# Pay as You Go Supercomputing: Does it Work?

*W T Hewitt,* CSAR, Manchester Computing, University of Manchester, Manchester M13 9PL, UK
+44 161 275 6095 Fax: +44 161 275 6800 w.t.hewitt@man.ac.uk
http://www.csar.cfs.ac.uk

**ABSTRACT:** *UK academia now gets its supercomputer service from a consortium of Silicon Graphics, Computer Sciences Corporation and the University of Manchester under the private finance initiative of the UK government. It is essentially a pay-per-use basis. This talk will review the mechanism in use for UK academia.*

## Introduction

UK academia now receives the main part of its supercomputer services from Computing Services for Academic Research (CSAR), which in turn is provided by a consortium of Computer Systems Corporation (CSC), SGI and the University of Manchester (UoM). Currently, its service is based around three machines: a Cray T3E-1200E, a Fujitsu VPP300 and an SGI Origin 2000. The service is operated under the 'Private Finance Initiative (PFI)'[1] of the UK Government whereby we are contracted to provide a basic service, but have the freedom to grow the service in a much more flexible way in response to user demands. This paper will describe the financing model of the service. Many aspects of the service will clearly be omitted from this paper.

This next section in this paper describes briefly how the system works particularly for the user. It is followed comparison with the administrative aspects and funding mechanisms before PFI. One of the main elements of PFI is the sharing of risk between the public and private sectors; this is covered next. The penultimate section describes how the consortium gets paid, and the final section provides a critical appraisal of the situation to date.

## How the System Works Today

### The Research Councils
An academic in the UK who needs resources (people and equipment) to undertake their research makes a proposal to one of the six UK research councils:

- Engineering and Physical Sciences Research Council (EPSRC)
- Natural Environment Research Council (NERC)
- Biotechnology and Biological Sciences Research Council (BBSRC)
- Economic and Social Research Council (ESRC)
- Medical Research Council (MRC)
- Particle Physics and Astronomy Research Council (PPARC)

The research councils fund a small number of central facilities and in particular the supercomputer service. One of the research councils, EPSRC for supercomputing, acts as managing agents on behalf of all.

### The Application Process
A proposal consists of three parts: a scientific case, the costs of the people resources requested and the amount of CSAR resource requested. The scientific case must include justification for all the resources requested. The CSAR resources are calculated via

---

[1] Now called "Public Private Partnership"

the resource calculator, a web form, and a simple example is shown in Figure 1. The next section will describe the elements of the resource. For now notice that the calculator produces two totals: a number of tokens and a notional cost.

| Resource Type | Months 1-6 | Months 7-12 | Months 13-18 | Months 19-24 | Months 25-30 | Months 31-36 | Total | Generic Service Tokens |
|---|---|---|---|---|---|---|---|---|
| T3E Hours (MPP PE Hours) | 1 | | | | | | 1 | 0.024 |
| Origin Hours (Origin CPU Hours) | 1 | | | | | | 1 | 0.025 |
| Guest Service (Generic Service Tokens) | 1 | | | | | | 1 | 1.000 |
| HPD (**Gbytes**) | 1 | | | | | | 1 | 3.876 |
| MPD (**Gbytes**) | 1 | | | | | | 1 | 2.146 |
| HSM/Tape (**Gbytes**) | 1 | | | | | | 1 | 0.298 |
| Training (Days) | 1 | | | | | | 1 | 3.000 |
| Optimisation Support (Days) | 1 | | | | | | 1 | 11.364 |
| Application Support (Days) | 1 | | | | | | 1 | 11.364 |
| | | | | | Total service tokens | | | 33.096 |
| | | | | | Notional Cost | | | £992.89 |

Recalculate

**Figure 1 A simple CSAR resource request**

The proposal is submitted to the appropriate research council who then send it to a number of referees who review the scientific content and we are asked to comment upon our capability to provide the requested resources and whether the application is suited to the facilities. Once the research council has received all the reports and comments they will accept or reject the application.

