

Cloud Programming Paradigms for Technical Computing Applications

Geoffrey Fox, Indiana University
Dennis Gannon, Microsoft

In the past four years cloud computing has emerged as an alternative platform for high performance computing. Unfortunately, there is still confusion about the cloud model and its advantages and disadvantages over traditional supercomputing based problem solving methods. In this note we characterize the ways in which cloud computing can be used productively in scientific and technical applications. As we shall see there is a large set of applications that can run on a cloud and a supercomputer equally well. There are also applications that are better suited to the cloud and there are applications where a cloud is a very poor replacement for a supercomputer. Our goal is to illustrate where cloud computing can complement the capabilities of a contemporary massively parallel supercomputer.

Defining the Cloud.

It would not be a huge exaggeration to say that the number of different definitions of cloud computing greatly exceeds the number of actual physical realizations of the key concepts. Consequently, if we wish to provide a characterization of what works “in the cloud”, we need a grounding definition and we shall start with one that most accurately describes the commercial public clouds from Microsoft, Amazon and Google. These public clouds consist of one or more large data centers with the following architectural characteristics

1. The data center is composed of containers of racks of basic servers. The total number of servers in one data center is between 10,000 and a million. Each server has 8 or more CPU cores and a 64GB of shared memory and one or more terabyte local disk drives. GPUs or other accelerators are not common.
2. There is a network that allows messages to be routed between any two servers, but the bisection bandwidth of the network is very low and the network protocols implement the full TCP/IP stack at a sufficient enough level so that every server can be a full Internet host. This is an important point of distinction between a cloud data center and a supercomputer. Data center networks optimize traffic between users on the Internet and the servers in the cloud. This is very different from supercomputer networks which are optimized to minimize interprocessor latency and maximize bisection bandwidth. Application data communications on a supercomputer generally take place over specialized physical and data link layers of the network and interoperation with the Internet is usually very limited.
3. Each server in the data center is host to one or more virtual machines and the cloud runs a “fabric controller” which schedules and manages large sets of VMs across the servers. This fabric controller is the operating system for the data center. Using a VM as the basic unit of scheduling means that an application running on the data center consists of one or more complete VM instances that implement a web service. This means that each VM instance can contain its own specific software library versions and internal applications such as databases and web servers. This greatly simplifies the required data

center software stack, but it also means that the basic unit of scheduling involves the deployment of one or more entire operating systems. This activity is much slower than installing and starting an application on a running OS. Most large scale cloud services are intended to run 24x7, so this long start-up time is negligible. Also the fabric controller is designed to manage and deal with faults. When a VM fails the fabric controller can automatically restart it on another server. On the downside, running a “batch” application on a large number of servers can be very inefficient because of the long time it may take to deploy all the needed VMs.

4. Data in a data center is stored and distributed over many spinning disks in the cloud servers. This is a very different model than found in a large supercomputer, where data is stored in network attached storage. Local disks on the servers of supercomputers are not frequently used for data storage.

Cloud Offerings as Public and Private, Commercial and Academic

As stated above, there are more types of clouds than is described by this public data center model. For example, to address a technical computing market, Amazon has introduced a specialized HPC cloud that uses a network with full bisection bandwidth and supports GPGPUs. “Private clouds” are small dedicated data centers that have various combinations of the properties 1 through 4 above. FutureGrid is the NSF research testbed for cloud technologies and it operates a grid of cloud deployments running on modest sized server clusters.

The major commercial clouds are those from Amazon, Google (App Engine), and Microsoft (Azure). These are constructed as a basic scalable infrastructure with higher level capabilities -- commonly termed Platform as a Service PaaS. For technical computing, important platform components include tables, queues, database, monitoring, roles (Azure), and the cloud characteristic of elasticity (automatic scaling). MapReduce, which is described in detail below, is another major platform service offered by these clouds. Currently the different clouds have different platforms although the Azure and Amazon platforms have many similarities. The Google Platform is targeted at scalable web applications and not as broadly used in technical computing community as Amazon or Azure, but it has been used by selected researchers on some very impressive projects.

Commercial clouds also offer IaaS Infrastructure as a Service with compute and object storage features and Amazon EC2 and S3 being the early entry in this field. There are four major open source (academic) cloud environments Eucalyptus, Nimbus, OpenStack and OpenNebula (Europe) which focus at the IaaS level with interfaces similar to Amazon. These can be used to build “private clouds” but the interesting platform features of commercial clouds are not fully available. Open source Hadoop and Twister offer MapReduce features similar to those on commercial cloud platforms and there are open source possibilities for platform features like queues (RabbitMQ, ActiveMQ) and distributed data management system (Apache Cassandra). However, there is no complete packaging of PaaS features available today for academic or private clouds. Thus interoperability between private and commercial clouds is currently only at IaaS level where it is possible to reconfigure images between the different virtualization choices and there is an active cloud standards activity. The major commercial virtualization products such as VMware and Hyper-V are also important for private clouds but also does not have built-in PaaS capabilities.

We expect more academic interest in PaaS as the value of platform capabilities become clearer from the ongoing application work such as that described in this paper for Azure.

Parallel Programming Models

Scientific Applications that are run on massively parallel supercomputers follow strict programming models that are designed to deliver optimal performance and scalability. Typically these programs are designed using a loosely synchronous or bulk synchronous model and the Message Passing Interface (MPI) communication library. In this model, each processor does a little computing and then stops to exchange data with other processors in a very carefully orchestrated manner and then the cycle repeats. Because the machine has been designed to execute MPI very efficiently, the relative time spent in communicating data can be made small compared to the time doing actual computation. If the complexity of the communication pattern does not grow as you increase the problem size, then this computation scales very well to large numbers of processors.

The typical cloud data center does not have a low latency high bandwidth network needed to run communication intensive MPI programs. However, there is a subclass of traditional parallel programs called MapReduce computations that do well on clouds architectures. MapReduce computations have two parts. In the first part you “map” a computation to each member of an input data set. Each of these computations has an output. In the second part you reduce the outputs to a single output. For example, finding the minimum of some function across a large set of inputs. First apply the function to each input in parallel and then compare the results to pick the smallest. Two variants of MapReduce are important. In one case there is no real reduce step. Instead you just apply the operation to each input in parallel. This is often called an “embarrassingly parallel” computation. At the other extreme there are computations that use MapReduce inside an iterative loop. A large number of linear algebra computations fit this pattern. Cluster and other machine learning computations also take this form. Figure 1 below illustrates these different forms of computation. Cloud computing works well for any of these MapReduce variations. In fact a small industry has been built around doing analysis (typically called data analytics) on large data collections using MapReduce on cloud platforms.

