# Selected Intel Xeon Phi Accelerated Libraries developed at IT4Innovations

**Lubomír Říha**      ESPRESO Library
**Michal Merta**      BEM4I Library
**Milan Jaroš**       CyclesPhi Render

# Approaches to Accelerate Application using Xeon Phi (KNC)

## ESPRESO

- Acceleration of the FETI solver
- FEM produces sparse matrices – new method converts sparse matrices to dense
- Memory bound code – utilizes high bandwidth of the KNC memory
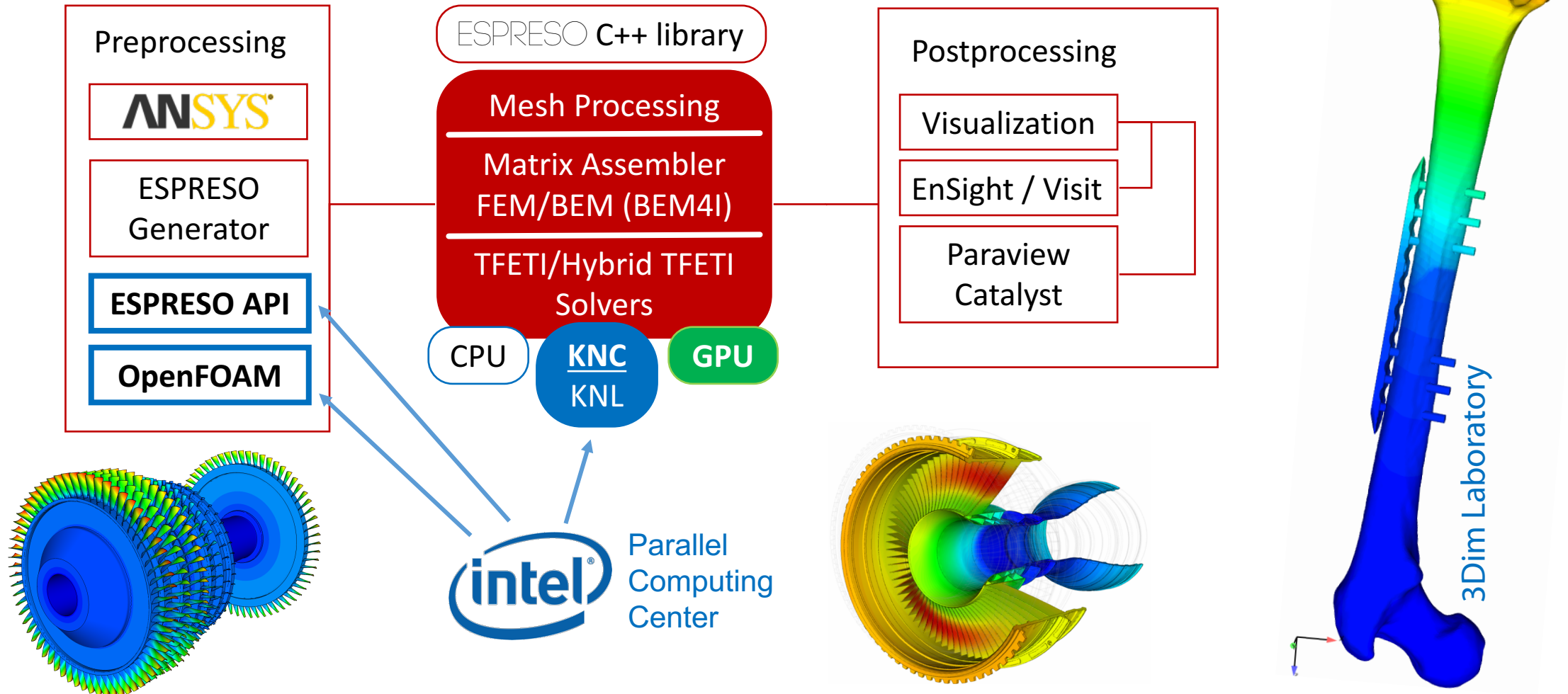
## BEM4I

- Acceleration of the BEM matrix assembler
- Compute bound code – relies on high performance (wide) vector units of the KNC core

## CyclesPhi

- Parallelization of the complex third party code using OpenMP and Offload pragmas to utilize KNC
- Developed MPI parallelization to support HPC environment
  - Enables real-time raytracing in interactive mode

# ESPRESO and Large Problems

## Weak Scalability Test

Up to **223 billion** DOF on 17576 Compute Nodes (281 216 cores)
Heat transfer

### 4th in TOP500 LIST

**18,688** AMD Opteron 6274 16-core CPUs
18,688 Nvidia Tesla K20X GPUs

2.7 million core hours dedicated to:
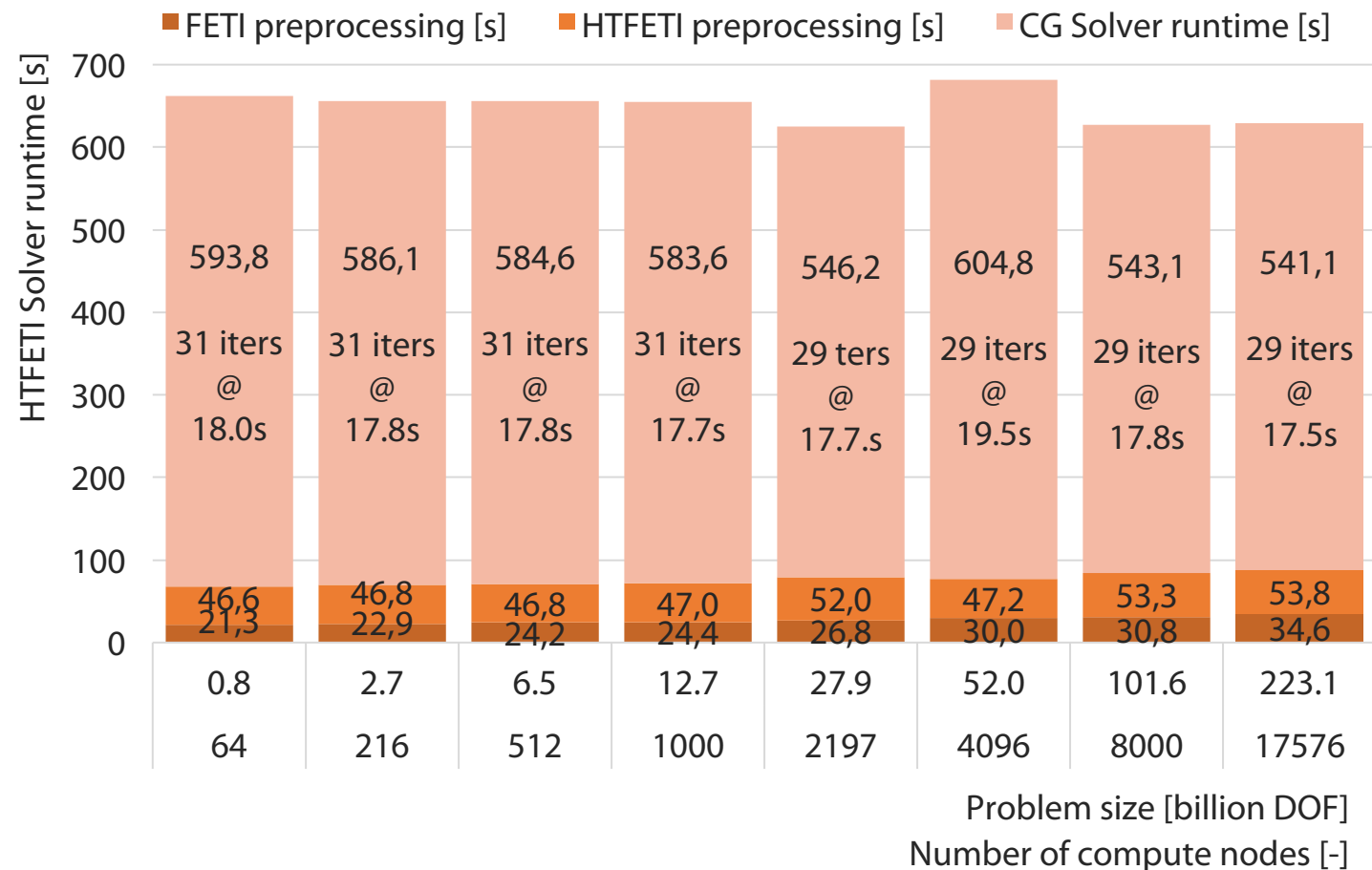- scalability optimization of ESPRESO
- optimization of GPU accelerated version for large scale problems

**OAK RIDGE National Laboratory**



Legend: ■ FETI preprocessing [s]   ■ HTFETI preprocessing [s]   ■ CG Solver runtime [s]

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 593,8 | 586,1 | 584,6 | 583,6 | 546,2 | 604,8 | 543,1 | 541,1 |
| 31 iters @ 18.0s | 31 iters @ 17.8s | 31 iters @ 17.8s | 31 iters @ 17.7s | 29 ters @ 17.7.s | 29 iters @ 19.5s | 29 iters @ 17.8s | 29 iters @ 17.5s |
| 46,6 | 46,8 | 46,8 | 47,0 | 52,0 | 47,2 | 53,3 | 53,8 |
| 21,3 | 22,9 | 24,2 | 24,4 | 26,8 | 30,0 | 30,8 | 34,6 |
| 0.8 | 2.7 | 6.5 | 12.7 | 27.9 | 52.0 | 101.6 | 223.1 |
| 64 | 216 | 512 | 1000 | 2197 | 4096 | 8000 | 17576 |

