

# Simulation-Based Performance Prediction of HPC Applications: A Case Study of HPL

Gen Xu\*, Huda Ibeid<sup>§</sup>, Xin Jiang\*, Vjekoslav Svilan<sup>§</sup>, and Zhaojuan Bian\*

\* Intel Corporation, Shanghai, China

<sup>§</sup> Intel Corporation, Santa Clara, US

{gen.xu, huda.ibeid, xin.jiang, vjekoslav.svilan, bianny.bian}@intel.com

**Abstract**—We propose a simulation-based approach for performance modeling of parallel applications on high-performance computing platforms. Our approach enables full-system performance modeling: (1) the hardware platform is represented by an abstract yet high-fidelity model; (2) the computation and communication components are simulated at a functional level, where the simulator allows the use of the components native interface; this results in a (3) fast and accurate simulation of full HPC applications with minimal modifications to the application source code. This hardware/software hybrid modeling methodology allows for low overhead, fast, and accurate exascale simulation and can be easily carried out on a standard client platform (desktop or laptop). We demonstrate the capability and scalability of our approach with High Performance LINPACK (HPL), the benchmark used to rank supercomputers in the TOP500 list. Our results show that our modeling approach can accurately and efficiently predict the performance of HPL at the scale of the TOP500 list supercomputers. For instance, the simulation of HPL on Frontera takes less than five hours with an error rate of four percent.

**Index Terms**—performance modeling, exascale systems, HPL

## I. INTRODUCTION

Currently, there are many efforts to evaluate the hardware and software bottlenecks of exascale designs to enable the development of applications that exploit the full performance of exascale computing platforms. However, the increasing complexity of modern computing architectures along with the exponentially growing configuration space and complex interactions among configuration options often make it difficult to develop accurate performance models. In recent years there have been several efforts to model the performance of HPC applications using simulation-based approaches. However, several challenges must be addressed to enable these approaches.

The full system stack consists of three layers: hardware infrastructure, middle layer libraries, and the application itself. Each layer can have a huge impact on the overall performance, which means that all layers should be modeled to achieve an acceptable accuracy. One of the main challenges is to determine which aspects are the most important to simulate when modeling each layer for large scale HPC applications. In terms of the hardware infrastructure layer, computation components, such as CPU, GPU, and memory, should be modeled. Similarly, the interconnect network is one of the essential parts. The computation and communication platforms are the most important to take into consideration for the distributed system.

Choosing which libraries to simulate is another important aspect. The basic principle is to choose the most widely used libraries. Science and engineering computations have been the dominant category of the applications running on HPC systems. In this area, Basic Linear Algebra Subprograms [1] (BLAS) is the most widely used mathematical library that forms the computational core of many HPC applications. BLAS operations very time-consuming as well as compute-intensive. Additionally, Message Passing Interface (MPI) has now emerged as the de-facto standard for node-to-node communication on supercomputers. MPI standards are used on all leading supercomputers of the TOP500 list [2]. Taking the characteristics of the software libraries is an essential requirement for accurate simulation-based modeling.

With the hardware infrastructure and software libraries models, our goal is to enable the modeling of HPC applications with minimal modification to the application source code. Among all HPC applications, the High-Performance LINPACK (HPL) Benchmark is the most widely recognized metric for ranking HPC systems, although other benchmarks such as HPGMG [3] and HPCG [4] have been proposed as either alternative or complementary benchmarks.

In this paper, we propose a simulation framework that employs a layered architecture to simulate HPC systems on standard client computers (desktop or laptop). We use HPL to demonstrate the capability and scalability of the simulation framework. The key contributions of this paper are as follows:

- We present a hardware platform model that includes the processing nodes and the interconnection network. The model employs a stream-level network model that balances the simulation speed and accuracy.
- We present abstracted library models for BLAS computations and MPI communications.
- We model HPL benchmark to demonstrate the capability and scalability of our simulation framework.
- We demonstrate that our modeling approach can accurately and efficiently predict the performance of HPL at the scale of the TOP500 list supercomputers.

The rest of the paper is organized as follows. In section II, we present a background on simulation-based approaches. We also describe related work in hardware infrastructure simulation and MPI modeling. In section III, we present an overview of our simulation framework and describe the design

of each of its layers. In section IV, we conduct extensive validation and performance studies. In section V we present some use cases. Finally, conclusions and future directions are presented in section VI.

## II. BACKGROUND AND RELATED WORK

In recent years there have been several efforts to model the performance of HPC applications using simulation-based approaches.

SimGrid [5] is an open-source simulation framework for large-scale distributed systems. It was originally designed to study the behavior of Grids but has been extended and applied to a wide range of distributed computing platforms, including Clouds and High Performance Computing systems. SimGrid uses a flow-level approach that approximates the behavior of TCP networks. Due to its use of a flow-level network simulation approach along with a coarse-grained CPU model for the computation, SimGrid can perform large numbers of statistically significant experiments on large TCP networks. However, SimGrid might result in an unacceptable accuracy when compared to packet-level simulators when the data sizes are small or when networks are highly contended [6]. In addition, the lack of detailed models for the processing components makes SimGrid unsuitable for several HPC applications.

