

# I/O Performance HOWTO

**Sharon Snider**

v1.1, 05/2002

## **Revision History**

Revision v1.1                    2002-05-01                    Revised by: sds

Updated technical information and links.

Revision v1.0                    2002-04-01                    Revised by: sds

Wrote and converted to DocBook XML.

This HOWTO covers information on available patches for the 2.4 kernel that can improve the I/O performance of your Linux" operating system.

---

# Table of Contents

<a href="#"><b>1. Distribution Policy</b></a> .....	<b>1</b>
<a href="#"><b>2. Introduction</b></a> .....	<b>2</b>
<a href="#"><b>3. Avoiding Bounce Buffers</b></a> .....	<b>3</b>
<a href="#">3.1. Memory and Addressing in the Linux 2.4 Kernel</a> .....	3
<a href="#">3.2. The Problem with Bounce Buffers</a> .....	3
<a href="#">3.3. Locating the Patch</a> .....	3
<a href="#">3.3.1. Configuring the Linux Kernel to Avoid Bounce Buffers</a> .....	3
<a href="#">3.3.2. Enabled Device Drivers</a> .....	4
<a href="#">3.4. Modifying Your Device Driver to Avoid Bounce Buffers</a> .....	4
<a href="#"><b>4. Raw I/O Variable–Size Optimization Patch</b></a> .....	<b>6</b>
<a href="#">4.1. Locating the Patch</a> .....	6
<a href="#">4.2. Modifying Your Driver for the Raw I/O Variable–Size Optimization Patch</a> .....	6
<a href="#"><b>5. I/O Request Lock Patch</b></a> .....	<b>7</b>
<a href="#">5.1. Locating the Patch</a> .....	7
<a href="#">5.2. Modifying Your Driver for the I/O Request Lock Patch</a> .....	7
<a href="#"><b>6. Additional Resources</b></a> .....	<b>8</b>

# **1. Distribution Policy**

The I/O Performance–HOWTO is copyrighted © 2002, by IBM Corporation

Permission is granted to copy, distribute, and/or modify this document under the terms of the GNU Free Documentation License, Version 1.1 or any later version published by the Free Software Foundation with no Invariant Sections, no Front-Cover text, and no Back-Cover text. A copy of the license can be found at <http://www.gnu.org/licenses/fdl.txt>.

---

## **2. Introduction**

This HOWTO provides information on improving the input/output (I/O) performance of the Linux operating system for the 2.4 kernel. Additional patches will be added as they become available.

Please send any comments, or contributions via e-mail to [Sharon Snider](#).

---

# 3. Avoiding Bounce Buffers

This section provides information on applying and using the bounce buffer patch on the Linux 2.4 kernel. The bounce buffer patch, written by Jens Axboe, enables device drivers that support direct memory access (DMA) I/O to high-address physical memory to avoid bounce buffers.

This document provides a brief overview on memory and addressing in the Linux kernel, followed by information on why and how to make use of the bounce buffer patch.

---

## 3.1. Memory and Addressing in the Linux 2.4 Kernel

The Linux 2.4 kernel includes configuration options for specifying the amount of physical memory in the target computer. By default, the configuration is limited to the amount of memory that can be directly mapped into the kernel's virtual address space starting at PAGE\_OFFSET. On i386 systems the default mapping scheme limits kernel-mode addressability to the first gigabyte (GB) of physical memory, also known as low memory. Conversely, high memory is normally the memory above 1 GB. High memory is not directly accessible or permanently mapped by the kernel. Support for high memory is an option that is enabled during [configuration of the Linux kernel](#).

---

## 3.2. The Problem with Bounce Buffers

When DMA I/O is performed to or from high memory, an area is allocated in low memory known as a bounce buffer. When data travels between a device and high memory, it is first copied through the bounce buffer.

Systems with a large amount of high memory and intense I/O activity can create a large number of bounce buffers that can cause memory shortage problems. In addition, the excessive number of bounce buffer data copies can lead to performance degradation.

Peripheral component interface (PCI) devices normally address up to 4 GB of physical memory. When a bounce buffer is used for high memory that is below 4 GB, time and memory are wasted because the peripheral has the ability to address that memory directly. Using the bounce buffer patch can decrease, and possibly eliminate, the use of bounce buffers.