HTFETI Solver runtime [s]
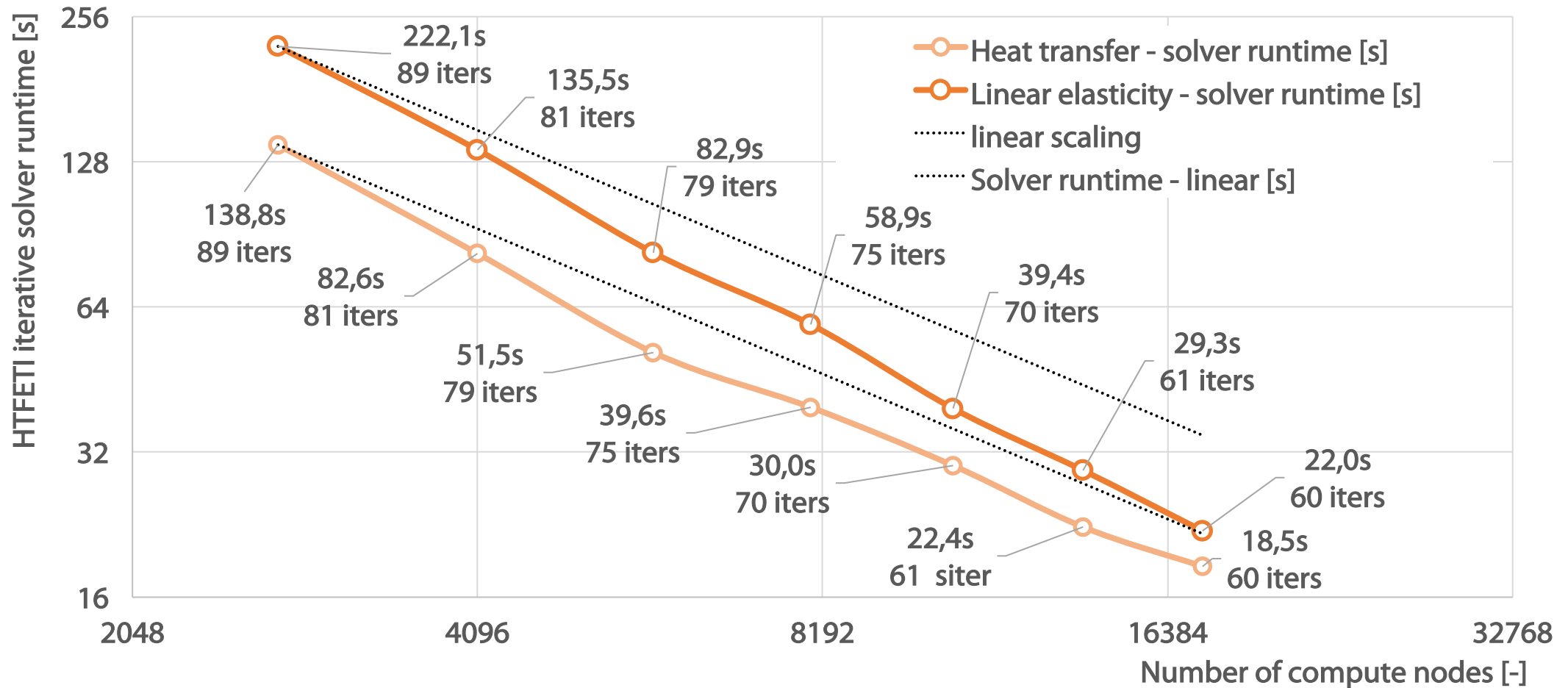
Problem size [billion DOF]
Number of compute nodes [-]

# Scalability of ESPRESO

**Heat transfer** 20 billion DOF on up to 17 576 Compute Nodes (281 216 cores)
**Linear elasticity** 11 billion DOF
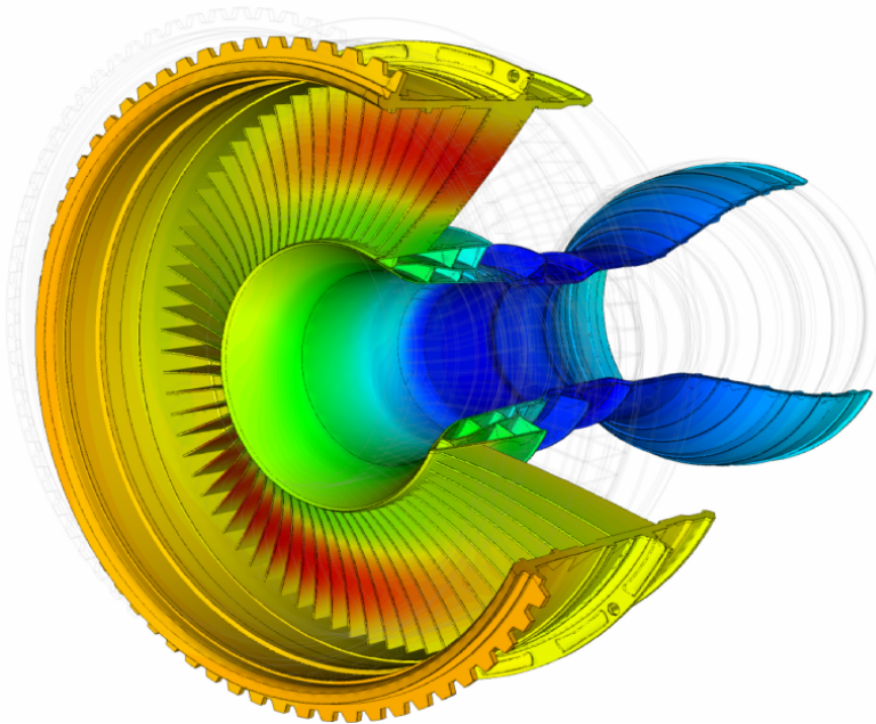
ORNL Titan 4[th] in TOP500 LIST

# Scalability for Real World Problems

300 million unknown - ANSYS Workbench real world problem

Linear elasticity – Hybrid FETI with Dirichlet preconditioner

IT4Innovations – SALOMON Supercomputer – only Haswell CPUs are used here
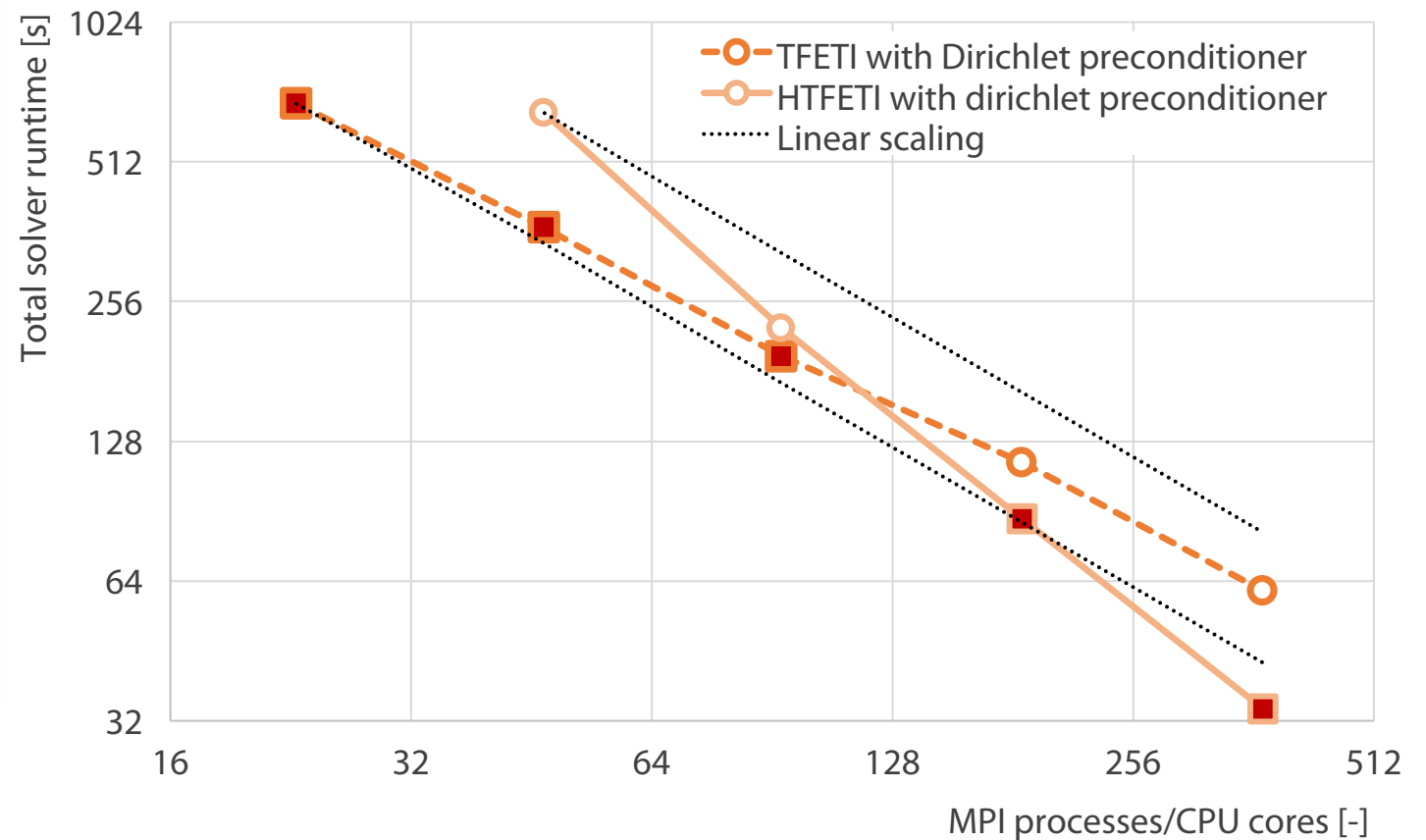


