

WHITE PAPER

THE HIGHEST AVAILABILITY FEATURES FOR PRIMEQUEST

Business continuity and cost-efficiency have become essential demands on IT platforms. Offering the best-blend of standard and high availability technologies, PRIMEQUEST is a rare enterprise system platform that fully maximizes uptime and greatly lowers TCO.

This whitepaper explains how business continuity can be maximized by the thorough data protection and fail-safe mechanisms in every PRIMEQUEST.

CONTENTS	
Introduction	2
Xeon processor RAS	4
Thorough memory protection	5
Automatic recovery from system board failure	8
Flexible I/O for swift I/O reconfiguration	9
Component redundancy maximizes business continuity	11
Intuitive management view	12
Cable-less design	13
Conclusion	15



INTRODUCTION

PRIMEQUEST is a mission-critical system server that supports up to eight Intel® Xeon® processor E7 family and 80 cores. It maximizes operational uptime by intimately combining Intel® processor technology, Microsoft® Windows Server® technology, Red Hat® Enterprise Linux® technology, and other standard technologies, with Fujitsu's high-availability technology. High levels of component redundancy in PRIMEQUEST include the very heart of the server - processors, memory, and system interconnects. In addition, should a problem occur in any system component, that component is automatically isolated, stopping the problem immediately and allowing uninterrupted system operation to continue with the remaining system resources.

There are many high-availability technologies in each PRIMEQUEST. Even faulty system boards can be automatically replaced with a reserved system board. This function significantly minimizes system down time.

System maintenance is also greatly simplified due to the use of a Mid-Plane which enables cable-less design. This frees engineers from the complexity of cable management as everything simply plugs into the solid-state Mid-Plane.

THE HIGHEST AVAILABILITY STANDARD SERVER - PRIMEQUEST

Both system operation uptime and efficient server use are concurrently maximized with PRIMEQUEST.

Redundant, hot-replaceable components and automatic recovery mechanisms are the core technologies for highest business continuity. In addition, hardware partitioning minimizes the effects of errors, and faulty components are automatically switched to standby components so that operation is resumed in the shortest possible time.

Further, PRIMEQUEST helps you reduce server related purchase, management, and maintenance, by its use of standard technologies and simple intuitive system administration interfaces. The GUI-based management view further facilitates early detection of problems.

In addition, simplified maintenance not only reduces required system downtime for component replacement, but also helps prevent problems otherwise caused by procedural complexity.

HARDWARE PARTITIONING

Even if a hardware failure occurs in a partition, that failure is totally isolated within that partition. It can never affect other partitions. PRIMEQUEST is able to support up to four of these highly robust hardware partitions.

AUTOMATIC RECOVERY FROM A SYSTEM BOARD FAILURE

Any faulty system board is automatically switched to a Reserved System Board. The automatic recovery feature of PRIMEQUEST means you will never suffer from prolonged system down time due to delays in problem detection and response.

With other vendors' industry standard servers, it takes longer to discover and manually resolve such problems. To automate recovery mechanisms and minimize system down time, they resort to the last resort of system clustering. But system clustering introduces greater complexity as well as costly integration and maintenance.

SIMPLE AND SWIFT I/O RECONFIGURATION MECHANISM

After a system board failure and the switch to the Reserved System Board, the I/O devices connected to the failed System Board can be simply redirected to the new System Board. This technology is called "Flexible I/O," and is specifically designed to achieve flexible connections between system boards and I/O devices.

HARDWARE REDUNDANCY FOR MAXIMIZING BUSINESS CONTINUITY

Most components are redundant and system interconnect paths are multiplexed. Even if a component becomes faulty, the system can continue its operation. Moreover, even if an error occurs on a system interconnect it can still work in degraded mode with the faulty link disconnected and available for maintenance.

INTUITIVE SYSTEM MANAGEMENT FOR SWIFT PROBLEM DETECTION

Swift problem resolution is a must-have for stable system operation.

The combination of PRIMEQUEST and ServerView Operations Manager lets you swiftly identify fault locations.

FOOL-PROOF MAINTENANCE BY CABLE-LESS DESIGN

To minimize down time and ensure stable operation, maintenance operations must be error free.

Component replacement with PRIMEQUEST is simple and swift. Server components have rigid plug-in form factors and no external cable connections. In the event of an error, the specific faulty unit is easily removed from either the front or rear of the cabinet.

XEON PROCESSOR RAS

CHARACTERISTICS COMMON TO THE PROCESSORS AND MEMORY

Xeon processor 7500 series and E7 family and the memory at the heart of PRIMEQUEST have thorough error recovery features. These correctly and quickly handle both recoverable and unrecoverable errors.

PROCESSOR RAS

■ Error handling outline

■ Recoverable errors

Energetic particles or electromagnetic waves can cause bit inversion in cache memory. The Xeon processor is able to detect and correct such corrupted bits.

■ Unrecoverable errors

The heart of the processor including the various cache memory levels can protect themselves from fatal errors. For instance, if the number of repeated errors exceeds a certain limit, the problem component can be off-lined.

■ Cache protection mechanisms

■ Handling of recoverable errors

Both data and tag fields in cache levels 1/2/3 can detect and correct bit errors. The data protection feature of level 3 cache is described below.

• Data array

Up to three bit errors can be detected, and up to two bit errors corrected. When a three-bit error is detected, the memory access operation that encountered the error is retried.

