

Master Of Science in Informatics at Grenoble  
*Parallel, Distributed, and Embedded Systems*

Parallel Systems

(1st Session)

Vincent Danjean, Arnaud Legrand

January 23rd, 2014

**Important information.  
Read this before anything else!**

- ▷ Any printed or hand-written document is authorized during the exam, even dictionaries. Books are not allowed though.
- ▷ Please write your answers to each problem on separate sheet of papers.
- ▷ The different problems are completely independent. You are thus strongly encouraged to start by reading the whole exam. You may answer problems and questions in any order but they have to be written in the order on your papers.
- ▷ All problems are independent and the total number of points for all problems exceeds 20. You can thus somehow choose the problems for which you have more interest.
- ▷ The number of points allotted to each question gives you an indication on the expected level of details and on the time you should spend answering.
- ▷ Question during the exam: if you think there is an error in a question or if something is unclear, you should write it on your paper and explain the choice you made to adapt.
- ▷ The quality of your writing and the clarity of your explanations will be taken into account in your final score. The use of drawings to illustrate your ideas is strongly encouraged.

## I. Parallel Code a Multi-core Machine ( $\approx$ 1h15)

The STL C++ library defines a template called `all_of` that takes as arguments, a array, a range within this array and a predicate `pred`. This template returns true if `pred` returns true for all the elements in the range `[first, last)` or if the range is empty, and false otherwise. In C, one would have to write an equivalent function with the following prototype:

```
int all_of(E* array, int first, int last, int (*pred)(E*))
```

where `E` is a previously defined type (e.g., a structure). In this problem, we look at various ways to parallelize such function.

The first listing (Figure 1) presents convenient functions that can be used in various implementations.

### Important remarks:

- ▷ The use of C codes is mainly to describe algorithms. Sometimes, to simplify presentation or notations, the proposed code is not a real working C code. You should not consider this as an error.
- ▷ In your answers, you can write code to show what you want to do but this is not required. In any case, you **have to explain** the goal of the proposed modifications. Once again, **the use of drawings is highly encouraged!**
- ▷ Remember there is not necessarily a single answer or cause to discuss. You should thus mention and discuss **any** flaw (involving hardware, algorithm, compiler, runtime, ...) you would think about.
- ▷ If you haven't done it yet, you should read the Important information paragraph on page 1.

```
typedef ... E;

#define N ... /* N is the number of processors on the machine */

int all_of_seq(E* array, int first, int last, int (*pred)(E*)) {
    int i;
    for (i=first; i<last; i++) {
        if (!pred(&(array[i]))) {
            return 0;
        }
    }
    return 1;
}

struct work_t {
    E* array;
    int first;
    int last;
    int (*pred)(E*);
    void* arg;
};

void fill_struct(struct work_t *w, E* array, int first, int last,
                int (*pred)(E*), void* arg) {
    w->array=array;
    w->first=first;
    w->last=last;
    w->pred=pred;
    w->arg=arg;
}

int all_of_struct(struct work_t *w) {
    return all_of_seq(w->array, w->first, w->last, w->pred);
}
```

Figure 1: Common functions available for all implementations

**Question I.1.** Quickly describe the rationale underlying the implementation proposed in Figure 2. There are at least two problems with this parallel implementation that may lead to a dramatic

```

static void* thread_work(void* arg) {
    return (void*)all_of_struct((struct work_t *)arg);
}

int all_of_1(E* array, int first, int last, int (*pred)(E*)) {
    struct work_t w[N];
    pthread_t tid[N];
    int t=(last-first)/N;
    int i,j;
    int ret=1;
    for(i=0, j=0; j<N; i+=t, j++) {
        fill_struct(&w[j], array, i, min(i+t, last), pred, NULL);
        pthread_create(&tid[j], NULL, thread_work, &w[j]);
    }
    for(j=0; j<N; j++) {
        void* retval;
        pthread_join(tid[j], &retval);
        ret &= (int)retval;
    }
    return ret;
}

```

Figure 2: Initial proposition for the all\_of function

slowdown. Describe the issues you can think of and the kind of workload where they would occur.

```

static void* thread_work(void* arg) {
    return (void*)all_of_struct((struct work_t *)arg);
}

#define T 1

int all_of_2(E* array, int first, int last, int (*pred)(E*)) {
    struct work_t w[N];
    pthread_t tid[N];
    int i,j;
    while (i<last) {
        int ret=1;
        for(j=0; j<N; i+=T, j++) {
            fill_struct(&w[j], array, i, min(i+T, last), pred, NULL);
            pthread_create(&tid[j], NULL, thread_work, &w[j]);
        }
        for(j=0; j<N; j++) {
            void* retval;
            pthread_join(tid[j], &retval);
            ret &= (int)retval;
        }
        if (ret == 0) { return 0 ; }
    }
    return 1;
}

```

Figure 3: Second proposition for the all\_of function

**Question I.2.** Quickly describe the rationale underlying the implementation proposed in Figure 3.

1. Discuss about the influence of the  $T$  value.
2. There are at least two other problems with this parallel implementation that may lead to a dramatic slowdown. Describe the issues you can think of and the kind of workload where they would occur. Without changing the main idea behind this implementation, propose how to fix these issues.