**Experiment setup:**
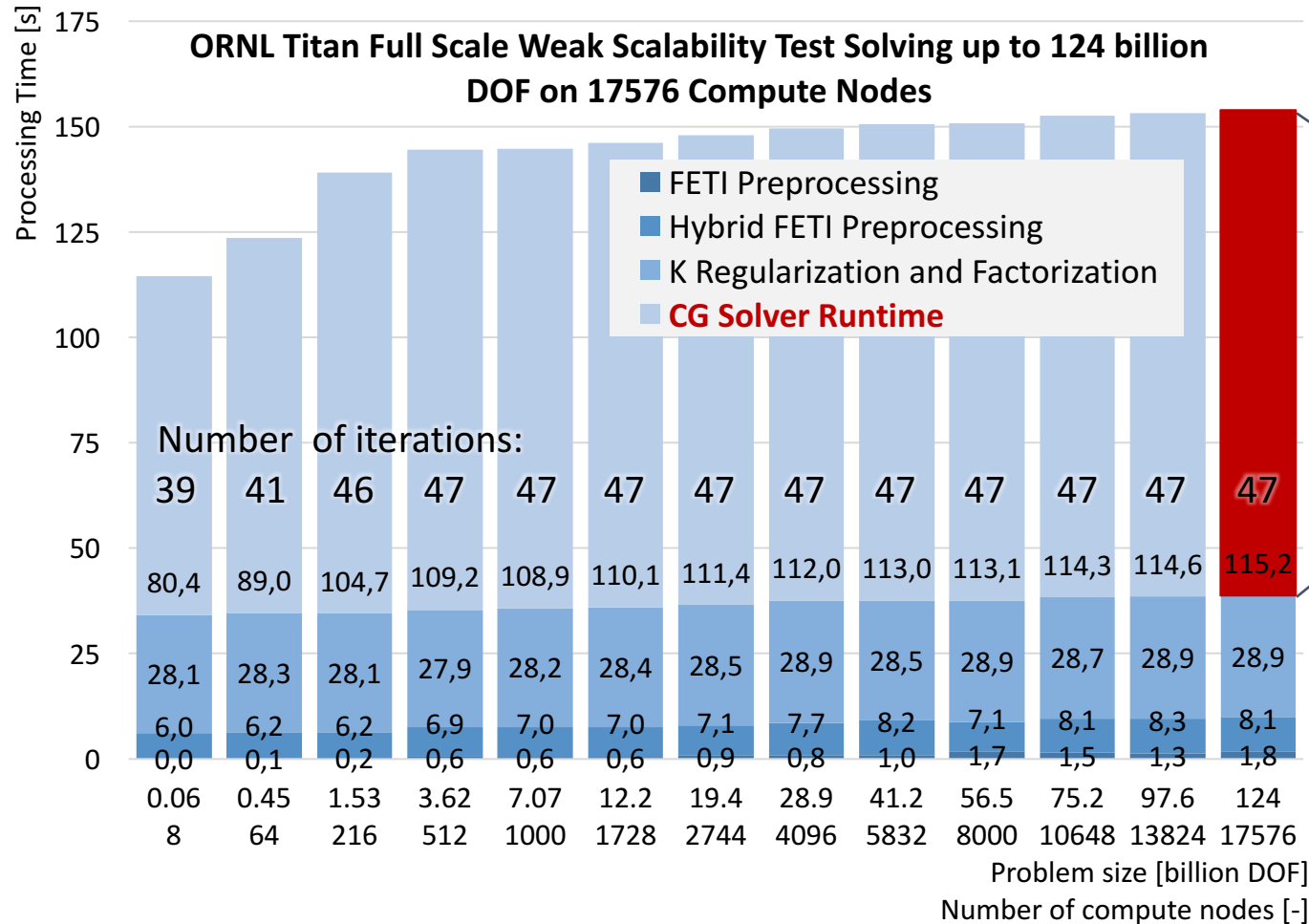Dirichlet preconditioner
regular CG solver
HTFETI method

# Motivation for Xeon Phi Acceleration

ORNL Titan Full Scale Weak Scalability Test Solving up to 124 billion DOF on 17576 Compute Nodes

This part of the solver is accelerated by the Intel Xeon Phi co-processor

**Legend:**
- FETI Preprocessing
- Hybrid FETI Preprocessing
- K Regularization and Factorization
- **CG Solver Runtime**

Number of iterations:

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 39 | 41 | 46 | 47 | 47 | 47 | 47 | 47 | 47 | 47 | 47 | 47 | 47 |
| 80,4 | 89,0 | 104,7 | 109,2 | 108,9 | 110,1 | 111,4 | 112,0 | 113,0 | 113,1 | 114,3 | 114,6 | 115,2 |
| 28,1 | 28,3 | 28,1 | 27,9 | 28,2 | 28,4 | 28,5 | 28,9 | 28,5 | 28,9 | 28,7 | 28,9 | 28,9 |
| 6,0 | 6,2 | 6,2 | 6,9 | 7,0 | 7,0 | 7,1 | 7,7 | 8,2 | 7,1 | 8,1 | 8,3 | 8,1 |
| 0,0 | 0,1 | 0,2 | 0,6 | 0,6 | 0,6 | 0,9 | 0,8 | 1,0 | 1,7 | 1,5 | 1,3 | 1,8 |
| 0.06 | 0.45 | 1.53 | 3.62 | 7.07 | 12.2 | 19.4 | 28.9 | 41.2 | 56.5 | 75.2 | 97.6 | 124 |
| 8 | 64 | 216 | 512 | 1000 | 1728 | 2744 | 4096 | 5832 | 8000 | 10648 | 13824 | 17576 |

Problem size [billion DOF]
Number of compute nodes [-]

Processing Time [s]

# Hardware Acceleration of FETI solvers
## Local Schur Complement (LSC) Method

1: $r_0 := b - Ax_0$; $u_0 := M^{-1}r_0$; $p_0 := u_0$
2: **for** $i = 0, \ldots, m-1$ **do**
3:    $s := Ap_i$
4:    $\alpha := \langle r_i, u_i \rangle / \langle s, p_i \rangle$
5:    $x_{i+1} := x_i + \alpha p_i$
6:    $r_{i+1} := r_i - \alpha s$
7:    $u_{i+1} := M^{-1}r_{i+1}$
8:    $\beta := \langle r_{i+1}, u_{i+1} \rangle / \langle r_i, u_i \rangle$
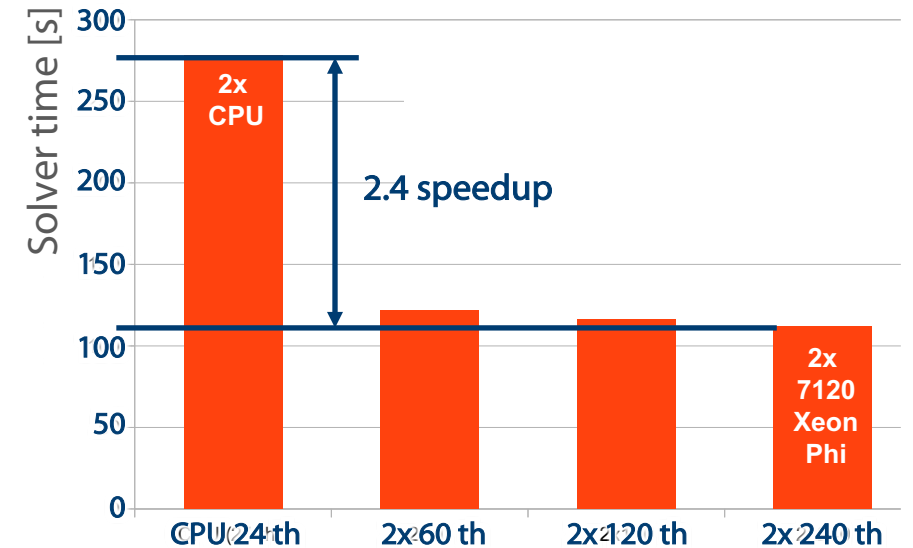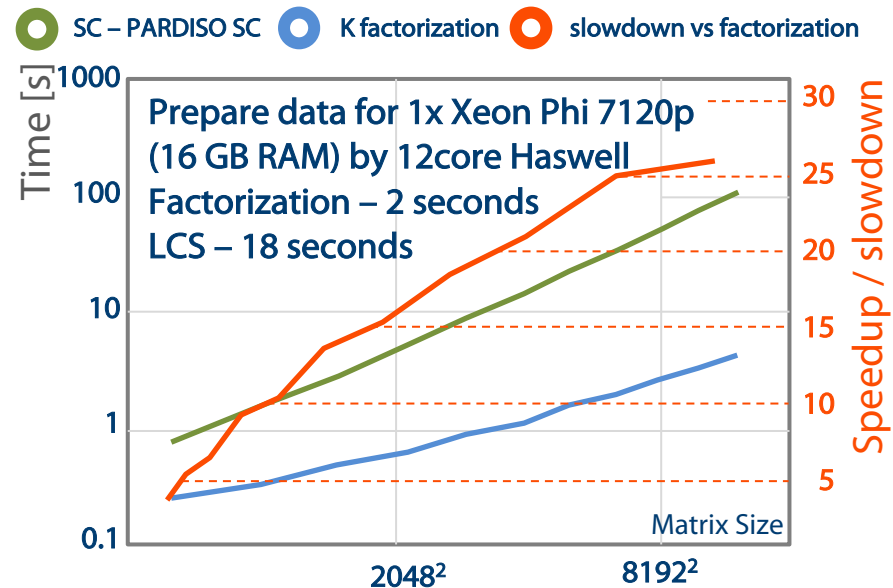9:    $p_{i+1} := u_{i+1} + \beta p_i$
10: **end for**

Pre-processing – K factorization
1.) $x = B_1^T \cdot \lambda$          - SpMV
2.) $y = K^{-1} \cdot x$          - solve
3.) $\lambda = B_1 \cdot y$          - SpMV
4.) stencil data exchange in $\lambda$
   - MPI – Send and Recv
   - OpenMP – shared mem. vec

Pre-processing - $S_c = B_1 K^{-1} B_1^T \to$ MIC
1.) $\lambda \to$ MIC   - PCIe transfer from CPU
2.) $\lambda = S_c \cdot \lambda$ - DGEMV, DSYMV on MIC
3.) $\lambda \leftarrow$ MIC   - PCIe transfer to CPU
4.) stencil data exchange in $\lambda$
   - MPI – Send and Recv
   - OpenMP – shared mem. vec

Using Schur complement we are still memory bounded, but fast Intel Xeon Phi memory can be now fully utilized due usage of dense matrices. This is the main factor that brings the speedup when compared to CPU which is two generations ahead.

