

NHSBT
UK National Health Service – Blood & Transplant
October 2008

NHS Blood & Transplant (NHSBT) is part of the UK's National Health Service and is responsible for blood services and organ transplantation throughout the UK. The organization manages the supply of blood to hospitals throughout England and North Wales, tracking blood from when it is first donated, through its testing and separation into various products, and, finally, to its dispatch to hospitals.

NHSBT has a particularly important role to play following any major incident where people might be injured and where blood is urgently needed to save lives. Therefore, the efficient operation of NHSBT's computing infrastructure is of paramount importance in ensuring that patients receive the blood they require.

To ensure the effective management and processing of data under all circumstances, NHSBT has two physically separate data centers using split-site OpenVMS clusters that provide multiple levels of redundancy. NHSBT has recently delivered a major upgrade to this system – an upgrade that was accomplished with minimal disruption to service. This article describes the development and delivery of the new NHSBT system.

NHSBT Facts & Figures

NHS Blood & Transplant's services include:

- promoting blood, tissue, and organ donation to the public.
- managing the supply of blood to hospitals in England and North Wales.
- working with hospital colleagues to promote the safe and appropriate use of blood.
- providing a range of tissues to hospitals.
- managing organ donation in the UK.
- managing the British Bone Marrow Register.

There are roughly 2.5 million blood donations per year. NHS Blood & Transplant (www.blood.co.uk) tracks and manages these blood donors and the blood they donate. NHSBT is the sole provider of blood and the majority of blood components to the National Health Service in England and North Wales and also to private healthcare providers.

Every blood sample is individually screened, identified, and tracked. A blood sample is usually split into one or more products, including red cells, plasma, and platelets. These products are stored in thirteen centers around the country.

The National Data-Centre Complex

The Birth of PULSE

The data-processing application that manages all of this activity is called PULSE. PULSE manages the core operations of NHSBT, including:

- administration of donors and donation records.
- testing of donations.
- manufacture of blood components into blood products.
- issue of blood products.

The origins of PULSE go back to the early 1980s when the blood supply at the various regional transfusion centers was managed by different and incompatible systems. There was no easy way to get an overall view of England's blood supply or to move blood from one center to another.

One of these systems, which ran at several centers, was called MITHRAS. MITHRAS ran on a DEC VAX system using the Mimer SQL database. When the National Blood Authority was created from the regional transfusion centers in 1994, it set out to select a single system from competing suppliers. In the 1997 to 1998 timeframe, it selected MITHRAS and renamed it PULSE.

Because of its MITHRAS legacy, the initial version of PULSE ran on HP (originally DEC, then Compaq) AlphaServers under OpenVMS with SCSI HSZ disk arrays and FDDI cluster interconnects. A second generation of PULSE upgraded to ES45 AlphaServers with a localized storage area network (SAN) using HSG80 disk arrays and Gigabit Ethernet cluster interconnects.

PULSE was installed in three independent data centers throughout England. They each served the blood centers in their region - Bristol in the southwest, London in the southeast, and Leeds in the north. Each of the data centers contained information only on the donors in its region.

The three independent databases were not fully integrated and reflected the social and organizational structures of the time, which created some operational problems:

- The transfer of products between data centers was clumsy and inefficient.
- Donors were not recognized in a region not serviced by their data centers.
- If an NHSBT center was compromised, blood collection, production, and blood product transfers were complex if they had to take place across regional data-center boundaries.
- Auditing of blood-product transfers and donor relocations was difficult.
- Duplicate reference tables for common data had to be maintained at each data center for each database. Change-control procedures were difficult to manage.
- The requirement to perform database joins across three disparate databases to provide management information added considerable complexity.
- National queries were available only for certain functions as each query had to access the three regional databases separately and join the returned results.

It was clear that the three databases should be consolidated into a single national database.

