

# Performance Evaluation of ParalleX Execution model on Arm-based Platforms

A Preprint - October 26, 2020

Nikunj Gupta\*<sup>§</sup>, Rohit Ashiwal\*, Bine Brank<sup>‡</sup>, Sateesh K. Peddoju\* Dirk Pleiter<sup>‡</sup>

\*Dept. of CSE, IIT Roorkee, Roorkee, India, Email: gnikunj@cct.lsu.edu;{rashiwal,sateesh}@cs.iitr.ac.in

<sup>‡</sup>JSC, Jülich Research Centre, 52425 Jülich, Germany, Email: {b.branc,d.pleiter}@fz-juelich.de

<sup>§</sup>The STE||AR Group, <http://stellar-group.org>

**Abstract**—The HPC community shows a keen interest in creating diversity in the CPU ecosystem. The advent of Arm-based processors provides an alternative to the existing HPC ecosystem, which is primarily dominated by x86 processors. In this paper, we port an Asynchronous Many-Task runtime system based on the ParalleX model, i.e., High Performance ParalleX (HPX), and evaluate it on the Arm ecosystem with a suite of benchmarks. We wrote these benchmarks with an emphasis on vectorization and distributed scaling. We present the performance results on a variety of Arm processors and compare it with their x86 brethren from Intel. We show that the results obtained are equally good or better than their x86 brethren. Finally, we also discuss a few drawbacks of the present Arm ecosystem.

**Index Terms**—Asynchronous Many-Task, HPX, Parallel Computing, ParalleX

## I. INTRODUCTION

The High-Performance ecosystem has shown a keen interest in shifting from the traditional x86 architecture to Arm-based processors. The Japanese exascale system called Fugaku [1] is based on the A64FX processor from Fujitsu. In Europe, the European Processor Initiative (EPI) is working on a processor that will (among others) include Arm-based cores and is positioned as an HPC technology. In the USA, Sandia National Laboratories has deployed a large-scale system based on the ThunderX2 processor from Marvell, which is based on the Armv8 ISA.

This development triggers the question of whether the HPC software ecosystem is ready for exploiting Arm. This concerns, in particular, a recent extension of the Arm ISA, which is called Scalable Vector Extension (SVE) [2]. SVE is vector-length agnostic, unlike AVX or AVX512 from Intel, where SIMD width is fixed to 256 and 512 bit, respectively. Size has some significant consequences when programming for SVE. For AVX and AVX512 (and similar SIMD ISAs), data types that are defined (e.g., `__m512d` for a vector of eight doubles) have a size known at compile-time. For SVE, this is a priori, not the case, as the vector length is only known at runtime.

These hardware complexities are matched with software complexities. Operating systems, compiler, and library support

are required to provide a functional environment that supports large-scale HPC applications and ensure they can both be easily ported to such new hardware and exploit it efficiently. One such class of applications is the one exploiting parallelism through task-based programming. Asynchronous Many-Task (AMT) runtime system models task-based programming and offers an alternative to the conventional programming models like Message Passing (MPI). In an AMT model, the program can be broken down into tasks, with each task having a dependency on some other task generating a data flow based process. During program execution, these tasks are launched arbitrarily based on the input data and the DAG generated, enabling multiple concurrent tasks running as computation kernels. The scheduler deals with the load imbalance. These characteristics of the AMT model poses it as a viable alternative in an era where future algorithms are expected to feature an increased dynamic behavior and low uniformity. In this paper, we explore an AMT model on Arm-based processors.

Section II talks about Related Work in porting and evaluating the Arm processors. Section III discusses the ParalleX execution model and HPX. Section IV touches on key concepts required to understand the benchmarks and the results. Section V describes the benchmarks and system setup, and Section VI discusses the results.

## II. RELATED WORK

S. McIntosh-Smith et al. [3] recorded the initial set of performance evaluations on mainstream Arm processors on large HPC systems. These evaluations showcased performance and cost benefits for a class of applications. The result considers performance on a single-node. They later released distributed application results in [4]. Jackson et al. [5] investigated the performance of distributed memory communications (MPI), as well as scientific applications utilizing MPI on ThunderX2. We also found that Mont-Blanc project [6] investigated energy consumptions during the execution of benchmarks and mini-apps.

While there has been decent research and performance evaluation on the conventional computation models (such as OpenMP + MPI), there are no performance numbers available for AMT runtime systems on mainstream Arm processors.

©2020 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

### A. Our contribution

In this paper, we execute several benchmarks on an AMT based on the ParalleX model, i.e. HPX. We investigate both distributed and shared memory models with a special emphasis on vectorization. Furthermore, we provide performance comparisons of mainstream Arm processors, such as Kunpeng 916 (HiSilicon), ThunderX2 (Marvell) and A64FX (Fujitsu), with their x86 brethren from Intel.

## III. BACKGROUND

### A. ParalleX Execution Model

ParalleX execution model [7] was devised to address the critical bottlenecks of exascale HPC systems namely:

- Starvation - insufficient parallelism
- Latencies - time-distance delay of remote resource accesses
- Overheads - extra work for management of parallel actions and resources on the critical path that are not necessary in the sequential variant
- Contention - delays due to lack of availability of resources

The ParalleX model offers an alternative to conventional computational models, such as MPI.

ParalleX improves the efficiency of the application by reducing synchronization and scheduling overheads. Resource utilization is achieved through increased asynchrony. Contention overheads are significantly reduced by employing adaptive scheduling and routing. This technique is instrumental in handling memory bank conflicts. Data directed computing using message-driven computation and lightweight synchronization mechanisms results in visible scalability improvements, at least for certain classes of problems. ParalleX achieves power reductions by reducing extraneous calculations and data movements.

### B. High Performance ParalleX Runtime System

High Performance ParalleX (HPX) [8]–[12] is the first open-source implementation of ParalleX execution model. HPX exposes an ISO C++ standard conforming API, which enables wait-free asynchronous parallel programming, including futures, channels, and other synchronization primitives. Being C++ standards conforming, HPX can run on a single machine as well as a cluster with thousands of nodes. Figure 1 gives us the architecture of HPX.