Prepare data for 1x Xeon Phi 7120p (16 GB RAM) by 12core Haswell
Factorization – 2 seconds
LCS – 18 seconds

SC – PARDISO SC   K factorization   slowdown vs factorization



2.4 speedup

2x CPU

2x 7120 Xeon Phi

CPU (24 th)   2x 60 th   2x 120 th   2x 240 th

# Architecture Comparison for LSC Method

| | Number of cores | Peak floating point performance SP/DP [GFLOPS] | Theoretical Memory Bandwidth [GB/s] | Memory type |
|---|---|---|---|---|
| Intel Xeon E5-2680v3 | 12 | 960 / 480 | 68 | DDR4 |
| Intel Xeon Phi 7120p | 61 | 2420 / 1210 | ~~352~~ **180** | GDDR5 |
| Intel Xeon Phi 7210 | 64 | 5325 / 2662 | 102 / 400 | DDR4/MCDRAM |
| NVIDIA Tesla K80 (1 chip) | 1597 SP / 533 DP | 4365 / 1455 | 240 | GDDR5 |
| NVIDIA Tesla K20X (Titan) | 2688 SP / 896 DP | 3950 / 1310 | 250 | GDDR5 |
| NVIDIA Tesla P100 (PCIe) | 3584 SP / 1792 DP | 8345 / 4217 | 720 | HBM2 |

| Method | Problem size and decomposition | | Solve time[s] | Speedup by LSC method for solve | | | | |
|---|---|---|---|---|---|---|---|---|
| | number of subdomains [-] | Subdomain size [DOF] | PARDISO CPU only | CPU only | KNC & CPU* | KNL only | K80 & CPU* | P100 & CPU* |
| Heat Transfer problem | | | | | | | | |
| TFETI | 512 | 6859 | 21.4 | 1.3 | **2.1** | 2.2 | 1.5 | 4.0 |
| HTFETI | 512 | 6859 | 26.0 | 1.2 | **2.4** | 2.5 | 1.6 | 3.5 |
| Linear elasticity problem | | | | | | | | |
| TFETI | 512 | 3993 | 20.4 | 1.1 | **1.8** | 2.5 | 1.3 | 3.5 |
| HTFETI | 512 | 3993 | 21.2 | 0.9 | **1.8** | 2.4 | 1.2 | 2.7 |

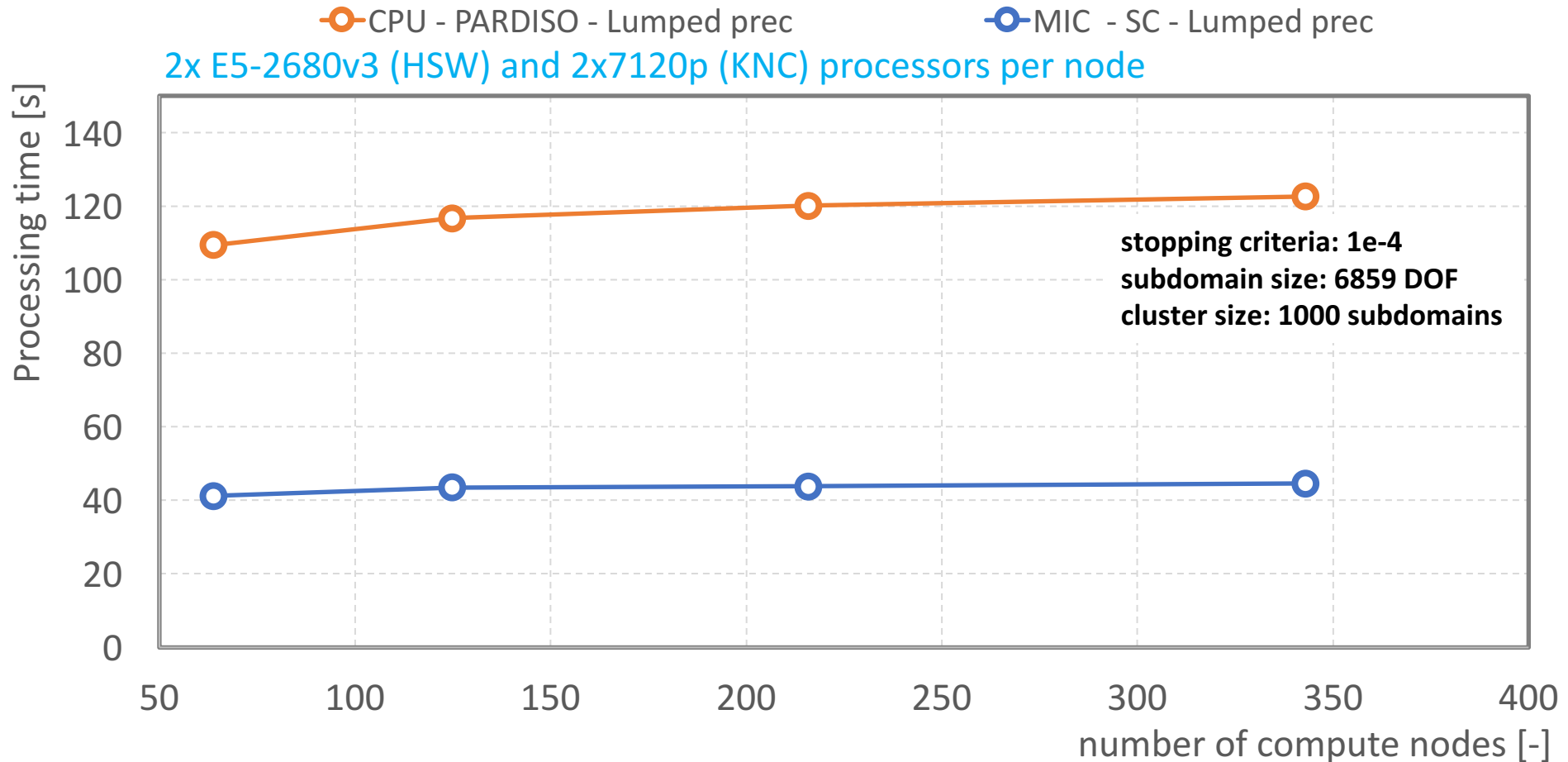**\* Note:** CPU does not process LSCs

# Large Scale Tests on Salomon

330 – 1776 million DOF  Hybrid FETI CG Solver Runtime

Heat Transfer – CG Solver Runtime w. Lumped Prec.

IT4Innovations Salomon Supercomputer

Speedup

⬤🔴  2.4



2x E5-2680v3 (HSW) and 2x7120p (KNC) processors per node

Legend:
- CPU - PARDISO - Lumped prec
- MIC - SC - Lumped prec

stopping criteria: 1e-4
subdomain size: 6859 DOF
cluster size: 1000 subdomains

y-axis: Processing time [s]
x-axis: number of compute nodes [-]

# Solving Real World Problems – Challenges

Input data in from **ELMER**, **OpenFOAM** or **Ansys Workbench**

Worse convergence – requires high number of iterations
- **requires efficient Dirichlet preconditioner**
- Full orthogonal CG solver with restarts for large coefficient jumps
- HTFETI method with clusterization by local kernels

Specific requirements for "smaller problems" up to 500 million unknown
- **multiple MPI processes share single node**
  - **resource sharing control for accelerators**
- **efficient use of all resources**

**Required enhancements in Xeon Phi Acceleration**

1. **Acceleration of Dirichlet preconditioner**

2. **Load balancing between CPU and Xeon Phi**

3. **Efficient sharing of one Xeon Phi by multiple MPI processes**

# Solving Real World Problems on Xeon Phi

## 1.) Acceleration of the Dirichlet Preconditioner

| | 2 MPI per MIC | | |
|---|---|---|---|
| | CPU | MIC | Speedup |
| MKL PARDISO / LSC action time [s] | 0.179 | 0.087 | 2.1 |
| Dirichlet preconditioner action time [s] | 0.105 | 0.060 | 1.8 |

Tire Rim benchmark with Dirichlet Preconditioner
Linear elasticity – 120 million unknown

## 2.) Dynamic Load balancing between CPU and accelerator

Minimizes runtime by automatically splitting the workload



64 Salomon's compute nodes
128 7120p accelerators
~2 million unknowns per node

## 3.) Efficient sharing of one Xeon Phi by multiple MPI processes

2 MPI ranks per Xeon Phi

| | CPU with PARDISO solver | | | MIC with LSC method | | | |
|---|---|---|---|---|---|---|---|
| # of subdomains per MPI rank [-] | Iteration time [s] | CG solver runtime [s] | number of iterations [-] | Iteration time [s] | CG solver runtime [s] | number of iterations [-] | speedup [-] |
| 128 | 0.298 | 88.0 | 295 | 0.146 | 42.9 | 294 | 2.0 |

# Transient Problem on Xeon Phi

Cube benchmark

Heat Transfer – 5 million unknowns per node

**A good use case for LSC and Xeon Phi acceleration**

1x Salomon's compute node
2x 7120p accelerators
~2.0 speedup which includes LSC and no. prec.

Benchmark execution:
`mpirun -n 2 ./espreso -c espreso-phi-bench.ecf 0  1 1 2  6 6 10  19 19 19`

Problem statistics:
Generate grid of hexahedrons
Coordinates loaded - total number of nodes: 5 051 950
Elements loaded - total number of elements: 4 938 480

| Time step | Transient solver stage | MIC LSC time [s] | CPU PARDISO time [s] | speedup [-] | savings [s] |
|---|---|---|---|---|---|
| preprocessing | FETI Preprocessing | 0.01 | 0.01 | 1.00 | 0.0 |
| preprocessing | Schur Complement asm. | 54.8 | 0.00 | --- | -54.8 |
| preprocessing | K  factorization | 0.00 | 2.9 | --- | 2.9 |
| 0 | CG Solver  runtime | 17.8 | 35.6 | 2.00 | 17.8 |
| 1 | RHS update | 0.01 | 0.01 | 1.0 | 0.0 |
| 1 | CG Solver  runtime | 17.2 | 34.7 | 2.02 | 17.5 |
| 2 | RHS update | 0.01 | 0.01 | 1.0 | 0.0 |
| 2 | CG Solver  runtime | 16.4 | 34.0 | 2.07 | 17.6 |
| 3 | RHS update | 0.01 | 0.01 | 1.0 | 0.0 |
| 3 | CG Solver  runtime | 16.1 | 33.5 | 2.08 | 17.4 |
| 4 | RHS update | 0.01 | 0.01 | 1.0 | 0.0 |
| 4 | CG Solver  runtime | 16.2 | 33.8 | 2.08 | 17.6 |
| | | | | | |
| 9 | RHS update | 0.01 | 0.01 | 1.0 | 0.0 |
| 9 | CG Solver runtime | 16.2 | 32.8 | 2.03 | 16.7 |
| ... | | | | | |

| 10 time steps | Total time | 219s | 339s | 1.5 | 120s |
| 100 time steps | Total time | 1694s | 3402s | 2.0 | 1708s |

# BEM4I Library

**Michal Merta**

# BEM4I library

## Parallel boundary element library



$$Vw(\boldsymbol{x}) = \frac{1}{2}u(\boldsymbol{x}) + Ku(\boldsymbol{x}) \quad \text{for } \boldsymbol{x} \in \partial\Omega$$

