiver  
MPI\_Send( to: 1 )      MPI\_Recv( from: 0 )  
MPI\_Recv( from: 1 )      MPI\_Send( to: 0 )

Helps to avoid deadlock,  
but increases time

### MPI\_Send

- Standard send operation, no extra synchronization
- Sender side handling implementation defined

message is  
- guaranteed to  
be sent eventually  
blocks until sent?

### MPI\_Bsend

- Buffered Send
- Force the use of a send buffer
- Returns immediately, but costs resources

Не блокируется,  
но это не гарантирует  
что сообщение будет  
получено в то же время  
когда было отправлено.  
Блокируется, пока  
приемник не получит  
(получить, что соответствует  
тому сообщению, которое было  
запрошено)

### MPI\_Ssend

- Synchronous send
- Only returns once receive has started
- Adds extra synchronization, but can be costly

### MPI\_Rsend

- Ready Send
- User must ensure that receive has been posted
- Enables faster communication, but needs implicit synchronization

Быстро, но требует чтобы приемник  
мог соответствовать  
тому сообщению, которое было  
запрошено и  
одновременно  
получил

## Bidirectional Send/Recv

```
int MPI_Sendrecv (void *sendbuf, int sendcount, MPI_Datatype sendtype,
 int dest, int sendtag,
 void *recvbuf, int recvcount, MPI_Datatype recvtype,
 int source, MPI_Datatype recvtag,
 MPI_Comm comm, MPI_Status *status)
```

### Combined blocking send and receive

- Can be matched with regular send/recv calls
- Can be matched with other Sendrecv calls with other destinations

```
int MPI_Sendrecv_replace(void *buf, int count, MPI_Datatype datatype,
 int dest, int sendtag,
 int source, int recvtag,
 MPI_Comm comm, MPI_Status *status)
```

### Same principle, but combined send and receive buffer

- Requires same message types

# Checking the UMQ with MPI\_Probe/Iprobe

Allow polling of incoming messages without actually receiving them.

- **MPI\_Probe**

- Blocks until a matching message was found
- Input: source, tag, communicator
- Output: status object

- **MPI\_Iprobe**

- Terminates after checking whether a matching messages is available
- Input: source, tag, communicator
- Output: flag (message found yes/no) and status object

## Usage

- Both operations inspect the Unexpected message queue
- Allows to wait for a message before issuing the receive
- Useful for wildcard receives
  - Issue probe call with wildcard
  - Call returns concrete source/tag pair
- MPI\_Iprobe can be used to overlap wait time with useful work.

*- size nrg. Monne corrigit noglog en niet  
reciever  
TO we canoes,  
no we don't keep.*

# Nonblocking P2P Operations in MPI

Initiation routines as counterparts to blocking P2P operations

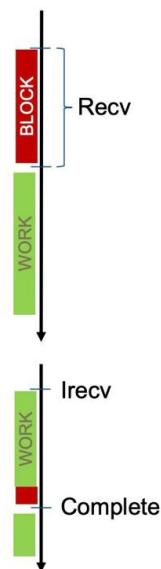
- Similar arguments, same semantics
- Can be matched with blocking operations
- Returns request object to track progress

Initiation call only starts operation

- Completion is indicated later
- Send and receive buffers are off limits before completion
- Semantics of operation

- **MPI\_Request**

- Request object
- User allocates object, but MPI maintains state
- Initiated operations must be completed



```
int MPI_Isend (void* buf, int count, MPI_Datatype dtype, int dest,
int tag, MPI_Comm comm, MPI_Request* request)
```

IN buf: Send buffer  
IN count: Number of elements  
IN dtype: Data type  
IN dest: Receiver  
IN tag: Tag  
IN comm: Communicator  
OUT request: Reference to request

No status param!  
As we do not know  
the status yet

```
int MPI_Irecv (void* buf, int count, MPI_Datatype dtype, int source,
int tag, MPI_Comm comm, MPI_Request* request)
```

