

DE LA RECHERCHE À L'INDUSTRIE



www.cea.fr

PORTING PARALLEL APPLICATIONS TO HETEROGENEOUS SUPERCOMPUTERS: LIBRARIES AND TOOLS CAN MAKE IT TRANSPARENT

Jean-Yves VET, DDN Storage
Patrick CARRIBAULT, CEA
Albert COHEN, INRIA

CATC 2016
September, 15

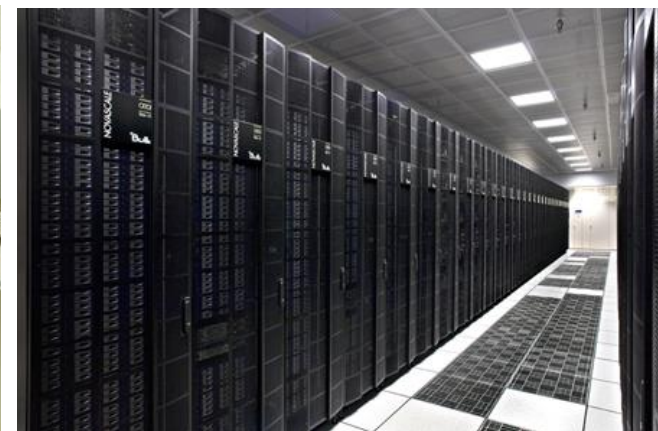
Since the 60s, about a new machine every 5 years...



CDC 7600 (1976)



CRAY 1S (1982)



Bull Tera-10 (2006)

Some legacy codes >100k lines. Maintaining and porting legacy codes is a huge amount of work.

A strong need for:

- Increasing portability
- Reaching decent compute efficiency (HPC)
- Porting code in a cost-efficient way (libraries or transparent mechanisms, incremental changes)

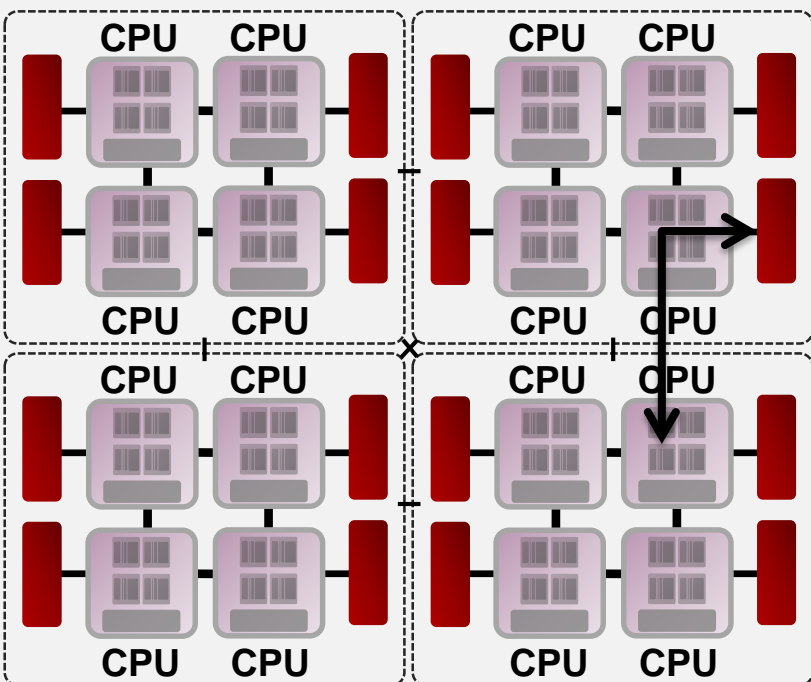


Tera-100 (2010)

Tera-100 (Bull)

- 1.254 PFLOP/s (ranked 6, November 2010 TOP500)
- Mainly homogeneous (Intel Xeon)
- Codes: MPI + OpenMP

Fat node (Bull Coherence System)



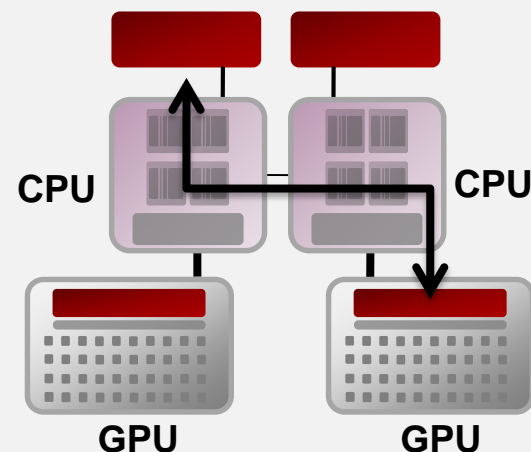
- Increased **Non-Uniform Memory Access (NUMA)** effects

- **Heterogeneous computing** (load balancing + programming models)

- **Non-hardware coherent memories** (GPU <-> CPU)

- **Non-Uniform IO Access (NUIOA⁽¹⁾)**

Heterogeneous node





H3LMS

*(Harnessing Hierarchy and Heterogeneity
with Locality Management and Scheduling)*

- **Bulk synchronous (multi-phase with barriers) task decomposition** to deal with heterogeneity.
- **NUMA aware scheduling**: mix data centric work distribution and hierarchical work stealing.
- **Transparent coupling with MPI and OpenMP** with an implementation into a single framework.

Common features
with **StarPU** ⁽¹⁾, **XKaapi** ⁽²⁾, or **OmpSs** ⁽³⁾

- Distributed Shared Memory (DSM) to handle data transfers automatically between non-hardware coherent memories in the compute node.
- Software caches to reduce memory transfers.



COMPAS

*(Coordinate and Organize
Memory Placement
and Allocation for Scheduler)*

- Keep track of data residency (NUMA node) to **guide scheduling**

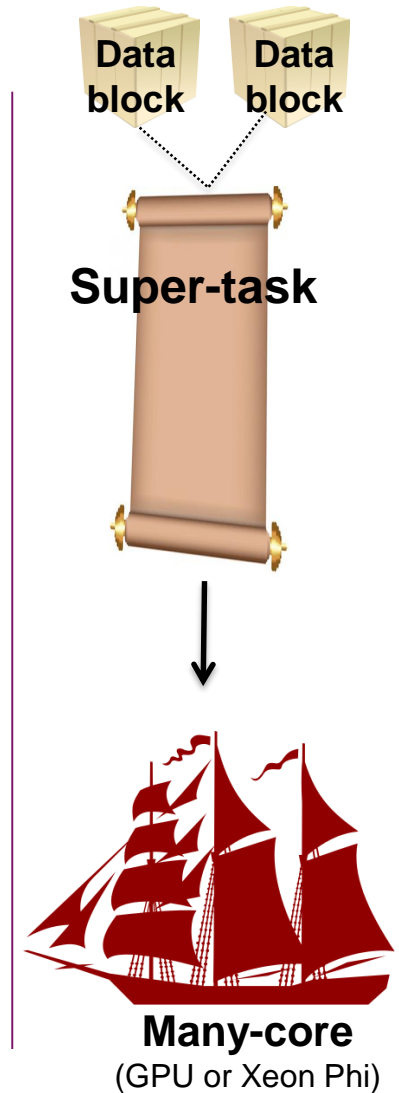
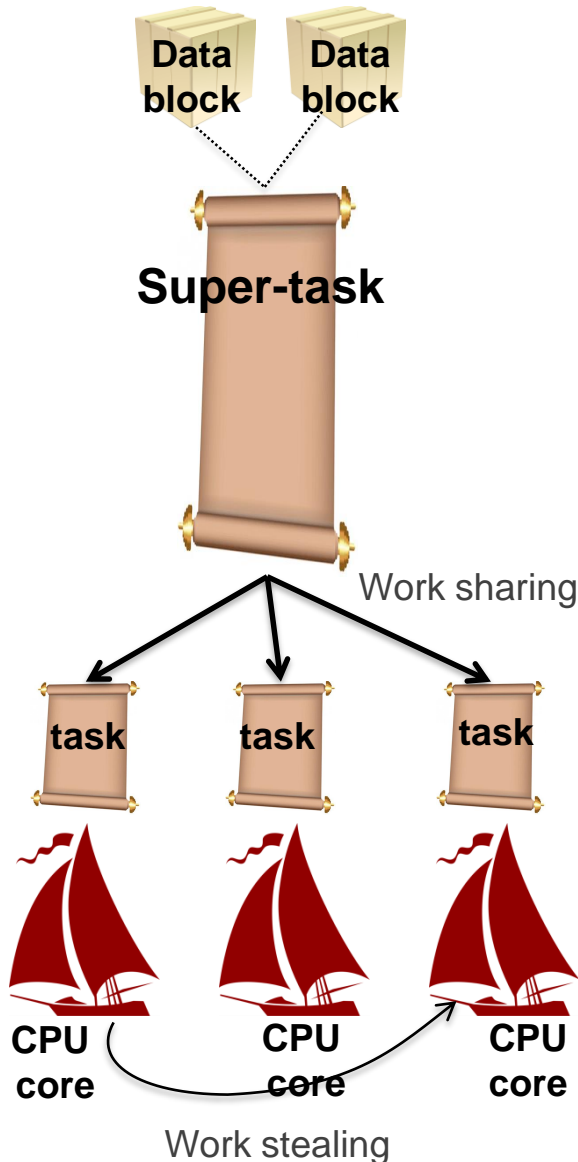
Common feature
with **Minas** ⁽⁴⁾

