

# the *Availability Digest*

## Fault-Tolerant Windows and Linux from Stratus

September 2007

### Introduction

ftServers from Stratus Technologies ([www.stratus.com](http://www.stratus.com)) provide plug-and-play fault tolerance for Windows and Red Hat Linux applications. Using Intel Xeon chips in a dual modular redundancy architecture, ftServers bring extremely high availability – five 9s and beyond – to the industry standard marketplace at affordable prices.



In contrast, average annual unplanned downtime for industry-standard servers has been measured to be about thirteen hours for Windows systems, sixteen hours for Red Hat Linux systems, and ten hours for UNIX systems. The availability of ftServers is continually monitored by Stratus, and a running availability is posted on its home page. This Uptime Meter indicates an average availability for the Stratus hardware and operating systems of greater than five 9s, or about two minutes per year of downtime.

Stratus' ftServer supports Windows, Linux, and Stratus' proprietary operating system, VOS. Their Continuum series of fault-tolerant platforms also support HP-UX and VOS.

### The Stratus ftServer Platform

The high availability achieved by the Stratus ftServer product line is achieved by running all applications on dual processors that are lockstepped at the memory access level. Should there be a disagreement between the processors, one of the processors has suffered a fault. If the faulty processor has detected its own fault, that processor is taken out of service. Otherwise, processing is paused as each processor enters a self-test mode; and the processor in error is taken out of service. Processing continues with the remaining good processor.<sup>1</sup>



The faulty processor can be replaced and synchronized with the operational processor while the system continues to run.

### The Processors

An ftServer comprises two logical processors. Each logical processor is packaged as a 2U module, similar to a blade, which Stratus calls a *slice*. The logical processors are interconnected via a small, passive backplane. A complete standalone ftServer, therefore, has a 4U form factor.

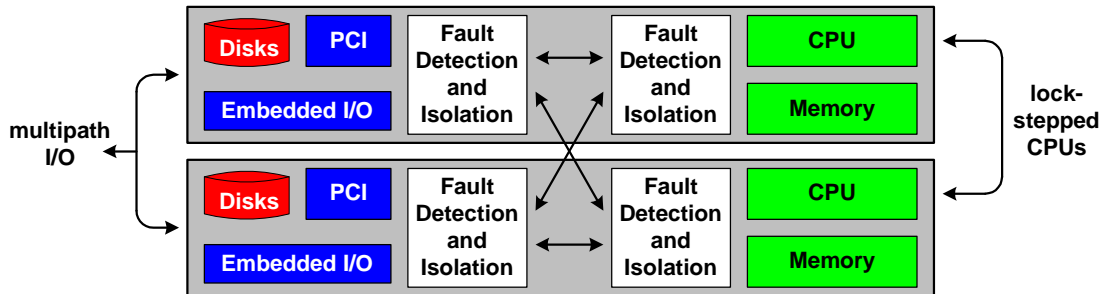
The ftServer logical processors use the Intel Xeon multicore microprocessors. A range of processing power is provided across the ftServer product line, from single dual-core logical

---

<sup>1</sup> This architecture differs from the early Stratus systems, in which there were two pairs of lockstepped processors. If one pair had a disagreement, it was taken out of service; and the other pair continued on.

processors to dual quad-core logical processors. The processing cores in a logical processor are organized as a symmetric multiprocessing (SMP) set.

However, the Xeon processing chips are not your everyday microprocessors. Their operation must be unfailingly deterministic in order for lockstepping to work. That means that any core must process instructions in exactly the same sequence.



Stratus works closely with its partner, Intel, to ensure that the Xeon processors used in the ftServers meet exacting deterministic standards. This determinism is now a standard feature of all delivered Xeon processors.

### ***I/O Subsystem***

Each 2U module, or slice, contains its own I/O subsystem. Each I/O subsystem also has its own fault detection and isolation logic. In normal operation, peripheral devices are driven by both I/O subsystems via a multi-path connection. However, if an I/O subsystem detects a malfunction, it will remove itself from service.

