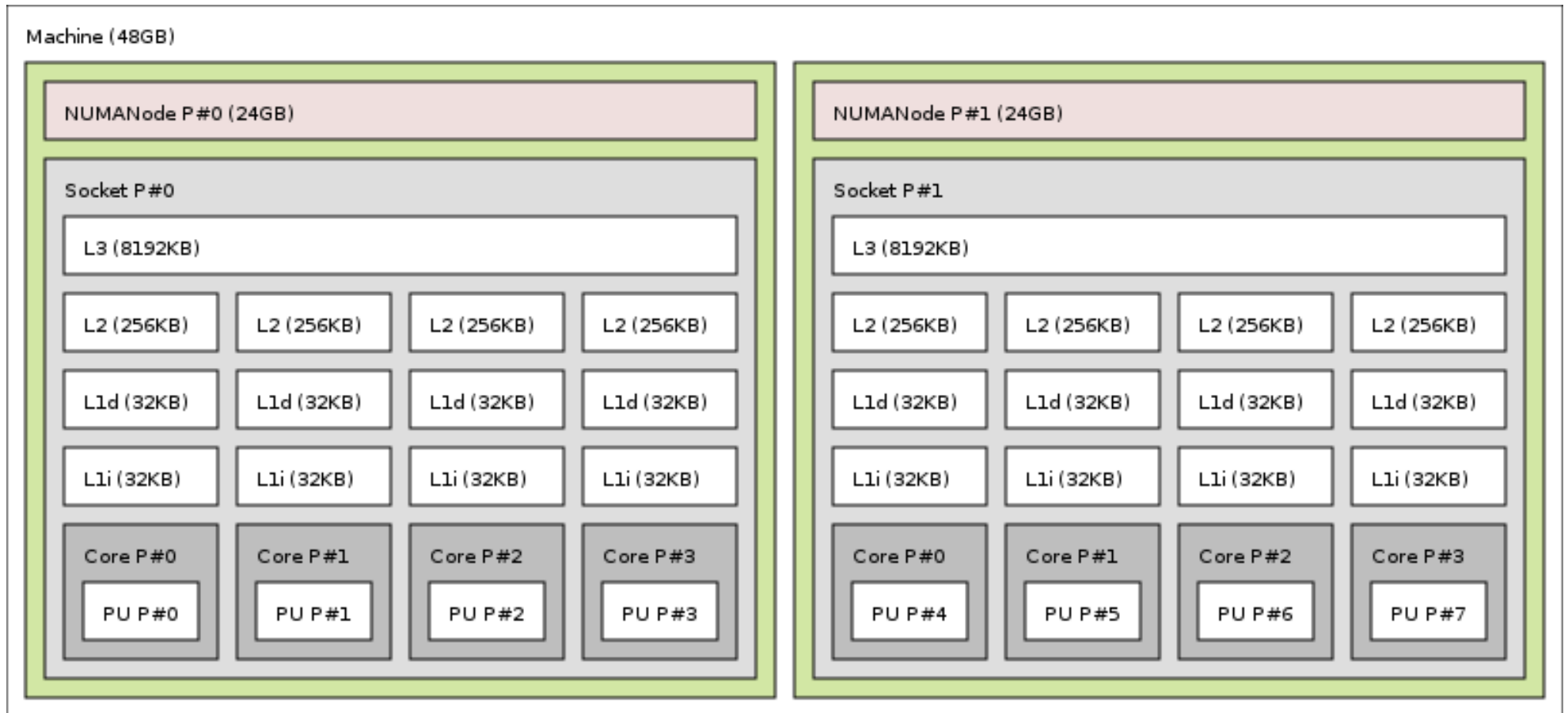# Locality and Physical Topology Support in MPI

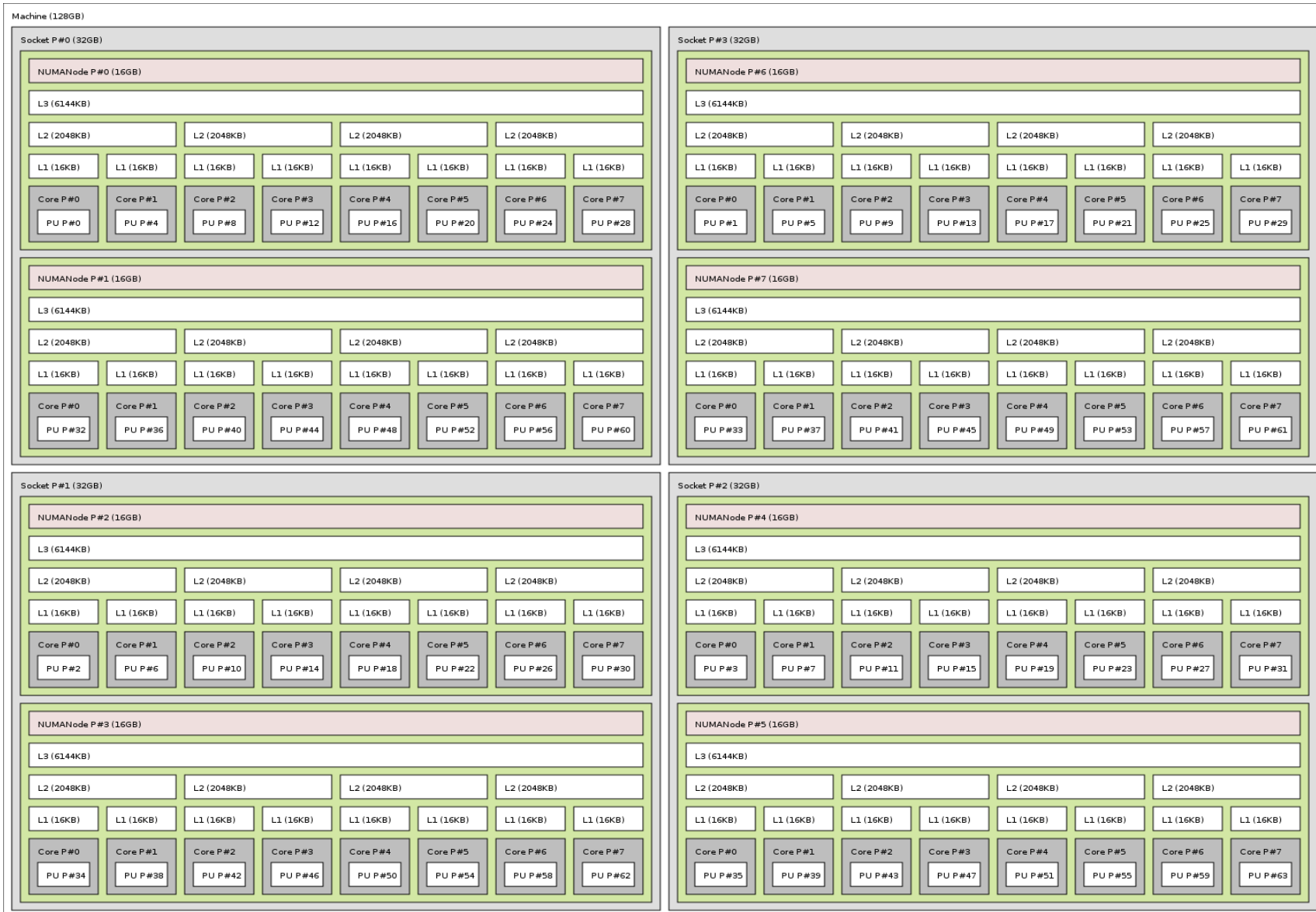Guillaume Mercier – TADaaM Team – Inria Bordeaux Sud-Ouest

# Machines are increasingly complex

- Multiple processors

- Multicore processors

- Simultaneous multithreading

- Shared caches

- NUMA effects


- We cannot expect users to understand all this ...

