



PIC codes in the HPC environment



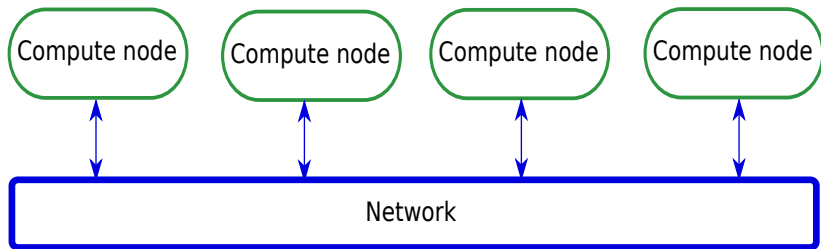


- 1 HPC environment, trends and perspectives
- 2 The PIC method and its parallelization
- 3 The load balancing issue



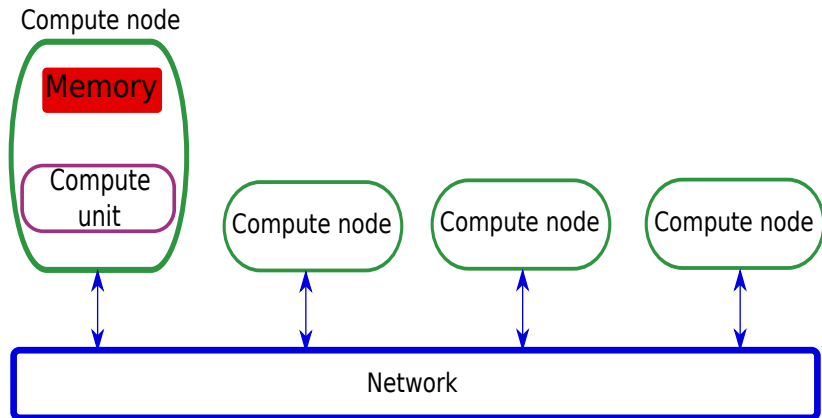
- 1 HPC environment, trends and perspectives
- 2 The PIC method and its parallelization
- 3 The load balancing issue

What is a super computer ?



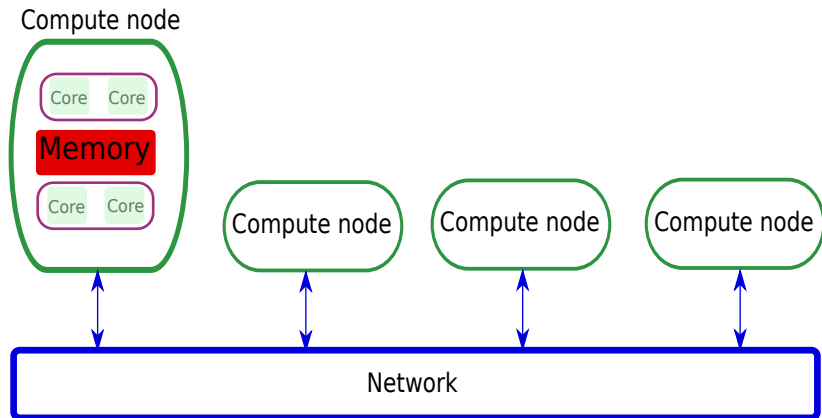
Distributed computing

What is a super computer ?



Distributed memory system

What is a super computer ?



Distributed {shared memory} system

Tianhe (China, June 2013) : 31 PFLOPS for 17 MW gives 1.85 GFLOPS/W.

Extrapolation : 1000 PFLOPS ==> 540 MW !

