

# Fault Tolerant Computer for Non-Stop Business

By Masahiko OKADA\*

**ABSTRACT** In recent years, intercorporate collaboration has been getting attention as a means to create new value to strengthen competitive power in market. In order to achieve efficient collaboration, both high availability and high connectivity are required at the same time for information systems. Express 5800/ft series based on standard IA server is one of the best products in this point. This paper describes the target and features of the Express 5800/ft series. It also describes various kinds of important technology to realize a redundant IA server. It also introduces several cases to explain how the fault tolerant computer is utilized.

**KEYWORDS** Fault tolerant, Mission critical, Server, Windows, IA (Intel Architecture)

## 1. INTRODUCTION

The surprising development of IT technologies has contributed to drastic changes in our daily life and business. For only ten-odd years, the Internet and personal computer have become popular at an explosive rate. Now everybody can communicate with all people throughout the world through e-mail and can buy goods through the Internet. This rapid change in market environment and the intensification of global competition has had a great effect on the way of business. In recent years, all enterprises have had no choice but to restructure their business and to concentrate on the core competence. In the future, in order to create a new value to strengthen the competitive power in a market, introduction of collaboration management will be required.

## 2. DYNAMIC COLLABORATION

In order to support intercorporate collaboration, NEC has been strengthening the products and services for supporting the collaboration while issuing a message “Dynamic Collaboration.” The target of intercorporate Dynamic Collaboration is the creation of new business through cost reduction, speedy management, and intercorporate cooperation and the creation of new business facing up to ubiquitous society. Enterprises in the future must consider what company they should have relations with and how they should contact by using property in addition to their own core competence. For this purpose, the enterprises should strengthen their core competence ultimately and understand the excellent technologies of other companies to unite their advantage and these

technologies organically to achieve their target. To implement fast and effective collaboration, utilization of information system will be an important factor.

To assist the purpose, NEC has been promoting OMCS (Open Mission Critical System) that includes a variety of products, technologies, and services. To enhance the connectivity of intercorporate cooperation, it is necessary to construct open products that conform to the business standards. For example, use of general-purpose software such as Windows and Linux and connection through Internet protocol are cited. In addition, high availability is required to support the security and continuity of the business.

## 3. FAULT TOLERANT COMPUTER

In order to bring secure computing to the business, NEC has developed a fault tolerant computer “Express 5800/ft series” which is based on IA server architecture in cooperation with Stratus Technologies, Ltd. The fault tolerant computer is one of the best platform products to meet the needs for high availability and robustness. Express 5800/ft series supports both Linux and Windows. For Windows model, NEC provides three categories: High-end, Mid-range, and Low-end model depends on the processor performance and I/O scalability. For Linux, only low-end model is available now due to market size, but we expect Linux market is expanding explosively for the future and plan to expand Linux products in response to rising needs. **Figure 1** shows the low-end Windows model as an example.

There are two options to increase system availability or robustness in general. One is fault tolerant system and the other is clustered system.

The fault tolerant system usually has two sets of hardware components and is able to work continuously by isolating the failed component in the event of hardware failure.

---

\*Client And Server Division



**Fig. 1 Fault tolerant computer product overview.**

The clustered system is configured with two or more systems and software-oriented architecture. In the event of system malfunction, the task is taken over from failed system to another system controlled by clustered software. The clustered system also has load balancing feature in addition to its high availability. These two kinds of systems are chosen to fit the purpose, as they have different features and advantages respectively. NEC is one of the few IT vendors that is able to provide various and secure solutions with these two kinds of high-availability systems. Here is the outline of the feature and the aim of the Express 5800/ft series.

### 3.1 Non-Stop Computing

The availability rate of Express 5800/ft series have achieved five nines (99.999%) according to our record. It has far exceeded standard IA servers and this means the downtime is just five minutes on annual average. It is an important issue to minimize downtime in the industry since downtime affects their business directly. Also, acceptable downtime is limited on the mission critical system. In the event of hardware failure, the fault tolerant computer is switched over from the failed component to another at once. Most operators will not realize the fact that the system was switched over, since the system keeps its performance. In the case of clustered system, it takes a long time and needs complex procedure to take over the tasks, because the failure is detected and controlled by software. This is one of the advantages of the fault tolerant computer.