**Question I.3.** Quickly describe the rationale underlying the implementation proposed in Figure 4. Does this new implementation solve the issues you previously mentioned? Does it create new issues?

```

static pthread_mutex_t m=PTHREAD_MUTEX_INITIALIZER;

static void* thread_work(void* arg) {
    struct work_t *w = (struct work_t *)arg;
    int i;
    while(i < w->last) {
        pthread_mutex_lock(&m);
        i=w->first++;
        pthread_mutex_unlock(&m);
        if (i >= w->last) { break ; }
        if (!w->pred(&w->array[i])) {
            return 0;
        }
    }
    return 1;
}

int all_of_3(E* array, int first, int last, int (*pred)(E*)) {
    pthread_t tid[N];
    struct work_t w;
    int j;
    int ret=1;
    w.first=first;
    w.last=last;
    w.pred=pred;
    for(j=0; j<N; j++) {
        pthread_create(& tid[j], NULL, thread_work, &w);
    }
    for(j=0; j<N; j++) {
        void* retval;
        pthread_join(tid[j], &retval);
        ret &= (int)retval;
    }
    return ret;
}

```

Figure 4: Third proposition for the all\_of function

Without changing the main idea behind this implementation, how can you improve this code?

**Question I.4.** Another classical technique to improve performances is to implement work stealing using a divide and conquer approach.

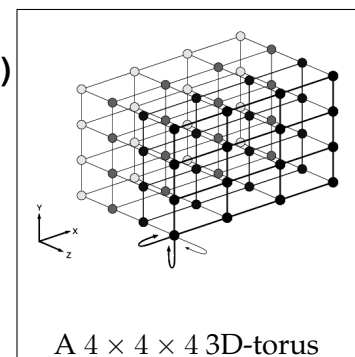
1. Write the corresponding (pseudo) code.
2. What are the issues solved by such an approach when compared to the simple work stealing strategy of the previous question?

**Question I.5.** What are the remaining issues? How could you further improve your code?

## II. Embedding Issues in a High Dimension Network ( $\approx 30$ min)

Stencil codes are a class of iterative kernels which update array elements according to some fixed pattern, called stencil. Classical parallel 2D stencil applications require that each cell be updated using its North, South, East, and West neighbors:  $\forall t, \forall i, j$ ,

$$A[i, j]^{(t+1)} = f \left( A[i-1, j]^{(t)}, A[i+1, j]^{(t)}, A[i, j-1]^{(t)}, A[i, j+1]^{(t)} \right)$$



**Question II.1.** Write the pseudo-code for such a stencil application using a  $n \times n$  domain on a  $q \times q$  2D processor grid (we assume  $q$  divides  $n$ ).

Give a performance model for your algorithm, assuming a 4-port bidirectional model.

**Question II.2.** Now, let us assume we have at our disposal a  $q \times q \times q$  3D-torus topology. It would obviously be possible to reuse the previous algorithm by restricting to a slice of the 3D-torus but we would be using only  $q^2$  processors instead of  $q^3$ ...

1. Propose a mapping of the  $n \times n$  matrix on the full  $q \times q \times q$  3D torus that preserves locality as much as possible.

**Hints:**

- ▷ **Draw pictures!!!**
  - ▷ Try to figure out how you would lay out a string onto a sheet of paper without breaking it
  - ▷ Try to figure out how you would fold a sheet of paper into a 3D box without stretching the paper too much
2. Explain how communications will conflict with each others and the kind of slowdown you may expect when compared with a  $q \times q$  2D grid.

### III. Performance Evaluation of a Recent Accelerator ( $\approx$ 1h15)

The article that can be found in the appendix is a technical report written by a Caltech research scientist (do not worry about the missing pages 7 and 8, which have been removed for clarity). This document reports his initial investigations (both results and experimental protocol) on the performance of the recently released Intel Xeon Phi. Unlike classical accelerators like GPU, this accelerator does not require major code rewriting. However, obtaining good performances can be challenging.

The goal of this section is to evaluate your ability to comment articles, figures and experimental protocol.

- ▷ We will denote by  $XP_1$  the set of experiments reported in the Figure 1 of this article. Likewise,  $XP_2$  and  $XP_3$  denote the set of experiments corresponding to Figure 2 and 3.
- ▷ Reading the code provided in the appendix of this article may provide with interesting information.

**Question III.1.** In Section 4, the author reports the performance he achieved to obtain on the Xeon Phi and on the Xeon E5 with the naive implementation.

1. The author explains that “the results are contrary to expectations reported by the vendor and others in the literature”. From the elements underlined in the text, compute the theoretical (peak) performance one could expect from the Xeon Phi and from the Xeon E5 (useful information has been underlined in the text).
2. What is the peak performance ratio between the Xeon Phi and the Xeon E5?
3. Compare the peak performance with the achieved performance. From the information given in the text, what can explain such disappointing performance?

**Question III.2.** In Section 4, the author perform a new series of experiments that are summarized on Figure 2. You will obviously justify any of your answer.

- ▷ List all the points that were changed between  $XP_1$  and  $XP_2$ .
- ▷ List all the points that should be changed/improved in Figure 1 and 2.
- ▷ What would you change/improve if you had to redo such experiments?

**Question III.3.** We consider now Figure 3 and  $XP_3$ . You will obviously justify any of your answer.

- ▷ If we denote by  $n$  the dimension of the square matrix used in the experiments, which formula corresponds the best to the information reported in Figure 3?