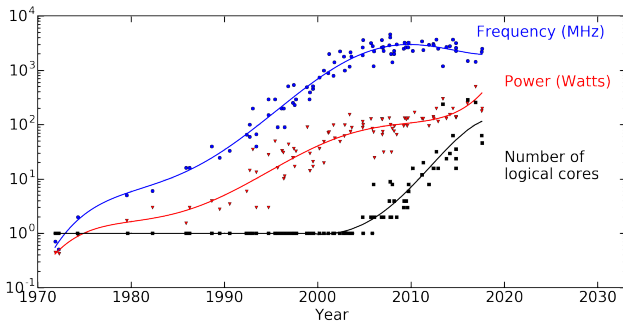


= ? =



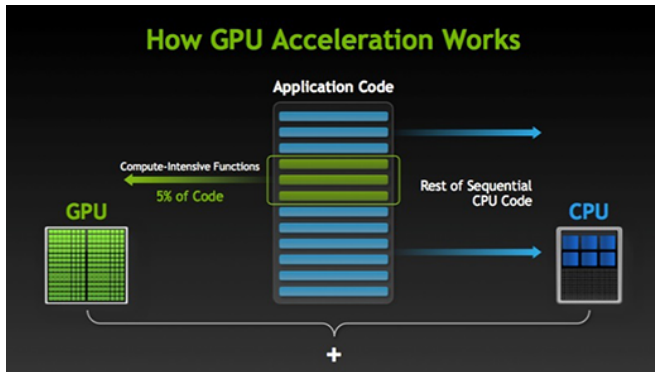
The objective is $P < 20 \text{ MW}$

The challenge for constructors is to increase both **total performance** and **energy efficiency** of computing nodes.



- Increased performances
- Reasonable energy budget

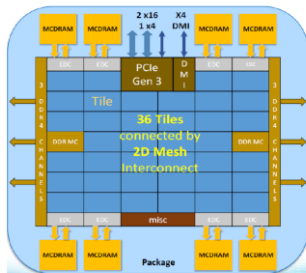




- Most energy efficient architecture today
- Difficult to adress :
 - Libraries : Cuda, OpenCL.
 - Directives programming : OpenMP 4 ou openACC.



KNL Overview



TILE



Chip: 36 Tiles interconnected by 2D Mesh

Tile: 2 Cores + 2 VPU/core + 1 MB L2

Memory: MCDRAM: 16 GB on-package; High BW

DDR4: 6 channels @ 2400 up to 384 GB

IO: 36 lanes PCIe® Gen3, 4 lanes of DMI for chipset

Node: 1-Socket only

Fabric: Intel® Omni-Path Architecture on-package (not shown)

Vector Peak Perf: 3+TF DP and 6+TF SP Flops

Scalar Perf: ~3x over Knights Corner

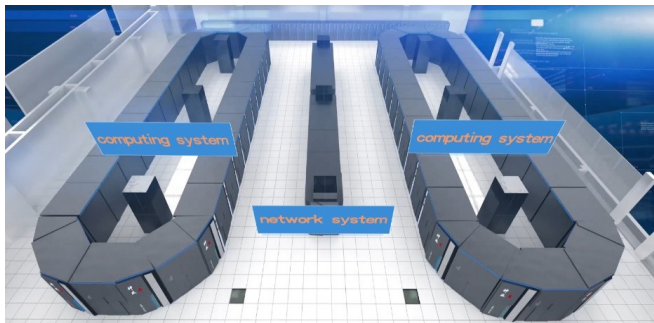
Streams Triad (GB/s): MCDRAM : 400+; DDR: 90+

MCDRAM
~5X Higher BW
than DDR

Source Intel. All products, computer systems, dates and figures specified are preliminary based on current expectations, and are subject to change without notice. KNL data are preliminary based on current expectations and are subject to change without notice. *Binary Compatible with Intel Xeon processors using Haswell Instruction Set (except TSX). 2D bandwidth numbers are based on STREAM like memory access pattern when MCDRAM used as flat memory. Results have been estimated based on internal Intel analysis and are provided for informational purposes only. Any difference in system hardware or software design or configuration may affect actual performance. *Other names and brands may be claimed as the property of others.



- Powers several top HPC systems.
- + Irene (France) - Aurora (U.S)
- Supposedly accessible through “Normal” programming but relies critically on the SIMD instruction set.



- Most powerful system in the world : 93 PFLOPS.
- 15 MW
- The SunWay architecture mimicks Xeon Phi.



Compiler Case Study

Introducing SIMD: Single Instruction, Multiple Data

• Scalar processing

- traditional mode
- **one operation produces one result**

X

+

Y

X + Y

X

• SIMD processing

- with SSE / SSE2
- **one operation produces multiple results**

x3 x2 x1 x0

+

y3 y2 y1 y0

X + Y

x3+y3 x2+y2 x1+y1 x0+y0

Intel Labs

Copyright © 2001 Intel Corporation.

- Excellent potential speed up, very good power budget.
- Heavy constraints on data structure and algorithm.
- Difficult to use at its full extent in a PIC code.



- U.S. : Exascale for 2021. No specifications.
- Japan : “Post K Supercomputer”. EFLOPS for 2020. Architecture ARM.
- China : 3 exascale systems for 2020.
- Europe : 2 Exascale systems for 2022. At least 1 powered by European technology (probably ARM).



