*informatics* *mathematics*

# A Topology-Aware Performance Monitoring Tool for Shared Resource Management in Multicore Systems
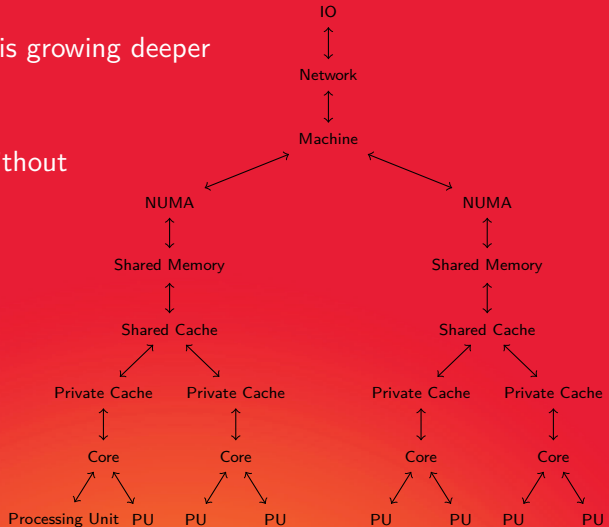
1. Context/Motivations

2. Fast presentation of the tool

3. Demonstration

4. How does it works ?

5. How is it made ?

6. Features & Future Works

# MOTIVATIONS

Memory hierarchy is growing deeper and larger.

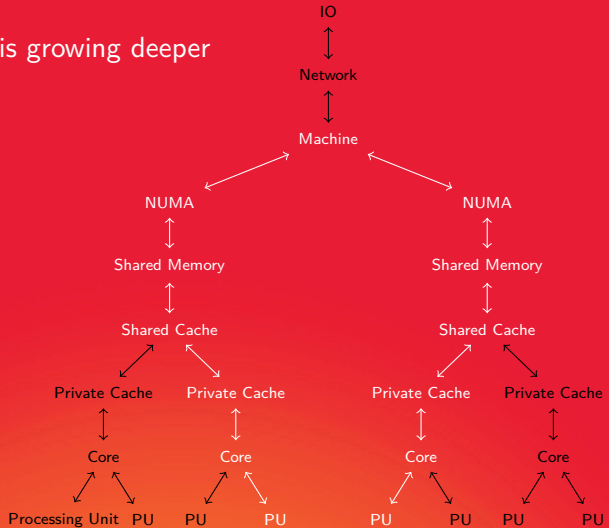No performance without a fair usage of the system topology

Batch schedulers, runtimes, applications themeselves . . . are getting topology aware.

# MOTIVATIONS

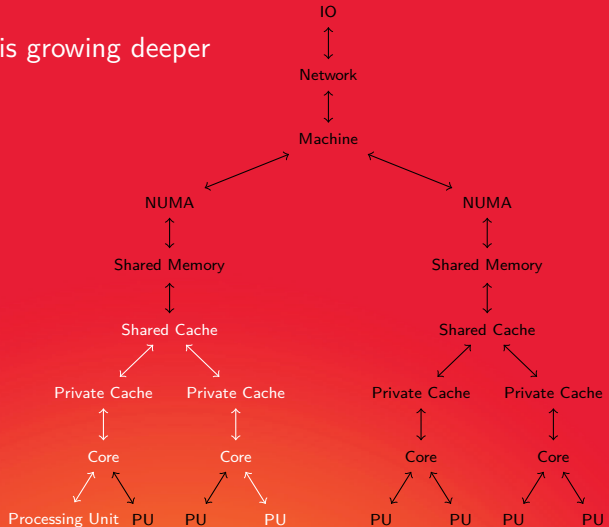Memory hierarchy is growing deeper and larger.

Hence, data management gives opportunities for performance improvements.

# MOTIVATIONS

Memory hierarchy is growing deeper and larger.

