

Delivering flexible High-Performance Computing (HPC) to enterprises for green and cost-effective improvement of risk management, a global and inclusive approach

Executive summary

Many organizations today require large quantities of IT resources to optimize risk analysis and management. Finance, energy, pharmaceuticals and large project management are examples of sectors running Monte Carlo algorithms to assess and manage risk, an approach for which High-Performance Computing (HPC) is required. Besides providing computing power, IT installations must also handle the cooling needed both cost-effectively and ecologically.

2CRSi, an international group founded in Strasbourg (France), designs, manufactures and supplies high-performance customized and environmentally friendly servers to meet these needs. The company has worked with a global banking organization to implement a solution for the pricing of complex financial derivatives. Making HPC available as a service by leveraging a hardware platform from a major IT partner, 2CRSi has added advanced cooling in a world-class data center to meet the bank's requirements.

The overall solution offers the bank significant advantages compared to its previous approach, including lower CAPEX and OPEX costs with enhanced data sovereignty. This solution can also benefit enterprises in other industries in a similar way.

Financial risk analysis use case

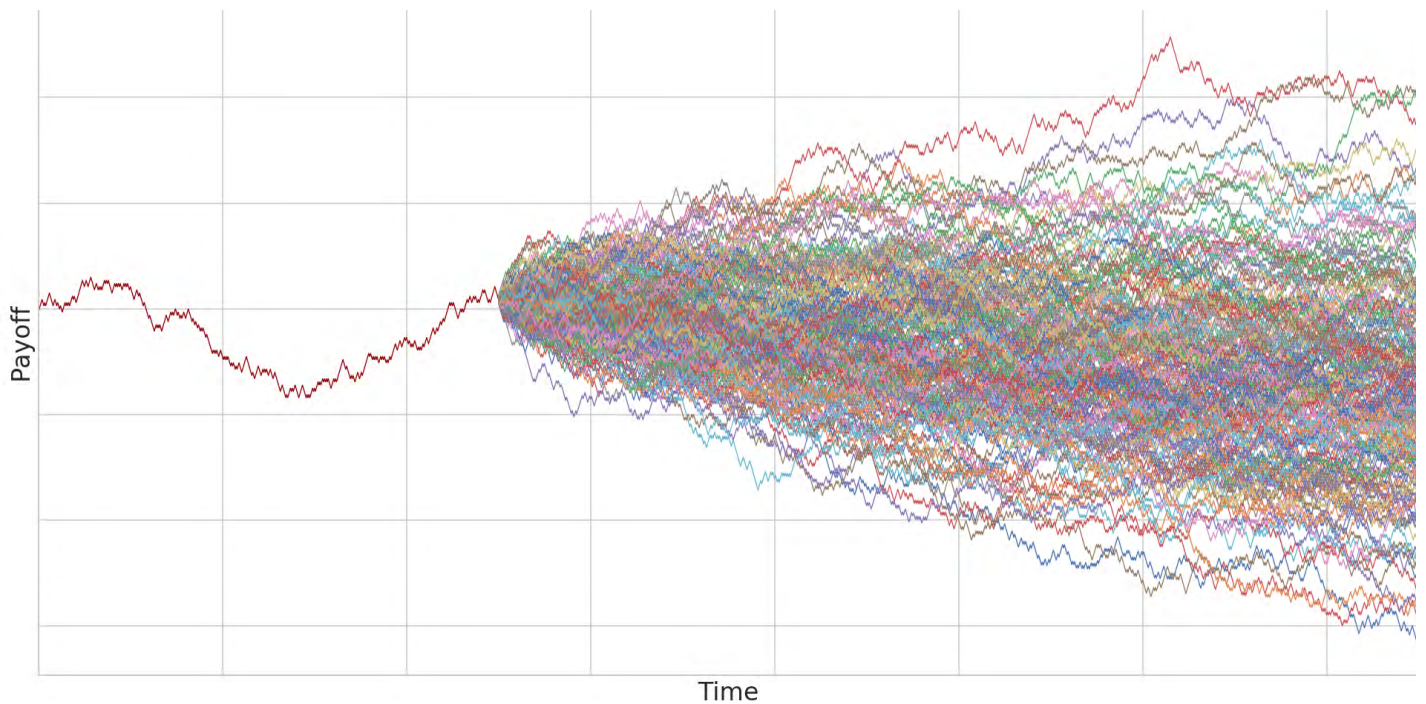
Massive computational resources are an essential feature of today's business and financial markets. A large part of these resources is used on the pricing and risk management of financial assets and their derivatives. Financial assets include ordinary shares, bonds and commodities. Financial derivatives are constructed on this panel of assets as contracts with a future payoff, such as options. The payoff of a derivative depends on the future price or the price trajectory of one or more underlying assets.