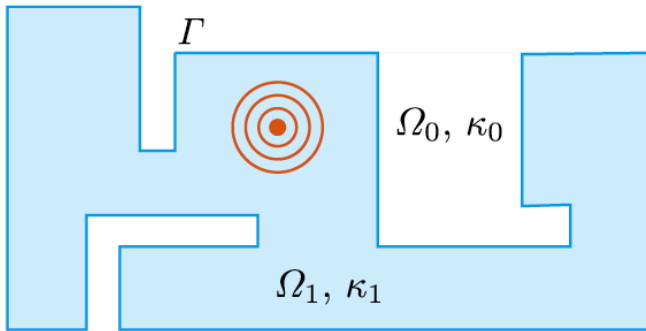$$\text{FEM} \quad \begin{cases} -\Delta u - \kappa^2 u = 0 & \text{in } \Omega \\ u = g & \text{on } \partial\Omega \end{cases}$$

- 3D boundary element solver in C++
- Laplace, Helmholtz, Lamé, wave equations
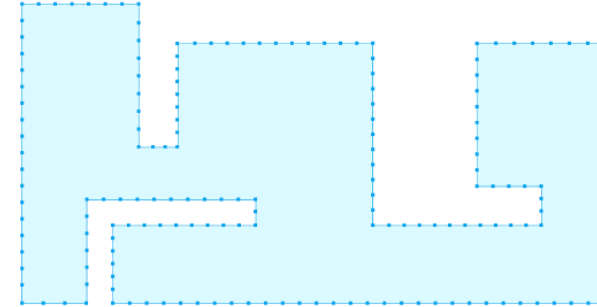- SIMD, OpenMP, MPI
- Offload to Intel Xeon Phi
- BEM4+ESPRESO = {BETI, FETI/BETI coupling}

# BEM4I library

## Parallel boundary element library



$$Vw(\boldsymbol{x}) = \frac{1}{2}u(\boldsymbol{x}) + Ku(\boldsymbol{x}) \quad \text{for } \boldsymbol{x} \in \partial\Omega$$

BEM

FEM

$$\begin{cases} -\Delta u - \kappa^2 u = 0 & \text{in } \Omega \\ u = g & \text{on } \partial\Omega \end{cases}$$

$\Gamma$

$\Omega_0, \kappa_0$

$\Omega_1, \kappa_1$

| Laplace | Lamé | Helmholtz | Wave |
|---|---|---|---|

Integrators and matrix assemblers

Common core: matrices, vectors, solvers

Wrappers to BLAS, LAPACK, Eigen, Metis

# BEM4I library

Porting BEM on modern Intel architectures

- BEM produces dense matrices
  - Regular, coalesced memory access pattern

- Large BEM matrices approximated using low-rank approximation
  - Large number of dense matrix-vector multiplications
  - Can furthermore benefit from cache hierarchy

- BEM is compute-intensive
  - Numerical/semi-analytical evaluation of singular integrals require large number of flops
  - Benefits from large number of threads and wide SIMD registers of Xeon and Xeon Phi (co)processors

# BEM4I library

## SIMD vectorization of numerical/semi-analytical quadrature

- SIMD instruction sets (AVX512, IMCI) enable concurrent operations on up to 8 DP operands

- Most beneficial optimization techniques
  - OpenMP SIMD pragmas
  - Data alignment and padding
  - AoS to SoA transition for spatial coordinates
  - Unit-strided memory loads and stores

- Tested using
  - 2 x Intel Xeon (Haswell)
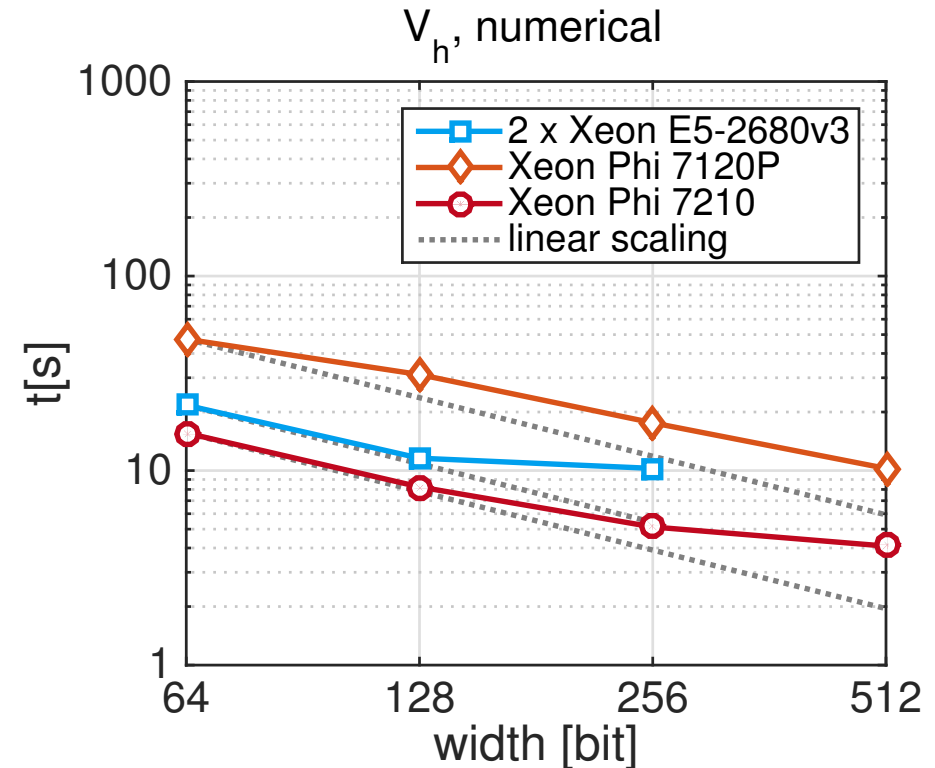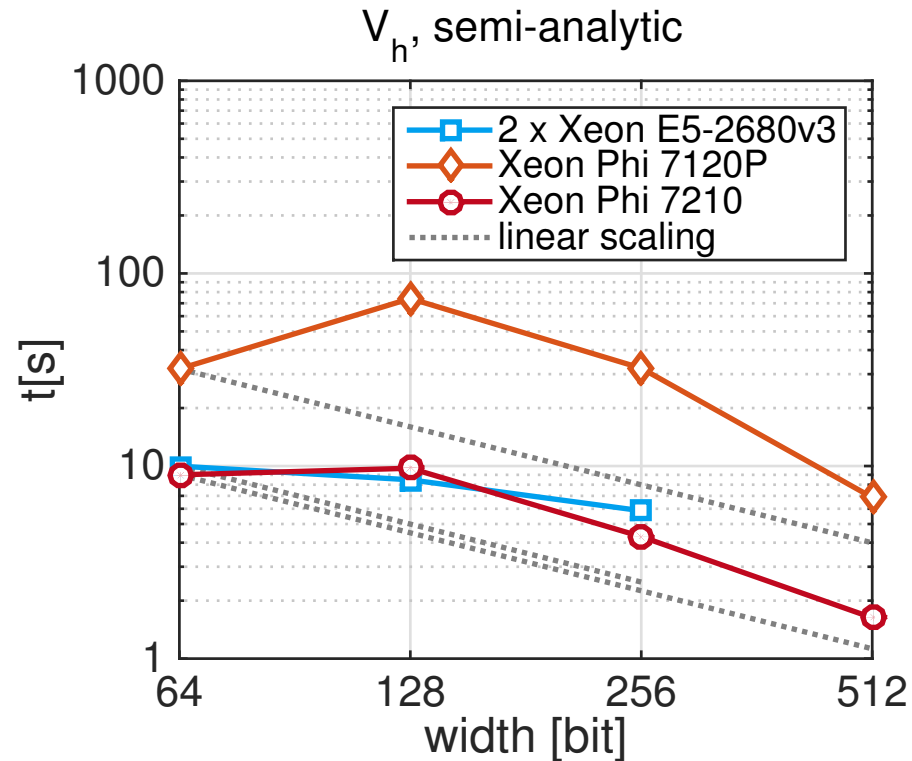  - Intel Xeon Phi 7120P (Knights Corner)
  - Intel Xeon Phi 7210 (Knights Landing)

```
1  #pragma omp declare simd simdlen( 8 )
2  evaluatePrimitive( double s, ... ) {
3  ...
4   // unmasked evaluation of sqrt
5   tmp1 = sqrt( tmp1 * tmp1 + q_sq );
6   // do not add to f in special case
7   if ( abs( s - sx ) > _EPS ) {
8    if ( tmp2 < 0.0 ) {
9     // masked division only
10    tmp3 = hh1 / ( tmp1 - tmp2 );
11
12   } else {
13
14    // masked addition only
15    tmp3 = tmp2 + tmp1;
16
17   }
18
19
20  } else {
21   tmp3 = 1.0;
22  }
23  // unmasked evaluation of log
24  f += ( s - sx ) * log( tmp3 );
25 ...}
```

# BEM4I library

## SIMD vectorization of numerical/semi-analytical quadrature



Performance of AVX2, AVX-512 and IMCI instruction sets on vectors with lengths 1, 2, 4, 8

# BEM4I library

SIMD vectorization of numerical/semi-analytical quadrature

|       | scalar | AVX512(1) | AVX512(2) | AVX512(4) | AVX512(8) |
|-------|--------|-----------|-----------|-----------|-----------|
| $V_h$ | 1.00   | 1.39      | 1.29      | 2.91      | 7.68      |
| $K_h$ | 1.00   | 1.62      | 1.72      | 3.41      | 8.25      |

Semi-analytical approach: scalar vs vectorized assembly on Xeon Phi 7210

|       | scalar | AVX512(1) | AVX512(2) | AVX512(4) | AVX512(8) |
|-------|--------|-----------|-----------|-----------|-----------|
| $V_h$ | 1.00   | 2.00      | 3.78      | 6.07      | 7.62      |
| $K_h$ | 1.00   | 1.20      | 2.26      | 3.89      | 5.53      |

Fully numerical approach: scalar vs vectorized assembly on Xeon Phi 7210

# BEM4I library

## Shared memory parallelization

- Using OpenMP parallel pragmas system matrix assembly is distributed among threads