IN buf: receive buffer  
IN count: number of entries in receive buffer  
IN dtype: data type  
IN dest: sender  
IN tag: tag  
IN comm: communicator  
OUT request: reference to request

### MPI\_Isend

- Standard send operation, no extra synchronization
- Sender side handling implementation defined

See MPI\_Send  
buffer, no  
completion -  
last parameter

### MPI\_Ibsend

- Buffered Send
- Force the use of a send buffer
- Completes immediately, but costs resources

### MPI\_Issend

- Synchronous send
- Only completes once receive has started
- Adds extra synchronization, but can be costly

### MPI\_Irsend

- Ready Send
- User must ensure that receive has been posted at Irsend time
- Enables faster communication, but needs implicit synchronization

All blocking collectives in MPI have a non-blocking counterpart

## Completion Operations

### Option 1: Blocking completion

```
int MPI_Wait (MPI_Request *request, MPI_Status *status)
```

INOUT request: *request*  
OUT status: *status of completed operation*

Returns if ...

- request is complete
- request is `MPI_REQUEST_NULL`

Sets request to `MPI_REQUEST_NULL`

Frees all resources associated with request

```
{
 MPI_Request req;
 MPI_Status status;
 int msg[10];
 ...
 MPI_Irecv(msg, 10, MPI_INT, MPI_ANY_SOURCE, 42,
 MPI_COMM_WORLD, &req);
 ...
 <do work>
 ...
 MPI_Wait(&req, &status);
 printf("Processing message from %i\n",
 status.MPI_SOURCE);
 ...
}
```

## Option 2: Nonblocking/Polling completion

Ke nogen, a request  
↳ bovenopges. gesloten  
↳ no cyclic bds

```
int MPI_Test (MPI_Request *request, int *flag, MPI_Status *status)
```

INOUT request:

*request*

OUT flag:

flag, whether complete or not

OUT status:

*status of completed operation*

Returns if immediately

flag set to True if request is complete

request is **MPI\_REQUEST\_NULL**

If completed,

Sets request to **MPI\_REQUEST\_NULL**

Frees all resources associated with request

```
{
 MPI_Request req;
 MPI_Status status;
 int msg[10], flag;
 ...
 MPI_Irecv(msg, 10, MPI_INT, MPI_ANY_SOURCE, 42,
 MPI_COMM_WORLD, &req);
 do
 {
 ...
 <do work>
 ...
 MPI_Test(&req, &flag, &status);
 } while (flag==0);

 printf("Processing message from %i\n",
 status.MPI_SOURCE);
}
```

↳ vanpuren?  
↳ execute

### MPI\_Wait

- Blocking wait for one request
- Output: status object

### MPI\_Waitall

- Blocking wait for several requests
  - Input: array of requests
  - Some can be MPI\_REQUEST\_NULL
- Returns once all are complete
  - Output: array of status objects

↙ array of  
requests

↙ also array

### MPI\_Waitany (MPI\_Waitsome)

- Blocking wait for several requests
  - Input: array of requests
  - Some can be MPI\_REQUEST\_NULL
- Returns once at least one is complete
  - Output: "completed" index (array)
  - Output: (array of) status object(s)

### MPI\_Test

- Checks completion of one request
- Output: flag, True if complete
- Output: status object

### MPI\_Testall

- Checks completion of several requests
  - Input: array of requests
  - Some can be MPI\_REQUEST\_NULL
- Returns immediately
  - Output: flag, True if all complete
  - Output: array of status objects

### MPI\_Testany (MPI\_Testsome)

- Checks completion of several requests
  - Input: array of requests
  - Some can be MPI\_REQUEST\_NULL
- Returns immediately
  - Output: flag, True if at least one complete
  - Output: "completed" index (array)
  - Output: (array of) status object(s)

Parallel Programming, IN.TUM, Prof. Martin Schulz (with material from Prof. Michael Gerndt)

## Terminology before MPI 4.0

### Blocking routines

- After return all buffers can be reused again
- E.g., after MPI\_Send, the data in the send buffer has been consumed by MPI