*Authorization, Registration and Trading*
Once a research council authorizes a grant, the applicant or principal investigator (PI) and CSAR are informed of the number of tokens allocated. Thus the PI receives tokens, and the research council knows how much it will cost them if they are all consumed. The CSAR service then creates accounts etc and using the original application form makes an initial trade of tokens into resources. This information is also used to form the basis of a computational plan that the PI is expected to keep up to date and which we use for planning ahead. Subsequently the PI can trade resources, see below.

Through another web form users associated with the project self register, and through and email to the PI and yet another web form the PI authorizes the user. The outcome of this process is that the user is provided with a username, password and resources on one of more or our systems. Eventually the users will run jobs, create disk files etc. Once a month the accounting software produces a 'bill' that is sent to the Research Councils, and within a few days the consortium receives its funds.

## The Previous System

For the users the proposal, review, allocation (or rejection), registration and use were essentially the same. The important differences are:
- Only CPU resource was quantified on the application. It was allocated in 6 monthly periods; if not used in that period the unused resources were lost
- The supercomputing resources were funded through top-slicing each research council's budget. Therefore individual research councils had already paid for the resource "up-front" and when a proposal was considered the cost of supercomputing was irrelevant. Now a research council must pay for the supercomputer resource and will balance such requests against its overall budget.

The major differences in the before and after PFI mechanisms are to the research councils and the service provider. Previously the RCs as a whole had purchased a machine and subsequent service and thus the RCs took all the risk associated with such a service. The supercomputer was a capital purchase and the maintenance and support service were paid from recurrent funds. Thus the service provider and technology supplier were well placed: they received their funded irrespective of whether the service was used or not. Now the RCs pay for all there supercomputing out of recurrent funds. The Consortium takes on board some of the risks.

The third major difference comes from the new providers. The RCs had to fund the resource they required, so conceptually there was no spare capacity on the system. So it was difficult to sell resources to industry and commerce. As the previous service providers had all been Universities who are not allowed to take the financial risk of investing in additional capacity for subsequent sale to the commercial sector. The consortium, being a private company, has invested in additional resources and is selling these to the commercial sector. Thus every customer on the service gets the capacity they require, but of course the capability available to any one customer is greater.

## Some Background

Early in 1997 the EPSRC set in motion a public procurement for its supercomputer service. Initially it was in two sections: technology suppliers and service providers. Once a short list of each type was established we were all encouraged to form consortium for the final round of bidding to secure the complete service (technology and support services). As a result SGI as a technology supplier and CSC and the University of Manchester as service providers formed a consortium. We believe it was an ad is an effective relationship bringing together a number of complementary skills. The responsibilities of each member of the consortium are summarized in Figure 2. SGI (who owned Cray at the time) brought the major elements of the technology together. CSC brought its experience of managing large systems and UoM brought its applications experience, and the experience of dealing with the scientific community (both users and the RCs). Most importantly of course was the ability of the consortium to take financial risk, which the RCs and the UoM are not allowed to take. This structure emphasizes that the UoM is the single point of contact for the users.



