# The Vinum Volume Manager

Greg Lehey

LEMIS (SA) Pty Ltd

PO Box 460

Echunga SA 5153.

`grog@lemis.com`

`grog@FreeBSD.org`

`grog@NetBSD.org`

*ABSTRACT*

The *Vinum Volume Manager* is a device driver which implements virtual disk drives. It isolates disk hardware from the device interface and maps data in ways which result in an increase in flexibility, performance and reliability compared to the traditional slice view of disk storage. Vinum implements the RAID-0, RAID-1, RAID-4 and RAID-5 models, both individually and in combination.

*Vinum* is an open source volume manager which runs under FreeBSD and NetBSD. It was inspired by the VERITAS® volume manager and implements many of the concepts of VERITAS®. Its key features are:

- Vinum implements many RAID levels:

    - RAID-0 (striping).

    - RAID-1 (mirroring).

    - RAID-4 (fixed parity).

    - RAID-5 (block-interleaved parity).

    - RAID-10 (mirroring and striping), a combination of RAID-0 and RAID-5.

    In addition, other combinations are possible for which no formal RAID level definition exists.

- Volume managers initially emphasized reliability and performance rather than ease of use. The results are frequently down time due to misconfiguration, with consequent reluctance on the part of operational personnel to attempt to use the more unusual features of the product. Vinum attempts to provide an easier-to-use non-GUI interface.

In place of conventional disk partitions, Vinum presents synthetic disks called *volumes* to

the user. These volumes are the top level of a hierarchy of objects used to construct volumes with different characteristics:

- The top level is the virtual disk or *volume*. Volumes effectively replace disk drives.

- Volumes are composed of one or more *plexes*, each of which represents the total address space of a volume. This level in the hierarchy thus provides redundancy.

- Since Vinum exists within the UNIX disk storage framework, it would be possible to use UNIX partitions as the building block for multi-disk plexes, but in fact this turns out to be too inflexible: UNIX disks can have only a limited number of partitions. Instead, Vinum subdivides a single UNIX partition (the *drive*) into contiguous areas called *subdisks*, which it uses as building blocks for plexes.

- Subdisks reside on Vinum *drives*, currently UNIX partitions. Vinum drives can contain any number of subdisks. With the exception of a small area at the beginning of the drive, which is used for storing configuration and state information, the entire drive is available for data storage.

With this structure, Vinum offers the following features:

## Unlimited disk size

Since Vinum plexes are composed of subdisks which can reside on any Vinum volumes can exceed the size of any single disk. There is no intrinsic limitation on the size of a plex or a volume.

## Faster access

Concurrent access to a disk severely limits the throughput: with modern disks, sequential disk transfer rates exceed 50 MB/s, but most transfers are less than 32 kB in size and thus take only in the order of 500 μs to complete. By contrast, seek latency takes about ten times this time, so every seek drops the aggregate throughput by a factor of ten or more, depending strongly on the length of the transfer.

The traditional and obvious solution to this bottleneck is "more spindles": rather than using one large disk, it uses several smaller disks with the same aggregate storage space. Each disk is capable of positioning and transferring independently, so the effective throughput increases by a factor close to the number of disks used.

The exact throughput improvement is, of course, smaller than the number of disks involved: although each drive is capable of transferring in parallel, there is no way to ensure that the requests are evenly distributed across the drives. Inevitably the load on one drive will be higher than on another.

The evenness of the load on the disks is strongly dependent on the way the data is shared across the drives. It's convenient to think of the disk storage as a large number of

data sectors which are addressable by number, rather like the pages in a book. The most obvious method is to divide the virtual disk into groups of consecutive sectors the size of the individual physical disks and store them in this manner, rather like taking a large book and tearing it into smaller sections. This method is called *concatenation* and has the advantage that the disks do not need to have any specific size relationships.

Concatenation does not correspond to any specific RAID level. It works well when the access to the virtual disk is spread evenly about its address space. When access is concentrated on a smaller area, the improvement is less marked. Figure 1 illustrates the sequence in which storage units are allocated in a concatenated organization.

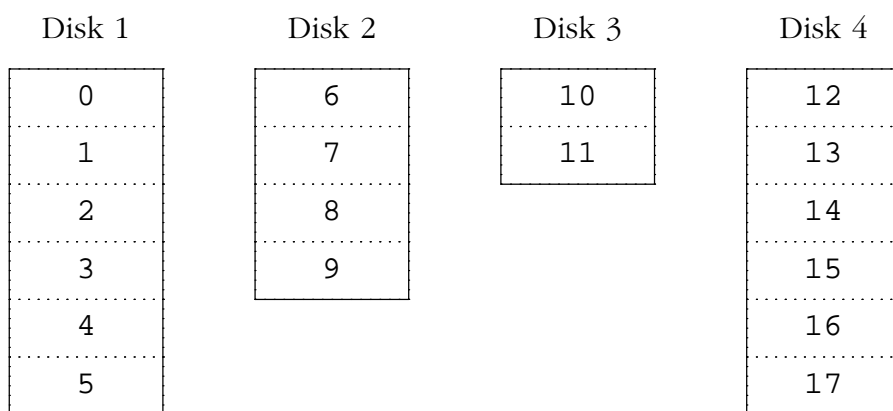| Disk 1 | Disk 2 | Disk 3 | Disk 4 |
|--------|--------|--------|--------|
| 0 | 6 | 10 | 12 |
| 1 | 7 | 11 | 13 |
| 2 | 8 | | 14 |
| 3 | 9 | | 15 |
| 4 | | | 16 |
| 5 | | | 17 |

Figure 1: Concatenated organization

An alternative mapping is to divide the address space into smaller, even-sized components and store them sequentially on different devices. For example, the first 256 sectors may be stored on the first disk, the next 256 sectors on the next disk and so on. After filling the last disk, the process repeats until the disks are full. This mapping is called *striping* or RAID-0, though the latter term is somewhat misleading: it provides no redundancy. Striping requires somewhat more effort to locate the data, and it can cause additional I/O load where a transfer is spread over multiple disks, but it can also provide a more constant load across the disks. Figure 2 illustrates the sequence in which storage units are allocated in a striped organization.

| Disk 1 | Disk 2 | Disk 3 | Disk 4 |
|:------:|:------:|:------:|:------:|
| 0 | 1 | 2 | 3 |
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 |
| 20 | 21 | 22 | 23 |

Figure 2: Striped organization

## Data integrity

Finally, Vinum offers improved reliability. Although disk drive reliability has increased tremendously over the last few years, they are still the most likely core component of a server to fail. When they do, the results can be catastrophic: replacing a failed disk drive and restoring data to it can take days.

The traditional way to approach this problem has been *mirroring*, keeping two copies of the data on different physical hardware. Since the advent of the RAID levels, this technique has also been called *RAID level 1* or *RAID-1*. Any write to the volume writes to both locations; a read can be satisfied from either, so if one drive fails, the data is still available on the other drive.

Mirroring has two problems:

- The price. It requires twice as much disk storage as a non-redundant solution.

- The performance impact. Writes must be performed to both drives, so they take up twice the bandwidth of a non-mirrored volume. Reads do not suffer from a performance penalty: it even looks as if they are faster. This issue will be discussed in more detail below.

## RAID-5

An alternative solution is *parity*, implemented in the RAID levels 2, 3, 4 and 5. Of these, RAID-5 is the most interesting. Vinum implements RAID-5 as a variant of striped plexes which dedicates one block of each stripe to parity of the other blocks. The location of this parity block changes from one stripe to the next. The numbers in the data blocks indicate the relative block numbers.

|   | Disk 1 | Disk 2 | Disk 3 | Disk 4 |
|---|--------|--------|--------|--------|
|   | 0      | 1      | 2      | Parity |
|   | 3      | 4      | Parity | 5      |
|   | 6      | Parity | 7      | 8      |
|   | Parity | 9      | 10     | 11     |
|   | 12     | 13     | 14     | Parity |
|   | 15     | 16     | Parity | 17     |

Figure 3: RAID-5 organization

Compared to mirroring, RAID-5 has the advantage of requiring significantly less storage space. Read access is similar to that of striped organizations, but write access is significantly slower, approximately 25% of the read performance. If one drive fails, the array can continue to operate in degraded mode: a read from one of the remaining accessible drives continues normally, but a read from the failed drive is recalculated from the corresponding block from all the remaining drives.

## RAID-4

Vinum also implements RAID-4, a variant of RAID-5 in which one disk contains all the parity information:

|   | Disk 1 | Disk 2 | Disk 3 | Disk 4 |
|---|--------|--------|--------|--------|
|   | 0      | 1      | 2      | Parity |
|   | 3      | 4      | 5      | Parity |
|   | 6      | 7      | 8      | Parity |
|   | 9      | 10     | 11     | Parity |
|   | 12     | 13     | 14     | Parity |
|   | 15     | 16     | 17     | Parity |

Figure 4: RAID-5 organization

The only advantage that RAID-4 offers over RAID-5 is marginally simpler code. As implemented in Vinum, it shares the code with the RAID-5 implementation. The only difference is the following code:

```
/* subdisk containing the parity stripe */
if (plex->organization == plex_raid5)
  m.psdno = plex->subdisks - 1
        - (*diskaddr / (plex->stripesize * (plex->subdisks - 1)))
        % plex->subdisks;
else                                          /* RAID-4 */
  m.psdno = plex->subdisks - 1;
```

On the other hand, RAID-4 has problems with the load on the parity subdisks when writing. The original motivation to implement RAID-4, which was added later, was to investigate whether improvements in drive cache hit rates might improve read performance and thus make up for the write performance issues. They did not: as will be discussed below, it is always a bad idea for an I/O request to span a drive, so the situation never occurred.

### RAID-3

RAID-3 is a very misunderstood RAID level. Many so-called implementations of RAID-3 are in fact implementations of RAID-4.

RAID-3 is like an implementation of RAID-4 with a stripe size of one byte. Each transfer involves all disks (with the exception of the parity disk for reads). Under certain circumstances (high-speed streaming data), this can be useful. Without spindle synchronization (where the corresponding sectors pass the heads of each drive at the same time), RAID-3 would be very inefficient. In a multiple-access system, it also causes high latency.

An argument for RAID-3 does exist where a single process requires very high data rates. With spindle synchronization, this would be a potentially useful addition to Vinum.

### RAID-2

RAID-2 uses two subdisks to store a Hamming code. Otherwise it resembles RAID-3. Compared to RAID-3, it offers a lower data density, higher CPU usage and no compensating advantages. As a result, it is effectively obsolete. Vinum does not implement it.

## Vinum plex organizations.

Vinum implements mirroring by attaching multiple plexes to a volume. Each plex is a representation of the data in a volume. A volume may contain between one and eight plexes.