### Nonblocking routines

- After return buffers are still under MPI control
- E.g., after MPI\_Isend, the data buffers cannot be written until MPI\_Wait

## MPI Blocking ≠ English Blocking

↙ makes sense  
for you  
me maybe

### Non-Local routines

- Requires some (specific) other execution to run in another process
- E.g., for MPI\_Recv to complete, another process must call MPI\_(X)send

### Local routines

- Not non-local

## MPI Local ≠ Executes only locally in one MPI process

All MPI routines are "weak local", i.e., assume progress on all MPI processes

# Terminology after MPI 4.0

## Non-Local procedures

- Requires some (specific) other execution to run in another process
- E.g., for MPI\_Recv to complete, another process must call MPI\_(X)send

## Local procedures

- Not non-local, now explicitly weak local

## Completing procedures (was blocking)

- After return all buffers can be reused again
- E.g., after MPI\_Send, the data in the send buffer has been consumed by MPI

## Incomplete procedures (was non-blocking)

- After return buffers are still under MPI control
- E.g., after MPI\_Isend, the data buffers cannot be written until MPI\_Wait

## Non-Blocking procedures

- incomplete and local

*isynchronous "non-blocking",  
nam kurewo ie  
meyas, a nges  
nge-ro ka opne  
zero-ro non-local-  
u for c kum*

MPI Blocking      ≠      English Blocking

## Blocking procedures

- Not “non-blocking”

```
MPI_Request req[2];
MPI_Status status[2];
int msg_in[10], msg_out[10];

MPI_Isend(msg_out, 10, MPI_INT, 1, 42, MPI_COMM_WORLD, &(req[0]));
MPI_Irecv(msg_in, 10, MPI_INT, 1, 42, MPI_COMM_WORLD, &(req[1]));

MPI_Waitall(2, req, status);

/* status[1] is the receive status */
```

# MPI\_Reduce

```
int MPI_Reduce (void* sbuf, void* rbuf, int count, MPI_Datatype dtype,
 MPI_Op op, int root, MPI_Comm comm)
```

|           |                                                          |
|-----------|----------------------------------------------------------|
| IN sbuf:  | <i>Send buffer</i>                                       |
| OUT rbuf: | <i>Receive buffer</i>                                    |
| IN count: | <i>Number of elements</i>                                |
| IN dtype: | <i>Data type</i>                                         |
| IN op:    | <i>Operation</i>                                         |
| IN root:  | <i>Rank of MPI process receiving the reduced results</i> |
| IN comm:  | <i>Communicator</i>                                      |

Combines the elements in the send buffer according to “op”

Delivers the result to root

## MPI Reductions

The data of the processes are combined via a specified operation

- Predefined operators, like ‘+’
- User defined operators (see `MPI_Op_create` and `MPI_Op_free`)

Often implemented in optimized way, e.g., using tree structures

### `MPI_Reduce`

- Reduce elements from all processes based on operation

### `MPI_Allreduce`

- Reduction followed by broadcast

### `MPI_Reducescatter_block`

- Reduction followed by a scatter operation (equal chunks)

### `MPI_Reducescatter`

- Reduction followed by a scatter operation (varying chunks)

`MPI_Scan` – *redukciū užduočių. X<sub>0</sub>, vėliau X<sub>0</sub>+X<sub>1</sub>, tada X<sub>0</sub>+X<sub>1</sub>+X<sub>2</sub> ...*

- Reduction using a prefix operation

# Available Reduction Operators

|          |                             |
|----------|-----------------------------|
| MPI_MAX  | maximum                     |
| MPI_MIN  | minimum                     |
| MPI_SUM  | sum                         |
| MPI_PROD | product                     |
| MPI_LAND | logical and                 |
| MPI_BAND | bit-wise and                |
| MPI_LOR  | logical or                  |
| MPI_BOR  | bit-wise or                 |
| MPI_LXOR | logical exclusive or (xor)  |
| MPI_BXOR | bit-wise exclusive or (xor) |

