

Assessing Cloud ROI for HPC and Enterprise

Rob Farber



The decision process cannot be couched in terms of traditional calculation

based on direct costs

A complicated decision: To purchase infrastructure or run remotely in the cloud? Bandwidth and data security issues provide the easiest gating factors to evaluate, because an inability to access data kills any chance of using remote infrastructure, be it the public cloud or at a remote HPC center. If running remotely is an option, then the challenge lies in determining the return on investment (ROI) for the remote and local options including:

1. using a hybrid local and remote solution
2. running entirely in a remote cloud
3. pursuing a traditional infrastructure procurement

Be aware of any technical biases that might affect the analysis. For example, there tends to be significant organizational pride associated with owning a large hardware deployment like a supercomputer or compute cluster, in which case legacy design decisions that justified previous hardware procurements might cause pushback against a decision to move to a remote solution, plus provide future headaches after a decision to run remotely has been made.

Estimating Costs

The ROI on direct costs provides a quantifiable estimate on the hard monetary return from the use of a product or service. For a remote solution, the first phases of the transition can be captured in a spreadsheet. For example, the cost difference between investing in a \$100,000 cluster versus a \$1,500-per-month fee for cloud computing, or a negotiated fee for a pre-determined number of hours of compute time at a supercomputing center versus a \$30M grant for a supercomputer procurement and associated building, manpower and maintenance fees.

Indirect costs require a judgment call, and can be the basis by which upper management is judged for making a “good” or “bad” decision. The unfolding of daily events will act as the forces on the teeter totter of opinion. So, a good day will be when the remote solution works well, and a bad day is when anything interferes with the production workflow.

Hybrid Private/Public Strategies

Rather than taking an “all or nothing” approach to run everything remotely, many situations allow the indirect costs to be evaluated by migrating the organization to a private cloud framework and gradually outsourcing more and more of the workload to a remote provider.

On the upside, remote infrastructure enables product/technology elasticity and provides access to resources for short time periods that can scale according to variable demands in workload. Hybrid cloud configurations provide the flexibility and cost benefits of the public cloud with the bandwidth, security and control of the private cloud. Combined, these factors contribute to overall cloud savings, but a hybrid private/public cloud strategy enables greater agility.

It is difficult to put a monetary value on agility that is not subjective and based to some degree on opinion and interpretation, but agility can contribute significantly to an organization’s competitive edge. The ability to move quickly to capitalize on new ideas can provide a big-time ROI, but one that is difficult to encapsulate in an ROI measure. Bottom line: agility benefits both commercial and scientific institutions.

Public Cloud

Market leaders like Google and Amazon Web Services have positioned themselves to be the marketplace that joins software as a service (SaaS) providers, the public cloud and private cloud providers. Whether an HPC or enterprise application, public exposure is the way to gain users, get funding and achieve market success. Effectively, public cloud providers act as “App Stores” for both enterprise and HPC applications. For example, it is straight-forward for an application provider to create a turnkey virtual machine to run an application — regardless of installation complexity, extensiveness of dependency chaining, or painful reliance package version numbers. Private and hybrid cloud users only need to start the remote instances (and, thus, see none of the complexity) to evaluate the application with minimal effort. If they like the functionality when running in the public cloud, customers can then pay to immediately start using the application, or download and install the software on their private cloud. Customers can expect higher performance because their private cloud is local to the organization and is not subject to the bandwidth, latency and data security concerns of the public cloud.

Basically, the public cloud provides a fantastic way to connect software developers, including scientists, HPC centers and commercial companies with customers and users and to “get the word out” about good work. Meanwhile, service providers like Amazon Web Services (AWS) and Google are shrewdly positioning themselves as the storefront and infrastructure providers that join application development efforts to both public and private cloud customers.

The public cloud is also good for the management team because the performance of cloud-based applications can be tested by individual technology teams with little or no cost on remote sites. Similarly, it is possible to move performance-friendly (e.g. latency and bandwidth tolerant) cloud-based applications from the private to a public cloud without impacting user workflows. This makes it possible for the management team to mitigate risk and plan for smaller procurements while still providing a way to handle unexpected peak loads and minimize headaches from institutional pride in hardware ownership. In short, moving to a cloud infrastructure, be it public or private, opens the door to risk management, planning flexibility, deployment agility, institutional publicity, and the potential to exploit the latest technology — such as GPUs and Knights Landing — without risking the creation of user- or application-specific hardware deployments.