- Allocate memory and **distribute pages across NUMA nodes according to provided pattern.**

(1) Cédric Augonnet et al., *StarPU : A unified platform for task scheduling on heterogeneous multicore architectures*, Euro-Par'09
 (2) Thierry Gautier et al., *Xkaapi : A runtime system for data-flow task programming on heterogeneous architectures*, IPDPS 2013
 (3) Alejandro Duran et al., *Ompss : a proposal for programming heterogeneous multi-core architectures*. Parallel Processing Letters 2011
 (4) Pousa Ribeiro (C.) et al., *Minas: Memory Affinity Management Framework*. INRIA 2009.

Corner stone of the H3LMS platform

- Multi-level load balancing. Decomposable tasks to adapt the workload to the target architecture.
- Increased spatial locality on targets + work stealing between units of the same type.
- 2 scheduling modes :
 - For compute intensive applications Dynamic heterogeneous load balancing by using shared queue of super-tasks.
 - Data centric to minimize transfers by directly using local queues.



MAGMA⁽¹⁾

- Linear-algebra library for a heterogeneous node
- StarPU task scheduler
- Kernels: Intel MKL + Nvidia CuBLAS

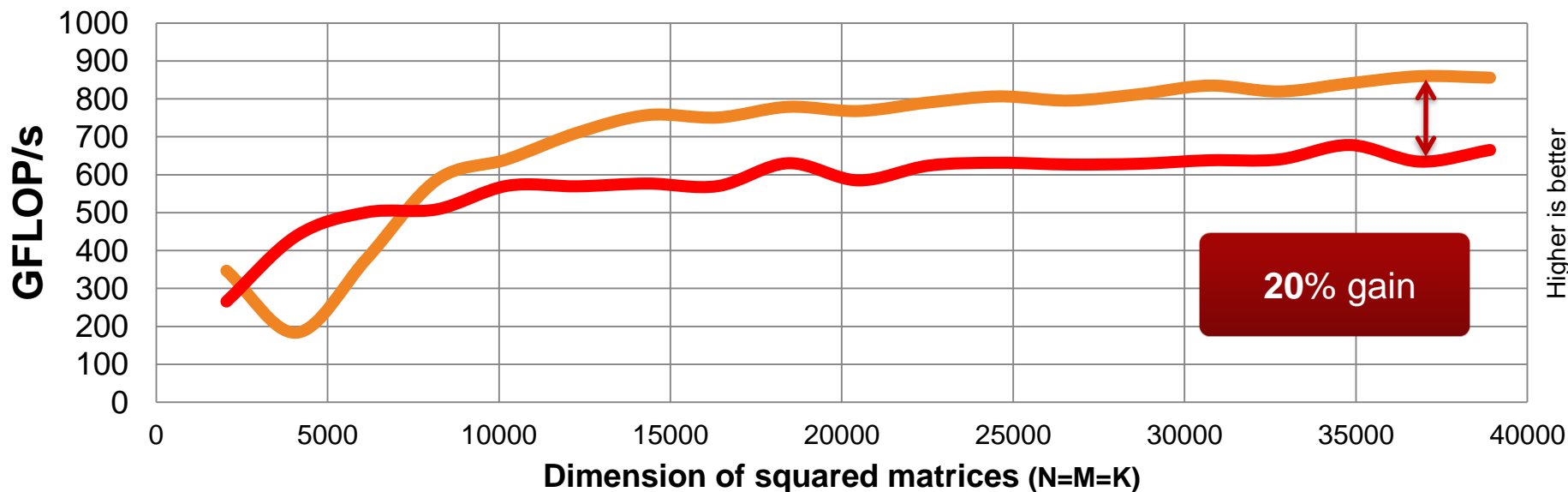
CPUs: 2 Intel Xeon X5660 (2x 4 cores)

GPUs: 2 Nvidia Tesla M2050

(Blocks: 1024x1024, Sub-blocks: 256x256)

Matrix multiply (SGEMM)

Comparison with MAGMA (including data transfers)



H3LMS (multi-granularity, single list of super-tasks)

MAGMA 1.3

(1) Agullo (E.) et al., Numerical linear algebra on emerging architectures : The PLASMA and MAGMA projects. Journal of Physics, 2009

Extended with an abstract organization

(based on the HWLOC ⁽¹⁾ library)

Worker Unit (constant granularity)

- Two kinds: small and accelerators
- Execute local tasks

Work Team (shared same physical memory)

- Execute super-tasks
- First level of work-stealing

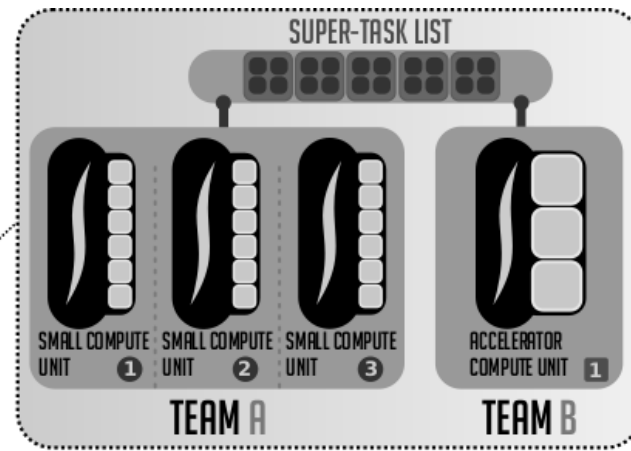
Work Pole (memory affinity)

- NUMA / NUIOA affinities
- As many lists of super-tasks as NUIOA nodes
- **COMPAS** to help selecting the pole and super-task-list

Choice of the list at the beginning of a bulk of super-tasks
(~owner compute rule⁽²⁾)