Although a plex represents the complete data of a volume, it is possible for parts of the representation to be physically missing, either by design (by not defining a subdisk for parts of a concatenated plex) or by accident (as a result of the failure of a drive). As long as at least one plex can provide the data for the complete address range of the volume, the volume is fully functional.

Conceptually, RAID-4 and RAID-5 are used for redundancy, but in fact the implementation is a kind of striping. This poses problems for the implementation of Vinum: should it be a kind of plex or a kind of volume? It would have been possible to implement it either way, but it proved to be simpler to implement RAID-4 and RAID-5 as a plex type. This means that there are two different ways of ensuring data redundancy: either have more than one plex in a volume, or have a single RAID-4 or RAID-5 plex. These methods can be combined.

Vinum implements both concatenation and striping at the plex level:

- A *concatenated plex* uses the address space of each subdisk in turn.

- A *striped plex* stripes the data across each subdisk. The subdisks must all have the same size, and there must be at least two subdisks in order to distinguish it from a concatenated plex.

- Like a striped plex, *RAID-4* and *RAID-5 plexes* stripe the data across each subdisk. The subdisks must all have the same size, and there must be at least three subdisks, since otherwise it would be more efficient to use mirroring.


## Which plex organization?

Each of plex organizations has its own advantages:

- Concatenated plexes are the most flexible: they can contain any number of subdisks, and the subdisks may be of different length. The plex may be extended by adding additional subdisks. They require less CPU time than striped, RAID-4 or RAID-5 plexes, though the difference in CPU overhead from striped plexes is not measurable. On the other hand, they are most susceptible to hot spots, where one disk is very active and others are idle.

- The greatest advantage of striped (RAID-0) plexes is that they reduce hot spots, at least in theory: by choosing an optimum sized stripe (empirically determined to be in the order of 256 kB), the load on the component drives can be made more even. In practice, since UFS also spreads the load across cylinder groups, it's not clear how much difference this advantage makes in practice. The disadvantages of this approach are (fractionally) more complex code and restrictions on subdisks: they must be all the same size. Extending a plex by adding new subdisks requires complete reorganization of the data in the subdisks. The issues regarding error recovery in this reorganization are so complicated that Vinum currently does not implement it. Vinum

imposes an additional, trivial restriction: a striped plex must have at least two subdisks, since otherwise it is indistinguishable from a concatenated plex.

- RAID-4 and RAID-5 plexes are effectively an extension of striped plexes. Compared to striped plexes, they offer the advantage of fault tolerance, but the disadvantages of higher storage cost and significantly higher CPU overhead, particularly for writes. The code is an order of magnitude more complex than for concatenated and striped plexes. Like striped plexes, RAID-4 and RAID-5 plexes must have equal-sized subdisks and cannot currently be extended. Vinum enforces a minimum of three subdisks for a RAID-4 or RAID-5 plex, since any smaller number would not make any sense.

  In comparison with RAID-5 plexes, RAID-4 plexes have no advantage. They should not be used.

Table 5 summarizes the advantages and disadvantages of each plex organization.

| Plex type | Minimum subdisks | can add subdisks | must be equal size | application |
|---|---|---|---|---|
| concatenated | 1 | yes | no | Large data storage with maximum placement flexibility and moderate performance. |
| striped | 2 | no | yes | High performance in combination with highly concurrent access. |
| RAID-4, RAID-5 | 3 | no | yes | Highly reliable storage, primarily read access. |

Figure 5: Vinum plex organizations

# Vinum features

Beyond object organization, Vinum includes a number of other features:

- Online configuration via the *vinum* utility program.

- Automatic error detection and recovery where possible.

- State information for each object. This enables Vinum to function correctly even if some objects are not accessible.

- Persistent configuration. Each Vinum drive stores two copies of the configuration, so the system can start up automatically. The configuration includes state information, so any degraded objects will remain so over a reboot, or even when moved to a new system.

- Support for Vinum root file systems.

- Online rebuild of objects.

# Vinum configuration

Vinum maintains a *configuration database* which describes the objects known to an individual system. Initially, the user creates the configuration database from one or more configuration files with the aid of the *vinum(8)* utility program. Vinum stores a copy of its configuration database on each drive (disk slice) under its control. This database is updated on each state change, so that a restart accurately restores the state of each Vinum object.

### The configuration file

The configuration file describes individual Vinum objects. The definition of a simple volume might be:

```
drive a device /dev/da0s4h
volume myvol
  plex org concat
    sd length 1g drive a
```

This file describes a four Vinum objects:

- The `drive` line describes a disk partition (*drive*) and its location relative to the underlying hardware. It is given the symbolic name *a*. This separation of the symbolic names from the device names allows disks to be moved from one location to another without confusion.

- The `volume` line describes a volume. The only required attribute is the name, in this case `myvol`.

- The `plex` line defines a plex. The only required parameter is the organization, in this case `concat`. No name is necessary: the system automatically generates a name from the volume name by adding the suffix `.p`$x$, where $x$ is the number of the plex in the volume. Thus this plex will be called *myvol.p0.*

- The `sd` line describes a subdisk. The minimum specifications are the name of a drive on which to store it, and the length of the subdisk. As with plexes, no name is necessary: the system automatically assigns names derived from the plex name by adding the suffix `.s`$x$, where $x$ is the number of the subdisk in the plex. Thus Vinum gives this subdisk the name *myvol.p0.s0*

The file is indented for legibility; it is not necessary for its correct interpretation.

After processing this file, *vinum(8)* produces the following output:

```
vinum -> create config1
1 drives:
D a                     State: up       /dev/da0s4h     A: 3070/4094 MB (74%)

1 volumes:
V myvol                 State: up       Plexes:       1 Size:      1024 MB

1 plexes:
P myvol.p0            C State: up       Subdisks:     1 Size:      1024 MB

1 subdisks:
S myvol.p0.s0           State: up       D: a            Size:      1024 MB
```

This output shows the brief listing format of *vinum(8)*. It is represented graphically in Figure 6.

This figure, and the ones which follow, represent a volume, which contains the plexes, which in turn contain the subdisks. In this trivial example, the volume contains one plex, and the plex contains one subdisk.

This particular volume has no specific advantage over a conventional disk partition. It contains a single plex, so it is not redundant. The plex contains a single subdisk, so there is no difference in storage allocation from a conventional disk partition. The following sections illustrate various more interesting configuration methods.
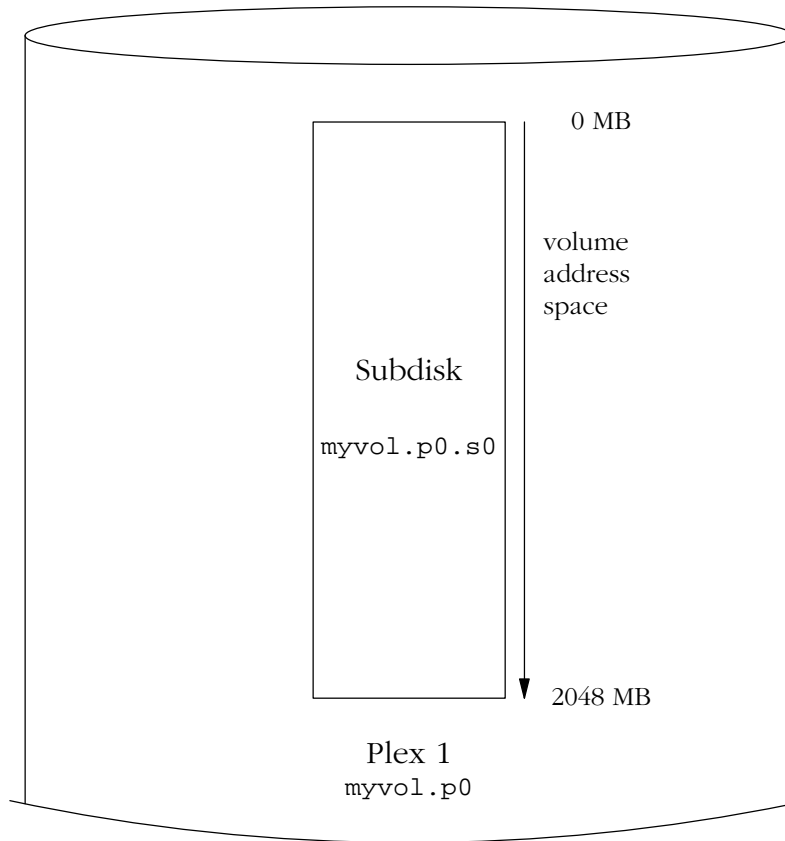
10

Figure 6: A simple Vinum volume

## Increased resilience: mirroring

The resilience of a volume can be increased either by mirroring or by using RAID-4 or RAID-5 plexes. When laying out a mirrored volume, it is important to ensure that the subdisks of each plex are on different drives, so that a drive failure will not take down both plexes. The following configuration mirrors a volume:

```
drive b device /dev/da1s4h
volume mirror
  plex org concat
    sd length 1g drive a
  plex org concat
    sd length 1g drive b
```

In this example, it was not necessary to specify a definition of drive *a* again, since Vinum keeps track of all objects in its configuration database. After processing this definition, the configuration looks like:

```
2 drives:
D a                      State: up        /dev/da0s4h      A: 2046/4094 MB (49%)
D b                      State: up        /dev/da1s4h      A: 3070/4094 MB (74%)

2 volumes:
V myvol                  State: up        Plexes:       1 Size:        1024 MB
V mirror                 State: up        Plexes:       2 Size:        1024 MB

3 plexes:
P myvol.p0        C State: up        Subdisks:     1 Size:        1024 MB
P mirror.p0       C State: up        Subdisks:     1 Size:        1024 MB
P mirror.p1       C State: faulty    Subdisks:     1 Size:        1024 MB

3 subdisks:
S myvol.p0.s0            State: up        D: a              Size:        1024 MB
S mirror.p0.s0           State: up        D: a              Size:        1024 MB
S mirror.p1.s0           State: empty     D: b              Size:        1024 MB
```