Structuring for High Availability and Disaster Tolerance

The availability requirements of the PULSE system are not as difficult to achieve as those of a real-time mission-critical system such as an air traffic control system. Though the PULSE system is both safety-critical and mission-critical, millisecond response times are not required. In fact, if the system is unavailable for a few minutes, this is generally acceptable. However, the system cannot suffer a loss of service for an extended period of time without extremely careful planning. In the event of a major incident where many people are injured, blood products must be immediately available and have to be moved quickly from their places of storage to where they are needed.

Consequently, the system must be designed to survive multiple failures. The system can be down for short periods (its Service Level Agreement is for three 9s or about eight hours of downtime per year), but it must be extremely resilient to failures of any kind, including a total site outage.

Structuring for Improved Efficiency

To accomplish the required database consolidation while meeting its availability requirements, PULSE was moved in 2002 to a new National Data Centre Complex (NDCC) comprising duplex data centers at two physically separate locations.

The three separate databases were collocated in the NDCC and ran on OpenVMS AlphaServer clusters. The PULSE service load was balanced between the two NDCC data centers.

The databases were still not consolidated as regional partitioning still applied. The databases for the three distinct regions were stored in separate databases, each running on its own OpenVMS cluster. However, the system could now provide limited national views of data to a client by using middleware to merge data from the three databases.

In addition, applications were becoming more complex; and there was a desire to move from the old character-cell interfaces to a modern GUI interface.

Therefore, in 2007-2008, the PULSE system went through another major upgrade to create a single national database and to improve the overall level of availability. This upgrade had to meet a series of challenges:

- The single-tiered architecture of the original PULSE system had to be replaced with a more flexible three-tiered architecture providing client, application, and database tiers.
- The design and specification of new hardware, a new version of the OpenVMS operating system, new storage subsystems, a new version of the Mimer SQL database, and updated PULSE applications had to be such that PULSE could run unchanged so far as the users were concerned.
- The migration to the new system had to be done with minimal interruption to PULSE services.
- The old and new systems had to work in parallel until the old system had been decommissioned.
- The three databases had to be merged into one with minimal impact on operations.

- The time to restore services to users in the event of a failure needed to be minimized. This led to the implementation of “fast restart” capability in the Mimer SQL database server.

OpenVMS Clusters

The new PULSE system runs on a system platform built using OpenVMS Integrity split-site, disaster-tolerant clusters.

Before describing the new PULSE architecture, let us briefly review OpenVMS clustering. OpenVMS clusters are the “founding fathers” of today’s cluster technologies. What’s more, even though they were first-to-market by decades, they still have several significant advantages over most, if not all, of today’s cluster technologies.¹

OpenVMS clusters are “shared everything” with a coordinated cluster-wide file system that provides cluster-wide record-level locking. Each node may have the same files open for true cluster-wide, shared read-write access. The distributed lock manager is the key component to making “shared everything” clusters work, and it ensures the synchronization and serialization of operations across the entire cluster.

The terms “cluster-aware” and “cluster-tolerant” are relevant here. Cluster-aware means that the software knows that it is running on a cluster and makes active use of the distributed lock manager for simultaneous access to the database files by all participating nodes. Cluster-tolerant means that the software will run on a cluster but may not use any or all of the cluster-locking facilities. The PULSE system fits into this latter category since the Mimer SQL database manager and the PULSE applications run in an active/hot-standby mode. Failover of the applications and the Mimer SQL database is very fast because of the Mimer SQL database’s “fast restart” capability described later.

The Upgraded PULSE System

The upgraded PULSE system achieves its high availability and its extreme tolerance to disasters through the extensive use of parallelism and equipment redundancy in its architecture, its data-storage systems, and its network infrastructure