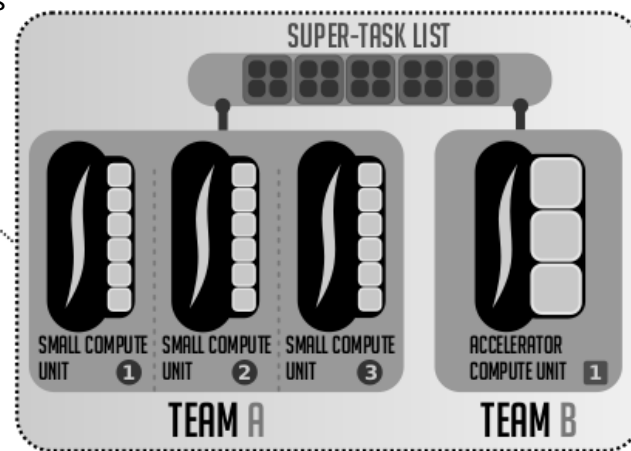
Pole hierarchy

- Poles organized in tree to map NUMA distances



POLE I

Abstract organization in a node:
2 NUMA nodes, 2x 4-core CPUs
and 2 accelerators



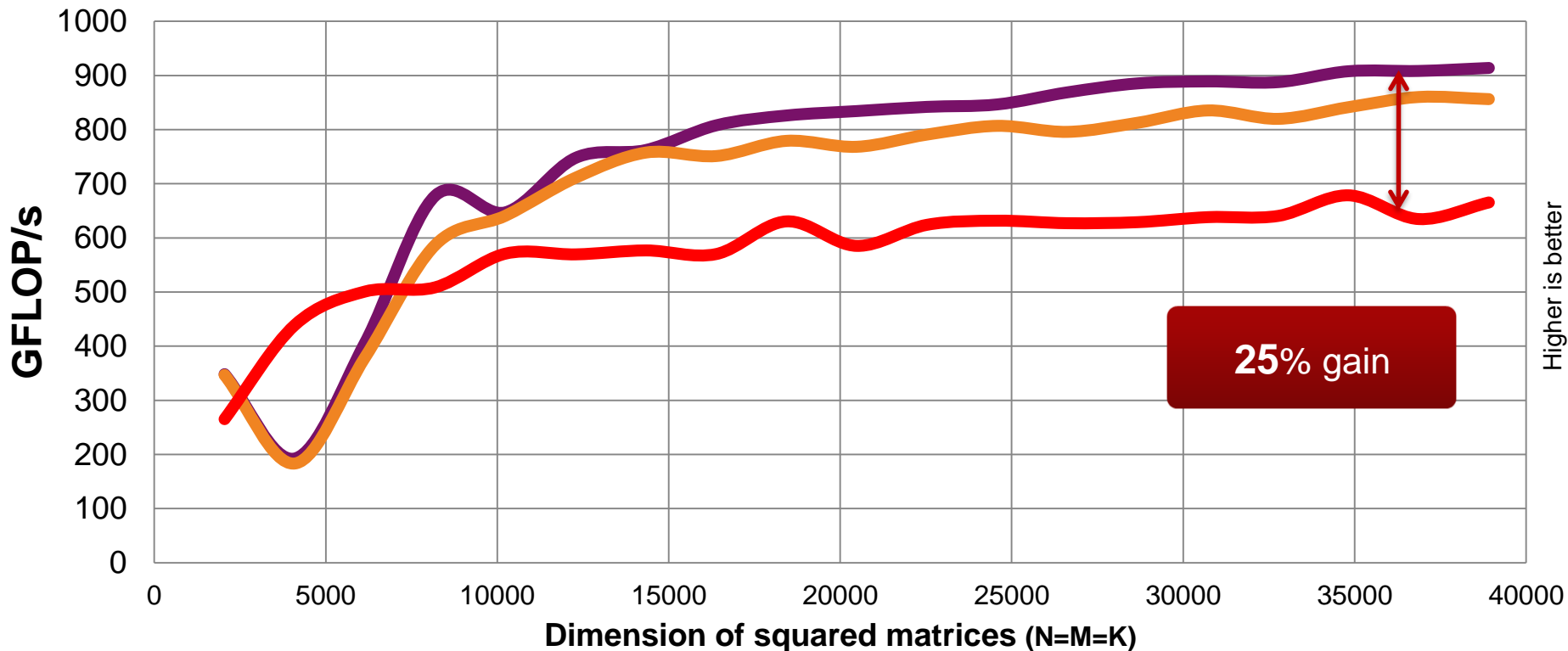
POLE II

included into

(1) F. Broquedis et al., HWLOC : A generic framework for managing hardware affinities in HPC applications. PDP '10.

(2) D. Callahan, et al., Compiling Programs for Distributed Memory Multiprocessors. The Journal of Supercomputing 1988.

Matrix multiply (SGEMM) Comparison with MAGMA (including data transfers)



— H3LMS (multi-granularity, hierarchical affinities) + **COMPAS**

— H3LMS (multi-granularity, single list of super-tasks)

— **MAGMA 1.3**

CPUs: 2 Intel Xeon X5660 (2x 4 cores)

GPUs: 2 Nvidia Tesla M2050

(Blocks: 1024x1024, Sub-blocks: 256x256)

5 SGEMMs accumulating in the same matrix ($C = A * B + C$) Comparison with MAGMA (including data transfers)

MAGMA with MORSE⁽¹⁾

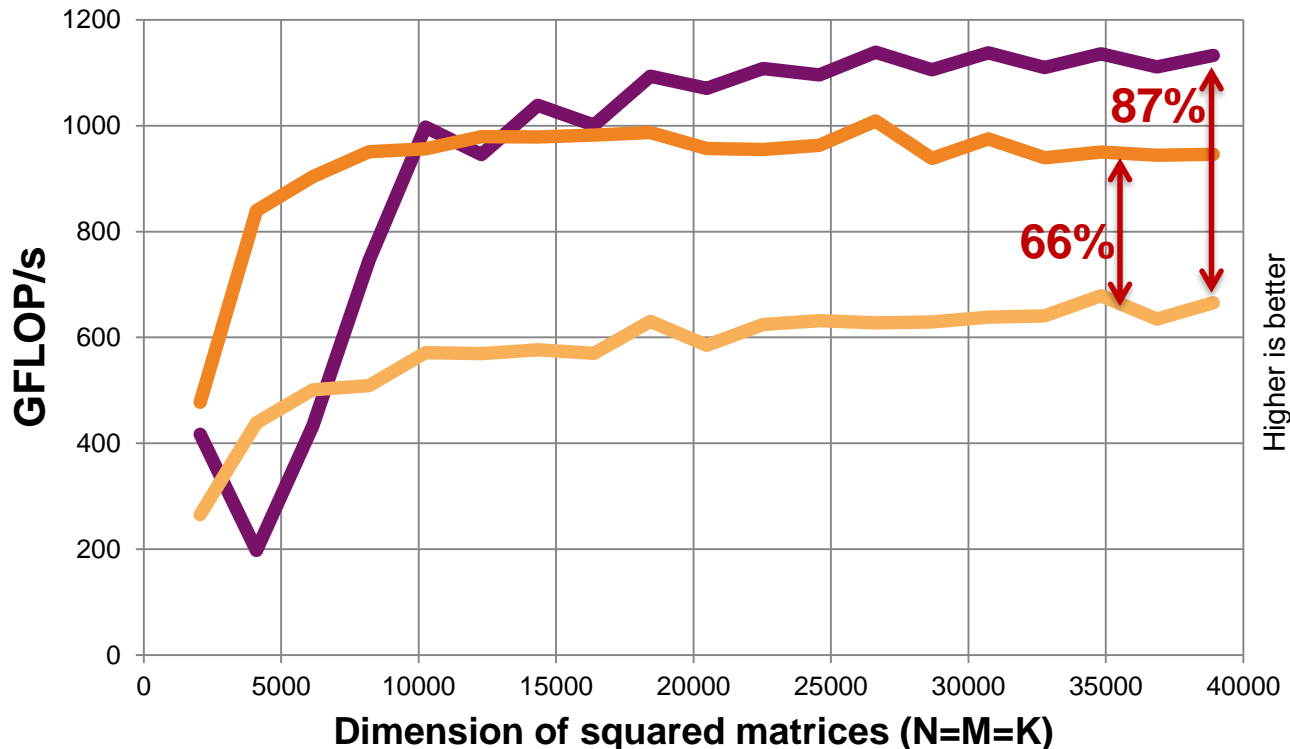
- Interface to link to linear-algebra libraries to runtime systems

- e.g. MAGMA on StarPU

With COMPAS:

- Keep same granularity of data blocks

- Turn off systematic flushing of software caches between library calls.