Note that the synchronization costs of clouds lie in between those of grids and supercomputers (HPC clusters) and applications suitable for clouds include all computational tasks appropriate for grids and HPC applications without stringent synchronization. Synchronization in clouds includes both effects of commodity networks but also overheads due to shared systems and virtualized (software) system. The latter is illustrated from Azure application reported by Thilina Gunarathne, Xiaoming Gao and Judy Qiu from Indiana University using Twister4Azure an early iterative MapReduce environment. The Multidimensional Scaling application is dominated by two linear algebra steps denoted by alternating blue and red task groups in figure 2. As in many classical parallel applications, each step consists of many tasks that should execute for essentially identical times as they are executing identical number of elemental matrix operations. The graph shows a job with striking fluctuations in execution time for individual tasks, which currently limits scale at which good efficient parallel speed up can be obtained. All tasks must wait until the slowest straggler finishes.

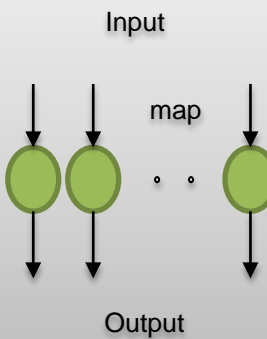
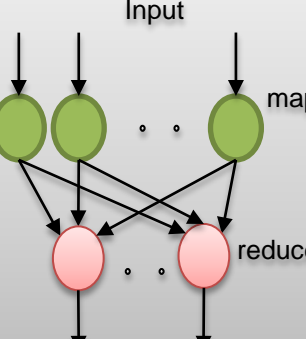
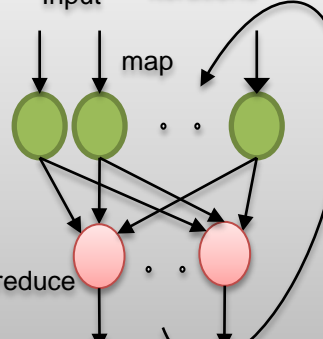
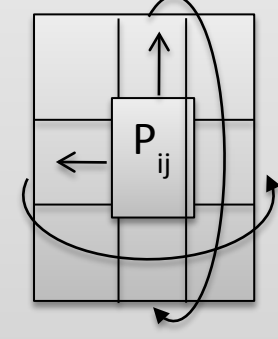
| (a) Map Only | (b) Classic MapReduce | (c) Iterative MapReduce | (d) Loosely or Bulk Synchronous |
|--|---|---|---|
|  |  |  |  |
| BLAST and other multi-protein, multi gene analyses Parametric sweeps Pleasingly Parallel | High Energy Physics Histograms Distributed search Distributed sorting Information retrieval | Expectation Maximization for datamining e.g. Clustering Linear Algebra Page Rank | Many MPI scientific applications such as solving differential equations and particle dynamics |
| ← Domain of MapReduce and Iterative Extensions → | | | MPI |

Figure 1: Forms of Parallelism and their application on Clouds and Supercomputers

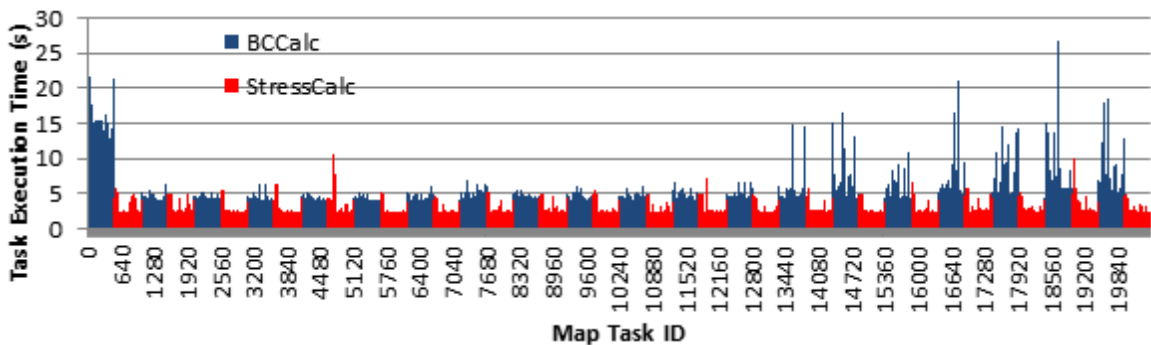


Figure 2: Execution times of individual tasks on Azure for an Iterative MapReduce implementation of Multidimensional Scaling

Classes of Cloud Applications

Limiting the applications of the cloud to classical scientific computation would miss the main reason that the cloud exists. Large scale data center are the backbone of many ubiquitous cloud services we use every day. Internet search, mobile maps, email, photo sharing, messaging, social networks are all dependent upon large data center clouds. These applications all depend upon the massive scale and bandwidth to the Internet that the cloud provides. Note clouds are executing jobs at scales much larger than those on supercomputers but that scale does not come from tightly coupled MPI. Rather scale can come from two aspects: Scale (or parallelism) from a multitude of users and scale from execution of an intrinsically parallel application within a single job (invoked by a single user).

Clouds naturally exploit parallelism from multiple users or usages. The Internet of things will drive many applications of the cloud. It is projected that there will soon be 50 billion devices on the Internet. Most will be small sensors that send streams of information into the cloud where it will be processed and integrated with other streams and turned into knowledge that will help our lives in a million small and big ways. It is not unreasonable for us to believe that we will each have our own cloud-based personal agent that monitors all of the data about our life and anticipates our needs 24x7. The cloud will become increasingly important as a controller of and resource provider for the Internet of Things. As well as today's use for smart phone and gaming console support, "smart homes" and "ubiquitous cities" build on this vision and we could expect a growth in cloud supported/controlled robotics.