The Structural Simulation Toolkit (SST) [7] enables the co-design of highly concurrent systems by allowing simulation of diverse aspects of the hardware and software. SST aims to simulate full-scale machines using a coarse-grained simulation approach for the processing and network components through the use of skeleton applications that replicate the full application control flow.

The work presented in this paper builds on our previous work, CSMMethod [8]. CSMMethod enables full-system performance modeling and prediction of big data clusters by simulating both the software stack (e.g. HDFS, OS, and JVM) and the hardware components (CPU, storage, and network). With CSMMethod, the computation and communication behaviors of the application are abstracted and simulated at a functional level. Software functions are then dynamically mapped onto hardware components. To achieve fast and accurate performance simulation, CSMMethod supports fine-grained analytical models for processor, memory, and storage. The timing of the hardware components is modeled according to payload and activities as perceived by the software. CSMMethod capabilities and accuracy have been demonstrated in [9]–[12]. However, CSMMethod is focused on big data applications and has not been applied to simulate HPC systems.

Cycle-accurate simulators are commonly used to evaluate next generation processors and system architectures. Traditionally, these simulators trade speed for accuracy. Similarly, packet-level or flit-level network simulators aim for a highly accurate representation of actual network behavior. Thus, large-scale simulations may be too time-consuming with packet-level simulation.

There are several different approaches to model MPI, ranging from analytical models to trace-based simulations. Some

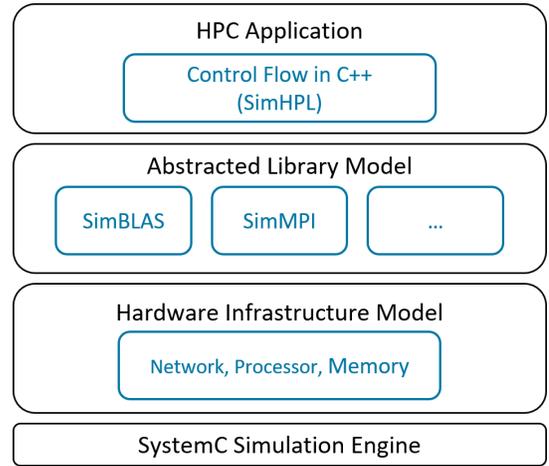


Fig. 1: Simulation framework architecture.

MPI modeling frameworks rely on the use of test environments based on “artificial communications” to perform synthetic tests of MPI performance. For example, LogGOPSim [13] replaces MPI collective operations by a set of point-to-point algorithms. While this approach is accurate on smaller systems, LogGOPSim ignores congestion in the network and assumes full effective bisection bandwidth, which may decrease the accuracy of the simulations on emerging large-scale systems. SMPI [14] simulates unmodified MPI applications on top of the SimGrid simulator. SMPI supports different performance modes through a generalization of the LogGPS model.

## III. SIMULATION FRAMEWORK

Our simulation framework employs a layered and configurable architecture to simulate the full stack of supercomputing systems, as shown in Figure 1. The top layer is the HPC application, where the application behavior is modeled. Underneath the top layer, computation and communication libraries are abstracted and simulated at a functional level. The library layer receives function calls from the top layer and dynamically connects to the hardware components. The hardware infrastructure layer beneath the library layer aims at defining the hardware components (processor, network, and storage) of the HPC system. In this framework, software behavior and hardware infrastructure are loosely coupled, which provides the flexibility to change the hardware platform without the need to modify the software behavior model and vice versa.

This paper discusses several techniques: (1) the hardware platform is modeled by an abstract yet high-fidelity model; (2) computation and communication components are simulated at a functional level, where the simulator allows the use of the component native interface; this results in a (3) fast and accurate simulation of HPC applications with minimal modifications to the application source code; and, at the bottom of these layers, (4) a simulation engine for SystemC-based discrete events. This is a low-overhead engine that enables fast simulations with good scalability. This hardware/software

hybrid modeling methodology allows for low overhead, fast, and accurate Exascale systems simulation and can be easily carried out on a standard client platform.

#### A. HPC hardware infrastructure simulation

The hardware model builds on our previous work, CSMethod [8]. Here, we extend CSMethod to enable the modeling of HPC applications. In particular, we implement an efficient CPU model for the computation operations as well as a GPU model. Moreover, a stream-level network model is implemented as an alternative to the original packet-level network model.

The hardware model simulates all the main components of the HPC platform, which includes the processing nodes and the interconnection network. In particular, the hardware infrastructure layer consists of models for the CPU, GPU, memory, and NIC. This section describes these models.

1) **Node architecture: CPU, GPU, and memory:** In this work, we extend [8] to support heterogeneous architectures. This new feature enables the simulation of accelerator-based architectures, such as CPU–GPGPU combinations. Our framework also utilizes analytical models to model compute-bound and bandwidth-bound operations, such as BLAS DGEMM operation and DSWAP described in section III-B1. Traditionally, compute-bound operations are modeled using an actual single-core execution time on real hardware scaled to the simulated processor core speed. In this work, we model the computation time of these operations analytically based on the theoretical peak performance and the efficiency of these operations on the CPU and GPU. The efficiency can be directly measured without complex computations. Similarly, modeling bandwidth-bound operations is based on the peak bandwidth and bandwidth efficiency.