— H3LMS + COMPAS

— MAGMA 1.3 (MORSE modified)

— MAGMA 1.3

CPUs: 2 Intel Xeon X5660 (2x 4 cores)

GPUs: 2 Nvidia Tesla M2050

(Blocks: 1024x1024, Sub-blocks: 256x256)

(1) Matrices Over Runtime Systems @ Exascale (MORSE)

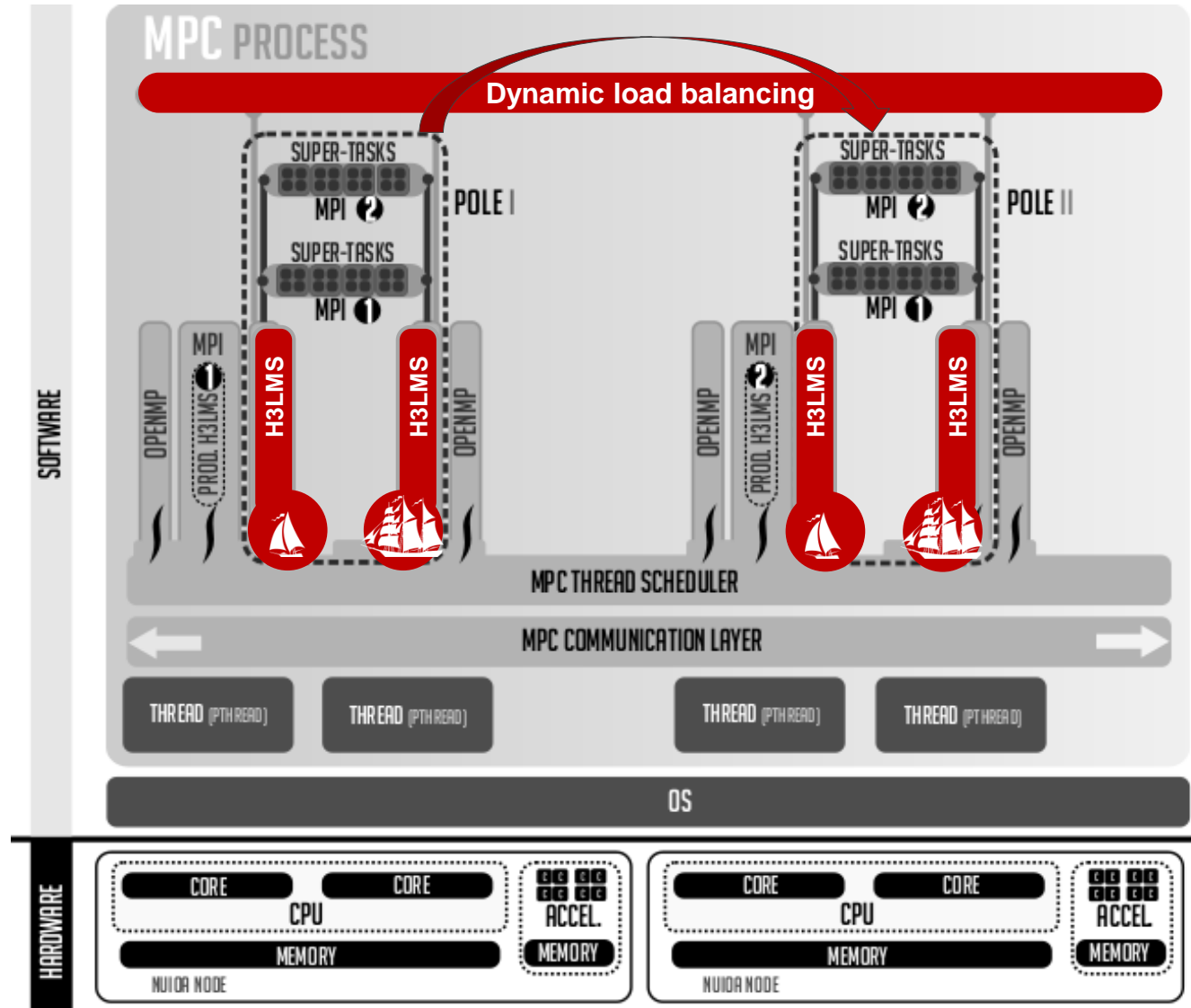
MPC⁽¹⁾

(Multi-Processor Computing)

- ➔ Framework developed by the CEA and the ECR
- ➔ Thread based MPI tightly coupled to an OpenMP implementation

H3LMS-MPC

- ➔ Rely on MPC for inter-node communications
- ➔ Load balancing with H3LMS, super-tasks generated from different MPI tasks in the same compute node



(1) M. Pérache et al., MPC: A unified parallel runtime for clusters of NUMA machines. Euro-Par '08.

EVALUATION WITH LEGACY CODES

LINPACK (HPL 2.0)

Need to modify 3 lines of code

1 Allocate page-locked memory with COMPAS.

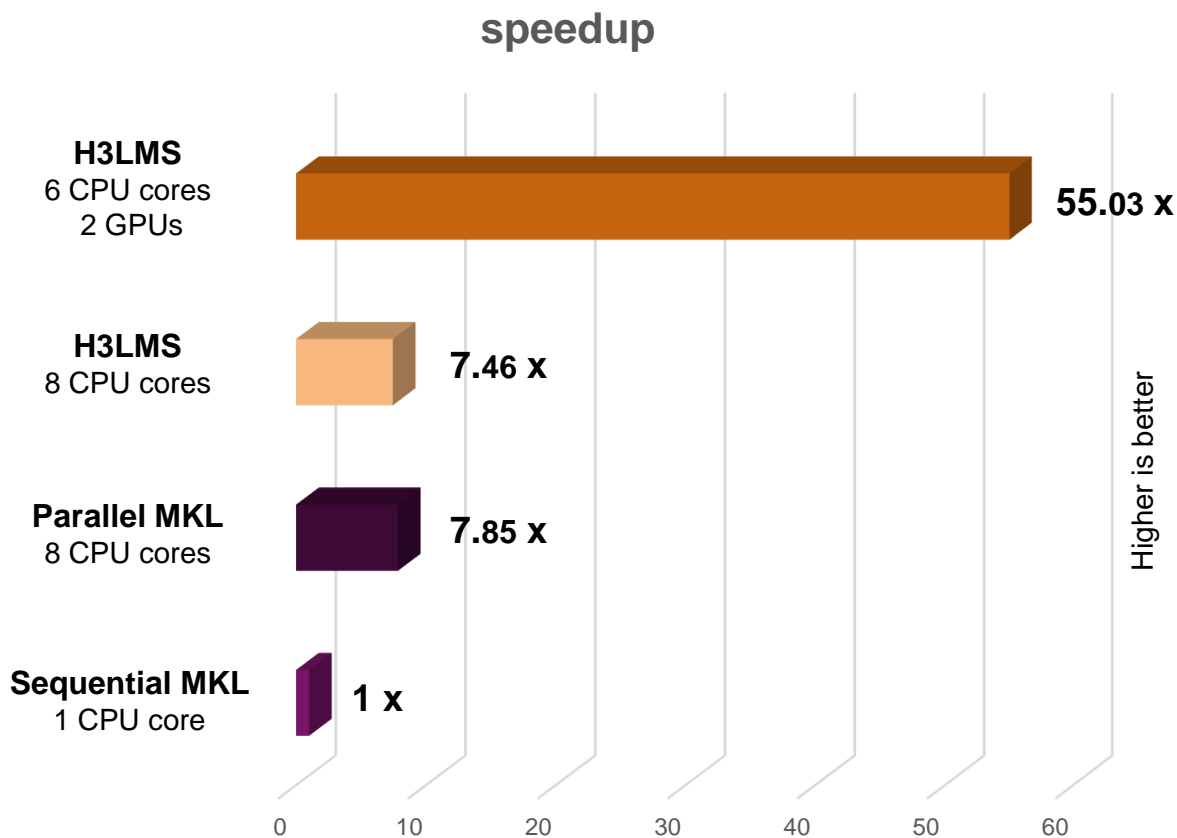
testing/ptest/HPL_pctest.c Before	testing/ptest/HPL_pctest.c After
164 vptr = (void*)malloc(165 ((size_t) (ALGO->align)+ 166 (size_t) (mat.ld+1)* 167 (size_t) (mat.nq)) * 168 sizeof(double));	164 vptr = (void*) (compas_malloc (NULL, COMPAS_PLOCKED, 165 ((size_t) (ALGO->align)+ 166 (size_t) (mat.ld+1)* 167 (size_t) (mat.nq)) * 168 sizeof(double));