Beyond the Internet of things, we have the commodity applications such as social networking, internet search and e-commerce. Finally we mention the "long tail of science" [1, 2] as an expected important usage mode of clouds. In some areas like particle physics and astronomy, i.e. "big science", there are just a few major instruments generating now petascale data driving discovery in a coordinated fashion. In other areas such as genomics and environmental science, there are many "individual" researchers with distributed collection and analysis of data whose total data and processing needs can match the size of big science. Clouds can provide scaling use for this important aspect of science.

Large scale Clouds also support parallelism within a single job with an obvious example of Internet search which has both "usage parallelism" listed above but also parallelism within each search as the summary of the web is stored on multiple disks on multiple computers and searched independently for each query. MapReduce was introduced to support successfully this type of parallelism. The examples given later will illustrate various forms of parallelism discussed above and summarized in figure 1. Today probably the "Map only" or Pleasingly Parallel mode is most common where a single job invokes multiple independent components. This is very similar to the parallelism from multiple users and indeed we will see examples where the different "map tasks" from a "single application" are generated outside the cloud which sees these tasks as separate usages. As well as Internet search, the second parallel category MapReduce is important in the Information Retrieval field with a good example being the well-known Latent Dirichlet Allocation algorithm for topic or "hidden/latent" feature determination. Basic data analysis can often be formulated as a simple MapReduce problem where independent event selection and processing (the maps) is followed by a reduction forming global statistics such as histograms. The final class of parallel applications explored so far is Iterative MapReduce implemented on Azure as Daytona from Microsoft or Twister4Azure from Indiana University. Here Page Rank is a well-known component of Internet search technology that corresponds to finding eigenvector of largest eigenvalue for the

sparse matrix of internet links. This algorithm illustrates those that fit well Iterative MapReduce and other algorithms with parallel linear algebra at their core have been explored on Azure.

Multiple usages or splitting the Internet summary over MapReduce nodes are all forms of a generalized data parallelism. However clouds also support well the functional parallelism seen where a given application breaks into multiple components which typically supported by workflow technology [3-6]. Workflow is an old idea and was developed extensively as part of Grid research. We will see from examples below that its use can be taken over directly by clouds without conceptual change. In fact the importance of Software as a Service for commercial clouds illustrates that another concept Services developed for science by the grid community is a successful key feature of cloud applications. Workflow as the technology to orchestrate or control clouds will continue to be a critical building block for cloud applications.

Important access models for clouds include portals, which are often termed Science Gateways and these can be used similarly to grids. Another interesting access choice seen in some cases is that of the queue either implemented using conventional publish-subscribe technology (such as JMS) or the built in queues of the Azure and Amazon platforms. Applications can use the advanced platform features of clouds (queues, tables, blobs, SQL as a service for Azure) to build advanced capabilities more easily than on traditional (HPC) environments. Of course the pervasive “on demand” nature of cloud computing emphasizes the critical importance of task scheduling where either the built-in cloud facilities are used or alternatively there is some exploration of technologies like Condor developed for grids and clusters.

The nature of the use of data [7] is another interesting aspect of cloud applications that currently is still in its infancy but is expected to become important as for example future large data repositories will need cloud computing facilities. A key challenge as the data deluge grows is how we avoid unnecessary data transport and if possible bring the computing to the data [8-12]. We need to understand the tradeoffs between traditional wide area systems like Lustre, Object stores which the heart of Amazon, Azure and OpenStack storage today and the “data-parallel” file systems popularized by HDFS, the Hadoop File System. We expect this to be a growing focus of future cloud applications.

The Process of Building a Cloud Application

Most attempts to directly port a conventional HPC application to a cloud platform fail¹. It is not unlike early attempts to move “vectorized” HPC application to massively parallel non-shared memory message passing clusters. The challenge is to think differently and rewrite the application to support the new computational and programming models. In the case of clouds the following practices lead to success

1. **Build the application as a service.** Because you are deploying one or more full virtual machines and because clouds are designed to host web services, you want your application to support multiple users or, at least, a sequence of multiple executions. If you are not using the application, scale down the number of servers and scale up with demand. Attempting to deploy 100 VMs to run a program that executes for 10 minutes is a waste of resources because the deployment may take more than 10 minutes. To minimize start up time one needs to have services running continuously ready to process the incoming demand.

¹ The exception to this is the Amazon HPC Cloud which seems to perform very well.

2. **Build on existing cloud deployments.** The cloud is ideal for large map reduce computations so use an existing map reduce deployment such as Hadoop or a similar service.
3. **Use PaaS if possible.** For platform-as-a-service clouds like Azure use the tools that are provided such as queues, web and worker roles and blob, table and SQL storage.
4. **Design for failure.** Applications that are services that run forever will experience failures. The cloud has mechanisms that automatically recover lost resources, but the application needs to be designed to be fault tolerant. In particular, environments like MapReduce (Hadoop, Daytona, Twister4Azure) will automatically recover many explicit failures and adopt scheduling strategies that recover performance "failures" from for example delayed tasks. One expects an increasing number of such Platform features to be offered by clouds and users will still need to program in a fashion that allows task failures but be rewarded by environments that transparently cope with these failures.
5. **Use as a Service where possible.** Capabilities such as SQLaaS (database as a service or a database appliance) provide a friendlier approach than the traditional non-cloud approach exemplified by installing MySQL on the local disk. We anticipate many prepackaged aaS capabilities such as Workflow as a Service for eScience will be developed and simplify the development of sophisticated applications.
6. **Moving Data is a challenge.** The general rule is that one should move computation to the data, but if the only computational resource available is in the cloud, you are stuck if the data is not also there. Moving a petabyte from a laboratory to the cloud over the Internet will take time. The ideal situation is when you can gradually stream the data to the cloud over time as it is being created. But if it exists in one place the best method of moving it is to physically ship the disk drives to the data center. This service is available from some cloud providers.