HPX utilizes Active Global Address Space (AGAS) for addressing any HPX object globally. Every object is assigned a Global Identifier (GID) that persists until object destruction. Furthermore, AGAS supports load balancing through object migration. HPX utilizes lightweight threads, called HPX threads, that are executed and scheduled on top of OS threads. Local Control Object (LCO) is a family of synchronization functions used to synchronize tasks generated by the application. HPX has an active-message networking layer that ships functions to the objects they operate on, i.e., the Parcel subsystem.

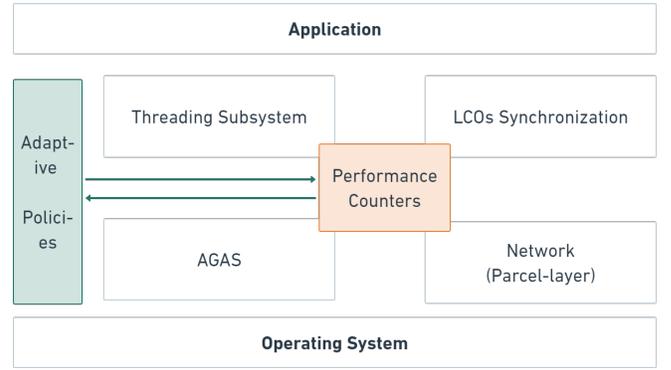


Fig. 1: HPX Architecture Overview

### C. Roofline Model

The Roofline Model [13]–[15] provides performance estimation keeping in mind the bottlenecks and bounds to predict a more realistic performance estimation. Usually, the performance is expressed as a function of peak Computational Performance (CP) of the architecture and the reachable peak I/O Bandwidth (BW). The CP is the maximum number of floating-point operations that the processor can achieve.

The Arithmetic Intensity (AI) is the number of operations executed per byte accessed from the main memory. It reveals the complexity of the algorithm. The model identifies the CP or the I/O BW as limiting factors. Ergo, the maximal attainable performance is always below the roofline obtained from both parameters:

$$\text{Attainable Performance} = \min(\text{CP}, \text{AI} \times \text{BW}) \quad (1)$$

## IV. STENCIL CODES

In this paper, we explore a 1D stencil solver for the heat equation in one dimension and a 2D stencil solver implementing a Jacobi solver.

### A. 1D Heat Equation

The diffusion equation, also known as heat equation in one dimension is given by:

$$\frac{\partial u}{\partial t} = \alpha \frac{\partial^2 u}{\partial x^2} \quad (2)$$

where  $u(x, t)$  is unknown and  $\alpha$  is the diffusion constant. On discretizing the domain and replacing the derivatives by finite differences, one can derive a 3-point stencil to update a single cell as:

$$\begin{aligned} T_{new}(x, y) &= T_{old}(x, y) \\ &+ \alpha \frac{\Delta t}{\Delta x^2} [T_{old}(x-1, y) - 2 * T_{old}(x, y) \\ &+ T_{old}(x+1, y)] \end{aligned} \quad (3)$$

## B. Jacobi Method for Linear Equation

Jacobi solvers are a class of iterative solvers for linear equations given by:

$$Ax = b, A \in \mathbb{R}^{N \times N}, x \text{ and } b \in \mathbb{R}^N$$

On a uniform grid with Dirichlet boundary conditions, the linear equation can be derived into a 5-point stencil to update a single cell as:

$$T_{new}(x, y) = [T_{old}(x, y + 1) + T_{old}(x, y - 1) + T_{old}(x + 1, y) + T_{old}(x - 1, y)]/4 \quad (4)$$

A detailed derivation is available at [16].

## V. BENCHMARKS

This Section discusses the implementation details of benchmarks, machine architectures and HPX configuration.

### A. 1D Stencil

For the 1D stencil, we implement a fully distributed 1D heat equation solver. This benchmark accurately measures the total execution time of the application. We report application execution time over kernel performance to investigate how the complete application scales on a distributed setting.

We run the application for two cases, i.e., weak and strong scaling. For weak scaling, we start with 480 million stencil points, and another 480 million stencil points are added for every node. For strong scaling, we benchmark with 1.2 billion stencil points. The benchmarks iterate over a hundred time steps.

```

1  hpx::parallel::for_each(policy, begin(range),
   end(range),
2  [&U, local_nx, nlp, t] (std::size_t i)
3  {
4    if (i == 0)
5      stencil_update(U, 1, local_nx, t);
6    else if (i == nlp-1)
7      stencil_update(U, i * local_nx,
8        (i + 1) * local_nx - 1, t);
9    else if (i > 0 && i < nlp-1)
10     stencil_update(U, i * local_nx,
11       (i + 1) * local_nx, t);
12  }
13 );
```

Listing 1: 1D stencil solver implemented in HPX

Listing 1 shows our implementation of the 1D stencil kernel. `stencil_update` takes in the two grids and applies the operation defined in Equation 3.

### B. 2D stencil

For the 2D stencil, we implement a shared-memory based 2D stencil implementing Jacobi method. We make use of explicit vectorization to compare performance with an auto vectorized scalar code from the compiler. Assuming that the cache size is large enough to accommodate three rows of the grid, three memory transfers have to be done for every iteration, implying that for a double, a total of 24 Bytes

are fetched from the main memory for every Lattice Site Update (LUP). Similarly, for a float, a total of 12 Bytes are fetched from the main memory. Thus, the Arithmetic Intensity (AI) for floats and doubles are 1/12 LUP/Byte and 1/24 LUP/Byte, respectively. The low arithmetic intensity makes the application memory bound for a broad class of processors.

```

1  template <typename Container>
2  void stencil_update(array_t<Container>& U,
   size_t ny, size_t t)
3  {
4  Grid<Container>& curr = U[t % 2];
5  Grid<Container>& next = U[(t + 1) % 2];
6
7  size_t row_length = curr.row_size();
8
9  #pragma unroll
10 for (size_t nx = 1; nx < row_length-1; ++nx)
11 {
12 // Stencil operation
13 next.in(nx, ny) = (curr.in(nx-1, ny) +
   curr.in(nx+1, ny) + curr.in(nx, ny-1)
   + curr.in(nx, ny+1)) * 0.25f;
14 }
15
16 // Maintain the halo in case of simd
17 if (std::is_same<typename Container::
   value_type, nsimd::pack<typename get_type
   <typename Container::value_type>::type
   >>::value)
18   helper<Container>::shuffle(next, ny);
19 }
20
21 // Call to stencil_update
22 hpx::util::high_resolution_timer t;
23 for (size_t t = 0; t < steps; ++t)
24 {
25   hpx::parallel::for_each(
26     policy, begin(range), end(range),
27     [&U, t] (size_t i)
28     {
29       stencil_update<Container>(U, i, t);
30     });
31 }
32 t.elapsed();
```

