

Chapter 2

Hardware in HPC

2.1 Introduction

The knowledge of hardware architecture is essential to reach performances through optimizations. Even if the nowadays software, API, framework or runtime already handle most of optimizations, the last percents of gain always need to be architecture dependent. In this chapter we describe the most important devices architectures from classical processors, General Purpose Graphics Processing Units (GPGPUs), Field Programmable Gate Arrays (FPGAs) and Application-specific integrated circuits (ASICs). in this study we will focus on multi-core processors and GPUs. We based our tests on those devices.

This chapter also details the architecture of supercomputers themselves. This has to go with the description of interconnection network with the most famous interconnection topologies.

We choose to present the architectures in a chronological order following the models presented in the previous chapter with: SISD, MIMD and SIMD/SIMT. We also present the optimizations of technologies through parallelism and memories.

2.2 Architectural improvements and parallel capabilities

The processors, as we know them today, begins theirs history around the 1970s. It is the reflection of the Von Neumann Machine we presented in Chapter I.

2.2.1 Single core processors

The first processors around 1970s were built using a single computation core like describe in the Von Neumann model. In order to reach more performances the first factor was to reach higher frequency in the computational cores reaching up to 4-5GHz for best nowadays processors.

Memory technologies

We first consider different memory technologies exists and the aim is always greater capacity, better speed and bandwidth while keeping the data integrity.

SRAM: The Static Random Access Memory is built using so called "flip-flop" circuits that can store the data as long as the machine is powered. This kind of memory is very expensive to produce due to the number of component needed and the size of the memory. Therefore it is usually limited for small amount of storage. The SRAM is mainly used for cache memory.

Cache is a memory mechanism that is useful to consider when targeting performance.

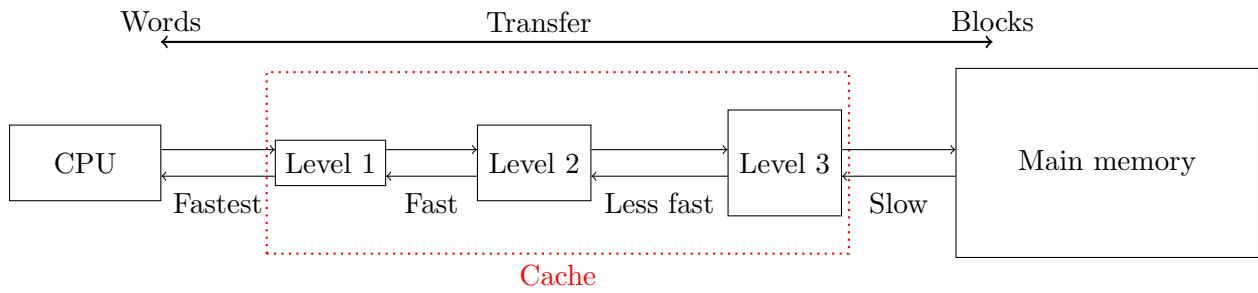


Figure 2.1: Cache memory technology on three levels L1, L2 and L3

Cache memory: This cache technology is presented on figure 2.1. This little memory is built over several levels. The closer to the CPU is L1, then L2 and generally no more than L3 except on specific architecture. When look for a data the CP will first check the L1 cache, otherwise L2 and L3 to get the data to higher level. This is based on the idea that if a data is used, it shall be use again in the near future. Many cache architectures exist like direct, associative, fully associative, etc.

DRAM: The Dynamic Random Access Memory is based on transistors and capacitors to store the binary information. This memory is less expansive to produce but needs to be refresh at a determined frequency however the data are lost. This refresh step is in fact a reading-writing operation on the whole memory at a specific frequency. There is several sub categories of DRAM used in different devices.

Depending on the way the bus are used we can find Single Data Rate, SDR, Double Data Rate, DDR and QDR, Quad Data Rates DRAM memories. The number of data carried can go from 1x to 4x but the limitation of those products is the price of memory constantly rising.

Historically more features were added to this simple core machine, going from 4 bits 1971 (Intel 4004), 8 bits in 1972 (Intel 8008), 12 bits, 16 bits, 32 bits and 64 bits bus size for recent CPUs. The cores/CPU's also get a huge performance gain based on the frequency acceleration, from the 100kHz to GHz nowadays. Plenty of other optimizations were added:

Multiple CPU cores: Multiple CPU cores on the same die. They can have independent or share part of the cache and access to the same main memory. The first machine were the IBM power4 with dual core.