2-3 Call BLAS function based on H3LMS.

include/hpl_blas.h Before	include/hpl_blas.h After
164 #define HPL_dgemm cblas_dgemm	164 #define HPL_dgemm (...) h3lms_blas_dgemm_ext_sync(__VA_ARGS__, 4096, -1, 256, -1)
165 #define HPL_dtrsm cblas_dtrsm	165 #define HPL_dtrsm h3lms_blas_dtrsm_sync

LINPACK – HPL 2.0

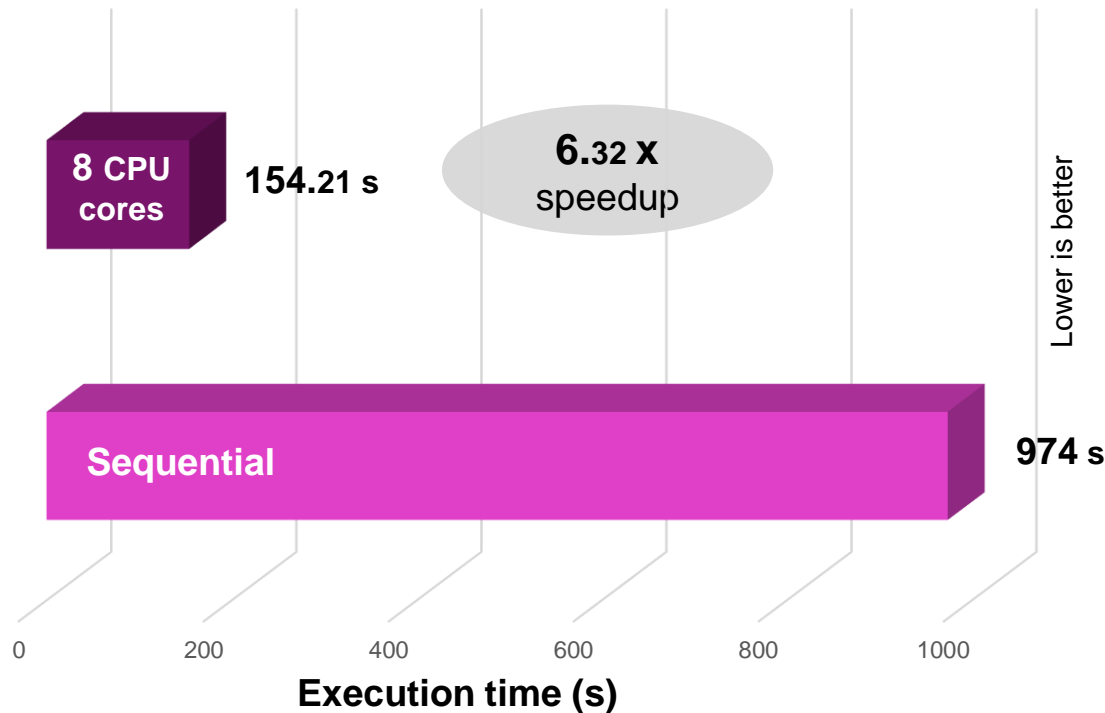
(N = 46080, N B = 512, P = 1, Q = 1, WC10L2L2)



- Based on optimized libraries Intel MKL 10.1 and NVIDIA CUBLAS 4.2
- Transparent for the user with synchronous function calls
- Internally decomposed into super-tasks and tasks
- Homogeneous performance close to parallel MKL (62.11 vs 68.94 GFLOP/s)
- Heterogeneous performance **482.4 GFLOP/s**

2x Intel Xeon Nehalem EP E5620 (2x 4 cores @ 2.4 GHz, peak: 2x 38.4 GFLOPS)
 2x NVIDIA Tesla Fermi M2090 (peak: 2x 665 GFLOPS)
 24 GB of DDR3 memory

- PN: solve the linear particle transport equation with deterministic resolution based on spherical harmonics approximation ⁽¹⁾
- Hybrid MPI-OpenMP code
- Focus on the *numerical_flux* function (~90% execution time).



PN: CPUs performance of *numerical_flux*

- **Double precision**
- Cartesian mesh : **1536x1536**, **N=15**, **36 iter.**
- Average on **20 runs**

Tera-100 heterogeneous compute nodes

CPUs: 2x 4 cores, Intel Xeon E5620

numerical_flux function on a $x * z$ Cartesian mesh

```

1  /* Part 1 */
2  DGEMM [read:A_X(136×136),B_X(x×z×136)] [write:C_X(x×z×136)]
3  DGEMM [read:A_Z(136×136),B_Z(x×z×136)] [write:C_Z(x×z×136)]
4
5  /* Part 2 */
6  DGEMM [read:A_X(136×136),B_L_X(x×136)] [write:C_L_X(x×136)]
7  DGEMM [read:A_Z(136×136),B_L_Z(z×136)] [write:C_L_Z(z×136)]
8  DGEMM [read:A_X(136×136),B_R_X(x×136)] [write:C_R_X(x×136)]
9  DGEMM [read:A_Z(136×136),B_R_Z(z×136)] [write:C_R_Z(z×136)]
10 BARRIER .....
11
12 /* Part 3 */
13 SUPER-TASK [read:D_X(1024),C_X(x×z×136)] [write:E_X(x×z×136)]
14 SUPER-TASK [read:D_Z(1024),C_Z(x×z×136)] [write:E_Z(x×z×136)]
15 BARRIER .....
16 SUPER-TASK [read:D_X(1024),C_L_X(x×136)] [read-write:E_X(x×z×136)]
17 SUPER-TASK [read:D_Z(1024),C_L_Z(z×136)] [read-write:E_Z(x×z×136)]
18 BARRIER .....
19 SUPER-TASK [read:D_X(1024),C_R_X(x×136)] [read-write:E_X(x×z×136)]
20 SUPER-TASK [read:D_Z(1024),C_R_Z(z×136)] [read-write:E_Z(x×z×136)]
21 BARRIER .....
22
23 /* Part 4 */
24 DGEMM [read:F_X(136×136),E_X(x×z×136)] [write:G_X(x×z×136)]
25 DGEMM [read:F_Z(136×136),E_Z(x×z×136)] [write:G_Z(x×z×136)]
26 BARRIER .....
27
28 FLUSH_SOFTWARE_CACHES