As a developer

- 1 Expose parallelism. Massive parallelization is key.
- 2 Focus on the algorithm and data structures. Not on architectures.
- 3 Reduce data movement : Computation is becoming cheaper, loads and stores not so much.
- 4 Be aware of the increasing gap between peak power and effective performances. The race to exascale is becoming a race to exaflops.

As a scientist

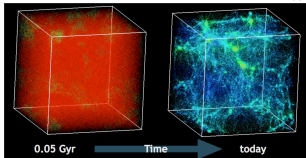
- 1 Collaborate with experts : complexity of HPC systems increases a lot !



- 1 HPC environment, trends and perspectives
- 2 The PIC method and its parallelization
- 3 The load balancing issue

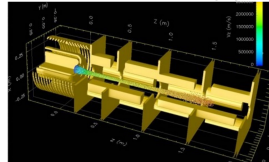
The Particle-In-Cell (PIC) method is a central tool for simulation over a wide range of physics studies

Cosmology



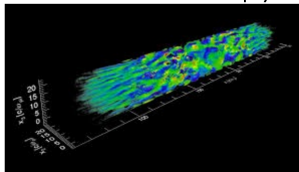
source: K. Heitmann, Argonne National Lab

Accelerator physics



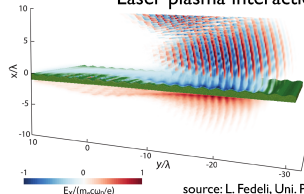
source: WARP, Berkeley Lab

Relativistic astrophysics



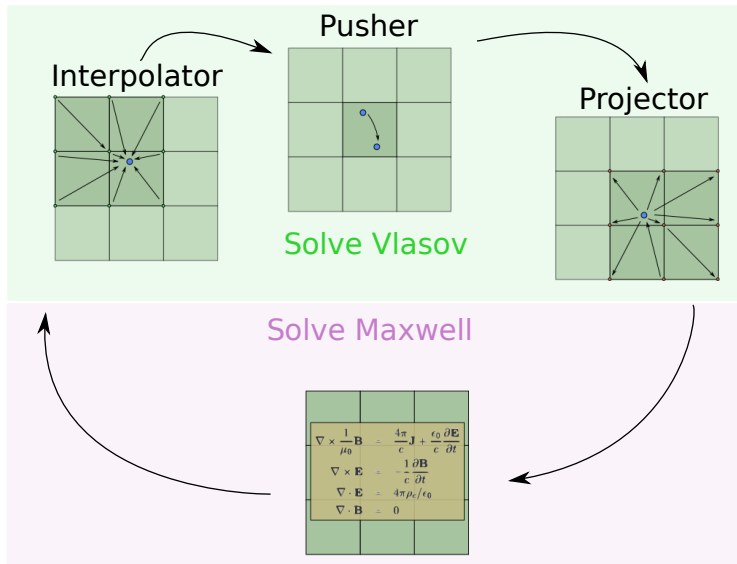
source: F. Fiuza, Livermore National Lab

Laser plasma interaction

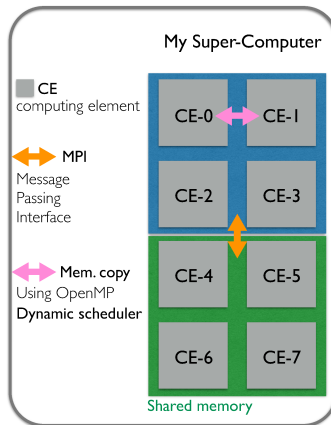
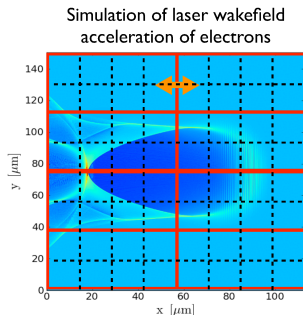


source: L. Fedeli, Uni. Pisa

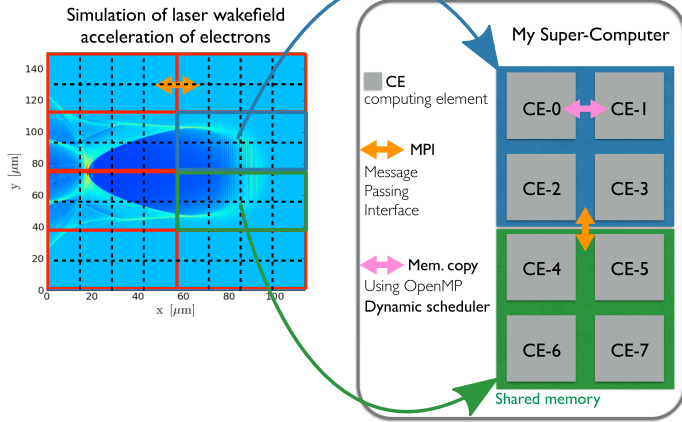
- Conceptually simple
- Efficiently implemented on (massively) parallel super-computers