Listing 2: Generic 2D stencil kernel implemented with HPX and NSIMD. Container can be an STL vector of scalar types, or an STL vector of vector types.

Listing 2 shows parts of the C++ code written in HPX to generate a generic 2D stencil kernel that supports all floating-point data types. `Grid` is our custom class that abstracts away the data layout of our stencil. We make use of the Virtual Node Scheme [17] to generate a SIMD data layout of the stencil. To maintain this data layout, we need to update the boundaries to keep consistent data. This is done by shuffling the halo vectors to update according to the changes brought after executing the current time step (see Line 18). At Line 17, we use C++ type traits and our custom `get_type` meta-class to identify if a type is scalar or vector.

We run the application with strong scaling since we are inclined to investigate the performance of the kernel. We report performance numbers for a grid size of  $8192 \times 131072$ . The

TABLE I: Specification of the Arm and x86 nodes utilised in the benchmarks.

	Intel Xeon E5-2660 v3	HiSilicon Kunpeng 916	Marvell ThunderX2	Fujitsu (FX1000) A64FX
<b>Processor Clock Speed</b>	2.6GHz	2.4GHz	2.4GHz	2.2Ghz
<b>Cores per processors</b>	10	64	32	48 (compute) + 4 (helper)
<b>Processors per node</b>	2	1	1	1
<b>Threads per core</b>	2	1	4	1
<b>Vectorization</b>	Double AVX2 Pipeline	Single NEON Pipeline	Double NEON Pipeline	Double SVE 512-bit
<b>Double Precision FLOPS per cycle</b>	16	4	8	32
<b>Peak Performance in GFLOP/s</b>	832	614	1228	3379

row size has been chosen such that it fits easily in caches for our described assumptions to be true. Furthermore, the grid size is chosen to be large enough in an attempt to keep all processing cores busy. The benchmark iterates over a hundred time steps.

## VI. SYSTEM SETUP

We use the following clusters to access processors:

- Juawei prototype cluster, JSC. We use the cluster to access Intel Xeon E5 2660 v3 and HiSilicon Kunpeng 916 nodes.
- Sage prototype cluster, JSC. We use the cluster to access Marvell ThunderX2 nodes.
- A64FX prototype cluster, Fujitsu. We use the cluster to access Fujitsu FX1000 A64FX nodes.

Table I lists the specification for all processors. Table II gives an overview of the dependencies of our benchmark. We choose GCC over Arm HPC compiler and Fujitsu Compiler. Arm compilers implement SVE data types using `__sizeless_struct`. `__sizeless_struct` determines vector length at runtime, thus making the SVE code portable in the sense that an executable can run all vector lengths supported by SVE (i.e., 128bits through 2048bits). However, one cannot wrap `__sizeless_struct` into a struct or a class unless the struct is also defined as `__sizeless_struct`. Since we make heavy use of the C++ STL vector and our custom class, we have to determine type and length at compile time. As SVE is bleeding-edge technology, there is no synergy between the implementation details of various compilers. Currently, only GCC provides the choice of passing SVE vector length at compile time. Therefore, we use GCC.

TABLE II: Benchmark dependencies Configuration

Package Name	Version
GCC	10.1
hwloc	2.1
jemalloc	5.2.1
boost	1.66
HPX	commit c62d992
NSIMD	commit d4f9fc5
PAPI	6.0.0

We make use of NSIMD<sup>1</sup> for explicit vectorization. We chose NSIMD over explicitly using vector intrinsics because

<sup>1</sup><https://github.com/agenium-scale/nsimd>

NSIMD provides a portable code supporting a wide array of vectorizations, including SVE. Portability allowed us to use the same code on all machine architectures and underlying vectors. To our knowledge, NSIMD and Inastemp [18] are the only C++ libraries for explicit vectorization that support SVE data types.

Memory Bandwidth using STREAM COPY Benchmark

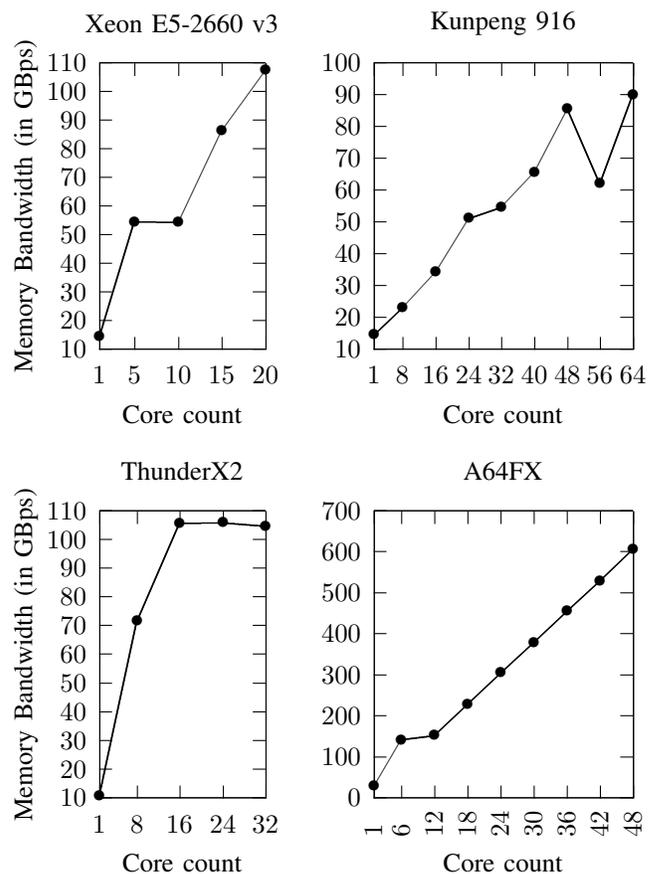


Fig. 2: Memory Bandwidth results using the STREAM COPY Benchmark with an array size of 128 million elements