Figure 7 shows the structure graphically. In this example, each plex contains the full 512 MB of address space. As in the previous example, each plex contains only a single subdisk.
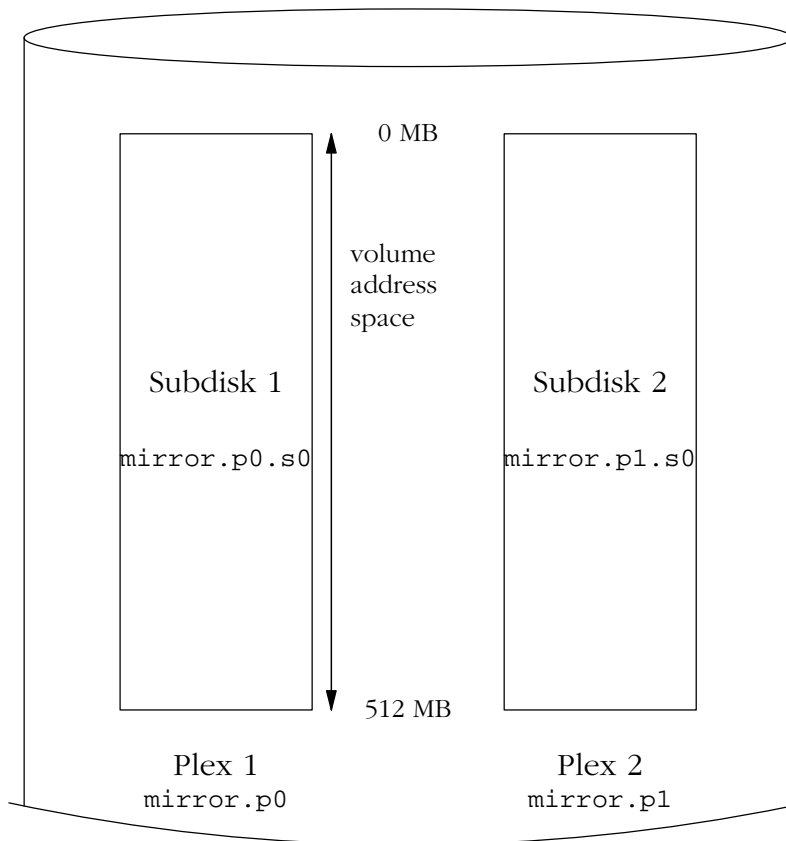


**Figure 7**: A mirrored Vinum volume

When creating a mirrored volume, Vinum should ensure that the contents of each plex are identical. This requires copying of data in some form, which can take a long time. Until this has occurred, the state of the second plex is set to `faulty`, and the corresponding subdisk is `empty`.

In practice, people are seldom interested in the contents of newly created plexes. Normally, data read from the volume has first been written to the volume, and Vinum ensures that this data is consistent across all the plexes. As a result, it is possible to ignore the inconsistencies which exist only in those parts of the plex which will never be read. The keyword `setupstate` tells to Vinum set all plexes and subdisks to the `up` state when creating mirrored volumes:

```
volume mirror1 setupstate
  plex org concat
    sd length 1g drive a
  plex org concat
    sd length 1g drive b
```

## Optimizing performance

The mirrored volume in the previous example is more resistant to failure than an unmirrored volume, but its performance is less: each write to the volume requires a write to both drives, using up a greater proportion of the total disk bandwidth. Performance considerations demand a different approach: instead of mirroring, the data is striped across as many disk drives as possible. The following configuration shows a volume with a plex striped across four disk drives:

```
drive c device /dev/da2s4h
drive d device /dev/da3s4h
volume stripe
  plex org striped 496k
    sd length 256m drive a
    sd length 256m drive b
    sd length 256m drive c
    sd length 256m drive d
```

As before, it is not necessary to define the drives which are already known to Vinum. After processing this definition, the configuration looks like:

```
4 drives:
D a                     State: up       /dev/da0s4h     A: 1790/4094 MB (43%)
D b                     State: up       /dev/da1s4h     A: 2814/4094 MB (68%)
D c                     State: up       /dev/da2s4h     A: 3838/4094 MB (93%)
D d                     State: up       /dev/da3s4h     A: 3838/4094 MB (93%)

3 volumes:
V myvol                 State: up       Plexes:       1 Size:        1024 MB
V mirror                State: up       Plexes:       2 Size:        1024 MB
V stripe                State: up       Plexes:       1 Size:        1023 MB
```

```
4 plexes:
P myvol.p0           C State: up       Subdisks:   1 Size:      1024 MB
P mirror.p0          C State: up       Subdisks:   1 Size:      1024 MB
P mirror.p1          C State: faulty   Subdisks:   1 Size:      1024 MB
P stripe.p0          S State: up       Subdisks:   4 Size:      1023 MB

7 subdisks:
S myvol.p0.s0          State: up       D: a          Size:      1024 MB
S mirror.p0.s0         State: up       D: a          Size:      1024 MB
S mirror.p1.s0         State: empty    D: b          Size:      1024 MB
S stripe.p0.s0         State: up       D: a          Size:       255 MB
S stripe.p0.s1         State: up       D: b          Size:       255 MB
S stripe.p0.s2         State: up       D: c          Size:       255 MB
S stripe.p0.s3         State: up       D: d          Size:       255 MB
```
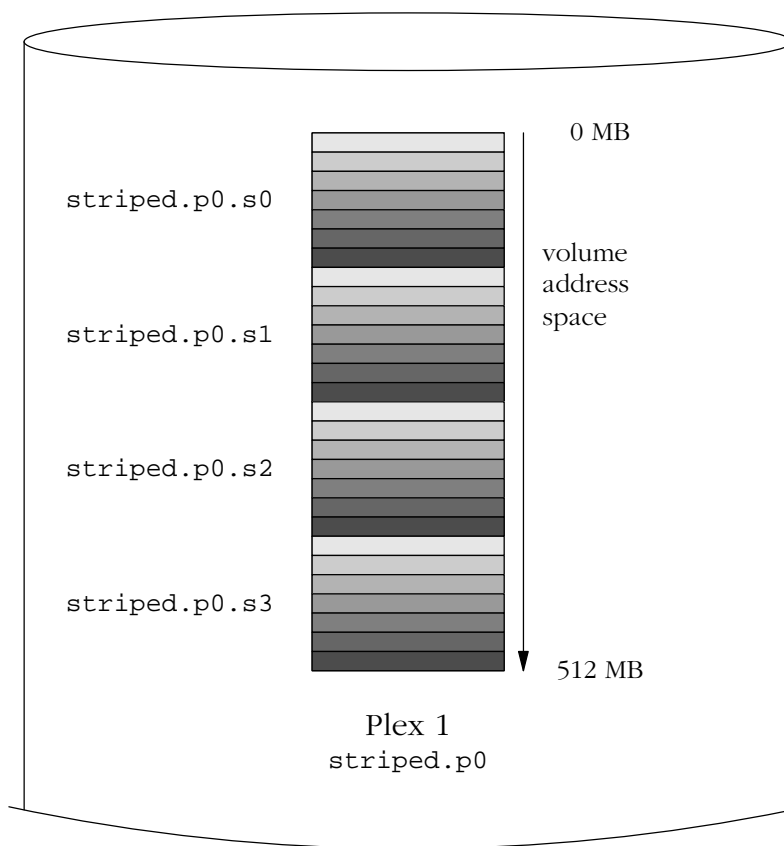


**Figure 8: A striped Vinum volume**

This volume is represented in Figure 8. The darkness of the stripes indicates the position within the plex address space: the lightest stripes come first, the darkest last.

14

# Increased resilience: RAID-5

The alternative approach to resilience is RAID-5. The following discussion also applies to RAID-4, but as discussed above, the use of RAID-4 is not recommended. A RAID-5 configuration might look like:

```
drive e device /dev/da4s4h
volume raid5
  plex org raid5 496k
    sd length 256m drive a
    sd length 256m drive b
    sd length 256m drive c
    sd length 256m drive d
    sd length 256m drive e
```

Although this plex has five subdisks, its size is the same as the plexes in the other examples, since the equivalent of one subdisk is used to store parity information. After processing the configuration, the system configuration is:

```
5 drives:
D a                     State: up        /dev/da0s4h    A: 1535/4094 MB (37%)
D b                     State: up        /dev/da1s4h    A: 2559/4094 MB (62%)
D c                     State: up        /dev/da2s4h    A: 3583/4094 MB (87%)
D d                     State: up        /dev/da3s4h    A: 3583/4094 MB (87%)
D e                     State: up        /dev/da4s4h    A: 3838/4094 MB (93%)

4 volumes:
V myvol                 State: up        Plexes:      1 Size:        1024 MB
V mirror                State: up        Plexes:      2 Size:        1024 MB
V stripe                State: up        Plexes:      1 Size:        1023 MB
V raid5                 State: down      Plexes:      1 Size:        1023 MB

5 plexes:
P myvol.p0         C State: up           Subdisks:    1 Size:        1024 MB
P mirror.p0        C State: up           Subdisks:    1 Size:        1024 MB
P mirror.p1        C State: faulty       Subdisks:    1 Size:        1024 MB
P stripe.p0        S State: up           Subdisks:    4 Size:        1023 MB
P raid5.p0        R5 State: init         Subdisks:    5 Size:        1023 MB

12 subdisks:
S myvol.p0.s0           State: up        D: a           Size:        1024 MB
S mirror.p0.s0          State: up        D: a           Size:        1024 MB
S mirror.p1.s0          State: empty     D: b           Size:        1024 MB
S stripe.p0.s0          State: up        D: a           Size:         255 MB
S stripe.p0.s1          State: up        D: b           Size:         255 MB
S stripe.p0.s2          State: up        D: c           Size:         255 MB
S stripe.p0.s3          State: up        D: d           Size:         255 MB
S raid5.p0.s0           State: empty     D: a           Size:         255 MB
S raid5.p0.s1           State: empty     D: b           Size:         255 MB
S raid5.p0.s2           State: empty     D: c           Size:         255 MB
S raid5.p0.s3           State: empty     D: d           Size:         255 MB
S raid5.p0.s4           State: empty     D: e           Size:         255 MB
```
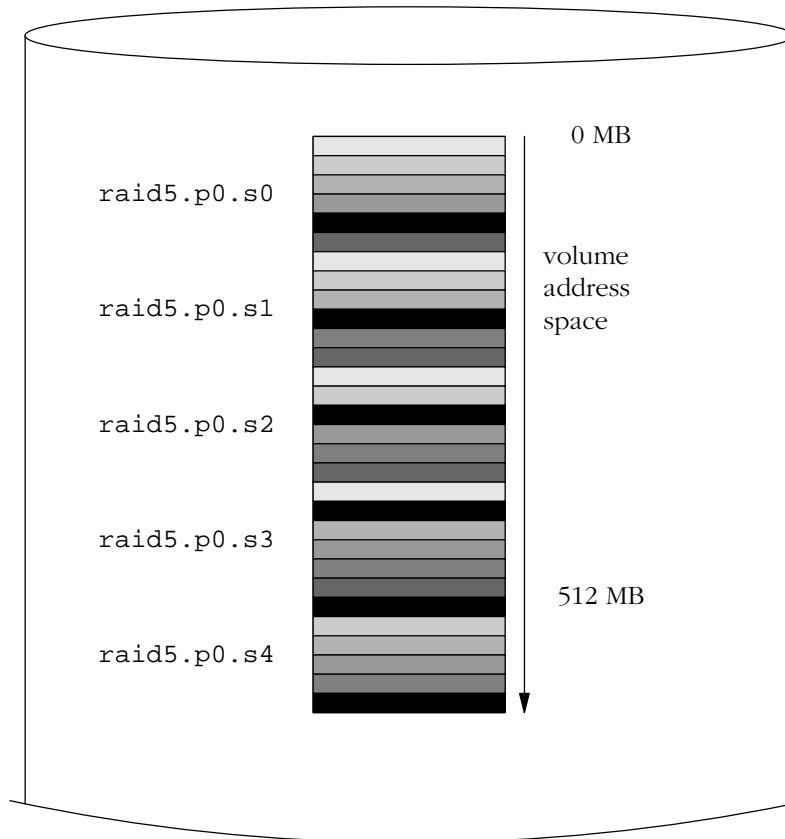
Figure 9 represents this volume graphically:



**Figure 9: A RAID-5 Vinum volume**

As with striped plexes, the darkness of the stripes indicates the position within the plex address space: the lightest stripes come first, the darkest last. The completely black stripes are the parity stripes.

