# Introduction to Scalasca

Gabriele Fatigati - g.fatigati@cineca.it
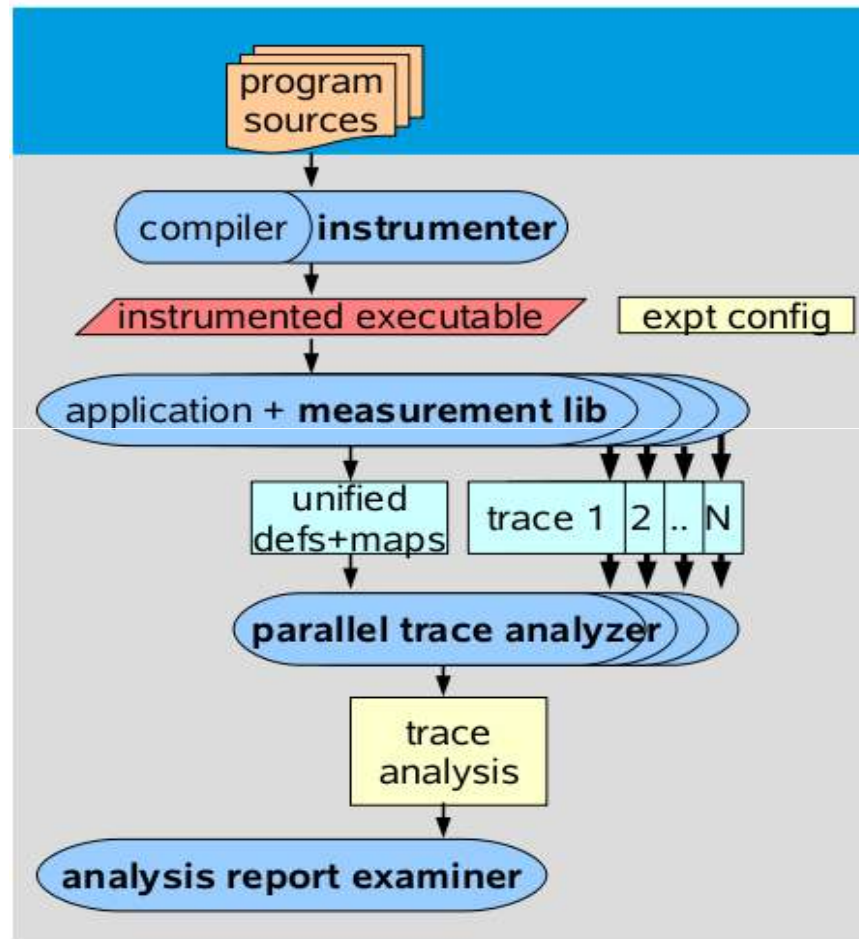
SuperComputing Group

CINECA

- SCalable performance Analysis of LArge SCale Applications
- Developed by Julich Supercomputer Centre
- Toolset for performance analysis of parallel applications on a large scale
- Manage programs MPI, OpenMP, MPI+OpenMP
- Latest releast 1.3
- www.scalasca.org

Event tracing

During the
measurement there is
a buffer for each
thread/process

Final collect of the
results

## Compilation

Original command:                    SCALASCA instrumentation command:

mpcc -c foo.c                        scalasca -instrument mpcc -c foo.c
mpxlf90  -o bar bar.f90               skin mpxlf90  -o bar bar.f90

```
#!/bin/bash
#
# @ job_name    = myjob
# @ output      = myjob.$(jobid)
# @ error       = myjob.$(jobid)
# @ wall_clock_limit = 0:10:00
# @ total_tasks = 8
# @ task_affinity=core(1)
# @ parallel_threads=1
# @ job_type    = parallel
# @ resources   = ConsumableMemory(320Mb)
# @ queue

module load profile/advanced
module load qt/4.5.2--xl--10.1
module load scalasca/1.2

scalasca -analyze poe ./c_example
```
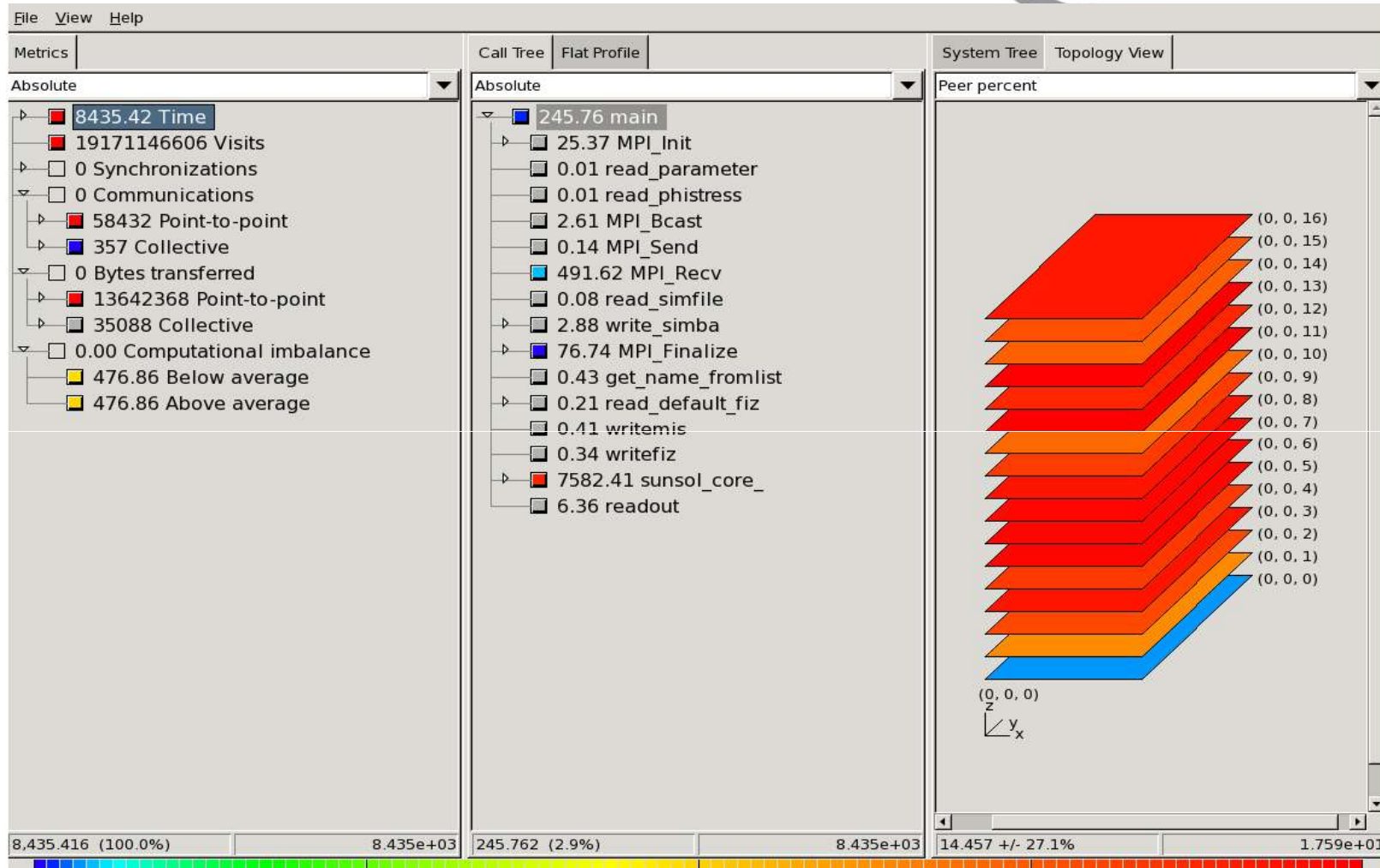
Results analysis:
scalasca -examine epik_...

**Log**

[00000]EPIK: Created new measurement archive ./epik_a
[00000]EPIK: Activated ./epik_a
 SWEEP3D – Pipelined Wavefront with Line-Recursion
 32 domains – 4 x 8 decomposition
 Iteration Monitor:
  its=1 err=1.000000 fixs=0

  ...
  its=12 err=5320.611978 fixs=19706584
 Balance quantities:

 ...
[00000]EPIK: Closing experiment ./epik_a
[00000]EPIK: 42 unique paths
[00000]EPIK: Unifying...done
[00000]EPIK: Collating...done
[00000]EPIK: Closed experiment ./epik_a

21st Summer
School of
**PARALLEL
COMPUTING**



Topology controls toolbar
(enable via 'Topology' menu)

What kind of
performance problem?

Where is it in the
source code?
In what context?

How is it distributed
across the system?
(graphical or tree-based view)

Select different
display modes

Colour coding according
to severity value and
display mode

**Context menus** via
right mouse button

Hierarchy minimum
(selected mode/absolute)

Selected value
(selected mode/absolute/
percentage of hierarchy total)

Hierarchy total
(selected mode/absolute)

CINECA

# Topology view

- Hardware (only on some systems, like Blue Gene)
- MPI topology ( eg: MPI_Cart_Create)
- Visual topology user-defined (next releases)

Currently supports cartesian topologies 1D,2D,3D
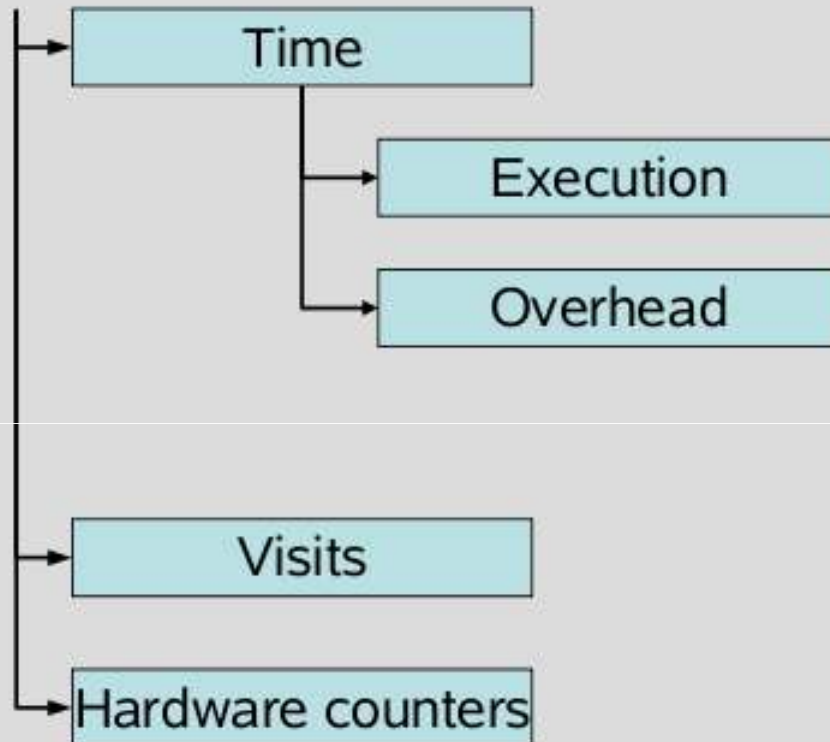
## Display modes