- … But we can help them to take advantage of  this complexity

  - Often seen as a hierarchy of resources

# A random dual-processor quad-core machine

# A random larger machine

# Current support of HW topology in MPI

- **MPI is hardware-agnostic**
  - And should remain so (I.e no assumptions about the HW)
  - Doesn't prevent from accessing the HW topology from MPI directly
- **Virtual topologies**
  - Machine-independent
  - Virtual to physical mapping "outside of the scope of MPI"
- **No standard behaviour (implementation-dependent)**
  - Reordering
    - MPI_Dist_graph_create, MPI_Graph_map, etc.
    - Side-effect of the function, implementation-dependent
  - Process Managers mapping and binding options
- **MPI Sessions?**

# Motivation

- Application developpers need abstract features to:
  - Deal with hardware characteristics (Caches, Interconnect, Cores, NUMA nodes, etc.)
  - Deal with low level tools (Hwloc, Lib_NUMA, Etc.)
- Expected performance improvements
  - Improved locality
  - Improved communication performance

# Basic Idea

- Use available abstractions in MPI: communicators

  - Well-known concept/object in MPI programming

  - "Natural" fit for our purpose:

    - Group MPI processes in communicators for each meaningful level in the hierarchy of the physical topology

    - Usable for collective communications

- Rather expand than add new features

  - Leverage existing mechanisms and abstractions

# Communicator creation functions

- MPI_Comm_create
- MPI_Comm_dup and friends
  - idup
  - with_info
- MPI_Comm_split
- MPI_Comm_split_type
  - MPI_COMM_TYPE_SHARED
  - MPI_COMM_TYPE_ADDRESS_SPACE (issue #31)
  - Implementations can also define their own values

# Proposal

- Add a new predefined value for the split_type arg
  - e.g MPI_COMM_TYPE_PHYSICAL_TOPOLOGY
  - Or any suitable (meaningful) name
- Property of the newly created communicator(s)
  - All newly created communicators should be a strict subset of the input communicator
    - MPI_Comm_compare(oldcomm,newcomm) yields MPI_UNEQUAL
    - **To ensure we don't create several useless communicators in case of physical levels redundancy**
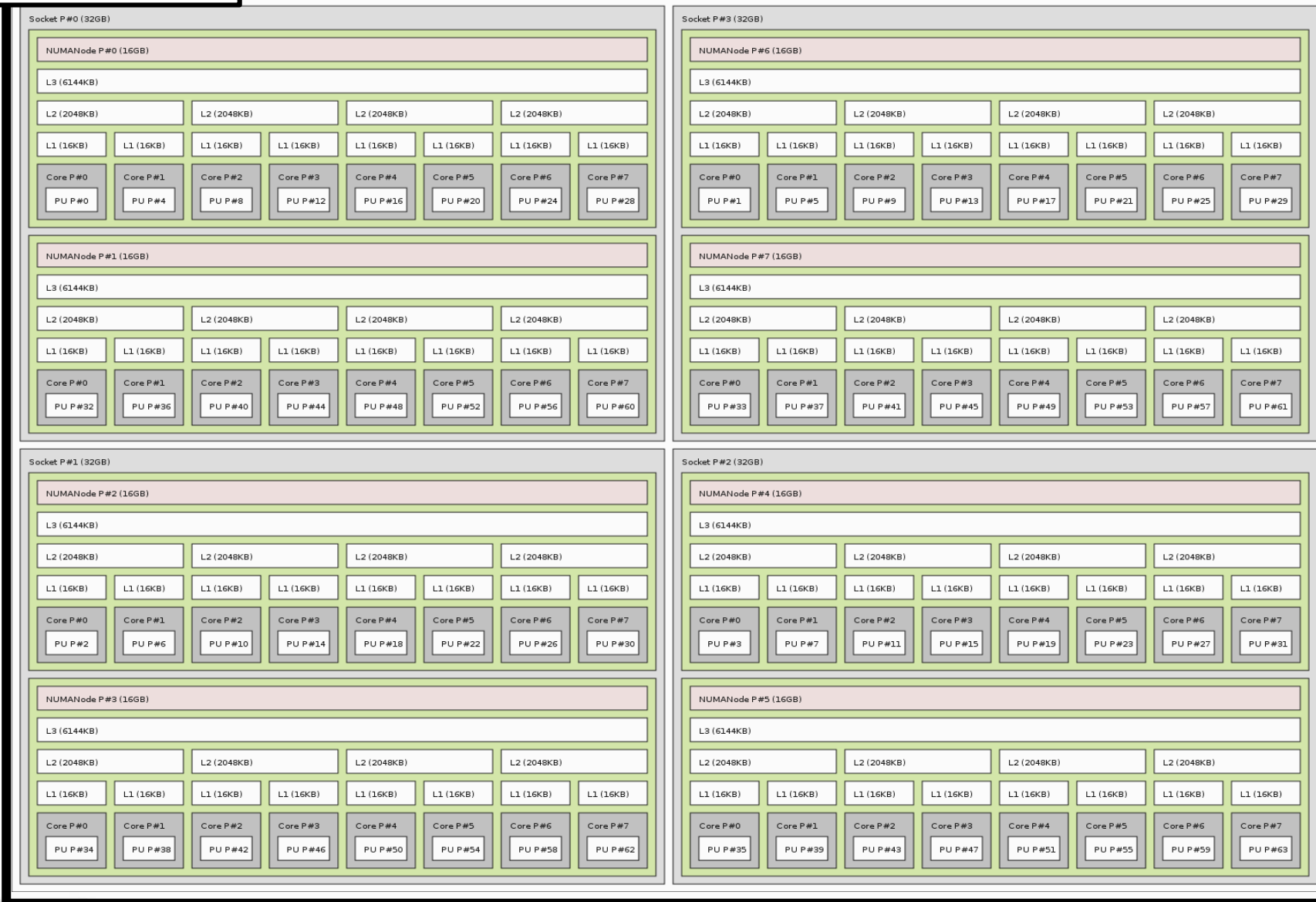
# Practical Use