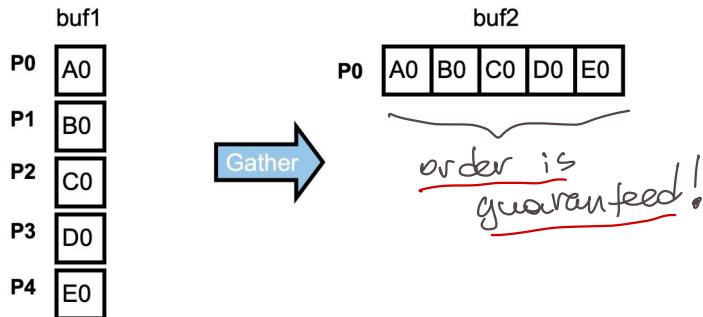
Standard has precise mapping of operations to allowable datatypes

Use with speciality datatype (tupel: <value> + <location>)

|            |                        |
|------------|------------------------|
| MPI_MAXLOC | max value and location |
| MPI_MINLOC | min value and location |

## Gather Operation

Get data from all MPI processes into one local array



## MPI\_Gather

```
int MPI_Gather (void *sendbuf, int sendcount, MPI_Datatype sendtype,
 void* recvbuf, int recvcount, MPI_Datatype recvtype,
 int root, MPI_Comm comm)
```

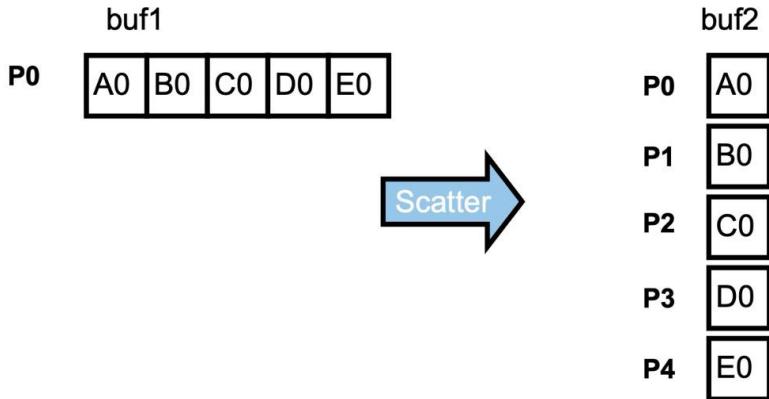
|               |                                                                 |
|---------------|-----------------------------------------------------------------|
| IN sendbuf:   | <i>Send buffer</i>                                              |
| IN sendcount: | <i>Number of elements to be sent to the root</i>                |
| IN sendtype:  | <i>Data type</i>                                                |
| OUT recvbuf:  | <i>Receive buffer</i>                                           |
| IN recvcount: | <i>Number of elements to be received from each MPI process.</i> |
| IN recvtype:  | <i>Data type</i>                                                |
| IN root:      | <i>Rank of receiving MPI Process</i>                            |
| IN comm       | <i>Communicator</i>                                             |

The root receives from all processes the data in the send buffer.

It stores the data in the receive buffer ordered by the process number of the senders.

## Scatter

with no broadcast, no user payette gather



## MPI\_Scatter

```
int MPI_Scatter(void *sendbuf, int sendcount, MPI_Datatype sendtype,
 void* recvbuf, int recvcount, MPI_Datatype recvtype,
 int root, MPI_Comm comm)
```

|               |                                                           |
|---------------|-----------------------------------------------------------|
| IN sendbuf:   | <i>Send buffer</i>                                        |
| IN sendcount: | <i>Number of elements sent to each MPI process</i>        |
| IN sendtype:  | <i>Data type</i>                                          |
| OUT recvbuf:  | <i>Receive buffer</i>                                     |
| IN recvcount: | <i>Number of elements to be received by a MPI process</i> |
| IN recvtype:  | <i>Data type</i>                                          |
| IN root:      | <i>Rank of sending/scattering MPI process</i>             |
| IN comm:      | <i>Communicator</i>                                       |

The root sends a part of its send buffer to each process.