Each logical processor is connected to each of the I/O subsystems. In this way, any combination of one logical processor failure and one I/O subsystem failure will not render the system inoperable.

### ***Integrated Disks***

Each slice contains not only its logical microprocessor cores, memory, and I/O subsystem but also three 500 gigabyte integrated disks. Thus, each logical processor has direct access to 1.5 terabytes of local storage. This storage is mirrored between the slices, providing fault-tolerant storage within the ftServer itself.

### ***Operating Systems***

The Windows and Red Hat Linux operating systems that Stratus supports are those that are commercially available to anyone. Therefore, the ftServers are application binary interface (ABI) compatible with Windows and Linux applications. Any application that can run under Windows or Red Hat Linux on an industry-standard server can run on an ftServer without modification. The installation and administration procedures are identical. The user should see no difference except for downtime – and that is the big difference.

Behind the scenes, Stratus spends a great deal of effort ensuring that these operating systems measure up to its high-availability requirements. Working with its partners, Microsoft and Red Hat, Stratus works diligently to harden the operating systems so that operating system faults are minimized.

When a new operating system version is about to be released, Stratus engineers do everything that they can to break it and the drivers that are supplied with it by the peripheral vendors. They

use a facility they call “Breaker” to induce a wide variety of faults, from device errors to power glitches. The Stratus engineers claim that they have not yet found a device driver that they could not break. Faults that involve the operating system are reported back to Microsoft and Red Hat, who make corrections to eliminate the sources of those failures. Actually, everyone benefits from this effort because it is the hardened version that is released for public use. The Windows and Red Hat Linux operating systems that Stratus runs are the standard commercially available versions.

However, Stratus has found that the weak links in the software systems are the device drivers. Therefore, Stratus develops its own hardened device drivers that will stand up to the stresses generated by Breaker. As a consequence, Stratus will only support certain peripheral devices. This limits the choice of peripherals such as external disk arrays and network interfaces. Since the operating systems are standard, anyone can install any device driver that they want. However, if they elect to do so, they are then on their own from an availability viewpoint.

### ***Hardware Fault Detection and Isolation***

In normal operation, both logical processors are processing the same instruction stream and are comparing their results at a memory-access level as they proceed. Providing that they agree, it is known that they are correct; and they continue on.

There are several processor failure modes:

- One failure mode occurs when one of the processors, through its own fault-detection logic, recognizes that its operation is erroneous. In this case, it takes itself out of service, and the other processor carries on with normal processing. This “fast-fail” action ensures that the faulty processor will not propagate the error to external interfaces or to the database.

However, it is possible that this error was a transient error. Therefore, the failed processor will run a self-check diagnostic. If it passes the diagnostic, it returns itself to service and is resynchronized with the operational processor so that fault tolerance is restored. A count is kept of transient errors for each logical processor. Should this count exceed a specified limit, the processor is taken out of service and must be replaced.

- A second failure mode occurs when the two logical processors disagree, but neither has declared itself at fault. In this case, processing is paused; and the two logical processors each run their self-test diagnostic. If one should determine that it is indeed at fault, it is taken out of service (since the fault was detectable, it was not a transient fault).
- A third failure mode occurs when the two logical processors disagree, and both pass their self-diagnostic tests. In this case, the problem is likely to be a timing problem, such as two simultaneous interrupts being processed in different order. One logical processor is declared the winner, and the other is resynchronized with it. Operation in fault-tolerant mode continues.

In any case, fault-tolerant operation continues on as long as there are two good logical processors. Should one processor fail, the system continues on as a single nonfault-tolerant system until the failed processor is replaced.

### ***System Recovery from a Hardware Fault***

When a logical processor fails, this event is transparent to the outside world. The system continues on with one processor.

When a replacement logical processor is available, it is slid into the 2U slice slot to replace the failed processor. Its memory and internal disk-resident database are then resynchronized with the operating processor while system operation continues, after which the new processor is put into service. Fault-tolerant operation is now restored.