2) **Interconnection network:** As discussed earlier, packet-level network models are not suitable for all scenarios. In this work, a stream-level network model is implemented as an alternative that offers latency and bandwidth restrictions. This work extends the capabilities of [8] network model in two ways. First, we include more network architectures, such as fat-tree and dragonfly, which are the most widely used networks in HPC systems. Second, traditionally, the implementation of routing policies calculates and stores all the routing paths during the initialization phase which uses a large amount of memory when simulating large-scale systems. Several routing algorithms, such D-mod-K for fat-tree [15] and minimal/non-minimal routing for dragonfly topology [16] can be dynamically calculated which reduces the memory consumption significantly.

To model the network communication, we divide large messages into smaller chunks and calculate the transmission time according to the currently allocated bandwidth. In addition, the network model supports communication primitives, such as send data and receive data, which enables the integration of external network simulators into our framework.

#### B. Computation and communication libraries simulation

When developing simulation models for large scale complex systems, it is important to consider which components to model. In HPC applications, computation and communication libraries are commonly utilized and tuned for optimal performance. In this work, BLAS and MPI libraries are simulated as modules on top of the infrastructure layer by leveraging dedicated APIs to access the hardware resources. These modules allow the use of the libraries native interface, thus easing the development of the simulation APIs.

In this section, a detailed discussion of the computation and communication libraries is presented.

1) **Performance modeling of BLAS library:** Many HPC applications rely heavily on BLAS kernels. The BLAS library implements fundamental dense vector and matrix operations, such as various types of multiplications and triangular linear system solvers. Since these kinds of kernels do not influence the control flow, the simulation time can be reduced by substituting the BLAS function calls with an analytical performance model for the respective kernel. The BLAS operation is data-independent, i.e., the data content does not affect the computation time. This means that all multiplications with zeros are explicitly performed no matter how sparse an operand is (i.e., how few non-zero entries it has).

BLAS functionality is categorized into three sets of levels according to the arithmetic density. Level 1 BLAS operations typically take linear time,  $\mathcal{O}(N)$ , Level 2 operations quadratic time, and Level 3 operations cubic time. Thus, we employ the same modeling approach but with different analytical performance models that are based on the Roofline model [17]. The Roofline model provides a simple way to estimate the performance based on the computation kernel and hardware characteristics. It relies on the concept of Arithmetic Intensity (in FLOPs/byte) and provides performance bounds for compute-bound and memory bandwidth-bound computations.

**Modeling Level-3 BLAS Kernels:** Here we describe in detail the methodology used to model the DGEMM operation. A similar approach is used to model the DTRSM kernel.

GEMM performs a matrix-matrix multiplication and an add operation

$$C \leftarrow \alpha AB + \beta C, \quad (1)$$

where  $C$  is  $m \times n$ ,  $A$  is  $m \times k$ ,  $B$  is  $k \times n$ , and  $\alpha$  and  $\beta$  are scalars.

For dense matrices, the total number of operations performed by GEMM is

$$O_{GEMM} = 2mnk + 2mn. \quad (2)$$

As the GEMM kernel is compute-bound, we use the following analytical model to estimate its compute time

$$E = \mu O_{GEMM} + \theta, \quad (3)$$

where  $\mu$  represents the computation cost of a single operation and  $\theta$  represents the overhead of each DGEMM function call. The Roofline model sets an upper bound on performance of a kernel depending on its arithmetic intensity. For a more

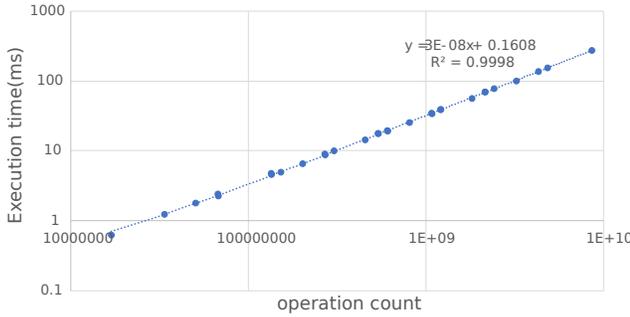


Fig. 2: Execution time of DGEMM kernel.

realistic estimates, we take into account the kernel efficiency on a given hardware. Let  $e$  be the GEMM efficiency on a given hardware, then  $\mu$  can be calculated as the inverse of the multiplication of  $e$  by the theoretical peak performance.

Both  $\mu$  and  $\theta$  in (3) are obtained through profiling and calibration. To calibrate and validate our model, we conduct a micro-test using MKL DGEMM kernel on a single core. The values of  $m$ ,  $n$ , and  $k$  range from 128 to 2048. Each case is executed 1000 times and then the average time is calculated. Figure 2 shows the impact of the total number of operations on the execution time along with the estimation model. The validation results show that the R-squared value is 0.9998. Here, the values of  $\mu$  and  $\theta$  are implementation and hardware dependent. This kind of analytical modeling speeds up the simulation by orders of magnitude, especially as the matrix size grows.