• Tag array, core valid array, and LRU (Least Recently Used)

Errors of up to two bits can be detected, and one-bit errors corrected. When a two-bit error is detected, the memory access operation that encountered the error is retried.

■ Handling of unrecoverable errors

If the above retry operations are successful, the application and operating system are not notified of the error. If recovery is unsuccessful, the application is stopped.

■ Other data protection features

One-bit errors are handled by Processor circuits: Registers, ALUs (Arithmetical and Logical Units), and TLBs (Translation-Lookaside Buffer). They can detect and correct such errors.

THOROUGH MEMORY PROTECTION

In general, memory is one of the most error-prone components. Memory errors are corrected, or isolated according to their error type such as recoverable or unrecoverable errors.

■ Handling of recoverable errors

Electromagnetic waves, energetic particles, and other external influences can cause signal bit inversion. A code to detect and correct this type of error is added to each signal. This means errors of this type can be corrected automatically or via a retry.

■ Handling of unrecoverable errors

Even if unrecoverable errors exist in memory, CPUs can correctly access memory using a memory redundancy mechanism. The failed memory is then isolated from further connection to CPUs.

REASONS FOR ELABORATE MEMORY PROTECTION MECHANISMS

To ensure stable memory accesses, there must be reliable methods for handling both recoverable and unrecoverable errors.

One paper on memory error experiments ^{*1} points out three important trends in memory errors.

- (1) The number of recoverable errors increases rapidly during the initial two-year period of memory use.
This means action must be taken to handle recoverable errors.
- (2) Positive correlation was observed between the unrecoverable error rate and the CPU and memory utilization rate.
This means systems with higher utilization rates require more robust memory RAS functions.
- (3) A positive correlation was also observed between recoverable errors and unrecoverable errors.
This means it is also necessary to provide robust memory RAS functions that can not only handle recoverable errors but also unrecoverable errors.

PRIMEQUEST can respond to the various memory errors by detecting and correcting recoverable errors, and isolating those components with unrecoverable errors.

*1 "DRAM Errors in the Wild: A Large-Scale Field Study" (research paper)

<http://www.cs.toronto.edu/~bianca/papers/sigmetrics09.pdf>

GUARANTEED READ AND WRITE OPERATIONS

PRIMEQUEST detects and corrects one-bit errors, detects two-bit errors and performs retry operations using ECC (Error Check and Correct).

If an error occurs on one SMI (Scalable Memory Interconnect) lane, which is an interface between processor and memory, memory access can continue using a spare lane.

RECOVERY FROM MULTI-BIT ERRORS

Even with an error occurring in a single DRAM module, an application can continue operating while the error is corrected. DRAM is assigned to each bit in 4-bit or 8-bit data chunks. ECC corrects this information so that CPU memory access can continue. This feature is called SDDC (Single Data Device Correction). Even if one DRAM module encounters a recoverable or unrecoverable error, SDDC can perform recovery from the error by correcting 1-bit errors as they occur (See lower part of FIGURE 1).

On the other hand, if SDDC is not used, one DRAM failure causes an unrecoverable 4-bit error (See upper part of FIGURE 1).

PRIMEQUEST 1800E2 with Xeon processor E7 family is able to recover dual DRAM failures using DDDC (Dual Data Device Correction) (Figure 2). PRIMEQUEST 1800E with Xeon 7500 processor series is able to recover single DRAM failures using SDDC.

FIGURE 1 DATA PROTECTION BY SDDC

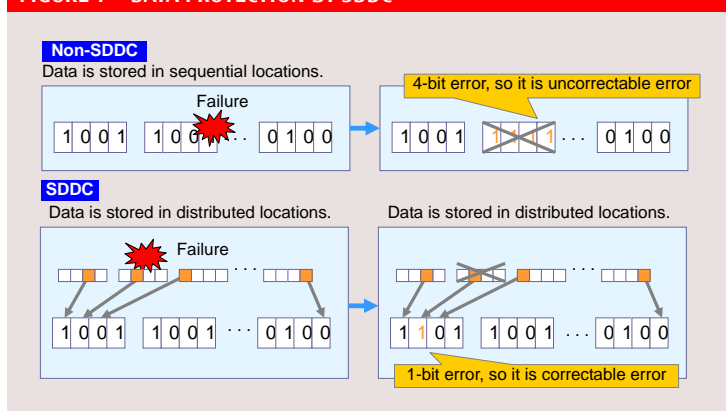
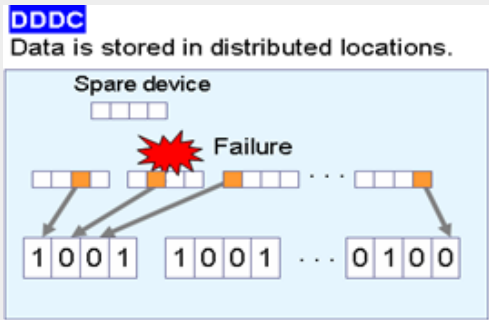
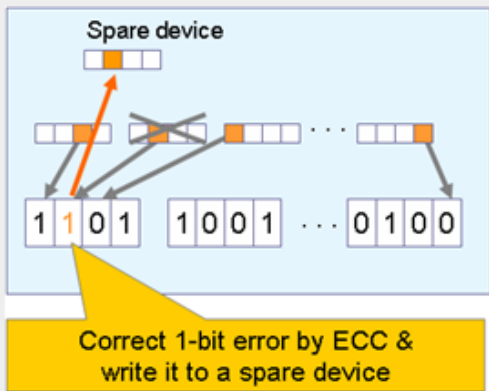