Process  $k$  receives  $sendcount$  elements starting with  $sendbuf + k * sendcount$ .

*form queue no context payload kon- $\rightarrow$*

### MPI\_Gatherv

- Gather operation with different number of data elements per process

### MPI\_Allgather

- Gather operation with broadcast to all processes

*gather & broadcast*

### MPI\_Allgatherv

- Gather operation with different number of data elements AND broadcast

### MPI\_Scatterv

- Scatter operation with different number of data elements per process

### MPI\_Alltoall

- Every process sends to every other process (different data)

### MPI\_Alltoallv

- All to all with different numbers of data elements per process

### MPI\_Alltoallw

- All to all with different numbers of data elements and types per process

## MPI\_Gatherv

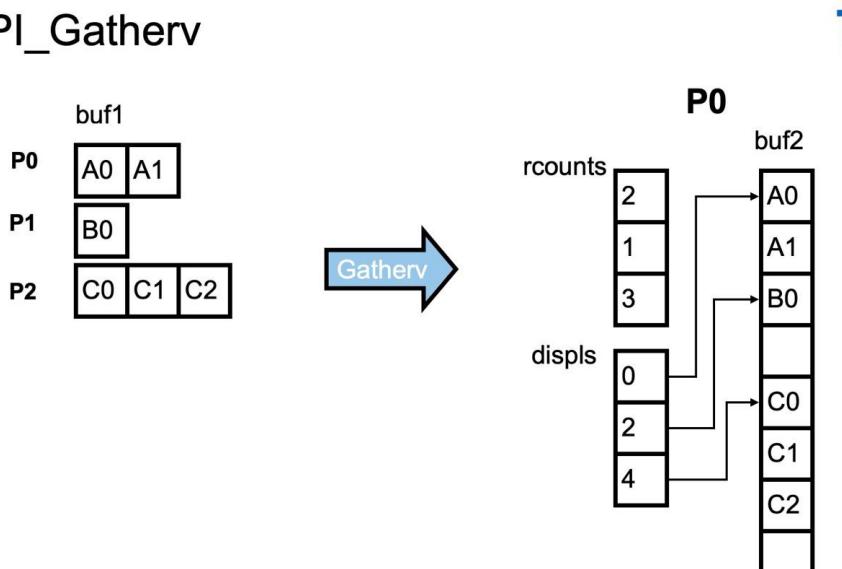
```
int MPI_Gatherv (void *sendbuf, int sendcount, MPI_Datatype sendtype,
 void* recvbuf, int *recvcount, int *displs,
 MPI_Datatype recvtype, int root, MPI_Comm comm)
```

|               |                                                                               |
|---------------|-------------------------------------------------------------------------------|
| IN sendbuf    | <i>Send buffer</i>                                                            |
| IN sendcount  | <i>Number of elements sent</i>                                                |
| IN sendtype   | <i>Send type</i>                                                              |
| OUT recvbuf   | <i>Receive buffer</i>                                                         |
| IN recvcounts | <i>Array with the number of elements to be received from each MPI process</i> |
| IN displs     | <i>Array with the first index in the receive buffer for each MPI process</i>  |
| IN recvtype   | <i>Receive data type</i>                                                      |
| IN Root       | <i>Receiving/Gathering MPI process rank</i>                                   |
| IN comm       | <i>Communicator</i>                                                           |

In contrast to MPI\_Gather, each process can send a different number of elements

The individual messages are stored according to *displs* in the receive buffer

## MPI\_Gatherv



```
int MPI_Gatherv (buf1, <2 or 1 or 3>, MPI_INT,
 buf2, rcounts, displs, MPI_INT,
 <rank of P0 in MPI_COMM_WORLD>,
 MPI_COMM_WORLD)
```