**Absolute**
Absolute value in seconds/number of occurrences

**Root Percent**
Percentage relative to the root ot the hierachy

**External percent**
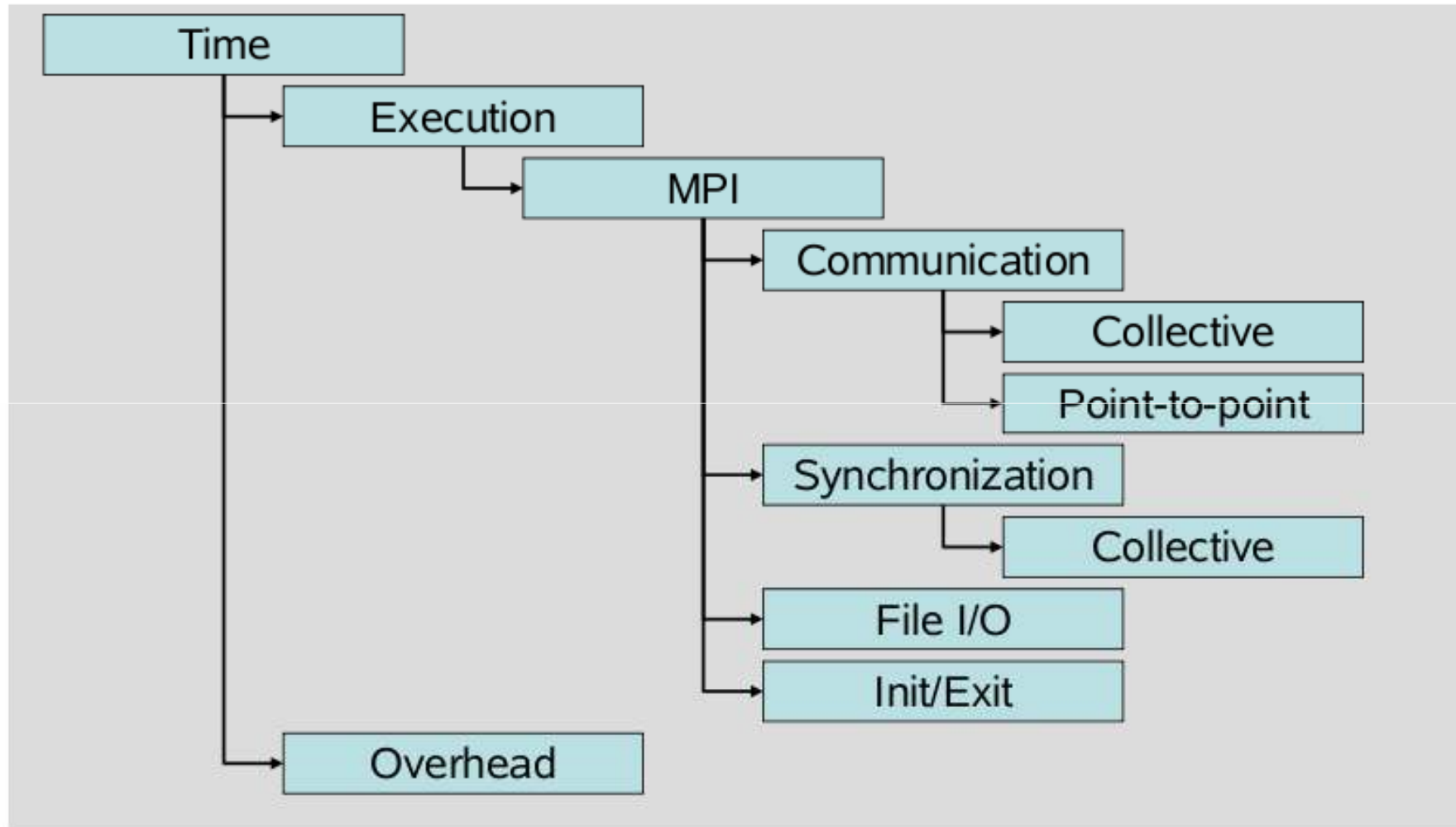Similar to "root percent", but for a different dataset

| | |
|---|---|
| **Time** | Total CPU allocation time |
| **Execution** | Execution time without overhead |
| **Overhead** | Time spent in activities related to measurement (not including dilation per instrumented routine/region!) |
| **Visits** | Number of times a routine/region was executed |
| **Hardware counters** | Aggregated counter values for each routine/region |

21st Summer
School of
**PARALLEL**
**COMPUTING**

| | |
|---|---|
| Time | Total CPU allocation time |
| Execution | Execution time without overhead |
| Overhead | Time spent in tasks related to measurement (not including dilation from instrumentation!) |
| MPI | Time spent in pre-instrumented MPI functions |
| Communication | Time spent in MPI communication calls, subdivided into collective and point-to-point |
| Synchronization | Time spent in calls to `MPI_Barrier()` |
| File I/O | Time spent in MPI file I/O functions |
| Init/Exit | Time spent in `MPI_Init()` and `MPI_Finalize()` |