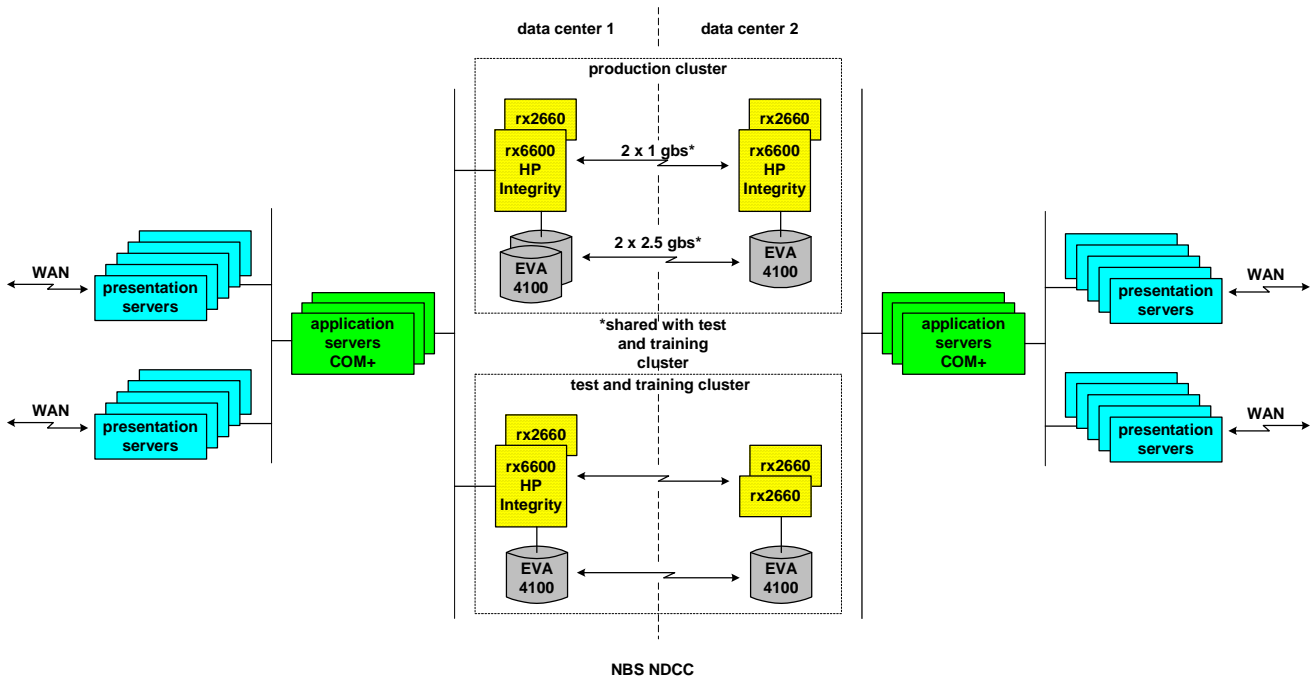
Production System

The upgraded PULSE system comprises a three-tiered architecture. The production configuration is similar in each data center with some differences:

- The Presentation tier (the client tier) comprises Win32 systems running a Citrix Presentation Server (now Citrix XenApp). Currently split between the data centers are over one hundred such servers, each serving thirty to forty thin clients in the field.
- The Application tier currently comprises two COM+ systems in each data center. They run as a load-balanced server farm. This configuration may soon be expanded to four systems per data center. The application servers are not clustered. Rather, they share the application load.
- The Database tier comprises redundant databases running on a four-node, split-site OpenVMS cluster. The subcluster configuration at each data center comprises two HP Itanium server nodes – an HP Integrity rx6600 server (a 7U machine with up to four dual-

¹ OpenVMS Active/Active Split-Site Clusters, *Availability Digest*, June, 2008.

core processors and 192GB memory) and a smaller HP Integrity rx2660 server (a 2U machine with up to two dual-core processors and 32GB memory).



- The database files reside on a number of “shadow sets,” which are OpenVMS HBVS (host-based volume shadowing) devices that synchronously replicate the data across three physically separate HP EVA 4100 storage arrays. HBVS ensures that all members of a shadow set are consistent and are correctly reconstructed with copy/merge operations if a node or a storage device should fail. Two of the EVA 4100 storage arrays are located at one data center, and one of the arrays is located at the other data center.
- The Mimer SQL database server maintains its own large buffer pools in 64-bit address space for performance, with up to 16GB memory being reservable for this space. The Mimer SQL database runs on a single cluster node at any one time, with the other database server nodes in “hot standby” mode and ready for failover if necessary.