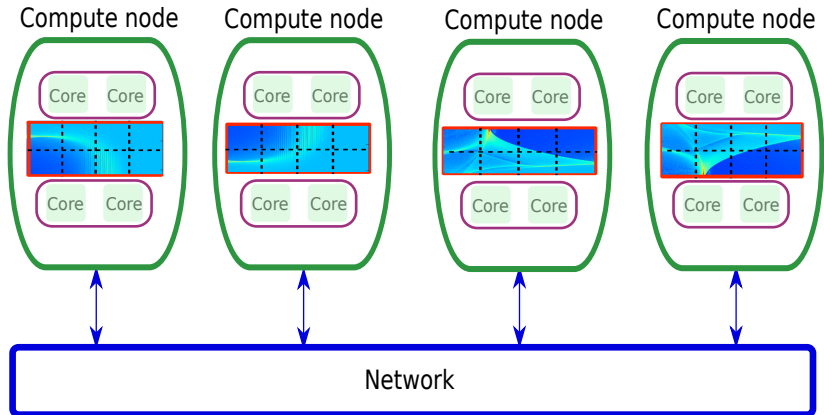


PIC code are 'easily' parallelized using **domain decomposition**



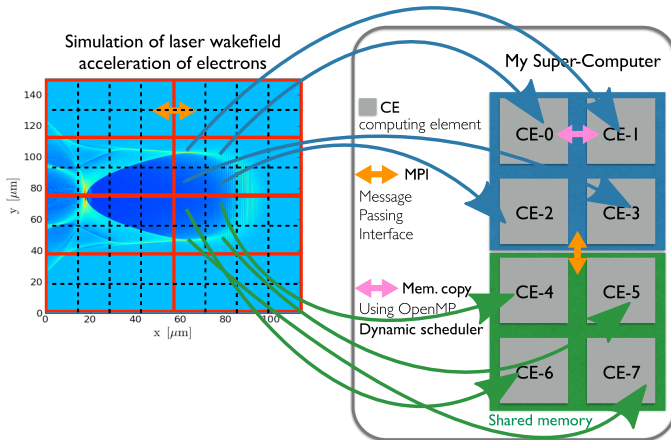
PIC code are 'easily' parallelized using domain decomposition

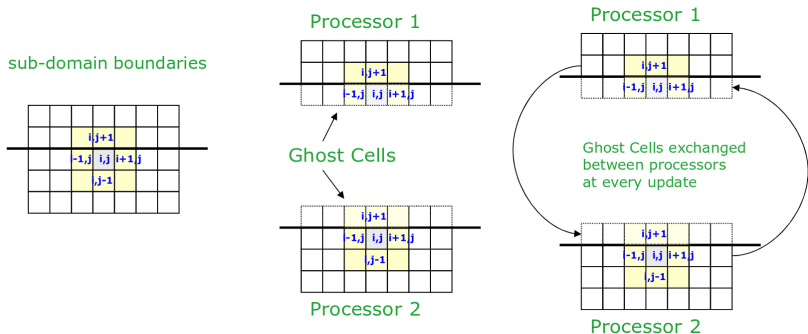






PIC code are 'easily' parallelized using **domain decomposition** + Patch





- If processors have a shared memory ==> OpenMP
- If processors have distributed memory ==> MPI
- Same logic for particles



Characteristics

- Library
- Coarse grain
- Inter node
- Distributed memory
- **Almost all HPC codes**

Issues

- Latency
- OS jitter
- Global communication scalability



Characteristics

- Compiler Directives
- Medium grain
- Intra node
- Shared memory
- **Many HPC codes**