In/Out-Of-Order: In-order-process is the one describes in previous chapter, the control unit fetches instruction in memory, then the operands and the ALU computes, and finally the results is stored in memory. In this model the time to perform an instruction, cumulation of instruction fetching + operand fetching + computation + store the result, can be high and the ALU itself is busy only one step for computation itself. The idea of Out-of-order is to compute the instructions without following the Program Counter order. Indeed, for independent tasks (this is know based on dependency graphs) while the process fetch the next instructions data, the ALU can perform another operation with already available data.

Pre-fetching: When a data is not available in L1 cache, it has to be moved from either L2 to L1 or L3 to L2 to L1 or in the worst case RAM to L3 to L2 to L1. Pre-fecthing technology is a way to, knowing the next instructions operands, pre-fetch the data in closer cache. The pre-fetch can either be hardware or software implemented and can concern data and even instructions.

Vectorization: Processors allows the instructions to be executed at the same time in a SIMD manner. If the same instruction is executed on coalescent data they can be executed in the same clock cycle. Of course this tool require specific care during coding.

Those optimizations can be found either in the classical processor model or accelerators. They are of first importance when we detail optimizations on our codes in part 2 and 3.

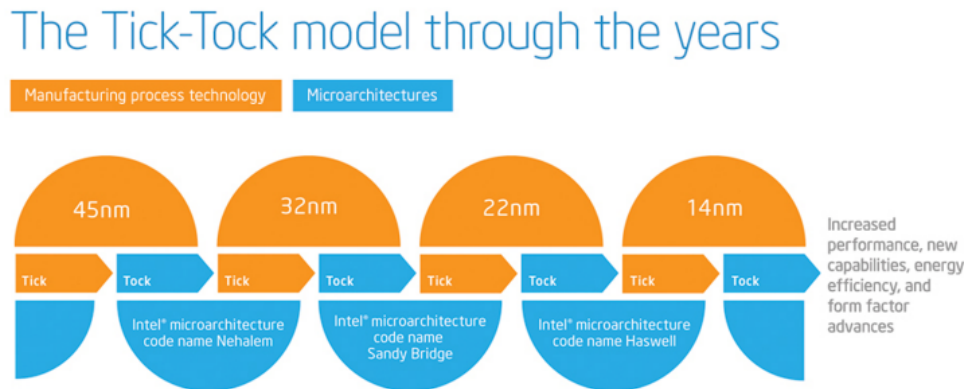


Figure 2.2: Intel Tick-Tock model

2.2.2 Multi-core

Nowadays processors share mostly the same architecture. They are called multi-cores and provide up to 2 to 16 cores and each constructor have its own specificities. Those processors are called "Host" because they are usually bootable and most of the accelerators need to be attached to them in order to work.

Intel

Intel was created in 1968 by a chemist and a physicists, Gordon E. Moore and Robert Noyce, in Mountain View, California. Nowadays processors are mostly Intel ones, this world leader equips around 90% of the supercomputers (November 2017 TOP500 list). In 2007 Intel adopted a production model called the "Tick Tock", presented on figure 2.2.

Since its creation this model followed the same fashion, a new manufacturing technology like shrink of the chip with better engraving on a "Tick" and a new micro-architecture delivered on a "Tock". The Intel processors for HPC are called Xeon and features ECC memory, higher number of cores, large RAM support, large cache-memory, Hyper-threading, etc. compared to desktop processors. Every new processor have a code name. The last generations are chronologically called Westemere, Sandy Bridge, Ivy Bridge, Haswell, Broadwell, Skylake and Kaby lake. Kaby Lake, the last architecture of processor, does not exactly fit the usual "Tick-Tock" process because it is just based on optimizations of the Skylake architecture. It is produce like Skylake in 14nm. This model seems to be hard to maintain due to the difficulties to engrave in less than 10nm with quantum tunneling. This leads to using more many-cores architecture and base next supercomputer generations on hybrid models.