```

1 « Large » matrix multiplies

2 Small matrix multiplies

3 Consecutive loops operating on each cell of the mesh

4 « Large » matrix multiplies

Reference

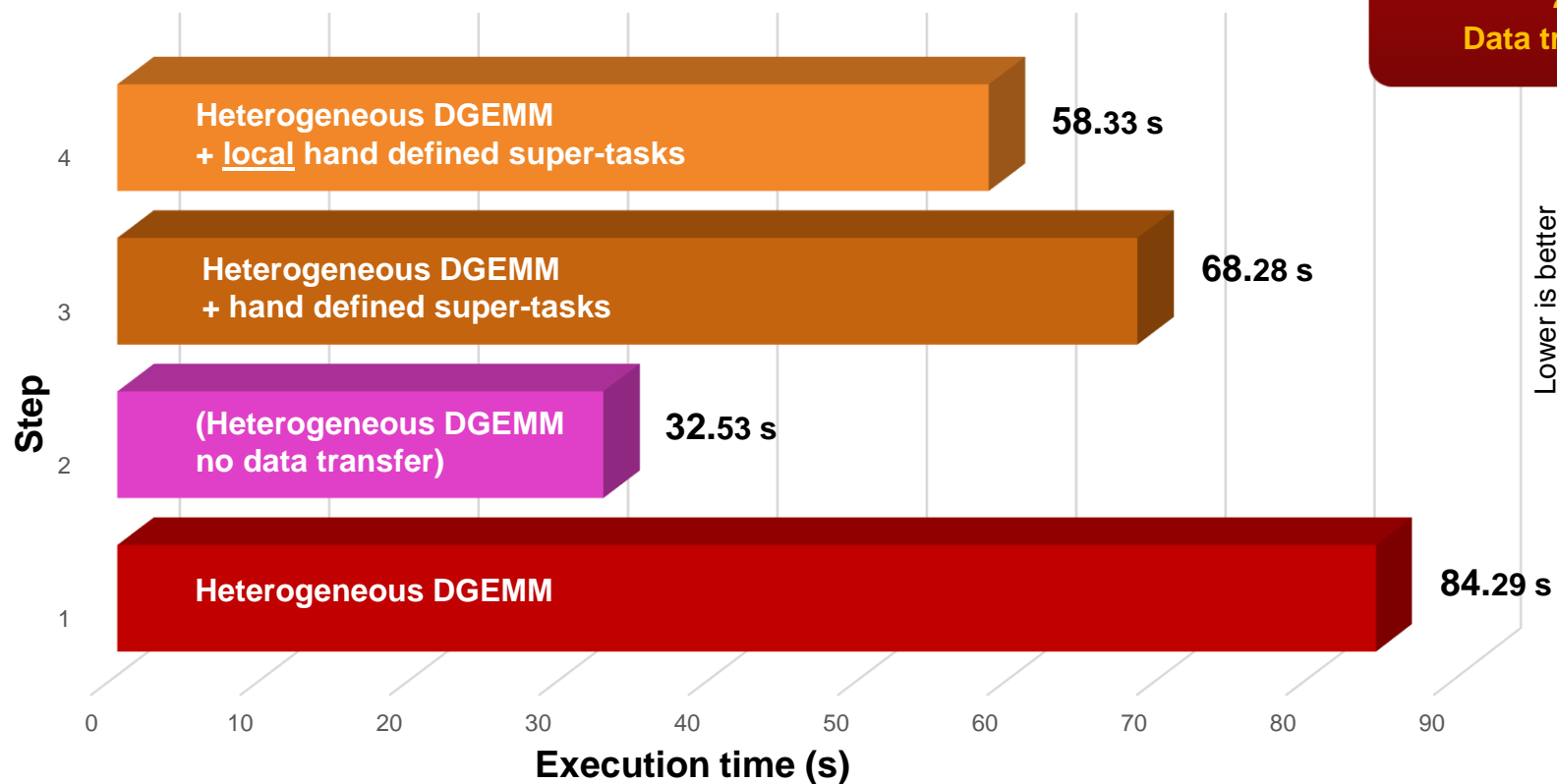
Homogeneous 8 cores: **154.21 s**

**PN: heterogeneous performance
of flux numerical_flux**

H3LMS+ COMPAS: multi-granularity,
spatial and temporal locality

Final speedup
Heterogeneous vs homogeneous

2.65 X
Data transfer bound



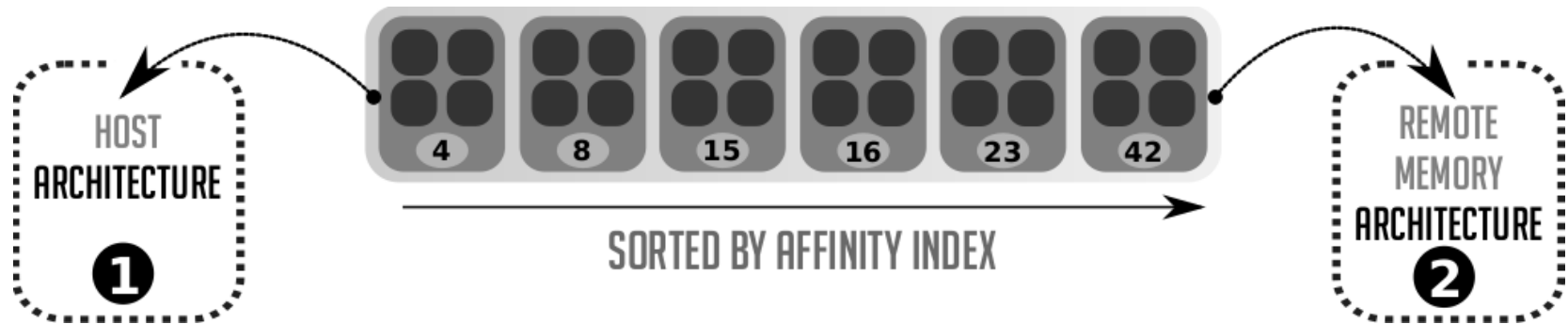
Lower is better

*double precision,
Cartesian mesh:
1536x1536,
N=15, 36 iter.
averaged on 20 runs*

THANK YOU

Backup

Software cache and execution order of the super-tasks⁽¹⁾



- Sort list of super-task at the bulk instantiation according to affinity with the corresponding memory.
- New scheduling policy based on cache statistics to reduce data transfers.

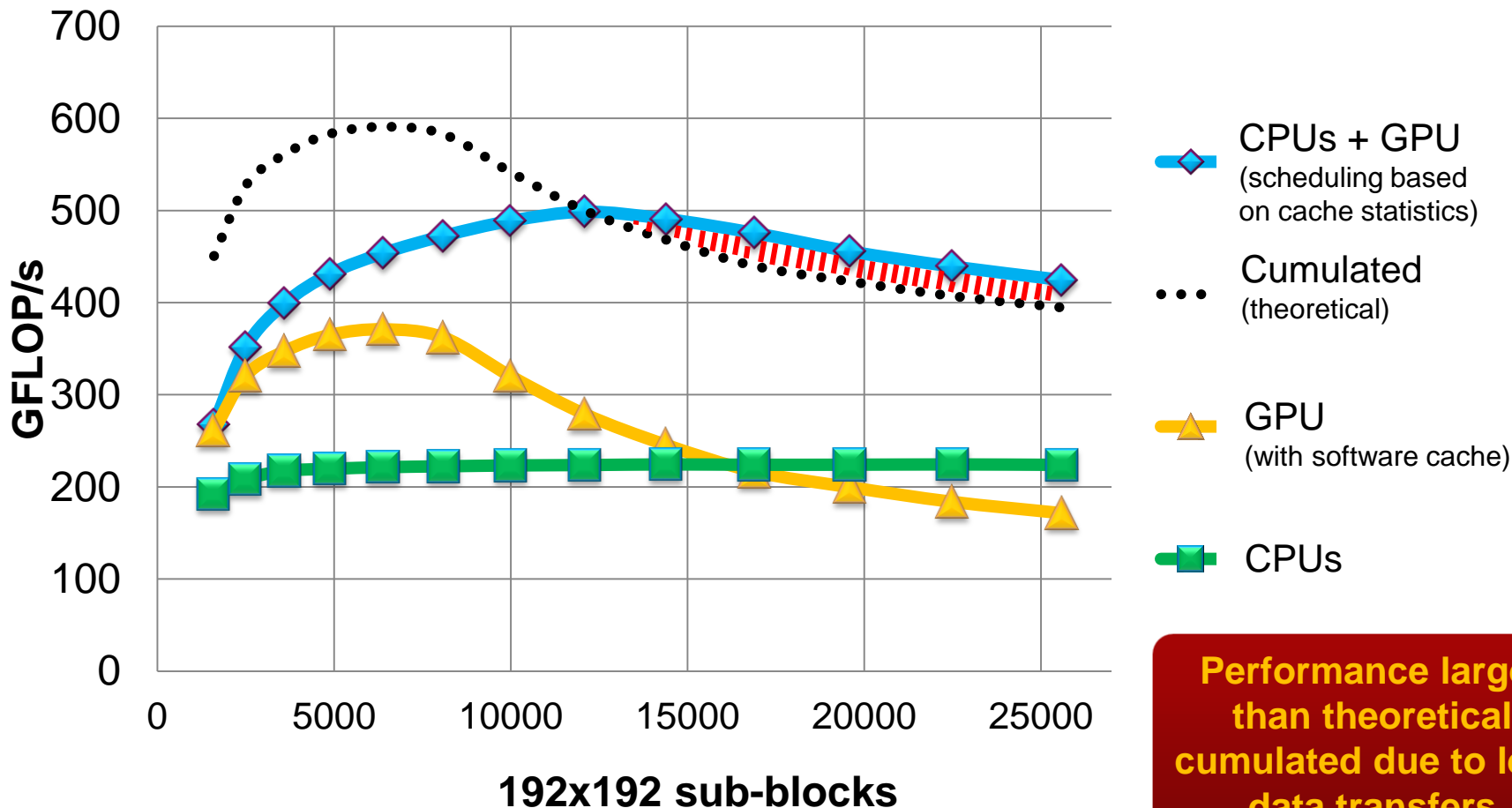
Affinity index depends on:

- Quantity of data
- If blocks of data are already held inside the cache

Backup

Sparse LU

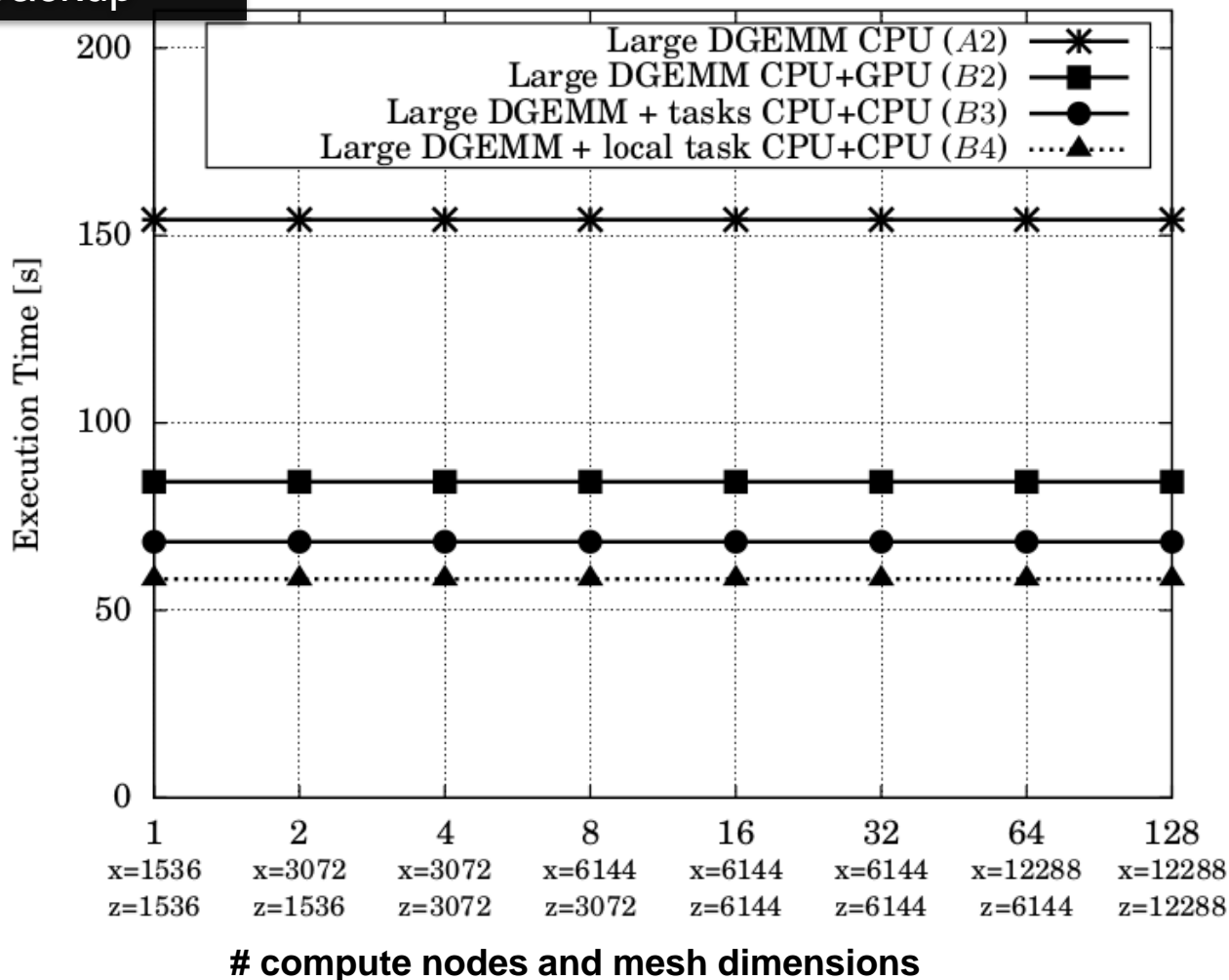
(step 3, simple precision, data transfers included)



Performance larger than theoretical cumulated due to less data transfers

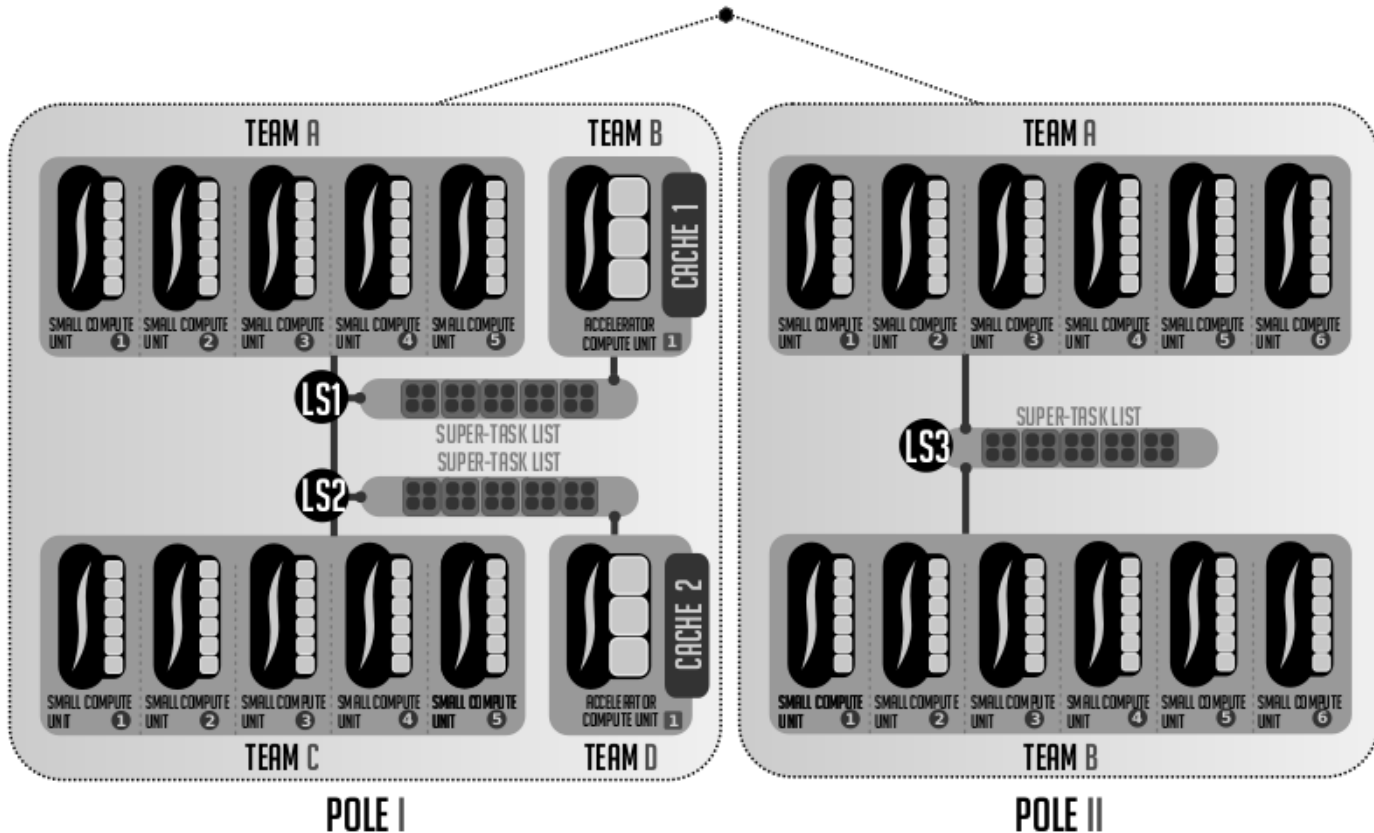
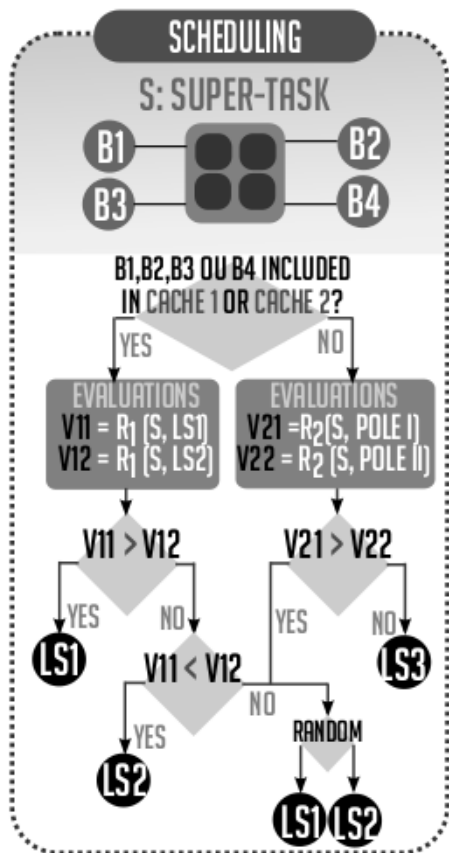
CPUs: 2x 12 cores AMD Opteron 6164 HE
Accel.: 1 GPU Nvidia Geforce GTX 470