Leadership-class Computing

In the HPC world, most leadership-class supercomputers spend most of their time running many smaller jobs. However, it is the ability of a leadership-class supercomputer to solve a single large job, or capability computing that is the *Raison d'être* (reason for being) for all leadership-class supercomputers. Capability computing means that the maximum computing power is applied to solve a single large problem in the shortest amount of time. (In contrast, capacity computing refers to workloads composed of a mix of smaller jobs that are queued to keep a machine running at capacity.)

The discussion about remote computing both facilitates and hinders the argument for leadership-class computers. On one hand, offloading cloud-infrastructure applications means that smaller remote devices can support the production computing workflows, thus leaving the leadership class supercomputers free to run the computationally groundbreaking capability workflows. The challenge with capability computing is that it effectively dedicates that single world-class resource to a single user, which is politically divisive in the user community and challenging to manage. In many cases, the leadership-class supercomputer can be so heavily utilized by half-, quarter-, or smaller-sized jobs that also accomplish wonderful science and set the stage for the dedicated and heroic capability runs, that it makes it very difficult to decide when and how much time should be dedicated to solo large jobs versus capacity workloads. Again, this highlights the challenges of making a subjective estimation of ROI; in this case, the ROI is based on the best use of the computing resource.

Analyzing Alternatives

Large procurements require extensive planning and analysis of alternative scenarios to ensure the best use of the funds. When planning for the Chinook supercomputer procurement at the US Pacific Northwest National Laboratory (PNNL), a viable alternative that had to be addressed considered running all the PNNL jobs at a remote HPC center, such as the National Energy Research Scientific Computing Center (NERSC). If this alternative won, there would be no need for supercomputer procurement or siting a big machine at PNNL. There were compelling advantages for the Department of Energy to have PNNL run remotely rather than spending \$30M for a new supercomputer — especially as much of those procurement funds would be spent on PNNL overhead and paid to the winning vendor for the machine plus a supporting maintenance contract.

In response to the “run remotely” alternative, PNNL argued that a separate procurement was required to develop and run large NWchem jobs — a focus of the PNNL computational research effort. In particular, the NWchem design team utilized a computational shortcut to accelerate a time-consuming $O(N^7)$ calculation, where the expensive computation was performed once at the start of a run for the user-defined input deck, and the results saved to disk, after which they were retrieved as needed during the remainder of the computational job. It was argued that this runtime characteristic is unique and required an I/O-intensive compute architecture that differed from any other supercomputer accessible to DOE users, hence the need for the PNNL Chinook supercomputer procurement.

The challenge in this argument lies in the motivation to move the NWchem kernels past the I/O-intensive computational shortcut. Once gone, it becomes a challenge to justify another supercomputer procurement. So, how much research effort should be dedicated to the development of newer technologies such as GPU accelerators, Intel Xeon Phi accelerators, and even DIMM-compatible flash memory such as the SanDisk ULLtraDIMM or Micron hybrid flash/DRAM memory? For example, NERSC will be building Cori, a new Intel Knights Landing supercomputer that will be operational in mid-2016. This system could swing the alternatives analysis in favor of PNNL running at NERSC or other remote resources, rather than having their own supercomputer.

It is almost a contradiction of terms to speak about calculating the ROI for the procurement and use of machines in an industry that strives to perform millions of billions of numerical calculations per second, because so many of the factors that affect the ROI are subjective. The best that people can do is realize that the decision process cannot be couched in terms of a traditional ROI calculation based on direct costs. Instead, the ROI must be based on circumstantial, indirect measurements of benefits over time. As a result, ROI should be viewed more as a noisy, quantitative measure that can reflect the progress the technology and management teams make as they work together to migrate the organization to a more flexible computing infrastructure over time.

Conclusion

At this time, a hybrid public/private cloud approach seems to offer risk mitigation and migration opportunities through the ability to evaluate, promote and sell applications on the public cloud while preserving the ability to perform capacity computing within a local, private cloud infrastructure that does not have bandwidth or data privacy concerns. When possible — and the direct-cost ROI indicates a benefit — the decision can be made to migrate user workloads to a remote provider. Further, the organization can start to pursue smaller, more frequent procurement and infrastructure refresh cycles to keep up with the latest in private cloud technology.

Rob Farber is an independent HPC expert to startups and Fortune 100 companies, as well as government and academic organizations. He may be reached at editor@ScientificComputing.com.

Source URL (retrieved on 05/24/2016 - 5:53am):

<http://www.scientificcomputing.com/articles/2014/06/assessing-cloud-roi-hpc-and-enterprise>