**Fourth Brazil-France Workshop**

On High Performance Computing and Scientific Data Management Driven by Highly Demanding Applications

15-18 September, Gramado, Brazil **2014**

# Emerging Challenges for EdgeCFD Simulations in Many-core Architectures

**Renato N. Elias**[1], José Camata[1], Albino Aveleda[1,] Alvaro Coutinho[1]

[1] High Performance Computing Center (NACAD)
{rnelias, camata, bino, alvaro}@nacad.ufrj.br
www.nacad.ufrj.br
Civil Engineering Department (PEC/COPPE)
www.pec.coppe.ufrj.br
Federal University of Rio de Janeiro (UFRJ)
www.ufrj.br

# *Outline*

- *Introduction & Motivation*

- *What's EdgeCFD*

  - *Main aspects and algorithms*

- *What's Xeon-Phi*

- *EdgeCFD on Xeon-Phi*

- *Test-case*

- *Results*

- *Closure and discussions*

NACAD

# *Introduction*

- ❑ *New HPC systems are getting*

  - – *Massive in slower cores = "Many-cores"*

  - – *Several parallelism and vectorization combinations/possibilities*

- ❑ *$1^{st}$ on top500.org (JUN14)*

  - – *Tianhe-2/China: 3.12M of cores (**Intel Xeon Phi**)*

  - – *33.86 Pflops in HPL*

NACAD

# *Questions/Concerns*

1. **Are applications ready for these new HPC systems?**

2. **How should we take advantage of them?**

3. **How do CoProcessors compare to "traditional" CPUs?**

    – *Speed?*

    – *Cost?*

    – *Complexity?*

4. **Are CoProcessors a viable solution for any kind of problem?**

NACAD

# *Introduction/Motivation:*

- ❑ ***It's an active and recent topic in HPC Community***

- ❑ ***Related works:***

  - *M. Vasquez et al @ BSC. "Xeon Phi Performance for HPC-based Computational Mechanics Codes", Partnership for Advanced Computing in Europe*

  - *Venetis et al @ NTUA/Greece and TUD/Germany. "Porting FEASTFLOW to the Intel Xeon Phi: Lessons Learned"*

  - *J. Waltz @ LANL. "Performance Analysis of a 3D Unstructured Mesh Hydrodynamics Code on Multi- and Many-Core Architectures to appear in IJNMF*

NACAD
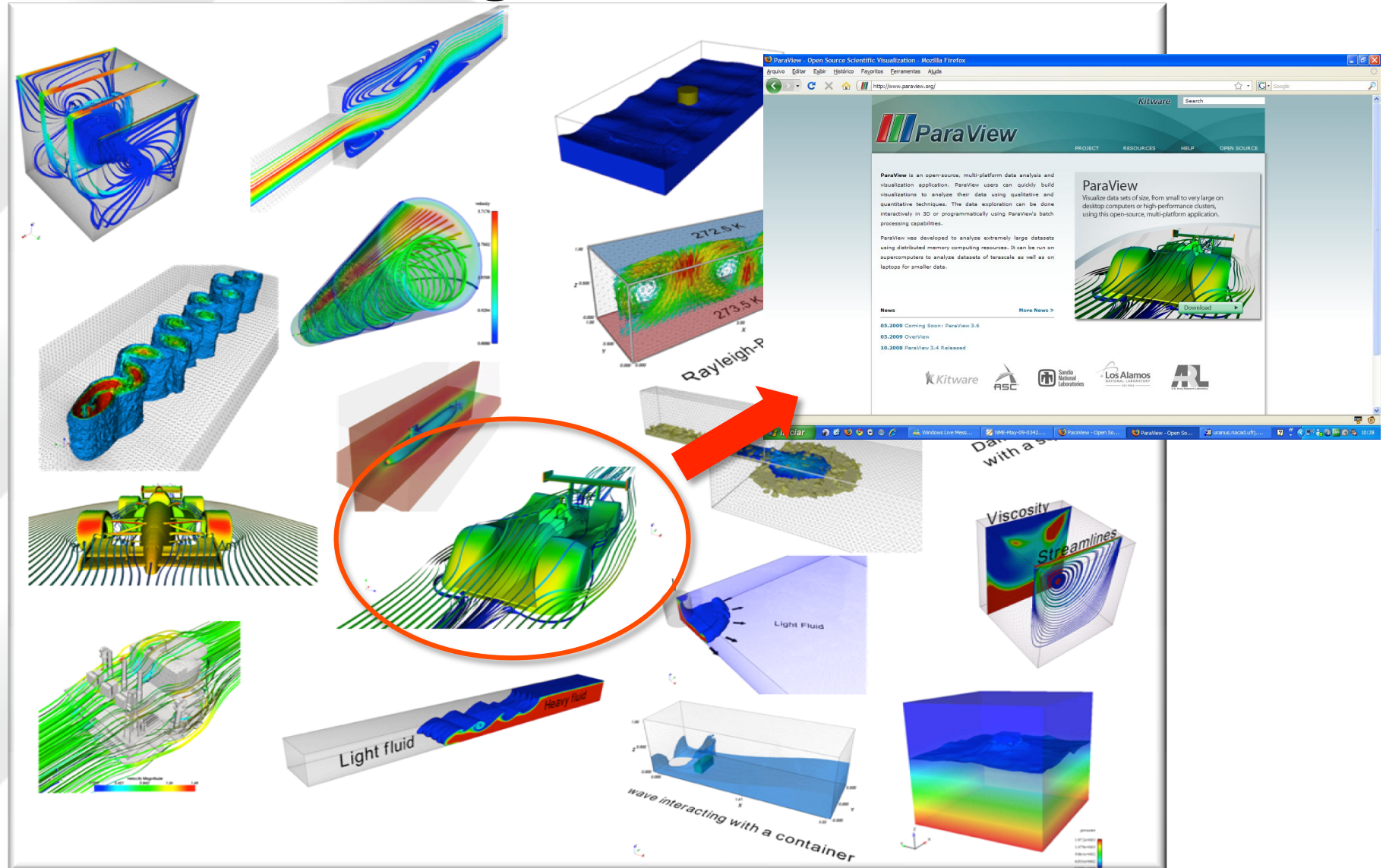
NACAD

# *What are we interested in?*

- ❑ *Prepare our software for these new architectures;*

- ❑ *EdgeCFD: A Stabilized Finite Element Simulator for Fluid Flow Problems.*

  - – *Inherit/incorporate results from theoretical and applied researches in our group*