Apart from the specifications in Table I, we measure the memory bandwidth using the STREAM COPY bench-

mark [19]. Results of the run are available in Figure 2<sup>2</sup>. The benchmark was run ten times, and the highest memory bandwidth for the core count is reported. Due to the memory bound nature of stencil codes, memory bandwidth results are important to determine how our results compare to the expected peak deriving from the STREAM benchmark.

We run each variant of the 1D stencil and 2D stencil for three and five times respectively. In case of 1D stencil, we report the least time consumed amongst all runs. For 2D stencil, we report the maximum performance achieved for a particular data type. While hyperthreading is enabled on all the cores, we pin to the physical PUs to ensure that the benchmark effectively uses L1 and L2 caches. In a hyperthreaded scenario, the pressure on the cache increases that may result in cache evictions leading to a possible loss in performance. All the benchmarks are run by pinning one thread per core using `hwloc-bind`. Finally, we turn all optimization flags on, i.e. `-O3`, `-ftree-vectorize`, `-ffast-math` for the best possible performances.

**Hardware Counters:** We use Linux `perf` and PAPI to get access to the hardware counters to better explain any aberration in results. All hardware counters were run on a single physical core on a smaller grid size of  $8192 \times 16384$  for a hundred iteration.

## VII. RESULTS

### A. 1D Stencil

Figure 3 succinctly presents the results of strong and weak scaling. The distributed application is implemented such that network latencies can be hidden under compute. Furthermore, the application is also NUMA aware. This is made possible by utilizing block allocators implemented within HPX. The allocator allocates memory based on Linux’s first touch data placement policy. This is similar to OpenMP’s `schedule(static)` policy. Combined with the block executor, we make sure that an HPX thread always spawns at a location of data. Thus, we are able to make up for the lack of bandwidth between chip-to-chip communications. For Fujitsu A64FX, we compiled all dependencies using the Fujitsu compiler. Furthermore, we build the HPX parcelport backend using Fujitsu MPI for easier integration with the system. For all other processors, results are considered according to the benchmark dependencies, as described in Table II.

For Intel Xeon E5 and Fujitsu A64FX under strong scaling, the application takes 28s and 18s respectively for a single node and 3.8s and 2.5s respectively involving eight nodes, which is close to linear scaling (the factor being 7.36 and 7.2, respectively). We expected a similar behavior as sufficient parallelism can be derived from the given number of stencil

<sup>2</sup>Higher STREAM benchmark results can be obtained on A64FX using Fujitsu compiler and special cache initialization techniques. We use the STREAM benchmark as a reference performance metric for our stencil application. For us to make assertions about the optimality of stencil codes, we need to make sure that we do not optimize STREAM benchmarks using techniques that cannot be applied to the 2D stencil code. We make sure that STREAM benchmarks are NUMA aware as we make our 2D stencil code NUMA aware.

### 1D Stencil: Distributed Results

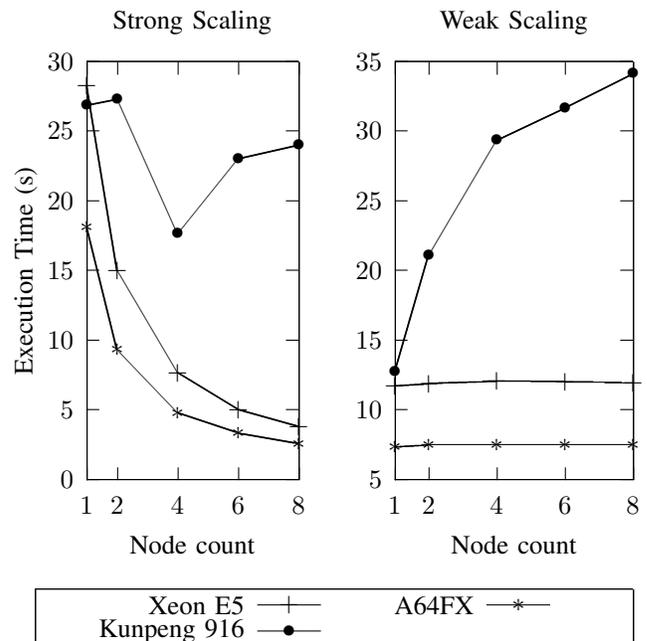


Fig. 3: 1D stencil: strong and weak scaling results. Strong scaling is done over 1.2 billion stencil points. Weak scaling is done by adding 480 million stencil points per node.

points to hide network latencies. Under weak scaling, the application takes 12s and 7.5s respectively irrespective of the number of nodes which proves that the network latencies are aptly hidden.

For Kunpeng 916, we do not observe linear scaling. On closer inspection, we observed that the network performance on the Hi1616 nodes is unsatisfactory and that the processor is not able to exploit the capabilities of the InfiniBand network making it difficult to hide network latencies and the results are transferred to the graph. This hypothesis is well proven under weak scaling, where we see a significant increase in execution times as we increase the number of nodes.

### B. 2D Stencil

The expected peak performance can be calculated with the recorded Memory Bandwidth, and calculated Arithmetic Intensity (See V-B) by putting the values in Equation 1. We do not aim to be optimal as further micro-optimizations can be made on our 2D stencil implementation. Our results will be considered optimal if we reach close to this expected peak performance. We make use of hardware performance counters to explain any aberrations in the results.