### 3.2 Maintenance without Downtime

The fault tolerant computer can be repaired while operation is maintained, even if failure occurred, since fault tolerant computer is configured mostly

duplex. The following describes the rack mounted system “Express 5800/320Lb” as an example. Please see also Fig. 1. The system is composed of two pieces of CPU modules and PCI modules. Each module is 1U sized and it is easy to insert the module into or remove it from the rack system. Usually, the work service engineer does this, but it is easy for the customer to exchange the module. Exchanged module is detected automatically and is duplicated once again. The system keeps working while this operation is done.

### 3.3 Open Architecture

In order to achieve collaboration between enterprises, it is necessary to interconnect their information systems so that they can share and use their information effectively. However, there was a problem regarding connection between different types of information systems in the past. This is because many of the existing high-availability systems such as mainframe computers or fault tolerant computers are designed as proprietary system. To solve this problem, general operating systems such as Windows and Linux are able to work on the fault tolerant computer and widely used application software and business package software are also able to work without any changes. There is no requirement to the application software unlike clustered system. On the contrary, you might need to use application for awareness of cluster or need to make script program for failover in case of system malfunction if you choose clustered system. Application softwares for standard IA servers are able to work on the fault tolerant computer without any changes in principle.

### 3.4 Total Cost

Recently, it has become more and more important

to reduce total cost and development term of information system. The fault tolerant computer is an appropriate product in this point. If you choose fault tolerant computer, you can reduce total cost and development term, because you can use general application softwares or existing business package softwares without any changes. Furthermore, it is easy to operate and maintain the fault tolerant computer. This means you do not need a system engineer with high skill who is able to handle clustered system. Thus, medium and small companies are able to bring high-available system into their business easily. As stated above, the fault tolerant computer is the best product in the point of low cost and it is easy to develop and maintain the system in addition to its high availability.

#### 4. TECHNOLOGY OF FAULT TOLERANT COMPUTER

##### 4.1 Key Technology

This section describes the key technology to achieve a mostly redundant IA server. **Figure 2** describes the configuration of the fault tolerant computer.

##### (1) Redundant Processor and Lock Stepped Execution

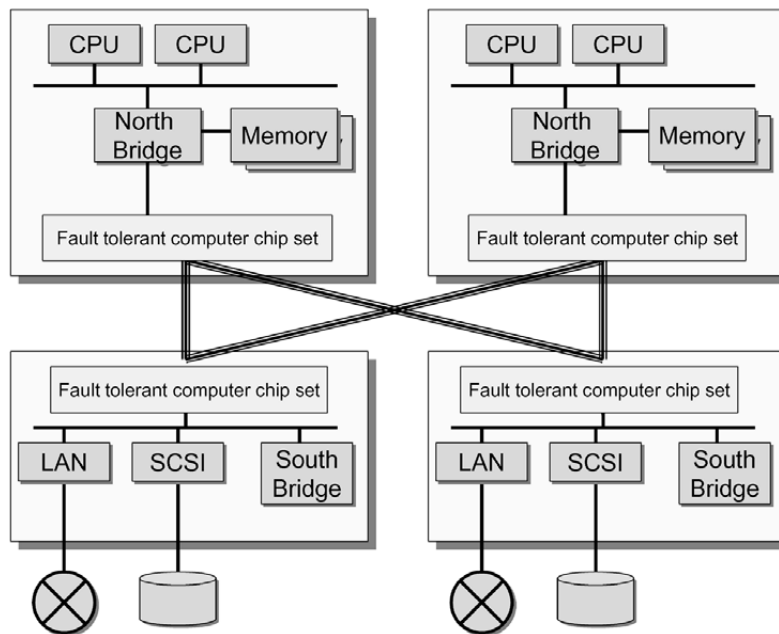
Although it has already been achieved to implement redundant component such as storage, network, and power supply for standard IA servers, it is diffi-

cult to implement redundant processor and main chip set from the point of cost and technique as yet.

There are two ways to configure redundant components in general. One of them is having two components and make one piece standby. In case of failure of the running component, the standby component takes over the following tasks which are supposed to be executed on the other component. Another way is also having two components. These two components are synchronized with each other. In case either of them fails, it is isolated and the other component works continuously. We adopted the latter way which synchronizes two processors for the fault tolerant computer. In concrete, CPU modules which is composed of processor, memory, and chip set are synchronized with each other. Both of the CPU modules are synchronized step by step based on a common clock signal. We call it lock stepped execution.