**Modeling Level-1 and Level-2 BLAS Kernels:** A similar approach is employed to model Level-1 and Level-2 BLAS kernels. On most architectures, Level-1 BLAS vector-vector operations, and Level-2 BLAS matrix-vector operations are memory-bound. As mentioned previously, we calibrate the models to take into account the memory efficiency of these operations.

Based on the methodology discussed, we present SimBLAS, a library to simulate and predict the performance of BLAS operations. Figure 3 shows a code snippet of Level-3 and Level-1 SimBLAS operations. There are different implementations of the BLAS library, for example, cuBLAS for GPUs, OpenBLAS, and Intel BLAS. Each implementation has different efficiency. Furthermore, these implementations can run on a single thread or with multi-threading. Hence, predicting efficiency analytically is a complicated task. In our simulations, we employ a microbenchmark to profile the efficiency and then use it as an input to SimBLAS.

SimBLAS library is coupled with the underlying hardware models, specifically, CPU, GPU, and memory models. As discussed earlier in this section, the execution time is determined by the operation complexity and hardware characteristics. The operation complexity is the operations count of a compute-bound operation or the memory access size of a bandwidth-bound operation. The hardware characteristics are obtained from the underlying hardware models.

```

1 void simblas_dgemm(const SIMBLAS_ORDER Layout,
2                  const SIMBLAS_TRANSPOSE TransA,
3                  const SIMBLAS_TRANSPOSE TransB,
4                  const int M, const int N,
5                  const int K, const double alpha,
6                  const double *A, const int lda,
7                  const double *B, const int ldb,
8                  const double beta,
9                  double *C, const int ldc) {
10
11     // Number of FLOPS
12     double op_count = M * N * (2 * K + 2);
13     // Achieved performance
14     double perf = getDgemmPerf();
15     waitns(op_count/getDgemmPerf() + getBlasLat());
16     return;
17 }
18
19 void simblas_dswap(const int N,
20                  double *X, const int incX,
21                  double *Y, const int incY) {
22
23     // Message size
24     double data_movement = 4.0 * N;
25     // Achieved performance
26     double perf = getDswapPerf();
27     waitns(data_movement/perf + getBlasLat());
28     return;
29 }

```

Fig. 3: Implementation of SimBLAS operations.

In summary, these performance models, in principle, balance simulation speed and accuracy to predict the performance of HPC systems.

2) **Performance modeling of MPI library:** In our previous work, a set of socket-like APIs are implemented to support TCP network transmission in big data environments. On HPC platforms, MPI is the de-facto standard for inter-node communication. This section details the MPI model in two aspects: peer-to-peer communication and collective communication.

First, all the peer-to-peer communication APIs, both synchronous and asynchronous operations, are implemented in the network model. The execution time of the MPI communication operations is independent of the message content. Hence, we model the performance based on the message size and the underlying network. Different communication protocols are used for different message sizes, such as “eager” or “rendezvous”. Many state-of-the-art MPI simulators, such as SMPI [14], have depicted this design methodology and proven good simulation accuracy for a wide range of settings without any application-specific tuning. Our approach is similar, a linear model is used to predict the MPI communication performance. This model is built on top of the hardware model discussed in section III-A. The network contention is simulated using the underlying network model. Figure 4 illustrates the implementation of MPI send operation. At first, the global server IDs of the sender and receiver are obtained. Then, a network function *SendData* is called to pass the communication request to the network model.

Second, we model collective communications. Previous studies show that the performance and scalability of MPI collective communication operations are critical to HPC applica-

```

1 extern "C" int MPI_Send(const void *buf, int count,
2                       MPI_Datatype datatype,
3                       int dest, int tag,
4                       MPI_Comm comm) {
5     int src_id, dest_id;
6     // Global ID of source and dest processes
7     MPI_Comm_globalID(dest, comm, dest_id, src_id);
8     double data_size = count * datatype;
9     // This function returns after the data is sent
10    SendData(src_id, dest_id, data_size);
11    return MPI_SUCCESS;
12 }

```

Fig. 4: Implementation of MPI Send.

tions. In major MPI implementations, each collective operation has several different algorithms to choose from depending on several factors, such as the message size and network topology. In some algorithms, collective communication is broken into a set of peer-to-peer operations. In our model, several algorithms for each operation are simulated mimicking the behavior of real implementations of OpenMPI and IntelMPI. In addition, optimized algorithms for specific network topologies, such as torus and dragonfly networks, are also available.

### C. Modeling applications behavior

In a previous section, we discussed several approaches to model application behavior. One traditional approach is to study and analyze the application source code, mimic its behavior at an abstract level, and model its critical components. While this method offers a high modeling accuracy, it is time-consuming and requires frequent follow-up model refinements.

With the hardware infrastructure and libraries models, our goal is to enable the modeling of HPC applications with few modifications to the application source code instead of mimicking applications behavior. To achieve this goal, several challenges need to be addressed. We use HPL as an example in this section.

**Parallel processes:** Our framework employs Intel CoFluent Studio (CoFluent) [18] which provides an easy to use graphical modeling tool in a SystemC simulation environment. Since SystemC is a sequential simulation engine, every MPI process of the application needs to be mapped onto a SystemC thread. [8] describes how to mimic an application parallel behavior in detail. As the native application source code is used in our approach, each MPI process is bound with a SystemC virtual thread. Using this approach, all the HPL processes are simulated with low overheads.