Hyper-threading

Another specificity of Intel processor is Hyper-threading (HT). This technology makes a single physical processor appearing as two logical processors for user's level. In fact a processor embedding 8 cores appears as a 16 cores for user. Adding more computation per node can technically allows the cores to switch context when data are fetched from the memory using the processor 100% during all the computation. A lot of studies have been released on HT from Intel itself [Mar02] to other studies [BBDD06, LAH⁺02]. This optimization does not fit to all the cases and can be disable for normal use of the processors.

ARM

Back in 1980s, ARM stood for Acorn RISC Machine in reference of the first company implementing this kind of architecture, Acorn Computers. This company later changed the name to Advanced RISC Machine (ARM). ARM is a specific kind of processor based on RISC architecture as its ISA despite usual processors using CISC. The downside of CISC machines makes them hard to create and they require way more transistor and thus energy to work. The ISA from the RISC is simpler and requires less transistors to operate. Therefore, the energy required and the heat dissipated is less important. It would then be easier to create massively parallel processors based on ARM. On the other hand, simple ISA impose more work on the source compilation to fit the simple architecture. That makes the instructions sources longer and therefore more single instructions to execute.

The ARM company provide several version of ARM processors named Cortex-A7X, Cortex-A5X and Cortex-A3X respectively balancing highest-performances, performances and efficiency and less power consumption. We find here the same kind of naming as Intel processors.

The new ARMv8 architecture starts to have the tools to target HPC context [RJAJVH17]. The European approach towards energy efficient HPC, Mont-Blanc project¹, already constructs ARM based supercomputers. For the exascale project in Horizon 2020 this project focus on using ARM-based systems for HPC with many famous contributors with Atos/Bull as a project coordinator, ARM, French Alternative Energies and Atomic Energy Commission (CEA), Barcelona Supercomputing Center (BSC), etc. The project is decomposed in several steps to finally reach exascale near 2020. The third step, Mont-Blanc 3, is about to work on a pre-exascale prototype powered by Cavium's ThunderX2 ARM chip based on 64-bits ARMv8.

2.2.3 Many-cores

Several architectures can be defined as many-cores. Those devices integrate thousands of cores that are usually control by a control unit. We can consider those cores as "simpler" since they have to work synchronously and under the coordination of a control unit. They are based on SIMD Flynn taxonomy. Some devices are specific like the Xeon Phi of Intel integrating a hundred of regular processor cores which can work independently.

GPU

GPUs are based on the SIMD model of the Flynn taxonomy presented previously, *Single Instruction, Multiple Data*. The specific execution model is called SIMT (*Single Instruction, Multiple Thread*). It enables the execution of millions of coordinated threads in a data-parallel mode. Two main companies provide GPGPUs for HPC: NVIDIA and AMD.

NVIDIA GPU architecture The NVIDIA company was founded in April 1993 in Santa Clara, Carolina, by three persons in which Jensen Huang, the actual CEO. The company name seems to come from *invidia* the Latin word for Envy and vision for the graphics generation.

Known as the pioneer in graphics, cryptocurrency, portable devices and now AI, it seems to be even the creator of the name "GPU". NVIDIA's GPUs, inspired from visualization and gaming at a first glance, are available as a dedicated device for HPC purpose since the company released the brand named *Tesla*. The public GPUs can also be use for dedicated computation but does not feature ECC memory, double precision or special functions/FFT cores. The different versions of the architecture are named following famous physicists, chronologically: Tesla, Fermi, Kepler, Pascal and Volta. We describe here the Kepler brand GPU on which we based our study.

As presented in figure 2.3, NVIDIA GPUs include many *Streaming Multiprocessors* (SM), each of which is composed of many *Streaming Processors* (SP). In the Kepler architecture, the