Resynchronization impacts system performance only to a small degree. At the end of the synchronization process, there is about a 100 millisecond pause as the new logical processor is put into service.

### ***Operating System Fault Recovery***

Even with the extensive efforts to harden the operating system, the operating system will crash on occasion. This will, of course, take down the ftServer; and it must be rebooted.

When this happens with an industry-standard server, it is good practice to obtain a dump of the operating system so that the fault can be tracked down and corrected. However, this is a timely process during which the system continues to be down.

Stratus makes use of its dual processors to alleviate this problem. Following a crash, only one of the logical processors is rebooted and becomes operational. The other is used to perform the dump. When the dump is finished, the processor will join the operating processor, restoring fault tolerance to the ftServer.

### **ftServer Models**

ftServers currently come in three models:

- ftServer 2500 – single socket dual core, 2 ghz processor, up to 6 gb of memory
- ftServer 4400 – single or dual socket dual core, 2 ghz processors, up to 12 gb of memory
- ftServer 6200 – dual socket quad core, 2.6 ghz processors, up to 24 gb of memory

### **Storage Arrays**

To augment the 1.5 terabyte mirrored internal storage, Stratus offers its ftScalable storage array. This array has a 2U form factor and includes dual RAID controllers with redundant cache. Each ftScalable array contains twelve 300 gb disks, yielding a capacity of 3.6 terabytes. Up to three ftScalable arrays can be attached to an ftServer and can provide up to 10.8 terabytes of data.



The disks can be arranged in any desired RAID array. RAID levels may be mixed within a single ftScalable array. The disks are hot-pluggable, and the dual RAID controllers provide a rolling upgrade capability for installing firmware upgrades without taking down the array.

Stratus also supports certain EMC disk storage units.

### **Virtualization**

Virtualization technology allows many virtual machines to run on a single server. Virtualization is an important technique for getting full utilization out of large server farms.

A problem faced with virtualization technology is that availability becomes far more important. One particular application may not have a high value and can suffer some downtime without serious consequence. However, run many of these as virtual machines on a single server, and a server failure becomes much more costly.

Stratus is integrating the ftServer with VMware's ESX Server to allow ftServers to host many virtual Windows and Linux machines in any combination. The ESX Server sits on top of the ftServer hardware and supports multiple instances of different operating systems running as if they were in their own physical server.

As a result, it will be very simple to add a fault-tolerant pool of servers to a virtualized server farm.

## VOS

VOS is Stratus' original proprietary operating system. Perfected by over two decades of engineering enhancements, VOS is still in common use today. As with many legacy systems, there are simply too many applications written for VOS for people to migrate to other platforms. In addition, the VOS Posix interface allows many open applications to run on this extremely reliable platform.

VOS is available on Stratus' line of V series ftServers.

## The Stratus Continuum Family

The Continuum product line is an older line that is still marketed and supported by Stratus. It is based on PA-RISC microprocessors. The PA-RISC chips are no longer made by HP, but Stratus states that it has an extensive inventory of these chips that will allow it to support the Continuum series into the foreseeable future.

The Continuum series of servers supports both VOS and HP-UX.

## Rolling Upgrades

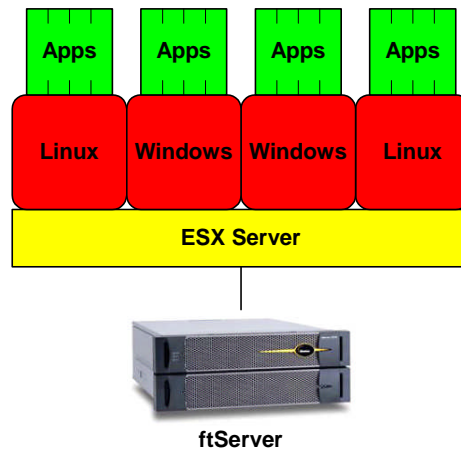
The Stratus ActiveUpgrade facility allows Windows updates to be installed without taking down the system. To do this, one of the logical processors is taken offline; and its operating system is updated. It is then put online, and the applications are restarted in order to be compatible with the new operating system. The other processor is then upgraded and is returned to service to restore fault tolerance.