Accordingly, option pricing is a major daily activity of banks. One of the most common types of options is the American option, which is a contract allowing the holder to exercise the option rights at any time before and including the day of expiration. Typically, the holder of an American option compares the payoff from immediate exercise with the payoff to be expected from continuing to hold the option. If the immediate payoff is higher, the holder exercises.

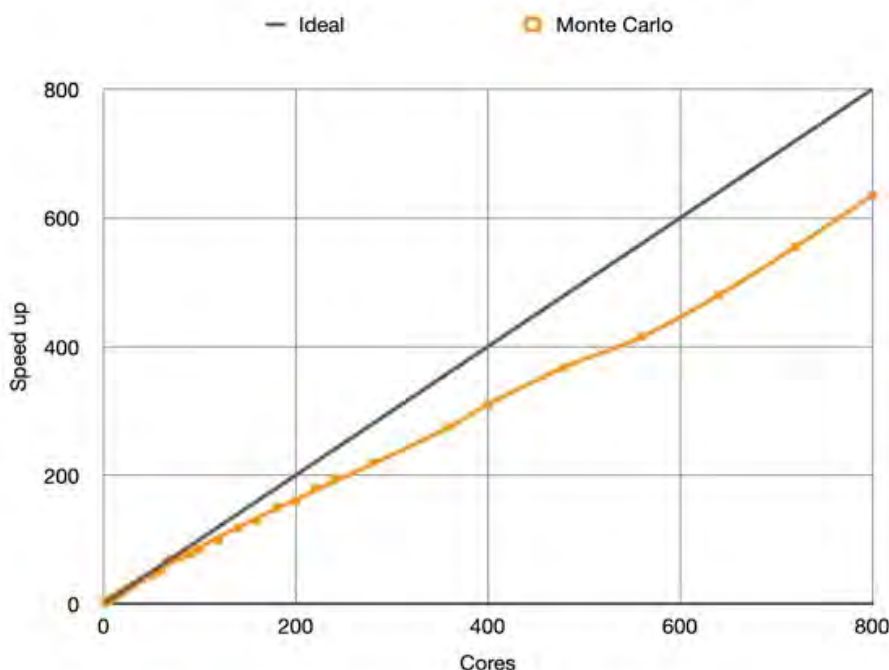
The challenge in correctly pricing derivatives like American options ⁽¹⁾ comes from the stochastic nature of the assets on which they are based. The inherent complexity of American option pricing is compounded when multiple underlying assets are used, each with its own price trajectory.

For financial derivatives, Monte Carlo was first used to solve the pricing problem for European options and later for Asian options. In 2001, Francis Longstaff and Eduardo Schwartz developed a practical technique for pricing American options that they named the least squares Monte Carlo (LSM) approach. Efficiently solving (i.e. ensuring scalability) such option pricing models involving large datasets requires high performance computing.

The graphs below show an example of trajectories resulting from a Monte Carlo computation (Graph 1) and an almost linear between increasing CPUs and acceleration of the execution of MC algorithms.