**Figure 2 The Consortium Structure and Roles**
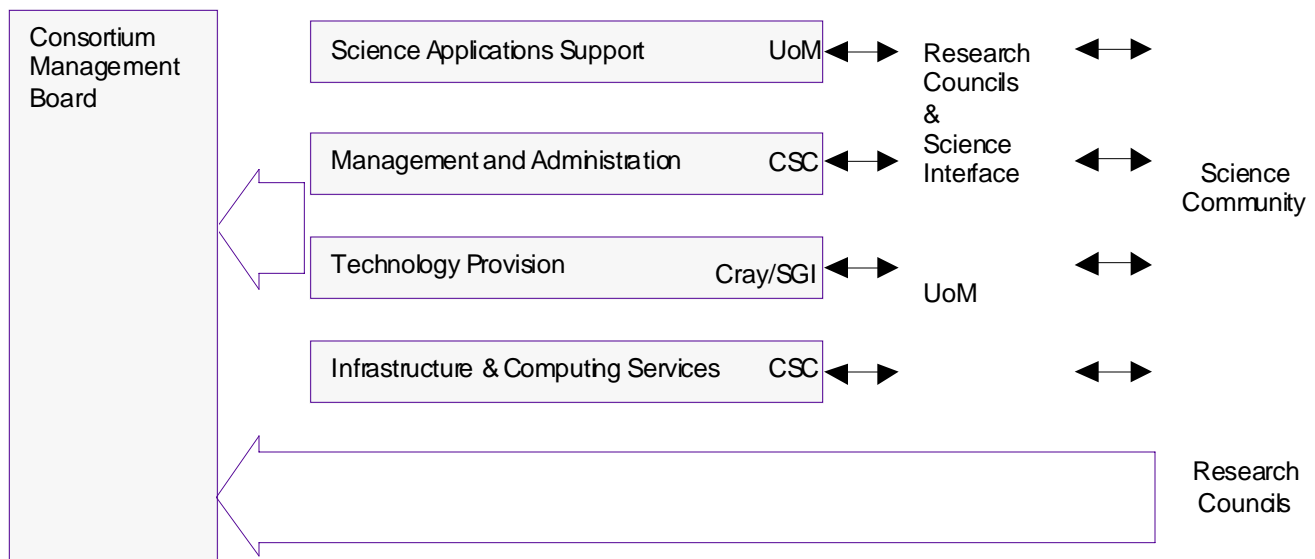
To cut a long story short in July 1998 the consortium was awarded a six year contract from the covering the period 1st January 1999 to 31st December 2004. The emphasis of the contract was that the consortium was supplying "a supercomputing service, which is capable of at least 0.1 TFlop sustained on a particular benchmark". In addition the service is required to supply other

resources, e.g., frontline and application staff. This capability is labeled the baseline service, and the research councils are committed to paying for this level of service for the six years of the contract. The RCs had also estimated that this level of service was insufficient for their projected needs. However the consortium takes the risk, based upon capacity plans on when and how to expand the service. Since the RCs approve the allocations they have control on how much funds they commit over the period of the contract. The Cray T3E-1200E which was installed was in excess of this baseline requirement, as can be seen in Figure 4. In addition all the other resources including memory, tape store, disc, staff resources were greater than that specified in the baseline.

To emphasize the CSAR service has both academic users funded by the RCs and a number of commercial users. The contract also recognizes that supercomputers rarely last six years and we are contracted to provide twice the compute capability before the end of December 2001. The contract indicates this will not be achieved by just enhancing the T3E-1200E. More information about the process and the services can be found in the CSAR newsletters that are available from our World Wide Web Server.

## CSAR Resources and Trading

Figure 1 shows the resources that a user can purchase from CSAR. The columns are filled in with the amount needed for each six-month period over three years. These are just convenient time scales, and it is unusual for grants to be longer than three years. The pen ultimate column contains the total and the final column the equivalent number of tokens. The exchange rates between tokens and resource were established in the initial contract. The resources are

- T3E Hours: this contains the total number of CPU hours required, e.g., 10 runs of a 20 PE job, each for 30 hours would mean 10*20*30 = 6000 hours.
- O2000 Hours: similar to above but for the SGI Origin2000 system. Note it is cheaper (or smaller number of tokens) per hour compared to the T3E-1200E since the T3E-1200E is a more powerful processor.
- Guest Services: The consortium has made other systems available, which are not part of the contract and this enable these boxes are filled in with the cost of these additional services.
- HPD stands for high performance disk, disk that is connected to the T3E.
- MPD: stands for medium performance, disk that is connected to the Origin2000 and is nominally slower access.
- HSM/Tape: amount of tape storage required.
- Training: number of days training required
- Optimization and Application Support for the purposes of this paper can be treated as the same thing. If a user needs some code developing by CSAR he/she needs to ensure they have purchased the necessary resources. They do not need to request tokens for day-to-day support, e.g. 'why hasn't my job run?' 'My code fails in 727, can you tell me what's wrong please?'