  - – *Helps other projects as a simulation tool*

> *Edge: Data structure used in the software*
>
> *+*
>
> *CFD: Computational Fluid Dynamics*

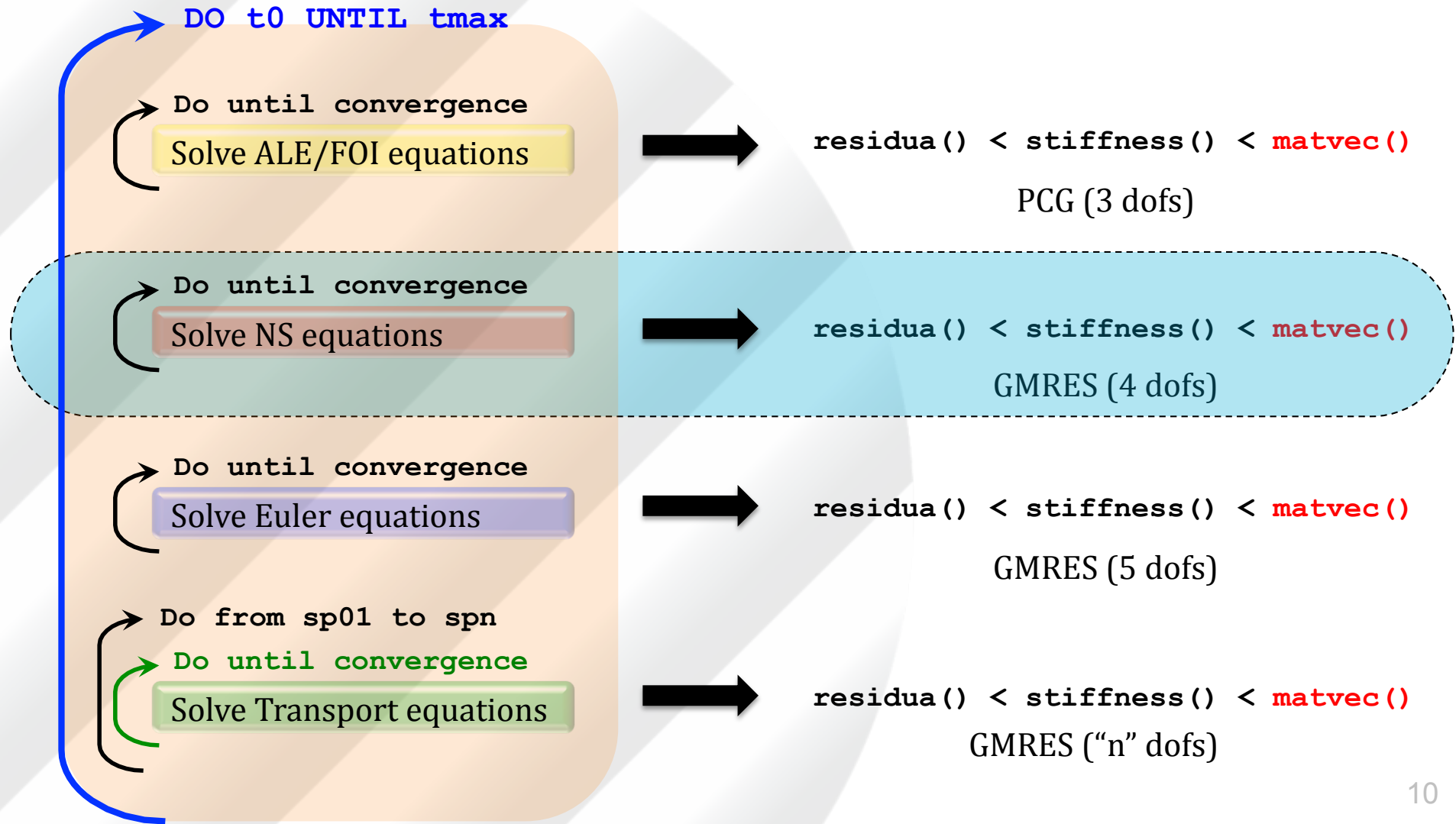# EdgeCFD in Action

NACAD

# Some "*Goodies* and *Baddies*"

❑ *Goodies:*

– *Built as a hybrid parallel code **from scratch***

– ***Highly portable:** *any* Fortran compiler is enough to compile…*

- *We extremely minimize external library dependencies;*

- *Only Metis and MPI is required;*

- *Build process is very simple for any platform (Mac Osx, Windows or *Nix)* ☺

– ***Efficient:** at least, Intel Vtune and TAU say that…* ☺

- *It's able to apply data ordering according to the different architecture.*

- *Employs fast non-linear and linear solvers*

- *Efficient data structure (by edges)*

- *Relies on good programming practices for HPC architectures*

❑ *Baddies*

– *No Open Source* ☹

– *No "fancy" **UI's*** ☹ *(pre-processing* → *GMSH and post-processing* → *ParaView)*