**Integration of SimBLAS and SimMPI libraries:** The original HPL source code supports several BLAS interfaces, for example, CBLAS and FBLAS. Here, we enable SimBLAS interfaces in HPL source code. Only three modifications to the HPL source code are needed, defining SimBLAS and including the new header file. SimMPI supports the same APIs as the standard MPI library. Hence, enabling SimMPI in HPL source code is simply achieved by including a header file.

**Simulation of other components:** In addition to the BLAS computations and MPI communications, HPL spends signif-

TABLE I: Hardware and Software configurations.

Category	Details
Node#	4
CPU#	Intel Xeon Broadwell E5-2699 v4 @ 2.2GHz
Socket#	2
Cores#/Socket	22
Memory/node	DDR4 256GB @ 2400MHz
Network	1 Port Intel OPA 100Gb
HPL version	Open HPL v2.3, Intel HPL v2.2
MPI version	Intel MPI 2019

icant time performing local copy and swap operations. In order to model HPL accurately, these HPL kernels, such as *HPL\_dlaswp\**, are simulated using the same approach used for BLAS Level-1 operations. Furthermore, min and max functions are simulated with random numbers as the content has no impact on HPL behavior.

**Privatization of global variables:** As the CoFluent kernels are implemented in SystemC, which uses a single process to simulate parallel MPI processes, global variables in the application code are shared between all MPI processes. In our framework, a private copy of the global variables is stored for each parallel process. CoFluent offers a simple API, *get\_container()*, which can be used by a virtual thread to get the corresponding MPI rank. A global array is used to store the privatized variables and can be accessed using a dedicated index.

The last challenge is to identify which components of the source code to modify. In this work, optimizations for simulation speed are used to identify the modifications. The two largest data structures in HPL are matrix *A* and the *panel* which stores the workplace. The total space allocated by the MPI processes on each node typically consumes most of the node memory while the content of *A* is irrelevant for the simulation. This memory allocation is removed with small modifications to code. The simulation results also indicate no impact on the execution flow and simulation accuracy.

Even though the matrix *A* can be removed, *panel* is used in every iteration of the factorization and, hence, must be stored. A possible workaround is to allocate and free *panel* structure at every iteration. However, this option is time-consuming. Alternatively, we use a global array to store *panel* structures for all MPI processes and *panel init/free* functions are reimplemented to map/demap corresponding spaces to private addresses.

## IV. PERFORMANCE VALIDATION AND SCALABILITY EVALUATION

In this section, we first discuss the accuracy of our framework. Then, we examine its scalability by performing simulations while changing the number of MPI processes from 2,000 to 10,000. Lastly, we demonstrate the fast simulation speed with different problem sizes and various configuration settings.

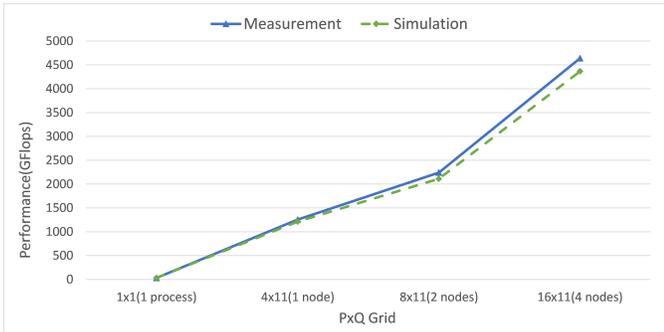


Fig. 5: OpenHPL performance.

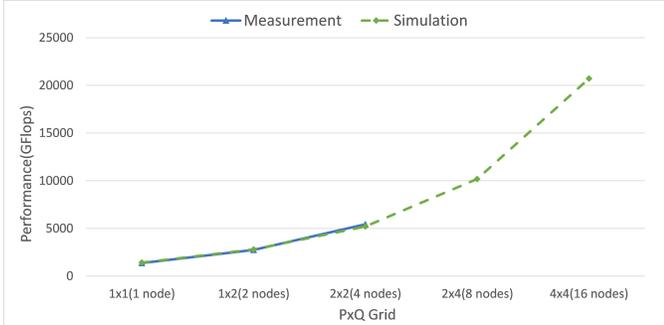


Fig. 6: Intel HPL performance.

#### A. Simulation accuracy

To validate the simulation accuracy, we conduct experiments on our local environment. Table I shows the configurations details of the environment. The cluster has 4 nodes, each node has a dual-socket of Intel Xeon CPU with 22 cores per socket. Each node has 256GB DDR4 memory operating at frequency 2.4GHz. All nodes are connected to the same switch with a single port of Intel 100Gb OPA. Software configurations are also shown in Table I. Two HPL versions, OpenHPL 2.3 and Intel HPL 2.2 are installed. We choose the two versions since they are both widely used in supercomputing systems as demonstrated in the TOP500 list. OpenHPL is compiled with Intel MKL 2019 and Intel MPI 2019. Intel HPL is based on Open HPL 2.2 and is available as a part of the Intel MKL library. Both HPL implementations use the same hardware and same Intel MPI library.