FIGURE 2. MEMORY ERROR CORRECTION BY DDDC

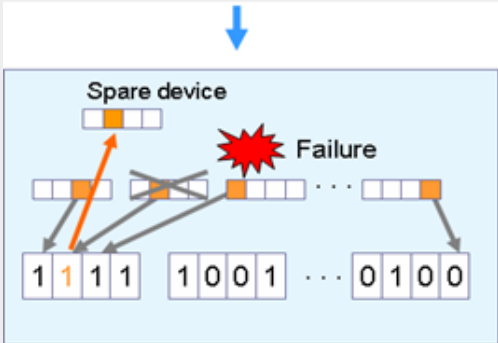
- (1) First DRAM failure
Due to one DRAM failure, one-bit is corrupted.



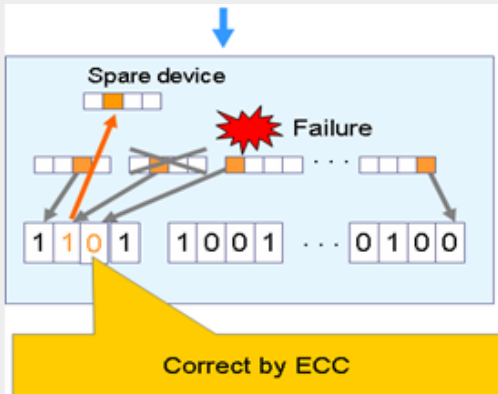
- (2) Recovery of first DRAM failure
The corrupted bit is corrected by ECC.



- (3) Second DRAM failure
Due to another DRAM failure, one-bit is corrupted.



- (4) Recovery of second DRAM failure
The corrupted bit is corrected by ECC.



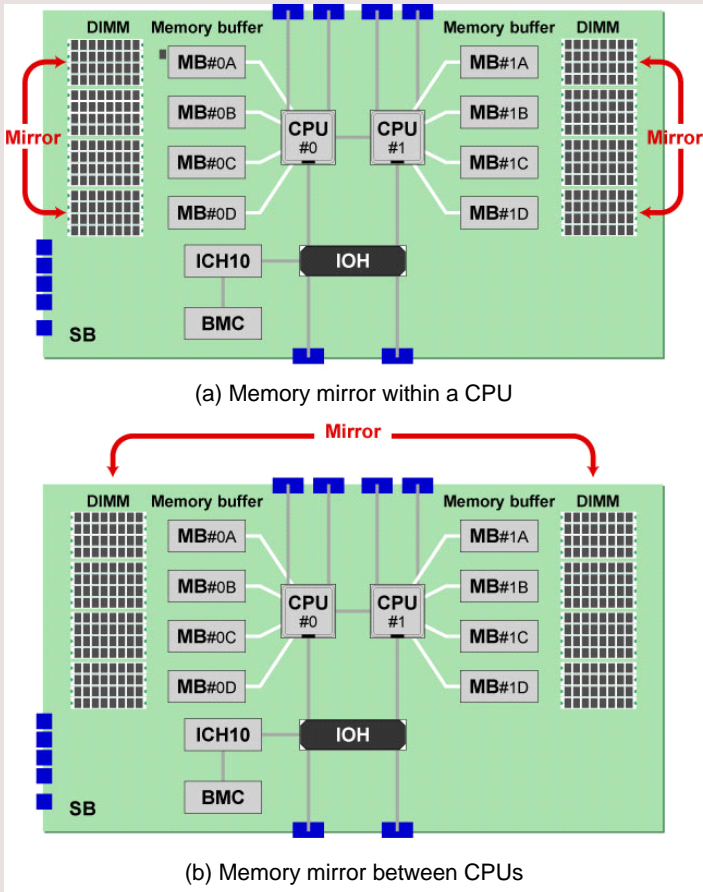
MEMORY MIRROR

Memory mirror is a memory redundancy function that allows each CPU to write and read to/from a memory pair. This means CPU-memory access can continue even if one DRAM module fails, as the other available DRAM module contains the correct data.

Memory mirror can be used to manage both recoverable and unrecoverable errors.

Memory mirror with PRIMEQUEST can be configured so that a memory pair is directly connected to the same CPU (FIGURE 3 (a) Memory mirror within a CPU), or also where a memory pair is connected to different CPUs. In the latter case, memory access can continue even if one Memory Buffer (MB) fails (FIGURE 3 (b) Memory mirror between CPUs).

FIGURE3 MEMORY MIRROR

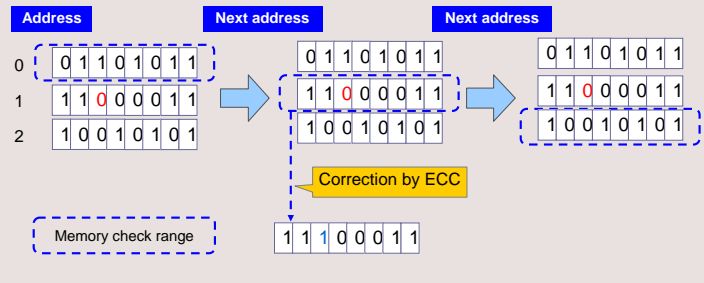


MEMORY SCRUBBING

Memory Scrubbing is designed for early detection and correction of memory errors by using ECC. This includes Demand Scrubbing error checking at memory read time, and periodic error checking by Patrol Scrubbing.