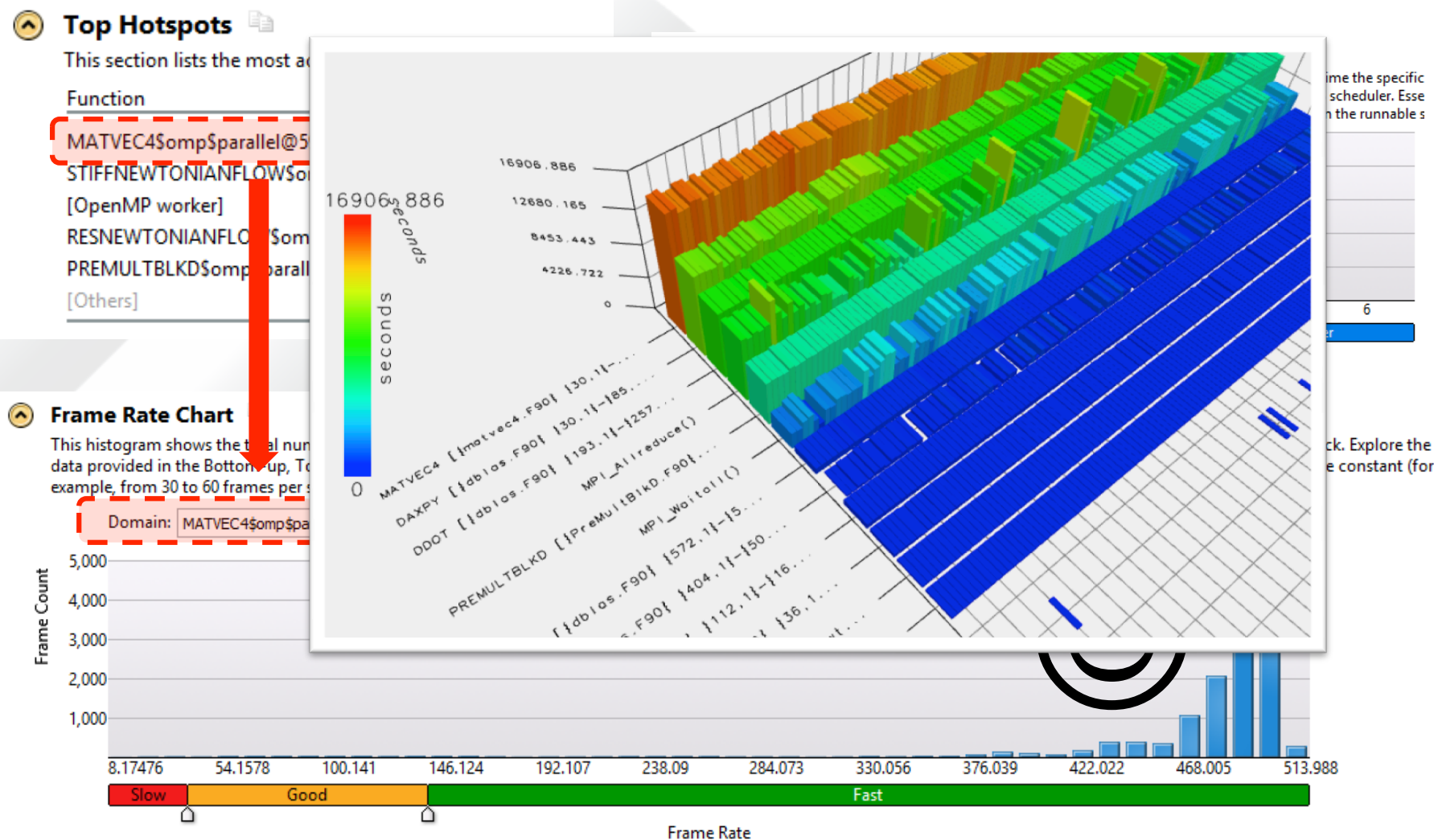– *Only **linear tetrahedra** elements are supported… (but any element can be dissassembled in edges)*

NACAD

# *Loops and Computational Costs...*

```
DO t0 UNTIL tmax
```

**Do until convergence**

Solve ALE/FOI equations

→ **residua() < stiffness() < matvec()**

PCG (3 dofs)

**Do until convergence**

Solve NS equations

→ **residua() < stiffness() < matvec()**

GMRES (4 dofs)

**Do until convergence**

Solve Euler equations

→ **residua() < stiffness() < matvec()**

GMRES (5 dofs)

```
Do from sp01 to spn
```
**Do until convergence**

Solve Transport equations

→ **residua() < stiffness() < matvec()**
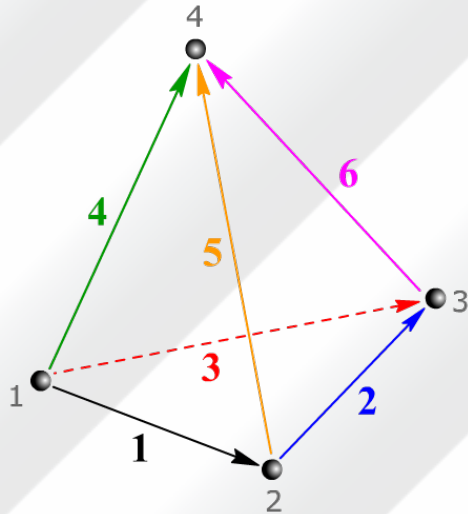
GMRES ("n" dofs)