OpenHPL uses one core per MPI process while Intel HPL uses all cores per node for each MPI process. Hence, the optimal  $P \times Q$  combination for each HPL implementation is different, where  $P$  and  $Q$  are the rows and columns of the MPI process grid in the benchmark. This allows for more validation scenarios while having no impact on the validation process as we are not comparing the variance of the two HPL implementations. For the given architecture, the HPL block size used is  $nb = 192$ . The efficiency of the BLAS operations is evaluated using the methodology discussed in section III-B1. The theoretical CPU peak performance and memory bandwidth are given as inputs to the simulator.

Figure 5 compares the simulated performance of OpenHPL

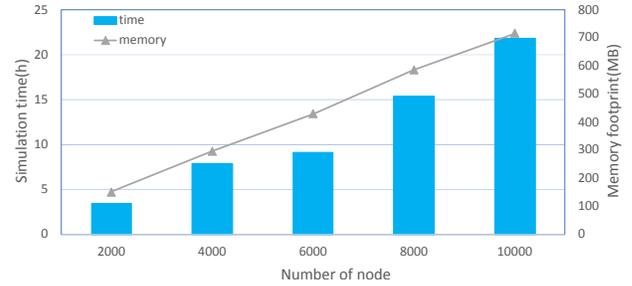


Fig. 7: HPL simulation time and memory footprint.

against the measured performance on 1 core to 4 nodes. Figure 6 shows the validation results of Intel HPL with node numbers scaling from 1 to 4. The performance on 8 and 16 nodes is predicted using the simulator. Overall, our framework achieves high accuracy at varying concurrency with an average of 3.7% discrepancy between the simulated and measured performance.

#### B. Simulation scalability

To evaluate the scalability of our framework, we simulate an HPC system consisting of 10,008 nodes. These nodes are connected using a two-level fat-tree topology. In total, 556 36-port switches are used at the edge level and 18 556-port switches are used at the core level. Each of the edge switches has 18 ports dedicated to connecting servers. The other 18 ports of each edge switch are connected to the core layer. In this scenario, the network of this hypothetical system may not be fully optimized as our goal is to evaluate the scalability of the simulator. The other hardware components are kept the same as those used for the experiments in the previous section.

The number of MPI processes and the matrix size are the two key factors impacting the HPL simulation time and memory consumption. In this section, we conduct a series of simulations where the matrix size is fixed to  $2 \times 10^7$  while the number of MPI processes varies from 2,000 to 10,000 with a step size of 2,000. The simulation results are shown in Figure 7. The bars in the figure represent the execution time. The largest simulation time is 21.8 hours which simulates 10,000 MPI processes with a matrix size of  $2 \times 10^7$ . The line in Figure 7 represents the memory footprint. The memory consumption grows linearly with the number of MPI processes. Simulating 10,000 MPI processes consumes about 720MB.

#### C. TOP500 HPC systems simulation

The TOP500 list ranks the most powerful supercomputing systems according to their performance on the HPL benchmark. Frontera [19] and PupMaya [20] supercomputers, which rank #5 and #25 on the TOP500 list, respectively, provide enough public information to allow the use of our simulator to predict their HPL performance.

TABLE II: TOP500 systems simulation.

	Real environment	Simulation	
Frontera	Node#	8,808	
	Core#	448,448	
	Memory	1,537,536 GB	550 MB
	Nmax	9,282,848	9,282,848
	Rmax	23,516 TFLOP/s	22,566 TFLOP/s
	Execute time	6.5 h (Estimated)	4.8 h
PupMaya	Node#	4,248	1
	Core#	169,920	1
	Memory	815,616 GB	300 MB
	Nmax	4,748,928	4,748,928
	Rmax	7,484 TFLOP/s	7,558 TFLOP/s
	Execute time	2.7 h (Estimated)	1.7 h

Table II shows the hardware configurations along with the performance reported in the TOP500 list. Frontera consists of 8,008 compute nodes, each node consists of a 2 socket Intel Xeon Platinum 8280 2.7GHz CPU with 28 cores per socket, and a 192GB DDR4 memory operating at frequency 2933 MHz. One thing to note here is that the Cascade Lake processor cannot operate at 2.7GHz continuously when running 512-bit Advanced Vector Extensions (AVX-512) unit and the actual running frequency is around 1.8 GHz. The peak CPU performance, memory bandwidth, and kernels efficiency are given as inputs to the simulator. Furthermore, we configure the simulator to use Frontera’s network topology which consists of six core switches, 182 leaf switches, and Mellanox HDR InfiniBand technology with 100Gbps and 90ns latency per routing hop [21], connected in a two-level fat-tree topology (Half of the nodes in a rack (44) connect to 22 downlinks of a leaf switch as pairs of HDR100 (100 Gb/s) links into HDR200 (200 Gb/s) ports of the leaf switch. The other 18 ports are uplinks to the six core switches). We assume that the routing algorithm is a non-blocking D-mod-K as it is commonly used in fat-tree networks [15]. We also assume default MPI configurations.

The simulation results are shown in Table II. The simulated performance of Frontera is 22,566 TFLOPs, while the Rmax performance reported in the TOP500 list is 23,516 TFLOPs. The error rate is around 4%. The simulator execution time is 4.8 hours with about 550MB memory consumption, which is faster than the actual running time of more than 6.5 hours on the full-system (we estimate the actual time based on the problem size).