- Tested on 2 x Intel Xeon (Haswell), Intel Xeon Phi 7120P (Knights Corner), Intel Xeon Phi 7210 (Knights Landing)
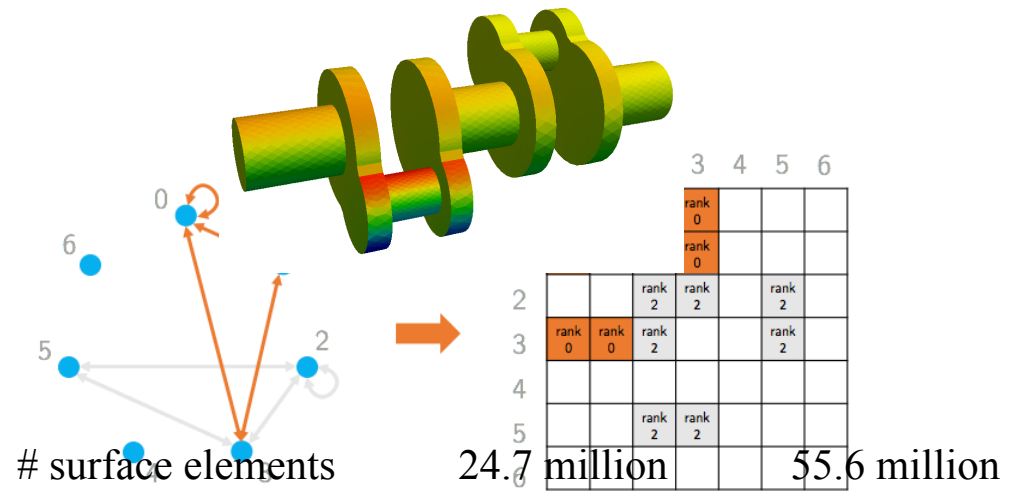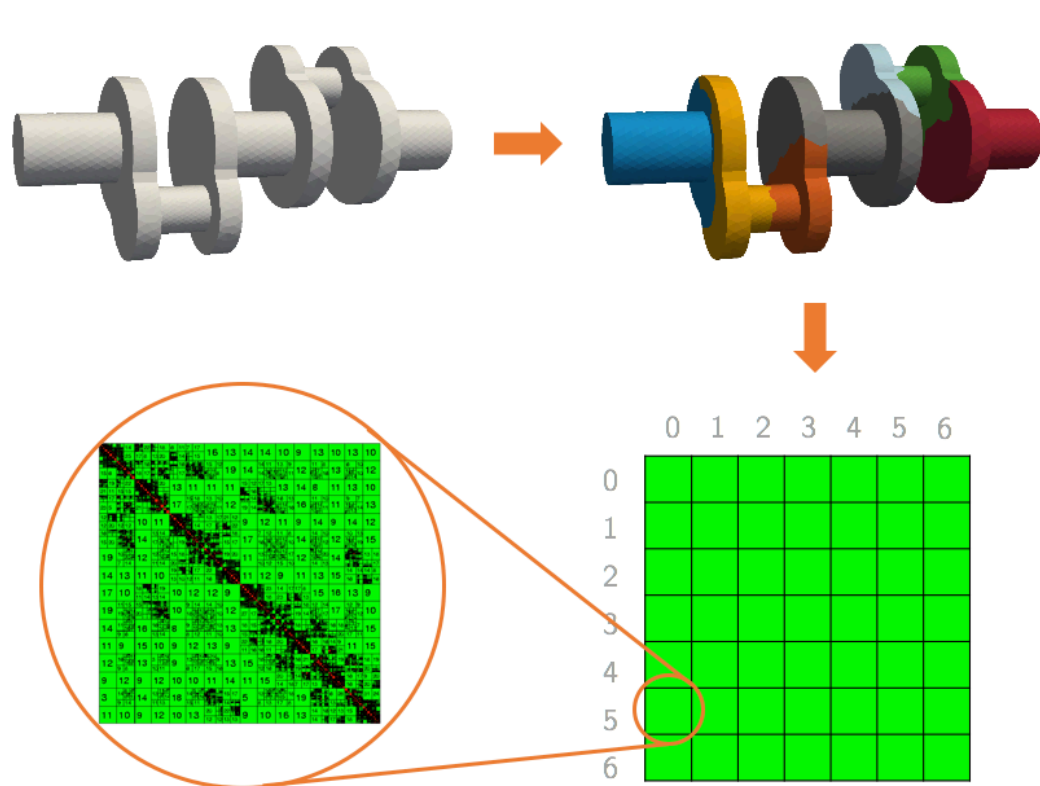


OpenMP scalability on Haswell, Knights Corner and Knights Landing

# BEM4I library

## Distributed memory parallelization using parallel ACA

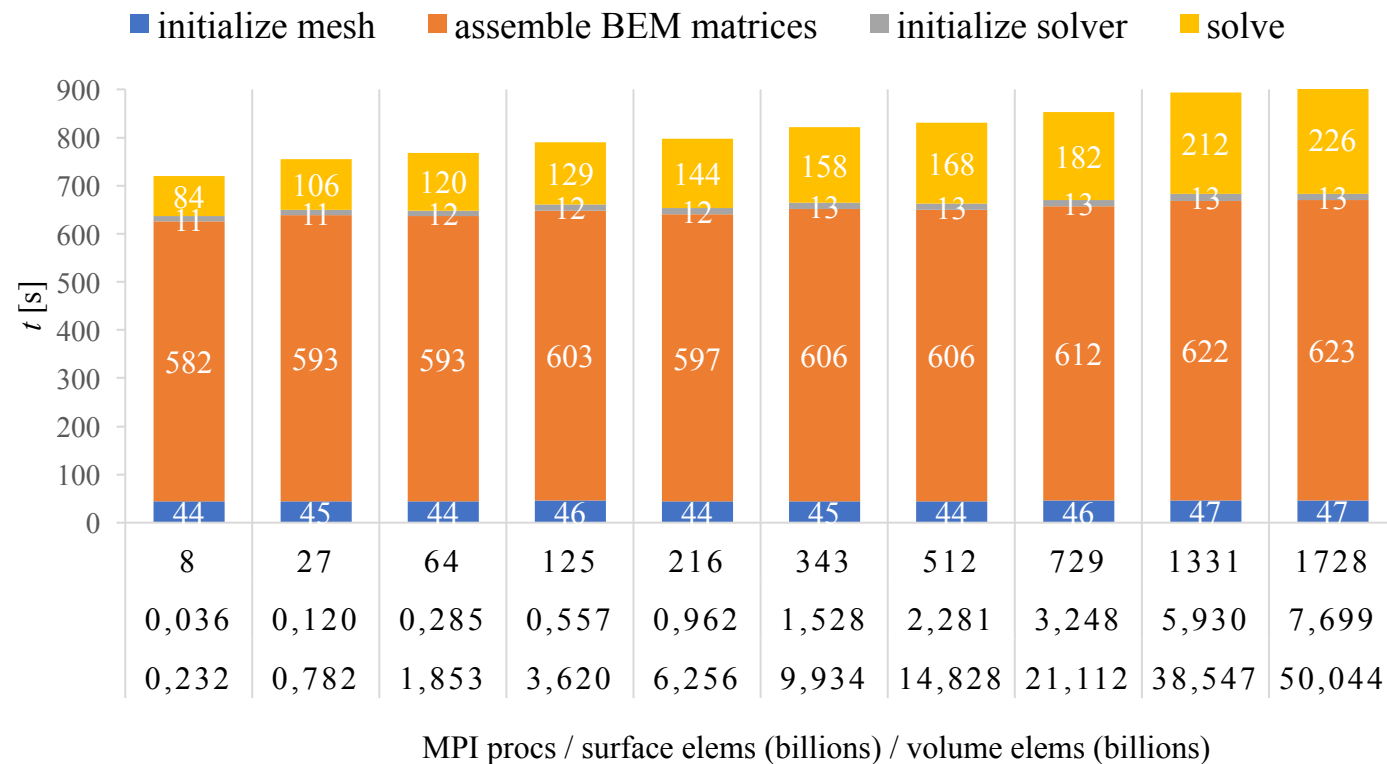- BEM4I implements parallel adaptive cross approximation (ACA)



| # surface elements | 24.7 million | 55.6 million |
| --- | --- | --- |
| # nodes | 64 | 256 |
| assembly $V_h/K_h$ [s] | 163/280 | 242/414 |
| compr. $V_h/K_h$ [%] | 0.08/0.16 | 0.09/0.18 |
| RAM $V_h/K_h$ [GB] | 3932/3987 | 22688/22873 |

# BEM4I library

## Distributed memory parallelization using BETI domain decomposition method

- Interface to the ESPRESO library enables us to use boundary element tearing and interconnecting method (BETI)

BETI on Salomon (Xeon E5-2680v3)



■ initialize mesh  ■ assemble BEM matrices  ■ initialize solver  ■ solve

MPI procs / surface elems (billions) / volume elems (billions)

BETI on HLRN TDS (Xeon Phi 7250)



■ initialize mesh  ■ assemble BEM matrices  ■ initialize solver  ■ solve

MPI procs / # of surface elems (billions) / # of volume elems (billions)

# BEM4I library

## Conclusion

- BEM4I provides BEM kernels optimized for Intel architectures

- SIMD and shared-memory parallelization using OpenMP

- Distributed memory parallelization by parallel ACA or BETI DDM

- Suitable for problems on unbounded domains (e.g., sound scattering), shape optimization, etc.