¹<http://montblanc-project.eu/>

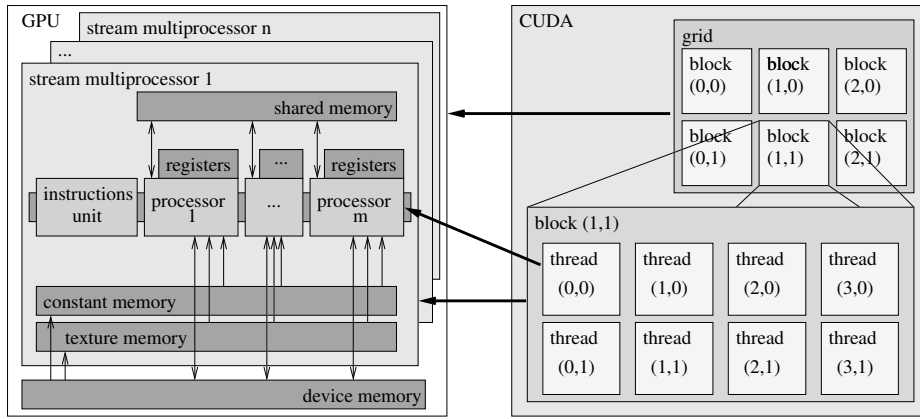


Figure 2.3: NVIDIA GPU and CUDA architecture overview

SM new generation is called SMX. Grouped into *blocks*, *threads* execute *kernels* functions synchronously. Threads within a block can cooperate by sharing data on an SMX and synchronizing their execution to coordinate memory accesses; inside a block, the scheduler organizes *warps* of 32 threads which execute the instructions simultaneously. The blocks are distributed over the GPU SMXs to be executed independently.

In order to use data in a device kernel, it has to be first created on the CPU, allocated on the GPU and then transferred from the CPU to the GPU; after the kernel execution, the results have to be transferred back from the GPU to the CPU. GPUs consist of several memory categories, organized hierarchically and differing by size, bandwidth and latency. On the one hand, the device's main memory is relatively large but has a slow access time due to a huge latency. On the other hand, each SMX has a small amount of shared memory and L1 cache, accessible by its SPs, with faster access, and registers organized as an SP-local memory. SMXs also have a constant memory cache and a texture memory cache. Reaching optimal computing efficiency requires considerable effort while programming. Most of the global memory latency can then be hidden by the threads scheduler if there is enough computational effort to be executed while waiting for the global memory access to complete. Another way to hide this latency is to use streams to overlap kernel computation and memory load.

It is also important to note that branching instructions may break the threads synchronous execution inside a warp and thus affect the program efficiency. This is the reason why test-based applications, like combinatorial problems that are inherently irregular, are considered as bad candidates for GPU implementation.

We give details on the K20X GPU we mainly used in this study in the ROMEO super-computer center. This NVIDIA Tesla Kepler GPU is based on the GK110 graphics processor describes in the white-paper[Nvi12] on 28nm process. The K20X comes in active and passive cooling mode with respectively K20Xc and K20Xm. This GPU embeds 2688 CUDA cores distributed in 14 SMX (we note that GK110 normally provides 15 SMX but only 14 are present on the K20X). In this model each SMX contains 192 single precision cores, 64 double precision cores, 32 special function units and 32 load/store units. In a SMX the memory provides 65536 32-bits registers, 64KB of shared memory L1 cache, 48KB of read-only cache The L2 cache is 1546KB shared by the SMX for a total of 6GB of memory adding the DRAM. The whole memory is protected using Single-Error Correct Double-Error Detect (SECDED) ECC code. The power consumption is estimated to 225W. This GPGPU is expected to produce 1.31 TFLOPS for double-precision and 3.95 TFLOPS of single-precision.

AMD Another company is providing GPUs for HPC, Advanced Micro Devices (AMD). In front of the huge success of NVIDIA GPU that leads from far the HPC market, it is hard for AMD to find a place for its GPGPUs in HPC. Their HPC GPUs are called FirePro. They are

targeted using a language near CUDA but not hold by a single company called OpenCL. An interesting creation of AMD is the Accelerated Processing Units (APUs) which embedded the processor and the GPU on the same die since 2011. This solution allows them to target the same memory.

In the race to market and performances, AMD found an accord with Intel to provide dies featuring Intel processor, AMD GPU and common HBM memory. The project is called Kaby Lake-G and announced for first semester of 2018 but for public, not HPC itself.