Our noticeable observation about GCC is its ability to auto vectorize our 2D stencil application. We observe that the instruction count with auto vectorized codes is similar to our explicitly vectorized codes, sometimes even beating it at instruction count. Visible differences come from cache-misses

2D Stencil: Intel Xeon E5-2660 v3

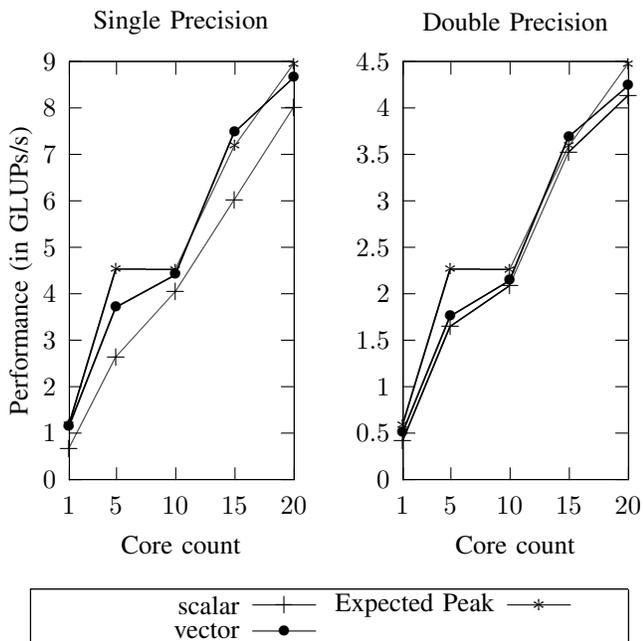


Fig. 4: 2D stencil: Results for Intel Xeon E5-2660 v3 with a grid size of  $8192 \times 131072$  iterated over 100 time steps

and CPU stalls arising from different data layout pattern design by GCC and our code. The next few paragraphs discuss processor specific results in detail.

Table III describes the major contributing hardware counters that directly contribute to performance differences for Intel Xeon E5. We observed a 2x difference in instruction count between scalar and vector types, i.e., GCC is not able to auto vectorize the code very well. The lower instruction count certainly helps an explicitly vectorized code to perform better. Interestingly, the cache friendly optimization from GCC is highly optimized for x86 architecture. This can be seen by the lower cache miss counts for auto vectorized codes.

TABLE III: Hardware Counters for Intel Xeon E5-2660v3

Data Type	Instruction	Cache Misses
Float	$3.153 \times 10^{10}$	$2.121 \times 10^8$
Vector Float	$1.783 \times 10^{10}$	$3.706 \times 10^8$
Double	$6.01 \times 10^{10}$	$4.74 \times 10^8$
Vector Double	$3.507 \times 10^{10}$	$8.751 \times 10^8$

Furthermore, PUs are able to keep only certain memory transactions in the flight. One can make use of CPU stall cycles to determine these numbers. Unfortunately, Intel Xeon E5 2660v3 doesn't support these counters. Nonetheless, we believe that the lower instruction counts relieve the memory controllers of excessive memory transactions improving the performance further. We hypothesize that the improvements

2D Stencil: Huawei Hi1616

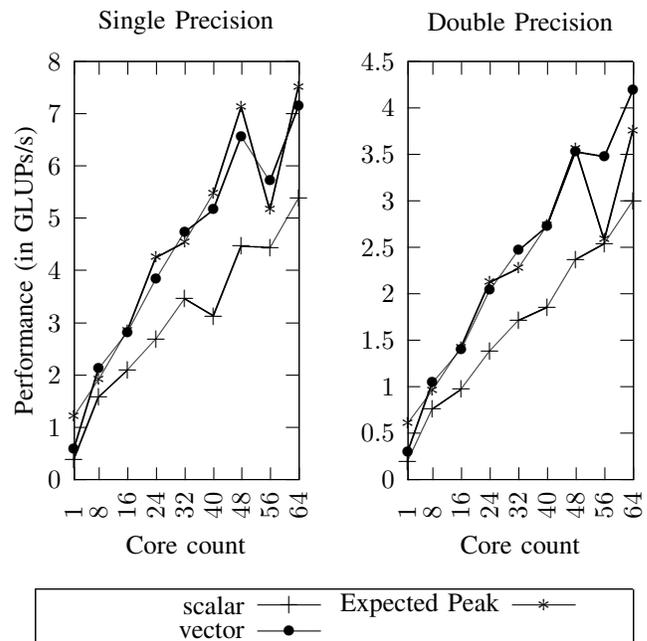


Fig. 5: 2D stencil: Results for Huawei Kunpeng 916 with a grid size of  $8192 \times 131072$  iterated over 100 time steps

of up to 50% with vectorized floats are a result of decreased memory transactions. Given doubles occupy eight bytes, we do not expect much performance improvements with decreased instruction count on an already busy memory bus, which can be easily visualized in the graphs with only up to 10% improvements in performances.

TABLE IV: Hardware Counters for HiSilicon Hi1616

Data Type	Instruction	Cache Misses
Float	$4.3 \times 10^{10}$	$3.148 \times 10^9$
Vector Float	$4.144 \times 10^{10}$	$2.512 \times 10^9$
Double	$8.321 \times 10^{10}$	$5.639 \times 10^9$
Vector Double	$8.236 \times 10^{10}$	$4.953 \times 10^9$

HiSilicon Hi1616 shows up to 80% improvements with explicit vectorization. Table IV describes differentiating hardware counter results. The table suggests that GCC is able to auto vectorize the results very well. Explicit vectorization resulted in a mere 5% improvement in instruction count. The auto vectorization does fail in exploiting caches that can be observed by a 10-20% decline in cache misses by moving to an explicitly vectorized code. We believe that the remaining difference arises from Backend stalls. Unfortunately, Hi1616 doesn't support CPU stall counters, but we hypothesize this based on our results on ThunderX2 (result described in the coming paragraphs) where we see similar behaviors. Backend stalls are generally caused either by many long-latency opera-

2D Stencil: Fujitsu A64FX (Compute cores only)