scalasca

21st Summer
School of
**PARALLEL**
**COMPUTING**

Communications

→ Collective
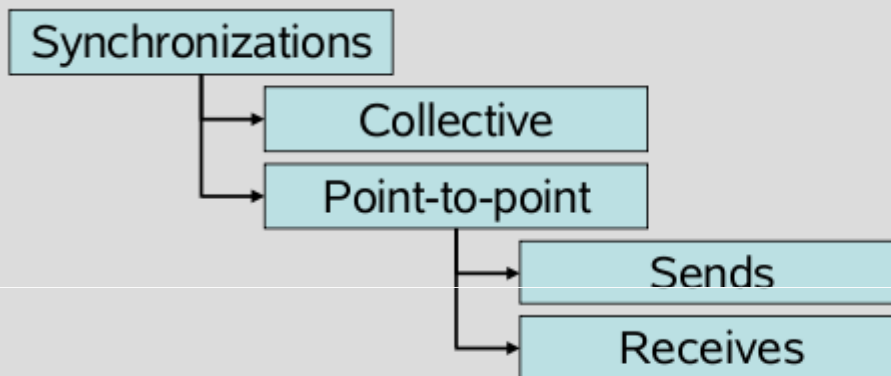
→ Exchange

→ As Source

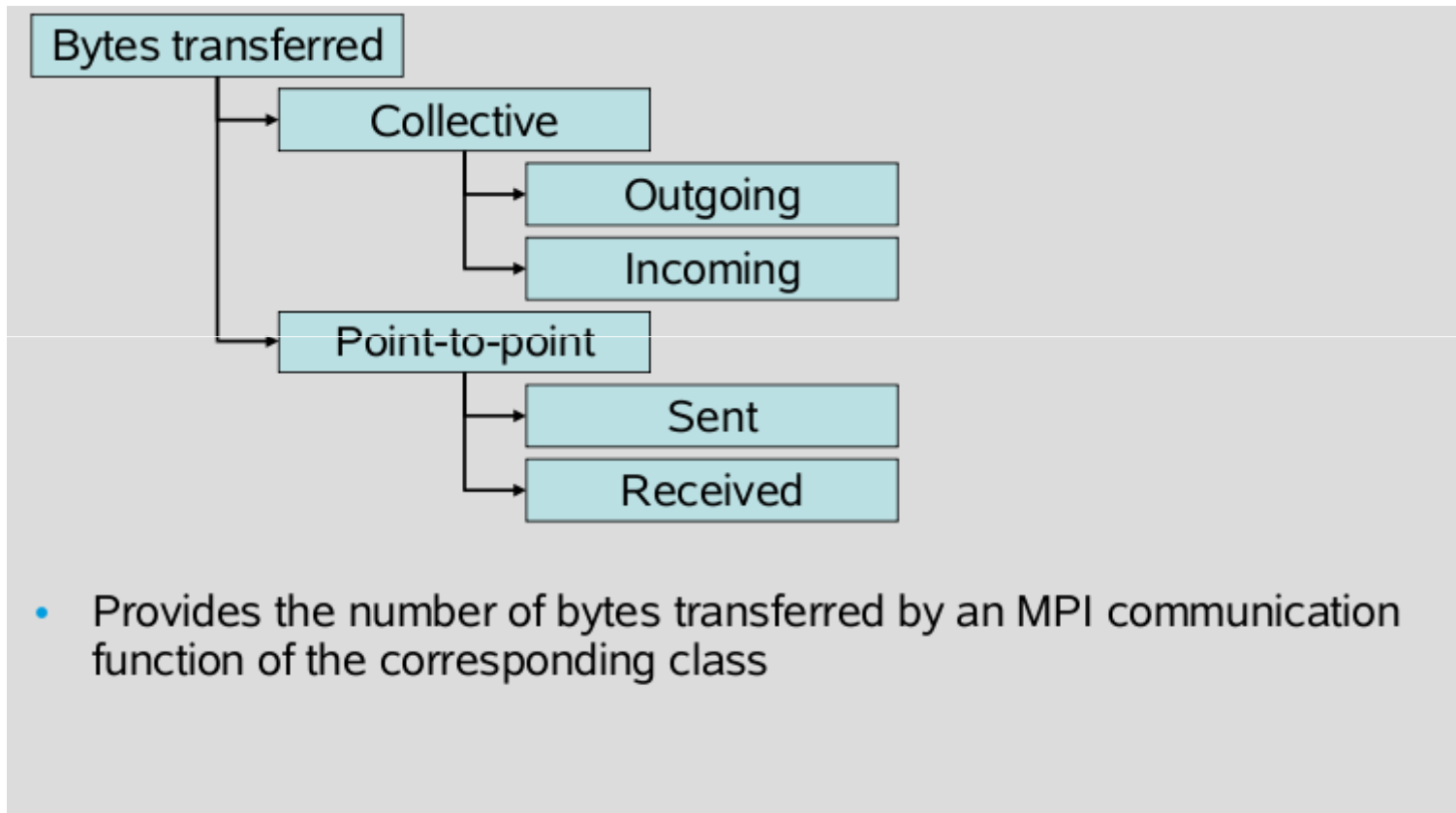→ As Destination

Point-to-point

→ Sends

→ Receives

- Provides the number of calls to an MPI communication function of the corresponding class
- Zero-sized message transfers are considered *synchronization*!
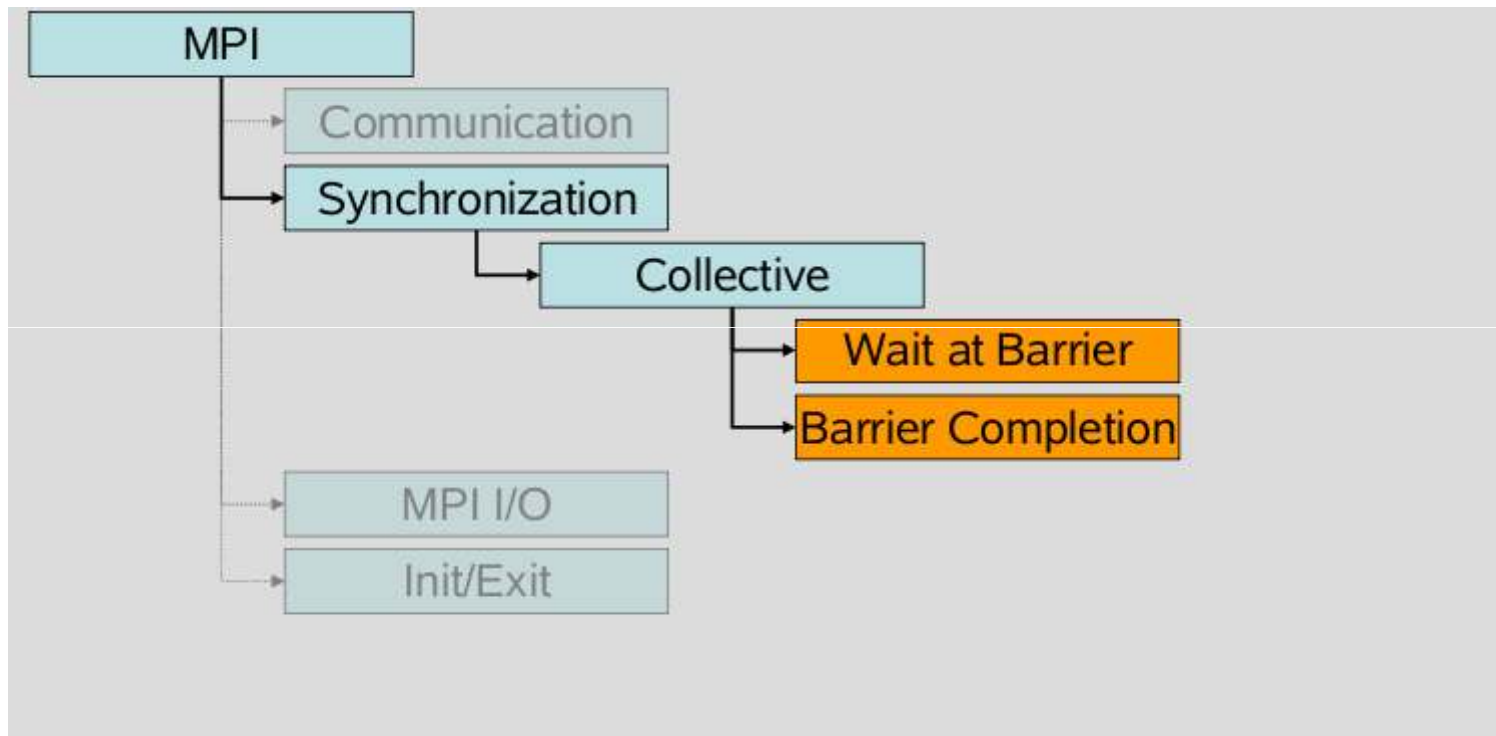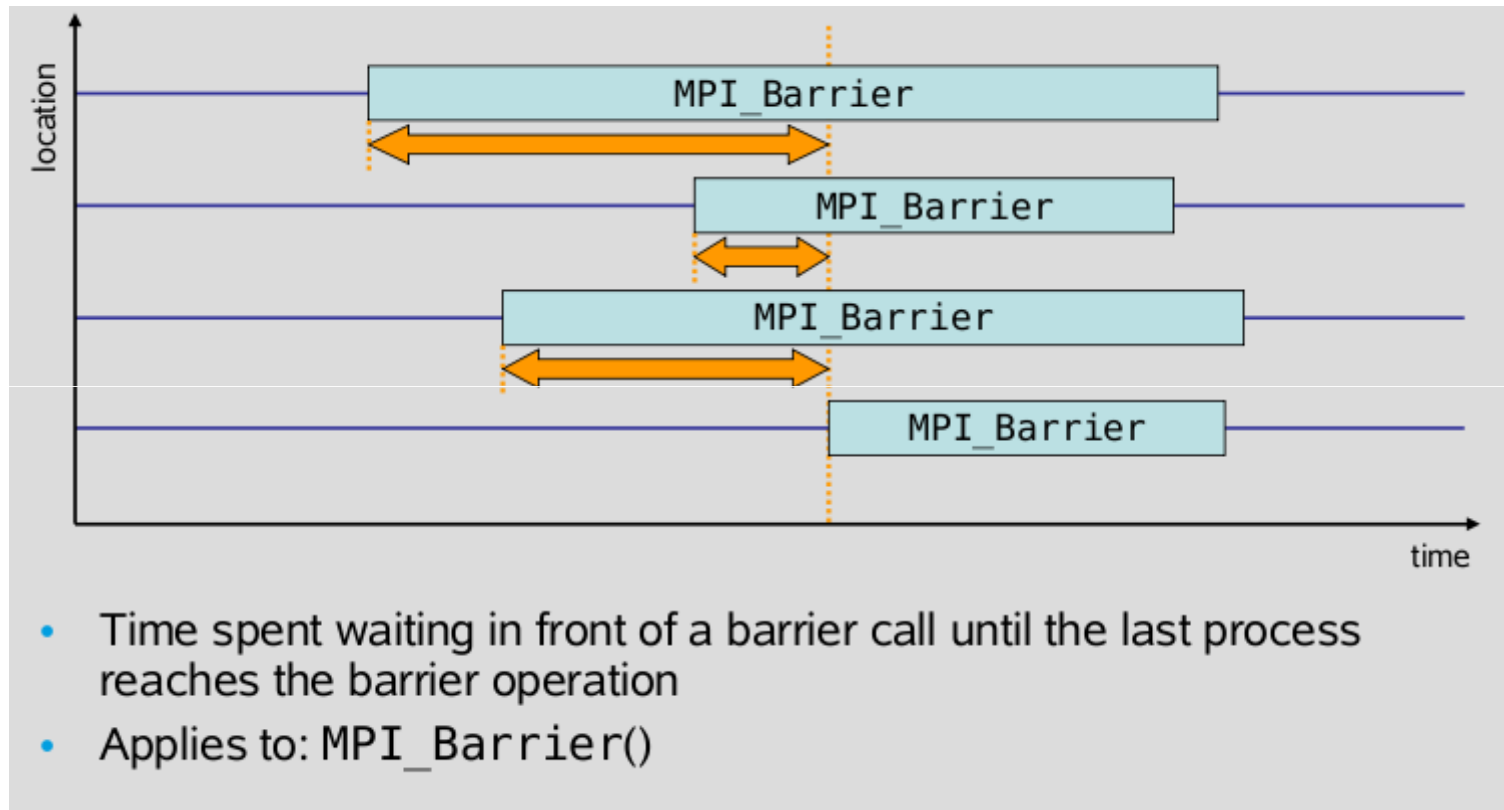
scalasca

- Provides the number of calls to an MPI synchronization function of the corresponding class
- MPI synchronizations include zero-sized message transfers!

# 21st Summer School of PARALLEL COMPUTING



```
Bytes transferred
        → Collective
                → Outgoing
                → Incoming
        → Point-to-point
                → Sent
                → Received
```