Whilst the monetary value appears to be a real price, it is only regarded as a not to exceed price. The resources are costed in the contract on a yearly basis so, for example, CPU hours become 'cheaper' in subsequent years. This version of the calculator uses the same costs for all years. A new version of the calculator will be made available soon that reflects the changes in costs in each year.

### Trading Pool

One of the problems that the new service was asked to overcome was the inflexibility in the allocation given to the user. Hence the trading pool. Once every three months a PI can trade up to 15% of his resources. Thus the total number of tokens is constant, the PI can convert unused disc space to CPU cycles for example. In addition to the total number of tokens being fixed the amount of resource that might be available might also be fixed, i.e., if we only had 1 PE we could only allocate 168 hours a week irrespective of how many tokens were available. Thus a trade is not necessarily assumed to be instantaneous. There may be a very long lead-time to the actual availability of resource. Since the consortium has invested in much more resource than required for the baseline service it is not yet a problem.

In some ways the users are surprised. The additional flexibility this service provides means they have to take more decisions and realize those through the trading pool. Thus they have had to pay a price for the flexibility; a small increase in the administration they must do on the system.

*Capacity Plans*

If the consortium is to take the risk about purchasing additional resource then it needs some estimate of how users will be using their resource over the remainder of their grant period. The initial allocation is no longer adequate as they may have traded resource, or may have decided for scientific reasons to use their resources in different time periods. To that end every PI is asked to keep up to date a capacity or computational plan. In this the PI projects how he or she might consume resources. For information the plan also contains summaries of resources used, so the PI cannot over allocate the tokens remaining. The plan is for planning purposes only, but it has proved extremely valuable in our planning the optimum time to upgrade the Cray T3E-1200E. The individual plans are aggregated into an overall plan and Figure 3 shows the T3E-1200E element. The demand is low in second half of 2001 and 2002 because at this stage few grants expand beyond that date. Those with allocations until mid-2001 are probably only know starting to think about their next proposal.
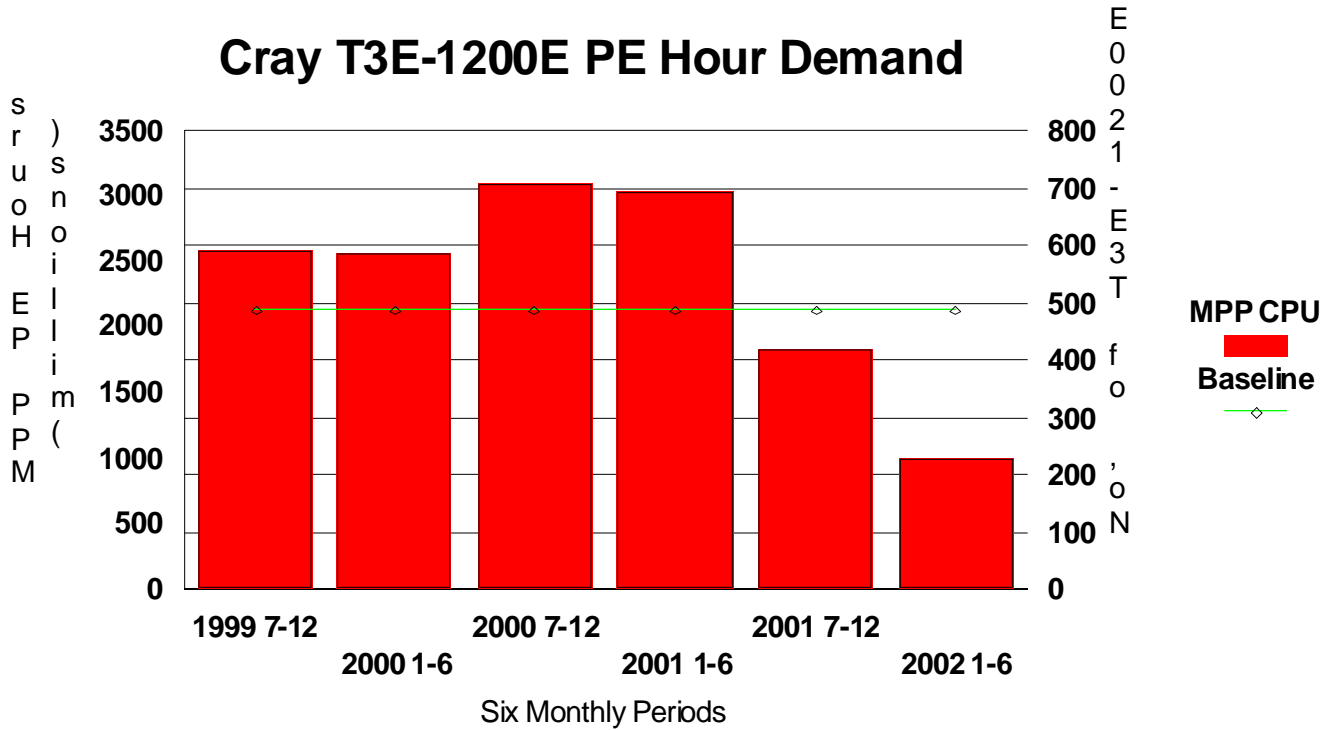