Graph 1: Trajectories plotted from Monte Carlo analysis.



Graph 2: Monte Carlo analysis scalability (CPU cores vs. execution speed of MC algorithm).

Monte Carlo compute requirements

Nowadays complex mechanisms are often described by non-closed-forms mathematical models. Monte Carlo is a common technique in many domains to solve such increasingly complex situations. As a probabilistic approach combining computing and scalability, Monte Carlo uses a stochastic “generator” to simulate the uncertainty of a parameter. It then processes the repeated sampling of the parameter to obtain a result or range of results. Depending on the context, MC (Monte Carlo) algorithms may need to be run with hundreds of thousands or even millions of samplings for a basic level of accuracy. From a CPU time point of view, the massive amount of floating operations involved in MC algorithms take benefit of parallel processing and vectorization (i.e AVX instructions). Also, the evolution of datasets and the dramatic increase of their sizes require more and more direct memory access. This changed the computer target of MC algorithm from commodity hardware to high-end HPC systems or SMP servers.

Examples of industries and applications using Monte Carlo methods include:



Financial institutions for pricing and risk management of financial assets and their derivatives.



Medical research and pharmaceuticals for in silico drug design such as High-throughput screening, prototyping of medical tests, and prediction and counteraction of viral infections, such as Covid-19. 2CRSi supports *dataagainstcovid19*, a French open data science initiative by sharing computing capacities to build predictive model.



Oil and gas operations for multiscale physics and fluid mechanics calculations for earth layer predictions for drilling and seismology for aftershocks and earthquake predictions.



Meteorology for weather forecasting, global warming simulation and the study of changes linked to pollution.

To meet the need for greater raw compute power, organizations have taken advantage of massive clustered server environment. These “Monte Carlo farms” as they are often known, may be housed on the premises of the organization or built as an outsourced configuration in the cloud or other external data center. They may also be constructed as a customized solution from standard servers or platforms from IT vendors.

Each approach has its advantages and disadvantages. On-premise installations offer greater control but demand greater effort from the user. They may make use of existing idle capacity or conversely require significant capital expenditure (CAPEX). Externally hosted configurations may lead to high operational expenditure (OPEX). Customized solutions can better meet specific needs, but prices may rise. Standard systems and platforms can offer cost-effective performance but must meet requirements.

In each case however, users will pay directly or indirectly for the costs of making the solution available and operating it. Energy required to cool IT installations can be more than 40% of the total energy consumed by a data center ⁽²⁾. This may be close to or even more than the cost of powering the IT systems themselves.

2CRSi, a different solution provider

Established in 2005, 2CRSi provides servers for verticals that require industrial, customized solutions and for data centers seeking high-performance yet dense and green systems. The company has a strong R&D task force with deep expertise in storage and compute equipment (servers, racks, components), and developed new services offerings (data center, on-demand compute power, software integration) in recent years.

In the last 10-15 years, the acceleration of technology introductions helped to face the demand of an increasingly online and real-time world : SSDs, accelerators such as FPGAs, and expanding multi-core processor offerings to address bigger datasets and diverse workloads. One underlying trend has become more prominent in recent years: all these technologies drive up power consumption, and with demand for density bring ever bigger cooling challenges to system builders.

Thanks to a strategic relationship with Intel, 2CRSi has been able to be on top of these developments and offer customers innovative, high-performance compute solutions in dense form factors.

2CRSi has also developed market-leading competences in **IT server cooling solutions**, helping customers to reduce their ecological footprint of their investments and their operating costs.



Air cooling

Classic air cooling can be optimized for higher performance and lower power consumption. By redesigning the implementation of fans and power supply at the rack level with smart fan speed control and power distribution monitoring, power consumption can be reduced by as much as 23%.