- Recursively split an input communicator until the bottom of the hierarchy is reached (MPI_COMM_NULL )

  - **Independent of the hierarchy depth**

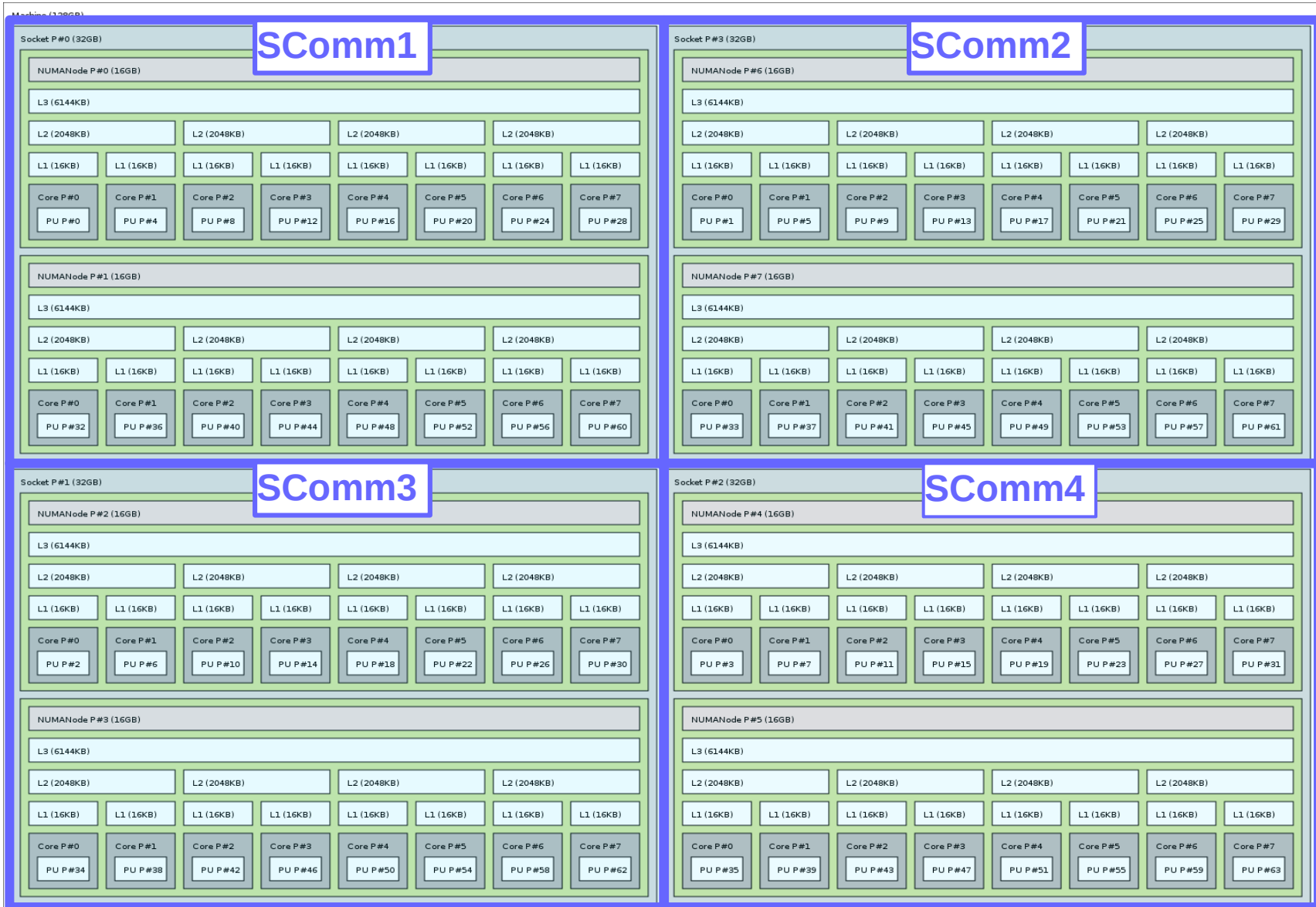  - **No (fixed) names for communicators**