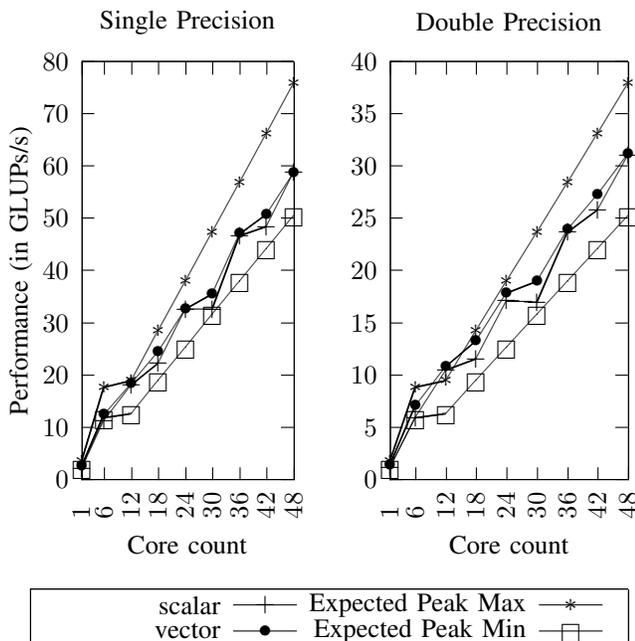


Fig. 6: 2D stencil: Results for Fujitsu A64FX with a grid size of  $8192 \times 131072$  iterated over 100 time steps. Expected Peak Max assumes two memory transfers per iteration and Expected Peak Min assumes three memory transfers per iteration.

tions such as multiply and divide, or by long-latency memory operations, or both. In our case, our stencil kernel uses simpler operations, so, a majority of such stalls are caused by memory operations. With explicit vectorization, we believe that we’ve reduced these memory transactions considerably allowing us for better performances.

Another interesting observation is the sudden decrease in performance while moving from 32 to 40 cores and 56 to 64 cores. It can be explained by understanding how memory controllers work. When the application runs on 40 cores, two of the NUMA domains are fully saturated with respect to memory bandwidth, while the third NUMA domain is only partially saturated. This uneven distribution leads to faster iteration in the fully saturated domains, leaving a trail of poor performance for the partially saturated domain. Therefore, a partially saturated NUMA domain becomes the critical path for the benchmark. With lower memory bandwidth available to the partially saturated domain, we see a loss in performance.

For Fujitsu A64FX, the execution time for the 2D stencil kernel is less than 2s for scalar and vector floats and about 3.5s for scalar and vector doubles while utilizing all 48 compute cores. Compared to the other processors, the execution time is significantly lower. This piqued our interest, and we investigated whether we have enough parallelism for HPX to

2D Stencil: Fujitsu A64FX (Compute cores only)

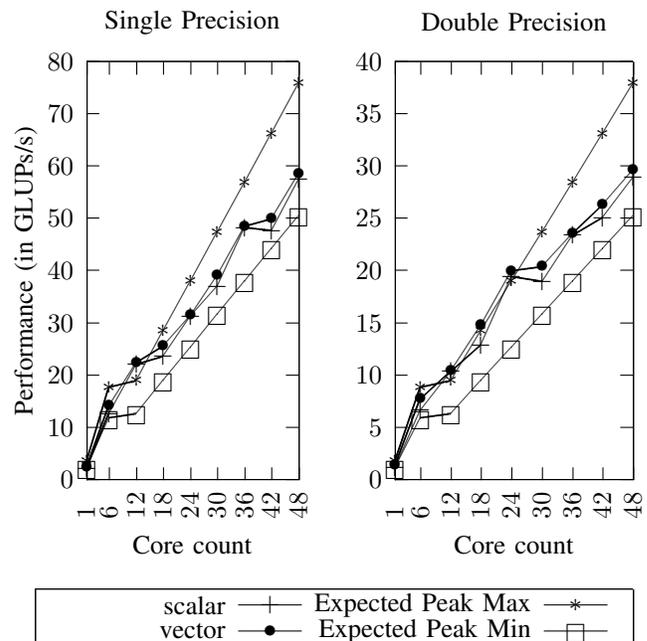


Fig. 7: 2D stencil: Results for Fujitsu A64FX with a grid size of  $8192 \times 196608$  iterated over 100 time steps. Expected Peak Max assumes two memory transfers per iteration and Expected Peak Min assumes three memory transfers per iteration.

take advantage of. Like every AMT model, HPX is known to have contention overheads when the grain size is too small.

We decided to investigate by increasing the grid size. We discovered that 32GB High Bandwidth Memory (HBM) can be a disadvantage for memory-intensive applications. For instance, our grid requires 9GB worth of DRAM. A 2D stencil code has two grids, i.e., 18GB worth of DRAM. We can, therefore, only test grid sizes of up to 1.5x the current size. Our investigation results suggest that there are no performance benefits (see Figure 7) in increasing grid size. This means that HPX is able to take advantage of the underlying parallelism.

TABLE V: Hardware Counters for Fujitsu FX1000 A64FX

Data Type	Instruction	Frontend Stalls	Backend Stalls
Float	$1.284 \times 10^{10}$	$3.801 \times 10^8$	$9.43 \times 10^9$
Vector Float	$1.496 \times 10^{10}$	$2.918 \times 10^8$	$8.003 \times 10^9$
Double	$2.299 \times 10^{10}$	$3.86 \times 10^8$	$1.871 \times 10^{10}$
Vector Double	$2.956 \times 10^{10}$	$3.56 \times 10^8$	$1.443 \times 10^{10}$

Another interesting observation is the better performance compared to the expected peak assuming three memory transfers per iteration. We claimed that the caches were large enough to accommodate three rows resulting in three memory transfer per iteration. With A64FX, we observe that the results