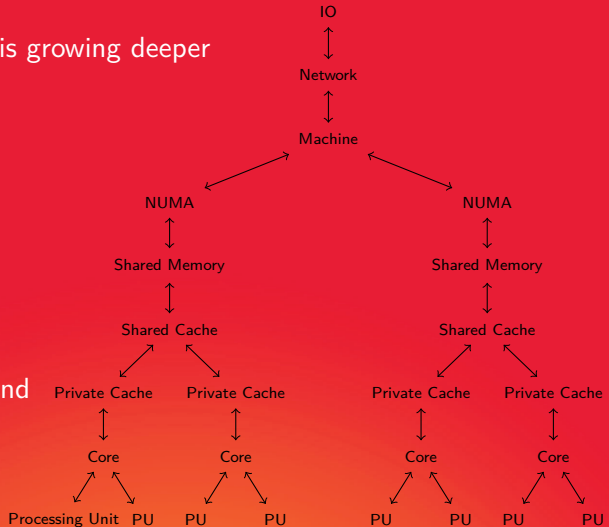Hence, data management gives opportunities for performance improvements.

# MOTIVATIONS

Memory hierarchy is growing deeper and larger.

Hence, data management gives opportunities for performance improvements.

It is a multi-level and a multi-criteria problem.

# MOTIVATIONS

- Need to match use cases, and relevant performance metrics for each level.
- **Need to match performance and topology.**



- Requires topology modeling skills.
- Requires adaptable performance monitoring.

# Yet Another Tool to Monitor Applications Performance

• Focus on data presentation to link the results with topology informations.

• Relies on two cornerstones of topology modeling (hwloc) and performance counter abstraction (PAPI) to map the latter on the former

• Minimal configuration and software requirements.

• Can help finding and caracterizing localized bottlenecks.
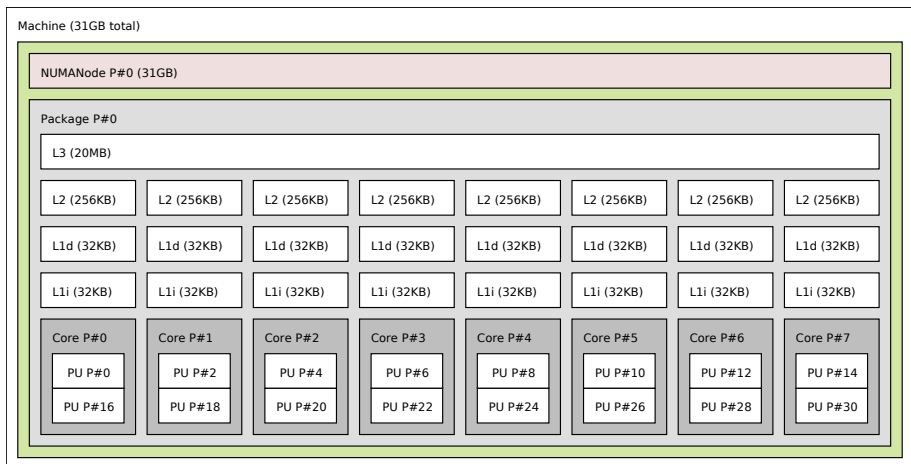
# Hardware Locality (hwloc)

Portable abstraction of hierarchical architectures for high-performance computing

- Performs topology discovery and extracts hardware component information.
- Provides tools for memory and process binding.
- Many operating systems supported
- ...
- lstopo utility to display the topology:

Developed at Inria Bordeaux.
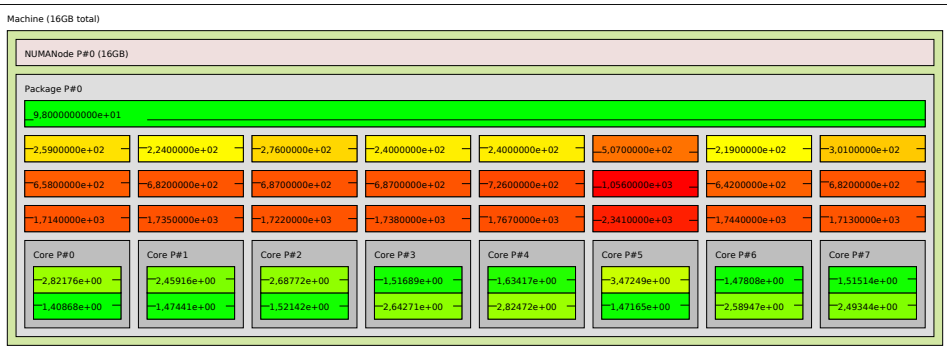
# Hardware Locality (hwloc)

# Performance Application Programming Interface (PAPI)

Consistent interface and methodology for use of the performance counter hardware.