Intel Xeon Phi

Another specific HPC product from Intel is the Xeon Phi. This device can be considered as a Host or Device/Accelerator machine. Intel describes it as "a bootable host processor that delivers massive parallelism and vectorization". This architecture embedded multiple multi-cores processors interconnected. This is called Intel's Many Integrated Core (MIC). The architectures names are Knights Ferry, Knights Corner and Knight Landing [SGC⁺16]. The last architecture, Knight Hill, was recently canceled by Intel due to performances and to focus the Xeon Phi for Exascale. The main advantage of this architecture compared to GPGPUs is the x86 compatibility of the embedded cores and the fact this device can boot and use to drive other accelerators. They also feature more complex operations and handle double precision natively. We considered the Xeon Phi in the many-cores architecture despite the fact that it is composed of completely independent processors. This is due to the number of cores that is very high and the fact it can be used as an accelerator instead of the host.

PEZY

Another many-core architecture just appears in the last benchmarks. The PEZY Super Computer 2, PEZY-SC2, is the third many-core microprocessor developed by the company PEZY. The three first machines ranked in the GREEN500 list are accelerators using this many-core die. We also note that in the November 2017 list the 4th supercomputer, Gyoukou, is also powered by PEZY-SC2 cards.

2.2.4 FPGA

Field Programmable Gate Array are devices that can be reprogrammed to fit the needs of the user after their construction. The leader was historically Altera with the Stratix, Arria and Cyclone FPGAs and is now part of Intel. With the FPGAs the user has access to the hardware itself and can design its own circuit. Nowadays FPGA can be targeted with OpenCL programming language. The arrival of Intel in this market promises the best hopes for HPC version of FPGAs. The main gap for users is the circuit building itself, perfect to respond to specific needs but hard to setup.

2.2.5 ASIC

ASICs are dedicated devices constructed for one purpose. An example of ASIC can be the Gravity Pipe (GRAPE) which is dedicated to compute gravitation given mass/positions. Google leads the way for ASIC and just created its dedicated devices to boost AI bots. We also find ASIC in some optimized communication devices like in fast interconnection networks in HPC.

2.3 Interconnection and clusters

2.3.1 Interconnection network

Interconnection network is the way the nodes of a cluster are connected together. Several topologies exist from point to point to multi dimensional torus.

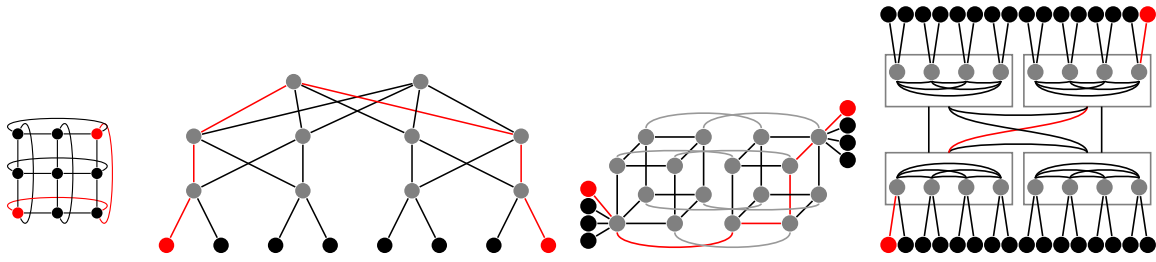


Figure 2.4: Torus, Fat-Tree, HyperX, DragonFly

| Name | Gbs | Year | Name | Gbs | Year |
|-----------|-----|------|-------------|-----|------|
| Single DR | 2.5 | 2003 | Enhanced DR | 25 | 2014 |
| Double DR | 5 | 2005 | Highg DR | 50 | 2017 |
| Quad DR | 10 | 2007 | Next DR | 100 | 2020 |
| Fourth DR | 14 | 2011 | | | |

Table 2.1: InfiniBand technologies name, year and bandwidth

The figure 2.4 is a representation of famous topologies. Each interconnect technology has its own specificity. These networks takes in account the number of nodes to interconnect and the targeted bandwidth/budget. Several declination of each network are not detailed here. the Mesh and the Torus are use as a basis in lower layers of others more complex interconnection networks. A perfect example is the supercomputer called K-Computer describe in the next section. The Fat Tree presented here is a k-ary Fat Tree, higher the position in the tree more connection are found and the bandwidth is important. The nodes are available as the leafs, on the middle level we find the switches and on top the routers. Another topology, HyperX[ABD⁺09], is base on Hyper-Cube. The DragonFly[KDSA08] interconnect is recent, 2008, and use in nowadays supercomputers.