Patrol Scrubbing scans the entire memory while incrementing memory addresses (FIGURE 4).

FIGURE 4 PATROL MEMORY SCRUBBING



AUTOMATIC RECOVERY FROM SYSTEM BOARD FAILURE

At the heart of each PRIMEQUEST are the System Boards containing processors, memory, and system interconnects.

System downtime is minimized as any errors in a system board are automatically detected and followed up by automatic switching of the faulty system board to a Reserved System Board.

Typically, in other vendors' industry standard servers, error recovery mechanisms require complex and time-consuming steps. On a failure you are required to stop system operation, identify the cause of the failure, and perform manual recovery operation. Such manual recovery causes delays in system restart. In addition, such manual operation also carries the risk of human error.

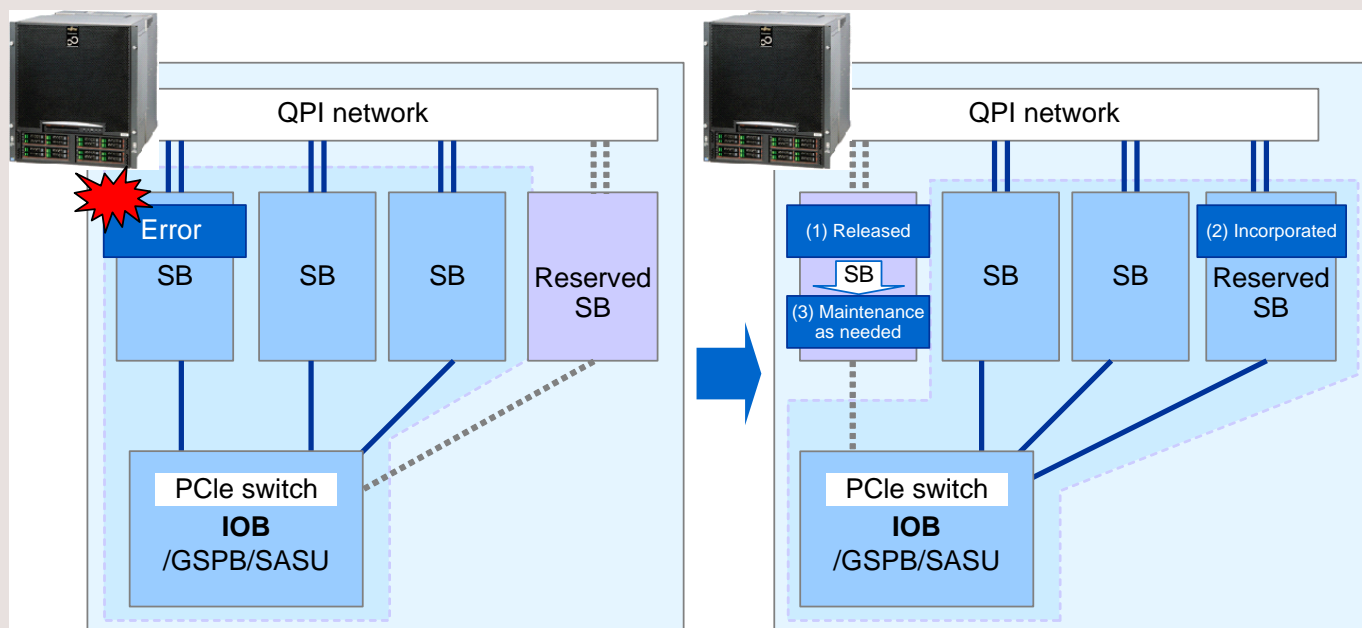
In contrast, the automatic recovery of PRIMEQUEST better ensures system recovery in a predictive time, due to its ability to detect errors and automatically switch System Boards.

PRIMEQUEST switches System Boards (SB) when triggered by the following error detection states.

- SB degradation
- CPU degradation (even if only one CPU is faulty)
- DIMM degradation (even if only one DIMM is faulty)
- Detection of a Memory mirror error
- Detection of QPI lane degradation
- Detection of a change of an SMI lane
- Detection of PCI Express lane or speed degradation (between a System Board and I/O Board)

The Reserved System Board is available for use during normal operation helping to improve system utilization. For example, you can use it for development work.

FIGURE 5 STANDBY SYSTEM BOARD ASSIGNMENT



FLEXIBLE I/O FOR SWIFT I/O RECONFIGURATION

On System Board failure, recovery procedures must include I/O reconfiguration. On PRIMEQUEST, this is greatly simplified by the "Flexible I/O" feature. This allows the Reserved System Board to take over the I/O configuration 'as is' from the failed System Board.

With some vendors' server products, I/O assignments to individual system boards cannot be changed. This means that the system must be clustered to enable swift operation resumption. The I/O configuration of the standby node must also have the same network configuration as the active node. You will find such duplicated configurations complex and cost-consuming.