---

## 3.3. Locating the Patch

The latest version of the bounce buffer patch is *block-highmem-all-18b.bz2*, and it is available from Andrea Arcangeli's -aa series kernels at <http://kernel.org/pub/linux/kernel/people/andrea/kernels/v2.4/>.

---

### 3.3.1. Configuring the Linux Kernel to Avoid Bounce Buffers

This section includes information on configuring the Linux kernel to avoid bounce buffers. The Linux Kernel-HOWTO at <http://www.linuxdoc.org/HOWTO/Kernel-HOWTO.html> explains the process of

re-compiling the Linux kernel.

The following kernel configuration options are required to enable the bounce buffer patch:

*Development Code* – To enable the configurator to display the High I/O Support option, select the Code maturity level options category and specify "y" to Prompt for development and/or incomplete code/drivers.

*High Memory Support* – To enable support for physical memory that is greater than 1 GB, select the Processor type and features category, and select a value from the High Memory Support option.

*High Memory I/O Support* – To enable DMA I/O to physical addresses greater than 1 GB, select the Processor type and features category, and enter "y" to the HIGHMEM I/O support option. This configuration option is a new option introduced by the bounce buffer patch.

---

### 3.3.2. Enabled Device Drivers

The bounce buffer patch provides the kernel infrastructure, as well as the SCSI and IDE mid-level driver modifications to support DMA I/O to high memory. Updates for several device drivers to make use of the added support are also included with the patch.

If the bounce buffer patch is applied and you configure the kernel to support high memory I/O, many IDE configurations and the device drivers listed below perform DMA I/O without the use of bounce buffers:

```
aic7xxx_drv.o
aic7xxx_old.o
cciss.o
cpqarray.o
megaraid.o
qllogicfc.o
sym53c8xx.o
```

---

## 3.4. Modifying Your Device Driver to Avoid Bounce Buffers

If your device drivers are not listed above in the [Enabled Device Drivers](#) section, and the device is capable of high-memory DMA I/O, you can modify your device driver to make use of the bounce buffer patch as follows. More information on rebuilding a Linux device driver is available at <http://www.xml.com/ldd/chapter/book/index.html>.

1. A.) For SCSI Adapter Drivers: set the *highmem\_io* bit in the *Scsi\_Host\_Template* structure.
- B.) For IDE Adapter Drivers: set the *highmembit* in the *ide\_hwif\_t* structure.
2. Call *pci\_set\_dma\_mask(struct pci\_dev \*pdev, dma\_addr\_t mask)* to specify the address bits that the device can successfully use on DMA operations.

If DMA I/O can be supported with the specified mask, `pci_set_dma_mask()` will set `pdev->dma_mask` and return 0. For SCSI or IDE, the mask value will also be passed by the mid-level drivers to `blk_queue_bounce_limit(request_queue_t *q, u64 dma_addr)` so that bounce buffers are not created for memory directly addressable by the device. Drivers other than SCSI or IDE must call `blk_queue_bounce_limit()` directly.

3. Use `pci_map_page(dev, page, offset, size, direction)`, instead of `pci_map_single(dev, address, size, direction)` to map a memory region so that it is accessible by the peripheral device. `pci_map_page()` supports both high and low memory.

The `address` parameter for `pci_map_single()` correlates to the `page` and `offset` parameters for `pci_map_page()`. Use the `virt_to_page()` macro to convert an `address` to a `page` and `offset`. The `virt_to_page()` macro is defined by including `pci.h`. For example:

```
void *address;
struct page *page;
unsigned long offset;
page = virt_to_page(address);
offset = (unsigned long) address & ~PAGE_MASK;
```

Call `pci_unmap_page()` after the DMA I/O transfer is complete to remove the mapping established by `pci_map_page()`.



`pci_map_single()` is implemented using `virt_to_bus()`. `virt_to_bus()` handles low memory addresses only. Drivers supporting high memory should no longer call `virt_to_bus()` or `bus_to_virt()`.

4. If your driver calls `pci_map_sg()` to map a scatter-gather DMA operation, your driver should set the `page` and `offset` fields instead of the `address` field of the `scatterlist` structure. Refer to step 3 for converting an `address` to a `page` and `offset`.