## 2D Stencil: Marvell ThunderX2

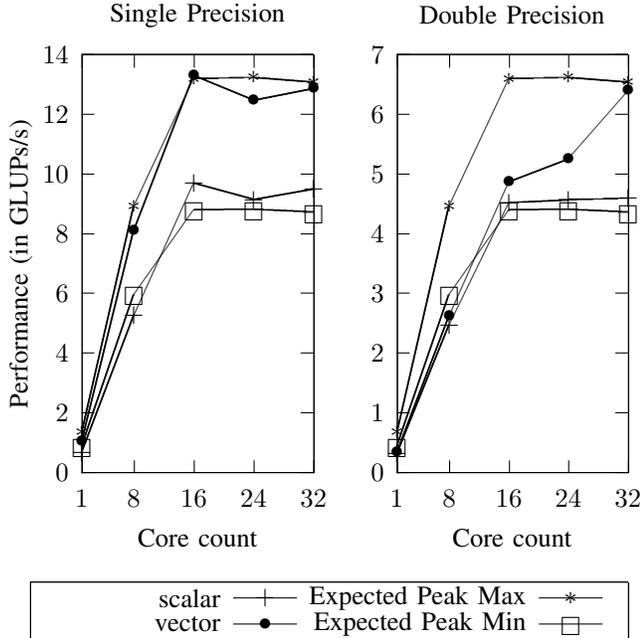


Fig. 8: 2D stencil: Results for Marvell ThunderX2 with a grid size of  $8192 \times 131072$  iterated over 100 time steps. Expected Peak Max assumes two memory transfers per iteration and Expected Peak Min assumes three memory transfers per iteration.

do not follow the scheme. Due to large sized cache lines, we see cache benefits resulting in better Arithmetic Intensity. We witness results equivalent to cache blocking version of 2D stencil up to 32 cores. A cache blocked version of 2D stencil essentially reduces the number of memory transfers per iteration, in our case, by one. This results in a 49% performance boost over the previously expected results.

The significant increase in performance going from half saturated NUMA domain to a fully saturated NUMA domain can be explained using the same principles as described for HiSilicon Hi1616. We use a 512-bit SVE vector length for benchmarking the 2D stencil. From the results, it is clear that no significant improvements are achieved by explicitly vectorizing the code. The improvements are anywhere from 5% to 15%. Table V provides some insights. Firstly, glancing at the instruction count, we observe that GCC does a better job of optimizing the instruction count than our explicitly vectorized code. The cache miss counts were very similar for auto vectorized and explicitly vectorized codes as well. What differs visibly are the CPU stall counts. This means that explicit vectorization is helping to relieve stress from the memory controllers. Due to significant reductions in CPU stalls for vectorized codes, we get marginally better results.

For Marvell ThunderX2, we see a similar behavior as

observed in A64FX. Single precision performance is taking complete advantage of a large cache line. Cache blocking behavior can be observed, leading to an additional 49% performance boost over an implementation assuming three memory transfers per iteration. With doubles, we see contradictory behavior. At lower core count, it behaves optimally according to our assumed arithmetic intensity. At 16 cores and above, the behavior changes to an arithmetic intensity of 1/8 and 1/16 for floats and doubles, respectively. This leads to an improvement in performance. Using hardware stall counters, we found that the number of backend and frontend stalls for explicitly vectorized code that shows this interesting switch reduced by about 40% from  $2.353 \times 10^{10}$  and  $1.144 \times 10^8$  (for auto vectorization) to  $1.577 \times 10^{10}$  and  $7.867 \times 10^7$  (for explicit vectorization), respectively, at a core count of 32, leading to a significant increase in performance. However, the total number of instructions and cache misses were similar for both auto vectorized and explicitly vectorized codes. Unfortunately, we were not able to identify the reason behind this interesting switch and it remains an open question.

TABLE VI: Hardware Counters for Marvell ThunderX2

Data Type	Instruction	L2 Cache Misses	Backend Stalls
Float	$4.039 \times 10^{10}$	$1.811 \times 10^9$	$1.522 \times 10^{10}$
Vector Float	$4.394 \times 10^{10}$	$1.69 \times 10^9$	$6.437 \times 10^9$
Double	$8.065 \times 10^{10}$	$5.716 \times 10^9$	$3.298 \times 10^{10}$
Vector Double	$8.756 \times 10^{10}$	$6.055 \times 10^9$	$2.826 \times 10^{10}$

Effects of adding explicit vectorization are significant. While GCC is able to auto vectorize the code well (see Table VI), we notice that the implementation is rather unoptimized. We observe a decrease in L2 cache misses of about 10% (No visible difference in L1 cache misses, hence not reported). The number of backend stalls observed in the case of GCC is considerably higher compared to an explicitly vectorized code. This means that outstanding load/store instructions with explicit vectorization is noticeably lower than an auto vectorized code. In this regard, we believe that ThunderX2 microarchitecture behaves similar to Cortex-A72 where we saw that cache misses differed by only 15%, but the resulting performance gap was upwards of 50%. We see significant improvements by utilizing explicit vectorization. These improvements were consistently within 50-60% for floats and up to 40% for doubles. The results also look nearly optimal for the given memory bandwidth.

## VIII. CONCLUSION

This paper gives a first overview of the performance of AMTs running at scale on a production HPC system that is based on Arm processors. Our experience of porting HPX and the benchmarks to Arm processors was mostly straightforward. We faced a few issues when building for SVE types. Arm Compilers implement SVE using `__sizeless_struct` making it impossible to wrap an SVE type to the custom container as they have no size at compile time. Currently, only GCC allows the commandline flag to pass SVE vector length with

the `-msve-vector-bits`. However, this comes at the cost of SVE type portability. Further development is required to integrate custom containers to work with `__sizeless_struct` to make it easier for application developers to port their application to Arm.

We demonstrate that the application scales both on-node and distributed. We found that performance on Arm processors is as good or better than their x86 brethren. For the 1D stencil, all processors except Kunpeng 916 showed good scaling results. In the case of Kunpeng 916, the poor interconnect network is to be blamed. For the 2D stencil, we observed that processors with large cache lines showed inherent cache blocking benefits (without explicit implementation). This resulted in about a 50% performance boost over the expected results. We also observed that explicit vectorization can improve the performance significantly for ThunderX2 and Kunpeng 916 due to considerably lower CPU stall counts when compared with auto vectorized codes. For A64FX, we did not observe any visible performance benefits by employing explicit vectorization.