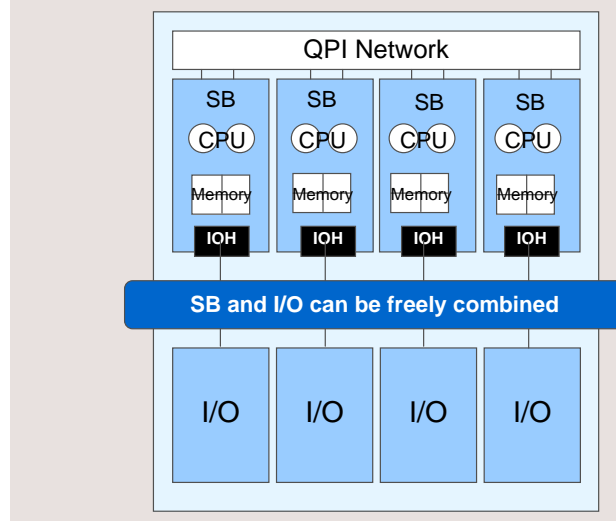
Using Flexible I/O, I/O reconfiguration takes a matter of seconds. Using the intuitive setup view, System Boards can be linked to I/O units such as I/O Boards^{*1} and Giga LAN SAS and PCI Box interface board (GSPB)^{*2}. You have no limitations on the combination of System Boards and I/O units.

- I/O devices can be efficiently used
You can add I/O Boards and GSPB to System Boards in response to I/O workload growth.

- Swift error recovery on System Board failure
As explained in the previous section, "Automatic recovery from system board failure", the I/O Boards and GSPB of the faulty system board in a partition can be switched to the Reserved System Board.

In PRIMEQUEST, each of four I/O Boards forms an independent PCI bus tree, connecting with a System Board, as shown in FIGURE 6. (For details of System Boards and I/O Boards and GSPB connections, refer to FIGURE 7.)

FIGURE 6 SYSTEM BOARD-TO-I/O CONNECTIONS



The following is a brief explanation of how Flexible I/O works.

Each hardware partition consists of one or more System Boards, I/O Boards, and GSPBs. In a hardware partition containing two or more System Boards, one predefined System Board is connected directly to the I/O resources within the partition. Any CPU on a System Board not directly connected to the I/O resources communicates with each I/O unit via this predefined system board (called the Home SB).

This simple scheme has the advantage for flexible change of partition configuration.

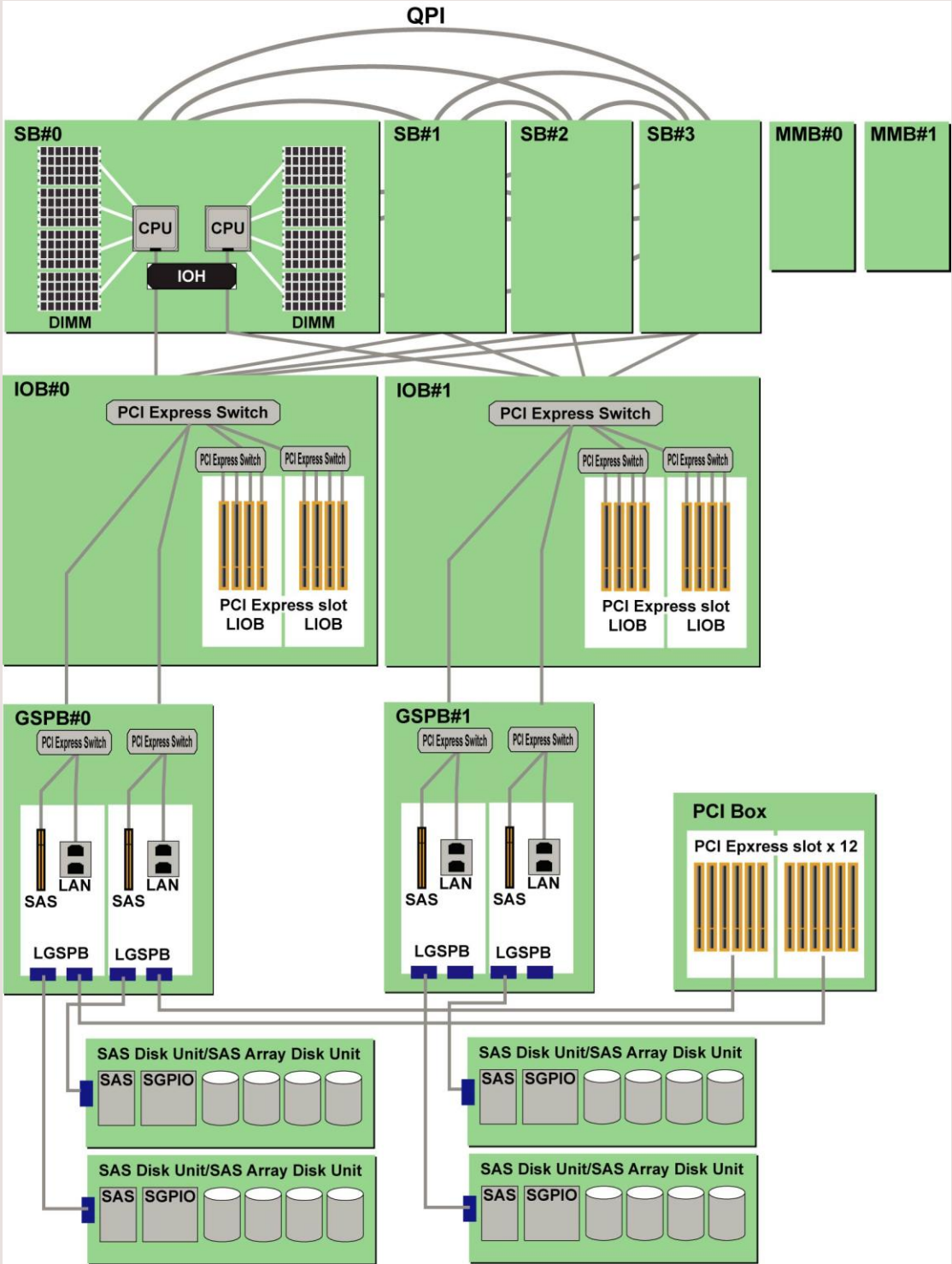
- Adding or removing a system board in a partition
Only the Home SB is connected directly to I/O resources. This completely eliminates the need to modify settings related to the I/O units.
- Changing I/O resources connected to a partition
Configuration can be simply changed by switching connections between the Home SB and I/O Boards/GSPBs. Should an error occur on a System Board, it is easy to change System Board connections with I/O units. After automatic switching from the faulty System Board to a Reserved System Board is completed, the I/O units are connected to the Home SB, and then the system is rebooted so it can resume the operation.