<input type="checkbox"/> $\frac{\alpha + \beta \times n}{n}$	<input type="checkbox"/> $\frac{n^2}{\alpha + \beta \times n^2}$	<input type="checkbox"/> $\frac{n^2}{\alpha \times n^3 + \beta \times n}$
<input type="checkbox"/> $\frac{\alpha + \beta \times n^2}{n^2}$	<input type="checkbox"/> $\frac{\alpha + \beta \times n^3}{n^3}$	<input type="checkbox"/> $\frac{\beta n^2}{\alpha \times n^2 \log(n) + n}$

What would be the interpretation of  $\alpha$  and  $\beta$ ? Note that if you think none of the above is a good model, you should explain why and provide your own model.

- ▷ What would you improve on Figure 3?
- ▷ What would you change/improve if you had to redo  $XP_2$  and  $XP_3$ ?

# Benchmarking the Intel<sup>®</sup> Xeon Phi<sup>™</sup> Coprocessor

Infrared Processing and Analysis Center, Caltech

F. Masci, 09/04/2013

## 1. Summary

This document summarizes our first experience with the Intel Xeon Phi. This is a coprocessor that uses Intel's Many Integrated Core (MIC) architecture to speed up highly parallel processes involving intensive numerical computations. The MIC coprocessor communicates with a regular Intel Xeon ("host") processor through its operating system. The Xeon Phi coprocessor is sometimes referred to as an "accelerator". In a nutshell, the Xeon Phi consists of 60 1.052 GHz cores each capable of executing four concurrent threads and delivering one teraflop of performance. For comparison, the host processor consists of 16 2.6 GHz cores, with two admissible threads per core. More details on the MIC hardware can be obtained from the references in Section 8.

Rather than present yet another guide on how to efficiently program for the Xeon Phi, our goal is to explore whether (and when) there are advantages in using the Xeon Phi for the processing of astronomical data, e.g., as in a production pipeline. Such pipelines are common-use at the Infrared Processing and Analysis Center (IPAC) and range from the instrumental calibration of raw image data, astrometry, image co-addition, source extraction, photometry, and source association. We also outline some lessons learned to assist future developers. Note: the findings and opinions reported here are exclusively the author's and do not reflect those of Intel or of any individual.

IPAC has recently acquired a single Xeon Phi card for preliminary benchmarking. We find in general that all existing heritage software based on C/C++/Fortran code can be made to run natively on the Xeon Phi with no recoding. However, whether it will run optimally to fully exploit the MIC architecture is a different question entirely. The answer is usually no. Even software that has been extensively multithreaded to utilize a multicore processor isn't guaranteed to run faster on the Xeon Phi than on a regular Intel Xeon machine. In fact, depending on memory and/or disk I/O usage, it can be much slower.

The key is to make efficient use of the Xeon Phi MIC architecture. This is not designed to handle jobs that are memory (primarily RAM) intensive. It is designed to utilize wide vector instruction units for floating point arithmetic (see below for details). Therefore, the types of problems the Xeon Phi is well suited for are intensive numerical computations with a low memory bandwidth. Additionally, the computations need to use one of the highly optimized vector math libraries that were implemented using assembly language constructs tuned specifically for the Xeon Phi architecture. Knowing this programming model beforehand can assist a developer to design software such that segments with intense numerical calculations can be offloaded to the Xeon Phi to be accelerated. The host processor then does most of the data I/O and memory management.

This document is organized as follows:

- Section 2 - Benchmarks conducted
- Section 3 - Programming philosophy, execution modes, and some lessons learned
- Section 4 - Preliminary "naïve" testing
- Section 5 - Reaching One Teraflop performance
- Section 6 - Conclusions and words of wisdom
- Section 7 - References and further reading
- Section 8 - Appendices: the author's "simple" benchmarking codes

## 2. Benchmarks conducted

To demonstrate that our Xeon Phi card performs according to the vendor's specifications and to explore its performance relative to the Xeon host, we ran two separate experiments :

1. A preliminary "naïve" test that multiplies two large matrices using explicit looping of matrix elements where loops were multithreaded using OpenMP pragmas. See Section 4. This makes heavy use of single-precision floating point arithmetic. A similar version with computations in double-precision yielded similar results. This code was compiled to run in "native-MIC" mode (see Section 3 for definition), "MIC-offload" mode, and "Host-exclusive" mode.
2. A test that also multiplies two large matrices but now using the highly (Phi-)optimized Math Kernel Library (MKL) to perform the matrix calculations. This also uses OpenMP to assist with multithreading. See Section 5. This code was compiled to run in "native-MIC mode" (see Section 3 for definition), "MIC-offload" mode, and "Host-exclusive" mode.

## 3. Programming philosophy, execution modes, and some lessons learned

The literature and various online documents broadly outline some "best practices" when programming for the Xeon Phi. But it's not all complete, nor collected in one place, or detailed enough to use in a practical sense, especially for a novice. Here's a summary of the programming models as well as some not-so-obvious ones that were (re)discovered from experimentation.