The Economics of Clouds

Comparing the cost of computing in the cloud to the cost of purchasing a cluster is somewhat challenging because of two factors. First, for most researchers, the Total Cost of Ownership (TCO) of a cluster they purchase is often not visible to them. Space and power are buried in the University's expense overall infrastructure bills. The systems administration for small clusters is often done by students, so it appears to be free of cost. For large, university managed supercomputer centers the situation is different and the costs are known but not widely published. However, there have been some TCO studies. Paterson, et. al. [13] report that the 3-year total cost of ownership for a cluster of servers for research varies between 8 and 15 times the purchase price of the hardware. In a study by IMEX Research, the 3 year TCO of a 512 node HPC cluster averages about \$7 million dollars. This is based on a 2005 configurations with 1U dual core servers with 2 GB memory [14], so it is by no means current technology. However if we compare this to the cost of 1024 cores on Windows Azure running 24x7 for 3 years the total is \$2.6M (with "committed use" discounts). For Linux on Amazon EC2 the cost is \$2.3M. These costs do not include persistent storage or network bandwidth. 100TBytes of triply replicated on-line cloud storage for 3 years is an additional \$400K. Network bandwidth and data transactions add additional costs

Based on these numbers one may conclude a modest advantage to computing in the cloud, but that would miss a critical point. The advantage of the cloud for researchers is the ability to use, and pay for it, on demand. If a

researcher has a need for 1000 cores for a week, the cost is \$21K. If the researcher needs only need a few cores between periods of heavy use the cost is a small increment. In this case there is a clear advantage to cloud computing over purchasing a cluster that is not fully utilized.

Case Studies from the Microsoft Cloud Computing Engagement Program

In the following paragraphs we describe a number of case studies from the Microsoft Cloud Computing Engagement Program which awarded time on Azure to 83 research teams. The examples here are abstracted from reports from the top 30 groups (based on resources consumed to date). These illustrate the models described above with services, workflow, parallelism from multiple users or repetitive internal tasks, and use of MapReduce. They also demonstrate the use of the cloud as a web-enabled “science gateway” where the application is built as a service that can be executed by remote users. Many of the projects in our “top 30” use variations on the MapReduce theme. We have selected this collection of 12 projects because they demonstrate the variety and creativity of research community using the Azure cloud.

1. University of Newcastle

Paul Watson and Jacek Cala of The University of Newcastle use an Azure based system to execute millions of workflows, each of which is a test of a target molecule for possible use as an anti-cancer drug. The scientists use a method known as Quantitative Structure-Activity Relationships (QSAR) to mine experimental data for patterns that relate the chemical structure of a drug to its kinase activity

The architecture of the solution is presented in Figure 3.

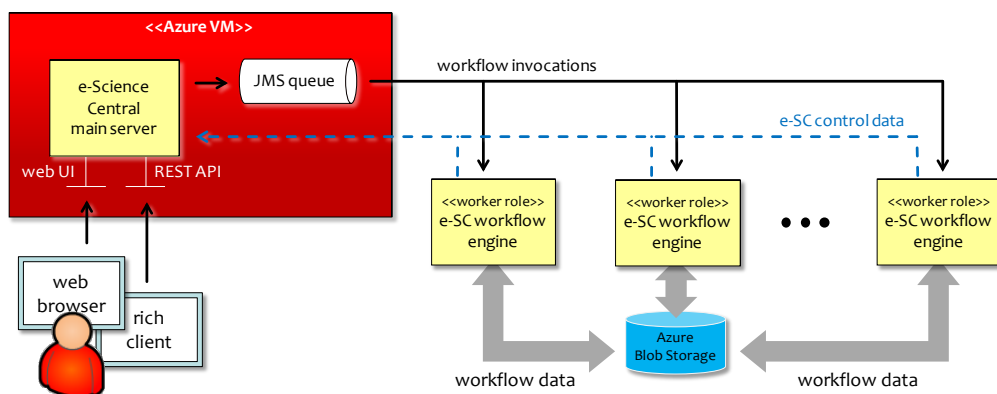


Figure 3. Architecture of the drug discovery system based on e-Science Central and Azure.

Workflows are modeled and stored within the e-Science Central main server which is the central coordination point for all workflow engines. The server dispatches work to a single JMS queue from which it is fetched by the engines. For every input data set the system issues a single top-level workflow invocation, which then results in a number of sub-workflow invocations. Altogether a single input dataset generates from 50 to 120 workflow invocations. The basic unit of work in their system — a workflow invocation — contains a number of tasks. After an engine accepts an invocation from the queue, it executes all included tasks. A task can be as simple as downloading data from blob storage or transposing a data matrix, or as complex as building a QSAR model with neural networks, which can consume over 1 CPU hour.

The Newcastle e-Science Central system also illustrates a two important characteristic of many cloud based scientific systems. This allows the tasks to be spread over 300 cores, with greater than 90% efficiency. The user input is through a web service that can allow multiple users to invoke the same instantiation of the service at the same time. This interactive model is far different from the traditional batch approach used in supercomputing facilities.

2. Georgia State University

This basic programming model which builds the system as a web service that takes user input from the web or rich client to a server in the cloud and then allows a pool of worker servers take work that is queued by the web server is natural for the cloud. The *Crayons* project lead by Sushil K. Prasad of Georgia State University uses a similar approach to doing vector overlay computations for geographic information systems. As shown in Figure 4, the system takes GML files and partitions them into the appropriate sub-domain tasks which are enqueued for worker servers to process. Data is stored in the cloud data storage and pulled to the workers as needed.