- Real time relation between software performance and processor events.
- Many operating systems supported too.
- Reliable and actively supported.
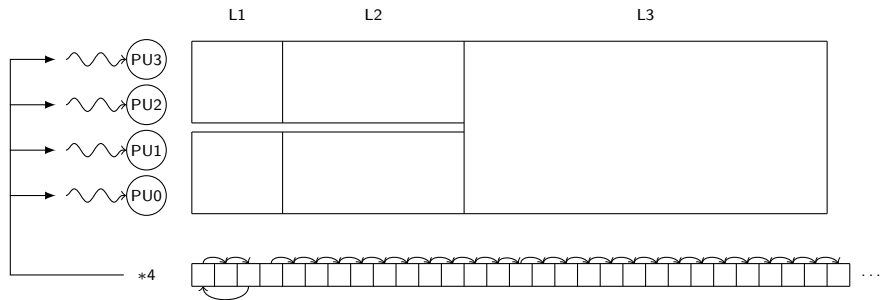- Used in a wide range of performance analysis applications.

An abstraction layer to plug some other performance library is under development.
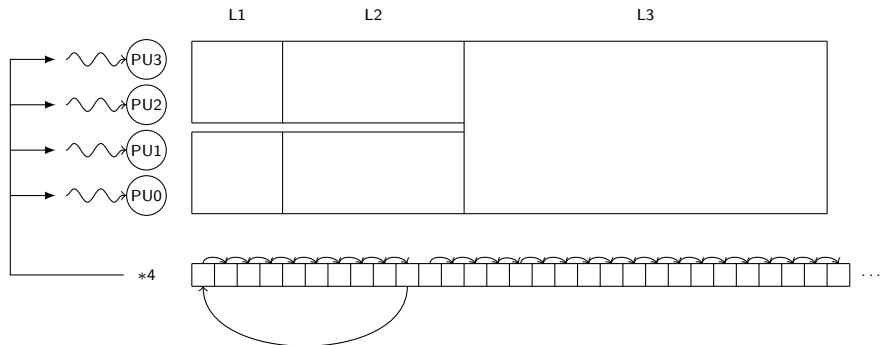
# Dynamic Lstopo (example)



Sample of hardware performance counters mapped on a single socket of an Intel Xeon E5-2650 CPU.
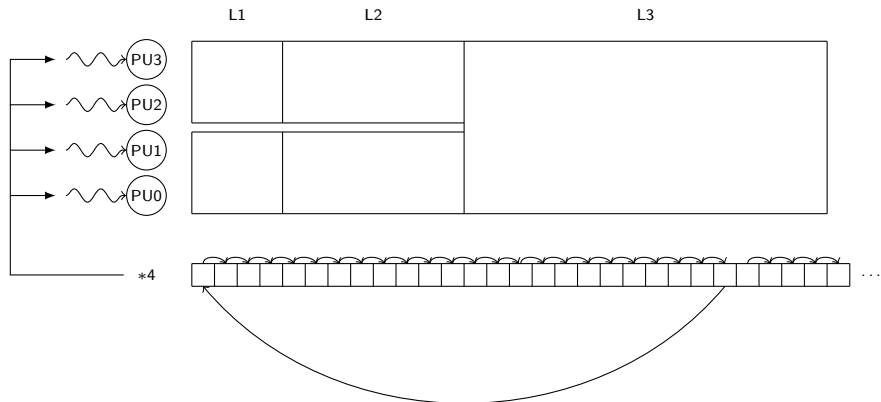
# A Demonstration Worth Thousand Words



Accesses to a linked list of variable size.

# A Demonstration Worth Thousand Words



Accesses to a linked list of variable size.
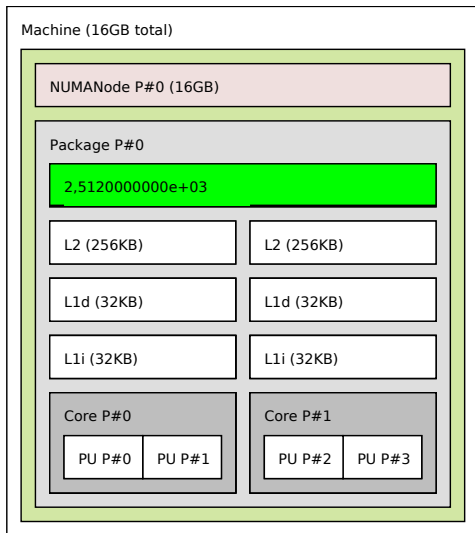
# A Demonstration Worth Thousand Words



Accesses to a linked list of variable size.