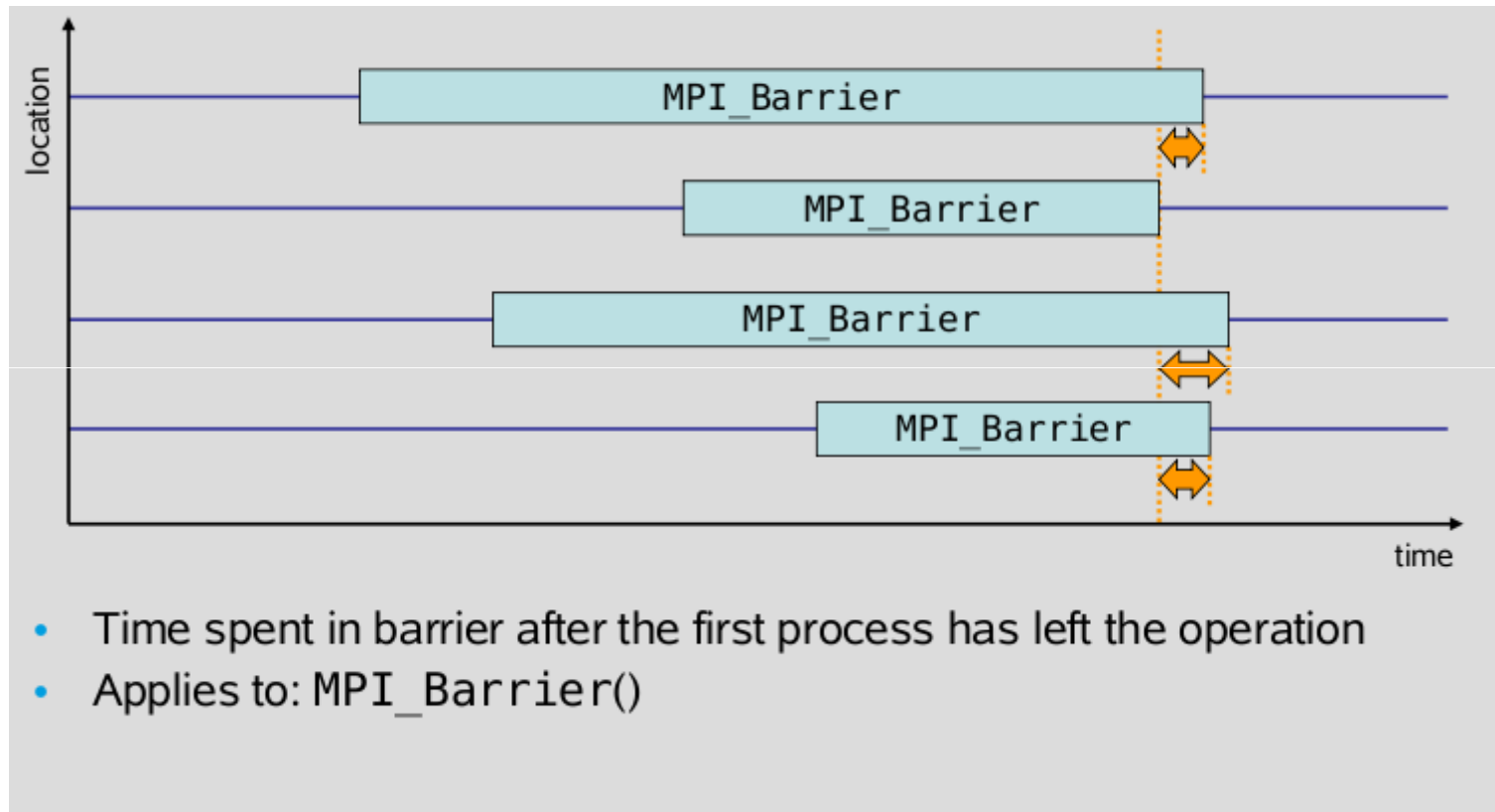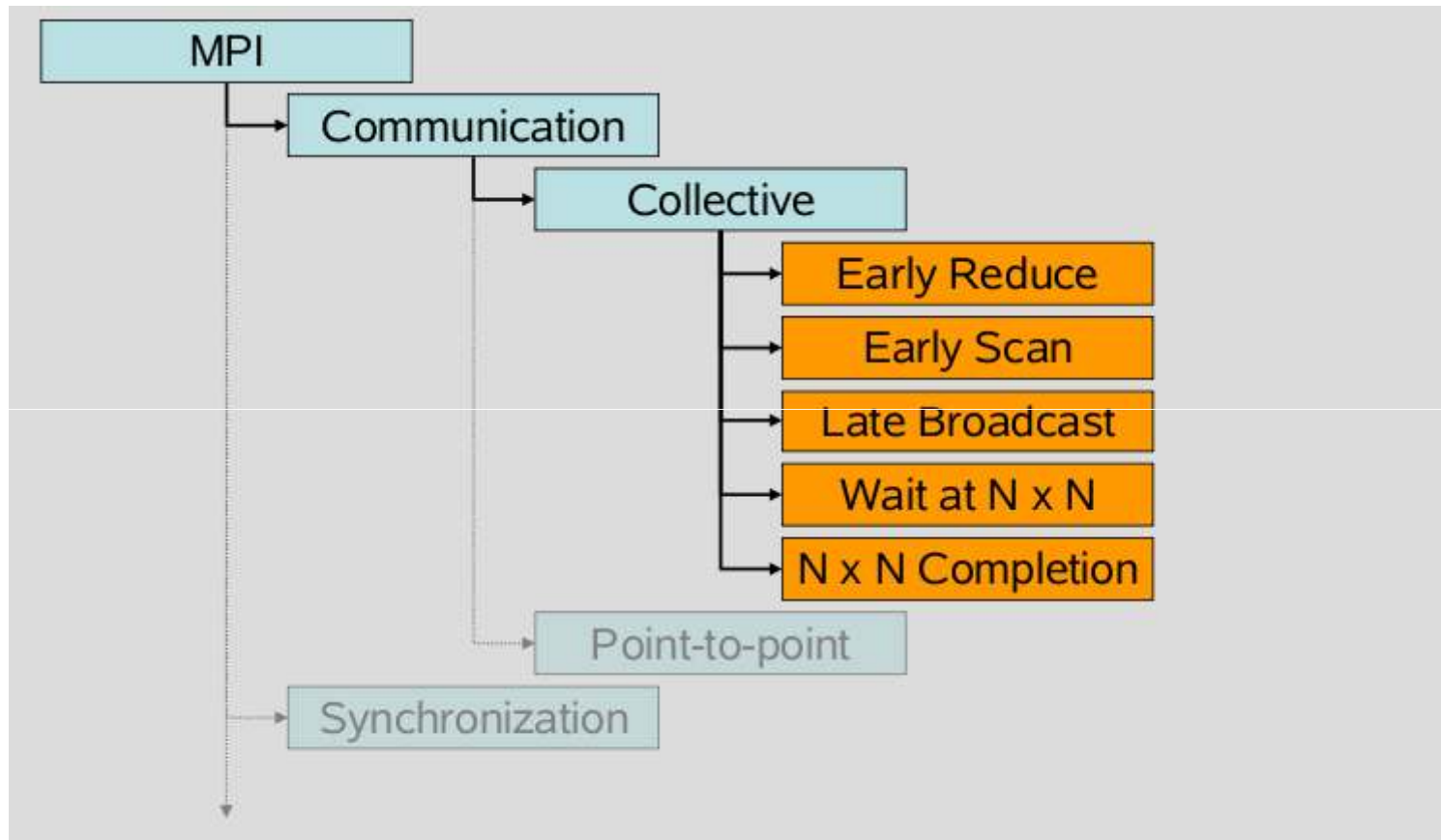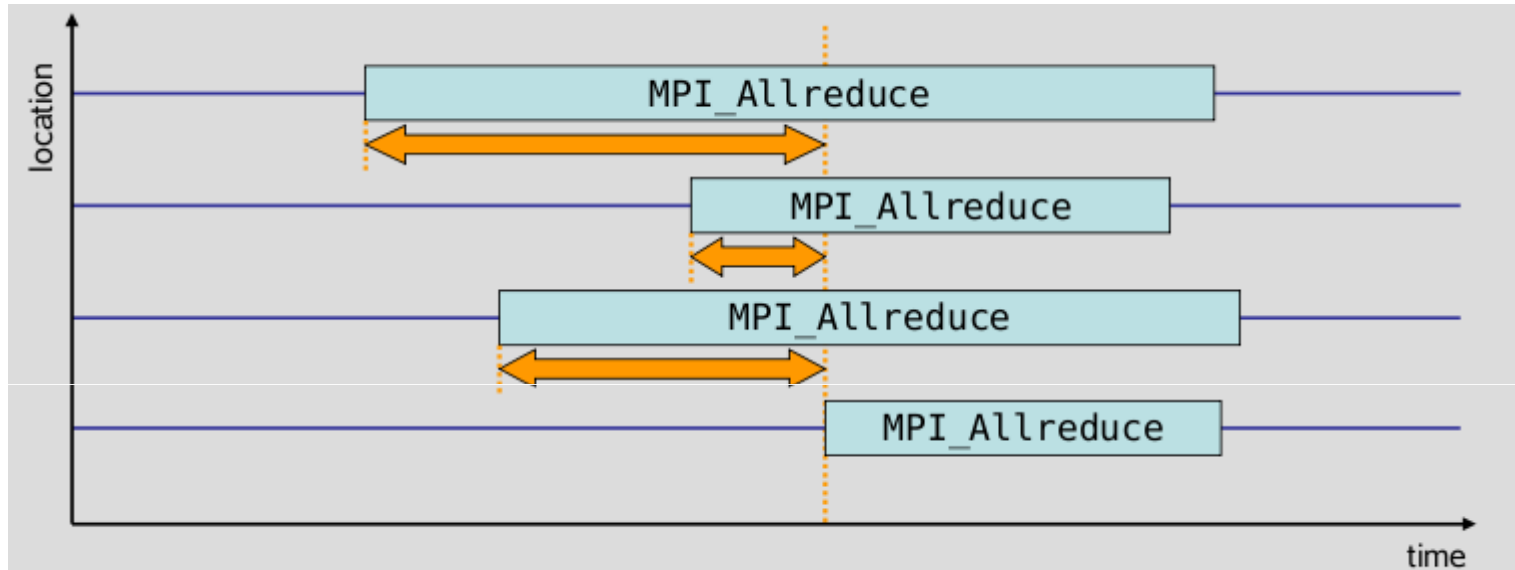- Provides the number of bytes transferred by an MPI communication function of the corresponding class

# MPI collective synchronization time

- Time spent waiting in front of a barrier call until the last process reaches the barrier operation
- Applies to: `MPI_Barrier()`