- The GNU compilers (e.g., gcc) are not suitable for compiling code for the Xeon Phi (for both native execution and offloading from the host). The Intel Compilers (e.g., icc) have the requisite optimization flags for code to execute efficiently on the Xeon Phi. This was discovered the hard way.
- The Xeon Phi is intended for highly parallelized numerical computations whose runtime scales linearly up to some maximum number threads that can be run concurrently on any processor, whether it's a regular Intel Xeon (host) or a Phi. For the Xeon Phi, the goal is to observe perfect linear scaling up to 120 threads (2 threads per core) before using additional Phi cards in a cluster setup. It is recommended that this scalability be demonstrated first on regular host processors before utilizing the Xeon Phi.
- Code that does heavy numerical work should also have a low memory bandwidth. Local CPU caches should be used as much as possible if there is frequent access to data, not the main RAM. This is referred to as maintaining "locality of reference". It is advised to keep all memory access within the L2 cache if possible. L2 is about 25 MB over all 60 cores on the Xeon Phi. I.e., there's little wiggle room - one of the challenges of optimizing code for the Xeon Phi.

- Make efficient use of the 64 byte-wide vector units (as opposed to 32 bytes on regular Xeon CPUs). I.e., 16 floats (or 8 doubles) can fit in the registers and be operated on simultaneously. Last, note that both the Xeon E5 and the Xeon cores have two vector units. The Single Instruction Multiple Data (SIMD) instruction set used for vector units on the Xeon Phi is not SSE compliant (Streaming SIMD Extension). This is applicable to Intel Xeon E5 processors only. Some time was spent using SSE-optimized compilation flags for the Xeon Phi that actually lead to a degradation in performance. The key is to use no SSE-specific flags when compiling for the Xeon Phi.
- If you are compiling to run natively on the Xeon Phi, you can get better performance if memory is allocated such that it's aligned on 64 byte cache-line boundaries. This can be achieved using the `posix_memalign()` or `malloc_aligned()` functions. For even better performance, you can make the array size a multiple of 16 for single-precision floating point, or 8 for double-precision, otherwise the unused registers will be internally masked and some overhead is incurred.
- Allow loops in the code to be auto-vectorized (unrolled) by the compiler. This optimizes the processing for the SIMD vector architecture on each individual core. At times, the icc compiler can be stubborn and refuse to auto-vectorize a loop because it thinks there's some variable dependence. You can use "#pragma simd" to force auto-vectorization by the compiler provided you're certain it's safe to unroll the loop in question.
- The compiler-assisted auto-vectorization step mentioned above is usually not enough to fully optimize code for the Xeon Phi. It's a necessary step for the SIMD 64-byte vector processing on each core, but more work on explicit multithreading to utilize all cores is needed if you have loops operating on large arrays of data with no dependencies between the variables. The latter can be achieved using either the openMP library via "omp pragmas" or via the pthreads library directly.
- Use the highly optimized vector/matrix functions in Intel's Math Kernel Library (MKL) or lower level Vector Math library (VML). These libraries use optimized assembly language that takes advantage of the SIMD vector instruction set for the Xeon Phi. So far, I've only been able to run things faster on the Xeon Phi (than on the host) when using MKL routines. In the end, all arithmetic operations on data arrays can be performed using vector/matrix operations so it is advised to use these whenever possible. One finding with MKL (see benchmarking below) is that computations run more efficiently for relatively large input vectors and matrices. This establishes the power of the Xeon Phi.

We explored two compilation/execution modes while testing on the Xeon Phi:

1. "Native-MIC" mode where a fully-native Xeon Phi binary is compiled for execution on the MIC O/S directly. This is a good choice when the code has low data I/O and/or memory needs and consists primarily of intensive numerical computations.
2. "MIC-offload" mode where a heterogeneous binary is compiled to run on the host CPU with selected code segments offloaded to the Xeon Phi during runtime. The code segments intended for execution on the Xeon Phi are designated using offload pragmas targeted for the MIC. The advantage here is that the host processor can manage any high data/memory I/O needs while segments involving heavy numerical work can be "accelerated" on the Xeon Phi.



#### 4. Preliminary “naïve” testing

We first conducted a simple test that multiplies two matrices with  $3000 \times 3000$  elements each. We call this test “naïve” since it’s programmed using nested loops over each matrix element and carrying out the arithmetic as we go along, i.e., as the textbooks say when multiplying matrices. The multiplication operation was repeated 10 times and the performance metrics were averaged at the end. We have also parallelized the loop computations using OpenMP pragmas (see code in Section 9.a.). We were suspicious from the outset that this approach was non-optimal for the Xeon Phi, but we wanted to see where this would take us. Such code is typical of heritage software lying around for performing matrix computations (without the multithreading of course).