The settings for these operations are easily made using a Web GUI.

^{*1} I/O Board : PCI Express slots, ports accessible with System Boards, ports accessible with Giga LAN SAS & PCI Box interface board (GSPB)

^{*2} Giga LAN SAS & PCI Box interface board (GSPB) :
On-Board Ethernet ports, PCI Express interface for PCI- Box, and Interface for SAS Disk Units/SAS Array Disk Units

FIGURE 7 I/O CONFIGURATION OUTLINE



COMPONENT REDUNDANCY MAXIMIZES BUSINESS CONTINUITY

Most PRIMEQUEST components can be made redundant. Memory can be redundant by Memory mirror. Power supplies, fans, disk drives, PCI cards, and service processors can also be duplicated.

Furthermore, the main components are hot-replaceable.

- Power supplies, fans, disk drives, PCI cards, service processors, and a DVD drive
- In addition, the systems interconnect (QPI) initiates retransmission in the event of an error, and places the system in degradation mode when the number of retries reaches the threshold. QPI uses CRC to detect errors, initiating retries to try and resolve them.

There is no sacrifice in performance with QPI because it requires no additional clock cycle for CRC control. Even if an error occurs on a QPI link, QPI guarantees fault tolerance by narrowing the transfer width or by bypassing the faulty route.

Not only are most components redundant, the Mid-Plane, which connects these components to each other, is structured to suppress the occurrence of errors.

QPI links between System Boards, and signal lines between System Boards and IOHs, are all wired on the Mid-Plane. Its simple and minimal design suppresses the occurrence of failures as the connectors are only the electronic parts. (See FIGURE 9). Moreover, because the components are solderless, the use of hazardous substances has also been eliminated.

FIGURE8 PRIMEQUEST COMPONENT DIAGRAM

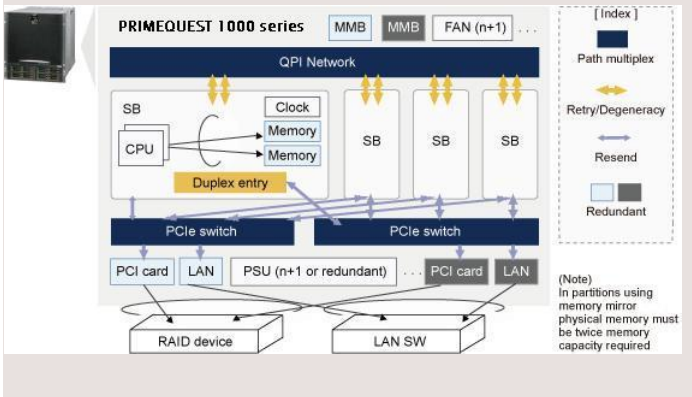


FIGURE 9 MID-PLANE



INTUITIVE MANAGEMENT VIEW

With the ServerView Operations Manager server management feature, you can quickly and correctly identify and solve server problems.

This feature displays diagrams showing faulty part locations so you can correctly identify the problem.

Moreover, it can remotely turn LEDs on or off on the server for fast identification of the failed part. The error notification most appropriate to your environment can be selected - possible choices include alarm notification by e-mail to a netbook PC or other mobile terminal, and error monitoring on the management console.

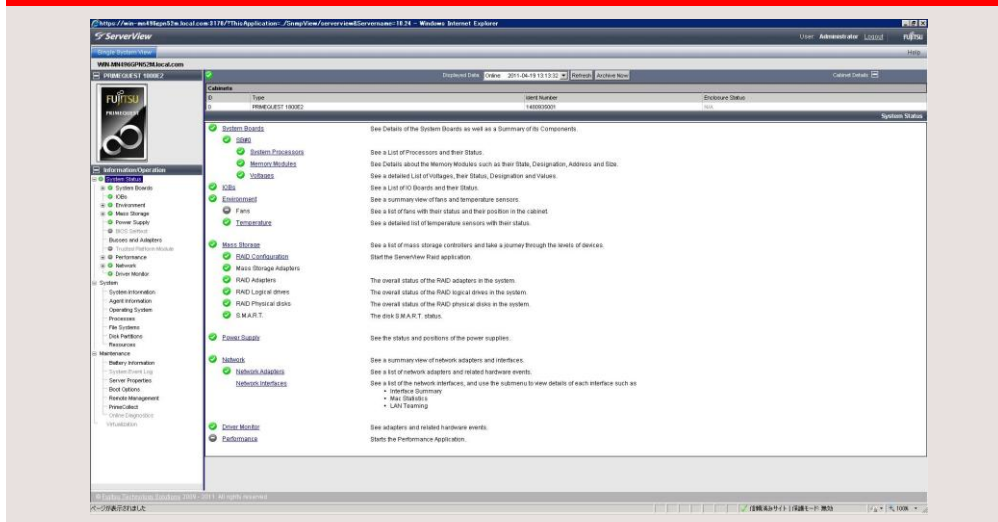
Each alarm can be assigned a priority and configured to specify the action to take when the alarm is received, ensuring alarms are

noticed, even in an emergency, and you can quickly respond to them.