PupMaya consists of 4,248 nodes, almost half the size of the Frontera supercomputer. We simulate the HPL performance on PupMaya using our framework and achieve good accuracy. Simulation results are shown in Table II.

## V. USE CASES

In this section, we use HPL as an example to demonstrate the simulation framework capabilities to perform what-if analysis.

In the previous section, the HPL performance on Frontera and PupMaya supercomputers is simulated. These two systems

use Mellanox InfiniBand 100Gbps as their interconnect. Here, we use the simulator to predict the HPL performance on a 200Gbps network. Our simulation results show that the performance of Frontera increases from 22,566 TFLOP to 23,143 TFLOPs, and that of PupMaya increases from 7,558 TFLOPs to 7,854 TFLOPs. The performance improvement rates are 2.6% and 3.9% for Frontera and PupMaya, respectively, which are very low. A closer look at the simulation results shows that network congestion occurs due to the non-blocking routing algorithm used in the fat-tree network. Therefore, in this scenario, the high cost of updating the network does not lead to significant performance improvement.

A large portion of HPC systems on the TOP500 list are equipped with accelerators, such as GPGPU. It is therefore of interest to simulate heterogeneous systems to predict and optimize the performance of scientific applications on emerging large scale systems. HPL CUDA [22] is an open-source HPL implementation for NVIDIA GPU. However, the code was last updated in 2011 and is based on HPL version 2.0. On our local server, this implementation achieved performance is about half the theoretical peak while both Summit [23] and Sierra [24] supercomputers achieve more than 75% efficiency. Unfortunately, although we can correlate the simulator with local measurements, the low compute efficiency is far from practical use for predicting the performance of modern HPC systems.

## VI. CONCLUSION

The exponential increase in core counts expected at exascale will lead to increases in the number of switches, interconnects, and memory systems. For this reason, modeling application performance at these scales and understanding what changes need to be made to ensure continued scalability on future exascale architectures is necessary.

This paper proposes a simulation approach to facilitate this process. Our approach enables full-system performance modeling: (1) the hardware platform is represented by an abstract yet high-fidelity model; (2) the computation and communication components are simulated at a functional level, where the simulator allows the use of the components native interface; this results in a (3) fast and accurate simulation of full HPC applications with minimal modifications to the application source code. This hardware/software hybrid modeling methodology allows for low overhead, fast, and accurate exascale simulation and can be easily carried out on a standard client platform (desktop or laptop). HPL is used to demonstrate the capability and scalability of the simulator. Two supercomputers from the TOP500, Frontera and PupMaya, are simulated with good simulation speed and accuracy. Specifically, the simulation of the HPL benchmark on Frontera takes less than 5 hours with an error rate of four percent.

We are extending our simulation framework in several ways to build a more comprehensive solution for modeling and exploiting the full performance of exascale computing platforms. Multithreading is widely used in HPC applications. In the current implementation, threads are extracted manually.

We are working on automating this process in CoFluent Virtual Thread by enabling the simulation of Linux Pthreads and C++ threads. We also plan to support an automatic privatizing of the global variables when mapping applications processes into virtual threads. Finally, power is a major challenge for exascale systems. We are planning to incorporate power models into the simulation framework to enable the design of energy-efficient hardware and software.