**Figure 3 The Aggregated Capacity Plan**

## Risk Transfer?

The function purpose of the public private finance initiative is to transfer risk from the government agency to the private sector, and I have mentioned briefly in the previous section how this has happened. The RCs have taken a conservative estimate of what they will use of the six-year period and have contracted to pay for this baseline service. However if there is demand from the users, and the proposal is scientifically sound then the RCs will allocate tokens. The consortium takes the risk on what upgrade or downgrade is needed, when it is needed and how to finance. This above baseline usage is on a pay as you go basis.

I have only identified one of the risks. As part of our original bid, and in conjunction with the RCs we created a risk register which identified which risks were being transferred from the RCs wholly or partly to the Consortium. Table 1 is an extract from the risk register; it contains over 40 rows.

The risk assumed by CfS is contractually under-written by the commercial members of CfS (CSC and Silicon Graphics) and the UoM is incentivised by the commercial members to fulfil the responsibilities assigned to it.

| Risk/Responsibilities Currently with Research Councils, from the Risk Register | Risk Description - e.g. the activities required to be undertaken in overall provision of the Service [WHAT] | Risk going to whom? [WHO] # = Lead organization | Mechanism for effecting the risk transfer [HOW] |
|---|---|---|---|
| **Responsibility for planning of service** | Plans for the service will include ordering the equipment, preparing the infrastructure, installing and configuring the HPC Services, setting up the User environment, accounting and job scheduling mechanisms, security and other Service support provisions | **CfS** | The CfS Operational Management Group will implement the plan in accordance with the Service Agreement containing the definition of the Service, the Service Levels and Transition Plan. CfS has strength in project management, which is a standard component of commercial contracts. In Baseline Service |
| **Responsibility for implementing appropriate hardware or service elements in response to the evolving user demand/requirements above the levels mandated by the Research Councils** | This covers predicting demand and responding in a cost-effective manner | **CfS (UoM#, CSC, SGI) With input from the Research Councils** | The Service model incorporates a proven capacity-planning model. It also incorporates schemes that will enable the CfS to monitor trends in an effective manner. CfS will deploy new resources in response to evolving demand. In Baseline Service. |
| **Risk of ensuring balanced hardware infrastructure and service components** | This covers configuration of the HPC Services for the efficient processing of the Research Council's workload | **CfS (SGI#, CSC)** | The proposed technology solutions are scaleable in all dimensions and will grow to meet demands. CfS will ensure that the overall HPC systems remain well balanced in order to enable high utilization levels with efficient usage. The Service Trading Pool will give the Users additional flexibility to adjust their resource allocation. In Baseline Service. |
| **Responsibility for technology refresh** | This covers changes to the HPC Services during the contract Term in order to provide newer technology for the Users, to give them competitive advantage for the throughput of science | **SGI#, CSC** | The proposed service includes technology refreshes in the 3rd year. The plan also incorporates a provision for extra resources to facilitate migration the newer technology. It should be noted that most of the costs for the Service would be in the Utilization band. |