- Example:

```
#define NLEVELS 16
int in_rank;
MPI_Comm out_comm[NLEVELS];
in_comm = MPI_COMM_WORLD;
idx = 0;
while(in_comm != MPI_COMM_NULL){
    MPI_Comm_rank(in_comm,&in_rank);
    MPI_Comm_split_type(in_comm,
                        MPI_COMM_TYPE_PHYSICAL_TOPOLOGY,
                        in_rank,MPI_INFO_NULL,
                        &out_comm[idx]);
  in_comm = out_comm[idx];
  assert(++idx < NLEVELS);
}
```
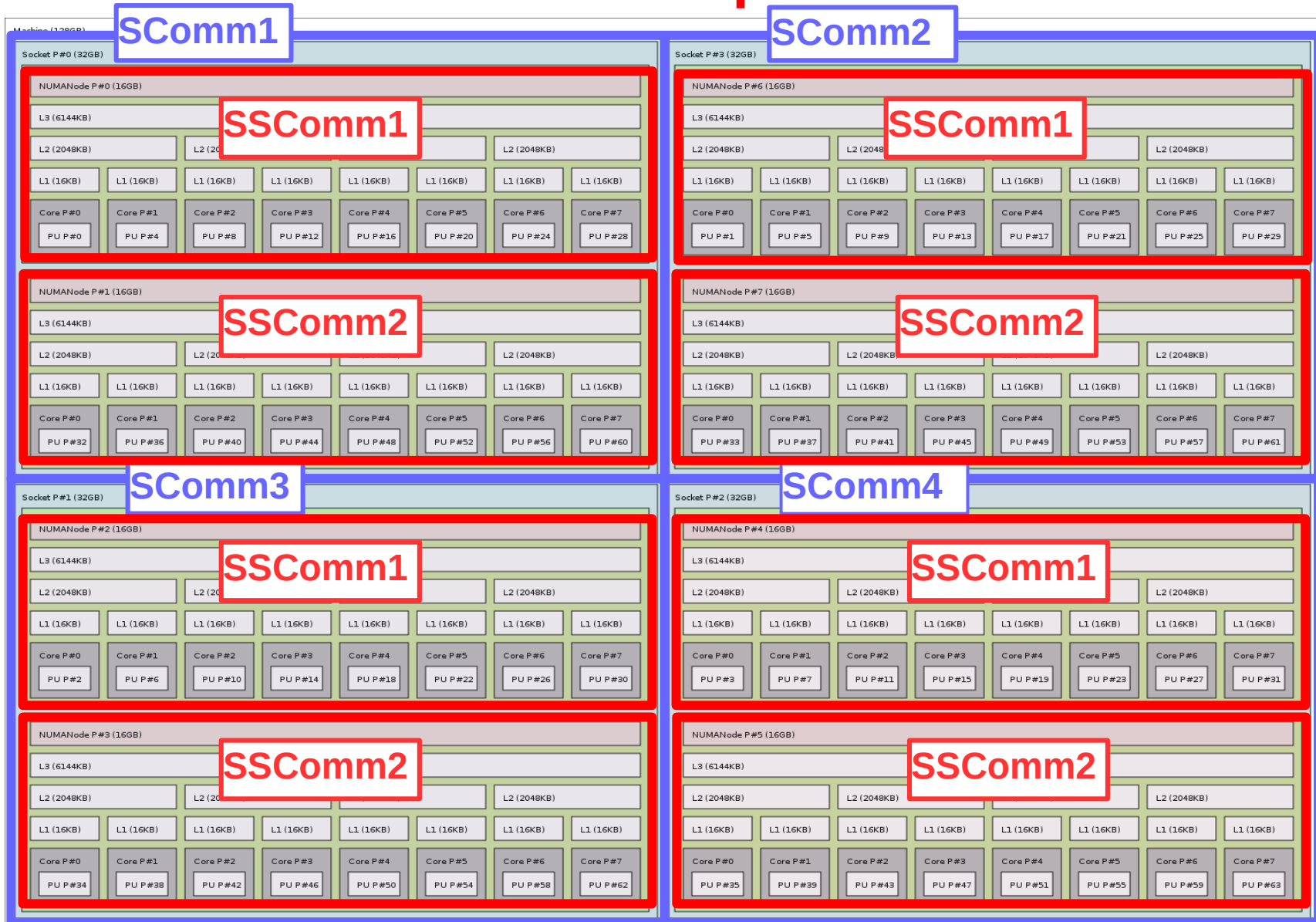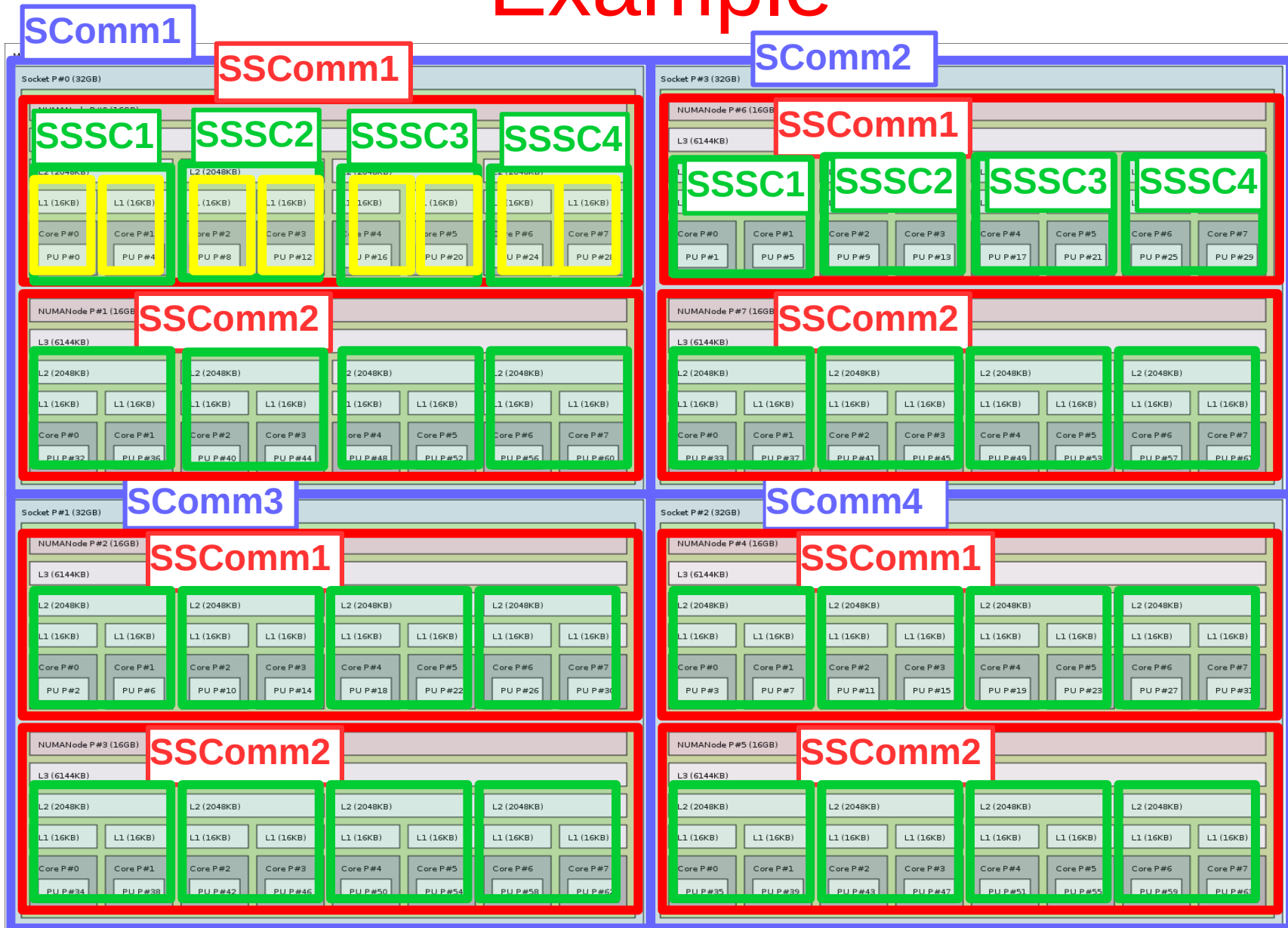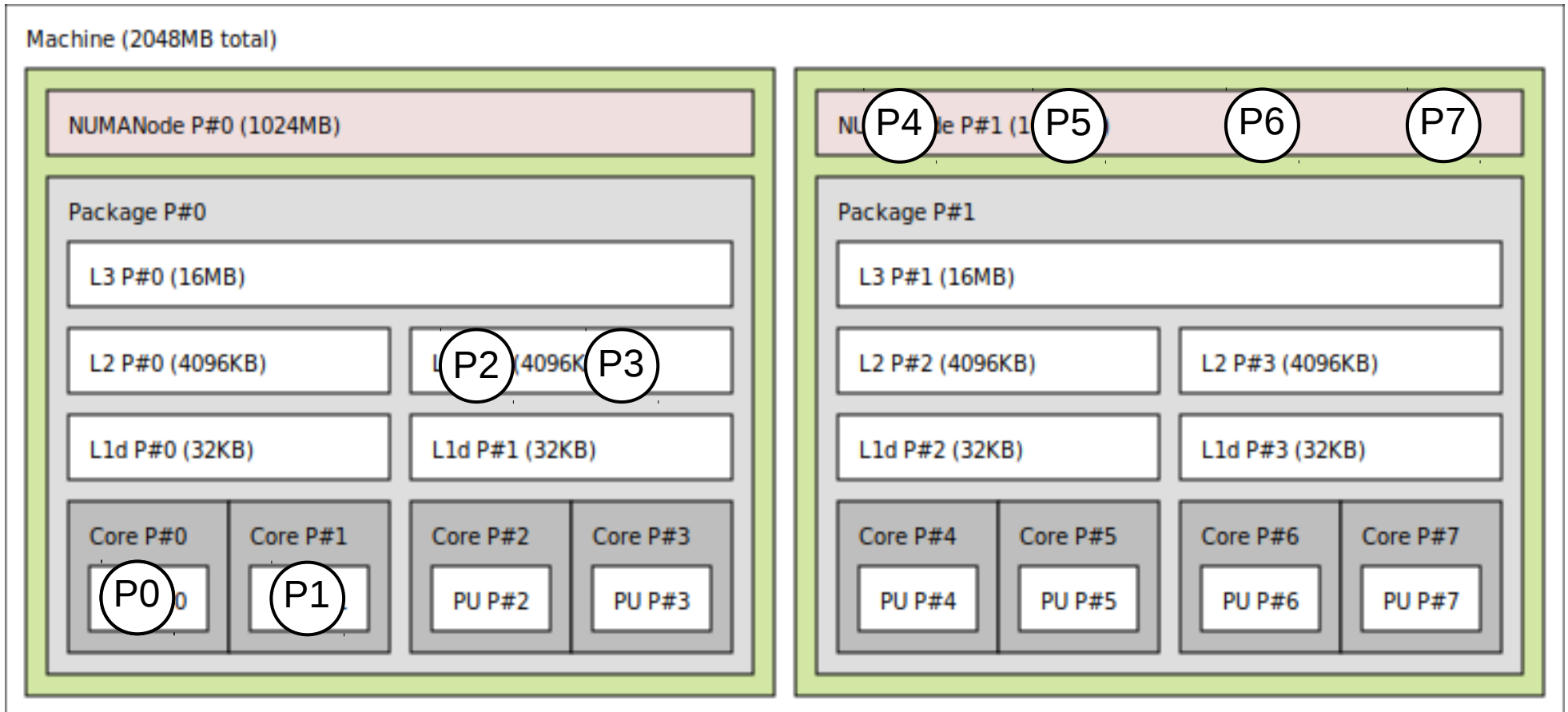
# Example