On creation, RAID-5 plexes are in the *init* state: before they can be used, the parity data must be created. Vinum currently initializes RAID-5 plexes by writing binary zeros to all subdisks, though a probably future alternative is to rebuild the parity blocks, which allows better recovery of crashed plexes.

16

## Resilience and performance

With sufficient hardware, it is possible to build volumes which show both increased resilience and increased performance compared to standard UNIX partitions. Mirrored disks will always give better performance than RAID-5, so a typical configuration file might be:

```
volume raid10
  plex org striped 496k
    sd length 102480k drive a
    sd length 102480k drive b
    sd length 102480k drive c
    sd length 102480k drive d
    sd length 102480k drive e
  plex org striped 496k
    sd length 102480k drive c
    sd length 102480k drive d
    sd length 102480k drive e
    sd length 102480k drive a
    sd length 102480k drive b
```
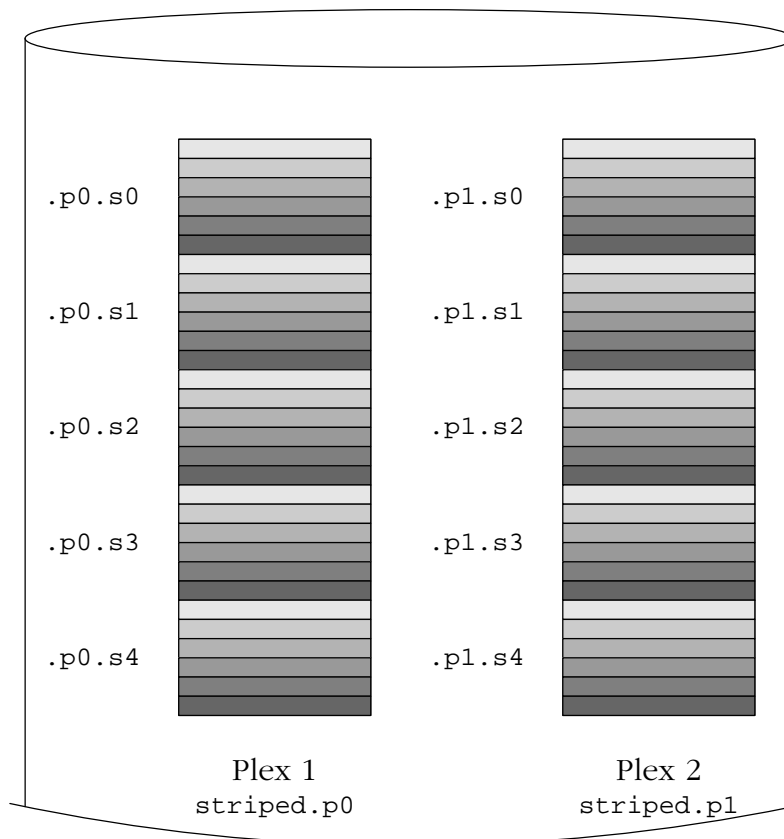


**Figure 10: A mirrored, striped Vinum volume**

Figure 10 represents the structure of this volume. The subdisks of the second plex are offset by two drives from those of the first plex: this helps ensure that writes do not go to the same subdisks even if a transfer goes over two drives.

After creating this volume, the configuration is:

```
5 drives:
D a                         State: up        /dev/da0s4h      A: 1335/4094 MB (32%)
D b                         State: up        /dev/da1s4h      A: 2359/4094 MB (57%)
D c                         State: up        /dev/da2s4h      A: 3383/4094 MB (82%)
D d                         State: up        /dev/da3s4h      A: 3383/4094 MB (82%)
D e                         State: up        /dev/da4s4h      A: 3639/4094 MB (88%)

5 volumes:
V myvol                     State: up        Plexes:       1 Size:        1024 MB
V mirror                    State: up        Plexes:       2 Size:        1024 MB
V stripe                    State: up        Plexes:       1 Size:        1023 MB
V raid5                     State: down      Plexes:       1 Size:        1023 MB
V raid10                    State: up        Plexes:       2 Size:         498 MB

7 plexes:
P myvol.p0         C State: up        Subdisks:      1 Size:        1024 MB
P mirror.p0        C State: up        Subdisks:      1 Size:        1024 MB
P mirror.p1        C State: faulty    Subdisks:      1 Size:        1024 MB
P stripe.p0        S State: up        Subdisks:      4 Size:        1023 MB
P raid5.p0        R5 State: init      Subdisks:      5 Size:        1023 MB
P raid10.p0        S State: up        Subdisks:      5 Size:         498 MB
P raid10.p1        S State: faulty    Subdisks:      5 Size:         498 MB

22 subdisks:
S myvol.p0.s0               State: up        D: a             Size:        1024 MB
S mirror.p0.s0             State: up        D: a             Size:        1024 MB
S mirror.p1.s0             State: empty     D: b             Size:        1024 MB
S stripe.p0.s0             State: up        D: a             Size:         255 MB
S stripe.p0.s1             State: up        D: b             Size:         255 MB
S stripe.p0.s2             State: up        D: c             Size:         255 MB
S stripe.p0.s3             State: up        D: d             Size:         255 MB
S raid5.p0.s0              State: empty     D: a             Size:         255 MB
S raid5.p0.s1              State: empty     D: b             Size:         255 MB
S raid5.p0.s2              State: empty     D: c             Size:         255 MB
S raid5.p0.s3              State: empty     D: d             Size:         255 MB
S raid5.p0.s4              State: empty     D: e             Size:         255 MB
S raid10.p0.s0             State: up        D: a             Size:          99 MB
S raid10.p0.s1             State: up        D: b             Size:          99 MB
S raid10.p0.s2             State: up        D: c             Size:          99 MB
S raid10.p0.s3             State: up        D: d             Size:          99 MB
S raid10.p0.s4             State: up        D: e             Size:          99 MB
S raid10.p1.s0             State: empty     D: c             Size:          99 MB
S raid10.p1.s1             State: empty     D: d             Size:          99 MB
S raid10.p1.s2             State: empty     D: e             Size:          99 MB
S raid10.p1.s3             State: empty     D: a             Size:          99 MB
S raid10.p1.s4             State: empty     D: b             Size:          99 MB
```

# Object naming

As described above, Vinum assigns default names to plexes and subdisks, although they may be overridden. Overriding the default names is not recommended: experience with the VERITAS volume manager, which allows arbitary naming of objects, has shown that this flexibility does not bring a significant advantage, and it can cause confusion.

Names may contain any non-blank character, but it is recommended to restrict them to letters, digits and the underscore characters. The names of volumes, plexes and subdisks may be up to 64 characters long, and the names of drives may up to 32 characters long.

Vinum objects are assigned device nodes in the hierarchy */dev/vinum*. The configuration shown above would cause Vinum to create the following device nodes:

- The control devices */dev/vinum/control* and */dev/vinum/controld*, which are used by *vinum(8)* and the Vinum daemon respectively.

- Character device entries for each volume. These are the main devices used by Vinum. The configuration above would include the devices */dev/vinum/myvol*, */dev/vinum/mirror, /dev/vinum/striped, /dev/vinum/raid5* and */dev/vinum/raid10*.

- The directories */dev/vinum/plex* and */dev/vinum/sd* which contain device nodes for each plex and subdisk.

For example, consider the following configuration file:

```
drive drive1 device /dev/da1s4h
drive drive2 device /dev/da2s4h
drive drive3 device /dev/da3s4h
drive drive4 device /dev/da4s4h
volume s64
 plex org striped 496k
    sd length 100m drive drive1
    sd length 100m drive drive2
    sd length 100m drive drive3
    sd length 100m drive drive4
```

After processing this file, *vinum(8)* creates the following structure in */dev/vinum*:

```
total 1
crw-------  1 root   wheel       91, 0x3fff00fe Dec 28 12:59 control
crw-------  1 root   wheel       91, 0x3fff00ff Dec 28 12:59 controld
dr-xr-xr-x  2 root   wheel          512 Dec 28 13:05 plex
crw-r-----  1 root   operator    91,    0 Dec 28 13:17 s64
dr-xr-xr-x  2 root   wheel          512 Dec 28 13:05 sd

/dev/vinum/plex:
total 0
crw-r-----  1 root   operator    91, 0x40000000 Dec 28 13:17 s64.p0

/dev/vinum/sd:
total 0
crw-r-----  1 root   operator    91, 0x80000000 Dec 28 13:17 s64.p0.s0
```

```
crw-r-----  1 root   operator   91, 0x80000001 Dec 28 13:17 s64.p0.s1
crw-r-----  1 root   operator   91, 0x80000002 Dec 28 13:17 s64.p0.s2
crw-r-----  1 root   operator   91, 0x80000003 Dec 28 13:17 s64.p0.s3
```

Although it is recommended that plexes and subdisks should not be allocated specific names, Vinum drives must have logical names separate from their device name. This makes it possible to move a drive to a different location and still recognize it automatically.

### Creating file systems

Volumes appear to the system to be identical to disks, with one exception: Vinum volumes do not contain a partition table. This has required modification to some disk utilities, notably *newfs*, which previously tried to interpret the last letter of a Vinum volume name as a partition identifier. For example, a disk drive may have a name like */dev/wd0a* or */dev/da2h*. These names represent the first partition (`a`) on the first (0) IDE disk (`wd`) and the eighth partition (`h`) on the third (2) SCSI disk (`da`) respectively. By contrast, a Vinum volume might be called */dev/vinum/concat*, a name which has no relationship with a partition name. Earlier versions of *newfs* interpreted the last character of the file name as the partition identifier and failed if the letter was not a valid identifier.

