

LVM Version 2.0 Volume Groups in HP-UX 11i v3



Abstract	3
Audience	3
Introduction	3
What is a Version 2.0 Volume Group?	3
How Do I Use a Version 2.0 Volume Group?.....	3
What Do I Need to Change to Manage my Version 1.0 Volume Groups?.....	3
New Limits	4
LVM Changes in the March 2008 Release of HP-UX 11i v3	4
Creation of a Volume Group.....	4
Creation of a Version 1.0 Volume Group.....	4
Creation of a Version 2.0 Volume Group.....	4
Display of LVM Information.....	6
vgdisplay	6
pvdisplay	7
lvdisplay	7
LVM Device Special Files	9
lvmtab, lvmtab_p	9
New Command (<i>lvmdm</i>).....	10
Summary List of New or Obsolete Command Options.....	11
Overview of Changes by Command	11
Other Differences Between Volume Group Versions 1.0 and 2.0	13
Sparing Not Supported on Version 2.0 Volume Groups	13
Boot, Dump, and Primary Swap Not Supported on Version 2.0 Volume Groups	13
Bad Block Relocation Not Supported on Version 2.0 Volume Groups	13
Cluster Lock Not Supported on Version 2.0 Volume Groups	13
Commands Not Supported on Version 2.0 Volume Groups	13
How to Provision a Version 2.0 Volume Group.....	14
Simplified Provisioning	14
What is the Disadvantage of Over-Provisioning?	15
Metadata Size on Disk versus Maximum Volume Group Size	15
What are the Advantages of Over-Provisioning with Version 2.0 Volume Groups?.....	16
Version 2.0 Can Leave as Much Space for User Data on Disk as Version 1.0	16

How to Pick an Extent Size	18
Memory Footprint	20
Minimum Memory Footprint	20
Maximum Memory Footprint	21
Comparison of Memory Footprints for the Same Number of Logical and Physical Volumes	22
How to Configure the New Limits	24
Number of Logical Volumes or Physical Volumes in a Volume Group	24
Maximum Size of the Volume Group.....	24
Size of Logical Volume or Number of Mirrors	24
Migration Between Version 1.0 and Version 2.0	24
Why Would You Migrate?	24
How to Migrate.....	24
When to Migrate.....	24
High Availability	25
When Does It Make Sense to Use More Than Two Mirrors?.....	25
Glossary	27
For More Information	27
Call to Action	27

Abstract

This whitepaper is an introduction to the LVM volume group version 2.0. These volume groups became available with the March 2008 release of HP-UX 11i v3 (11.31). Prior to this release, only Version 1.0 volume groups were available.

Audience

The document is intended for system administrators, operators, and customers who wish to utilize and know about the LVM 2.0 volume groups. It is assumed that the reader has a basic knowledge of LVM.

Introduction

LVM and MirrorDisk/UX now support two versions of volume groups. Version 1.0 is the version supported on all current and previous versions of HP-UX 11i. The procedures and command syntax for managing Version 1.0 volume groups are unchanged from previous releases, except for the enhancements described later in this paper. When creating a new volume group, `vgcreate` defaults to Version 1.0.

What is a Version 2.0 Volume Group?

A Version 2.0 volume group is a volume group whose metadata layout is different from the one used for Version 1.0 volume groups. A Version 2.0 volume group extends the limits of a Version 1.0 volume group.

How Do I Use a Version 2.0 Volume Group?

A Version 2.0 volume group is managed the same way as a Version 1.0 volume group using the same user interface. This interface is the same as in the HP-UX 11i v3 (11.31) initial release but has been extended for Version 2.0 volume groups. These additions are listed in the section "[LVM Changes in the March 2008 Release of HP-UX 11i v3](#)".

What Do I Need to Change to Manage my Version 1.0 Volume Groups?