##### (2) Virtual Devices

We adopted standby method as redundant I/O subsystem as contrasted to redundant processor. I/O devices are switched over from failed one to the other at the time of hardware failure. This is achieved by configuring virtual single device from two real devices. It is controlled by customized hardware and device driver. In the case of hardware failure including power failure or module exchanges, it is detected and the working component is switched automatically by the customized hardware and the device



**Fig. 2** Fault tolerant computer architecture.

driver. The I/O commands are taken over and retried by the virtual device driver if necessary. The series of fail over procedure is carried out transparently to application software. The software does not realize the device is switched even if it is in use. Thanks to the mechanism, existing application can be used without any changes. With regard to the implementation for each individual device, this is explained in detail in the following section.

### (3) The Fault Tolerant Computer Chip Set

This section describes the most important chip set to implement fault tolerant computer architecture. The fault tolerant computer is composed of two sets of CPU modules and PCI modules. Each module is interconnected via crossbow link to each other. The crossbow link is controlled by the fault tolerant computer chip set. In other words, four modules are connected via the fault tolerant computer chip set which is mounted on each module. The fault tolerant computer chip set has the following important functions.

The crossbow link of the fault tolerant computer supports not only peer to peer connection but also peer to multiple nodes connection at the same time. For example, two CPU modules and one PCI module can be connected. And also one CPU module and two PCI modules can be connected. CPU modules and PCI modules can be connected, disconnected, or duplicated flexibly by this mechanism. In terms of logical aspect, crossbow link works as if they are transparent to the software. In this way, general operating system for standard IA server such as Windows and Linux are able to run on fault tolerant computer with only a few changes.

The fault tolerant computer chip set always checks whether both the CPU modules are synchronized or not after duplicated while comparing the transactions which are going through in the crossbow link. In the event of hardware failure, the fault tolerant computer chip set detects the event and isolates the failed module. Once the CPU module is isolated, accesses from the CPU module to PCI module are blocked by the fault tolerant computer chip set to avoid negative effect from the failed CPU module.

### (4) Modular Structure and Hot Swap

The policy regarding maintenance is not to stop the system. Availability does not improve if the system needs to stop the operation during maintenance, even though the system has duplicated configuration. To improve maintenance ability, most of the elements are configured as modular structure and we made it easy to exchange. Figure 1 shows an image of the

modular structure.

IA servers which have hot swap function such as storage device, power supply, and PCI card have increased recently. Furthermore, we enhanced the hot swap function to repair most of the elements so that we can reduce the down time. As mentioned above, the fault tolerant computer is composed of CPU module and PCI module. Most of the elements are gathered and mounted on either of these modules. Processor, memory and main chip set are mounted on the CPU module. Disk drive and various kinds of I/O adapter such as video, SCSI, and network controller are mounted on the PCI module. Even the power supplies are also mounted on the modules.

In order to implement hot swapping of the CPU module, the CPU module needs to be able to connect, disconnect, or duplicate the system dynamically. To achieve that, working state of the CPU module (processor context, memory, cache, bus status, and so on) is copied to another CPU module completely. In addition, some processes are done transparently to software in order to keep its operation.

The hot swapping of the PCI module is implemented by expanding the coverage of the existing hot swapping function of storage device and PCI card. **Figure 3** shows how it works in case of failure as an example.