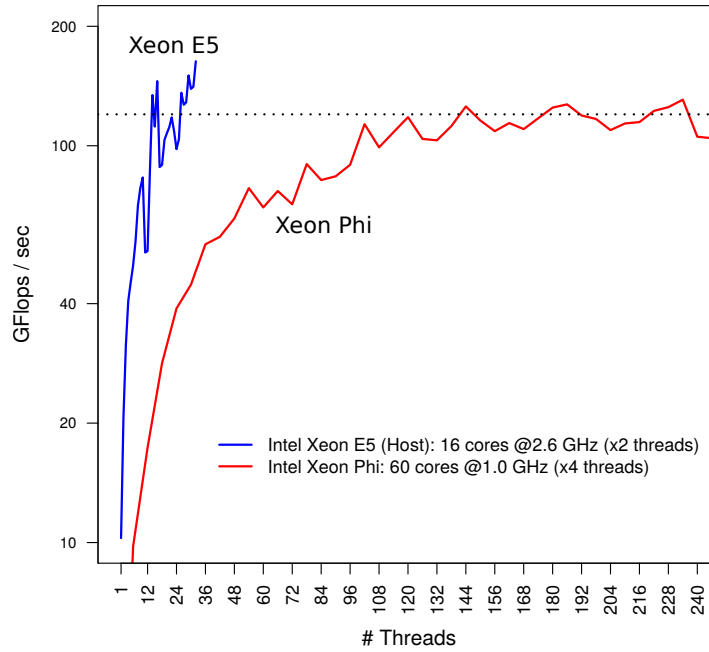
Test code “test1.c” (Section 9.a.) was compiled to run in “native-MIC” mode, “MIC-offload” mode, and “Host-exclusive” mode (see Section 3 for definitions). Figure 1 shows the number of Floating point Operations per second, or FLOPs/sec versus the number of threads spawned. All computations used single-precision floating point. Only the native-MIC (red) and Host-exclusive (blue) runs are shown. We found that the MIC-offload mode performed ~10% worse (lower flops) than the native-MIC run. This is probably due to the additional overhead in transferring (offloading) the large data arrays from the host’s main memory to the Xeon Phi.

Figure 1 shows that the maximum achievable performance on the Xeon Phi (~120 GFlops/sec) is comparable to that on the host (within the measured variance). This is contrary to expectations reported by the vendor and others in the literature.

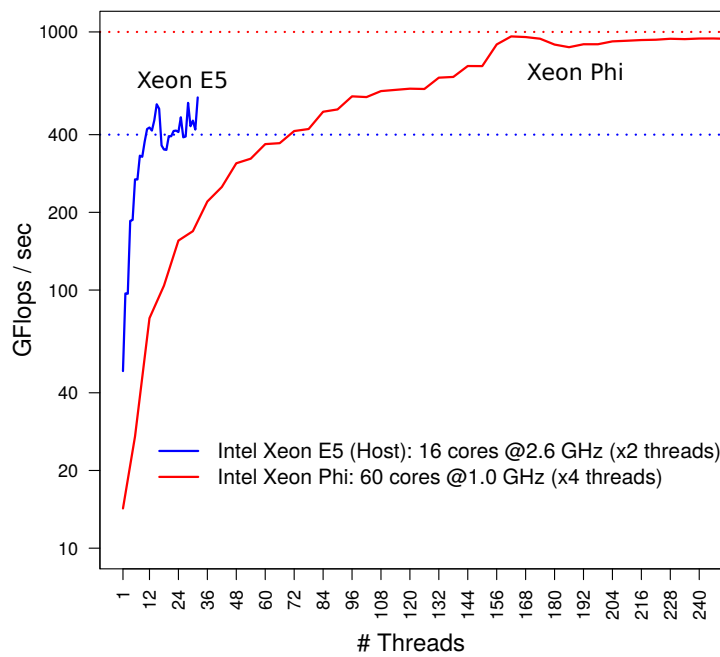
#### 5. Reaching One Teraflop performance

We decided to replace the “naïve” multiplication step in test1.c (Section 4) with the `s g e m m f ( )` matrix multiplication function from Intel’s Math Kernel Library (MKL). As outlined in Section 3, this library uses optimized assembly language constructs to take advantage of the SIMD vector instruction set for Xeon Phi coprocessors. Results for test code “test2.c” (Section 9.b.) are shown in Figure 2. The long sought-after one TFlops/sec performance metric is now achieved. Benchmarks advertised by the vendor and others peak at ~1.6 - 1.8 TFlops/sec, but no guidance is given on the set-up or circumstances under which these are achieved. Furthermore, relative to the host processor, we find the Xeon Phi performs ~2.5 better! It appears the use of MKL routines for floating-point arithmetic is pivotal in exploiting the Xeon Phi.

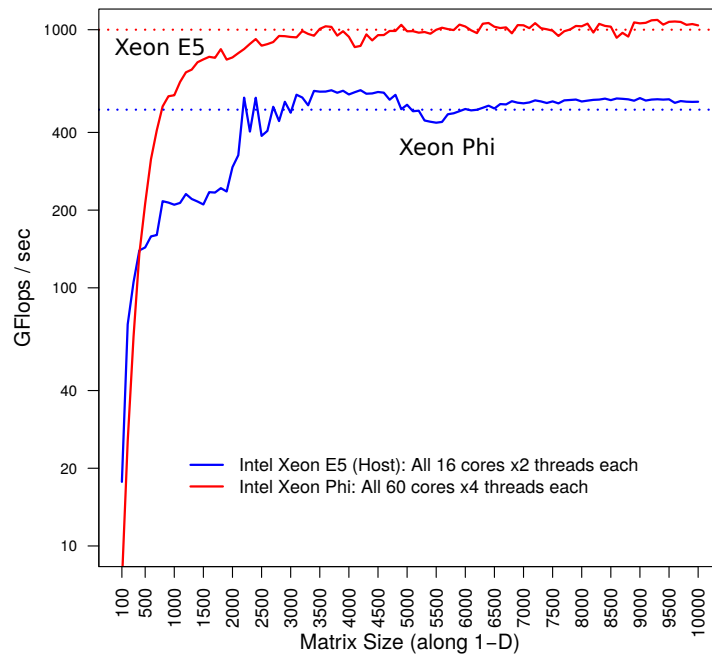
As a further test, we fixed the number of threads to the maximum possible on the Xeon Phi (240) and the host processor (32) and explored the performance as a function of matrix size using test code test2.c. Figure 3 shows the result. It appears that we reach maximum performance when the matrix size exceeds  $2500 \times 2500$  elements. That is, both the Xeon host and Phi perform best (with better throughput and efficiency) when the number of computations is large. This shows the Xeon Phi is best suited for large/heavy numerical problems and not worth the effort for small ones where data/memory I/O is likely to dominate (relatively speaking).