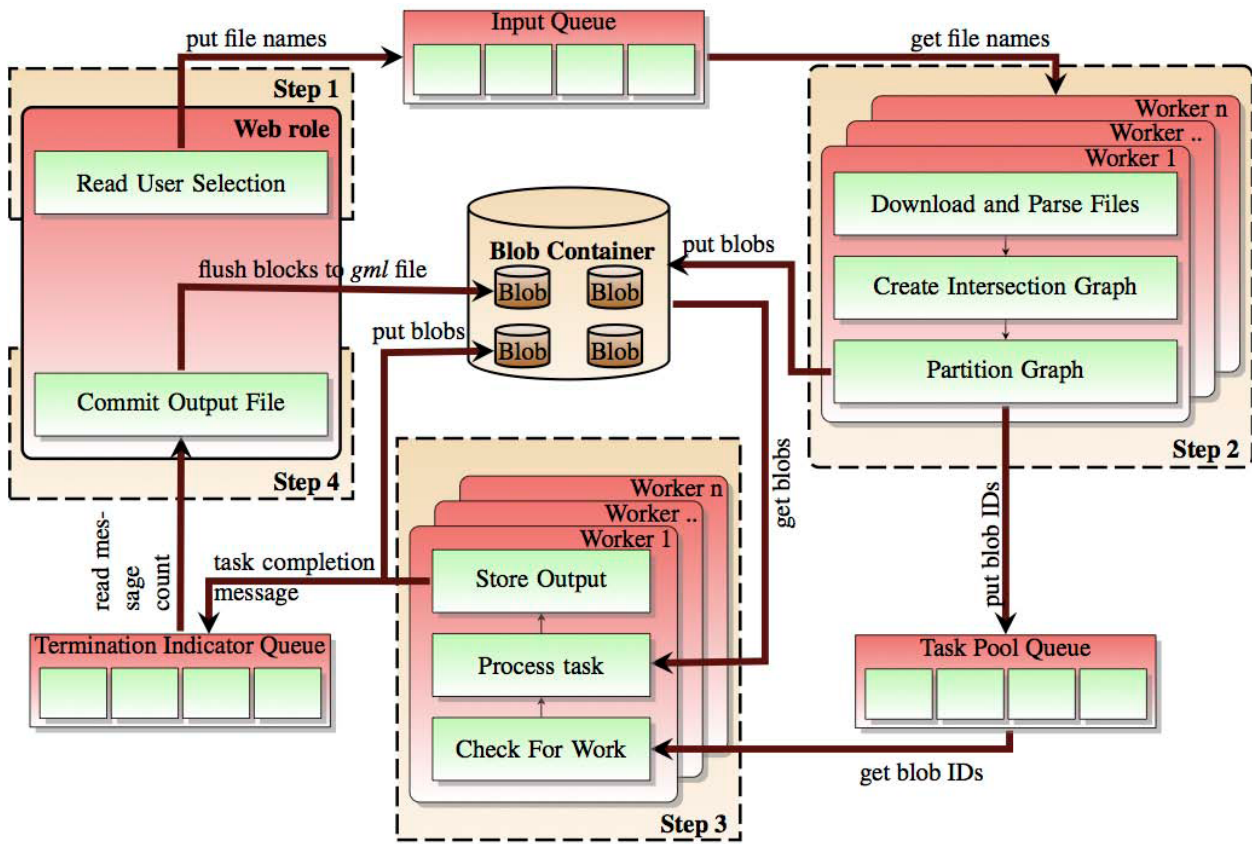


Figure 4. *Crayons* GIS vector overlay processing architecture.

The Crayons project has a fixed workflow in which the tasks are distributed over a pool of workers. *Crayons* is the first distributed GIS system over cloud capable of end-to-end spatial overlay analysis. It scales well for sufficiently large data sets, achieving end-to-end speedup of over 40-fold employing 100 Azure processors.

3. The Microsoft Research - University of Trento Centre for Computational and Systems Biology

Angela Sanger, Michele Di Cosmo and Corrado Priami at COSBI have been investigating the behavior of p53, one of the most important transcription factors in the genome. They have developed a model that includes all known reactions between p53 and DNA, and are fitting this complex model using experimental data obtained in different conditions and using different mutants of p53. COSBI has developed BetaSIM, a simulator, driven by BlenX - a stochastic, process algebra based programming language for modeling and simulating biological systems as well as other complex dynamic systems which they have ported to the cloud. AzureBetaSIM is a data-flow driven parallelization on Windows Azure of the BetaSIM simulator. The aim of AzureBetaSIM is to provide researchers with the ability to quickly run (despite the length of the job queue) a large number of concurrent simulations and, in return, quickly gather research data. More importantly, this approach enables the execution of complex flows based on the aggregation of a large number of identical or slightly different models that, because of the stochastic nature of BetaSIM, will produce different outputs. In AzureBetaSIM, the user can alter the flow dynamically based on the previous outputs and ask the researcher remotely to continue, stop or alter manually the model before a new step by providing a “sparkle” python script to control the high-level flow of execution as in Figure 5. All jobs in this use case are scripted, and the scripts can be provided by the user or generated automatically by another script.

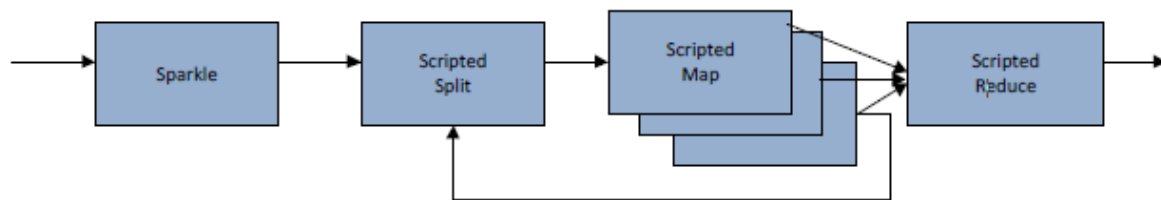
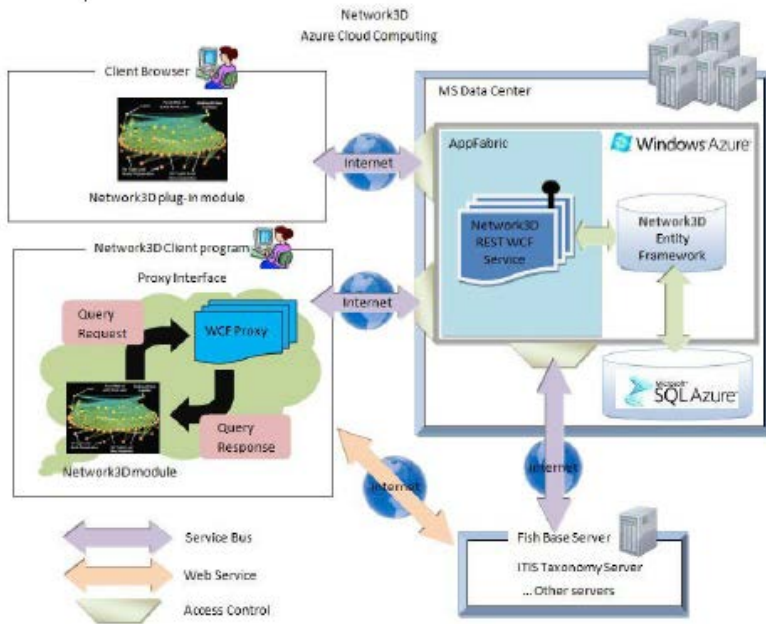


Figure 5. Scripted dataflow execution in COSBI's AzureBetaSIM.