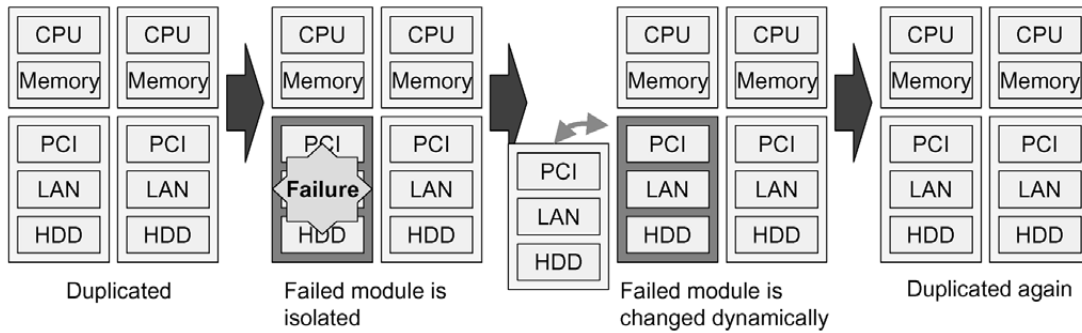
## 4.2 Implementation for Each Device

### (1) Storage Device

The internal storage devices are duplicated by software mirroring, which the general operating system already has. The hardware RAID that is used with general IA server is a reasonable solution, but if RAID controller fails or SCSI bus is disturbed by failed device, the system is not able to work continuously. There is a solution that has fully duplicated controllers and buses among high-end storage products, but this is still expensive to use for our fault tolerant computer. In such a present situation, we have chosen software mirroring, because its cost, performance, feasibility, and availability are well balanced. All of the components including disk drive, controller, and busses are separated in order to work continuously even if a failure occurs somewhere. **Figure 4** shows the mirroring configuration. All the disk drives and buses are separated. Hardware RAID is required to improve performance in the future.

### (2) Network

Redundant network technique has already been used with general IA server. There are two kinds of



**Fig. 3** Fault tolerant operation.

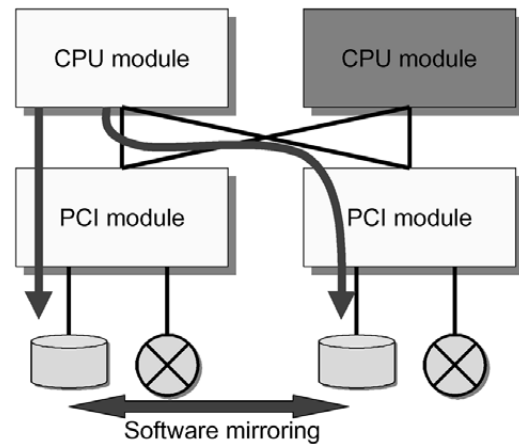
methods with regard to redundant network: standby and load balancing. Both the methods are enough from the viewpoint of availability. Then we have chosen AFT (Adapter Fault Tolerance) and ALB (Adaptive Load Balancing) function which existing network controller and device driver already have. The device driver creates virtual single network devices and acts in the same manner as an actual single network device. Switching is done by device driver transparently to application software in the event of failure, including problems occurring in the public network.

(3) Power Supply

In general, the failure rate of the power supply is relatively high, so redundant power supply is also used for general IA server. It is also important to ensure the quality of power supply for the fault tolerant computer. However, these kinds of power supply are supposed to supply to common devices, so they have shared parts as a necessity. Thus, failed power supply might have a negative effect on another power supply. To avoid this, the fault tolerant computer has a fully separated power supply. It is mounted on each CPU and PCI modules and its power feed line is also separated.

(4) Legacy Devices

General IA servers are composed of lots of devices. For example Video (VGA), Keyboard, Mouse, Serial interface, Timer, FDD, CD-ROM, Management controller and so on. Although some devices are not necessarily required redundant configuration, all of these devices must be made redundant to use general OS such as Windows or Linux, because failure of these devices might cause serious problems. These devices, called legacy devices, are furnished to maintain the compatibility with conventional PC/AT machine. They are lacking in flexibility and hot plug function, and some resources are even located at fixed



**Fig. 4** Redundant storage.

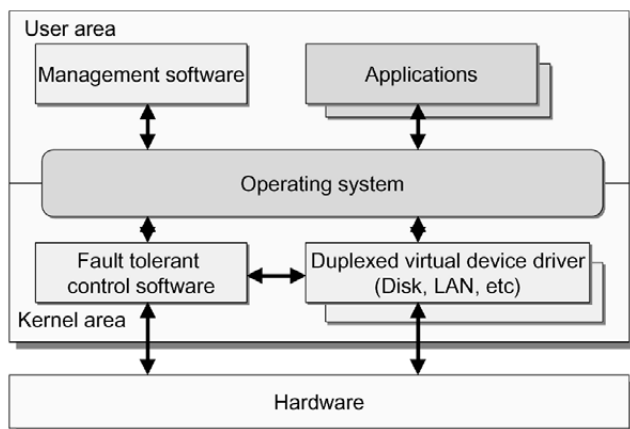
address and cannot be re-assigned. The legacy devices are duplicated in combination with some technique depending on its features. As for some devices such as keyboard and mouse, we are able to choose USB type device. It is difficult to create a virtual PS2 device but it is easy to switch USB device, as USB devices have hot plug function originally.

4.3 Software

This section describes software architecture. **Figure 5** shows simple software architecture.