Test and Training

Like the production system, NHSBT has implemented a four-node, split-site OpenVMS cluster for testing and training. In this case, the subcluster nodes at one site comprise an HP Integrity rx6600 server and an rx2660 server, just like the production cluster. The subcluster nodes at the other site are a pair of HP Integrity rx2660 servers. Both sites have an additional test and training EVA storage array managed by the test and training cluster.

The test and training nodes are provided to train blood service staff and systems-management staff, to test PULSE upgrades, and to test failover procedures. They also provide the facilities to test new releases of the operating-system software and equipment firmware.

The rx6600 test and training node at one of the sites provides two functions. It matches the power of the production database server nodes and can therefore be used for scalability and performance testing. It also serves as a cold standby backup for a production node should one node suffer an extended outage.

Each site also has an application server to complete the test and training configuration. These nodes also serve as cold-standby backups for the equivalent production nodes.

Management Information

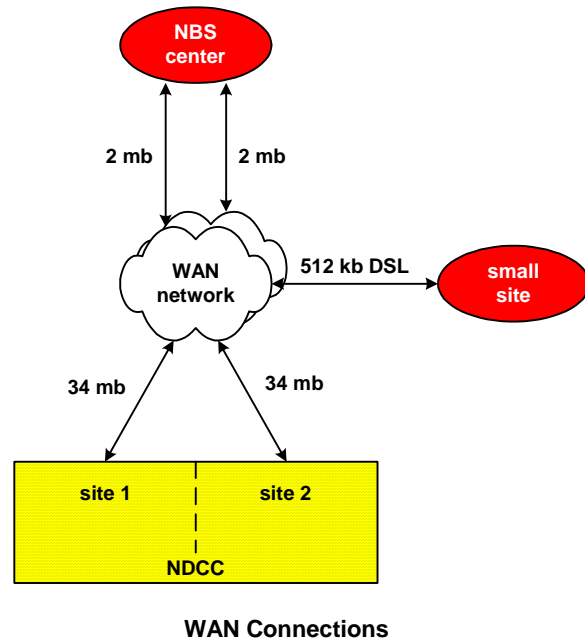
Business Objects Servers at each site operate in a load-balanced configuration to drive an Oracle Operational Database System (ODS) for management information. Web access is provided to this system for management reports.

In addition, two of the EVA 4100 storage arrays (one per site) are configured with additional capacity to support archiving. The Archive Server is an identical rx2660 node that is configured as a single node cluster.

Data and Storage Networks

In order to minimize risk due to contention for network resources, there are private links for cluster communication and for data replication as well as a WAN network to interconnect the users in the field to the presentation servers. All network links are redundant to ensure the uninterrupted provision of services to PULSE users in the event of any network path or switch failures.

Cluster interconnect traffic is sent over two physically separate Gigabit Ethernet circuits carried on the DWDM² intersite links. The private cluster-interconnect data network is implemented using HP ProCurve switches with ProCurve meshing. The cluster nodes and storage arrays are interconnected via dual SAN fabrics, with each cluster and each storage array connected to both fabrics. SAN-fabric traffic is carried over two physically separate pairs of 2 gigabit/second fibre channel circuits (a total of 4 gbs per SAN fabric) transmitted on physically separate DWDM intersite links.



The user data network is Cisco-based. All cluster nodes are connected to both the HP ProCurve private cluster interconnect network and to the Cisco user data network. This completely isolates the two types of traffic with no risk of traffic contention.

Client connections to the fifteen major NHSBT distribution sites are via a redundant WAN network that interfaces with the Cisco user data network. Each remote distribution site is connected to both WANs, and each NDCC data center has connections to both WANs. Each data center and each distribution site connects to the redundant WAN network via separate POPs (points of presence) to mitigate against exchange failures. In the event that both WAN connections should fail, DSL backup is provided between the NDCC and the distribution site.