## MPI\_Alltoall

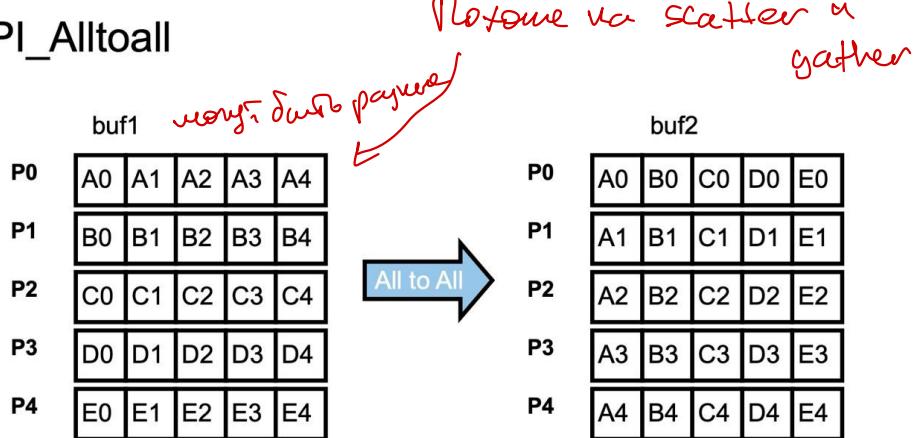
```
int MPI_Alltoall(void* sendbuf, int sendcount, MPI_Datatype sendtype,
 void* recvbuf, int recvcount, MPI_Datatype recvtype,
 MPI_Comm comm)
```

|              |                                                           |
|--------------|-----------------------------------------------------------|
| IN sendbuf   | <i>Send buffer</i>                                        |
| IN sendcount | <i>Number of elements sent</i>                            |
| IN sendtype  | <i>Send type</i>                                          |
| OUT recvbuf  | <i>Receive buffer</i>                                     |
| IN recvcount | <i>Number of elements to be received by a MPI process</i> |
| IN recvtype  | <i>Data type</i>                                          |
| IN comm      | <i>Communicator</i>                                       |

Sends data from each MPI process to each other MPI process

Warning: this is often a scaling bottleneck

## MPI\_Alltoall



```
int MPI_Alltoall (buf1, 1, MPI_INT,
 buf2, 1, MPI_INT,
 MPI_COMM_WORLD)
```

# Nonblocking Collectives (since MPI 3.0)

Same concept as non-blocking P2P

- Defined for all collective operations
  - Include `MPI_Ibarrier`, the non-blocking barrier
  - Enables overlap with extra computation
- Returns request object
- Request can be used in `MPI_Wait` and `MPI_Test` operations

**Cannot** be mixed with blocking collectives

The initiation routines have to be executed in the same order on all processes

Example:

```
int MPI_Ibcast(void* buffer, int count, MPI_Datatype datatype,
 int root, MPI_Comm comm, MPI_Request *request)
```

INOUT buffer starting address of buffer (choice)

IN count number of entries in buffer (non-negative integer)

IN datatype data type of buffer (handle)

IN root rank of MPI process acting as broadcast root (integer)

IN comm communicator (handle)

OUT request communication request (handle)

## MPI\_Ialltoall

```
int MPI_Ialltoall(void* sendbuf, int sendcount, MPI_Datatype
 sendtype,
 void* recvbuf, int recvcount, MPI_Datatype recvtype,
 MPI_Comm comm, MPI_Request *req)
```

|              |                                                    |
|--------------|----------------------------------------------------|
| IN sendbuf   | Send buffer                                        |
| IN sendcount | Number of elements sent                            |
| IN sendtype  | Send type                                          |
| OUT recvbuf  | Receive buffer                                     |
| IN recvcount | Number of elements to be received by a MPI process |
| IN recvtype  | Data type                                          |
| IN comm      | Communicator                                       |
| OUT req      | MPI Request to complete operation                  |

Starts the send of data from each MPI process to each other MPI process  
 Had to be complete later with Test or Wait

Warning: this is still often a scaling bottleneck, therefore a good target for overlap

# Derived MPI Datatypes

Several MPI functions available to create arbitrary layouts

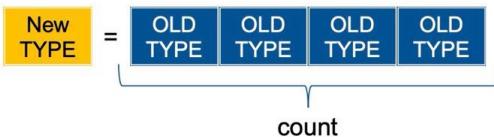
- Recursive specification possible
- Declarative specification of data-layout