# Example

# Example

# Example

# Example

# Another example



Machine (2048MB total)

NUMANode P#0 (1024MB)

Package P#0

L3 P#0 (16MB)

L2 P#0 (4096KB)　　L2 P#1 (4096KB) P2 P3

L1d P#0 (32KB)　　L1d P#1 (32KB)

Core P#0　　Core P#1　　Core P#2　　Core P#3

P0　　P1　　PU P#2　　PU P#3

NUMANode P#1 (1024MB) P4 P5 P6 P7

Package P#1

L3 P#1 (16MB)

L2 P#2 (4096KB)　　L2 P#3 (4096KB)

L1d P#2 (32KB)　　L1d P#3 (32KB)

Core P#4　　Core P#5　　Core P#6　　Core P#7
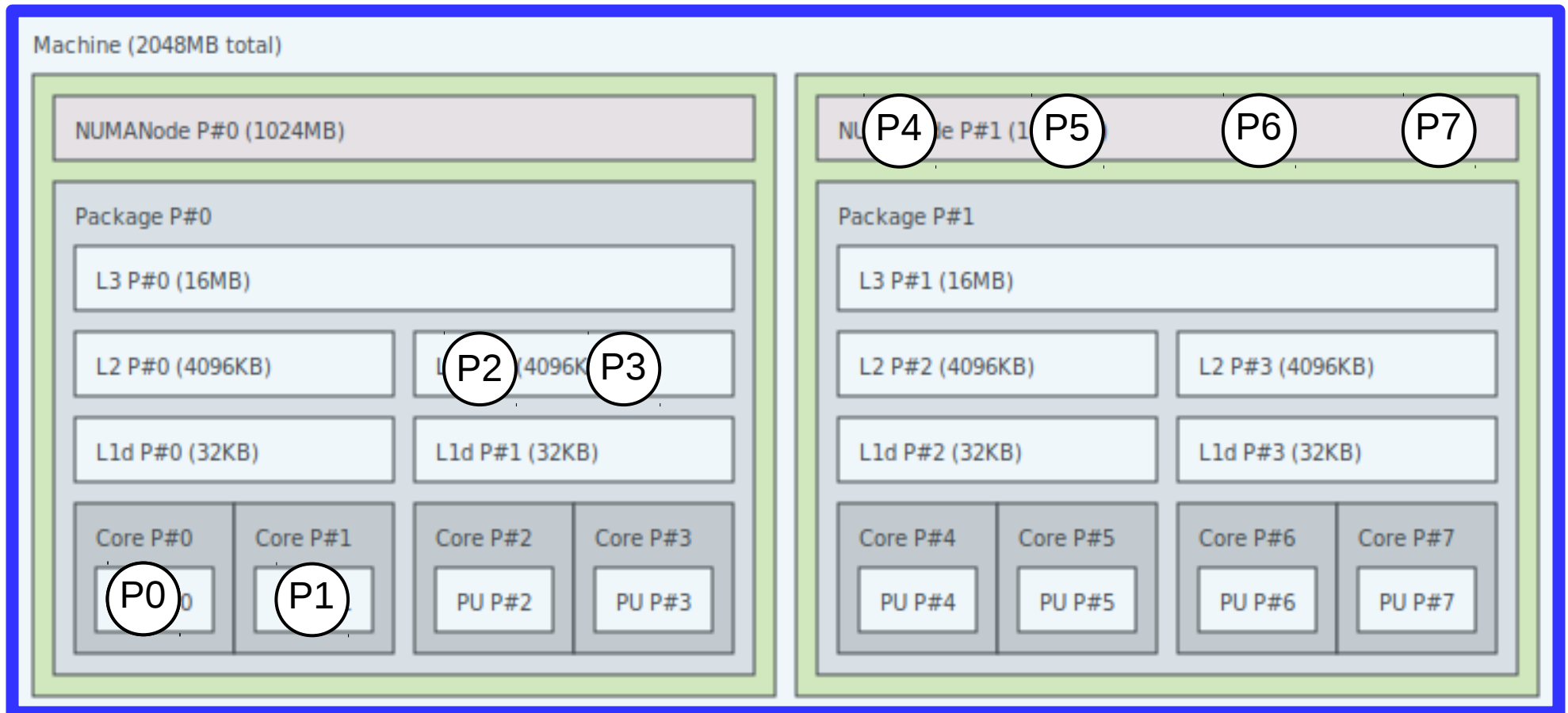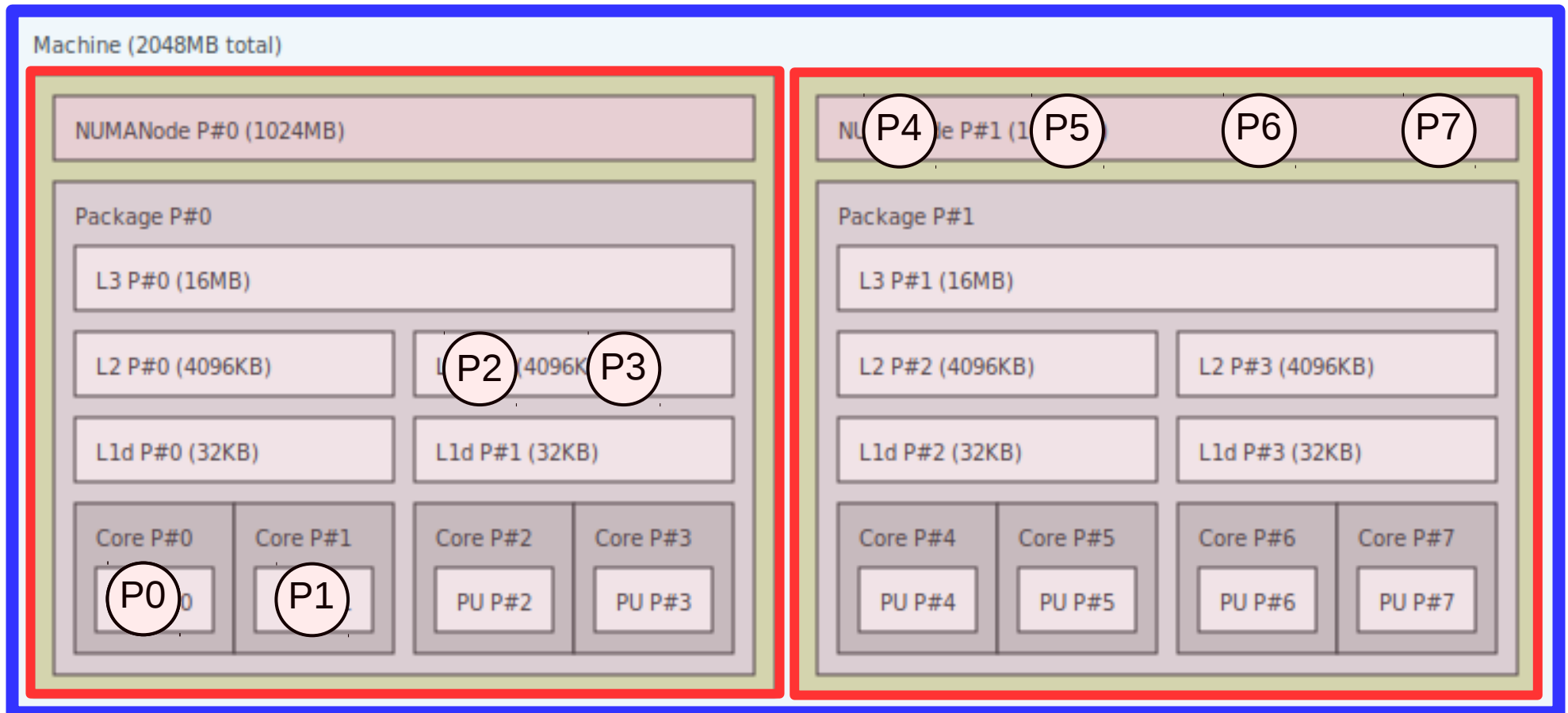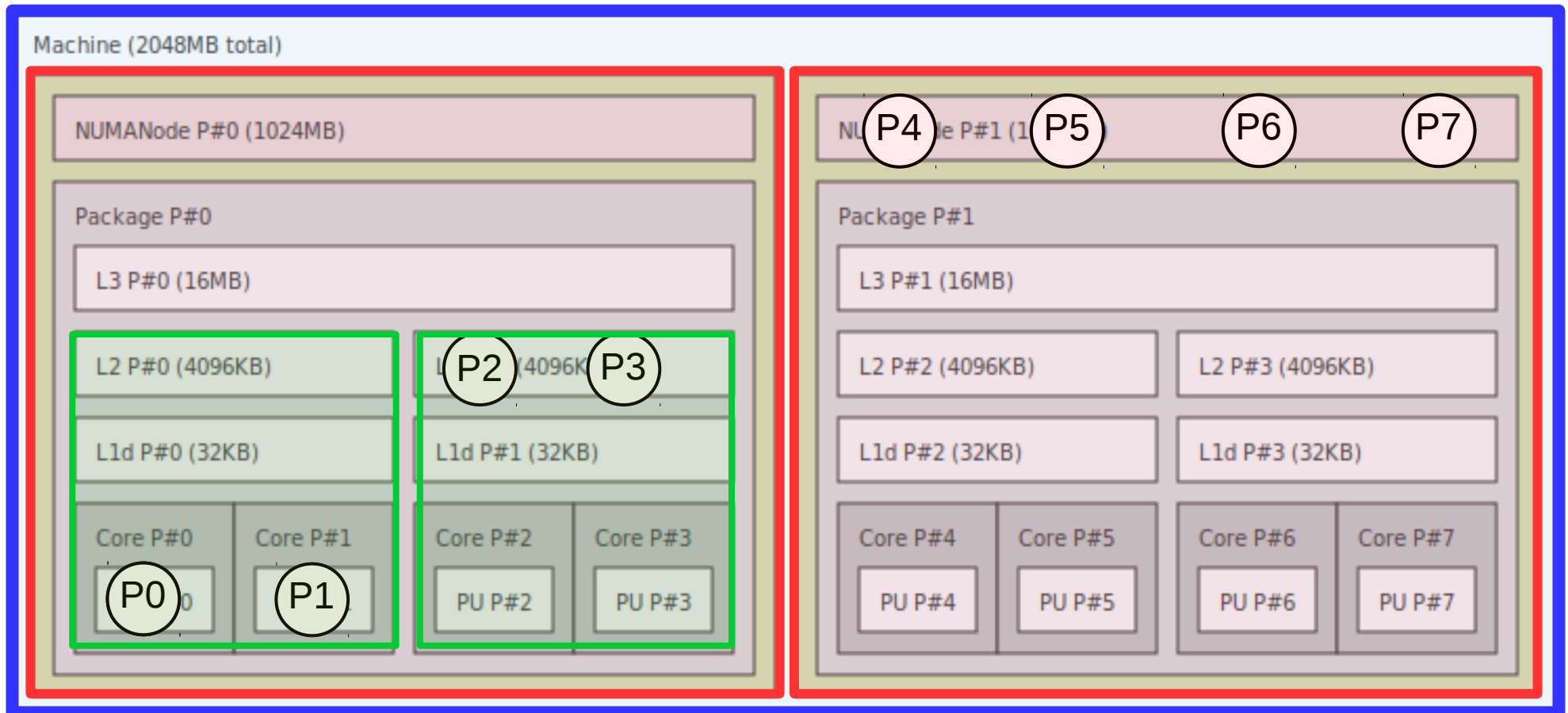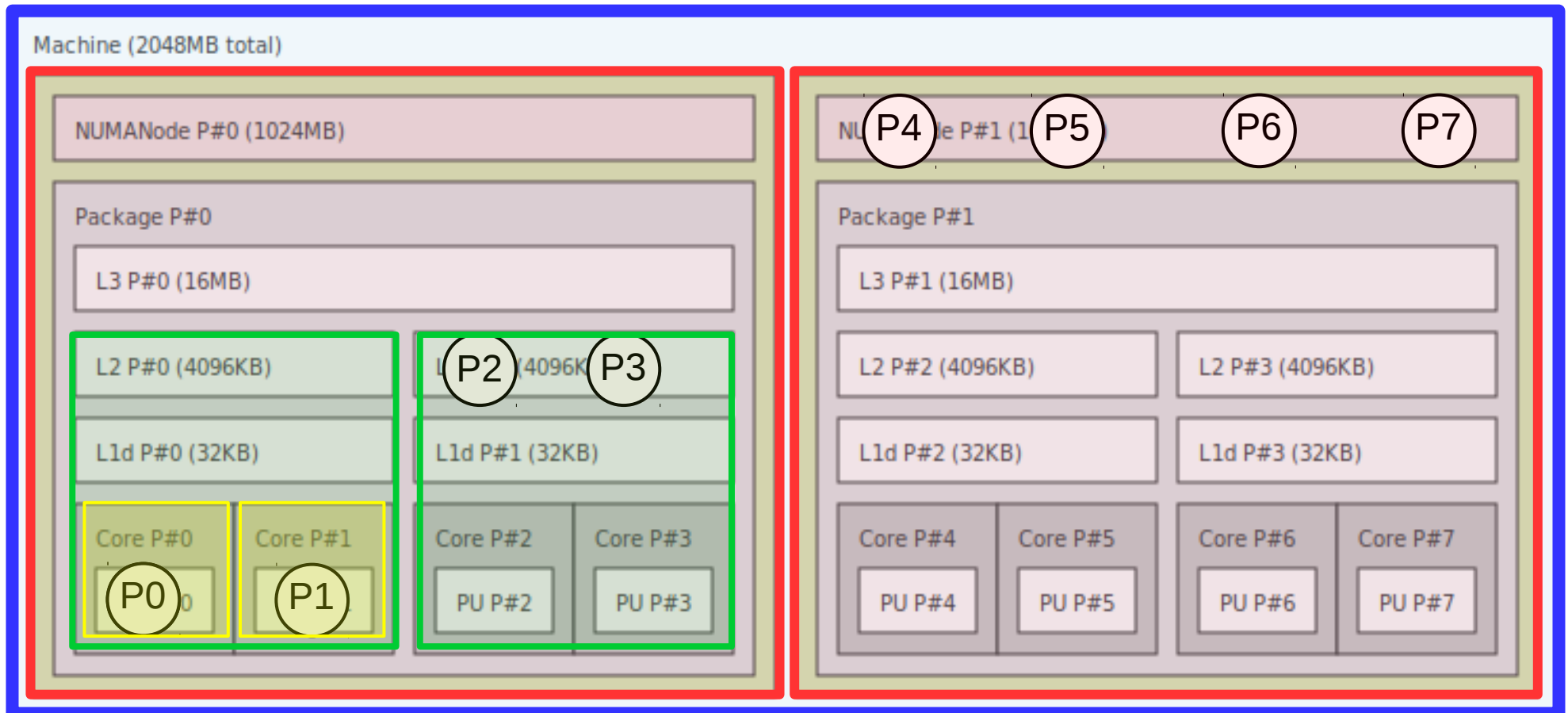
PU P#4　　PU P#5　　PU P#6　　PU P#7