The WAN speed is 34 megabits/sec. at the data centers and two megabits/sec. at the sites. Traffic is load-balanced across the dual WANs. Should a WAN fail, the other WAN carries all of the traffic.

² Dense wavelength division multiplexing. DWDM is an optical technology that increases the bandwidth capability of a strand of fibre optic

Minor NHSBT sites are connected to the data centers via 512 kilobit/sec. DSL connections.

System Management

The NDCC OpenVMS clusters for the PULSE applications are monitored and managed with HP's DTCS (Disaster Tolerant Computer Services) products. They provide detailed alerting of any issues or deviations from the expected operational state of any of the nodes or storage subsystems.

Cutover From Old To New

Migration of the application clearly had to be done as speedily as possible. However, some outage time was essential for this to happen. The time-consuming aspect was the data migration, where data had to be unloaded from the old systems, moved to the new systems, and then merged into a single database. Because the data was previously held in three almost identical databases, the plan was to migrate in three stages, one for each database.

Stage 1 involved proving that the software functioned properly on the HP Integrity server platform. Data from one of the three regions was moved to the new system, which provided confidence that the hardware and software was correctly configured and which demonstrated compatibility between the old and new systems.

Stage 2 involved unloading the data from the second regional database, copying it to the new platform, and then merging it into the new database. This provided confidence in the database merge process prior to the final regional database being moved across.

Stage 3 completed the data unload/merge process, at the end of which the single national database was operational.

Several trial migrations were completed prior to the final live migration.

The phased cutover also provided confidence in the performance behavior and capacity of the new system platform. The platform provides excellent response-time performance with minimal latency and sufficient capacity to absorb intermittent spikes in workload with little impact on response times.

Disaster Tolerance

Redundant Data Centers

The NDCC's disaster-tolerance capability starts with a pair of fully configured and geographically separated data centers sharing the load. Each data center is equipped with dual uninterruptable power supplies (UPS), automatic fire-suppressant devices, and full environmental monitoring that is linked to automatic notification devices.

Redundant Processors

All processing capacity is redundant and is distributed between the two data centers.

During normal operation, the primary rx6600 production cluster node runs the Mimer SQL production database server. Consistent data replication across the three storage arrays is provided by OpenVMS host-based volume shadowing. Other nodes in the cluster provide ancillary services such as running historical query requests against a copy of the database, tape backup processing, etc.

Should the primary node fail, its workload is swiftly passed to the rx6600 cluster node in the opposite data center. In the event of a serious failure such as the long-term loss of a data center, the test and training rx6600 can be brought into the production cluster. In reality, production can survive on a single node in a single data center. It would take the simultaneous failure of three rx6600 cluster nodes or all of the EVA storage controllers to take down the NDCC, a highly unlikely scenario.

Should an rx2660 node in the production cluster fail, its processing load is taken over by the surviving nodes in the cluster. The loss of an rx2660 node is simply an inconvenience as much of its functionality can be absorbed by the remaining nodes.

Presentation servers and application servers are not clustered. However, the load is balanced across them by utilizing Citrix load-balancing techniques for the presentation servers and by explicit mapping of COM+ application servers to specific "failsafe IP" addresses to spread the load across all available NICs. Should a presentation server fail, the users connected to that server can reconnect to a surviving server and continue receiving services. Should an application server fail, a surviving server can take over its load. The GUI applications are written in such a way that each interaction between the COM+ servers and the database and presentation programs are stateless.

Since all nodes are actively participating in the application, failover to a surviving node following a node failure is virtually instantaneous.

Redundant Storage

The NHSBT database is three-way, host-based volume shadowed (HBVS) and is distributed between the two data centers. Since replication is synchronous, no file system data is lost following a processor or storage-array failure. Shadow-set reconstruction following restoration of equipment to service is invoked by the DTCS products during the node boot process and can also be invoked manually if necessary. Rapid shadow-set rebuild is provided using the OpenVMS HBVS features known as "mini-copy" and "mini-merge."