*Dof = degrees of freedom*

# *Performance analysis*

*Typical incompressible flow simulation running in OpenMP*

NACAD

# EDS: Edges-by-Edge (cont.)

*"…algebraic disassembling of a tet into edges…"*

$$T_1^e = \begin{bmatrix} T_{11}^1 & T_{12}^1 & 0 & 0 \\ T_{21}^1 & T_{22}^1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \qquad T_2^e = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & T_{22}^2 & T_{23}^2 & 0 \\ 0 & T_{32}^2 & T_{33}^2 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$T_4^e = \begin{bmatrix} T_{11}^4 & 0 & 0 & T_{14}^4 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ T_{41}^4 & 0 & 0 & T_{44}^4 \end{bmatrix} \qquad T_3^e = \begin{bmatrix} T_{11}^3 & 0 & T_{13}^3 & 0 \\ 0 & 0 & 0 & 0 \\ T_{31}^3 & 0 & T_{33}^3 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$
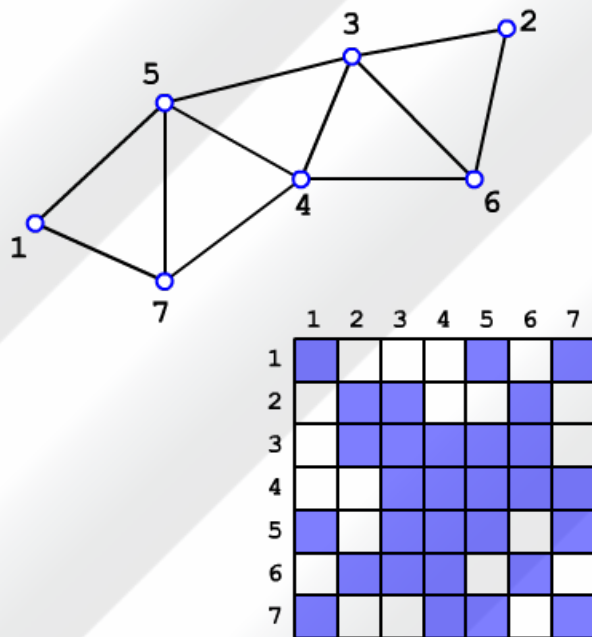
$$J^e = \begin{bmatrix} J_{11} & J_{12} & J_{13} & J_{14} \\ J_{21} & J_{22} & J_{23} & J_{24} \\ J_{31} & J_{32} & J_{33} & J_{34} \\ J_{41} & J_{42} & J_{43} & J_{44} \end{bmatrix}$$

$$J^e = T_1^e + T_2^e + T_3^e + T_4^e + T_5^e + T_6^e$$

$$T_5^e = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & T_{22}^5 & 0 & T_{24}^5 \\ 0 & 0 & 0 & 0 \\ 0 & T_{42}^5 & 0 & T_{44}^5 \end{bmatrix} \qquad T_6^e = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & T_{33}^6 & T_{34}^6 \\ 0 & 0 & T_{43}^6 & T_{44}^6 \end{bmatrix}$$
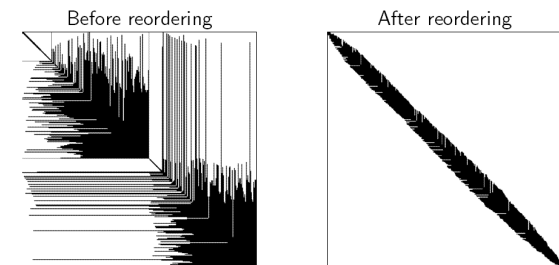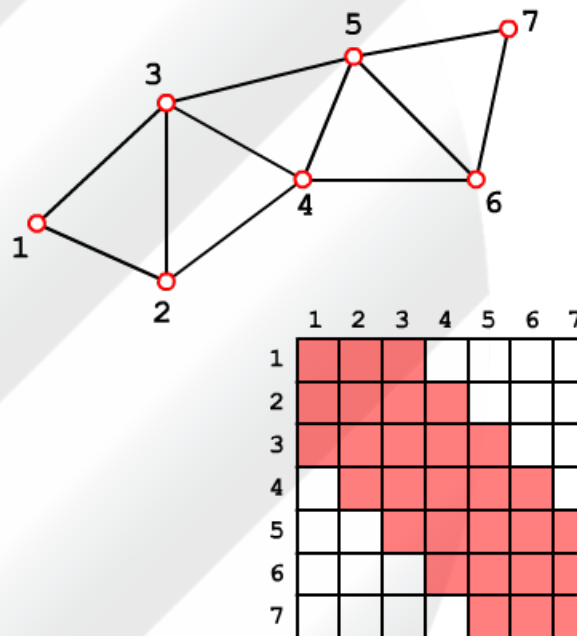
**In practice, nodal block-diagonals are also used to build preconditioner**