(1) Use of General OS and Application

In order to achieve collaboration between enterprises, it is necessary to interconnect their information systems flexibly. To improve connectivity, it is efficient to use general products and technology. At this point, fault tolerant computer supports general OS such as Windows and Linux to respond to the needs. The mechanism is implemented by fault tolerant control software and a virtual device driver.



**Fig. 5 Software architecture.**

Thanks to the mechanism, most general application softwares and existing business package software are able to work. There is no need to change to those application softwares. Thus, it is now possible to develop flexible collaboration system easily.

#### (2) Fault Tolerant Control Software and Virtual Device Driver

The fault tolerant control software controls the whole system, for example, by duplicating and isolating each module in the event of hardware failure or module exchange. If hardware failure occurs during duplex operation, the failed module is isolated by the fault tolerant computer chip set and the fault tolerant control software checks it and tries to duplicate the module again if necessary.

On the other hand, the virtual device driver also has two important major functions regarding duplicating. One of them is creating a virtual device combined with two devices and acts like a single device. Application software is able to use it like a single device. Another function is failover. In case any device fails, the driver switches from the failed one to another and the command is retried. Thanks to this mechanism, there is no need to be aware that two devices are behind it for application software.

#### (3) Software Trouble

The principal purpose of the fault tolerant computer is improving the availability of hardware. However, we are considering to minimize the downtime making good use of duplicated hardware even if a sudden software trouble has happened.

Quick dumping function is one of them. Normally, it takes a very long time to dump memory in addition to restart when OS has crashed. In this case, fault

tolerant computer breaks lock stepped execution once. One of the CPU module tries to restart the OS, and the other CPU module keeps its memory and dumps them to storage device after the system boot up. After that the system is duplicated again. As a result, down time can be minimized.

#### (4) Management Software

We provide two kinds of management software: fault tolerant computer utility and ESMPRO (Fig. 6). To use the fault tolerant computer utility, you can monitor the status and diagnose the fault tolerant computer. Also, you can update firmware such as BIOS and BMC keeping software operation. This eliminates shutting down the fault tolerant computer to update them. The ESMPRO is common management software for NEC's IA server series. It is able to consolidate all Express 5800 servers including fault tolerant computer via network. You can manage your system stably and continuously using these management software to monitor and manage your servers.

#### (5) Platform Firmware

The fault tolerant computer works as simplex configuration until OS boots. Needless to say, it has the mechanism to recover when it fails during simplex operation. The BMC (Baseboard Management Controller) mounted on each PCI module has the role to do that. The BMC is an independent component which has its own processor and firmware and it controls power supply and monitors sensors such as temperature or voltage. The BMC also monitors whether the system boots normally by watch dog timer. If the BMC detects failure while booting, it attempts to reboot using another modules. In this way, the system can be recovered even if failure occurs during simplex booting. There is a possibility of failure of BMC. Both PCI module's BMC communicate with each other on a regular basis to solve the problem.

## 5. CASE STUDY

Nowadays, information systems are widely used as infrastructure, and high availability is expected. The following describes how fault tolerant computer is used and what is expected to the fault tolerant computer while taking up several cases.

Recently, we often hear the word "e-Government" or "electronic medical chart." It becomes popular to introduce to the field of the government and municipal offices or medical institute that was more passive to introduction of the information system compared

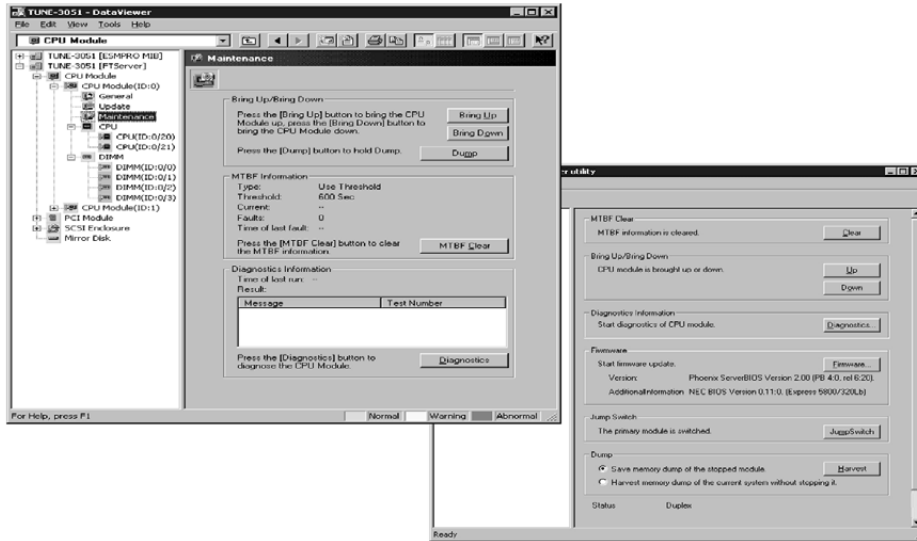


Fig. 6 Management software.