# A Demonstration Worth Thousand Words

```
L3_MISS{
  OBJ = L3;
  CTR = PAPI_L3_TCM;
  LOGSCALE = 1;
}

L2_MISS{
  OBJ = L2;
  CTR = PAPI_L2_TCM;
  LOGSCALE = 1;
}
```

```
L1_MISS{
  OBJ = L1d;
  CTR = PAPI_L1_DCM;
  LOGSCALE = 1;
}

SINGLE_L3_MISS{
  OBJ = PU;
  CTR = PAPI_L3_TCM;
  LOGSCALE = 1;
}
```

# Dynamic Lstopo (Usage)



Counters input:
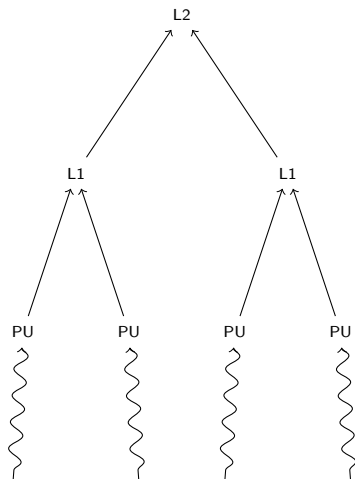
```
SINGLE_L3_MISS{
  OBJ = L3;
  CTR =
  PAPI_L2_TCM/PAPI_L2_TCA;
  LOGSCALE = 1;
  MAX=1000000;
  MIN=0;
}
```
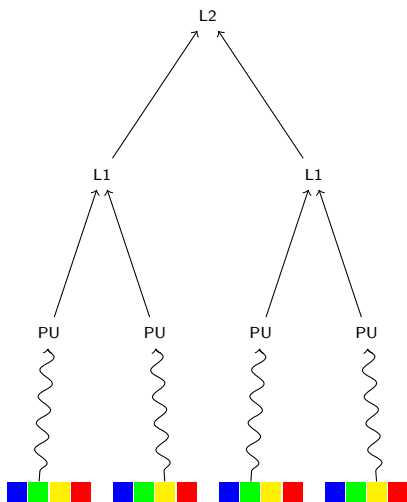
Command line:
lstopo –perf-input counters_input

# Dynamic Lstopo (Theory)

1. Spawn one pthread per hardware thread (PU#0, ..., PU#3).

# Dynamic Lstopo (Theory)

1. Spawn one pthread per hardware thread (PU#0, ..., PU#3).

2. For each timestamp, each thread collects a local set of performance counters.

# Dynamic Lstopo (Theory)

1. Spawn one pthread per hardware thread (PU#0, . . . , PU#3).

2. For each timestamp, each thread collects a local set of performance counters.

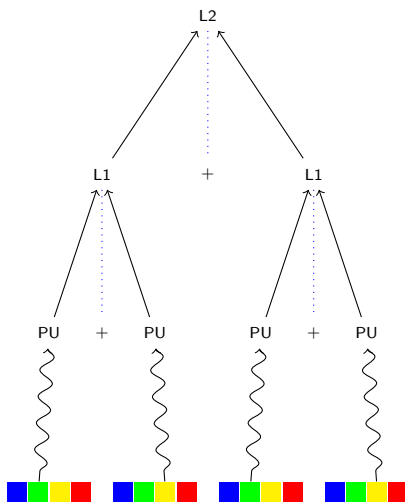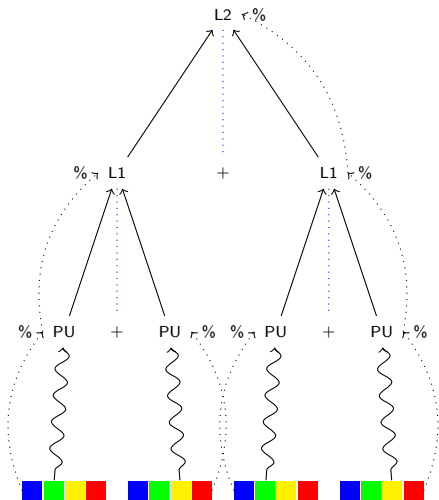3. Counters are accumulated in each upper level.