# Nodal Degrees-of-Freedom Reordering

- *Matrix profile is directly related to the node (degrees-of-freedom) order*

- *Reordering algorithms such as Reverse Cuthill Mckee reorder the unknowns to reduce bandwidth (profile)*
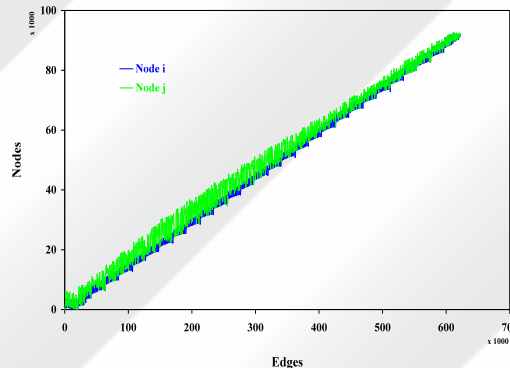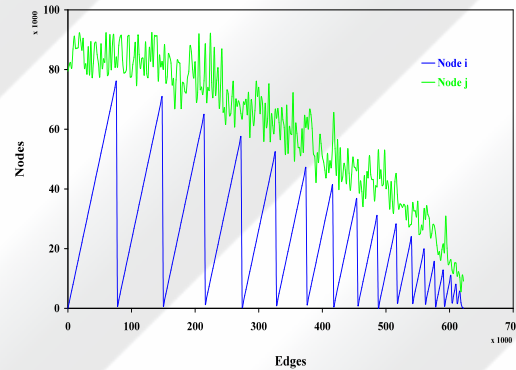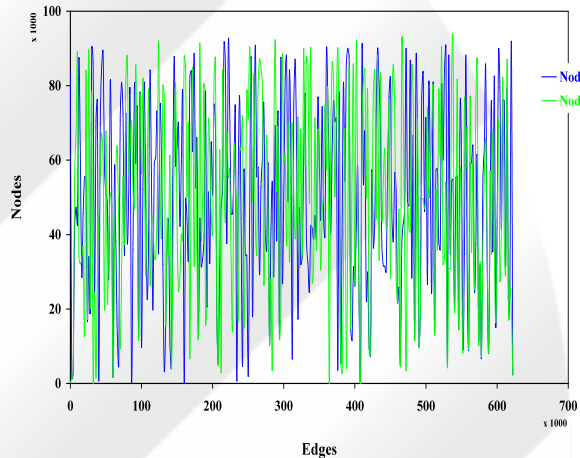
- *NP-complete graph problem*



real problem

basic example

NACAD

# *Memory Access Patterns*

**Reordering for minimizing i/a**

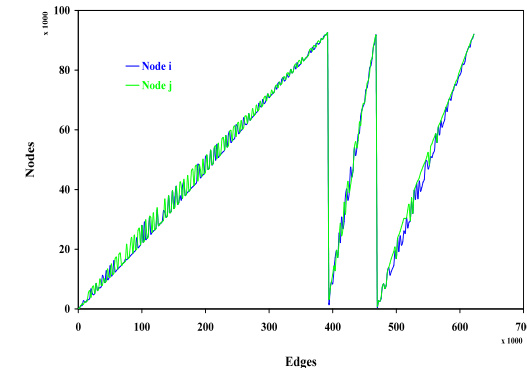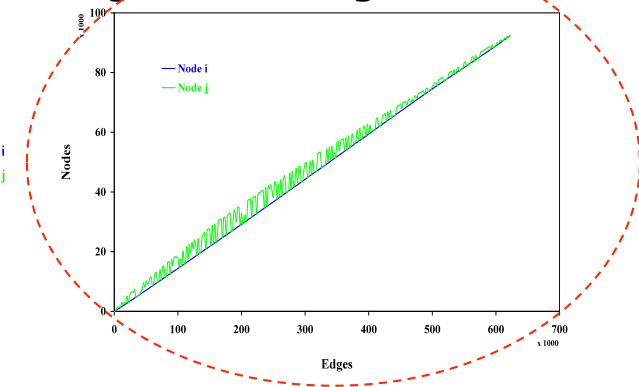**Improving data locality:   RCM vertex reordering and edge reordering in ascending vertex order**



**Improving data locality with NO memory dependencies**

**Original edge order**
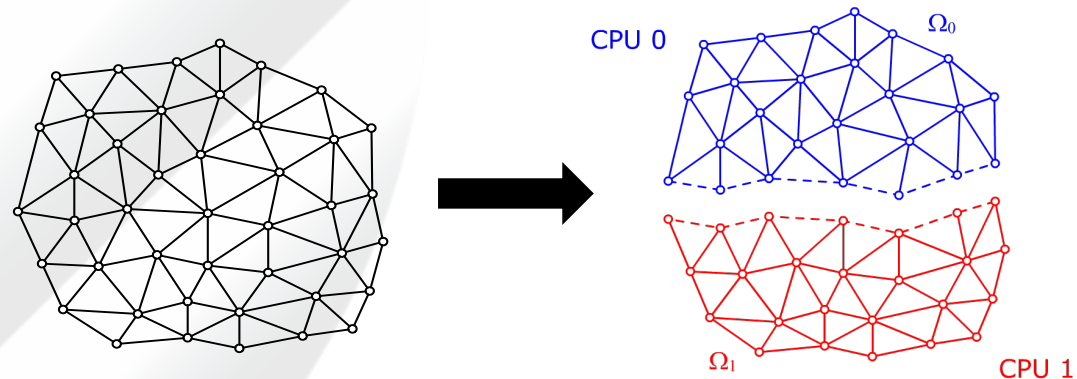
**Data reordering schemes provided by EdgePack®**

**Improving data locality for each superedge type. Superedge improves the use of CPU registers**

NACAD
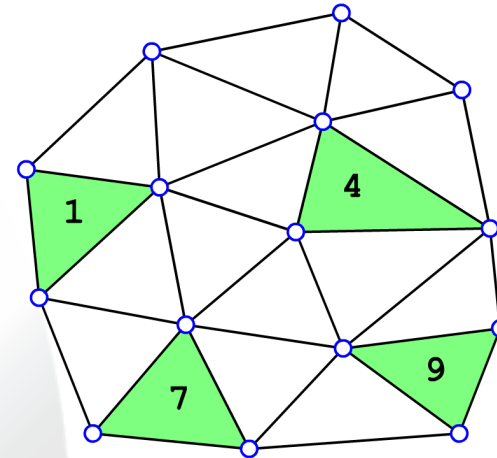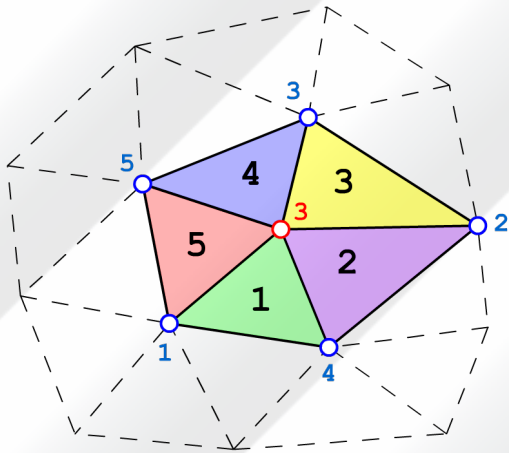
# Distributed Memory Parallelism