Constructor functions

- Additional arguments to describe new type
- Input: existing datatype (basetype for this operation)
- Output: newly created derived datatype

```
int MPI_Type_contiguous(int count, MPI_Datatype oldtype,
 MPI_Datatype *newtype)
```

New datatype is “count” contiguous elements of “oldtype”



Datatype handle

```
MPI_Datatype
```

Datatype construction

Set of routines (see previous and next slides)

Datatype commit

```
int MPI_Type_commit(MPI_Datatype *datatype)
```

necessary before type can be used,

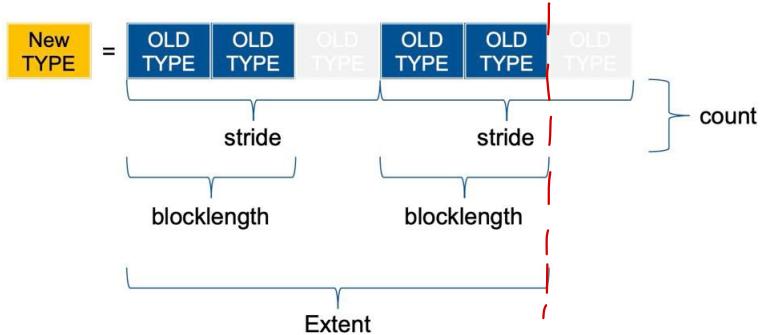
but uncommitted types can be used to create new types

Datatype free

```
int MPI_Type_free(MPI_Datatype *datatype)
```

## Constructing a Vector Datatype

```
int MPI_Type_vector(int count, int blocklength, int stride,
 MPI_Datatype oldtype, MPI_Datatype *newtype)
```

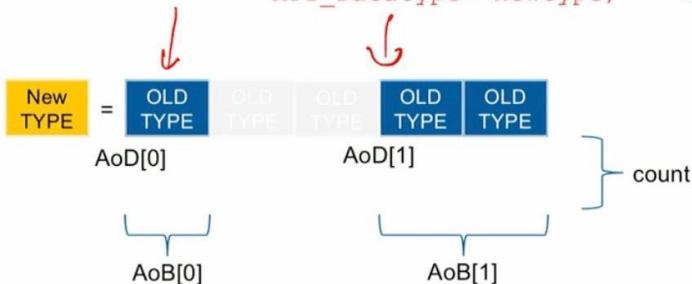


## Constructing an Indexed Datatype

```
int MPI_Type_indexed(int count,
 const int array_of_blocklengths[],
 const int array_of_displacements[],
 MPI_Datatype oldtype,
 MPI_Datatype *newtype)
```

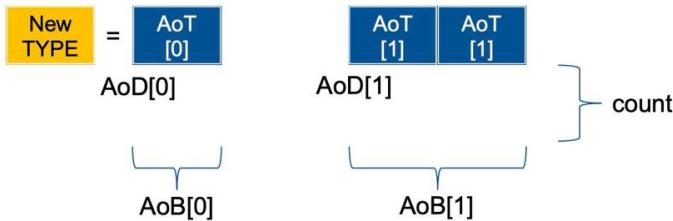
Annotations in red:

- $(1, 1)$  above the first element of the new type array.
- $(0, 3)$  below the second element of the new type array.
- $(AoS)$  to the right of the first element.
- $(AoD)$  to the right of the second element.