# Startup

The configuration information which Vinum stores the drives has the same form as in the configuration files. When reading from the configuration database, Vinum recognizes a number of keywords which are not allowed in the configuration files. For example, a disk configuration might contain to following text, which has been wrapped to fit on the page:

```
volume myvol state up
volume bigraid state down
plex name myvol.p0 state up org concat vol myvol
plex name myvol.p1 state up org concat vol myvol
plex name myvol.p2 state init org striped 512b vol myvol
plex name bigraid.p0 state initializing org raid5 512b vol bigraid
sd name myvol.p0.s0 drive a plex myvol.p0 state up len 1048576b driveoffset 265b plex
 offset 0b
sd name myvol.p0.s1 drive b plex myvol.p0 state up len 1048576b driveoffset 265b plex
 offset 1048576b
sd name myvol.p1.s0 drive c plex myvol.p1 state up len 1048576b driveoffset 265b plex
 offset 0b
sd name myvol.p1.s1 drive d plex myvol.p1 state up len 1048576b driveoffset 265b plex
 offset 1048576b
sd name myvol.p2.s0 drive a plex myvol.p2 state init len 524288b driveoffset 1048841b
 plexoffset 0b
sd name myvol.p2.s1 drive b plex myvol.p2 state init len 524288b driveoffset 1048841b
 plexoffset 524288b
sd name myvol.p2.s2 drive c plex myvol.p2 state init len 524288b driveoffset 1048841b
 plexoffset 1048576b
```

```
sd name myvol.p2.s3 drive d plex myvol.p2 state init len 524288b driveoffset 1048841b
 plexoffset 1572864b
sd name bigraid.p0.s0 drive a plex bigraid.p0 state initializing len 4194304b driveof
fset 1573129b plexoffset 0b
sd name bigraid.p0.s1 drive b plex bigraid.p0 state initializing len 4194304b driveof
fset 1573129b plexoffset 4194304b
sd name bigraid.p0.s2 drive c plex bigraid.p0 state initializing len 4194304b driveof
fset 1573129b plexoffset 8388608b
sd name bigraid.p0.s3 drive d plex bigraid.p0 state initializing len 4194304b driveof
fset 1573129b plexoffset 12582912b
sd name bigraid.p0.s4 drive e plex bigraid.p0 state initializing len 4194304b driveof
fset 1573129b plexoffset 16777216b
```

The obvious differences here are the presence of explicit location information and naming (both of which are also allowed, but discouraged, for use by the user) and the information on the states (which are not available to the user). Vinum does not store information about drives in the configuration information: it finds the drives by scanning the configured disk drives for partitions with a Vinum label. This enables Vinum to identify drives correctly even if they have been assigned different drive IDs.

At system startup, Vinum reads the configuration database from each of the Vinum drives in reverse order of last modification time. Under normal circumstances, each drive contains an identical copy of the configuration database, so it does not matter which drive is read. After a crash, however, it is possible that the drives contain inconsistent configuration information. After reading the most recent configuration, Vinum adds only those objects which were not previously defined, thus avoiding changing the state of existing objects.

# Performance issues

Performance measurements show that the performance is very close to what could be expected from the underlying disk driver performing the same operations as Vinum performs: in other words, the overhead of Vinum itself is negligible. This does not mean that Vinum has perfect performance: the choice of requests has a strong impact on the overall subsystem performance, and there are some known areas which could be improved upon. In addition, the user can influence performance by the design of the volumes.

The following sections examine some factors which influence performance.

> Note: The performance measurements in this section were done on old disk drives. The absolute performance is correspondingly significantly poorer than that of modern drives. The intention of the following graphs is to show relative performance, not absolute performance. Other tests indicate that the performance relationships also apply to modern high-end hardware.

## The influence of stripe size

In striped plexes, including RAID-4 and RAID-5 implementations, the stripe size has a significant influence on performance. In all plex structures except a single-subdisk plex (which by definition is concatenated), the possibility exists that a single transfer to or from a volume will be remapped into more than one physical I/O request. This is never desirable, since the average latency for multiple transfers is always larger than the average latency for single transfers to the same kind of disk hardware. Spindle synchronization does not help here, since there is no deterministic relationship between the positions of the data blocks on the different disks. Within the bounds of the current BSD I/O architecture (maximum transfer size 128 kB) and current disk hardware, this increase in latency usually offsets any speed increase in the transfer.

In the case of a concatenated plex, this remapping occurs only when a request overlaps a subdisk boundary. In a striped or RAID-5 plex, however, the probability is an inverse function of the stripe size. For this reason, a stripe size between 256 kB and 512 kB appears to be optimum: it is small enough to create a relatively random mapping of file system hot spots to individual disks, and large enough to ensure than 95% of all transfers involve only a single data subdisk.
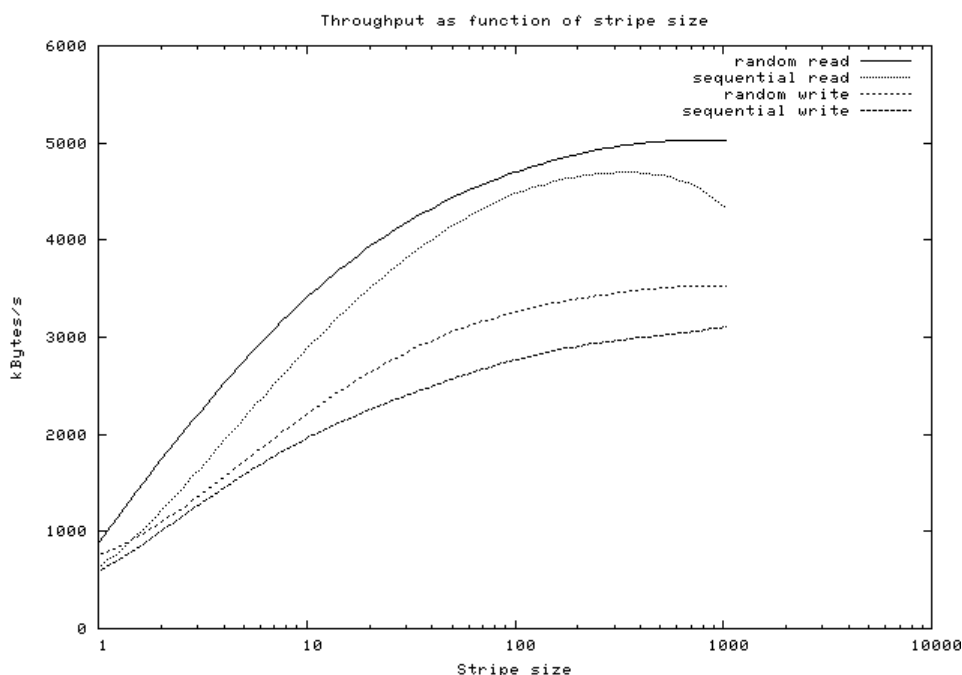


Figure 11: The influence of stripe size and mirroring

22

Figure 11 shows the effect of stripe size on read and write performance, obtained with *rawio*. This measurement used eight concurrent processes to access volumes with striped plexes with different stripe sizes. The graph shows the disadvantage of small stripe sizes, which can cause a significant performance degradation even compared to a single disk.

Other factors which influence the choice of stripe size are:

- The file system structure. UFS structures file systems into fixed sized cylinder groups with an inode table at the beginning of the cylinder group. It is easy for the relationship between cylinder group and stripe size to place the beginning of each cylinder group on the same subdisk. This occurs when both the stripe size and cylinder group size are a power of 2.

- The file system block size. While UFS does not always transfer entire blocks, there is a tendency for it to do so.

- The disk cache size. This factor remains to be investigated.

For these reasons, the current recommendation is to have a stripe size which is a power of two plus or minus a block size; 496 kB is a good example for a file system with the default block size of 16 kB.

## The influence of RAID-1 mirroring

Mirroring has different effects on read and write throughput. A write to a mirrored volume causes writes to each plex, so write performance is less than for a non-mirrored volume. A read from a mirrored volume, however, reads from only one plex, so read performance can improve.

There are two different scenarios for these performance changes, depending on the layout of the subdisks comprising the volume. Two basic possiblities exist for a mirrored, striped plex.

### One disk per subdisk

The optimum layout, both for reliability and for performance, is to have each subdisk on a separate disk. An example might be the following configuration, similar to the configuration on page 17:

```
volume raid10
  plex org striped 512k
    sd length 102480k drive a
    sd length 102480k drive b
    sd length 102480k drive c
    sd length 102480k drive d
  plex org striped 512k
    sd length 102480k drive e
    sd length 102480k drive f
    sd length 102480k drive g
    sd length 102480k drive h
```

In this case, the volume is spread over a total of eight disks. This has the following effects:

- Read access: by default, read accesses will alternate across the two plexes, giving a performance improvement close to 100%.

- Write access: writes must be performed to both disks, doubling the bandwidth requirement. Since the available bandwidth is also double, there should be little difference in througput.

At present, due to lack of hardware, no tests have been made of this configuration.

### Both plexes on the same disks

An alternative layout is to spread the subdisks of each plex over the same disks:

```
volume raid10
  plex org striped 512k
    sd length 102480k drive a
    sd length 102480k drive b
    sd length 102480k drive c
    sd length 102480k drive d
  plex org striped 512k
    sd length 102480k drive c
    sd length 102480k drive d
    sd length 102480k drive a
    sd length 102480k drive b
```

This has the following effects:

- Read access: by default, read accesses will alternate across the two plexes. Since there is no increase in bandwidth, there will be little difference in performance through the second plex.

- Write access: writes must be performed to both disks, doubling the bandwidth requirement. In this case, the bandwidth has not increase, so write throughput will decrease by approximately 50%.

Figure 11 also shows the effect of mirroring in this manner. The results are very close to the theoretical predictions.

### The influence of request size

As seen above, the throughput of a disk subsystem is the sum of the latency (the time taken to position the disk hardware over the correct part of the disk) and the time to transfer the data to or from the disk. Since latency is independent of transfer size, overall throughput is strongly dependent on the size of the transfer, as Figure 12 shows. Unfortunately, there is little that can be done to influence the transfer size. In FreeBSD, it tends to be closer to 10 kB than to 30 kB.

**Figure 12: Throughput as function of transfer size**

## The influence of concurrency

Vinum aims to give best performance for a large number of concurrent processes performing random access on a volume. Figure 13 shows the relationship between number of processes and throughput for a raw disk volume and a Vinum volume striped over four such disks with between one and 128 concurrent processes with an average transfer size of 16 kB. The actual transfers varied between 512 bytes and 32 kB, which roughly corresponds to UFS usage.

**Figure 13:** Concurrent random access with 32 sector transfers

## The influence of request structure

For concatenated and striped plexes, Vinum creates request structures which map directly to the user-level request buffers. The only additional overhead is the allocation of the request structure, and the possibility of improvement is correspondingly small.