Nothing, unless you have scripts that parse the output of `lvdisplay`, `vgdisplay`, `pvdisplay` or `vgscan`. The output of these commands is slightly changed and that may impact parsing. The user interface to manage the Version 1.0 volume group didn't change. The same LVM commands with the same options (as in HP-UX 11i v3 (11.31) initial release) are used to handle Version 1.0 volume groups. There are some additional options that may be used on Version 1.0 volume groups.

New Limits

The limits remain the same for Version 1.0 volume groups. The table below illustrates the new upper limits with Version 2.0 volume groups.

	Maximum Supported Version 1.0 volume group limits	Maximum Supported Version 2.0 volume group limits
VG size	510TB	2PB
LV size	2^{46} (64TB) 16TB supported	2^{48} (256TB)
PV size	2^{41} (2TB)	2^{44} (16TB)
Number of VGs	256	512
Number of LVs per VG	255	511
Number of PVs per VG	255	511
Number of mirror copies	2	5
Number of extents per VG	2^{16} (64K)	2^{25} (32M)
Extent size	1 to 256MB	1 to 256 MB
Stripe width	255	511

LVM Changes in the March 2008 Release of HP-UX 11i v3

Creation of a Volume Group

Automatic creation of the volume group directory and group file is available starting with the March 2008 release of HP-UX 11i v3. It is available for Version 1.0 and 2.0 volume groups. Below are examples of volume group creation using the legacy method (*mkdir*, *mknod*) and examples of automatic creation.

Creation of a Version 1.0 Volume Group

To avoid manually creating the volume group directory and the group file, just use *vgcreate*. In this case *vgcreate* automatically creates the directory and group file if they don't exist for this volume group.

Example

```
# vgcreate -s 8 -l 3 -p 16 -e 63535 /dev/vg01 /dev/dsk/c3t4d0
```

To select a particular volume group number or to create the volume group the same way as in releases before March 2008, first create the volume group directory and the group file.

Example

```
# mkdir /dev/vg01
# mknod /dev/vg01/group c 64 0x010000
# vgcreate -s 8 -l 3 -p 16 -e 63535 /dev/vg01 /dev/dsk/c3t4d0
```

Creation of a Version 2.0 Volume Group

- As with Version 1.0 volume groups, for Version 2.0 volume groups, you can manually create the volume group directory and group file or *vgcreate* automatically creates them if they don't already exist.

- A new option `-V` specifies the version of the volume group to create. To create a Version 2.0 volume group, specify `"-V 2.0"` as an option on the command line. Version 1.0 volume groups are the default and `-V 1.0` is also valid to specify a Version 1.0 volume group.
- A new option `-S` must be used when creating a Version 2.0 volume group. It specifies the maximum size the volume group may reach. LVM uses this size to reserve enough space on disk to accommodate the metadata for a volume group of that size. Note this does not need to be the actual size of the volume group but how large the volume group can grow over time. The "How to Provision a Version 2.0 Volume Group" section later in this document covers provisioning strategies and tradeoffs.
- For a Version 2.0 volume group, if the sum of the physical volume sizes passed on the command line is greater than the maximum size specified with the `-S` option, `vgcreate` automatically increases the maximum size of the volume group to the sum of the physical volumes' sizes and displays an informational message.
- The maximum size specified requires a unit qualifier. The units are in Megabytes (2^{20}), Gigabytes (2^{30}), Terabytes (2^{40}) and Petabytes (2^{50}), represented by `m`, `g`, `t` and `p` respectively on the command line.
- A new option `-E` displays which volume group size can be reached for a given extent size, or which minimum extent size must be used for a given volume group size. This can be used in determining the best fit of volume group size and physical extent size.
- The option `-s` (to specify the extent size) is mandatory for Version 2.0 volume groups.
- The options `-e`, `-l`, `-p`, and `-f` are invalid when applied to a Version 2.0 volume group. If used for a Version 2.0 volume group, `vgcreate` fails.
 - `-e max_pe` is not needed for Version 2.0 volume groups because the Version 2.0 metadata format is provisioned so that the number of physical extents can always grow to the maximum size specified with `-S`.
 - `-l max_lv` and `-p max_pv` are not needed for Version 2.0 volume groups because any Version 2.0 volume group is provisioned to handle the maximum supported number of logical volumes and physical volumes.
 - `-f` is not needed for Version 2.0 volume groups because bad block handling is handled by all currently supported physical disk drives.