□ **Standard Approach:**

– *Takes a mesh and give it to a partitioner (**Metis**)*

– *Once partitioned, **keeps the data parallel***

– *Global IDs → **Local IDs***

– *Each process records its own files for **parallel I/O***

– *Shared information is synchronized by **p2p non-blocking communications***

– *No ghost/halo information employed*

NACAD

# *Shared Memory Parallelism*

❑ *Standard blocked loops to remove memory dependency*



```
!$OMP PARALLEL DO
do i=1,nel
    ! retrieve element nodes
    x(no) = x(no) + a
enddo
!$OMP END PARALLEL DO
```

```
ielm = 0
do icor = 1, ncores
    nvec = ielblk(icor)
!$OMP  DO
    do i = ielm+1, ielm+nvec
        ! Retrieve element nodes
        x(no) = x(no) + a
    enddo
!$OMP END DO
    ielm = ielm+nvec
enddo
```

NACAD

# *Hybrid Matrix-Vector Product*

**Edge-by-Edge**

```
iside = 0
DO iblk = 1, nedblk
nvec  = ia_edblk(iblk)
!dir$ ivdep
!$OMP DO
    DO ka = iside+1, iside+nvec, 1
       ...MATVEC computations...
    ENDDO
!$OMP END DO
ENDDO
```
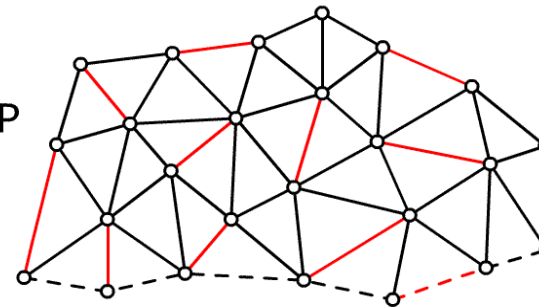
```
...over interface nodes...
#ifdef MPICODE
call MPI_AllReduce
#endif
```

OpenMP

OpenMP

NACAD

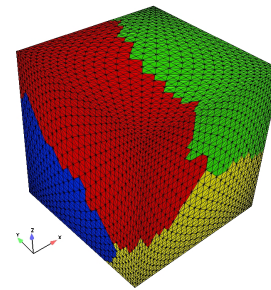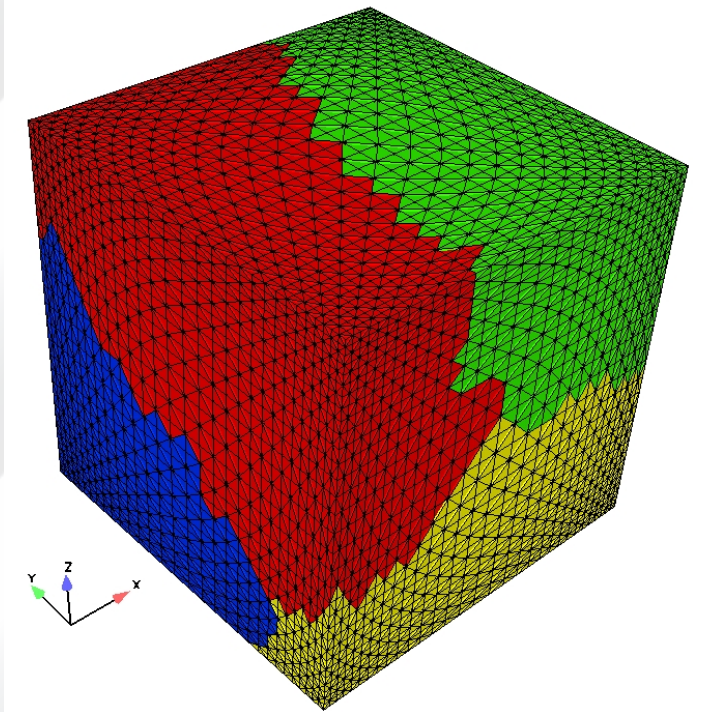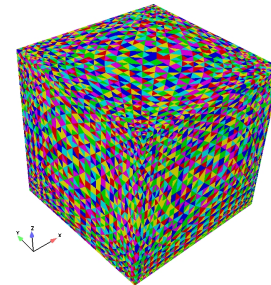# *Data Structure for Hybrid Paralelism*



**Metis partitioning**

Distributed (or shared) memory with MPI

**Mesh coloring**

Vectorization and/or shared memory parallelism (OpenMP)

**Data partitioning + Mesh coloring**

NACAD

# Intel Xeon Phi

❑ *Intel's Many-Integrated-Core (MIC) architecture*

❑ *In few words:*

– *"A PCI express board with 57-61 cores supporting up to 4 threads each (228-244 threads per board)*

❑ *Dense and simplified collection of (old) processors*

– *Based on Pentium P54c x86 architecture, with 64-bits, vector capabilities and cache coherent.*