Direct liquid cooling

In ultra-high-density configurations where air cooling alone is no longer enough, direct liquid cooling capability is added by 2CRSi to air-cooled racks. Higher power server components like CPUs and GPUs benefit from the liquid cooling for the greatest overall cooling impact. The combined cooling design offers hot-swappable liquid distribution and avoids single points of failure. Energy consumption is reduced, and heat can be reused for other (non-IT) applications.



Immersion cooling

For the highest level of data center efficiency, the latest emerging technology involves that servers are submerged in non-conductive fluid to cool all components uniformly. Data centers using this technology no longer need fans to cool their servers and can reduce or eliminate investment in air conditioning systems. Server component life is extended, energy consumption is reduced, and heat can be reused for other (non-IT) applications. Immersion cooling can contribute as much as 25% savings on data center implementation (CAPEX) and 40% savings on its operating costs (OPEX).

As a subsidiary of 2CRSi, Green Computing (founded in 2019) provides on-demand housing and hosting services in its two French data centers. The company specializes in the rental of high-density compute power backed up by green IT cooling solutions, like direct liquid cooling and immersion cooling.

The Intel® Server System S9200WK product family

Among its strategic relationships with IT vendors, 2CRSi is Platinum Intel® Technology Provider and HPC Data Center Specialist. This high-level status, conferred by Intel on select partners only, gives 2CRSi privileged access to Intel high-performance computing solutions with a range of advantages that it can then pass on to its own customers.

One of these solutions is the **Intel® Server System S9200WK Product Family platform**, featuring Intel® Xeon® Platinum 9200 Processors, which is designed for both HPC and artificial intelligence (AI) application performance. The platform offers:

- Up to 112 powerful processor cores per 2-socket node (448 cores per 2U system)
- 24 memory channels per node delivering ~400GB/sec memory bandwidth (2x bandwidth compared to prior generation platforms)
- Built-in AI capabilities to boost inference performance by up to 30x⁽³⁾
- A choice of an air-cooled or a liquid-cooled data center block.

Further advantages of the S9200WK product family for 2CRSi customers to run Monte Carlo applications include:

- Intel® Advanced Vector Extensions 512 (Intel® AVX512 on CPU) for optimizing floating-point operations in Monte Carlo workloads
- High density and data-centric technology that keeps networking costs down and reduces the number of machines required while being suited to large parameter Monte Carlo models
- Efficiency in sequential tasks and high scalability in parallel framework.

HPC as a service: a specific solution for a global bank

2CRSi has worked with a global bank headquartered in France to provide a new, high-performance yet dense and energy-efficient Xeon-based solution to the bank's computational requirements for Monte Carlo pricing of American options and other financial instruments such as credit derivatives. The specific customer requirements regarding data governance and its green IT strategy, has been the starting point of our client-centric approach.

The bank had an existing risk calculation platform of 120,000 cores, hosted in its own data centers in France and in Iceland with a third-party provider. Its objective was to renew the equivalent of 30,000 cores running an internal application with Monte Carlo algorithm calculations. CPU performance was the overriding consideration, with associated memory, local data storage, and network bandwidth requirements. The new solution had to be turnkey, integrating servers, local networking, hardware hosting, and hardware maintenance.

2CRSi's response to the bank's requirements is based on the Intel® Server System S9200WK product family using Intel® Xeon® Platinum 9242 dual processors with 96 cores per node for a total capacity of 30,000 cores. The hosting is set up in the Green Computing data center with direct liquid cooling for a dedicated solution of MC farms within a green IT context.

As an optional service, a further 20,000 cores are made available as on-demand extra capacity, billed on an hourly basis. 2CRSi brings the bank the possibility of "one-click" deployment with customizable templates and API/extranet management of hundreds of additional servers for these additional compute resources.



Intel® Xeon®
Scalable processors

The Green Computing data center selected for this solution offers tier III+ availability. With power input of 7 megawatts from the main French electricity provider, the data center provides 1.5 megawatt of direct liquid cooling and 2.5 megawatts of immersion cooling.