Example

To create a Version 2.0 volume group using 32 MB extents and a maximum size of 1 PB:

```
# vgcreate -V 2.0 -s 32 -S 1p /dev/vg01 /dev/disk/disk50
```

To select a particular volume group number or to create the volume group the same way as in releases before March 2008, first create the volume group directory and the group file.

Example

To create a Version 2.0 volume group using 4 MB extents and a maximum size of 8 TB:

```
# mkdir /dev/vg01
# mknod /dev/vg01/group c 128 0x001000
# vgcreate -V 2.0 -s 4 -S 8t /dev/vg01 /dev/disk/disk49
```

Display of LVM Information

With the introduction of the new volume group version, the display commands display additional information. In addition, because Version 2.0 volume groups support increased sizes, some of the display information may contain much larger numbers than previously.

This section illustrates the difference for each display operation. As a general practice, when writing scripts to gather volume manager information, use the `-F` option introduced in 11i v3 to reduce the impact of future volume manager changes.

`vgdisplay`

In the March 2008 release of HP-UX 11i v3, the `vgdisplay` command displays three more lines than `vgdisplay` from the HP-UX 11i v3 initial release. In addition, the displayed number of physical extents may be larger on Version 2.0 volume groups.

The new lines appear for both Version 1.0 and Version 2.0 volume groups:

- The volume group version, as "VG version".
- The maximum size of the volume group, as "VG Max Size".
- The maximum number of physical extents the volume group can contain, as "VG Max Extents". This value is the ratio of volume group maximum size to extent size.

Example

```
# vgdisplay vgtest12
--- Volume groups ---
VG Name                /dev/vgtest12
VG Write Access        read/write
VG Status              available
Max LV                 511
Cur LV                1
Open LV               1
Max PV                 511
Cur PV                70
Act PV                 70
Max PE per PV         2097152
VGDA                   140
PE Size (Mbytes)      8
Total PE               1123123
Alloc PE               40
Free PE               1123083
Total PVG              0
Total Spare PVs       0
Total Spare PVs in use 0
VG Version             2.0
VG Max Size            16t
VG Max Extents         2097152
```

} New lines for March 2008
release

pvdisk

In the March 2008 release of HP-UX 11iV3, the *pvdisk* command displays one more line than *pvdisk* from the HP-UX 11i v3 initial release.

When the *-d* option is used with *pvdisk*, the new line appears for both Version 1.0 and Version 2.0 volume groups. The new line is "Data End". "Data Start" and "Data End" are block numbers from the start of the disk. In this display, a block is always 1024 bytes, no matter what the disk sector size is. "Data End" is the block number of the last block on the disk that may be used by LVM to store user data.

LVM only uses storage in multiples of extent size. In some cases, "Data End" – "Data Start" may not be a multiple of the extent size. In this case, the space between the last extent and the "Data End" will not contain any user data.

Example

```
# pvdisk -d /dev/disk/disk148
--- Physical volumes ---
PV Name                /dev/disk/disk148
VG Name                /dev/vgtest12
PV Status              available
Allocatable            yes
VGDA                   2
Cur LV                0
PE Size (Mbytes)       8
Total PE               16401
Free PE                16401
Allocated PE           0
Stale PE                0
IO Timeout (Seconds)  default
Autoswitch              On
Data Start              1024
Data End                134358016 } New line for March 2008
Boot Disk               no      } release
Relocated Blocks        0
Proactive Polling      On
```

lvdisplay

For Version 1.0 volume groups, the output of *lvdisplay* is the same as the HP-UX 11i v3 initial release. For Version 2.0 volume groups, *lvdisplay* output changes if the number of mirrors is greater than 2. In that case when *lvdisplay -v* is used, *lvdisplay* displays extent mapping data for the additional mirrors on additional lines. This is illustrated in the example below.