– *It was made to support existing HPC applications*

❑ *Each board runs a simplified Linux*

– *Has ip address and is fully functional. In other words, it's a computing node inside a computing node*

NACAD

# *EdgeCDF on Xeon-Phi*

- ❑ *As simple as call the compiler with a new compilation flag (`-mmic`)* ☺

- ❑ *Ok, but does it scale? Let's see on next slides...*

- ❑ *Why not GPGPU's?*

  - – *It would require **deep structural code changes** (data structure, reordering schemes, etc...)*

  - – *EdgeCFD is not focused on specific platforms. We seek for a **good and overall** performance to keep portability!*

  - – *Planning some tests with PGI CUDA Fortran compiler @ CRAY*

NACAD

# *Xeon-Phi Execution Models*

# *EdgeCFD Tests on Xeon Phi*

- ❑ *"Blind"/Naive test*
    - – *Upload, compile and run! Nothing changed and/or tuned besides compiler optimizations (-O3)*

- ❑ *Execution Models*
    - – *Host only*
    - – *Mic only*

- ❑ *Parallel Models*
    - – *OpenMP and MPI on Host*
    - – *OpenMP and MPI on Mic*
    - – *Hybrid on Mic*

- ❑ ***Only strong scalability!***
    - – *How solution time behaves as the number of cores increase for a fixed effort/ problem.*

NACAD

# *Test Problem*

❑ ***Cavity Flow Problem (Incompressible Flow)***

- – **50 Time steps**

- – **Edges: 133,849**

- – **Tets.: 108,104**

- – **Nodes: 20,589**





*NOTE*

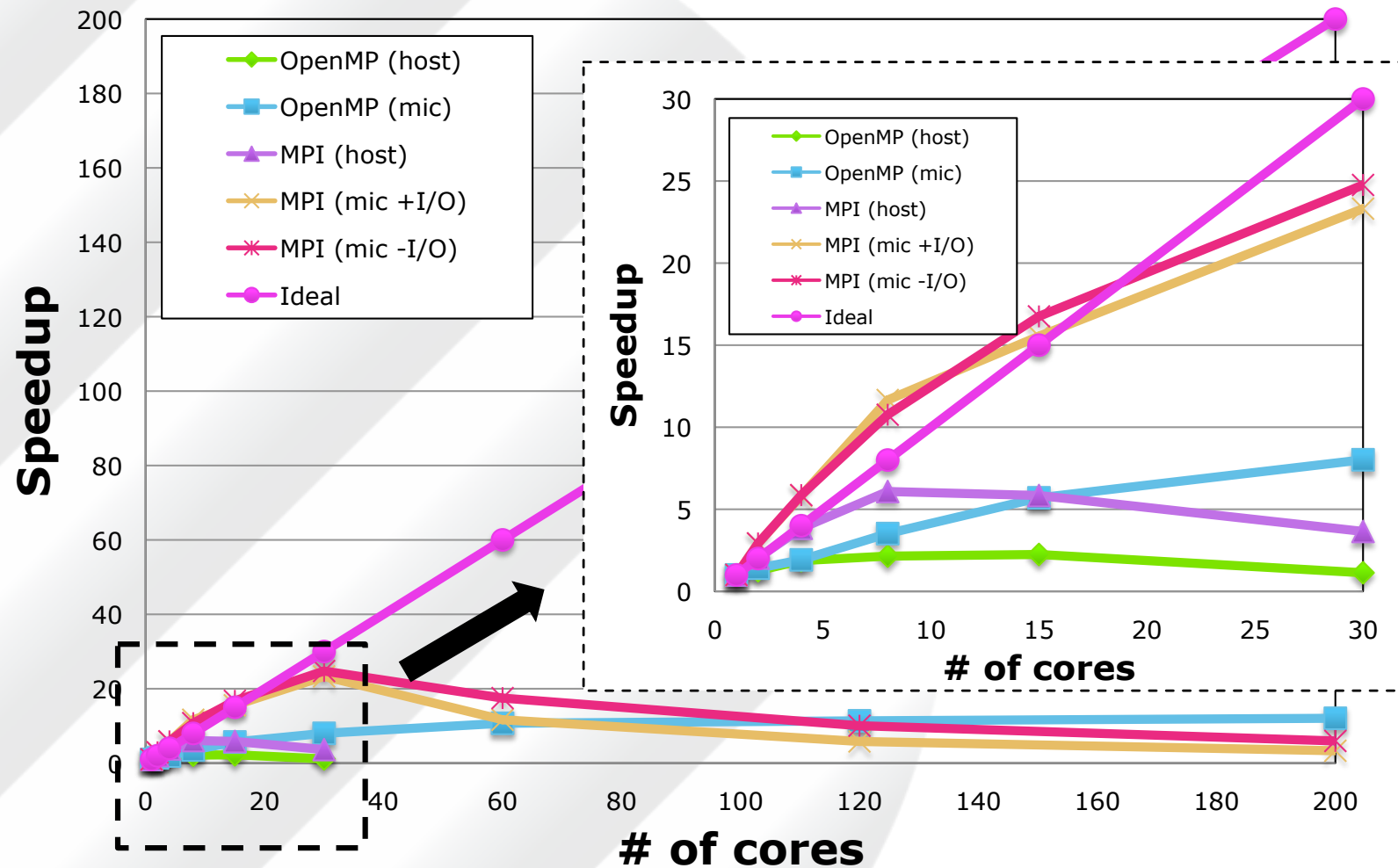- *It's a small problem! Could be easily ran on a laptop!*

*BUT...*

- *Can be considered as a MPI partition of a bigger problem placed on a computing node...* ***we're just trying to measure Xeon-Phi performance alone***