4. The University of North Carolina at Charlotte

Zhengchang Su, Srinivas Arkela and Youjie Zhou, from the Department of Bioinformatics & Genomics and Computer Science at The University of North Carolina at Charlotte are using a similar mapreduce style workflow to annotate regulatory sequences in sequenced bacterial genomes using comparative genomics-based algorithms. Regulatory sequences specify when, how much, and where the genes should be expressed in the cell through their interactions with proteins called transcription factors (TFs). They have built the system using the basic Web Role – Worker Role programming model native to Azure. Web roles are used as interfaces between the system and the users. After jobs are submitted by the users, web roles build job messages and send them through Azure's queues. Worker roles automatically pull messages from the queues and perform real tasks. They have also implemented this same system on top of Hadoop.

5. Pacific Ecoinformatics and Computational Ecology Lab (Berkeley, CA) and the Santa Fe Institute

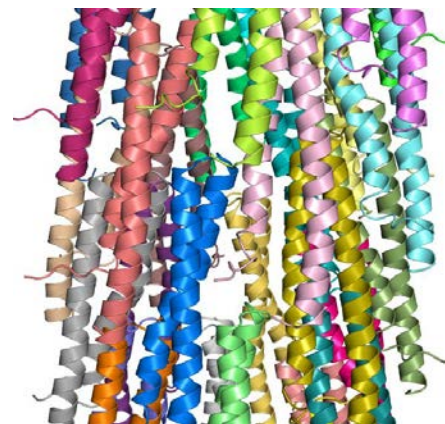


Jennifer Dunne, Sanghyuk Yoon and Neo Martinez from the Pacific Ecoinformatics and Computational Ecology Lab (Berkeley, CA) and the Santa Fe Institute are looking at the grand challenge in the science of ecology of explaining and predicting the behavior of complex ecosystems comprised of many interacting species. The long-term scientific goal is the development of theory that accurately predicts the response of different species and whole ecosystems to physical, biological and chemical changes. The team has developed a tool called Network3D which is used to simulate the complex non-linear systems that characterize these problems. They have ported Network3D to

Azure. The Network3D engine uses Windows Workflow Foundation to implement long-running processes as workflows. Given requests which contain a number of manipulations, each manipulation is delivered to a worker instance to execute. The result of each manipulation is saved to SQL Azure. A web role provides the user with an interface where the user can initiate, monitor and manage their manipulations as well as web services for other sites and visualization clients. Once the request is submitted through the web role interface, the manipulation workflow starts the task and a manipulation is assigned to an available worker to process. The Network3D visualization client communicates through web services and visualizes the ecological network and population dynamics results.

6. University of Washington Baker Lab

A classic example of embarrassingly parallel computation in science is the Seti@home. This is based on volunteer computing where thousands of people make their pc available for an external agent to download tasks to it when it is not being heavily used. A standard framework for these applications is BOINC from Berkeley. In David Baker's lab at the University of Washington, they have built a protein folding application (Rosetta@home) based on BOINC. The problem with traditional volunteer computing is that volunteers are not very reliable. On the other hand, the cloud can be considered a very large pool of capability that can easily be turned into a "high



throughput” BOINC service. To demonstrate this we used 2000 Azure cores to run a substantial folding challenge provided by Dr. Nikolas Sgourakis, a postdoc in the lab. Specifically the challenge was to elucidate the structure of a molecular machine called the needle complex, which is involved in the transfer between cells of dangerous bacteria, such as salmonella, e-coli, and others. One of the main advantages of using Azure is that they didn't have to handle support of the system, a very common problem with Rosetta@home. Of course on a real volunteer system the price is free. Consequently one has a choice. You can get the job done quickly, or you can get it for free.

7. University of Nottingham

A novel use of the cloud for science was provided by Dominic Price at the University of Nottingham as part of the Horizon Digital Britain project. Horizon research focuses on the role of “always on, always with you” ubiquitous computing technology. Their use of Azure was to build support for crowd-sourcing. Specifically they are building a “marketplace” for crowd-sourcing activities. This takes the form of a toolkit that provides an infrastructure in which crowd-sourcing modules that support a particular activity can be developed and then different crowd-sourcing workflows constructed by combining different modules. These modules can then be shared with other users to facilitate their crowd-sourcing activities. This enables non-programmers to reuse existing modules in creating new crowd-sourcing applications. At the heart of the toolkit is the crowd-sourcing factory which is an application running in the cloud, providing the front-end interface for crowd-sourcing administrators (the user group which requests crowd-sourcing activities) and developers (the user group which develops crowd-sourcing modules). This allows the administrators and developers to log in and create crowd-sourcing modules and create an activity from selected modules. Once the workflow has been defined, the factory generates a crowd-sourcing instance, a separate application that is sandboxed from the factory and all other crowd-sourcing instances. This instance contains all of the necessary functionality for recruiting and managing crowd-sourcing participants as well as some storage for storing the results of the activity to be retrieved by the administrator once the activity ends.

8. Old Dominion University

Harris Wu, Kurt Maly and Mohammad Zubair of Old Dominion University are developing a web-based system (FACET) that allows users to collaboratively organize multimedia collections into an evolving faceted classification. The FACET system includes a wiki-like interface that allows users to manually classify documents into their personal document hierarchies as well as the global faceted classification schema, and backend algorithms that automatically classify documents. Evaluating FACET with millions of documents and thousands of users allows them to answer research questions specific to large-scale deployment of social systems that harness and cultivate collective intelligence. For example, how to merge thousands of users' individual document hierarchies into a global schema? How to build a knowledge map of thousands of experts in different domains?