Example

```
# lvdisplay -v /dev/vgtest12/lvol1
--- Logical volumes ---
LV Name                /dev/vgtest12/lvol1
VG Name                /dev/vgtest12
LV Permission          read/write
LV Status              available/syncd
Mirror copies          3
Consistency Recovery   MWC
Schedule               parallel
LV Size (Mbytes)       80
Current LE             10
```

```

Allocated PE          40
Stripes              0
Stripe Size (Kbytes) 0
Bad block            NONE
Allocation            strict
IO Timeout (Seconds) default

```

--- Distribution of logical volume ---

```

PV Name      LE on PV  PE on PV
/dev/disk/disk79  10      10
/dev/disk/disk80  10      10
/dev/disk/disk81  10      10
/dev/disk/disk82  10      10

```

--- Logical extents ---

```

LE          PV1          PE1          Status 1 PV2          PE2
Status 2 PV3          PE3          Status 3
LE          PV4          PE4          Status 4
000000000 /dev/disk/disk79  00000000 current /dev/disk/disk80
000000000 current /dev/disk/disk81  00000000 current
000000000 /dev/disk/disk82  00000000 current
000000001 /dev/disk/disk79  00000001 current /dev/disk/disk80
000000001 current /dev/disk/disk81  00000001 current
000000001 /dev/disk/disk82  00000001 current
000000002 /dev/disk/disk79  00000002 current /dev/disk/disk80
000000002 current /dev/disk/disk81  00000002 current
000000002 /dev/disk/disk82  00000002 current
000000003 /dev/disk/disk79  00000003 current /dev/disk/disk80
000000003 current /dev/disk/disk81  00000003 current
000000003 /dev/disk/disk82  00000003 current
000000004 /dev/disk/disk79  00000004 current /dev/disk/disk80
000000004 current /dev/disk/disk81  00000004 current
000000004 /dev/disk/disk82  00000004 current
000000005 /dev/disk/disk79  00000005 current /dev/disk/disk80
000000005 current /dev/disk/disk81  00000005 current
000000005 /dev/disk/disk82  00000005 current
.
.
.

```

There is another difference regarding *lvdisplay*.

For Version 1.0 volume groups, the order in which the physical extents are displayed by *lvdisplay -v* may change across activation cycles. This is because when a Version 1.0 volume group is activated, LVM will order the physical extents of a logical extent (LE) in order of increasing physical volume number order.

For example the display:

```

--- Logical extents ---
LE          PV1          PE1          Status 1          PV2          PE2          Status 2
000000003 /dev/disk/disk79  00000003 current /dev/disk/disk80  00000014 current

```

may change to:

```

--- Logical extents ---
LE          PV1          PE1          Status 1          PV2          PE2          Status 2
000000003 /dev/disk/disk80  00000014 current /dev/disk/disk79  00000003 current

```

after deactivation and activation again.

For Version 2.0 volume groups, the order in the display does not change across activation.

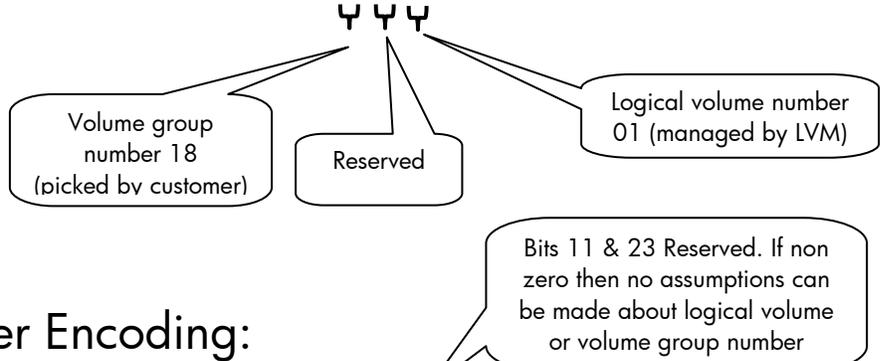
LVM Device Special Files

With Version 2.0 volume groups, LVM device special files have changed. For Version 1.0 volume groups, procedures will not need to change. For Version 2.0 version volume groups, there are a few differences.