If your driver is already using the PCI DMA API, continue to use `pci_map_page()` or `pci_map_sg()` as appropriate. However, do not use the `address` field of the `scatterlist` structure.

---

# 4. Raw I/O Variable–Size Optimization Patch

This section provides information on the raw I/O variable–size optimization patch for the Linux 2.4 kernel written by Badari Pulavarty. This patch is also known as the RAW VARY or PAGESIZE\_io patch.

The raw I/O variable–size patch changes the block size used for raw I/O from *hardsect\_size* (normally 512 bytes) to 4 kilobytes (K). The patch improves I/O throughput and CPU utilization by reducing the number of buffer heads needed for raw I/O operations.

---

## 4.1. Locating the Patch

You can download the patch from one of the following locations:

- Andrea Arcangeli has made the patch available at <http://www.kernel.org/pub/linux/kernel/people/andrea/kernels/v2.4/2.4.18pre7aa2/>. The name of the file is *10\_rawio-vary-io-1*.
  - Alan Cox has included the patch in the 2.4.18pre9–ac2 kernel patch. The patch is available at <http://www.kernel.org/pub/linux/kernel/people/alan/linux-2.4/2.4.18/>.
  - The patch is available from SourceForge at <http://sourceforge.net/projects/lse/io>. The latest version is *PAGESIZE\_io-2.4.17.patch*.
- 

## 4.2. Modifying Your Driver for the Raw I/O Variable–Size Optimization Patch

In previous versions of this patch, changes were enabled for all drivers. However, the 2.4.17 and later versions of the patch enable the changes only for the Adaptec, Qlogic ISP1020, and IBM ServerRAID drivers. All other drivers for version 2.4.17 and later must be modified to make use of the patch by setting the *can\_do\_varyio* bit in the *Scsi\_Host\_Template* structure.



Drivers that have the raw I/O patch enabled must support buffer heads of variable sizes (*b\_size*) in a single I/O request because *hardsect\_size* is used until the data buffer is aligned on a 4 K boundary.

Additional information is available on rebuilding Linux device drivers at <http://www.xml.com/ldd/chapter/book/index.html>.

---

# 5. I/O Request Lock Patch

This section provides information on the I/O request lock patch, also known as the scsi concurrent queuing patch (*sior1*), written by Johnathan Lahr.

The I/O request lock patch improves SCSI I/O performance on Linux 2.4 multi-processor systems by providing concurrent I/O request queuing. There are significant I/O performance and CPU utilization improvements possible by enabling multi-processors to concurrently drive multiple block devices.

Before the patch is applied block I/O requests are queued one at a time holding the global spin lock, *io\_request\_lock*. Once the patch is applied, SCSI requests are queued while holding the lock specific to the queue associated with the request. Requests that are made to different devices are queued concurrently, and requests that are made to the same device are queued serially.

---

## 5.1. Locating the Patch

You can download the I/O request patch from Sourceforge at <http://sourceforge.net/projects/lse/io>. The latest version is *sior1-v1.2416*. Patches that enable concurrent queuing for specific drivers are also available at SourceForge. The patch for the Emulex SCSI/FC is *lpfc\_sior1-v0.249* and the patch for Adaptec SCSI is *aic\_sior1-v0.249*.

---

## 5.2. Modifying Your Driver for the I/O Request Lock Patch

The I/O request lock patch installs concurrent queuing capability into the SCSI midlayer. Concurrent queuing is activated for each SCSI adapter device driver. To activate the driver, the *concurrent\_queue* field in the *Scsi\_Host\_Template* structure must be set when the driver is registered.



Drivers that activate concurrent queuing must ensure that any access of the *request\_queue* by the driver is protected by the *request\_queue.queue\_lock*.

Additional information is available on rebuilding device drivers at <http://www.xml.com/ldd/chapter/book/index.html>.

---

## 6. Additional Resources

The following list of Web sites provides additional information on modifying device drivers and configuring the Linux kernel.

- Information on Dynamic DMA mapping is available at  
<http://lwn.net/2001/0712/a/dma-interface.php3>.
- Kernel-HOWTO is available from the Linux Documentation Project at  
<http://www.linuxdoc.org/HOWTO/Kernel-HOWTO.html>.
- Linux Device Drivers, 2nd Edition published by O'Reilly is available online at  
<http://www.xml.com/ldd/chapter/book/index.html>.