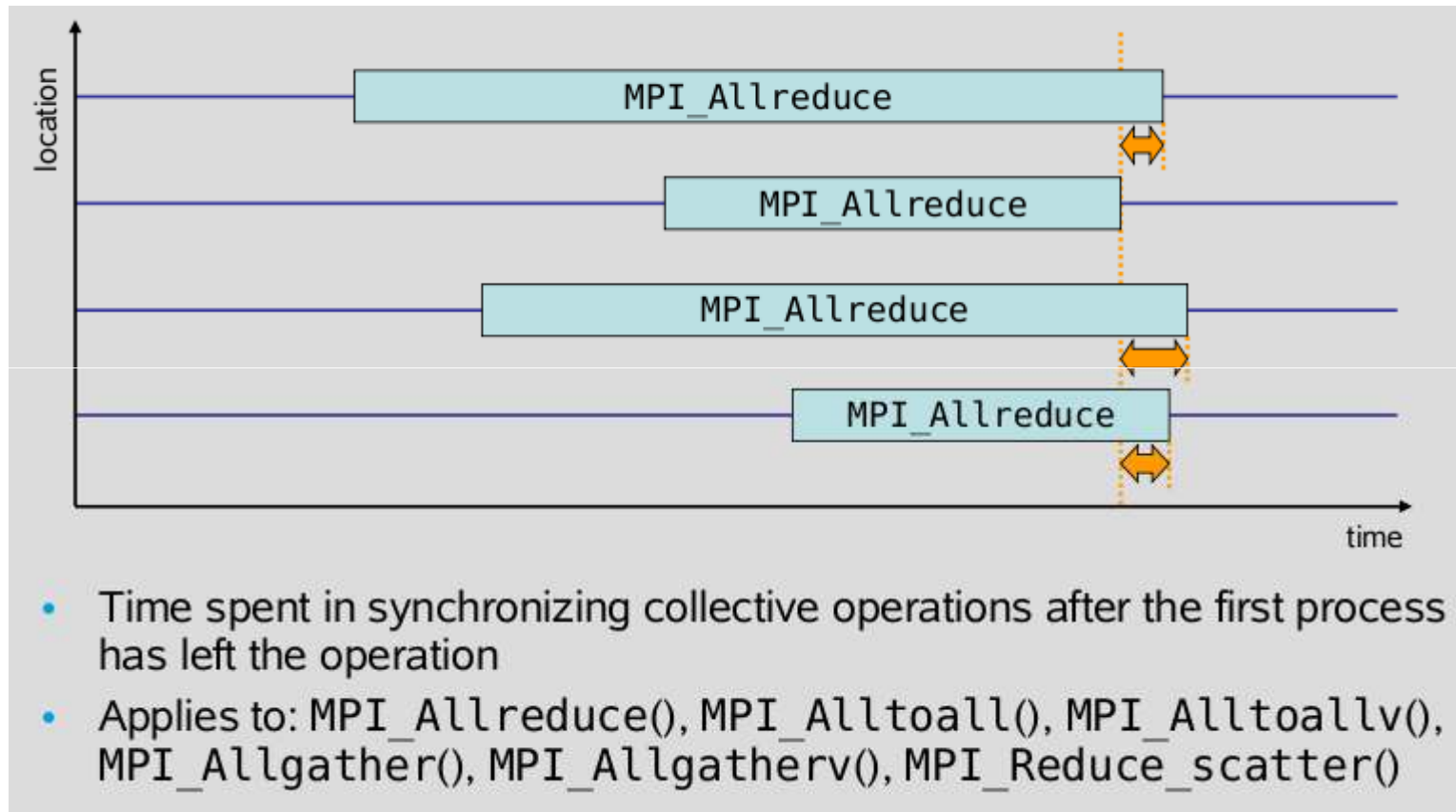- Time spent in barrier after the first process has left the operation
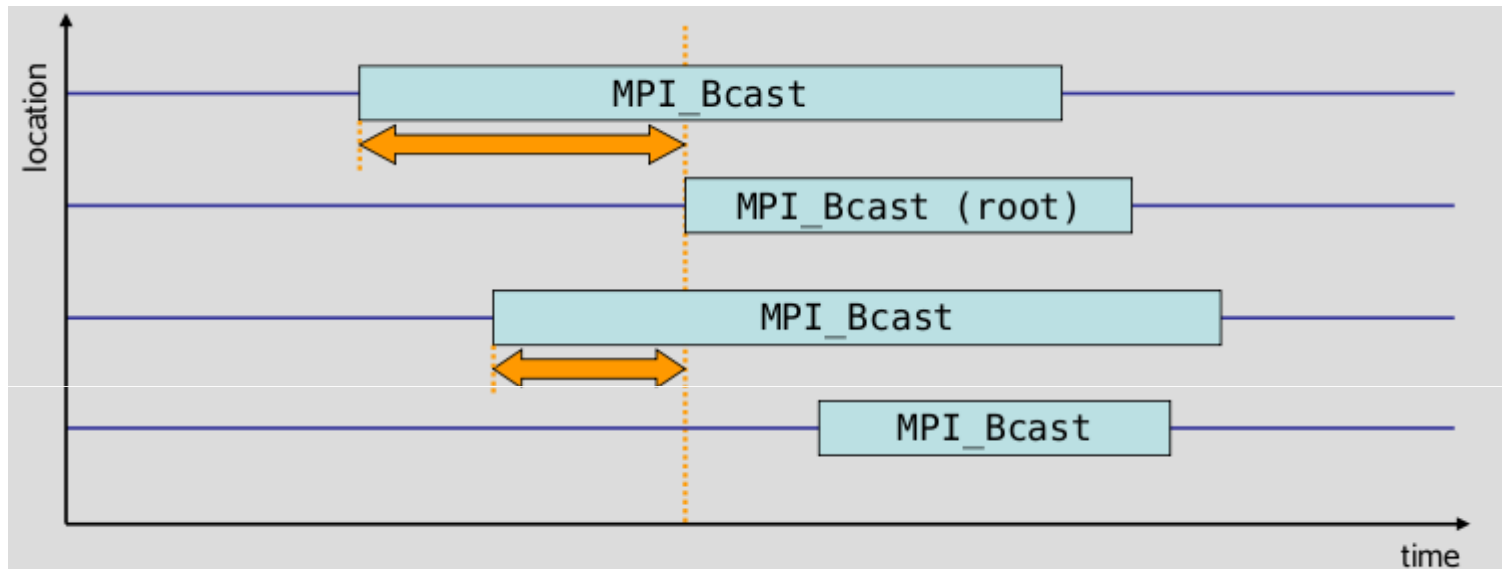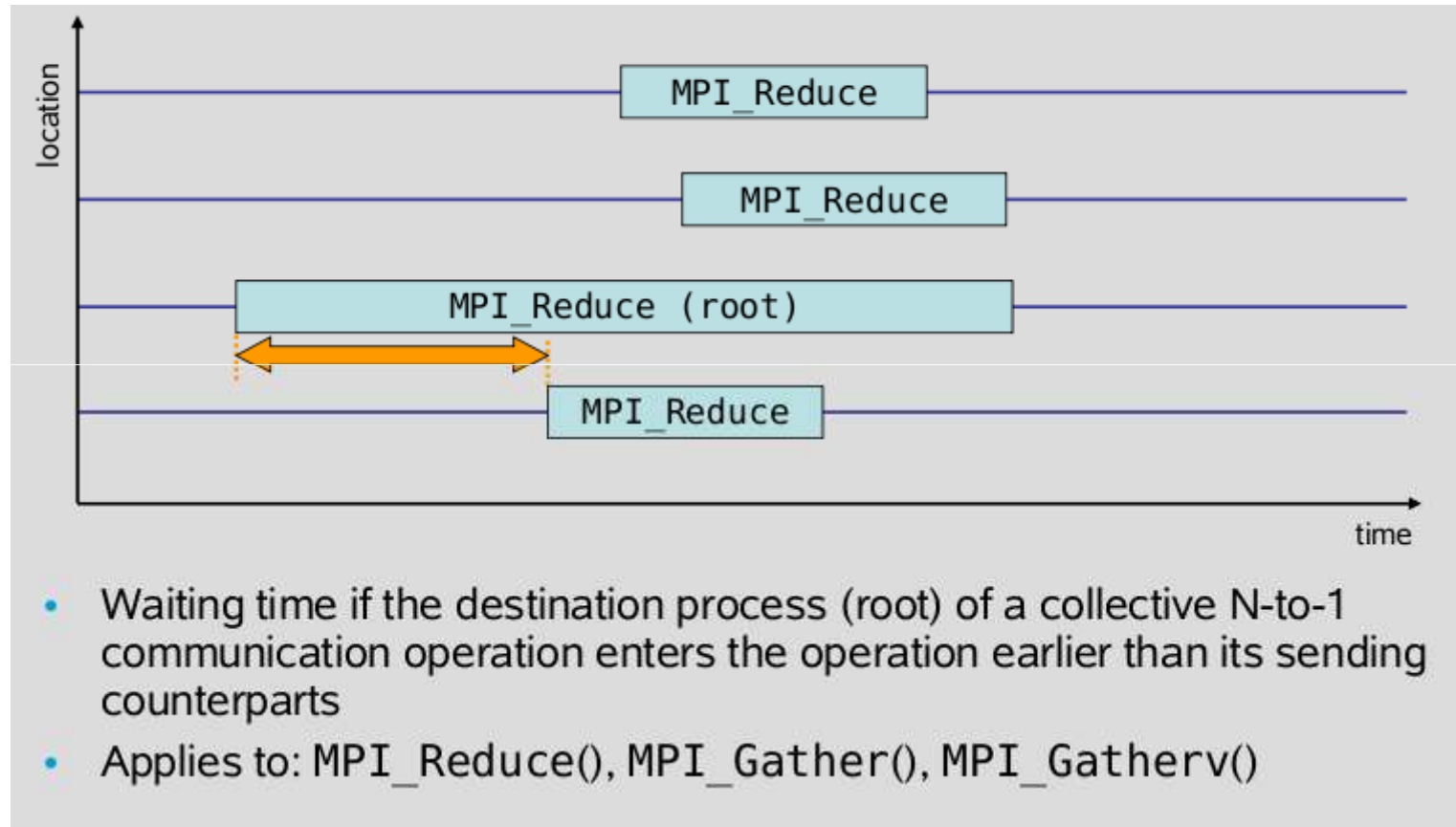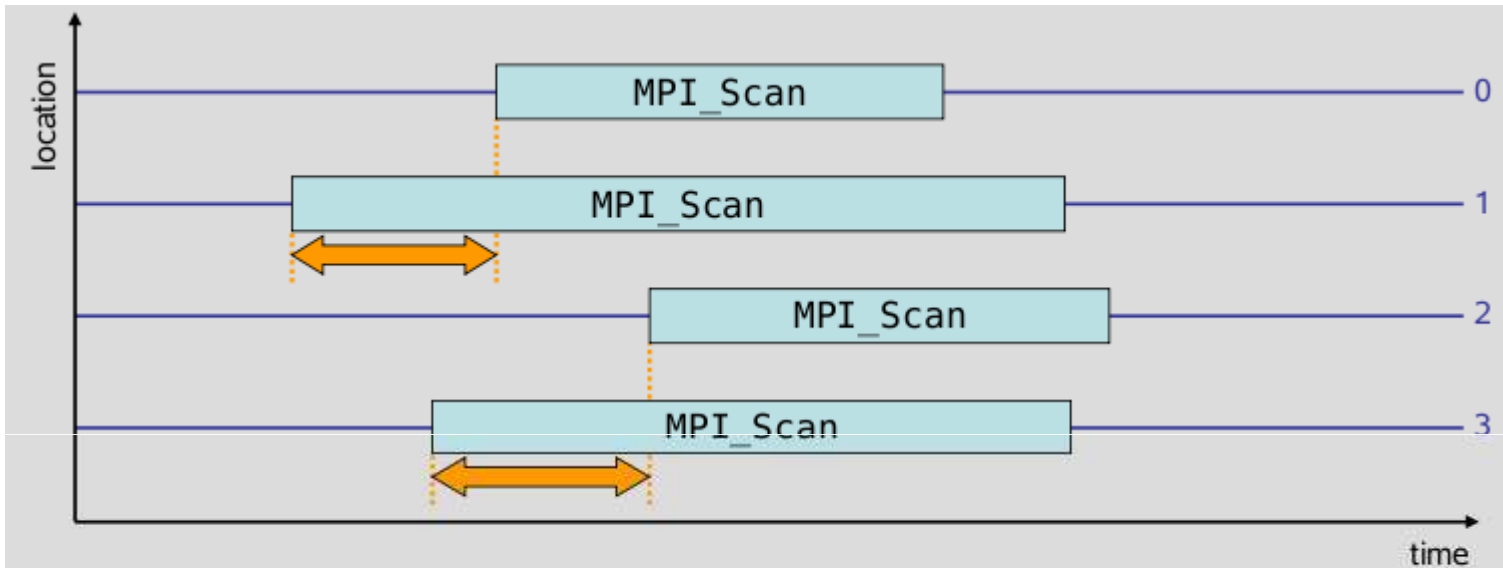- Applies to: `MPI_Barrier()`

- Time spent waiting in front of a synchronizing collective operation call until the last process reaches the operation
- Applies to: MPI_Allreduce(), MPI_Alltoall(), MPI_Alltoallv(), MPI_Allgather(), MPI_Allgatherv(), MPI_Reduce_scatter()