**Figure 1: Single-precision floating-point performance for the simple “naive” matrix multiplication test code (test1.c) involving two  $3000 \times 3000$  matrices. Blue curve is for the “Host-exclusive” runs and red curve is for the Xeon Phi runs using “native-MIC” mode. Dotted line is at 120 GFlops/sec.**



**Figure 2: Single-precision floating-point performance for our optimized matrix multiplication test code (test2.c) involving two  $3000 \times 3000$  matrices. Blue curve is for the “Host-exclusive” runs and red curve is for the Xeon Phi runs using “native-MIC” mode. Horizontal dotted lines (at 400 and 1000 GFlops/sec) are to guide the eye.**



**Figure 3: Single-precision floating-point performance for our optimized matrix multiplication test code (test2.c) as a function of matrix size. Blue curve is for the “Host-exclusive” runs and red curve is for the Xeon Phi runs using “native-MIC” mode. The number of concurrent threads used was fixed at the maximum values of 32 and 240 for the host and Xeon Phi respectively.**

## 7. Conclusions and words of wisdom

While most heritage software will run on the Xeon Phi with little effort, there's no guarantee it will run optimally to fully exploit its architecture, even if the segments that carry out intensive numerical computations have been highly parallelized. Some recoding using offload pragmas targeted for the MIC and use of vector/matrix libraries optimized for numerical work is inevitable. This could be expensive for existing codes that are not structured in a manner to trivially separate the heavy data/memory I/O steps from the pure computational ones.

Approximately one teraflop can be obtained on the Xeon Phi when using optimized vector-math libraries. Even then, we find this is only a factor of 2 to 2.5 times greater than that achieved on regular Intel Xeon E5 processors with 16 2.6 GHz cores. Is a factor of 2 to 2.5 improvement really worth the human labor to recode existing "numerical-heavy" software for optimal execution on Xeon Phi cards? This depends on the developer's experience with the software and its complexity, but more importantly, on the type of software being run. As it stands, much of the software executed in IPAC's mission-production pipelines (*Spitzer*, WISE, and PTF for example) are primarily data I/O limited and at the mercy of file servers cooperating and responding in sync with the CPUs.

However, this doesn't mean there's no place for Xeon Phi cards at IPAC. Future developers should use the Xeon Phi when it makes sense to do so. Existing heritage code can be judiciously recycled with an eye for "accelerating" heavy computational steps with the Xeon Phi. There's always a gain, but as mentioned above, the factor of 2 to 2.5 improvement over regular processors won't buy you much if most of the time is spent moving data off/onto disks and/or in/out of RAM. The compute-heavy steps (until now) are usually a small fraction (<~ 30%) of pipeline processing budgets for astronomical applications that go from raw-image data to source-catalogs. Regular Intel Xeon processors are getting faster and more efficient at managing memory. Nonetheless, there will be scientists who would benefit from the Xeon Phi by running customized code involving heavy numerical work, e.g., N-body simulations, gravitational lensing shear calculations, or radiative transfer models.

## 8. References and further reading

The following links and documents contain good examples and guidelines for a novice when programming for the Intel Xeon Phi.

- <http://software.intel.com/sites/default/files/article/335818/intel-xeon-phi-coprocessor-quick-start-developers-guide.pdf>
- <http://software.intel.com/en-us/articles/building-a-native-application-for-intel-xeon-phi-coprocessors>
- <http://www.drdoobs.com/parallel/programming-intels-xeon-phi-a-jumpstart/240144160>
- <http://www.prace-project.eu/IMG/pdf/Best-Practice-Guide-Intel-Xeon-Phi.pdf>
- <http://software.intel.com/sites/default/files/article/366893/offload-runtime-for-the-intel-xeon-phi-m-coprocessor.pdf>
- [http://research.colfaxinternational.com/file.axd?file=2013%2F5%2FCofax\\_Static\\_Libraries\\_Xeon\\_Phi.pdf](http://research.colfaxinternational.com/file.axd?file=2013%2F5%2FCofax_Static_Libraries_Xeon_Phi.pdf)
- <https://hpcforge.org/plugins/mediawiki/wiki/pracewp8/images/6/68/XeonPhi.pdf>
- <http://software.intel.com/en-us/articles/getting-started-with-openmp>
- [http://d3f8ykwhia686p.cloudfront.net/1live/intel/An\\_Introduction\\_to\\_Vectorization\\_with\\_Intel\\_Compiler\\_021712.pdf](http://d3f8ykwhia686p.cloudfront.net/1live/intel/An_Introduction_to_Vectorization_with_Intel_Compiler_021712.pdf)

## 9. Appendices

Below are the C codes used in the “simple” benchmarking tests presented in Sections 4 and 5. These were written by combining various code snippets (with much experimentation) from the links in Section 8. Further below are the compilation, environment variables and execution command lines used.