Users of the FACET system are served by multiple Web Role instances. Browsing and classification are supported by queries to SQL Azure. Evaluation has shown that a single Web Role instance of medium size (2 core processors, 3.5GB RAM) can support ~200 concurrent users. The backend classification algorithms run on Worker Role instances. One of the most computing intensive backend procedures is to compute the similarities (both textual and structural) among user-created categories. By dividing the work into 2 extra-

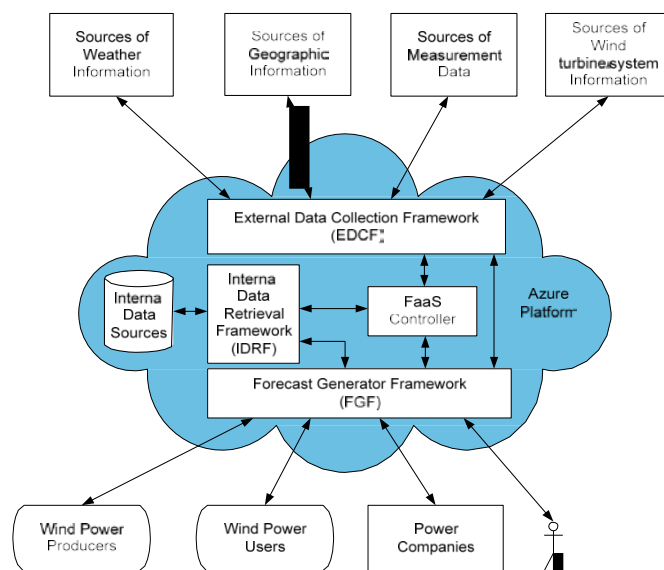
large instances, the pair-wise similarity computation of over 10,000 user categories can be computed with 24 hours.

The FACET system is a research prototype built on Joomla, a popular open-source content management system, and originally on the LAMP stack: Linux, Apache, MySQL and PHP. To deploy the FACET system on Azure, we had to move the metadata repository and session management from MySQL to SQL Azure. To continue to benefit from evolving features in the Joomla community, however, we chose to keep MySQL to support Joomla's core non-data intensive features.

9. Virginia Tech

Kwa-Sur Tam of Virginia Tech have been developing a "Forecast-as-a-Service (FaaS) Framework for Renewable Energy Sources" such as wind and solar. To generate wind power forecasts at specific locations, additional data such as orography, land surface condition, wind turbine characteristics, etc., need to be obtained from multiple sources. In

addition to the diversity of the types of data and the sources of data, there are different forecasting models that have been developed using different approaches. A goal of the project is to support on-demand delivery of forecasts of different types and at different levels of detail for different prices. The FaaS framework consists of the Forecast Generation Framework (FGF), the Internal Data Retrieval Framework (IDRF), the External Data Collection Framework (EDCF) and the FaaS controller. The FGF, IDRF and EDCF all adopt service-oriented architecture (SOA) and the activities of these three



frameworks are orchestrated by the FaaS controller. Since Windows Azure and the associated .NET technologies support the implementation of service-oriented architecture and its design principles, this project can focus on achieving its goals rather than dealing with the underlying support infrastructure.

10. The University South Carolina and the University of Virginia

Jon Goodall at the University of South Carolina and Marty Humphrey at the University of Virginia are creating a cloud-enabled hydrologic model and data processing workflows to examine the Savannah River Basin in the Southeastern United States. Understanding hydrologic systems at the scale of large watersheds is critically important to society when faced with extreme events, such as floods and droughts, or with concern about water quality. This project will advance hydrologic science and our ability to manage water resources. Today, models that are used for engineering analysis of hydrologic systems are insufficient for addressing current water resource challenges such as the impact of land use change and climate change on water resources. The challenge being faced extends beyond modeling and includes the entire workflow from data collection to decision making. The project is being built in three stages. First they will create a cloud-enabled hydrologic model. Second, they will improve the process of hydrologic model

parameterization by creating cloud-based data processing workflows. Third, in Windows Azure, they will apply the model and data processing tool to a large watershed.

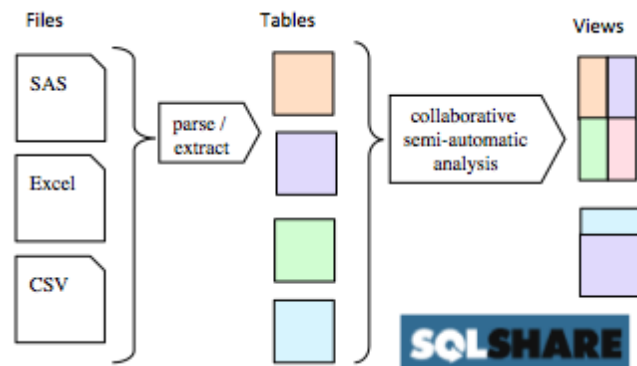
They plan to use Windows Azure for data preparation, model calibration, and large-scale model execution. To date, they have focused on model calibration, whose goal is to search the space of potential model parameters for the best match against observed data. In their system the user uploads her model and then specifies the set of parameters and range of values for each parameter, effectively defining the search space. Their architecture is based on cloudbursting, which uses the Microsoft HPC support extensively. When a user submits a job, the first resources sought are the Microsoft HPC Cluster at the University of Virginia. When there is too much work in the queue, they “burst” onto Windows Azure. This architecture has shown to be very flexible, and they are able to spin up new compute nodes in Windows Azure (and shut them down) whenever they desire. Within Windows Azure, they use Blob storage to prepare the nodes – i.e., when they boot, they automatically load our generic watershed modeling code. They use the Windows Azure Virtual Private Network (VPN) support to make the Windows Azure nodes appear to be local to their enterprise – they have found that this has greatly simplified the use of the cloud.