- Time spent in synchronizing collective operations after the first process has left the operation
- Applies to: `MPI_Allreduce()`, `MPI_Alltoall()`, `MPI_Alltoallv()`, `MPI_Allgather()`, `MPI_Allgatherv()`, `MPI_Reduce_scatter()`
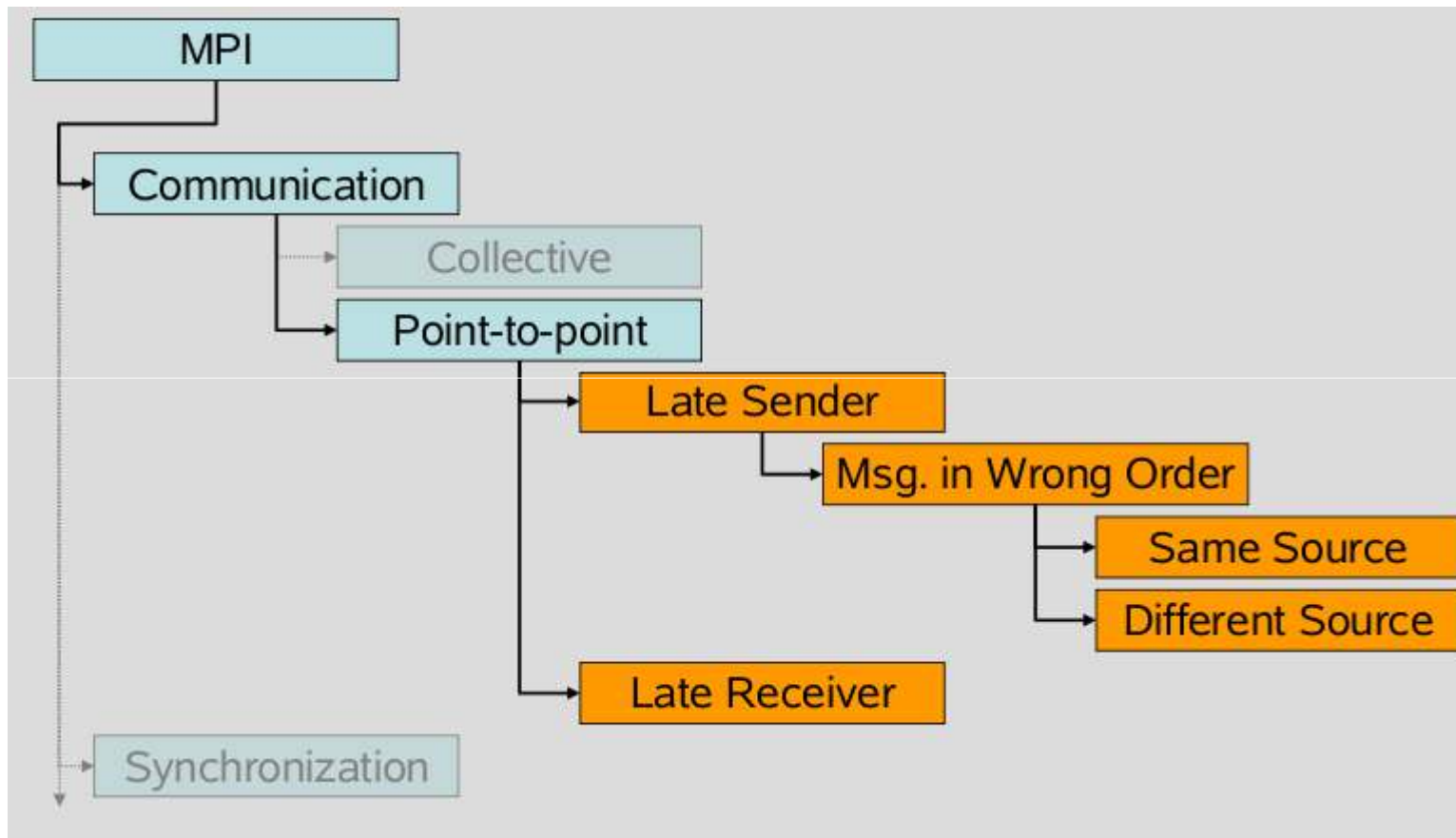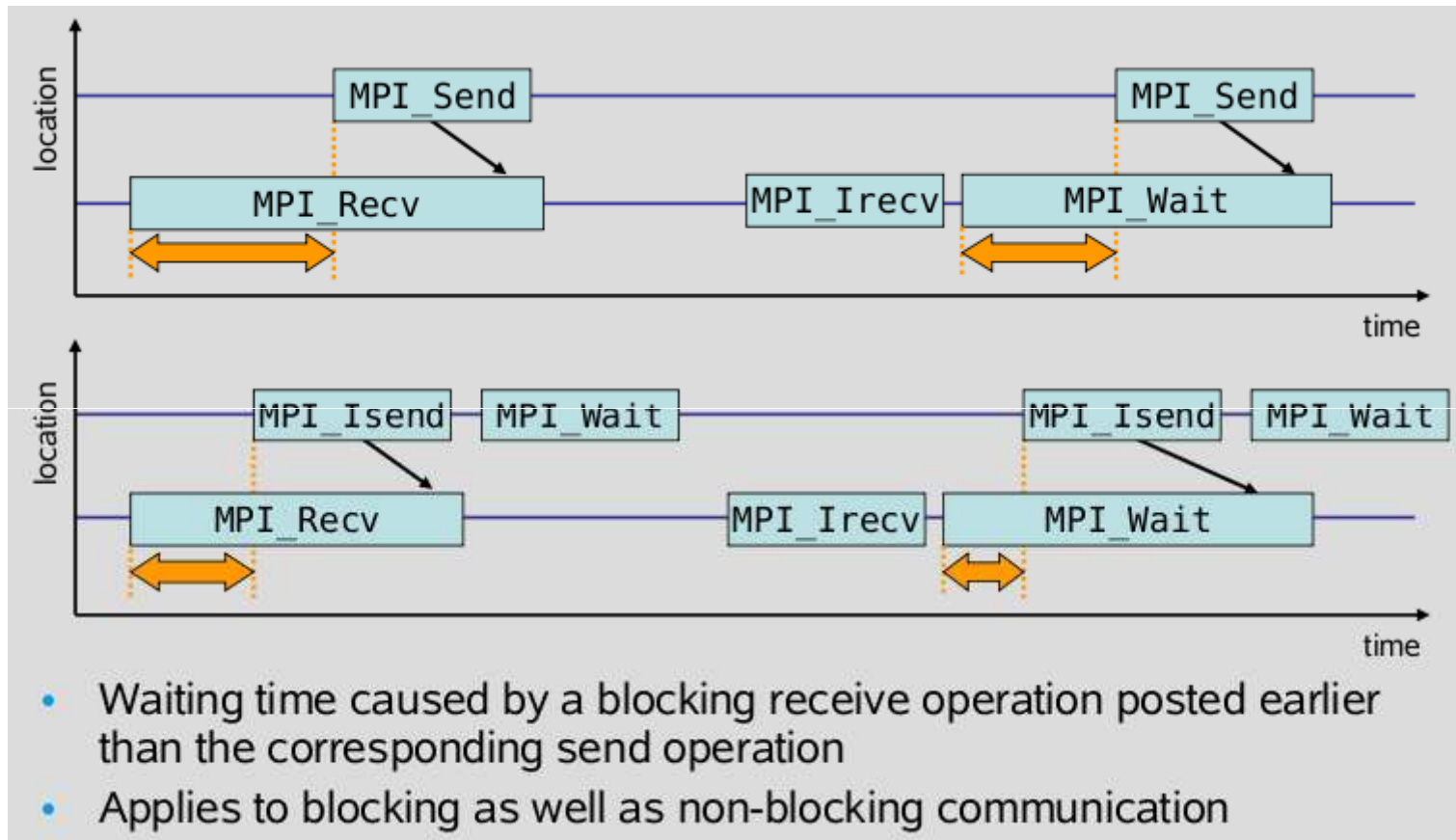
- Waiting times of the destination processes of a collective 1-to-N communication operation which enter the operation earlier than the source process (root)
  - Late Broadcast by source = Early Broadcast by destinations
- Applies to: MPI Bcast(), MPI Scatter(), MPI Scatterv()

- Waiting time if the destination process (root) of a collective N-to-1 communication operation enters the operation earlier than its sending counterparts
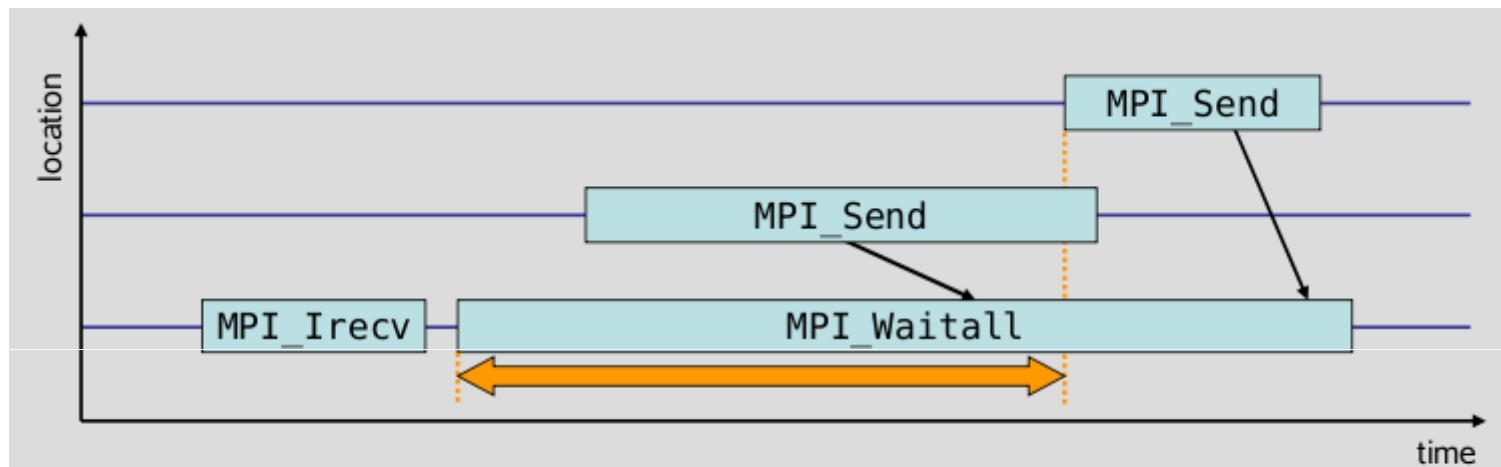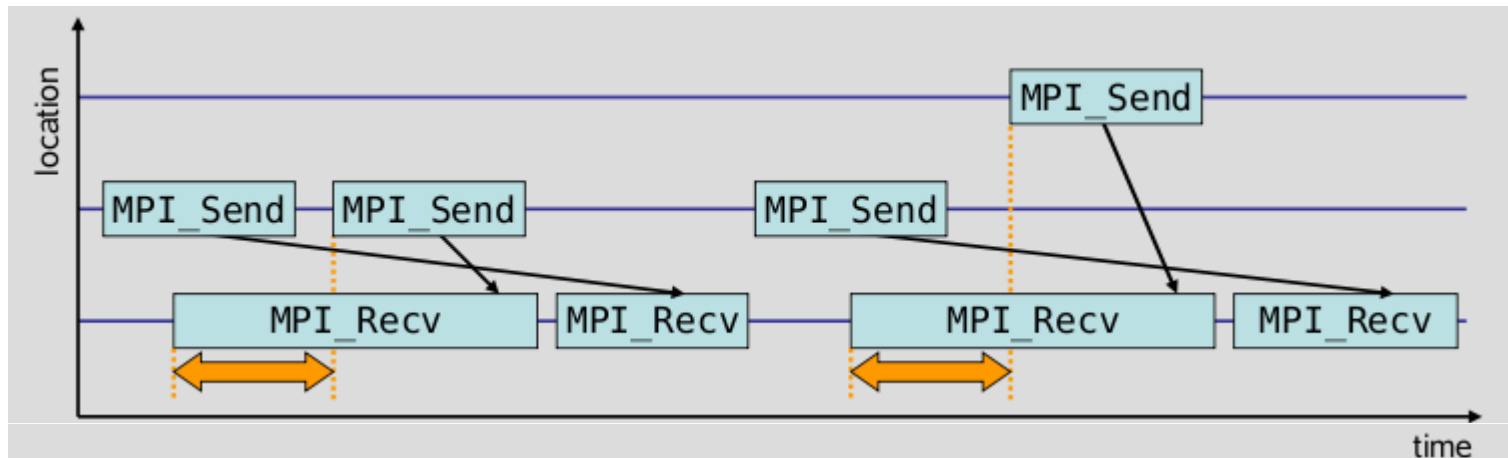- Applies to: `MPI_Reduce()`, `MPI_Gather()`, `MPI_Gatherv()`