InfiniBand (IB) is the most spread technology used for interconnect with different kind of bandwidth presented in figure 2.1. It provides high bandwidth and small latency and companies like Intel, Mellanox, etc provide directly adapters and switches specifically for IB.

2.3.2 Remarkable supercomputers

The TOP500 is the reference benchmarks for the world size supercomputers. Most of the TOP10 machines have specific architectures and, of course, the most efficient ones. In this section we give details on several supercomputers about their interconnect, processors and specific accelerators.

Sunway Taihulight

Sunway Taihulight is the third Chinese supercomputer to be ranked in the first position of the TOP500 list. A recent report from Jack J. Dongarra, a figure in HPC, decrypt the architecture of this supercomputer[Don16]. The most interesting point is the conception of this machine, completely done in China. The Sunway CPUs were invented and built in China, the Vendor is the Shanghai High Performance IC Design Center.

The SW26010, a many core architecture processor, features 260 cores based on RISC architecture and a specific conception depicted on figure 2.5. The processor is composed of the master core, a Memory Controller (MC), a Management Processing Element (MPE) that manages the Computing Processing Elements (CPE) which are the slaves cores.

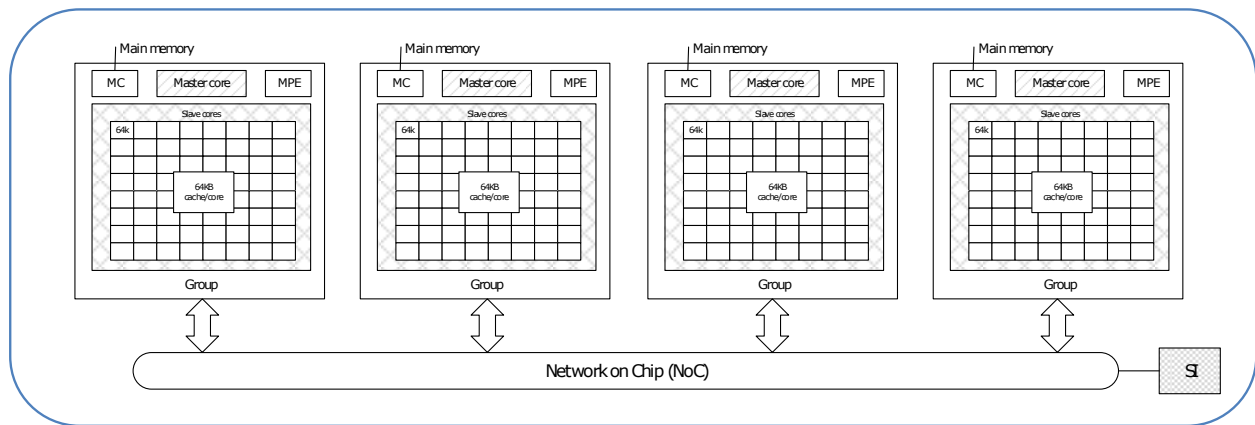


Figure 2.5: Sunway Taihulight node architecture from *Report on the Sunway TaihuLight System*, Jack Dongarra, June 24, 2016.

The interconnect network is called Sunway Network and connected using Mellanox Host Channel Adapter (HCA) and switches. This is a five level interconnect going through computing nodes, computing board, super-nodes and cabinets to the complete system. The total memory is 1.31 PB and the number of cores available is 10,649,600. The peak performance is 125.4 PFLOPS and the Linpack is 93 PFLOPS which induce 74.16% of efficiency.

Piz Daint

The supercomputer of the CSCS, Swiss National Supercomputing Center, is currently ranked 2nd of the November 2017 TOP500 list. This GPUs accelerated supercomputer is a most powerful representative of GPU hybrid acceleration. This is also the most powerful European supercomputer. He is composed of 4761 hybrids and 1210 multi-core nodes. The hybrids nodes embedded an Intel Xeon E5-2690v3 and an NVIDIA Tesla Pascal P100 GPGPU. The interconnect is based on a Dragonfly network topology and Cray Aries routing and communications ASICs. The peak performance is 25.326 TFLOPS using only the hybrid nodes and the Linpack gives 19.590 TFLOPS. The low power consumption rank Piz Daint as 10th in the GREEN500 list.