At the request of NHSBT, Mimer implemented a significant improvement known as "fast restart" to the Mimer SQL database manager. If the database server terminates abnormally (e.g., a system failure), the database can be left in an inconsistent state. Thus, the databases must be scanned to check their integrity during the database server startup process. Previously, users could not access PULSE until this process was completed, which could take on the order of an hour.

The "fast restart" improvement was made to the Mimer SQL database manager so that users can connect to PULSE immediately upon the restart of the database server. Database integrity and consistency checking proceed in parallel while users start to use the database. This reduces application failover time experienced by the users following the failure of the primary node from about one hour to a few seconds. System availability is therefore increased by greatly reducing the time required to restore service.

Redundant Networks

Every fibre channel and Gigabit Ethernet connection is redundant.

Each rx6600 has four 4-gbs fibre channel paths (i.e., a total of 16 gbs of bandwidth) to each local EVA 4100 storage array via two entirely separate FC fabrics. The intersite fibre channel links take diverse paths so that any intersite link disruption is confined to path switching between the fibre channel switches and does not impact the production system. HBVS balances the read load to the fastest responding disks, which results in the local EVA storage arrays being used for the majority of read requests.

The ProCurve private cluster connections are “dual rail” for cluster traffic (SCS protocol³) and use “LAN failover” (equivalent to NIC teaming on Proliant Servers) for DTCS monitoring, DECnet interconnects, and EVA scripting functions. This is implemented using ProCurve “meshing” and VLANs for the different types of traffic, where meshing provides a shortest-path connection between ProCurve switches. Any intersite link disruption is confined to path selection between the ProCurve switches and does not impact the production system.

This design has already proven its worth when one of the intersite DWDM pipes was severed by road construction work shortly after the national system was fully operational. There was no disruption to normal operation, and the DTCS monitoring subsystem alerted the operations staff to the problem.

The Cisco user data network IP connections use “Failsafe IP” that allows the active IP addresses to be automatically moved to other NICs in the same failsafe IP group in the event of a NIC, cable or switch failure.

Every WAN connection is redundant or has an alternate backup. There are two separately routed connections to the WAN from each data center. Should both fail, a DSL connection is established to continue communications.

Each major NHSBT center is served by two independent WAN connections. Small sites are served by DSL links that are inherently redundant in the carriers’ networks.

Summary

NHS Blood & Transplant provides a life-saving service. It cannot fail, especially in the face of major incidents. NHSBT has achieved an extremely high level of disaster tolerance through the use of dual data centers and split-site OpenVMS clusters.

Every component within the National Data Center Complex is redundant, from data centers to processors, storage arrays, and networks. There are no standby backup components. Every component is active and is used in a load-sharing configuration. Therefore, should there be a component failure, its load can be shifted to a surviving component very quickly. In addition, NHSBT stocks on-site spares for high probability of failure components such as disks and power supplies.

Though NHSBT’s contractual requirement is to achieve an availability of three 9s, the NDCC configuration should far exceed this requirement for unplanned downtime. It is only planned downtime that will have an impact on the availability of blood services to those serviced by NHSBT.

In summary, this system demonstrates considerable in-depth strength to deliver extremely high-availability blood-product services to NHSBT by using the PULSE software, the Mimer SQL database, and the OpenVMS clusters running on Integrity Server systems and EVA storage subsystems.

³ OpenVMS cluster System Communication Services.

Acknowledgements

We at the Availability Digest would like to thank the following firms for their efforts in implementing this system and in the preparation of this paper:

- NHS Blood & Transplant (www.blood.co.uk)
- Savant (application provider and operational support) www.savant.co.uk
- Mimer (database supplier) www.mimer.com
- OCSL (hardware reseller and installer) www.ocsl.co.uk
- HP C&I (project management and delivery)
- HP DTCS (system and infrastructure monitoring)
- XDelta (platform design and implementation on behalf of HP C&I) www.xdelta.co.uk