Stratus has not made this facility available for VOS, as VOS is now very stable with few upgrades. They also have not made it available for Linux as they feel that virtualization will make this unnecessary.

## Stratus Call Home

All Stratus servers monitor themselves for faults. If a problem of any kind is detected, the system will automatically call a Stratus support facility (provided customer permission has been granted to do so); and action is immediately taken to diagnose the problem.

If it is decided that a component needs to be replaced, Stratus will immediately send the component to the customer's site. There are many cases in which the first sign of a problem to the customer is when he receives the replacement part in the mail.



## Disaster Recovery

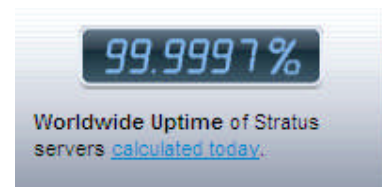
Stratus supports remote standby sites that are kept synchronized with the primary site via data replication. For Windows environments, the Double-Take asynchronous replication engine is used. Double-Take will also replicate virtual Windows environments.

GoldenGate is used to replicate databases running under Linux, and DataWise is used to back up VOS systems.

## The Stratus Uptime Meters

The success of all of the availability features described above is reflected in the Stratus Uptime Meters, one for ftServers and one for Continuum servers. The ftServer meter is updated daily, whereas the Continuum meter is updated weekly. The meters are accessible from Stratus' home page (the ftServer Uptime Meter is shown there).

Calculated availability is based on all reported service incidents over the preceding six-month period that impacted production. Both hardware-related and Stratus-supplied software-related incidents are included.



The Uptime Meters typically show availabilities of 99.9996% or better, equivalent to about two minutes per year of downtime.

## Market Positioning

Stratus' mantra is availability, simplicity, and cost:

- Availability - ftServers achieve availabilities of five 9s and beyond.
- Simplicity - The ftServers are plug-and-play. They are installed and are administered just like any other industry-standard server.
- Cost - Server prices, ranging from the low- to mid-five figures, are affordable.

Stratus focuses on the industry-standard server community by offering fault-tolerant systems that can replace normal industry-standard servers.

What about the other fault-tolerant systems – HP NonStop servers? Both HP and Stratus marketing people say that they rarely run into each other. Stratus focuses on the low-end industry-standard server market, and the HP NonStop folks focus on the very large systems requiring massive scalability, a feature Stratus does not support. They each say that if they run into each other, probably one of them is in the wrong place.

Nevertheless, the low end is pervasive. Stratus is found in eight of the top ten banks, 12 of the top 15 pharmaceuticals, fourteen of the top twenty telecommunication companies, and in over 225 public agencies.

## Stratus – the Phoenix Risen from the Ashes

Stratus Computers was founded in 1980 by Bill Foster. It became a publicly held company but was then acquired by Ascend in 1998 for Stratus' telecommunication products. Shortly after that, Ascend was acquired by Lucent. Lucent spun off all but the telecommunication products to former Stratus management in 1999, and Stratus Technologies was born. In 2003, Stratus Technologies reacquired its telecommunications products from Lucent.

The new Stratus Technologies now has over 5,000 customers in 47 countries. Its 700+ employees generate over a quarter billion dollars in revenue. Stratus, like the Phoenix, has come full circle from birth to disappearance to rebirth.

## Summary

Stratus today is the predominant player in the fault-tolerant, industry-standard server marketplace. With its lockstep hardware technology, its hardened failsafe software, and its call-home and availability services, it provides affordable, plug-and-play fault-tolerant solutions for Windows, Red Hat Linux, VOS, and HP-UX systems. Via its port of VMWare, it is now entering the fault-tolerant server market for virtualized data centers.

It is a tribute to the short-sightedness of today's IT managers that fault-tolerant solutions such as this and others have not yet become pervasive in modern data centers.