By coupling this infrastructure to the heat recovery system, Green Computing can heat the building in winter and thus obtain a PUE of less than 1. This solution for heating the building with the heat from the supercomputer prevents the emission of around 200 tons of CO2 per year.

Green Computing world-class density and cooling

The Green Computing data center has an industry-leading cooling infrastructure in several respects: reuse of heat given off by servers in the data center (see below) or installation of large immersion cooling units (Immersion cooling technology achieves extremely high densities, over 100 kW per tank).

The Power Usage Effectiveness (PUE) of the system in the Green Computing data center is less than 1.1 in summer. In winter, the data center achieves a PUE of less than 1, a feat unheard of in many other data centers, by transforming the heat from servers into heating for the building of 65,000m².

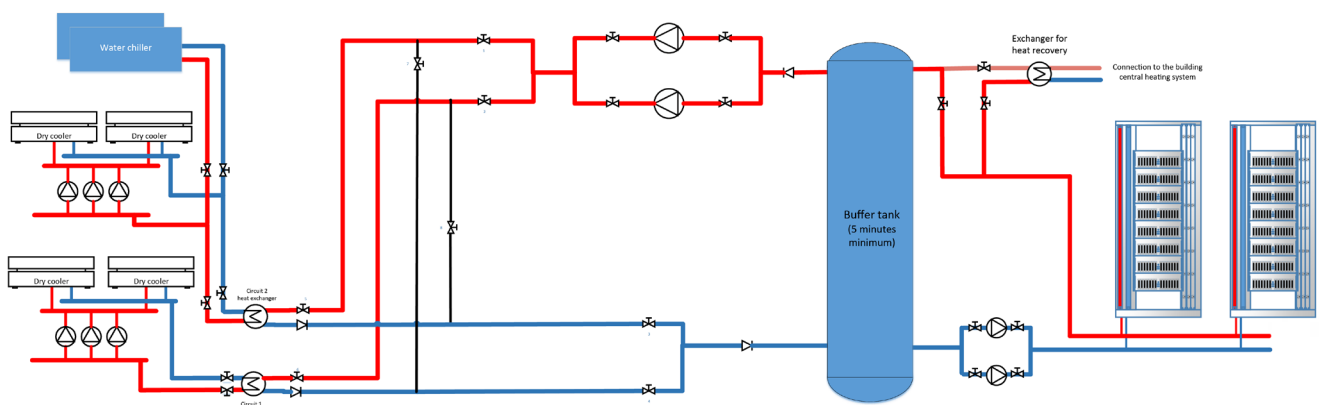


Image 1: Direct Liquid Cooling water flow

The dedicated room in the data center is fully equipped with 3 Direct Liquid Cooling systems with a unit power of 500 kW, for a total of 1.5MW to power the supercomputers. Direct Liquid Cooling technology achieves extremely high densities, over 70 kW per rack.

The combination of the Intel® Server System S9200WK platform and the Direct Liquid Cooling system allows free cooling all year round. The overall PUE of the installation being less than 1.1.

The investment in the Direct Liquid Cooling system in the Green Computing data center is less than half the outlay for comparable commercial systems and components to achieve the same cooling performance. These other systems and components do not necessarily offer the features listed above either. Green Computing can thus pass on operational, ecological, and economic benefits to its customers.

Summary of solution advantages

By switching to the solution from 2CRSi, the banking customer not only moved to leading edge servers but also reduced its costs. 2CRSi, through its close collaboration with Intel as a Platinum Partner and HPC Data Center Specialist, has been able to optimize the TCO for the bank of the Intel® Server System S9200WK platform (both CAPEX and OPEX savings). Moving to the Green Computing data center then gave the bank significant OPEX cost savings compared to the remote rental of data center facilities in Iceland, the bank's previous solution. The bank also improved its data sovereignty and lessened its overall environmental impact by moving its risk calculation and data back to France.