- Device special files for Version 1.0 volume groups are unchanged
- Device special files for Version 2.0 volume groups have a new major number (128)
- Device special files for Version 2.0 volume groups have a different minor number encoding scheme

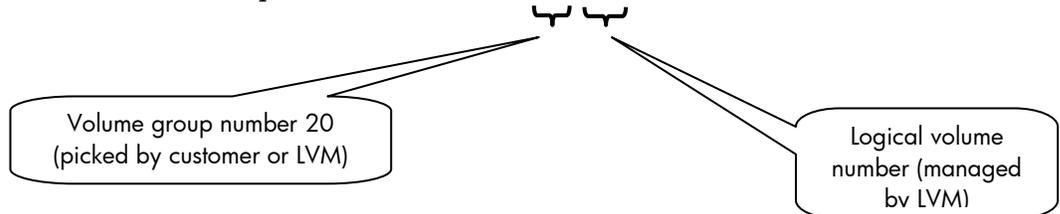
1.0 Minor Number Encoding:

```
brw-r----- 1 root sys 64 0x120001 Jun 23 2006 lv011
```



2.0 Minor Number Encoding:

```
brw-r----- 1 root sys 128 0x014000 Jun 23 2006 lv011
```



lvmtab, lvmtab_p

LVM has a well known configuration file “/etc/lvmtab”. With the introduction of Version 2.0 volume groups, LVM uses a new configuration file.

- /etc/lvmtab is a private binary file. It contains information related to Version 1.0 volume groups only. It is compatible with all supported HP-UX releases.
- /etc/lvmtab_p is a new private binary file. It contains information related to Version 2.0 volume groups only. It has a different internal structure from /etc/lvmtab and contains additional information.

Example

```
# strings /etc/lvmtab_p
/dev/vgtest12
A00000000000000004Sun Mar 16 03:02:19 2007241fa9de-dba3-11da-9cd2-23efa80dc3e7
/dev/disk/disk79
/dev/disk/disk80
/dev/disk/disk81
/dev/disk/disk82
/dev/disk/disk83
/dev/disk/disk84
/dev/disk/disk85
```

New Command (*lvmadm*)

lvmadm (1M)

The *lvmadm* command displays the supported limits for Version 1.0 and 2.0 volume groups. It is not possible to create a volume group that exceeds these limits.

Example

```
# lvmadm -t

--- LVM Limits ---
VG Version                1.0
Max VG Size (Tbytes)     510
Max LV Size (Tbytes)     16
Max PV Size (Tbytes)     2
Max VGs                   256
Max LVs                   255
Max PVs                   255
Max Mirrors               2
Max Stripes               255
Max Stripe Size (Kbytes) 32768
Max LXs per LV            65535
Max PXs per PV            65535
Max Extent Size (Mbytes) 256

VG Version                2.0
Max VG Size (Tbytes)     2048
Max LV Size (Tbytes)     256
Max PV Size (Tbytes)     16
Max VGs                   512
Max LVs                   511
Max PVs                   511
Max Mirrors               5
Max Stripes               511
Max Stripe Size (Kbytes) 262144
Max LXs per LV            33554432
Max PXs per PV            16777216
Max Extent Size (Mbytes) 256
```

Summary List of New or Obsolete Command Options

The following table summarizes the command option changes with Version 1.0 or Version 2.0 volume groups.

“No change” means that the option works as in the HP-UX 11i v3 (11.31) initial release.

“Optional” means that the option is not necessary to create or use the volume group.

“Invalid” means that the LVM command fails.

“Ignored” means that the LVM command ignores the option, but does not automatically fail.

“Obsolete” means that the LVM command ignores the option and displays a warning message, but does not automatically fail.