With RAID-5 plexes, the picture is very different. The strategic choices described above work well when the total request size is less than the stripe width. By contrast, consider the following transfer of 32.5 kB:

| Offset | Subdisk 1 | Subdisk 2 | Subdisk 3 | Subdisk 4 | Subdisk 5 |
|---|---|---|---|---|---|
| 0x0000 | 0x0000 | 0x1000 | 0x2000 | 0x3000 | Parity |
| 0x1000 | 0x4000 | 0x5000 | 0x6000 | Parity | 0x7000 |
| 0x2000 | 0x8000 | 0x9000 | Parity | 0xa000 | 0xb000 |
| 0x3000 | 0xc000 | Parity | 0xd000 | 0xe000 | 0xf000 |
|  | Parity | 0x10000 | 0x11000 | 0x12000 | 0x13000 |

☐ Parity block
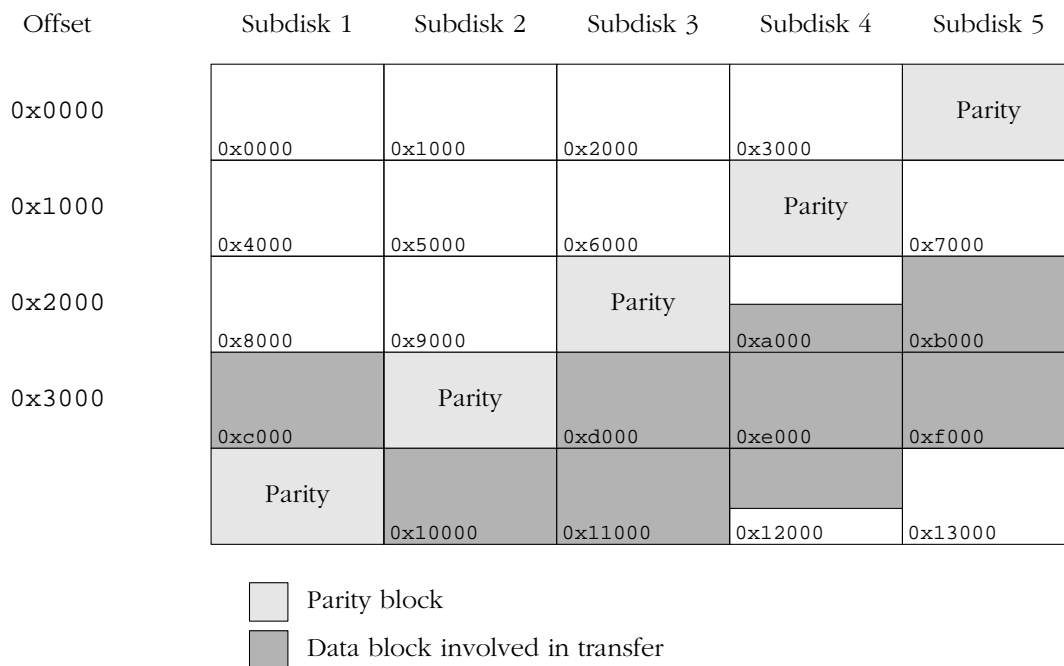▨ Data block involved in transfer

**Figure 14: A sample RAID-5 transfer**

An optimum approach to reading this data performs a total of 5 I/O operations, one on each subdisk. By contrast, Vinum treats this transfer as three separate transfers, one per stripe, and thus performs a total of 9 I/O transfers.

In practice, this inefficiency should not cause any problems: as discussed above, the optimum stripe size is larger than the maximum transfer size, so this situation does not arise when an appropriate stripe size is chosen.

These considerations are shown in figure 15, which clearly shows the RAID-5 tradeoffs:

- The RAID-5 write throughput is approximately half of the RAID-1 throughput in figure 11, and one-quarter of the write throughput of a striped plex.

- The read throughput is similar to that of striped volume of the same size.

Although the random access performance increases continually with increasing stripe size, the sequential access performance peaks at about 20 kB for writes and 35 kB for reads. This effect has not yet been adequately explained, but may be due to the nature of the test (8 concurrent processes writing the same data at the same time).

Figure 15: RAID-5 performance against stripe size

# The implementation

The implementation of Vinum required a number of tradeoffs. This section looks at some of the more interesting ones.

### Where the driver fits

To the operating system, Vinum looks like a block device, so it is normally be accessed as a block device. Instead of operating directly on the device, it creates new requests and passes them to the real device drivers. Conceptually it could pass them to other Vinum devices, though this usage makes no sense and would probably cause problems. The following figure  is © 1996 Addison-Wesley, and is reproduced with permission. It shows the standard 4.4BSD I/O structure:

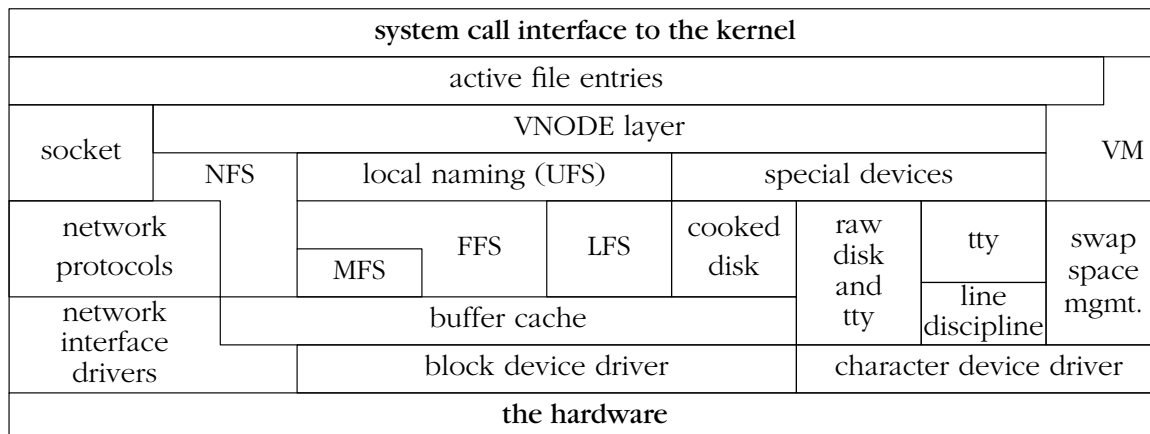| system call interface to the kernel | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| active file entries | | | | | | | | |
| socket | VNODE layer | | | | | | | VM |
| | NFS | local naming (UFS) | | | special devices | | | |
| network protocols | | MFS | FFS | LFS | cooked disk | raw disk and tty | tty | swap space mgmt. |
| | | | | | | | line discipline | |
| network interface drivers | buffer cache | | | | | | | |
| | block device driver | | | | character device driver | | | |
| the hardware | | | | | | | | |

Figure 16: Kernel I/O structure, after McKusick et. al.

With Vinum and some other changes, it becomes:

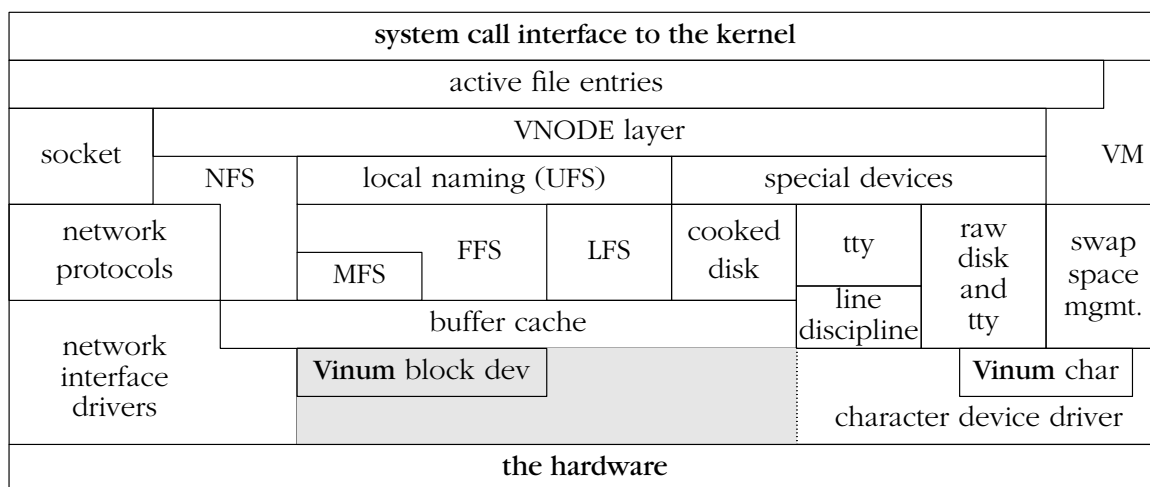| system call interface to the kernel | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| active file entries | | | | | | | | |
| socket | VNODE layer | | | | | | | VM |
| | NFS | local naming (UFS) | | | special devices | | | |
| network protocols | | MFS | FFS | LFS | cooked disk | tty | raw disk and tty | swap space mgmt. |
| | | | | | | line discipline | | |
| network interface drivers | buffer cache | | | | | | | |
| | **Vinum** block dev | | | | | | **Vinum** char | |
| | | | | | | character device driver | | |
| the hardware | | | | | | | | |

Figure 17: Kernel I/O structure with Vinum

Apart from the effect of Vinum, it shows the gradual lack of distinction between block and character devices that has occurred since the release of 4.4BSD. NetBSD implements disk block and character devices in the same driver. FreeBSD has completely dispensed

with disk block devices, the shaded area in the figure.

## Design limitations

Vinum was intended to have as few arbitrary limits as possible consistent with an efficient implementation. Nevertheless, a number of limits were imposed in the interests of efficiency, mainly in connection with the device minor number format, which currently persists even in the device file system.

The original release of Vinum was significantly more restrictive in the number of object that it supported. The current limits are:

- 16382 volumes in NetBSD, 262160 volumes in FreeBSD. The highest numbered two volume numbers are reserved for the control devices.

- 16384 plexes in NetBSD, 262162 plexes in FreeBSD.

- 32768 subdisks in NetBSD, 524324 subdisks in FreeBSD.

The difference between NetBSD and FreeBSD is due to the different width of the minor device number in each system.

In addition, Vinum requires a minimum device size of 1 MB. This assumption makes it possible to dispense with some boundary condition checks. Vinum requires 133 kB of disk space to store the header and configuration information, so this restriction does not appear serious.

## Memory allocation

In order to perform its functionality, Vinum allocates a large number of dynamic data structures. Currently these structures are allocated by calling kernel `malloc`. This is a potential problem, since `malloc` interacts with the virtual memory system and may trigger a page fault. The potential for a deadlock exists if the page fault requires a transfer to a Vinum volume. It is probable that Vinum will modify its allocation strategy by reserving a small number of buffers when it starts and using these if a `malloc` request fails.

## To cache or not to cache

Traditionally, UNIX block devices are accessed from the file system via caching routines such as *bread* and *bwrite*. It is also possible to access them directly, but this facility is seldom used. The use of caching enables significant improvements in performance.

Vinum does not cache the data it passes to the lower-level drivers. It would also seem counterproductive to do so: the data is available in cache already, and the only effect of caching it a second time would be to use more memory, thus causing more frequent cache misses.