Sovereignty implies a robust data governance considering data mechanism and network engineering. Through expert matchmaking of its group competencies, 2CRSi can provide guidance to clients in the construction of their own dedicated cloud computing solution.

The advantages of this 2CRSi solution can also apply to any organization that requires high-performance computing for risk management or other applications.

To find out more about 2CRSi solutions that can help your business run Monte Carlo algorithms whilst managing risk better and more cost-effectively, contact us today.

www.2crsi.com → contact@2crsi.com →

[Find out more about Intel solutions in the Financial Services Industry](#) →

(1) Longstaff-Schwartz paper citation. <https://people.math.ethz.ch/~hjfurrer/teaching/LongstaffSchwartzAmericanOptionsLeastSquareMonteCarlo.pdf>
(2) https://www.researchgate.net/publication/317308758_Cooling_Energy_Consumption_Investigation_of_Data_Center_IT_Room_withVertical_Placed_Server
(3) <https://www.intel.com/content/www/us/en/benchmarks/server/xeon-scalable/platinum-9200-performance.html>

Using Intel® Deep Learning Boost (Intel® DL Boost) combined with Intel® Optimization for Caffe*, new breakthrough levels of performance can be achieved. Up to 30x improvement in inference performance on Intel® Xeon® Platinum 9282 processor (56 cores) w/ Intel® Deep Learning Boost (Intel® DL Boost) for ResNet-50 (image classification workload) vs. Intel® Xeon® Platinum 8180 processor at launch.

Max Inference Throughput

Intel® Xeon® Platinum 9282 processor: Tested by Intel as of 3/04/2019. DL Inference: Platform: Dragon rock 2S Intel® Xeon® Platinum 9282 processor (56 cores per socket), HT ON, turbo ON, Total Memory 768 GB (24 slots/ 32 GB/ 2933 MHz), BIOS: SE5C620.86B.0D.01.0241.112020180249, Centos* 7 Kernel 3.10.0-957.5.1.el7.x86_64, Intel® Deep Learning Framework: Intel® Optimization for Caffe* version: <https://github.com/intel/caffe> Commit id: 362a3b3, ICC 2019.2.187 for build, MKL DNN version: v0.17 (commit hash: 830a10059a018cd2634d94195140cf2d8790a75a), model: https://github.com/intel/caffe/blob/master/models/intel_optimized_models/int8/resnet50_int8_full_conv.prototxt, BS=8, synthetic Data:3x224x224, 28 instance/2 socket, Datatype: INT8. BKMs for running multi-stream configurations on Xeon: https://www.intel.ai/wp-content/uploads/sites/69/TensorFlow_Best_Practices_Intel_Xeon_AI-HPC_v1.1_Q119.pdf.

Notice

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products. For more complete information visit www.intel.com/benchmarks.

Performance results are based on testing as of dates shown in configurations and may not reflect all publicly available updates. See backup for configuration details. No product or component can be absolutely secure. Your costs and results may vary. Results have been estimated or simulated. You may not use or facilitate the use of this document in connection with any infringement or other legal analysis concerning Intel products described herein. You agree to grant Intel a non-exclusive, royalty-free license to any patent claim thereafter drafted which includes subject matter disclosed herein. No license (express or implied, by estoppel or otherwise) to any intellectual property rights is granted by this document.

All product plans and roadmaps are subject to change without notice. The products described may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Intel disclaims all express and implied warranties, including without limitation, the implied warranties of merchantability, fitness for a particular purpose, and non-infringement, as well as any warranty arising from course of performance, course of dealing, or usage in trade.

Intel technologies may require enabled hardware, software or service activation.

© Intel Corporation. Intel, the Intel logo, Xeon and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

Copyright © 2020 2CRSi. All Rights reserved. Other names and brands may be claimed as the property of others. All specifications are subject to change without notice.