“New option” means that the option is first delivered in the March 2008 release of HP-UX 11i v3 (11.31).

LVM command	Option	New option	Effect on Version 1.0 volume group	Effect on Version 2.0 volume group
<i>vgcreate</i>	-V	Y	Optional	Mandatory
	-S	Y	Invalid	Mandatory
	-E	Y	Invalid	Optional
	-s	N	No change	Mandatory
	-e	N	No change	Invalid
	-l	N	No change	Invalid
	-p	N	No change	Invalid
<i>vgextend</i>	-f	N	No change	Obsolete
	-z y	N	No change	Invalid
	-z n	N	No change	Obsolete
<i>vgcfgrestore</i>	-F	N	No change	Obsolete
	-v	Y	Optional	Optional
<i>vgremove</i>	-X	Y	Optional	Optional
<i>pvchange</i>	-z y	N	No change	Invalid
	-z n	N	No change	Obsolete
<i>lvcreate</i>	-r	N	No change	Ignored
<i>Pvcreate</i>	-s	N	No change	Obsolete
<i>lvchange</i>	-r	N	No change	Ignored

Overview of Changes by Command

vgcreate(1M)

- *vgcreate* is changed. It can be used as before to create Version 1.0 volume groups. To create 2.0 volume groups, new options have to be used. See [“Creation of a Volume Group”](#) and [“How to Provision a Version 2.0 Volume Group”](#) for more details.

vgextend(1M)

- The option -f is obsolete on a Version 2.0 volume group. If used on a Version 2.0 volume group, *vgextend* displays a warning.
- The option -z y is invalid on a Version 2.0 volume group. If used on a Version 2.0 volume group, *vgextend* fails.
- The option -z n is obsolete on a Version 2.0 volume group. If used on a Version 2.0 volume group, *vgextend* displays a warning.

- *vgextend* on a Version 2.0 volume group fails if the volume group is already at the maximum size specified with *vgcreate*.
- *vgextend* on a Version 2.0 volume group displays a warning if only part of the physical volume can be added to the volume group. If adding the whole physical volume results in exceeding the maximum size of the volume group, only a part of the physical volume is added.
- It is not possible to add bootable disks to a Version 2.0 volume group.

vgcfgrestore(1M)

- The option *-F* is obsolete on a Version 2.0 volume group. If used on a Version 2.0 volume group, *vgcfgrestore* displays a warning.

vgimport(1M)

- Automatically creates the volume group directory (under */dev*) and the group file if they do not already exist. This applies to Version 1.0 and 2.0 volume groups.

vgremove(1M)

- For both Version 1.0 and 2.0 volume groups by default *vgremove* does not delete the volume group directory and group file. A new option *-X* has been added to *vgremove*. If used, *-X* option deletes the volume group directory and group file. This option applies to Version 1.0 and 2.0 volume groups.

pvchange(1M)

- The option *-z y* is invalid on a Version 2.0 volume group. If used on a Version 2.0 volume group, *pvchange* fails.
- The option *-z n* is obsolete on a Version 2.0 volume group. If used on a Version 2.0 volume group, *pvchange* displays a warning.

lvcreate(1M)

- Larger limits are allowed on 2.0 logical volumes.
- The option *-r* is ignored on 2.0 logical volumes.

lvextend(1M), lvreduce(1M)

- Larger limits are allowed on 2.0 logical volumes.

pvcreate(1M)

- The option *-s* is obsolete on a Version 2.0 volume group. When a physical volume that was created with *-s* is added to a Version 2.0 volume group, *pvcreate* displays a warning and the *-s* is ignored.
- It is not possible to add bootable disks to a Version 2.0 volume group.

lvchange(1M)

- The option *-r* is obsolete on a Version 2.0 volume group. If used on a Version 2.0 volume group, *lvchange* displays a warning.

vgdisplay (1M)

- Change in the display. See "[Display of LVM Information.](#)"

pvdisplay(1M)

- Change in the display. See "[Display of LVM Information.](#)"

lvdisplay(1M)

- Change in the display if using more than two mirrors. See "[Display of LVM Information.](#)"

Other Differences Between Volume Group Versions 1.0 and 2.0

Sparing Not Supported on Version 2.0 Volume Groups

The sparing function that is available on Version 1.0 volume groups is not available on Version 2.0 volume groups.