# Another example

# Another example

# Another example

# Another example

# Other routines

- Query routine:

int MPI_Get_min_hierarchy(MPI_Comm comm,int size, int *ranks, MPI_Info info);

- Returns the deepest level in the hierarchy encompassing the ranks
- Result is a string (key: MPI_HIER_LEVEL)
- In practice, an Hwloc name (e.g : L1,L2,Package)
- Really useful?

# Potentially useful other things

- For data distribution, one might want:
  - Subcommunicators number
  - Subcommunicators ranks
- This can be stored in an info objet attached to the subcommunicator
  - MPI_Comm_set_info
  - MPI_Comm_get_info
- Really needed?
  - See Roots comms
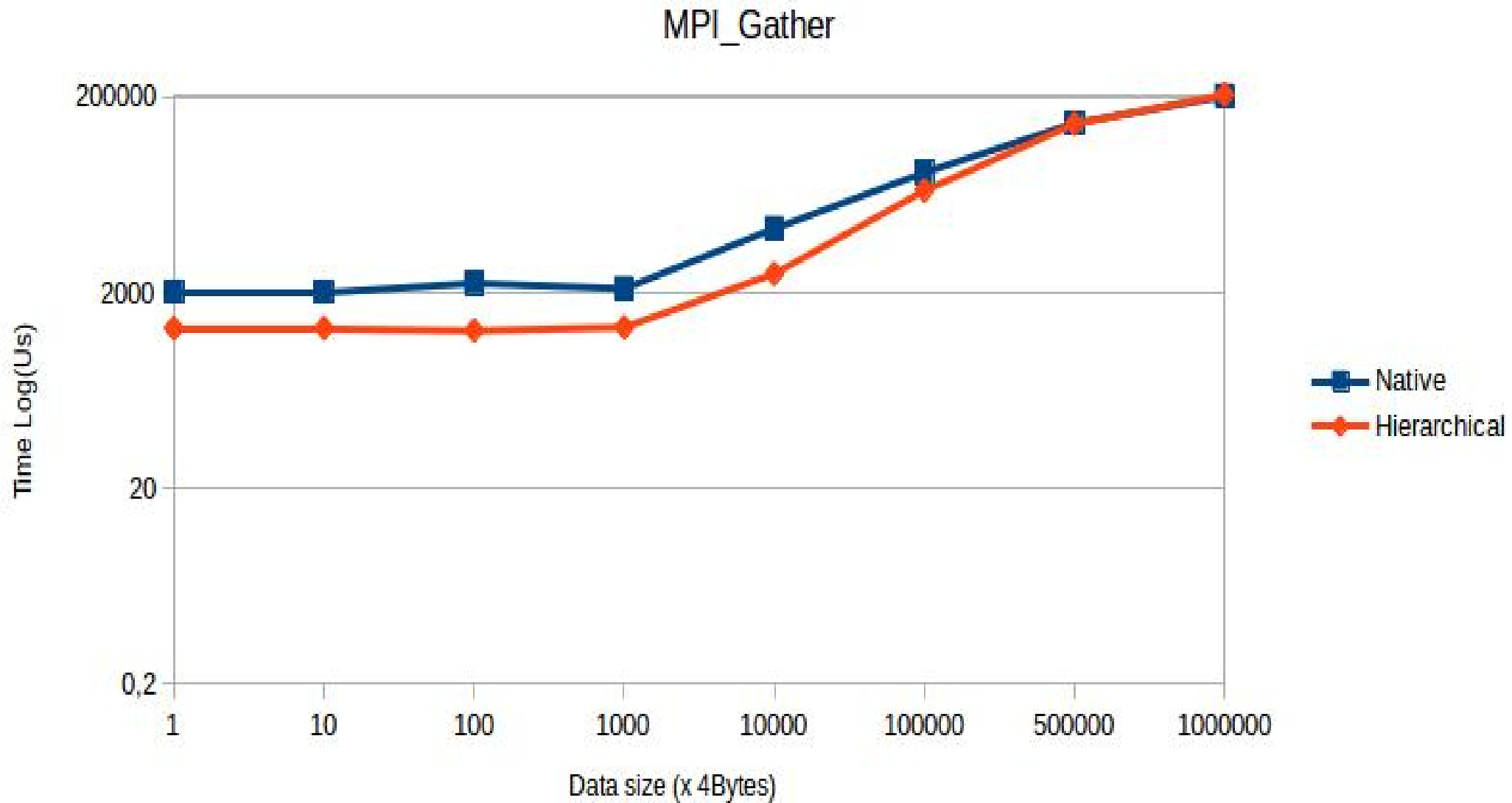
# Roots Communicators