**Table 1 Extract from the Risk Register**


## How do We Get Paid?

*Pricing Model*

The pricing model has a simple structure with two components:

- There is a Baseline Fee that decreases annually, though the Baseline Service increases each year. Even during the first year, the Baseline Fee is set at a level below the price corresponding to the minimum service specified in the contract. The Baseline Fee entitles the Research Council users to utilise a specified level of computing resources and all the services required to support such use. The computing capacity included in the baseline is set at a constant level throughout the lifetime of the contract. It essentially specifies the minimum amounts of resource in Figure 1 that must be available.

- All usage above the baseline is charged based on the actual utilization.

The contract has agreed rates for all the resources identified in Figure 1 both within the baseline and for above baseline usage. However the model assumes that CPU, Tape and disk become cheaper over the six-year period, and that staff costs increase over the same period. We wish to remain cost effective so the pricing model includes reduction/rises in costs over the contract term. However the user does not see the complexity of the finance model, he or she only uses tokens and the exchange rates between tokens and resources are constant.
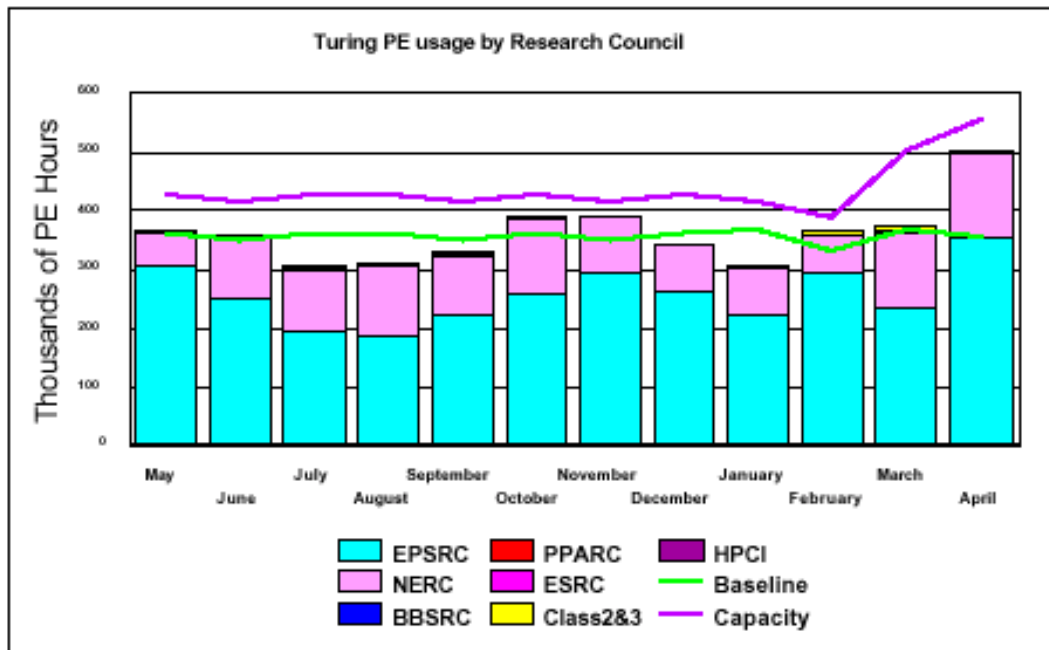


**Figure 4 Usage on the T3E (turing)**

Figure 4 may help understand it. The bars indicate the amount of resource consumed in that month, T3E time in this case. The green line represents the baseline capacity. It goes up and down according to the number of days in a month. The purple line represents the capacity of the system. So for example for July 1999 we will have received the baseline fee only, whereas in April 2000 we will receive the baseline fee plus a fee for being approximately 45% above baseline. The jump in the capacity corresponds to the additional 200PEs that were installed in March 2000. The decision about increasing the capacity was taken wholly by the consortium on the basis of usage and capacity plans. If we do not sustain usage close to this capacity the consortium may well elect to remove them. You may know like to compare actual usage in Figure 4 to a similar period in capacity plan of Figure 3. Users have not be able to match their capacity plans!