### 9.a. test1.c: code used to generate Figure 1 (Section 4)

```

#ifndef MIC_DEV
#define MIC_DEV 0
#endif

#include <stdio.h>
#include <stdlib.h>
#include <omp.h>
#include <mkL.h> /* needed for the dsecnd() timing function. */
#include <math.h>

/* Program test1.c: multiply two matrices using explicit looping of elements. */
/*-----*/
/* Simple "naive" method to multiply two square matrices A and B
to generate matrix C. */

void myMult(int size,
            float (* restrict A)[size],
            float (* restrict B)[size],
            float (* restrict C)[size])
{
    #pragma offload target(mic:MIC_DEV) in(A:length(size*size)) \
                                in( B:length(size*size)) \
                                out(C:length(size*size))
    {
        /* Initialize the C matrix with zeroes. */

        #pragma omp parallel for default(none) shared(C,size)
        for(int i = 0; i < size; ++i)
            for(int j = 0; j < size; ++j)
                C[i][j] = 0.f;

        /* Compute matrix multiplication. */

        #pragma omp parallel for default(none) shared(A,B,C,size)
        for(int i = 0; i < size; ++i)
            for(int k = 0; k < size; ++k)
                for(int j = 0; j < size; ++j)
                    C[i][j] += A[i][k] * B[k][j];
    }
}

/*-----*/
/* Read input parameters; set-up dummy input data; multiply matrices using
the myMult() function above; average repeated runs. */

int main(int argc, char *argv[])
{
    if(argc != 4) {
        fprintf(stderr, "Use: %s size nThreads nIter\n", argv[0]);
        return -1;
    }

    int i, j, nt;
    int size=atoi(argv[1]);
    int nThreads=atoi(argv[2]);
    int nIter=atoi(argv[3]);

    omp_set_num_threads(nThreads);

    /* when compiled in "mic-offload" mode, this memory gets allocated on host,
when compiled in "mic-native" mode, it gets allocated on mic. */

    float (*restrict A)[size] = malloc(sizeof(float)*size*size);
    float (*restrict B)[size] = malloc(sizeof(float)*size*size);
    float (*restrict C)[size] = malloc(sizeof(float)*size*size);

    /* this first pragma is just to get the actual #threads used

```

```

    (sanity check). */
#pragma omp parallel
{
    nt = omp_get_num_threads();

    /* Fill the A and B arrays with dummy test data. */
#pragma omp parallel for default(none) shared(A,B,size) private(i,j)
    for(i = 0; i < size; ++i) {
        for(j = 0; j < size; ++j) {
            A[i][j] = (float)i + j;
            B[i][j] = (float)i - j;
        }
    }

    /* warm up run to overcome setup overhead in benchmark runs below. */
    myMult(size, A,B,C);

    double aveTime,minTime=1e6,maxTime=0.;

    /* run the matrix multiplication function nIter times and compute
       average runtime. */
    for(i=0; i < nIter; i++) {
        double startTime = dsecnd();
        myMult(size, A,B,C);
        double endTime = dsecnd();
        double runtime = endTime-startTime;
        maxTime=(maxTime > runtime)?maxTime:runtime;
        minTime=(minTime < runtime)?minTime:runtime;
        aveTime += runtime;
    }
    aveTime /= nIter;

    printf("%s nThreads %d matrix %d maxRT %g minRT %g aveRT %g GFlop/s %g\n",
           argv[0],nt,size,maxTime,minTime,aveTime, 2e-9*size*size*size/aveTime);

    free(A);
    free(B);
    free(C);

    return 0;
}

```

### 9.b. test2.c: Code used to generate Figures 2 and 3 (Section 5)

```

#include <stdio.h>
#include <stdlib.h>
#include <omp.h>
#include <mkl.h>
#include <math.h>

/* Program test2.c: multiply two matrices using a highly thread-optimized
   routine from the Intel Math Kernel Library (MKL). */

/*-----*/
/* Multiply two square matrices A and B to generate matrix C using the
   optimized sgemm() routine (for single precision floating point) from MKL. */

float fastMult(int size,
              float (* restrict A)[size],
              float (* restrict B)[size],
              float (* restrict C)[size],
              int nIter)
{
    float (*restrict At)[size] = malloc(sizeof(float)*size*size);
    float (*restrict Bt)[size] = malloc(sizeof(float)*size*size);
    float (*restrict Ct)[size] = malloc(sizeof(float)*size*size);

    /* transpose input matrices to get better sgemm() performance. */

#pragma omp parallel for
    for(int i=0; i < size; i++)
        for(int j=0; j < size; j++) {
            At[i][j] = A[j][i];
            Bt[i][j] = B[j][i];
        }
}