K-Computer

K-Computer was the top 1 supercomputer of TOP500 2011 list. The TOFU interconnect network makes the K-Computer unique [ASS09] and stands for TORus FUSion. This interconnect presented in figure 2.6 mixes a 6D Mesh/Torus interconnect. The basic units are based on a mesh and are interconnected together in a 3 dimensional torus. In this configuration each node can access to its 12 neighbors directly. It also provide a fault tolerant network with many routes to reach distant node.

AJouter MIRA/SEQUOIA pour parler de IBM, pour le Graph500 = meilleurs supercomputers

Sequoia/Mira

Sequoia supercomputer was top 1 of the TOP500 2012 list. It is based on BlueGene from IBM. The BlueGene project made up to three main architectures with BlueGene/L, BlueGene/P and BlueGene/Q.

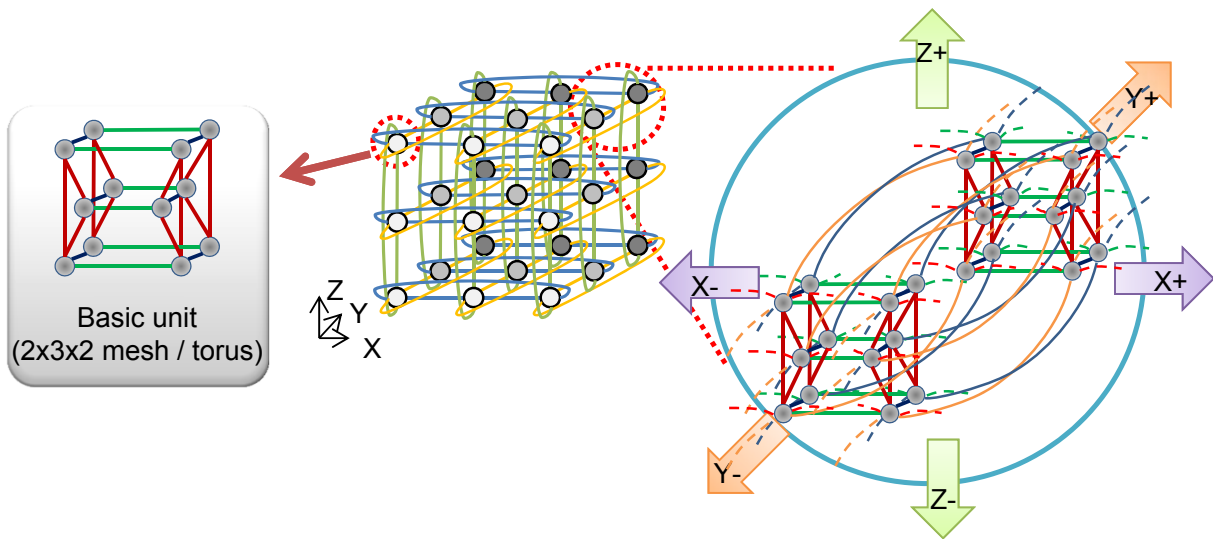


Figure 2.6: TOFU Interconnect schematic from *The K-Computer: System Overview*, Atsuya Uno, SC11

2.4 ROMEO Supercomputer

The ROMEO supercomputer center is the computation center of the Champagne-Ardenne region in France. Hosted since 2002 by the University of Reims Champagne-Ardenne, this so called meso-center (French name for software and hardware architectures) is used for HPC for theoretic research and domain science like applied mathematics, physics, biophysics and chemistry.

This project is support by the Champagne-Ardenne region and the CEA (French Alternative Energies and Atomic Energy Commission), aim to host research and production codes of the region for industrial, research and academics purposes.

We are currently working on the third version of ROMEO, installed in 2013. As many of our tests in this study have been done on this machine, we will carefully describe its architecture.

This supercomputer was ranked 151st in the TOP500 and 5th in the GREEN500 list.

2.4.1 ROMEO hardware architecture

ROMEO is a Bull/Atos supercomputer composed of 130 BullX R421 computing nodes.