Issues

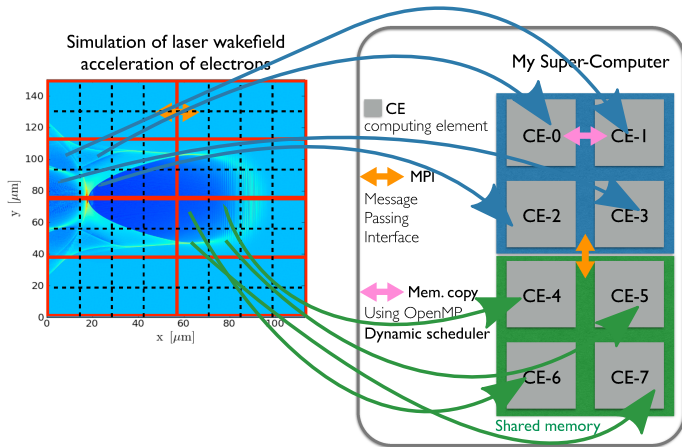
- Thread creation overhead
- Memory/core affinity
- Interface with MPI
(`MPI_THREAD_MULTIPLE`)



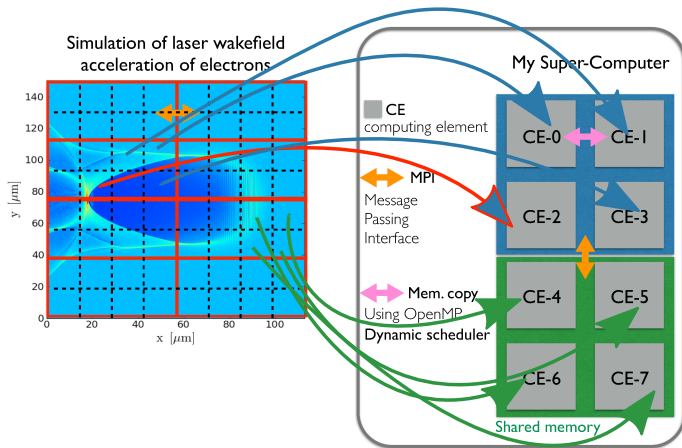
- 1 HPC environment, trends and perspectives
- 2 The PIC method and its parallelization
- 3 The load balancing issue



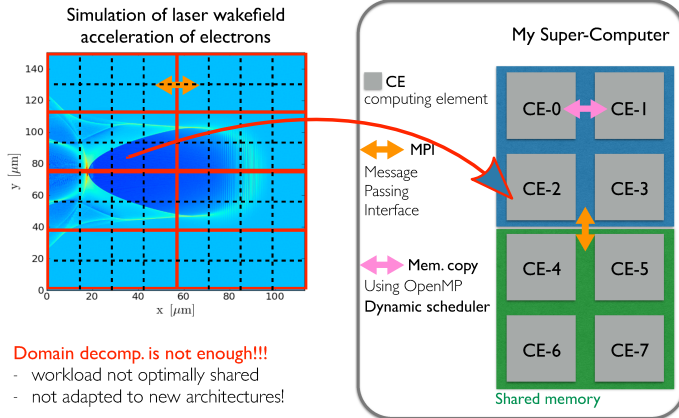
PIC code are 'easily' parallelized using **domain decomposition**



PIC code are 'easily' parallelized using **domain decomposition**

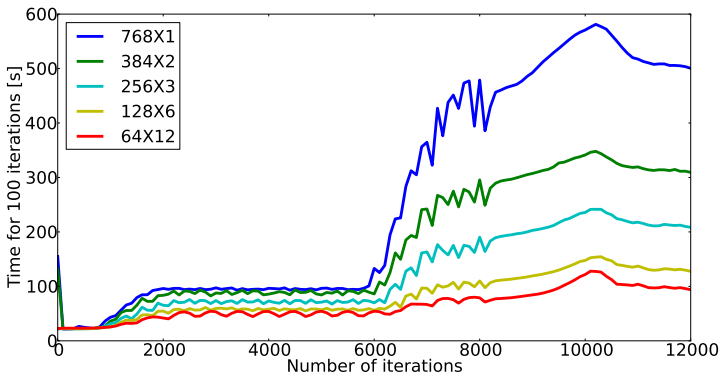


PIC code are 'easily' parallelized using **domain decomposition**

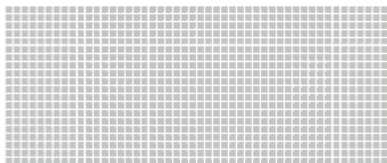




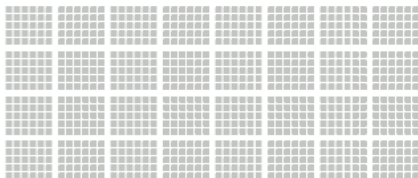
$MPI \times OpenMP$



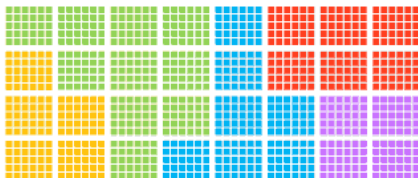
OpenMP dynamic scheduler is able to smooth the load but only at the node level.