- Waiting time if process $n$ enters a prefix reduction operation earlier than its sending counterparts (i.e., ranks $0..n-1$)
- Applies to: `MPI_Scan()`

- Waiting time caused by a blocking receive operation posted earlier than the corresponding send operation
- Applies to blocking as well as non-blocking communication

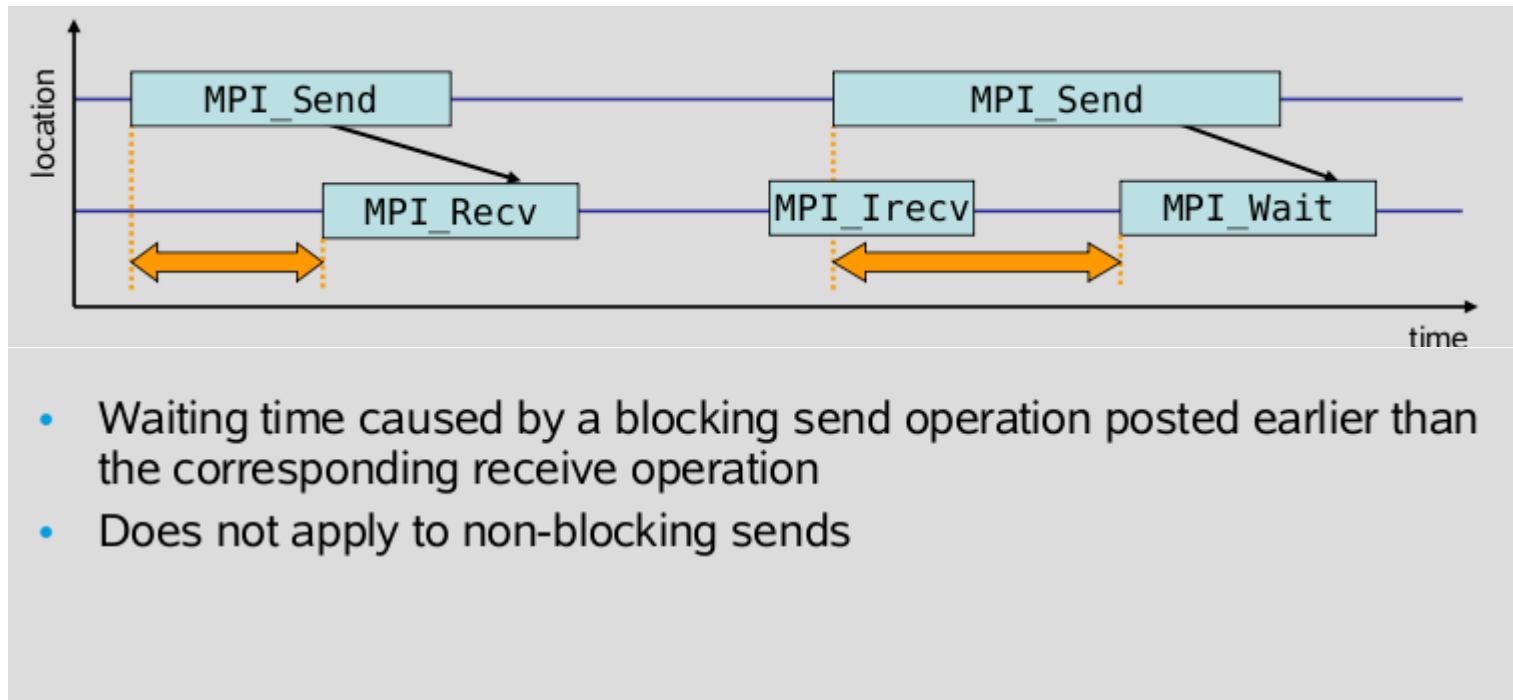21st Summer School of **PARALLEL COMPUTING**
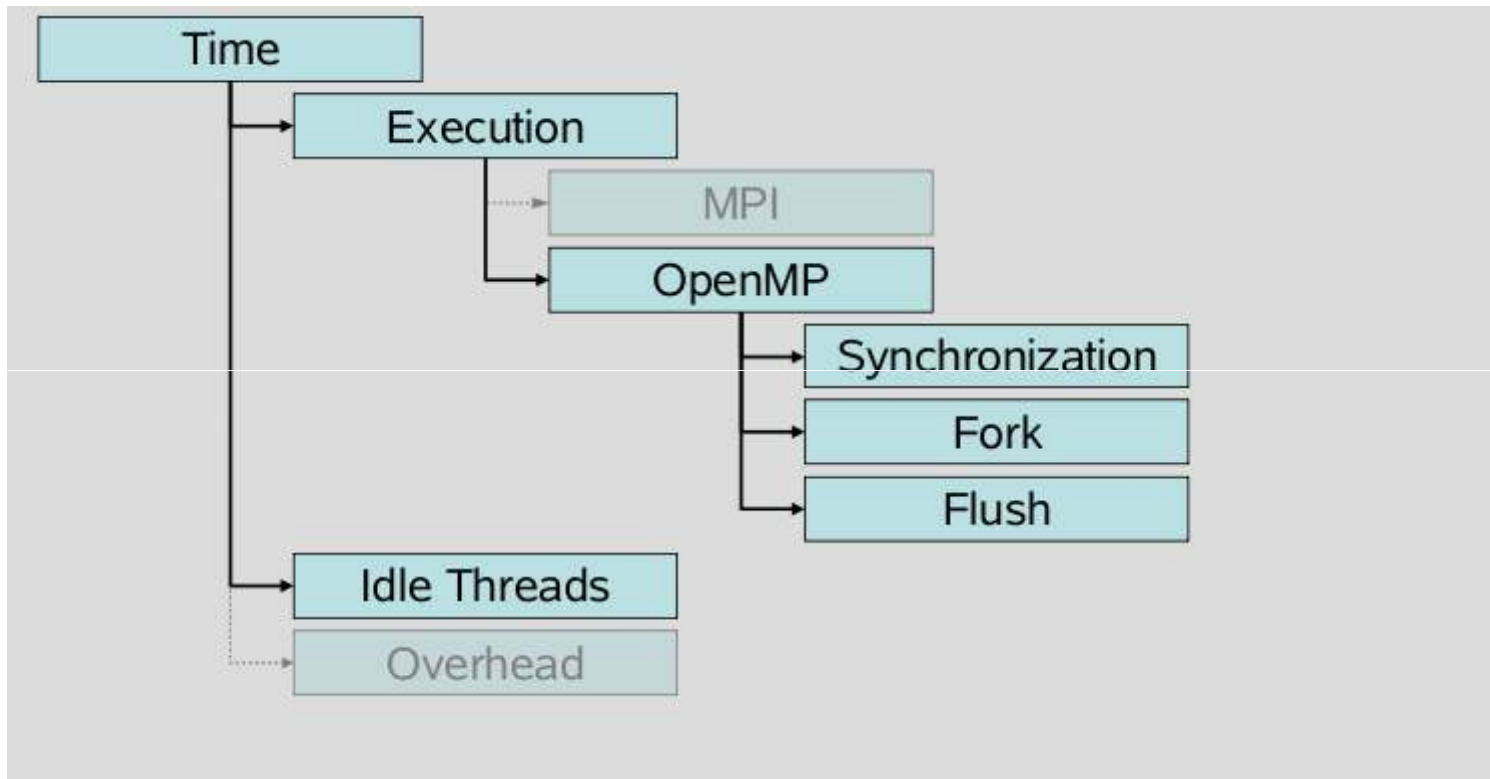


- While waiting for several messages, the maximum waiting time is accounted
- Applies to: `MPI_Waitall()`, `MPI_Waitsome()`

- Refers to Late Sender situations which are caused by messages received in wrong order
  - Early receive of message out of order
- Comes in two flavours:
  - Messages sent from same source location
  - Messages sent from different source locations

- Waiting time caused by a blocking send operation posted earlier than the corresponding receive operation
- Does not apply to non-blocking sends