# Creating a Struct Datatype

```
int MPI_Type_create_struct(int count,
 const int array_of_blocklengths[], (AoB)
 const MPI_Aint array_of_displacements[], (AoD)
 const MPI_Datatype array_of_types[], (AoT)
 MPI_Datatype *newtype)
```



Note: Displacements in Bytes (MPI\_Aint) *addresses integers*

## Example



```
#include <stdio.h>
#include <stdlib.h>
#include <mpi.h>

#define INSPECT(dt) {MPI_Type_get_extent(dt,&lb,&ex); printf("%li\n",ex);}

int main(int argc, char **argv)
{
 MPI_Datatype dt1[2],dt2,dt3;
 int blens[2];
 MPI_Aint ex,lb,disps[2];

 MPI_Init(&argc,&argv);

 dt1[0]=MPI_DOUBLE; INSPECT(dt1[0]);
 dt1[1]=MPI_CHAR; INSPECT(dt1[1]);
 blens[0]=1; blens[1]=2;
 disps[0]=0; disps[1]=sizeof(double);
 MPI_Type_create_struct(2,blens,disps,dt1,&dt2);
 INSPECT(dt2);

 MPI_Type_vector(2,1,2,dt2,&dt3);
 INSPECT(dt3);

 MPI_Finalize();
}
```

Query extent

16 Double C C

48 Double C C

# MPI Threading

New MPI initialization routine

```
int MPI_Init_thread(int *argc, char ***argv,
 int required, int *provided)
```

Application states needed level of thread support  
MPI returns which one it supports

1) call\_init\_thread

3) finalize

2) check

MPI\_THREAD\_SINGLE - oguz u wce, go konus

- Only one thread exists in the application

MPI\_THREAD\_FUNNELED - outside parallel regions, during master region

- Multithreaded, but only the main thread makes MPI calls master region

MPI\_THREAD\_SERIALIZED - critical & single regions

- Multithreaded, but only one thread at a time makes MPI calls

MPI\_THREAD\_MULTIPLE

- Multithreaded and any thread can make MPI calls at any time

## Consequences of using MPI\_THREAD\_MULTIPLE

Each call to MPI completes on its own

- Effect is as if they would have been called sequentially in some order
- Blocking calls will only block the calling thread

T1: MPI\_Send( to: 1 )  
T2: MPI\_Recv( from: 1 )

T1: MPI\_Send( to: 0 )  
T2: MPI\_Recv( from: 0 )



User is responsible for making sure racing calls are avoided

- Example: access to an already freed object
- Collective operations must be ordered correctly among threads

Collective operations  
should be called in the  
same order!

T1: MPI\_Bcast( comm )  
T2: MPI\_Barrier( comm )

T1: MPI\_Bcast( comm )  
T2: MPI\_Barrier( comm )  
order is not determined



Possible performance impact, e.g., caused by locking in MPI

# MPI - SPMD (single process, multiple data)

## MPI Persistent Communication

### — setup communication

```
int MPI_Send_init(const void *buf, int count, MPI_Datatype datatype,
 int dest, int tag, MPI_Comm comm, MPI_Request *request);
```

int MPI\_Start(MPI\_Request \*request); — every loop iteration

int MPI\_Request\_free(MPI\_Request \*request); — after loop

Also available for collectives

when send to

Sender: fill partition and mark as ready

Receiver: notification for each arriving partition

Partitions do not need to be finished in order!

Each thread fills one or more partitions.

### Sender:

1. MPI\_Psend\_init sets up the communication.
2. MPI\_Start indicates the coming transfers to MPI.
3. MPI\_Pready indicates that a partition is ready. — filled one partition
4. MPI\_Wait completes the operation.
5. MPI\_Request\_free frees the resources.

## Receiver:

1. `MPI_Precv_init` sets up the communication.
2. `MPI_Start` indicates the coming transfers to MPI.
3. `MPI_Parrived` indicates whether a partition arrived.
4. `MPI_Wait` completes the operation.
5. `MPI_Request_free` frees the resources.

true or  
false

```
if (rank == 0) {
 MPI_Psend_init(message, sendbuf, arraysize, MPI_DOUBLE, 1, 0, MPI_COMM_WORLD, MPI_INFO_NULL, &request);

 for (int k = 0; k < num_steps; k++) {
 MPI_Start(&request);

 #pragma omp parallel
 {
 int i = omp_get_thread_num();
 // Do computation on i-th block of data
 MPI_Pready(i, &request);
 }

 MPI_Wait(&request);
 }

 MPI_Request_free(&request);
} else if (rank == 1) {
```