Backup



PN: heterogeneous performance of flux numerical_flux on multiple compute nodes

Backup



Abstract organization in a node:

2 NUMA nodes,

2x 12-core AMD Magny-Cours CPUs (Dual package, 1 memory controller per CPU)
and 2 accelerators attached to the same NUMA node

Backup

COMPAS API

```

9  struct compas_pages_s pages;
10 compas_pages_init (&pages);
11 pages.geometry = COMPAS_2D;
12 pages.pattern = COMPAS_CYCLIC;
13 pages.xsize = NB*BS;
14 pages.ysize = NB*BS;
15 pages.xblocksize = BS;
16 pages.yblocksize =
    (BS*sizeof(double))/PSIZE;
17 distri.elsize =
    (BS*sizeof(double))/PSIZE;
18
19 double *A = compas_malloc(&distri,
    COMPAS_PLOCKED, NB*BS*NB*BS
    +sizeof(double));
20 init_matrix(A);

```

Proposal based on pragma

```

9  #pragma compas malloc
    elsize(sizeof(double))
    size(NB*BS, NB*BS)
    bsize((BS*sizeof(double))/PSIZE,
    (BS*sizeof(double))/PSIZE)
    pattern(cyclic)
    page(locked)
10
11 double *A = malloc(NB*BS*NB*BS +
    sizeof(double));
12 init_matrix(A);

```

Backup

H3LMS API

```

1 H3LMS_BEGIN_SUPERTASK_FUNCTION(l_bmod_cpu)
2 double *A1 = GETDATA(0, double*);
3 double *A2 = GETDATA(1, double*);
4 double *A3 = GETDATA(2, double*);
5
6 const int subblocks = 16;
7 const int subblocksize = BS/subblocks;
8
9 for (int i=0; i<subblocks; i++) {
10     double *SA1 = A1;
11     double *SA2 = &A2[i*subblocksize
12                 +BS+NB];
13     double *SA3 = &A3[i*subblocksize
14                 +BS+NB];
15     H3LMS_BEGIN_SUBTASK_DECLARATION
16     H3LMS_REGISTER_ARGS(SA1, SA2, SA3,
17                         subblocksize)
18     H3LMS_CPU_LAUCHER_FUNCTION(
19         bmod_cpu)
20     H3LMS_END_SUBTASK_DECLARATION
21 }
22 H3LMS_END_SUPERTASK_FUNCTION
23 ...
24
25 for (int ii=kk+1; ii<NB; ii++) {
26     for (int jj=kk+1; jj<NB; jj++) {
27         H3LMS_BEGIN_SUPERTASK_DECLARATION
28         H3LMS_REGISTER_NB_DEP(3)
29         H3LMS_REGISTER_ARGS( &A[ii*MROW +
30                             kk*BS], &A[kk*MROW + jj*BS],
31                             &A[ii*MROW + jj*BS])
32         H3LMS_REGISTER_DEP_2D( &A[ii*MROW +
33                               kk*BS], IN, sizeof(double),
34                               BS, BS+NB, BS)
35         H3LMS_REGISTER_DEP_2D( &A[kk*MROW +
36                               jj*BS], IN, sizeof(double),
37                               BS, BS+NB, BS)
38         H3LMS_REGISTER_DEP_2D( &A[ii*MROW +
39                               jj*BS], INOUT, sizeof(double),
40                               BS, BS+NB, BS)
41         H3LMS_CPU_LAUCHER_FUNCTION(
42             l_bmod_cpu)
43         H3LMS_ACCEL_LAUCHER_FUNCTION(
44             l_bmod_accel)
45         H3LMS_COMPUTE_RATIO(5)
46         H3LMS_END_SUPERTASK_DECLARATION
47     }
48 }

```

Proposal based on pragmas

```

1
2 for (int ii=kk+1; ii<NB; ii++) {
3     for (int jj=kk+1; jj<NB; jj++)
4     {
5         #pragma h3lms supertask
6         depend(in:
7             A[ii*BS:ii*BS+BS-1,
8               kk*BS:kk*BS+BS-1],
9             A[kk*BS:kk*BS+BS-1,
10              jj*BS:jj*BS+BS-1], inout:
11             A[ii*BS:ii*BS+BS-1,
12              jj*BS:jj*BS+BS-1])
13         accel(l_bmod_cuda)
14         cratio(5)
15     {
16         const int subblocks = 16;
17         const int subblocksize = BS
18             /subblocks;
19         for (int i=0; i<subblocks;
20             i++) {
21             double *SA1 = &A[ii*MROW
22                 + kk*BS];
23             double *SA2 = &A[kk*MROW
24                 + jj*BS + i*
25                 subblocksize+BS+NB];
26             double *SA3 = &A[ii*MROW
27                 + jj*BS + i*
28                 subblocksize+BS+NB];
29
30             #pragma h3lms subtask
31             bmod_cpu (SA1, SA2, SA3,
32                     subblocksize);
33         }
34     }
35 }
36 #pragma h3lms barrier

```

Glossary

BCS	– Bull Coherence System
BLAS	– Basic Linear Algebra Subprograms
COMPAS	– Coordinate and Organize Memory Placement and Allocation for Scheduler
CPU	– Central Processing Unit
DGEMM	– Double precision matrix matrix multiplication
DSM	– Distributed Shared Memory
DTRSM	– Solves one of the matrix equations $\text{op}(A) * X = \alpha * B, \text{ or } X * \text{op}(A) = \alpha * B$
GFLOPS	– Giga FLoating-point Operations Per Second
GPU	– Graphics Processing Unit
H3LMS	– Harnessing Hierarchy and Heterogeneity with Locality Management and Scheduling
HPC	– High Performance Computing
HPL	– High Performance Linpack
LRU	– Least Recently Used
NUMA	– Non-Uniform Memory Access
NUIOA	– Non-Uniform Input/Output Access
PFLOPS	– Peta FLoating-point Operations Per Second