- International cooperation
  - Cooperation with TU Graz (parallel ACA, time-domain BEM, BETI, etc.)
  - Cooperation with Elmer (alternating Dirichlet-Neumann method, proof of concept)
  - Cooperation with UPMC Paris VI (DD based on MTF for the Helmholtz equation)

Agent 327: Operation Barbershop

Cosmos Laundromat: First Cycle

The Daily Dweebs

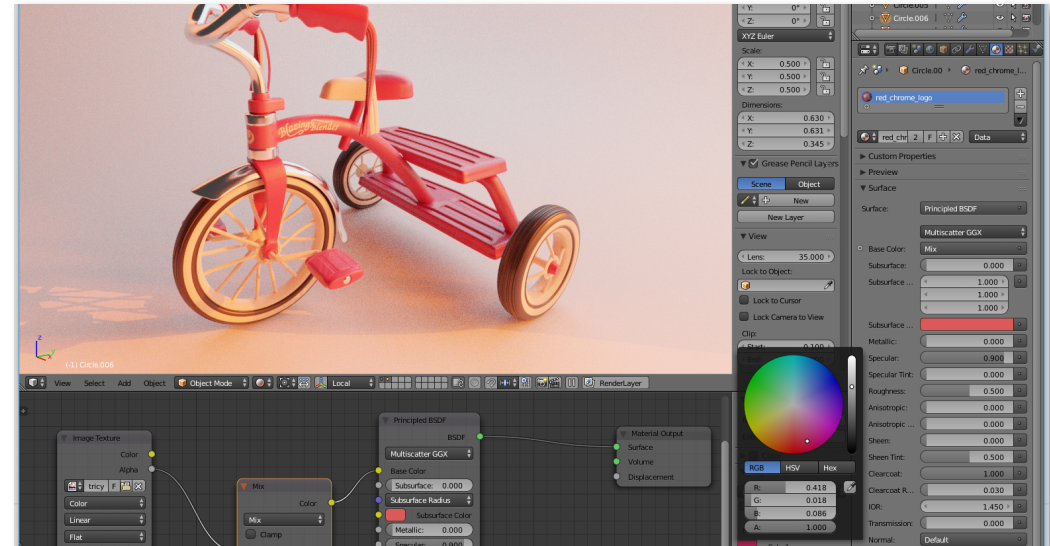# CyclesPhi Render

**Milan Jaroš**

# Blender Cycles
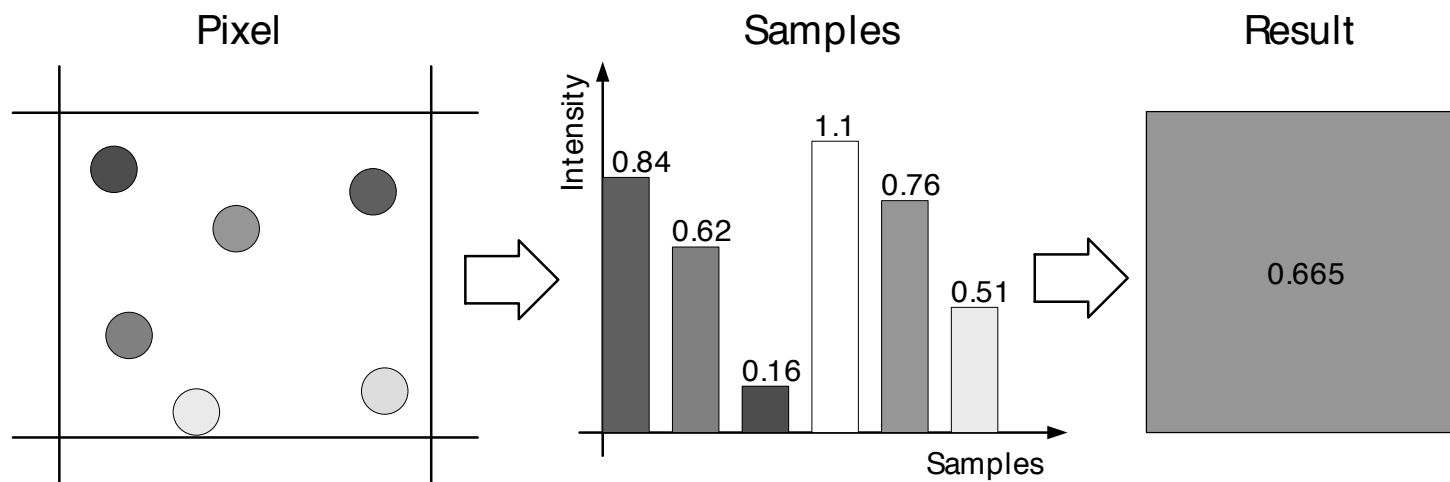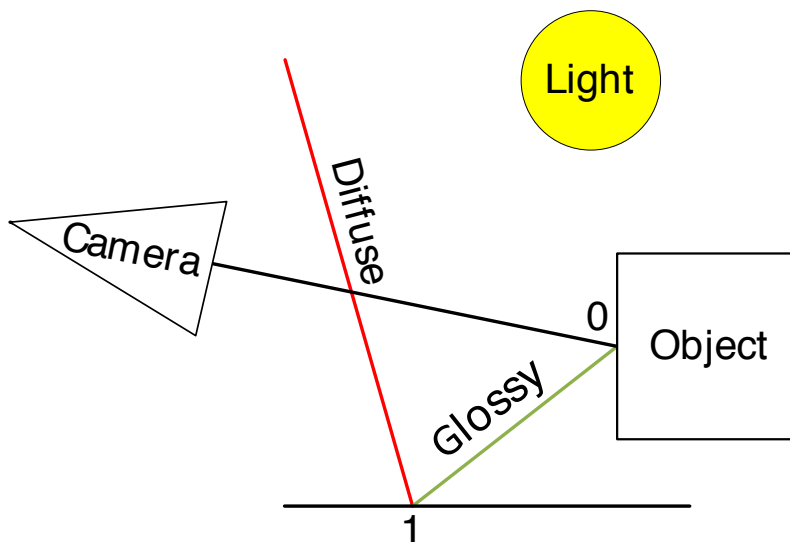
- **Blender**
  - open source 3D creation
  - render engines: Blender Internal and Cycles.

- **Cycles**
  - path-tracing based render engine
  - contains shading node system, texture workflow
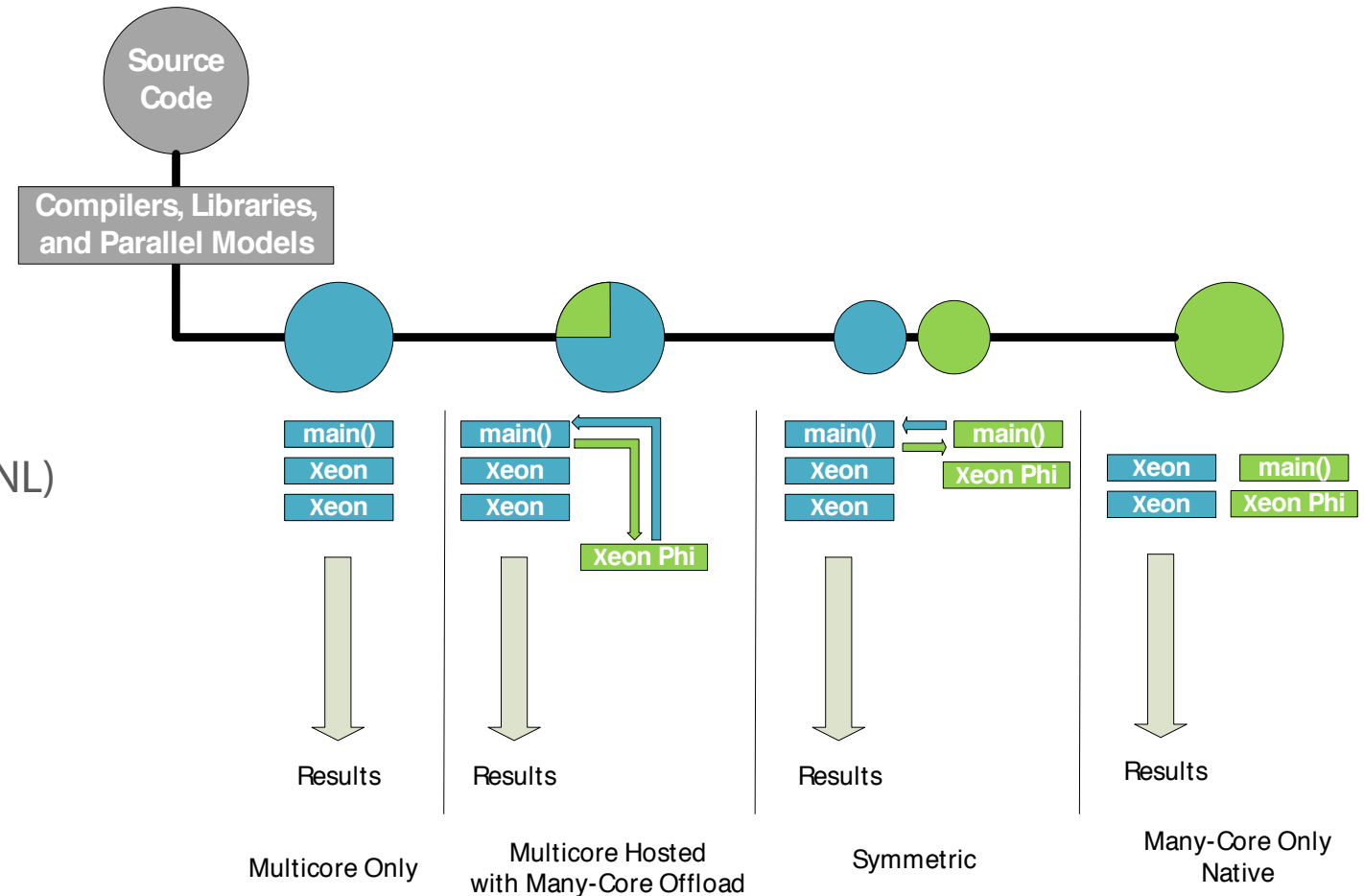  - interactive rendering
  - offline renderer

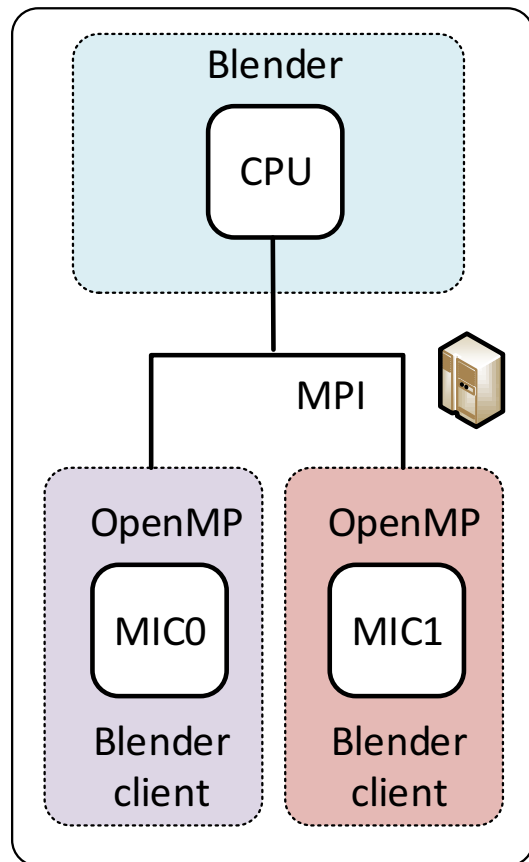# Path Tracing

# CyclesPhi

**Blender Cycles rendering engine was extended to support HPC environment:**

- OpenMP

- MPI

- Intel® Xeon Phi™
  - with Offload concept (KNC)
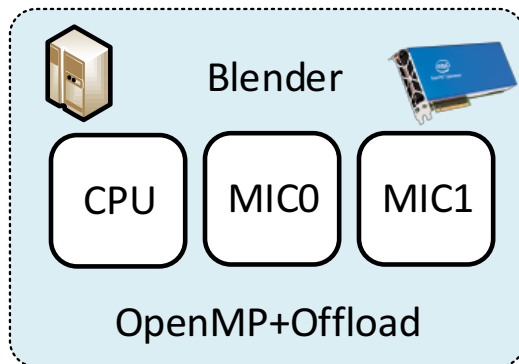  - with Symmetric mode (KNC, KNL)