## REFERENCES

- [1] L. S. Blackford, A. Petit, R. Pozo, K. Remington, R. C. Whaley, J. Demmel, J. Dongarra, I. Duff, S. Hammarling, G. Henry, and others, "An updated set of basic linear algebra subprograms (BLAS)," *ACM Transactions on Mathematical Software*, vol. 28, no. 2, pp. 135–151, 2002.
- [2] W. Gropp, W. D. Gropp, E. Lusk, A. Skjellum, and A. D. F. E. E. Lusk, *Using MPI: portable parallel programming with the message-passing interface*. MIT press, 1999, vol. 1.
- [3] M. Adams, J. Brown, J. Shalf, B. V. Straalen, E. Strohmaier, and S. Williams, "Hpgmg 1.0: A benchmark for ranking high performance computing systems," 5 2014.
- [4] J. Dongarra, M. A. Heroux, and P. Luszczyk, "Hpcg benchmark: a new metric for ranking high performance computing systems," Tech. Rep. ut-eecs-15-736, 2015-01 2015. [Online]. Available: <http://www.eecs.utk.edu/resources/library/file/1047/ut-eecs-15-736.pdf>
- [5] H. Casanova, A. Legrand, and M. Quinson, "SimGrid: A Generic Framework for Large-Scale Distributed Experiments," in *Tenth International Conference on Computer Modeling and Simulation (uksim 2008)*, Apr. 2008, pp. 126–131, iSSN: null.
- [6] K. Fujiwara and H. Casanova, "Speed and accuracy of network simulation in the simgrid framework," in *Proceedings of the 2nd international conference on Performance evaluation methodologies and tools*, 2007, pp. 1–10.
- [7] A. F. Rodrigues, K. S. Hemmert, B. W. Barrett, C. Kersey, R. Oldfield, M. Weston, R. Risen, J. Cook, P. Rosenfeld, E. Cooper-Balis, and others, "The structural simulation toolkit," *ACM SIGMETRICS Performance Evaluation Review*, vol. 38, no. 4, pp. 37–42, 2011, publisher: ACM New York, NY, USA.
- [8] Z. Bian, K. Wang, Z. Wang, G. Munce, I. Cremer, W. Zhou, Q. Chen, and G. Xu, "Simulating Big Data Clusters for System Planning, Evaluation, and Optimization," in *Parallel Processing (ICPP), 2014 43rd International Conference on*. IEEE, 2014, pp. 391–400. [Online]. Available: <http://ieeexplore.ieee.org/abstract/document/6957248/>
- [9] K. Wang, Z. Bian, Q. Chen, R. Wang, and G. Xu, "Simulating Hive Cluster for Deployment Planning, Evaluation and Optimization," in *Cloud Computing Technology and Science (CloudCom), 2014 IEEE 6th International Conference on*. IEEE, 2014, pp. 475–482. [Online]. Available: <http://ieeexplore.ieee.org/abstract/document/7037705/>
- [10] K. Wang, Z. Bian, and Q. Chen, "Millipedes: Distributed and Set-Based Sub-Task Scheduler of Computing Engines Running on Yarn Cluster," in *High Performance Computing and Communications (HPCC), 2015 IEEE 7th International Symposium on CyberSpace Safety and Security (CSS), 2015 IEEE 12th International Conferen on Embedded Software and Systems (ICESSE), 2015 IEEE 17th International Conference on*. IEEE, 2015, pp. 1597–1602. [Online]. Available: <http://ieeexplore.ieee.org/abstract/document/7336396/>
- [11] J. Liu, B. Bian, and S. S. Sury, "Planning Your SQL-on-Hadoop Deployment Using a Low-Cost Simulation-Based Approach," in *Computer Architecture and High Performance Computing (SBAC-PAD), 2016 28th International Symposium on*. IEEE, 2016, pp. 182–189. [Online]. Available: <http://ieeexplore.ieee.org/abstract/document/7789339/>
- [12] Q. Chen, K. Wang, Z. Bian, I. Cremer, G. Xu, and Y. Guo, "Simulating spark cluster for deployment planning, evaluation and optimization," in *2016 6th International Conference on Simulation and Modeling Methodologies, Technologies and Applications (SIMULTECH)*. IEEE, 2016, pp. 1–11.
- [13] T. Hoefler, T. Schneider, and A. Lumsdaine, "LogGOPSim: simulating large-scale applications in the LogGOPS model," in *Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing*, ser. HPDC '10. Chicago, Illinois: Association for Computing Machinery, Jun. 2010, pp. 597–604. [Online]. Available: <https://doi.org/10.1145/1851476.1851564>
- [14] A. Degomme, A. Legrand, G. S. Markomanolis, M. Quinson, M. Stillwell, and F. Suter, "Simulating MPI Applications: The SMPI Approach," *IEEE Transactions on Parallel and Distributed Systems*, vol. 28, no. 8, pp. 2387–2400, Aug. 2017.
- [15] E. Zahavi, "D-Mod-K routing providing non-blocking traffic for shift permutations on real life fat trees," *CCIT Report*, vol. 776, p. 840, 2010.
- [16] J. Kim, W. J. Dally, S. Scott, and D. Abts, "Technology-Driven, Highly-Scalable Dragonfly Topology," in *2008 International Symposium on Computer Architecture*, Jun. 2008, pp. 77–88, iSSN: 1063-6897.
- [17] S. Williams, A. Waterman, and D. Patterson, "Roofline: an insightful visual performance model for multicore architectures," *Communications of the ACM*, vol. 52, no. 4, pp. 65–76, 2009, publisher: ACM New York, NY, USA.
- [18] Intel, "Product Brief: Intel® CoFluent™ Studio," library Catalog: [www.intel.com](http://www.intel.com). [Online]. Available: <https://www.intel.com/content/www/us/en/cofluent/cofluent-studio-brief.html>
- [19] TOP500, "Frontera - Dell C6420, Xeon Platinum 8280 28C 2.7GHz, Mellanox InfiniBand HDR | TOP500 Supercomputer Sites." [Online]. Available: <https://www.top500.org/system/179607>
- [20] —, "PupMaya - Apollo 2000, Xeon Gold 6148 20C 2.4GHz, Infiniband EDR | TOP500 Supercomputer Sites." [Online]. Available: <https://www.top500.org/system/179604>
- [21] Frontera, "System Architecture - TACC Frontera User Guide." [Online]. Available: <https://frontera-portal.tacc.utexas.edu/user-guide/system/>
- [22] avidday, "avidday/hpl-cuda," Jan. 2020, original-date: 2011-05-09T19:12:49Z. [Online]. Available: <https://github.com/avidday/hpl-cuda>
- [23] TOP500, "Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband | TOP500 Supercomputer Sites." [Online]. Available: <https://www.top500.org/system/179397>
- [24] —, "Sierra - IBM Power System AC922, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband | TOP500 Supercomputer Sites." [Online]. Available: <https://www.top500.org/system/179398>