Earlier I wrote that the resource calculator produces only an approximate price. The actual price paid by the RCs depends upon:
- As CPU usage gets cheaper in each time period, if a user consumes his resources in different time periods than in the proposal then a higher or lower price will be paid, depending on whether the resources are used earlier or later than specified in proposal.
- As above and below Baseline are charged at different rates it is not clear that a particular job is charged at above or below the baseline rate.

### Performance Criteria
If it were that simple! The quality of the CSAR service is also measured and the performance above or below target, on a monthly basis is converted via agreed exchange rates into tokens. At the end of each year that is converted into money for the consortium or additional resources to the RCs in subsequent years depending upon which side is in credit. Tables 1 and 2 taken from the CSAR quarterly report to the managing agents shows how this assessment takes place. It gives the measure by which the quality

of the CSAR Service is judged. It identifies the metrics and performance targets, with colour coding so that different levels of achievement against targets can be readily identified. Unsatisfactory actual performance will trigger corrective action.

**CSAR Service - Service Quality Report - Performance Targets**

| Service Quality Measure | Performance Targets | | | | | |
|---|---|---|---|---|---|---|
| | White | Blue | Green | Yellow | Orange | Red |
| HPC Services Availability | | | | | | |
| Availability in Core Time (% of time) | > 99.9% | > 99.5% | > 99.2% | > 98.5% | > 95% | 95% or less |
| Availability out of Core Time (% of time) | > 99.8% | > 99.5% | > 99.2% | > 98.5% | > 95% | 95% or less |
| Number of Failures in month | 0 | 1 | 2 to 3 | 4 | 5 | > 5 |
| Mean Time between failures in 52 week rolling period (hours) | >750 | >500 | >300 | >200 | >150 | otherwise |
| Help Desk | | | | | | |
| Non In-depth Queries - Maximum Time to resolve 50% of all queries (working days) | < 1/4 | < 1/2 | < 1 | < 2 | < 4 | 4 or more |
| Non In-depth Queries - Maximum Time to resolve 95% of all queries (working days) | < 1/2 | < 1 | < 2 | < 3 | < 5 | 5 or more |
| Administrative Queries - Maximum Time to resolve 95% of all queries (working days) | < 1/2 | < 1 | < 2 | < 3 | < 5 | 5 or more |
| Help Desk Telephone - % of calls answered within 2 minutes | >98% | > 95% | > 90% | > 85% | > 80% | 80% or less |
| Others | | | | | | |
| Normal Media Exchange Requests - average response time in month (working days) | < 1/2 | < 1 | < 2 | < 3 | < 5 | 5 or more |
| New User Registration Time (working days) | < 1/2 | < 1 | < 2 | < 3 | < 4 | otherwise |
| Management Report Delivery Times (working days) | < 1 | < 5 | < 10 | < 12 | < 15 | otherwise |
| System Maintenance - no. of scheduled sessions taken per system in the month | 0 | 1 | 2 | 3 | 4 | otherwise |

**Table 2 Service Quality Report – Performance Targets**

So for example, in any month we measure the HPC service availability. According to what we achieve we colour code an entry in Table 3 appropriately. Green corresponds to the Baseline or adequate service. White and blue are above average performance and yellow brown and red correspond to below average. This chart is also a useful management tool, making it is straightforward to identify problem areas or trends.