- And their combinations



Source Code

Compilers, Libraries, and Parallel Models

| main() | main() | main() | main() | Xeon | main() |
| Xeon | Xeon | Xeon | Xeon Phi | Xeon | Xeon Phi |
| Xeon | Xeon | Xeon | | | |

Xeon Phi

Results        Results        Results        Results

Multicore Only        Multicore Hosted with Many-Core Offload        Symmetric        Many-Core Only Native

# Native Mode, Offload Mode and Symmetric Mode

Symmetric mode (KNC, KNL)

Blender
CPU
MPI
OpenMP — MIC0 — Blender client
OpenMP — MIC1 — Blender client

Offload mode (KNC)

Blender
CPU | MIC0 | MIC1
OpenMP+Offload

Native mode (KNL)

Blender
MIC0

# Benchmarks

- **Hardware setup**
  - Salomon supercomputer
    - 2x Intel Xeon E5-2680v3 CPUs and 2x Intel Xeon Phi 7120P per node
  - Intel Xeon Phi 7250 at HLRN-III Cray System
  - NVIDIA GeForce GTX 970

**Classroom**

**Fishy cat**

**Victor**



Classroom by Christophe Seux



Fishy Cat by Manu Jarvinen



Victor by Blender Foundation

# Performance Comparison

IT4Innovations national supercomputing center

Rendering time

| | GTX 970 Win | KNL (OMP272) | OMP24+MIC2 | OMP24 | KNC (OMP240) |

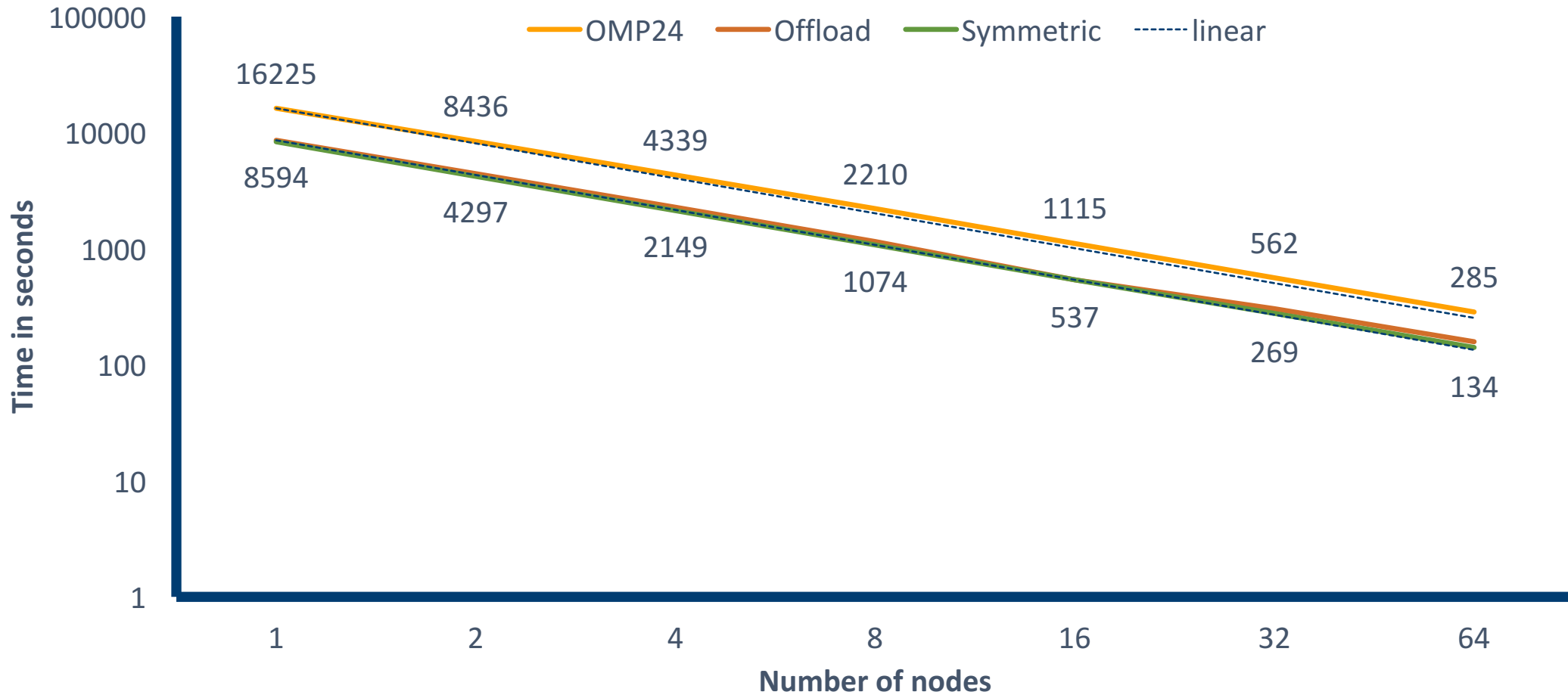# Benchmark Worm: Strong Scalability MPI Test (offline)

- Up to 64 computing nodes of the Salomon supercomputer

- 13.2 million triangles

- Resolution: 4096x2048

- 1024 samples:



Cosmos Laundromat - First Cycle

# Benchmark Tatra T87: Strong Scalability MPI Test (interactive)

- Up to 64 computing nodes of the Salomon supercomputer

- 1.2 million triangles

- HDRI lighting
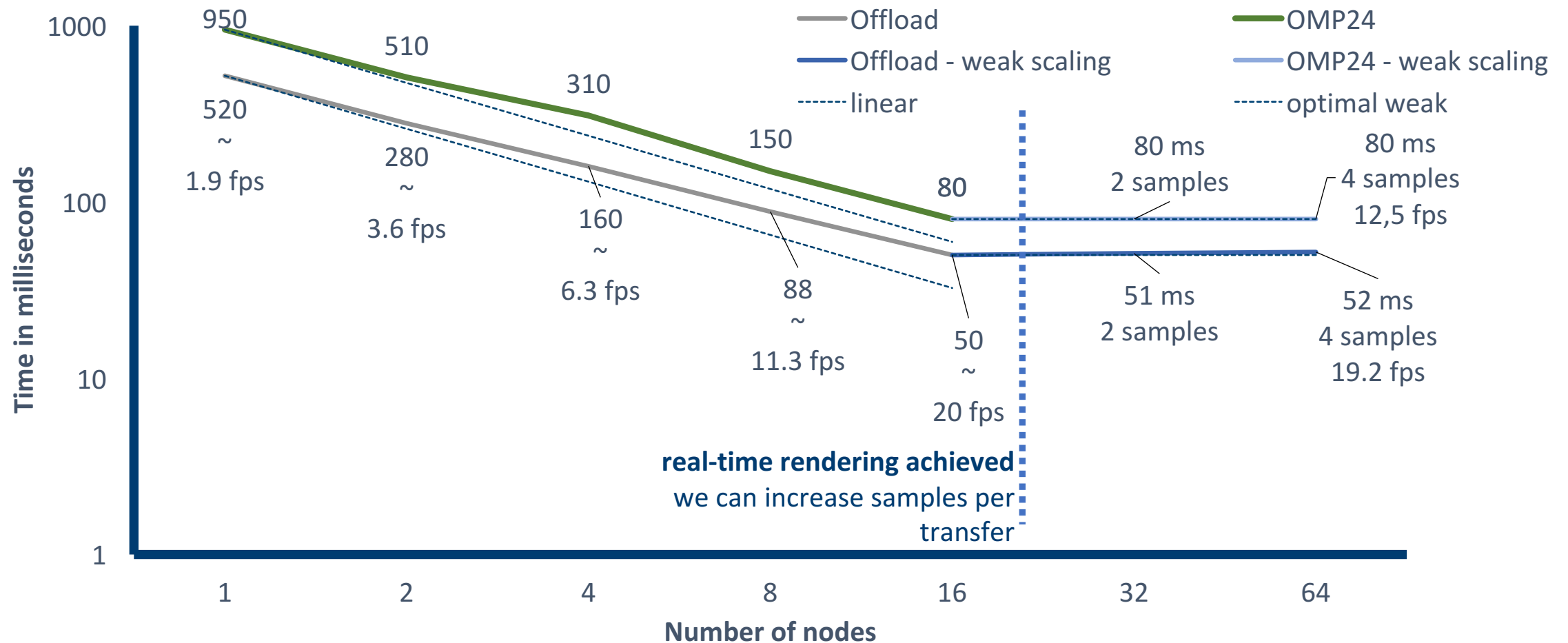
- Resolution: 1920x1080,

- 1 sample



Tatra T87 by David Cloete

# Benchmark Tatra T87: Strong Scalability MPI Test (interactive)

**Legend:**
- Offload
- OMP24
- Offload - weak scaling
- OMP24 - weak scaling
- linear
- optimal weak

Data labels (Offload/OMP24 line):
- 950
- 510
- 310
- 150
- 80
- 80 ms / 2 samples
- 80 ms / 4 samples / 12,5 fps

Offload line:
- 520 ~ 1.9 fps
- 280 ~ 3.6 fps
- 160 ~ 6.3 fps
- 88 ~ 11.3 fps
- 50 ~ 20 fps
- 51 ms / 2 samples
- 52 ms / 4 samples / 19.2 fps

Y-axis: Time in milliseconds (1, 10, 100, 1000)

X-axis: Number of nodes (1, 2, 4, 8, 16, 32, 64)

**real-time rendering achieved**
we can increase samples per transfer

THANK YOU

THANK YOU