```

```

    }

    /* scaling factors needed for sgemm(). */
    float alpha = 1.0f;
    float beta = 0.0f;

    /* warm up run to overcome setup overhead in benchmark runs below. */
    sgemm("N", "N", &size, &size, &size, &alpha,
          (float *)At, &size, (float *)Bt, &size, &beta, (float *) Ct, &size);

    double StartTime=dsecnd();

    for(int i=0; i < nIter; i++)
        sgemm("N", "N", &size, &size, &size, &alpha,
              (float *)At, &size, (float *)Bt, &size, &beta, (float *) Ct, &size);

    double EndTime=dsecnd();

    float tottime = EndTime - StartTime;
    float avgttime = tottime / nIter;
    printf("tot runtime = %f sec\n", tottime);
    printf("avg runtime per vec. mult. = %f sec\n", avgttime);
    float GFlops = (2e-9*size*size*size)/avgttime;

    free(At);
    free(Bt);
    free(Ct);

    return ( GFlops );
}

/*-----*/
/* Read input parameters; set-up dummy input data; multiply matrices using
   the fastMult() function above; average repeated runs therein. */

int main(int argc, char *argv[])
{
    if(argc != 4) {
        fprintf(stderr,"Use: %s size nThreads nIter\n",argv[0]);
        return -1;
    }

    int i,j,nt;
    int size=atoi(argv[1]);
    int nThreads=atoi(argv[2]);
    int nIter=atoi(argv[3]);

    omp_set_num_threads(nThreads);
    mkl_set_num_threads(nThreads);

    /* when compiled in "mic-offload" mode, this memory gets allocated on host,
       when compiled in "mic-native" mode, it gets allocated on mic. */

    float (*restrict A)[size] = malloc(sizeof(float)*size*size);
    float (*restrict B)[size] = malloc(sizeof(float)*size*size);
    float (*restrict C)[size] = malloc(sizeof(float)*size*size);

    /* this first pragma is just to get the actual #threads used
       (sanity check). */

    #pragma omp parallel
    {
        nt = omp_get_num_threads();

        /* Fill the A and B arrays with dummy test data. */
        #pragma omp parallel for default(none) shared(A,B,size) private(i,j)
        for(i = 0; i < size; ++i) {
            for(j = 0; j < size; ++j) {
                A[i][j] = (float)i + j;
                B[i][j] = (float)i - j;
            }
        }
    }

    /* run the matrix multiplication function nIter times and average
       runs therein. */

    float Gflop = fastMult(size,A,B,C,nIter);

```

```

printf("size = %d x %d; nThreads = %d; #GFlop/s = %g\n",
      size, size, nt, Gflop);

free(A);
free(B);
free(C);

return 0;
}

```

### 9.c. Compilation and runtime environment for codes above

All environment variables and example command lines below were set/executed directly on the Host (Xeon E5) processor. The environment variables containing “MKL” and compiler option “-mkl” can be omitted for the “naïve” test code of Section 9.a. Furthermore, for the code in Section 9.b., the MKL\_MIC\_ENABLE environment variable was used to control the offloading of MKL-specific routines to the Xeon Phi. This avoids explicit use of offload “target(mic)” pragmas for the MKL routines.

Depending on the execution mode, three different binaries were generated for each test?.c code above: *test\_mic*, *test\_mic\_offload*, *test\_host*. The input arguments are “size of square matrix along one dimension”, “number of concurrent threads”, “number of iterations for runtime averaging”.

#### Full native execution on the Xeon Phi:

```

source /opt/intel/bin/compilervars.csh intel64

icc -mkl -O3 -mmic -openmp -L /opt/intel/lib/mic -Wno-unknown-pragmas -std=c99
-vec-report2 -liomp5 -o test_mic test.c

setenv SINK_LD_LIBRARY_PATH
/opt/intel/composer_xe_2013/lib/mic:/opt/intel/mkl/lib/mic;
setenv MKL_MIC_ENABLE 1;
setenv MIC_ENV_PREFIX MIC;
setenv MIC_KMP_AFFINITY "granularity=thread,balanced";
setenv MIC_USE_2MB_BUFFERS 32K
setenv MIC_MKL_DYNAMIC false

/opt/intel/mic/bin/micnativeloadex ./test_mic -a "3000 240 10";

```

#### Heterogeneous binary with offloading of selective code to the Xeon Phi:

```

source /opt/intel/bin/compilervars.csh intel64

icc -offload-option,mic,compiler,"-mP2OPT_hpo_vec_check_dp_trip=F -fimf-
precision=low -fimf-domain-exclusion=15 -opt-report 1" -
mP2OPT_hlo_pref_issue_second_level_prefetch=F -
mP2OPT_hlo_pref_issue_first_level_prefetch=F -vec-report2 -O3 -openmp -
intel-extensions -opt-report-phase:offload -openmp-report -mkl -Wno-unknown-
pragmas -std=c99 -o test_mic_offload test.c

setenv MKL_MIC_ENABLE 1;
setenv MIC_ENV_PREFIX MIC;
setenv MIC_KMP_AFFINITY "granularity=thread,balanced";
setenv MIC_USE_2MB_BUFFERS 32K
setenv MIC_MKL_DYNAMIC false;
setenv KMP_AFFINITY "granularity=thread,scatter";

./test_mic_offload 3000 240 10

```



**Host-only execution (no use of the Xeon Phi):**

```
source /opt/intel/bin/compilervars.csh intel64

icc -xhost -mkl -O3 -no-offload -openmp -Wno-unknown-pragmas -std=c99 -vec-
report2 -o test_host test.c

setenv MKL_MIC_ENABLE 0;
setenv KMP_AFFINITY "granularity=thread,scatter";
setenv USE_2MB_BUFFERS 32K;
setenv MKL_DYNAMIC false;

./test_host 3000 32 10
```