11. The University of Washington

Bill Howe, Garret Cole, Alicia Key, Nodira Khossainova and Leilani Battle of the University of Washington have developed “SQLShare: Database-as-a-Service for Long Tail Science”. This project addresses a growing problem in science involving the ability of researchers to save, share and query the data results from scientific research. Spreadsheets and ASCII files remain the most popular tools for data management, especially in the long tail. But as data volumes continue to explode, cut-and-paste manipulation of spreadsheets cannot scale, and the relatively cumbersome development cycle of scripts and workflows for ad hoc, iterative data manipulation becomes the bottleneck to scientific discovery and a fundamental barrier to those without programming experience. SQLShare (<http://sqlshare.escience.washington.edu>) is cloud-based relational data sharing and analysis platform that allows users to upload their data and immediately query it using SQL — no schema design, no reformatting, no DBAs. These queries can be named, associated with metadata, saved as views, and shared with collaborators.

The SQLShare platform is implemented as an Azure Web Role that issues queries against a SQL Azure back end. The Azure Web Role implements a REST API and manages communication with the database, enforcing SQLShare semantics when they differ from conventional databases. In

particular, the Web Role manages fault-tolerant and incremental upload of large datasets, analyzes the uploaded data to infer types and recommend example queries, manages authentication with external authentication services, operates on the system catalog, parses and formats data for interoperability with external systems, provides asynchronous semantics for all operations that may operate on large datasets, and handles all REST requests. Windows Azure was essential to the success of the project by empowering a



single developer to build, deploy, and manage a production-quality web service. This work was supported in part by the Gordon and Betty Moore Foundation.

12. Kyoto University

Daisuke Kawahara and Sadao Kurohashi of Kyoto University have been developing a search engine infrastructure, TSUBAKI, which is based on deep Natural Language Processing. While most of conventional search engines register only words to their indices, TSUBAKI provides a framework that indexes synonym relations, hypernym-hyponym relations, dependency/case/ellipsis relations and so forth. These indices enable TSUBAKI to capture the semantic matching between a given query and documents more precisely and flexibly. Case/ellipsis relations have not been indexed in a large scale because the speed of these analyses is not fast enough due to the necessity of referring to a large database of predicate-argument patterns (case frames). To apply case/ellipsis analysis to millions of Web pages of TSUBAKI in a practical time, it is necessary to use 10,000 CPU cores. Because of limits on the Azure fabric controller, it was necessary to divide this into 29 hosted services of 350 CPUs each. This was the largest experiment of any of the research engagement projects.

References

1. P. Bryan Heidorn, *Shedding Light on the Dark Data in the Long Tail of Science*. Library Trends, 2008. 57(2, Fall 2008): p. 280-299. DOI:10.1353/lib.0.0036
2. Peter Murray-Rust. *Big Science and Long-tail Science*. 2008 January 29 [accessed 2011 November 3]; Available from: <http://blogs.ch.cam.ac.uk/pmr/2008/01/29/big-science-and-long-tail-science/>.
3. Taylor, I.J., E. Deelman, D.B. Gannon, and M. Shields, *Workflows for e-Science: Scientific Workflows for Grids*. 2006: Springer. ISBN:1846285194
4. Gannon, D. and G. Fox, *Workflow in Grid Systems, Editorial of special issue of Concurrency&Computation:Practice&Experience based on GGF10 Berlin meeting*. Concurrency and Computation: Practice & Experience, August 2006, 2006. 18(10): p. 1009-1019. DOI:<http://dx.doi.org/10.1002/cpe.v18:10>.
<http://grids.ucs.indiana.edu/ptliupages/publications/Workflow-overview.pdf>
5. Deelman, E., D. Gannon, M. Shields, and I. Taylor, *Workflows and e-Science: An overview of workflow system features and capabilities*. Future Generation Computer Systems, May 2009, 2009. 25(5): p. 528-540. DOI:<http://dx.doi.org/10.1016/j.future.2008.06.012>
6. Yolanda Gil, Ewa Deelman, Mark Ellisman, Thomas Fahringer, Geoffrey Fox, Dennis Gannon, Carole Goble, Miron Livny, Luc Moreau, and Jim Myers, *Examining the Challenges of Scientific Workflows*. Computer, 2007. 40(12): p. 24-32. DOI:10.1109/mc.2007.421
7. Geoffrey Fox, Tony Hey, and Anne Trefethen, *Where does all the data come from?*, Chapter in *Data Intensive Science* Terence Critchlow and Kerstin Kleese Van Dam, Editors. 2011. <http://grids.ucs.indiana.edu/ptliupages/publications/Where%20does%20all%20the%20data%20come%20from%20v7.pdf>.
8. Faris, J., E. Kolker, A. Szalay, L. Bradlow, E. Deelman, W. Feng, J. Qiu, D. Russell, E. Stewart, and E. Kolker, *Communication and data-intensive science in the beginning of the 21st century*. Omics: A Journal of Integrative Biology, 2011. 15(4): p. 213-215. <http://www.ncbi.nlm.nih.gov/offcampus.lib.washington.edu/pubmed/21476843>
9. Gordon Bell, Tony Hey, and Alex Szalay, *COMPUTER SCIENCE: Beyond the Data Deluge*. Science, 2009. 323(5919): p. 1297-1298. DOI:10.1126/science.1170411

10. Gray, J., *Jim Gray on eScience: A transformed scientific method*, Chapter in *The Fourth Paradigm: Data-intensive Scientific Discovery*. 2009, Microsoft Research. p. xvii-xxxi.
11. Gordon Bell, Jim Gray, and Alex Szalay, *Petascale Computational Systems: Balanced CyberInfrastructure in a Data-Centric World (Letter to NSF Cyberinfrastructure Directorate)*. IEEE Computer, January, 2006. 39(1): p. 110-112. <http://research.microsoft.com/en-us/um/people/gray/papers/Petascale%20computational%20systems.doc>
12. Jim Gray, Tony Hey, Stewart Tansley, and Kristin Tolle. *The Fourth Paradigm: Data-Intensive Scientific Discovery*. 2010 [accessed 2010 October 21]; Available from: <http://research.microsoft.com/en-us/collaboration/fourthparadigm/>.
13. Paterson, et. al. "Recovery Oriented Computing (ROC): Motivation, Definition, Techniques, and Case Studies", Computer Science Technical Report UCB//CSD-02-1175, U.C. Berkeley, March 15, 2002.
14. IMEX TCO Analysis High Performance Computing.
http://internet.ziffdavis.com/creatives/appro/Imex_Study_Appro_HyperBlades_Excell_in_HPC.pdf