Each node is composed of two processors Intel Ivy Bridge 8 cores @ 2,6 GHz. Each processor have access to 16GB of memory for a total of 32GB per node, the total memory if 4.160TB. Each processor if linked, using PCIe-v3, to an NVIDIA Tesla K20Xm GPGPU. This cluster provide then 260 processors for a total of 2080 CPU cores and 260 GPGPU providing 698880 GPU cores. The computation nodes are interconnected with an Infiniband QDR non-blocking network structured as a FatTree. The Infiniband is a QDR providing 10GB/s.

The storage for users is 57 TB and the cluster also provide 195 GB of Lustre and 88TB of parallel scratch file-system.

In addition to the 130 computations nodes, the cluster provides a visualization node NVIDIA GRID with two K2 cards and 250GB of DDR3 RAM. The old machine, renamed Clovis, is always available but does not features GPUs.

The supercomputer supports MPI with GPU Aware and GPUDirect.

2.4.2 New ROMEO supercomputer, June 2018

[Avoir les info et decrire le nouveau ROMEO](#)

2.5 Conclusion

In this chapter we reviewed the most important nowadays hardware architectures and technologies. In order to use the driver or API in the most efficient way we need to keep in mind the way the data and instructions are proceed by the machine.

As efficiency is based on computation power but also communications we showed different interconnection topologies and their specificities. We presented perfect use cases of the technologies in nowadays top ranked systems. They also show that every architecture is unique in its construction and justify the optimization work dedicated to reach performance.

We can see through the new technologies presented here that every one is moving toward hybrids architectures featuring multi-core processors accelerated by one or more devices, many-core architectures. The exascale supercomputer of 2020 will be shape with hybrid architectures and they represent the best of nowadays technology for purpose of HPC. Combining CPU and GPUs or FPGA on the same die, sharing the same memory space can also be the solution.

Bibliography

- [ABD⁺09] Jung Ho Ahn, Nathan Binkert, Al Davis, Moray McLaren, and Robert S Schreiber. Hyperx: topology, routing, and packaging of efficient large-scale networks. In *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*, page 41. ACM, 2009.
- [ASS09] Yuichiro Ajima, Shinji Sumimoto, and Toshiyuki Shimizu. Tofu: A 6d mesh/torus interconnect for exascale computers. *Computer*, 42(11), 2009.
- [BBDD06] Luciano Bononi, Michele Bracuto, Gabriele D’Angelo, and Lorenzo Donatiello. Exploring the effects of hyper-threading on parallel simulation. In *Distributed Simulation and Real-Time Applications, 2006. DS-RT’06. Tenth IEEE International Symposium on*, pages 257–260. IEEE, 2006.
- [Don16] Jack Dongarra. Report on the sunway taihulight system. *PDF*). *www.netlib.org*. Retrieved June, 20, 2016.
- [KDSA08] John Kim, William J Dally, Steve Scott, and Dennis Abts. Technology-driven, highly-scalable dragonfly topology. In *Computer Architecture, 2008. ISCA’08. 35th International Symposium on*, pages 77–88. IEEE, 2008.
- [LAH⁺02] Tau Leng, Rizwan Ali, Jenwei Hsieh, Victor Mashayekhi, and Reza Rooholamini. An empirical study of hyper-threading in high performance computing clusters. *Linux HPC Revolution*, 45, 2002.
- [Mar02] Deborah T Marr. Hyperthreading technology architecture and microarchitecture: a hyperhext history. *Intel Technology J*, 6:1, 2002.
- [Nvi12] C Nvidia. Nvidias next generation cuda compute architecture: Kepler gk110. *Technical report, Technical report, Technical report, 2012.[28]j*, 2012.
- [RJAJVH17] Alejandro Rico, José A Joao, Chris Adeniyi-Jones, and Eric Van Hensbergen. Arm hpc ecosystem and the reemergence of vectors. In *Proceedings of the Computing Frontiers Conference*, pages 329–334. ACM, 2017.
- [SGC⁺16] Avinash Sodani, Roger Gramunt, Jesus Corbal, Ho-Seop Kim, Krishna Vinod, Sundaram Chinthamani, Steven Hutsell, Rajat Agarwal, and Yen-Chen Liu. Knights landing: Second-generation intel xeon phi product. *Ieee micro*, 36(2):34–46, 2016.