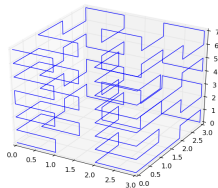
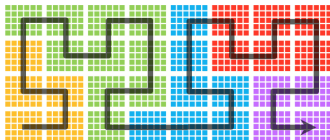
960 cells



32 patches

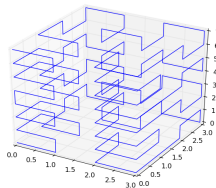
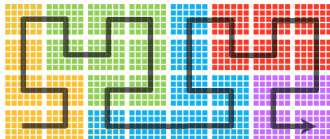


5 MPI regions



We need a policy to assign patches to MPI processes. To do so, patches are organized along a one dimensional **space-filling curve**.

- 1 Continuous curve which goes across all patches.
- 2 Each patch is visited only once.
- 3 Two consecutive patches are neighbours.
- 4 In addition we want compactness !

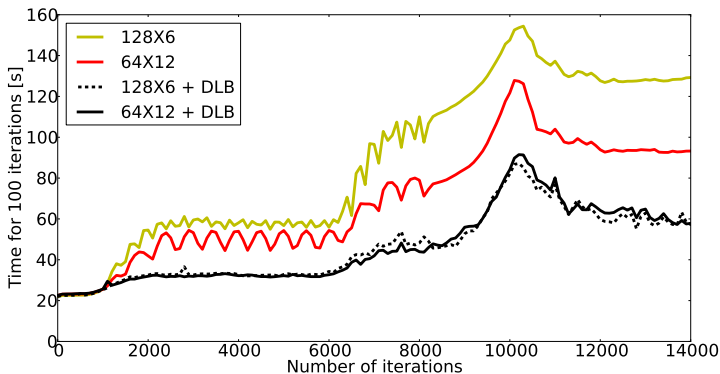


We need a policy to assign patches to MPI processes. To do so, patches are organized along a one dimensional **space-filling curve**.

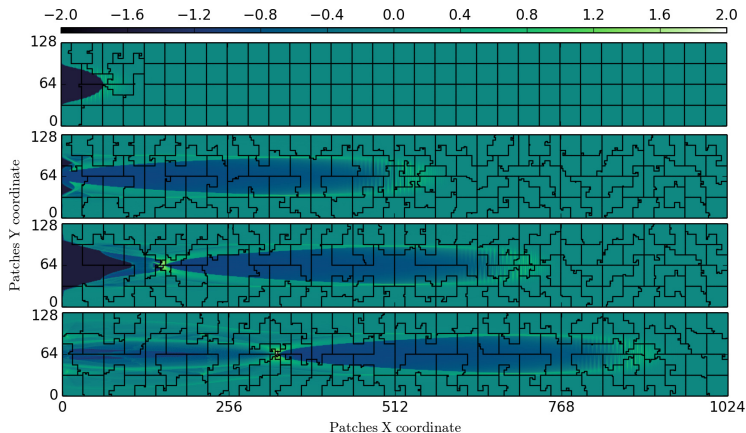
- 1 Continuous curve which goes across all patches.
- 2 Each patch is visited only once.
- 3 Two consecutive patches are neighbours.
- 4 In addition we want compactness !



MPI \times OpenMP

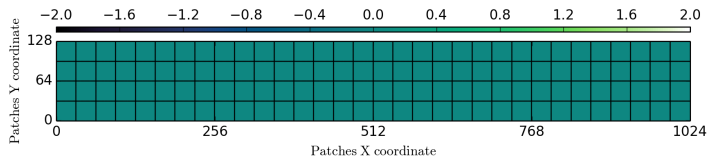


Yellow and red are copied from previous figure.



Color represents the local patch computational load imbalance

$$I_{\text{loc}} = \log_{10} (L_{\text{loc}}/L_{\text{av}})$$



Color represents the local patch computational load imbalance

$$I_{\text{loc}} = \log_{10} (L_{\text{loc}}/L_{\text{av}})$$