# Dynamic Lstopo (Theory)

1. Spawn one pthread per hardware thread (PU#0, ..., PU#3).

2. For each timestamp, each thread collects a local set of performance counters.

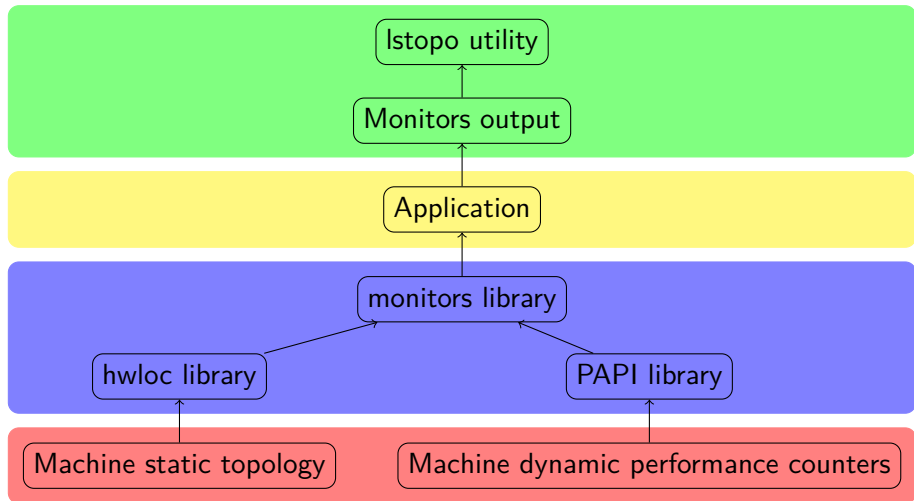3. Counters are accumulated in each upper level.

4. For each level, a leaf computes an arithmetic expression of the performance counters in the set.

# Dynamic Lstopo Software Architecture in Brief

# Dynamic Lstopo Software Architecture in Brief

# API

```
monitors = load_Monitors_from_config(NULL,
    "my_perf_file", "my_output_file", 0);

Monitors_watch_pid(monitors, getpid());

Monitors_start(monitors);

/* ... */

Monitors_update_counters(monitors);

delete_Monitors(monitors);
```

# Dynamic Lstopo (Output paje trace)

```
%EventDef val 0                1 3 PU 0 SINGLE_L3_MISS 1
%    Id        int             1 2 PU 1 SINGLE_L3_MISS 1
%    Phase     int             1 1 PU 2 SINGLE_L3_MISS 1
%    Time_us   date            1 0 PU 3 SINGLE_L3_MISS 1
%    Value     double
%EndEventDef                   0 2 0 962832762224 67,00
                               0 1 0 962832762224 58,00
%EventDef container 1          0 3 0 962832762225 77,00
%    Id        int             0 0 0 962832762236 64,00
%    Level     string          0 3 0 962832860676 94514,00
%    Sibling   int             0 0 0 962832860676 121746,00
%    Name      string          0 2 0 962832860676 205170,00
%    Logscale  int             0 1 0 962832860717 200931,00
%EndEventDef
```

# Features

- Record and/or Display live machine performance counters and match them with topology.

- Several settings: counters accumulation, sampling rate, attach to a process. . .

- Replay any trace with a topology file (for external display)..

- Sample specific parts of an application with the monitor library.

- Support legacy lstopo options (restrict topology, change display format. . . ).

# Future works

- Match code and performance informations

- Accept user defined aggregation operator.

- Provide performance abstraction layer.

- Be able to delimit phases during execution.

- Find and give explicit hints on bottlenecks.

# Conclusion

Data locality becomes a main criterion for high performance.

We built a tool based on a topology model and a performance library to help taking up the challenge.

It maps performance values to machine objects.

It is a visual tool, fast and easy to use.

It is lightweight and causes less than 1% CPU overhead.

Let you build topology aware performance models.

Dynamic lstopo is into the process of beeing merged with hwloc project.

# Thank you