To analyse application performance changes, you can narrow scope problems by system resources utilization changes using the intuitive graphical views.

In this way, the PRIMEQUEST, with ServerView Operations Manager, helps maximize your server operating time through early resolution of problems.

FIGURE 10 SERVERVIEW MANAGEMENT SCREEN



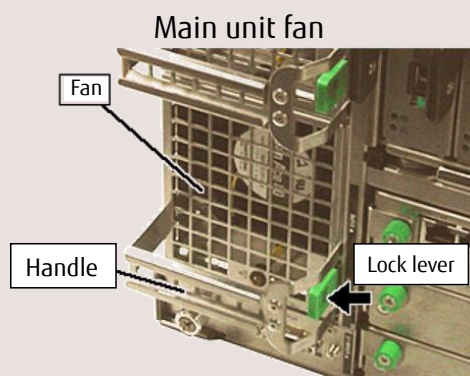
CABLE-LESS DESIGN

Simplifying maintenance operations can help prevent human error and reduce system down time. Other vendor's industry standard servers allow hot replacement of components, but hot replacement may actually be difficult to do. Especially if the servers are not designed with simple maintenance in mind.

PRIMEQUEST is trouble free in such situations as it is designed with an emphasis on convenience of actual maintenance work (FIGURE 11 (1) and (2)). Faulty units can be easily and individually removed from the front or rear of the cabinet.

Each unit uses only a minimum number of connectors, and they are interconnected through the Mid-Plane. Mistakes are avoided due to the easy installation and removal of units and the lack of separate cables.

FIGURE 11 UNIT SWAPPING FOR PRIMEQUEST 1800E (1)



To remove the fan, grasp the fan handle and pull it while pushing on the lock lever.
The handle is located either above or below the fan, depending on the fan mounting location (right or left side).

*** A Fujitsu certified field engineer performs this work.**

In contrast, component replacement with some other vendor servers can require complex procedures, including removal of the top panel of the server cabinet. Replacing components in an IBM X3950 X5 for example requires the removal of the cabinet's top panel (FIGURE 12)^{*1}.

- PCI card (addition, removal, or replacement)
- Fan (replacement)
- PSU
- The top panel must be put back into place within two minutes to maintain normal cooling and air flow.
- Hard disk drive

Although IBM asserts that PSUs, fans, and hard disks are all hot-replaceable^{*2}, hot replacement of these components seems to be very difficult in practice.

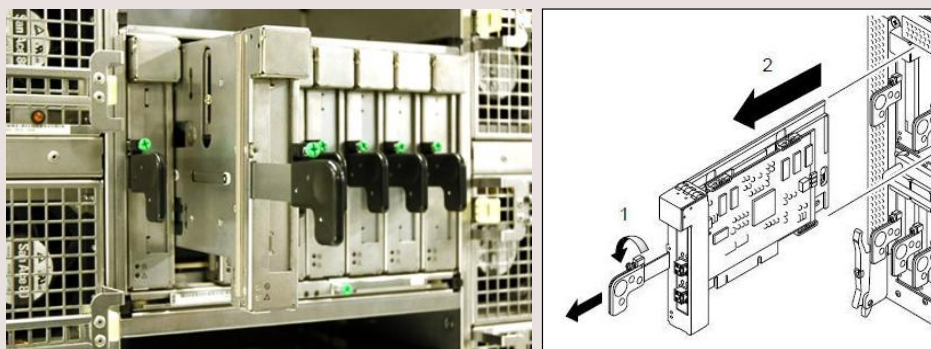
PRIMEQUEST frees you from complex work and avoids mistakes by its cable-less design that simplifies replacement work as much as possible.

*1 "IBM System x3950 Type 8878 and System x3950 E Type 8879 Problem Determination and Service Guide"

*2 IBM System x3950 specifications

<http://www-03.ibm.com/systems/x/hardware/enterprise/x3690x5/specs.html>

FIGURE 11 UNIT SWAPPING FOR PRIMEQUEST 1800E (2)