21st Summer School of **PARALLEL COMPUTING**

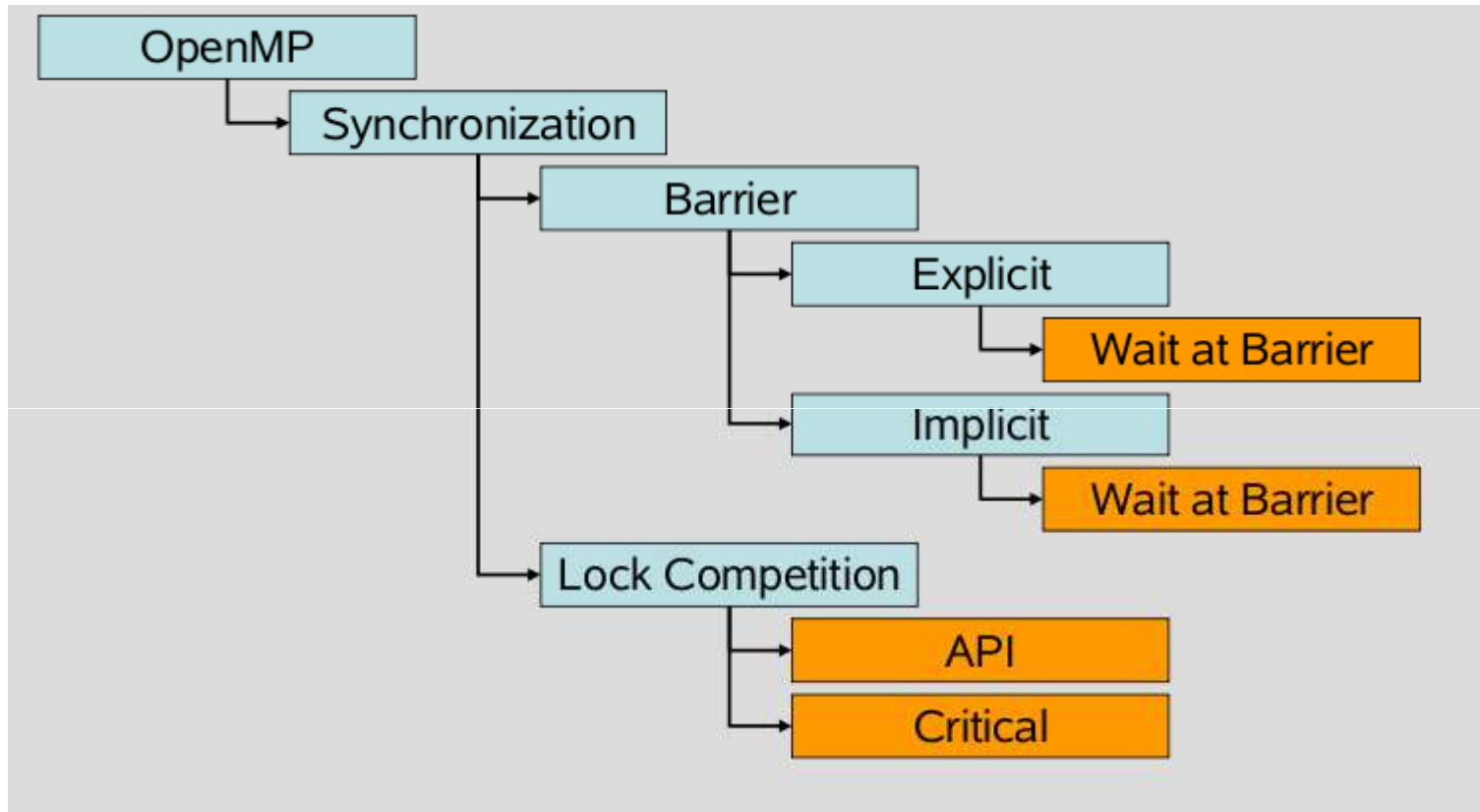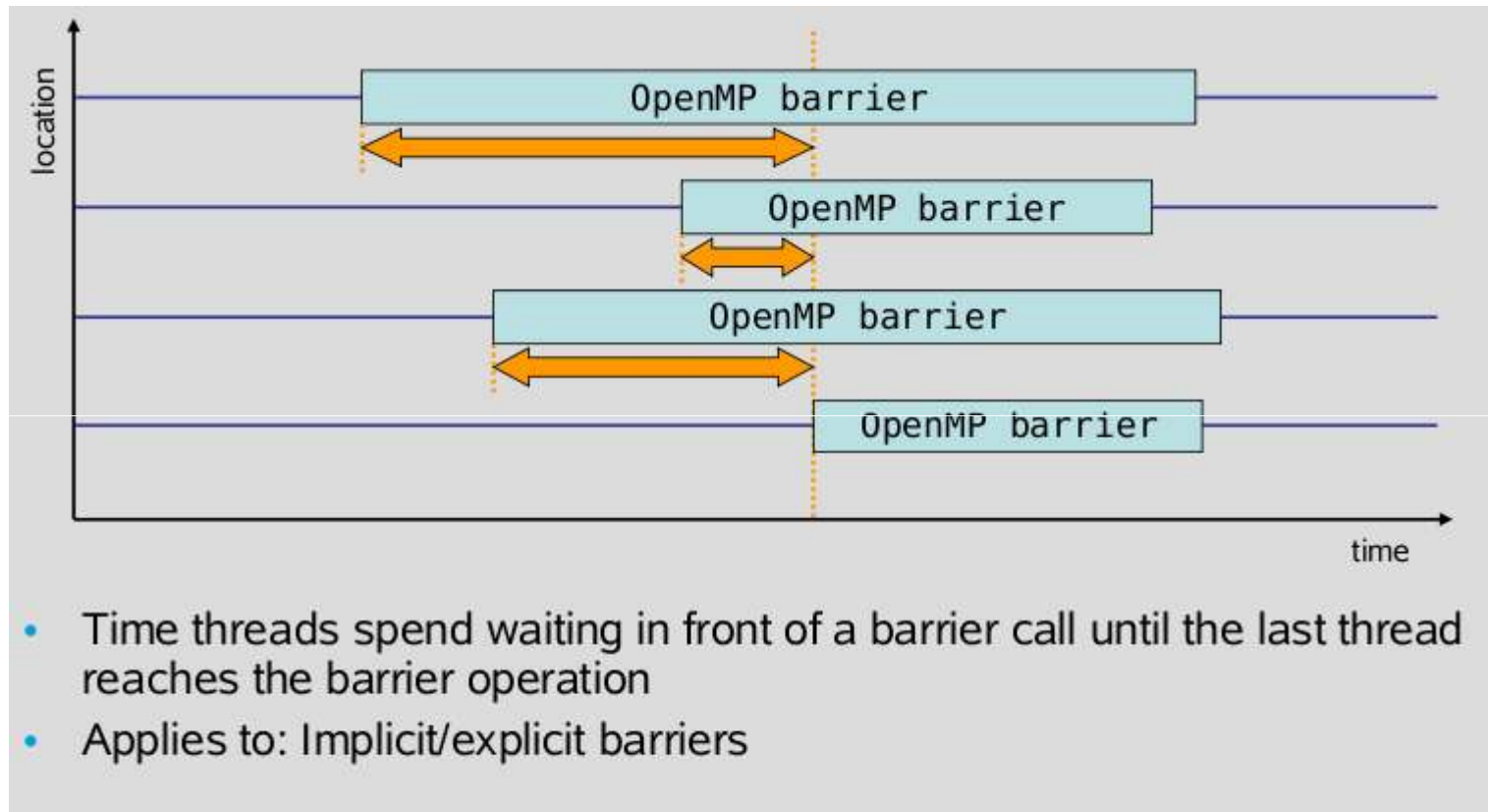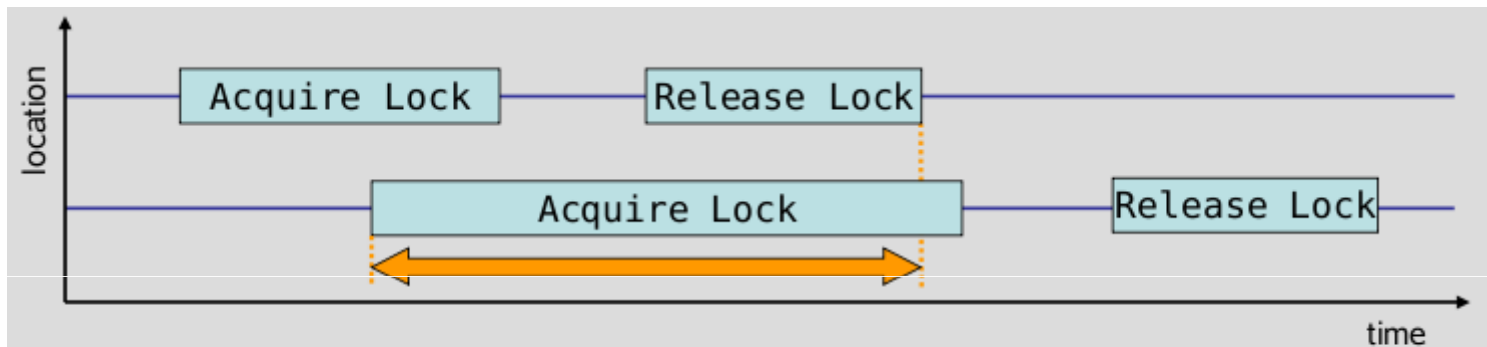| | |
|---|---|
| OpenMP | Time spent for all OpenMP-related tasks |
| Synchronization | Time spent synchronizing OpenMP threads |
| Fork | Time spent by master thread to create thread teams |
| Flush | Time spent in OpenMP flush directives |
| Idle Threads | Time spent idle on CPUs reserved for slave threads |

- Time threads spend waiting in front of a barrier call until the last thread reaches the barrier operation
- Applies to: Implicit/explicit barriers

- Time a thread spends waiting for a lock that is held by other threads until it is released and can be acquired by this thread
- Applies to: critical sections, OpenMP lock API

## Code instrumentation

C/C++:

```
#include "epik_user.h"
...
void foo() {
   ...  // local declarations
   ...  // more declarations
   EPIK_FUNC_START();
   ...  // executable statements
   if (...)  {
     EPIK_FUNC_END();
     return;
   } else {
     EPIK_USER_REG(r_name,"region");
     EPIK_USER_START(r_name);

     ...
     EPIK_USER_END(r_name);
   }
   ...  // executable statements
   EPIK_FUNC_END();
   return;
}
```

Fortran:

```
#include "epik_user.inc"
...
subroutine bar()
  EPIK_FUNC_REG("bar")
   ...  ! local declarations
  EPIK_FUNC_START()
   ...  ! executable statements
  if (...)  then
    EPIK_FUNC_END()
    return
  else
    EPIK_USER_REG(r_name,"region")
    EPIK_USER_START(r_name)

    ...
    EPIK_USER_END(r_name)
  endif
   ...  ! executable statements
  EPIK_FUNC_END()
  return
end subroutine bar
```

C++:

```
#include "epik_user.h"
...
{
  EPIK_TRACER("name");
  ...
}
```

- EPIK_FUNC_START, EPIK_FUNC_END mark the entry and exit from the piece of code

- The regions should be initialized with EPIK_USER_REG

- Each exit/break/continue/return must have EPIK_FUNC_END

- Need **-user** flag to decode instrumentations