RAID-5 plexes pose a problem to this reasoning. A RAID-5 write normally first reads the parity block, so there might be some advantage in caching at least the parity blocks. This issue has been deferred for further study.

## Access optimization

The algorithms for RAID-5 access are surprisingly complicated and require a significant amount of temporary data storage. To achieve reasonable performance, they must take error recovery strategies into account at a low level. A RAID 5 access can require one or more of the following actions:

- *Normal read.* All participating subdisks are up, and the transfer can be made directly to the user buffer.

- *Recovery read.* One participating subdisk is down. To recover data, all the other subdisks, including the parity subdisk, must be read. The data is recovered by exclusive-oring all the other blocks.

- *Normal write.* All the participating subdisks are up. This write proceeds in four phases:

  1. Read the old contents of each block and the parity block.

  2. "Remove" the old contents from the parity block with exclusive or.

  3. "Insert" the new contents of the block in the parity block, again with exclusive or.

  4. Write the new contents of the data blocks and the parity block. The data block transfers can be made directly from the user buffer.

- *Degraded write* where the data block is not available. This requires the following steps:

  1. Read in all the other data blocks, excluding the parity block.

  2. Recreate the parity block from the other data blocks and the data to be written.

  3. Write the parity block.

- *Parityless write*, a write where the parity block is not available. This is in fact the simplest: just write the data blocks. This can proceed directly from the user buffer.

## Combining access strategies

In practice, a transfer request may combine the actions above. In particular:

- A read request may request reading both available data (normal read) and non-available data (recovery read). This can be a problem if the address ranges of the two reads do not coincide: the normal read must be extended to cover the address range

of the recovery read, and must thus be performed out of malloced memory.

- Combination of degraded data block write and normal write. The address ranges of the reads may also need to be extended to cover all participating blocks.

An exception exists when the transfer is shorter than the width of the stripe and is spread over two subdisks. In this case, the subdisk addresses do not overlap, so they are effectively two separate requests.

## Examples

The following examples illustrate these concepts:

| Offset | Subdisk 1 | Subdisk 2 | Subdisk 3 | Subdisk 4 | Subdisk 5 |
|--------|-----------|-----------|-----------|-----------|-----------|
| 0x0000 | 0x0000 | 0x1000 | 0x2000 | 0x3000 | Parity |
| 0x1000 | 0x4000 | 0x5000 | 0x6000 | Parity | 0x7000 |
| 0x2000 | 0x8000 | 0x9000 | Parity | 0xa000 | 0xb000 |
| 0x3000 | 0xc000 | Parity | 0xd000 | 0xe000 | 0xf000 |
|  | Parity | 0x10000 | 0x11000 | 0x12000 | 0x13000 |

<div style="margin-left:2em">
▢ Parity block

▨ Data block involved in transfer
</div>

**Figure 18**: A sample RAID-5 transfer

Figure 18 illustrates a number of typical points about RAID-5 transfers. It shows the beginning of a plex with five subdisks and a stripe size of 4 kB. The shaded area shows the area involved in a transfer of 4.5 kB (9 sectors), starting at offset `0xa800` in the plex. A read of this area generates two requests to the lower-level driver: 4 sectors from subdisk 4, starting at offset `0x2800`, and 5 sectors from subdisk 5, starting at offset `0x2000`.

Writing this area is significantly more complicated. From a programming standpoint, the simplest approach is to consider the transfers individually. This would create the following requests:

- Read the old contents of 4 sectors from subdisk 4, starting at offset `0x2800`.

- Read the old contents of 4 sectors from subdisk 3 (the parity disk), starting at offset `0x2800`.

- Perform an exclusive OR of the data read from subdisk 4 with the data read from subdisk 3, storing the result in subdisk 3's data buffer. This effectively "removes" the old data from the parity block.

- Perform an exclusive OR of the data to be written to subdisk 4 with the data read from subdisk 3, storing the result in subdisk 3's data buffer. This effectively "adds" the new data to the parity block.

- Write the new data to 4 sectors of subdisk 4, starting at offset `0x2800`.

- Write 4 sectors of new parity data to subdisk 3 (the parity disk), starting at offset `0x2800`.

- Read the old contents of 5 sectors from subdisk 5, starting at offset `0x2000`.

- Read the old contents of 5 sectors from subdisk 3 (the parity disk), starting at offset `0x2000`.

- Perform an exclusive OR of the data read from subdisk 5 with the data read from subdisk 3, storing the result in subdisk 3's data buffer. This effectively "removes" the old data from the parity block.

- Perform an exclusive OR of the data to be written to subdisk 5 with the data read from subdisk 3, storing the result in subdisk 3's data buffer. This effectively "adds" the new data to the parity block.

- Write the new data to 5 sectors of subdisk 5, starting at offset `0x2000`.

- Write 5 sectors of new parity data to subdisk 3 (the parity disk), starting at offset `0x2000`.

This approach is clearly suboptimal. The operation involves a total of 8 I/O operations and transfers 36 sectors of data. In addition, the two halves of the operation block each other, since each must access the same data on the parity subdisk. Vinum optimizes this access in the following manner:

- Read the old contents of 4 sectors from subdisk 4, starting at offset `0x2800`.

- Read the old contents of 5 sectors from subdisk 5, starting at offset `0x2000`.

- Read the old contents of 8 sectors from subdisk 3 (the parity disk), starting at offset `0x2000`. This represents the complete parity block for the stripe.

- Perform an exclusive OR of the data read from subdisk 4 with the data read from subdisk 3, starting at offset `0x800` into the buffer, and storing the result in the same place in subdisk 3's data buffer.

- Perform an exclusive OR of the data read from subdisk 5 with the data read from subdisk 3, starting at the beginning of the buffer, and storing the result in the same place in subdisk 3's data buffer offset.

- Perform an exclusive OR of the data to be written to subdisk 4 with the modified parity block, starting at offset `0x800` into the buffer, and storing the result in the same place in subdisk 3's data buffer.

- Perform an exclusive OR of the data to be written to subdisk 5 with the modified parity block, starting at the beginning of the buffer, and storing the result in the same place in subdisk 3's data buffer offset.

- Write the new data to 4 sectors of subdisk 4, starting at offset `0x2800`.

- Write the new data to 5 sectors of subdisk 5, starting at offset `0x2000`.

- Write the 8 sectors of new parity data to subdisk 3 (the parity disk), starting at offset `0x2000`.

This is still a lot of work, but by comparison with the non-optimized version, the number of I/O operations has been reduced to 6, and the number of sectors transferred is reduced by 2. The larger the overlap, the greater the saving. If the request had been for a total of 17 sectors, starting at offset `0x9800`, the unoptimized version would have performed 12 I/O operations and moved a total of 68 sectors, while the optimized version would perform 8 I/O operations and move a total of 50 sectors.

### Degraded read

Figure 19 illustrates the situation where a data subdisk fails, in this case subdisk 4. In this case, reading the data from subdisk 5 is trivial. Recreating the data from subdisk 4, however, requires reading all the remaining subdisks. Specifically,

- Read 4 sectors each from subdisks 1, 2 and 3, starting at offset `0x2800` in each case.

- Read 8 sectors from subdisk 5, starting at offset `0x2800`.

- Clear the user buffer area for the data corresponding to subdisk 4.

- Perform an "exclusive or" operation on this data buffer with data from subdisks 1, 2, 3, and the last four sectors of the data from subdisk 5.

- Transfer the first 5 sectors of data from the data buffer for subdisk 5 to the corresponding place in the user data buffer.
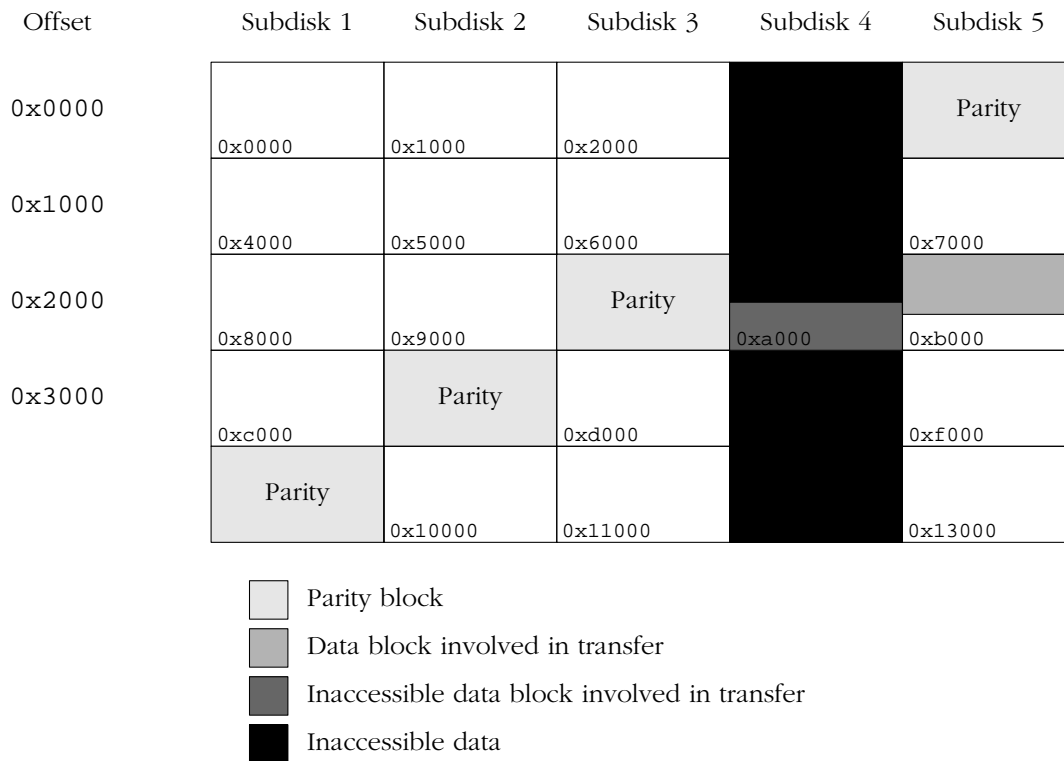
| Offset | Subdisk 1 | Subdisk 2 | Subdisk 3 | Subdisk 4 | Subdisk 5 |
|---|---|---|---|---|---|
| 0x0000 | 0x0000 | 0x1000 | 0x2000 | *Inaccessible data* | Parity |
| 0x1000 | 0x4000 | 0x5000 | 0x6000 | *Inaccessible data* | 0x7000 |
| 0x2000 | 0x8000 | 0x9000 | Parity | 0xa000 | 0xb000 |
| 0x3000 | 0xc000 | Parity | 0xd000 | *Inaccessible data* | 0xf000 |
|  | Parity | 0x10000 | 0x11000 | *Inaccessible data* | 0x13000 |

- Parity block
- Data block involved in transfer
- Inaccessible data block involved in transfer
- Inaccessible data

**Figure 19: RAID-5 transfer with inaccessible data block**

## Degraded write

There are two different scenarios to be considered in a degraded write. Referring to the previous example, the operations required are a mixture of normal write (for subdisk 5) and degraded write (for subdisk 4). In detail, the operations are:

- Read 4 sectors each from subdisks 1 and 2, starting at offset 0x2800, into temporary storage.