- In practice, the subcommunicators are not enough to write applications

  - Need to access easily the roots of each subcommunicator

  - Hierarchy of Roots subcommunicators

    - The root of a communicator should also be the root of a newly created communicator in the hierarchy

    - Easy with the right key in Comm_split

- Lastovtesky papers

  - Hierarchical decomposition of collectives can improve performance

  - Topology-oblivious implementations (so far)
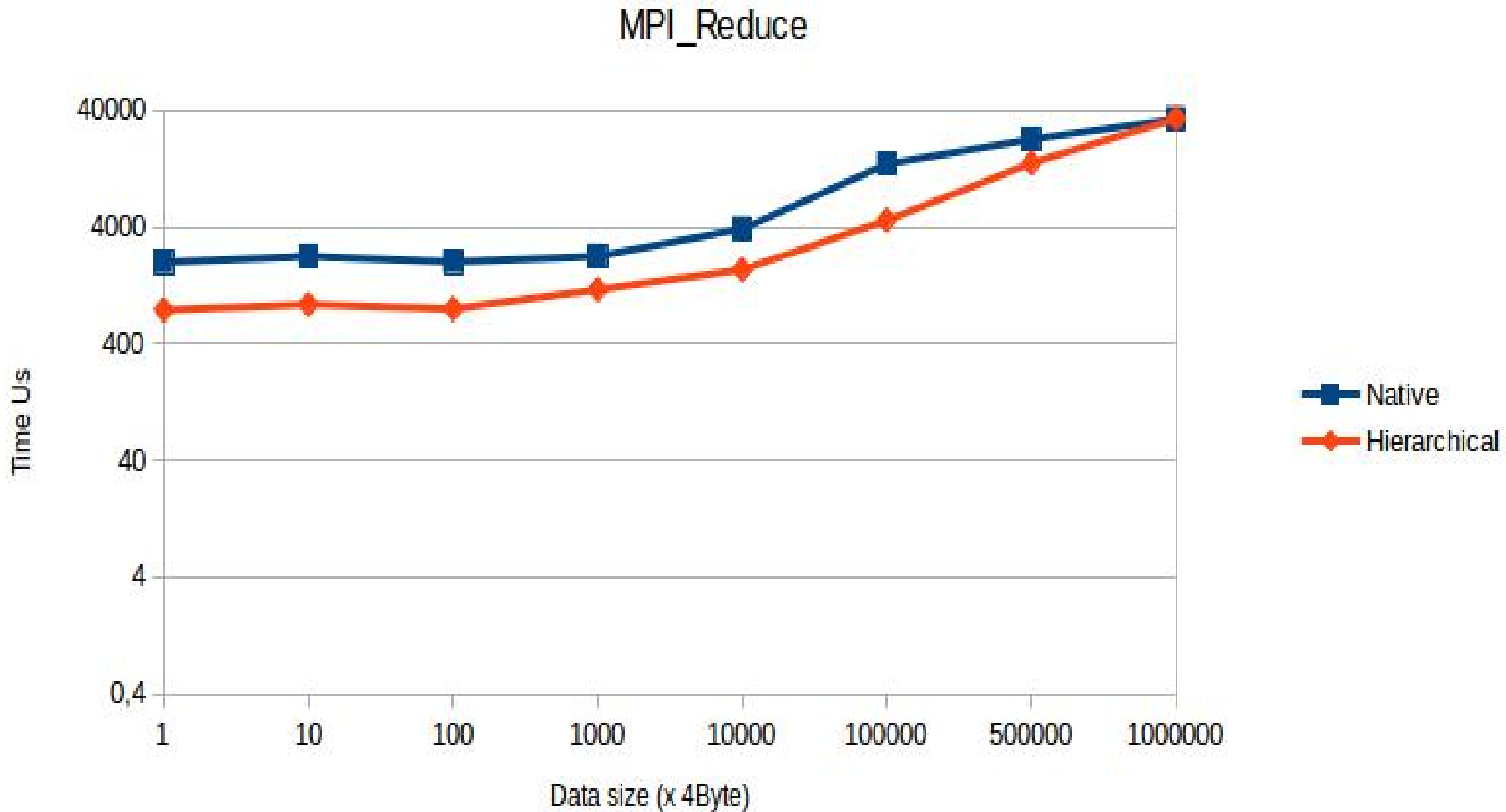
# Implementation

- Prototype available

- Implemented as an external library

- Hwloc-based implementation

  - Does not (yet) address network topology

  - But Netloc should help :)

- Use of MPI_Comm_split

- Roots Communicators creation

  - Easy with hwloc: use the logical_index of the parent level as the color for the split

  - Possible with MPI calls, but less efficient (involves more collectives)

# Preliminary Results



MPI_Gather

# Preliminary Results



MPI_Reduce

# Acknowledgements

- Brice Goglin

- Emmanuel Jeannot

- Farouk Mansouri

- Work funded by ELCI Project :

    http://elciproject.fr/