with the civilian sector such as the manufacturing industry and service industry. Naha city hall is one of them. They aim to achieve e-Government preceding other municipality. They are promoting computerized public service and sharing information among office staff using groupware. The system is used not only inside the city hall but also many public service agency such as branch office, hospital, waterworks department, boards of education and so on. The reason to choose fault tolerant computer is that this kind of public service is not allowed to stop its services even if in a few minutes. In case of general clustered system, it takes a long time to failover and public service is stopped while the failover process is done. Therefore fault tolerant computer is chosen.

Aozora Bank, Ltd. replaced a part of its information system from conventional mainframe computer to fault tolerant computer. Needless to say, although the information system is not so severe as the accounting system, it is mission critical for financial institutes. It is obvious to require high availability because of replacing a part of mainframe computer. In addition, the reason to choose fault tolerant computer is development term and operating cost. Clustered system is one of the options to improve availability, but it needs additional development and a skilled operator to handle the system. On the contrary, the fault tolerant computer can be handled as if it were a standard IA server. As a result, it can reduce total cost compared with the clustered system. For this reason, fault tolerant computer is being chosen

increasingly.

## 6. PROSPECTS

According to iDC report (September 2003, IDC#30028, Volume 1), high availability server market place is expected to grow at an annual rate of 10% or more for the future. The fault tolerant computer is used actually in various kinds of industry, including the above example, such as government and public offices, autonomy, medical, financial, media, manufacturing, distribution, services and so on. We can understand the growing needs for high availability from this fact. In my opinion, there are two trends for the future.

Although high availability products are mainly used in the server field so far, they are expected to spread gradually into the client products market from now on. It will become more and more important not only for servers but also client terminals in commercial transactions. ATM terminal and store computer are good examples of this. “e-procurement system” used for intercorporate transaction is also one of the mission critical systems. With such a perspective, high availability market place will grow in the future.

Clustered system and fault tolerant computer are used according to their characteristics at present. Both systems will improve while covering the weak point of each method and getting the advantages of the other method. Clustered technology will be installed in general OS as a basic function some day

and it will be easy to configure clustered system.

On the other hand, there are technical problems. Processor performance is going to progress more and more and asynchronous interface such as PCI Express will increase. These factors make it difficult to synchronize processors. In any case, it is expected that high availability market place will grow more and more while overcoming the difficulty.

### 7. CONCLUSION

To achieve effective collaboration between enterprises, it is important to interconnect their information systems so that they can make good use of their information for each other. In order to improve connectivity, it is efficient to use general products and technology. At this point, fault tolerant computer support general OS such as Windows and Linux to respond to needs, and this makes it easy to construct intercorporate information systems. As stated above, the fault tolerant computer is one of the best solutions from the viewpoint not only of its extremely high-availability but also its high connectivity and total cost.

NEC is going to provide high availability products including fault tolerant computer for the future continuously in response to rising demand.

### ACKNOWLEDGMENTS

The authors would like to thank the following individual: Messrs. Hajime Fukuzawa, Akihiko Obayashi.

### REFERENCES

- [1] R. Mittal, "Fault-tolerant Systems and Software," May 1996.
- [2] September 2003, IDC#30028, 1.
- [3] NEC Technical Journal, 55, 7, 2002-7.
- [4] <http://www.necsam.com/servers/products/>

\*Windows is a registered trademark or a trademark of Microsoft Corporation.

†Linux is a registered trademark of Linus Torvalds.

*Received March 16, 2004*

\* \* \* \* \*



Masahiko OKADA graduated from Oyama National College of Technology in 1987. He joined NEC Ibaraki, Ltd. in 1987 and was transferred to NEC Corporation in 2003. He is now Manager of the 1st Engineering Department, Client And Server Division. He is engaged in the development of Express 5800/ft series.

\* \* \* \* \*