- Read 5 sectors from subdisk 3 (parity block), starting at offset 0x2000, into the beginning of an 8 sector temporary storage buffer.

- Clear the last 3 sectors of the parity block.

- Read 8 sectors from subdisk 5, starting at offset 0x2000, into temporary storage.

- "Remove" the first 5 sectors of subdisk 5 data from the parity block with exclusive or.

- Rebuild the last 3 sectors of the parity block by exclusive or of the corresponding data from subdisks 1, 2, 5 and the data to be written for subdisk 4.

- Write the parity block back to subdisk 3 (8 sectors).

- Write 5 sectors user data to subdisk 5.

## Parityless write

Another situation arises when the subdisk containing the parity block fails, as shown in figure 20.



| Offset | Subdisk 1 | Subdisk 2 | Subdisk 3 | Subdisk 4 | Subdisk 5 |
|--------|-----------|-----------|-----------|-----------|-----------|
| 0x0000 | 0x0000 | 0x1000 | | 0x3000 | Parity |
| 0x1000 | 0x4000 | 0x5000 | | Parity | 0x7000 |
| 0x2000 | 0x8000 | 0x9000 | | 0xa000 | 0xb000 |
| 0x3000 | 0xc000 | Parity | | 0xe000 | 0xf000 |
| | Parity | 0x10000 | | 0x12000 | 0x13000 |

Parity block

Data block involved in transfer

Inaccessible data

**Figure 20: RAID-5 transfer with inaccessible parity block**

This configuration poses no problems on reading, since all the data is accessible. On writing, however, it is not possible to write the parity block. It is not possible to recover from this problem at the time of the write, so the write operation simplifies to writing only the data blocks. The parity block will be recreated when the subdisk is brought up again.

36

# Driver structure

Vinum can issue multiple disk transfers for a single I/O request:

- As the result of striping or concatenation, the data for a single request may map to more than one drive. In this case, Vinum builds a request structure which issues all necessary I/O requests at one time. This behaviour has had the unexpected effect of highlighting problems with marginal SCSI hardware by imposing heavy activity on the bus.

- As seen above, many RAID-5 operations require a second set of I/O transfers after the initial transfers have completed.

- In case of an I/O failure on a resilient volume, Vinum must reschedule the I/O to a different plex.

The second set of RAID-5 operations and I/O recovery do not match well with the design of UNIX device drivers: typically, the "top half"[1] of a UNIX device driver issues I/O commands and returns to the caller. The caller may choose to wait for completion, but one of the most frequent uses of a block device is where the virtual memory subsystem issues writes and does not wait for completion.

This poses a problem: who issues the second set of requests? The following possibilities, listed in order of increasing desirability, exist:

1.  The top half can wait for completion of the first set of requests and then launch the second set before returning to the caller. This approach can seriously impact system performance and possibly cause deadlocks.

2.  In a threaded kernel, the strategy routine can create a thread which waits for completion of the first set of requests and starts the second set without impacting the main thread of the process. At the moment this approach is not possible, since FreeBSD currently does not provide kernel thread support. It also appears likely that it could cause a number of problems in the areas of thread synchronization and performance.

3.  Ownership of the requests can be "given" to another process, which will be awakened when they complete. This process can then issue the second set of requests. This approach is feasible, and it is used by some subsystems, notably NFS. It does not pose the same severe performance penalty of the previous possibility, but it does require that another process be scheduled twice for every I/O.

---

1. UNIX device drivers run in two separate environments. The "top half" runs in the process context, while the "bottom half" runs in the interrupt context. There are severe restrictions on the functions that the bottom half of the driver can perform.

4.  The second set of requests can be launched from the "bottom half" of the driver. This is potentially dangerous: the interrupt routine must call the `start` routine. While this is not expressly prohibited, the `start` routine is normally used by the top half of a driver, and may call functions which are prohibited in the bottom half.

Initially, Vinum used the fourth solution. This worked for most drivers, but some drivers required functions only available in the "top half", such as `malloc` for ISA bounce buffers. Current FreeBSD drivers no longer call these functions, but it is possible that the situation will arise again.

On the other hand, this method does not allow I/O recovery. Vinum now uses a daemon process for I/O recovery and a couple of other housekeeping activities, such as saving the configuration database. The additional scheduling overhead for these activities is negligible, but it is the reason that the RAID-5 second stage does not use the daemon.

## Vinum root file systems

Initially Vinum did not support the root file system. The main problem was that the FreeBSD boot loader does not understand the generality of a Vinum volume. The small space reserved for the boot loader also makes it difficult to modify it to understand Vinum. This is a situation paralleled in many commercial operating systems.

Based on prior attempts at solving this issue, the following alternatives were considered:

- Teach the bootstrap code about Vinum so that it could start Vinum and load directly from a Vinum volume. This was considered an unlikely alternative, since Vinum needs to know too much about the kernel environment, and there's not enough space for the code in conventional bootstraps.

- Create a separate boot file system and put the kernel there, then start Vinum and the root file system. Many System V implementations used this method, but it's untidy.

- Create an root file system on ramdisk. This effectively lives in swap, but there's no problem there.

- Boot normally, start Vinum and then mount the root file system on top of the old root file system. It appeared that there would be a number of problems with this approach, since the system must access */etc/fstab* to determine where to mount root.

In practice, however, it turned out that the last alternative was not so difficult after all, and Vinum on FreeBSD now supports the root file system. The same considerations apply to NetBSD, but at the time of writing the coding is not yet complete.

The changes required were:

- Load the Vinum module at boot time. Changes in the boot loader had already made this possible.

- Start Vinum automatically at boot time, which required some changes in Vinum itself and also minor modifications to the FreeBSD startup procedure.

- Mount the Vinum volume instead of the underlying partition. This proved to be trivial.

Some restrictions remain for Vinum volumes for the root file system:

- They must contain at least one concatenated plex with only one subdisk.

- These subdisks must correspond to a UFS partition.

- The system can only boot from one of these subdisks.

Theoretically it would be possible to have other plexes in the root file system, but since it's not possible to boot from them, there's no obvious reason to do so.

# Future directions

A number of additional features have been proposed for Vinum:

- *Hot spare* capability: on the failure of a disk drive, the volume manager automatically recovers the data to another drive. Some work has been done in this direction, but it is currently not yet complete.

- *Logging* changes to a degraded volume. Rebuilding a plex usually requires copying the entire volume. In a volume with a high read to write, if a disk goes down temporarily and then becomes accessible again (for example, as the result of controller failure), most of the data is already present and does not need to be copied. Logging pinpoints which blocks require copying in order to bring the stale plex up to date.

- *Snapshots* of a volume. It is often useful to freeze the state of a volume, for example for backup purposes. A backup of a large volume can take several hours. It can be inconvenient or impossible to prohibit updates during this time. A snapshot solves this problem by maintaining *before images*, a copy of the old contents of the modified data blocks. Access to the plex reads the blocks from the snapshot plex if it contains the data, and from another plex if it does not.

  Implementing snapshots in Vinum alone would solve only part of the problem: there must also be a way to ensure that the data on the file system is consistent from a user standpoint when the snapshot is taken. This task involves such components as file systems and databases and is thus outside the scope of Vinum. A snapshot facility

already exists for UFS, so the concept of snapshots at the volume manager level has become less interesting.

- A *SNMP interface* for central management of Vinum systems.

- A *GUI* interface is currently *not* planned, though it is relatively simple to program, since no kernel code is needed. As the number of failures testify, a good GUI interface is apparently very difficult to write, and it tends to gloss over important administrative aspects, so it's not clear that the advantages justify the effort. On the other hand, a graphical output of the configuration could be of advantage.

- *Remote data replication* is of interest either for backup purposes or for read-only access at a remote site. From a conceptual viewpoint, it could be achieved by interfacing to a network driver instead of a local disk driver.

- *Extending striped and RAID-5 plexes* is a slow complicated operation, but it is feasible.

Currently there are no further plans for any of these features.

## Vinum under other operating systems

Vinum first appeared in FreeBSD in 1998. In 2003, a couple of independent efforts ported it to NetBSD (where it is now in the source tree) and OpenBSD (where it is not). Vinum currently does not run under Linux. Currently some interest has been expressed on the developer list `vinum-developers@auug.org.au`. A port would require rewriting the driver interface layer, affecting perhaps 15% of the total code. The question remains whether there's enough interest in doing so in view of the multitude of other similar products available.

## Vinum: the future

Vinum is no longer new: it has been around for nearly six years. In the meantime, many other products have become available, both for FreeBSD and for other systems. In particular, since release 5.0 FreeBSD includes the GEOM subsystem, which addresses some of the same issues that Vinum addresses in a different, usually more general manner. Inevitably, there are conflicts. For example, at the time of writing it is no longer possible to use a Vinum volume as a swap partition on FreeBSD -CURRENT, the development version of FreeBSD.

This problem will be resolved. In a more general sense, though, it does not make sense to have two different and conflicting implementations on a single platform. Since GEOM is more general, it makes sense to adapt Vinum to GEOM.

Desirable features in GEOM include:

- You can define a new mapping and insert it into the GEOM stack. For example, some users have asked for the ability to build a single plex out of multiple mirrored subdisks, which is not possible in Vinum. This could be done with GEOM. In this respect, GEOM resembles STREAMS.

- Autodiscovery allows "consumers" to find about new devices as they arrive. This would make it unnecessary for Vinum to scan every disk in the system. When a disk of type Vinum is discovered, whether at boot time or later, the Vinum code will be informed and will have the opportunity to add additional configuration information.

At present, this issue is under discussion in the FreeBSD mailing lists. It's possible that an initial implementation will simply convert Vinum to work as a single, complex GEOM transformation, and that later the individual objects will be split up as independent transforms.

# References

FreeBSD project: *http://www.FreeBSD.org/*

Marshall Kirk McKusick, Keith Bostic, Michael J. Karels, John S. Quarterman. *The Design and Implementation of the 4.4BSD Operating System*, Addison Wesley, 1996.

NetBSD project: *http://www.NetBSD.org/*

Rawio: a raw disk I/O benchmark. *ftp://ftp.lemis.com/pub/rawio.tar.gz*

The VERITAS Volume manager, *http://www.veritas.com/*.

Vinum web site: *http://www.vinumvm.org/*

Vinum developers list: *http://www.auug.org.au/mailman/listinfo/vinum-devel/*