**CSAR Service - Service Quality Report - Actual Performance Achievement**

| Service Quality Measure | 1999 | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Jan | Feb | March | April | May | June | July | Aug. | Sept | Oct | Nov | Dec |
| HPC Services Availability | | | | | | | | | | | | |
| Availability in Core Time (% of time) | 99.70% | 100% | 100% | 97.10% | 98.50% | 99.70% | 99.70% | 100% | 100% | 100% | 100% | 100% |
| Availability out of Core Time (% of time) | 100% | 99.40% | 98.51% | 99.10% | 99.71% | 99.40% | 99.40% | 99.40% | 99.5% | 100% | 100% | 99.70% |
| Number of Failures in month | 1 | 3 | 1 | 1 | 3 | 2 | 2 | 1 | 1 | 0 | 0 | 1 |
| Mean Time between failures in 52 week rolling period (hours) | 744 | 354 | 432 | 480 | 463 | 396 | 391 | 418 | 437 | 498 | 534 | 563 |
| Fujitsu Service Availability | | | | | | | | | | | | |
| Availability in Core Time (% of time) | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A | 99.30% | 100% |
| Availability out of Core Time (% of time) | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A | 100% | 100% |
| Help Desk | | | | | | | | | | | | |
| Non In-depth Queries - Max Time to resolve 50% of all queries | <0.25 | <0.25 | <0.25 | <0.25 | <0.25 | <0.25 | <0.25 | <0.25 | <0.25 | <0.25 | <0.25 | <0.25 |
| Non In-depth Queries - Max Time to resolve 95% of all queries | <1 | <2 | <2 | <1 | <3 | <3 | <2 | <2 | <1 | <3 | <2 | <1 |
| Administrative Queries - Max Time to resolve 95% of all queries | <1 | <5 | <2 | <2 | <2 | <1 | <1 | <1 | <1 | <2 | <1 | <0.5 |
| Help Desk Telephone - % of calls answered within 2 minutes | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% |
| Others | | | | | | | | | | | | |
| Normal Media Exchange Requests - average response time | <0.5 | 0 | <0.5 | <0.5 | <0.5 | <0.5 | 0 | 0 | 0 | 0 | 0 | 0 |
| New User Registration Time (working days) | <2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Management Report Delivery Times (working days) | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 | 10 |
| System Maintenance - no. of sessions taken per system in the month | 2 | 2 | 2 | 0 | 1 | 2 | 2 | 2 | 1 | 2 | 2 | 2 |

**Table 3 Service Quality Report – Actual Performance**

The monthly reports also contain another table giving the Service Credit values for each month to date. In fact the report is full of tables and figures showing how each resource is being used. All the monthly service and quarterly management reports are available from the CSAR World Wide Web Server.

Some of the targets must be met on a month-by-month basis, but others we have been more inventive. New Users must be registered within 3 working days for Green. We have in place the web based registration software whereby the new user is registered "instantaneously" so through development of the appropriate infrastructure we achieve the white band every month. (Due to a misunderstanding it was only coded green for the first nine months.)

## Is it working?

Yes and no! In terms of the hardware resources we have invested in additional equipment a Fujitsu VPP300 to satisfy a vector-processing requirement of one of the RCs. We have increased the capability of the T3E-1200E through the addition of 200 PEs in March 2000.

However we have not collected sufficient tokens to equate to baseline usage in applications and optimization and training. We only need to collect 60% of the tokens corresponding to scientific support; 40% of the baseline service is assumed to be for day-to-day support of the users. The reasons for the lack of uptake are complex:

- Users are not yet used to including such resource requests in grant applications
- If CSAR improves the code efficiency it shows the code was in efficient in the first place!
- There are other places where such resources can be obtained without appearing as a direct cost to the grant
- If they report a problem they think will get help fixing it anyway.

We have had an initiative to give each project 3 days free optimization support and all projects that took up the support have been able to see benefit. It is interesting that those projects that have purchased such resource have been satisfied; in some cases we have be able to improve performance of their codes by a factor of ten.

## Conclusions

The new regime is working. Users get a more flexible service, the Research Councils have transferred risk to the consortium, and therefore capital spend to recurrent spend. The system has grown in response to user demand that would have been extremely difficult to realize without the private finance initiative mechanism.

## Acknowledgements