NACAD

# *Hardware*

- ❑ ***HOST: Dell server ( ~$4.2k @ Amazon 2T14)***

  - *Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz*

  - *2 x 8 cores, 2 threads/core, up to 32 threads or **64 hyperthreaded***

  - *64GB memory*

- ❑ ***Xeon Phi 3120A CoProcessor (~ $1.55k @ Amazon 2T14)***

  - *6GB memory (12 memory channels @ 240GB/s)*

  - *28,5 MB cache memory*

  - *1.053 GHz Clockspeed*

  - *57 Cores / 4 threads per core = up to 228 Threads*

> *But…, hey, Xeon-Phi does not live alone*
> *It needs a host anyway… ☹*

NACAD

# *SpeedUp*

# *Host x Mic Performance*

NACAD

# *Money x Time Comparison*



Host: ~$4.2k @ Amazon 2T14
Xeon-Phi: ~$1.55 @ Amazon 2T14
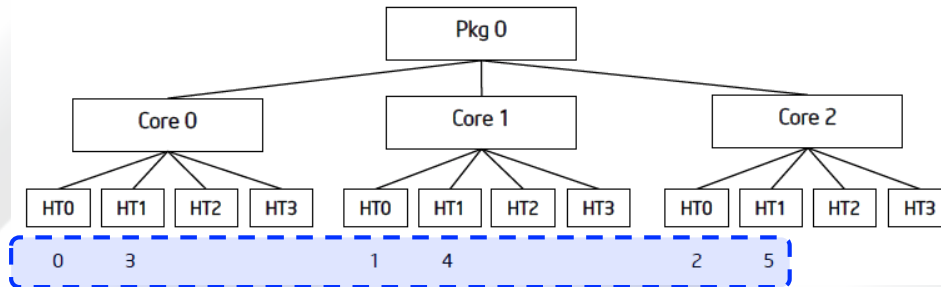
NOTE: Mics need a host but it was considered alone in this comparison

NACAD

# *Thread Affinity on Xeon-Phi*

*Scatter*



*Balanced*

*Compact*

NACAD

# How does thread affinity affect performance?

# How does block size affect EdgeCFD Performance

❑ *Block sizes => Bigger blocks means denser inner loops and less barrier synchronization*

**Outer/Inner colored loops
(compact affinity)**



Legend:
- 14/8192
- 27/4096
- 53/2048
- 106/1024
- 212/512
- 1690/64

X-axis: Number of threads
Y-axis: Time (sec)

```
iside = 0
DO iblk = 1, nedblk
nvec  = ia_edblk(iblk)
!$OMP DO
    DO ka = iside+1, iside+nvec, 1
        ...MATVEC computations...
    ENDDO
!$OMP END DO  ← BARRIER!!!!
iside = iside + nvec
ENDDO
```

NACAD

# Hybrid Performance on Xeon-Phi

- *Parallel I/O cost is inversely proportional to the number of MPI process;*

- *60 threads trying to write on the same disk at once* ☹

- *Solution: CoProcessing could help to overcome this problem , but computational mechanics software still needs to save some data on disk...*



Time (sec) / MPI processes

Legend: ■ 1 thread  ■ 2 threads  ■ 4 threads

NACAD

# *Summary and Conclusions (1/3)*

| | Best @ Host | | Best @ Mic | |
|---|---|---|---|---|
| Overall | 24.0 | MPI / 8 procs | 113.0 | MPI / 30 procs -I/O |
| OpenMP | 65.0 | 15 threads | 232.0 | 200 threads |
| MPI | 24.0 | 8 procs | 113.0 | 30 procs |
| cost x time | 1.0 | $ seconds | 5.0 | $ seconds |

❑ *Xeon-Phi could be **5x cheaper or 5x faster** from this comparison*

  – ***Cost-wise**, for our application, **regular processors are the best investment** (yet!)*

❑ ***New algorithms** for improving many-cores parallelism are necessary **for unstructured methods** running on newer architectures based on Mics;*

  – *Or we should leave Mics for other more specific applications (codes with predictable memory access pattern...)*

NACAD

# *Conclusions (2/3)*

- ❑ ***Xeon-Phi is indeed a flexible*** *alternative to get into microprocessors world (but **not the most efficient** at first... ☹)*

  - – *Microprocessors pose other options and complexity about parallel and vectorization combinations;*

  - – *Flexibility comes with a (high) price.*

- ❑ *Some **parameters** on EdgeCFD may affect **OpenMP performance***

  - – *outer/inner blocked loops ratio;*

  - – *Memory access pattern (not deeply investigated yet...);*

NACAD

# *Conclusions (3/3)*

- ❑ ***Thread affinity*** *is also important when squeezing the best* ***performance on Xeon-Phi****.*

- ❑ *Xeon-Phi, although support MPI* ***natively****, it must be used with care due to I/O costs...*

  - – *Offloads to host might mitigate this issue...*

- ❑ *In very few words, the* ***main issues using Xeon-Phi natively*** *are:*

  - – *Memory access pattern for unstructured mesh methods;*

  - – *Thread synchronization in OpenMP;*

  - – *Disk use in MPI.*

NACAD

# *Trying to Answer Our Questions:*

- ❑ *Are old applications ready for these new HPC systems?*

  – *Definitely not!*

- ❑ *How could/should we take advantage of them?*

  – *We don't know (yet). But we already have some guesses….*

- ❑ *How do CoProcessors compare to "traditional" CPUs?*

  – *They are (much) slower and expensive as well*

- ❑ *Are CoProcessors a viable solution for any kind of problem?*

  – *Our results pointed that they're more suitable for applications with a predictable memory access pattern;*

  – *We have room for improvements in new algorithms here…*

NACAD

# *Thanks for your attention!*

**NACAD**