## IX. ACKNOWLEDGMENT

We would like to thank Huawei for making the JUAWEI cluster hardware available at JSC as well as Fujitsu for allowing us to participate in the A64FX early access program. Funding for parts of this work is received from the European Commission H2020 program under Grant Agreement 779877(Mont-Blanc 2020).

We are grateful for the insights provided by Dr. Thomas Heller (Exasol) and Prof. Hartmut Kaiser (LSU, USA) in analyzing and optimizing our benchmarks.

## REFERENCES

- [1] Fujitsu. (2019) Supercomputer fugaku. [Online]. Available: <https://www.fujitsu.com/global/Images/supercomputer-fugaku.pdf>
- [2] N. Stephens, S. Biles, M. Boettcher, J. Eapen, M. Eyole, G. Gabrielli, M. Horsnell, G. Magklis, A. Martinez, N. Premillieu *et al.*, “The arm scalable vector extension,” *IEEE Micro*, vol. 37, no. 2, pp. 26–39, 2017.
- [3] S. McIntosh-Smith, J. Price, T. Deakin, and A. Poenaru, “Comparative benchmarking of the first generation of hpc-optimised arm processors on isambard,” 2018.
- [4] S. McIntosh-Smith, J. Price, A. Poenaru, and T. Deakin, “Scaling results from the first generation of arm-based supercomputers.”
- [5] A. Jackson, A. Turner, M. Weiland, N. Johnson, O. Perks, and M. Parsons, “Evaluating the arm ecosystem for high performance computing,” in *Proceedings of the Platform for Advanced Scientific Computing Conference*, ser. PASC ’19. New York, NY, USA: Association for Computing Machinery, 2019. [Online]. Available: <https://doi.org/10.1145/3324989.3325722>
- [6] F. Banchelli, M. Garcia, M. Josep, F. Mantovani, J. Morillo, K. Peiro, G. Ramirez, G. Valenzano, and J. W. Weloli, “Mb3 d6. 9–performance analysis of applications and mini-applications and benchmarking on the project test platforms. version 1.0.”
- [7] H. Kaiser, M. Brodowicz, and T. Sterling, “ParalleX an advanced parallel execution model for scaling-impaired applications,” in *2009 International Conference on Parallel Processing Workshops*, 2009, pp. 394–401.
- [8] M. Anderson, M. Brodowicz, H. Kaiser, and T. Sterling, “An application driven analysis of the parallex execution model,” 2011.
- [9] T. Heller, H. Kaiser, and K. Iglberger, “Application of the parallex execution model to stencil-based problems,” *Computer Science - Research and Development*, vol. 28, no. 2, pp. 253–261, May 2013. [Online]. Available: <https://doi.org/10.1007/s00450-012-0217-1>
- [10] H. Kaiser, T. Heller, B. Adelstein-Lelbach, A. Serio, and D. Fey, “Hpx: A task based programming model in a global address space,” in *Proceedings of the 8th International Conference on Partitioned Global Address Space Programming Models*, ser. PGAS ’14. New York, NY, USA: Association for Computing Machinery, 2014. [Online]. Available: <https://doi.org/10.1145/2676870.2676883>
- [11] T. Heller, P. Diehl, Z. Byerly, J. Biddiscombe, and H. Kaiser, “HPX – An open source C++ Standard Library for Parallelism and Concurrency,” in *Proceedings of OpenSuCo 2017, Denver, Colorado USA, November 2017 (OpenSuCo’17)*, 2017, p. 5.
- [12] H. Kaiser, B. A. L. aka wash, T. Heller, M. Simberg, A. Bergé, J. Biddiscombe, aurianer, A. Bikineev, G. Mercer, A. Schäfer, K. Huck, A. S. Lemoine, T. Kwon, J. Habraken, M. Anderson, M. Copik, S. R. Brandt, M. Stumpf, D. Bourgeois, D. Blank, S. Jakobovits, V. Amaty, rstobaugh, L. Viklund, Z. Khatami, P. Diehl, T. Pathak, D. Bacharwar, S. Yang, and E. Schnetter, “STELLAR-GROUP/hpx: HPX V1.4.1: The C++ Standards Library for Parallelism and Concurrency,” Feb. 2020. [Online]. Available: <https://doi.org/10.5281/zenodo.3675272>
- [13] S. Williams, A. Waterman, and D. Patterson, “Roofline: An insightful visual performance model for multicore architectures,” *Commun. ACM*, vol. 52, no. 4, p. 65–76, Apr. 2009. [Online]. Available: <https://doi.org/10.1145/1498765.1498785>
- [14] S. Williams, D. Patterson, L. Oliker, J. Shalf, and K. Yelick, “The roofline model: A pedagogical tool for auto-tuning kernels on multicore architectures,” in *Hot Chips*, vol. 20, 2008, pp. 24–26.
- [15] S. W. Williams, “Auto-tuning performance on multicore computers,” Ph.D. dissertation, USA, 2008.
- [16] R. Barrett, M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, and H. van der Vorst, *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*. Society for Industrial and Applied Mathematics, 1994. [Online]. Available: <https://epubs.siam.org/doi/abs/10.1137/1.9781611971538>
- [17] P. Boyle, A. Yamaguchi, G. Cossu, and A. Portelli, “Grid: A next generation data parallel c++ qcd library,” *arXiv preprint arXiv:1512.03487*, 2015.
- [18] B. Bramas, “Inastemp: A novel intrinsics-as-template library for portable simd-vectorization,” *Scientific Programming*, vol. 2017, p. 5482468, Sep 2017. [Online]. Available: <https://doi.org/10.1155/2017/5482468>
- [19] J. D. McCalpin, “Memory bandwidth and machine balance in current high performance computers,” *IEEE Computer Society Technical Committee on Computer Architecture (TCCA) Newsletter*, pp. 19–25, Dec. 1995.