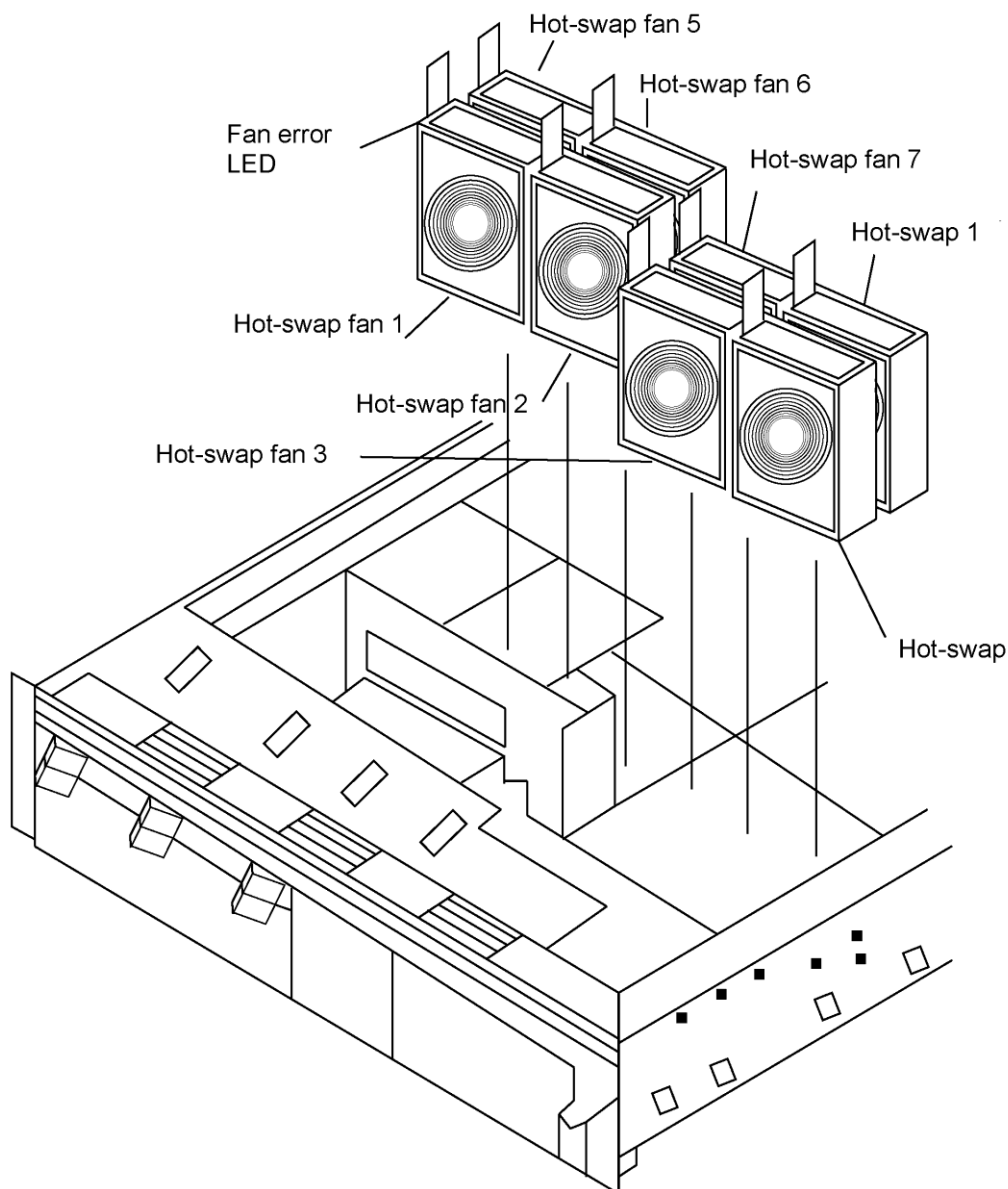


To remove the PCI cassette, first loosen the screw on top of the black handle, and pull out this lever as far as it will go while pressing on the top part of the PCI cassette. Then, remove the PCI cassette by pulling on the lever again.

*** A Fujitsu certified field engineer performs this work.**

FIGURE 12 EXAMPLE OF SERVER FAN REPLACEMENT

To remove hot-swap fan, complete the following steps.



CONCLUSION

HIGHEST BUSINESS CONTINUITY WITH EXCELLENT COST-EFFICIENCY

There is a "rumor" in IT industry that standard technology is unsuited to business critical applications. This whitepaper proves that is no longer true. Thorough data protection mechanisms in Intel® Xeon®E7 processor family and in memory are now sufficient for the handling of both recoverable and unrecoverable errors. In particular, error-prone areas of server memory are protected by multiple data protection mechanisms including error detection/correction by ECC, redundancy by Memory mirror, and memory health checking by Scrubbing.

Plus with Fujitsu's unique high-availability technologies in PRIMEQUEST platform uptime demanded by business critical applications is maximized. For example a Reserved System Board in combination with Flexible I/O can automate error recovery procedures, much reducing downtime.

With the obvious cost efficiency of using Xeon® processors, the intuitive administrative interface by Server View and the simple maintenance interface, and TCO can be reduced. In particular the Cable-less design and Mid-Plane dramatically simplify maintenance work, with most components replaceable by simple plug-in operation.

With its supreme simplicity and elegance of design PRIMEQUEST is one of the most robust and cost-efficient server platforms on the market today.

CONTACT

FUJITSU LIMITED
Website: www.fujitsu.com

© Copyright 2012 Fujitsu Limited

Fujitsu, the Fujitsu logo, PRIMEQUEST are trademarks or registered trademarks of Fujitsu Limited in Japan and other countries. Microsoft, MS, Windows, and Windows Server are trademarks or registered trademarks of Microsoft Corporation in the United States and, or other countries. LINUX is a trademark or registered trademark of Linus Torvalds in the United States and other countries. Intel and Xeon are trademarks or registered trademarks of Intel Corporation. Other company names and product names are the trademarks or registered trademarks of their respective owners.