Boot, Dump, and Primary Swap Not Supported on Version 2.0 Volume Groups

A Version 2.0 volume group can not be the root volume group nor can it contain a boot disk. As a consequence, Ignite/UX does not allow the creation of a Version 2.0 root volume group during a cold install and the commands *lvrmboot* and *lvlnboot* do not apply to Version 2.0 volume groups.

A Version 2.0 volume group can not be used to save crashdumps. The *crashconf* command displays an error if used on a Version 2.0 volume group. For example:

```
# crashconf -s /dev/vgtest2_12/lvol1
/dev/vgtest2_12/lvol1: error: unsupported disk layout
warning: All dump devices may not be marked persistent
```

A logical volume in a Version 2.0 volume group can't be used for primary swap but may be used for secondary swap. The *swapon* command displays an error if used on a Version 2.0 volume group. For example:

```
# swapon -s /dev/vgtest2_12/r1vol1
swapon: /dev/vgtest2_12/lvol1: Invalid argument
```

Bad Block Relocation Not Supported on Version 2.0 Volume Groups

Version 2.0 volume groups do not perform bad block relocation in software because modern disks and disk arrays handle such relocation in their own hardware.

Cluster Lock Not Supported on Version 2.0 Volume Groups

It is not possible to create a cluster lock on a Version 2.0 volume group. When using cluster lock with Serviceguard, a Version 1.0 volume group must be used.

Commands Not Supported on Version 2.0 Volume Groups

vgmodify(1M), lvrmboot(1M), lvlnboot(1M), pvck(1M)

How to Provision a Version 2.0 Volume Group

In the context of this document “provisioning” is the set of parameters used at the time of the volume group creation (*vgcreate*).

Provisioning is defined as how the volume manager reserves space on disk for future growth. Getting the provisioning correct is important because it is a tradeoff between ease of future growth and efficient disk space usage.

Simplified Provisioning

With Version 1.0 volume groups, you provision the volume group with three parameters: max PVs, max extents, and max LVs. Furthermore, Version 1.0 volume groups had to keep their disk metadata within one physical extent. Both of these factors made it challenging to configure Version 1.0 volume groups so that they could easily grow over time. With Version 2.0 volume groups, you have to give only one maximum when provisioning. This is the maximum size of the volume group (new option *-S*). The size entered with *-S* is the size of the user data. LVM guarantees that for any Version 2.0 volume group, you can later add physical volumes and logical volumes up to the Version 2.0 supported limits. In addition, LVM disk metadata can be larger than one physical extent, thus giving much more flexibility on how to configure and provision a volume group.

To make room for LVM metadata, the actual size of the volume group is larger than what is specified with *-S*. For a Version 2.0 volume group, you don't need to think if the default maximum number of extents per physical volumes or the default maximum number of logical volumes or the default maximum number of physical volumes is sufficient. These are automatically managed by LVM.

Tip: When planning your volume group needs, consider how fast your storage requirements have been growing over time. Estimate how fast and how long particular volume groups will exist. For example, if you have a database that doubles in size every two years, and you expect the application environment to last ten years, then provisioning five times the current amount may be a good place to start.

Example

- Creation of 2.0 volume group provisioned for 1 petabyte.

```
# vgcreate -V 2.0 -s 32 -S 1p myvg /dev/disk/disk149
```

There is a relationship between the maximum number of extents and the extent size when selecting the maximum volume group size. To help in selecting the extent size, you can preview via *vgcreate* extent size or maximum volume group size.

Once you know the volume group size you want to provision for, you can determine the minimum extent size required to achieve it. For that use *vgcreate* with the option *-E*.

Example

- What is the minimum extent size to provision a volume group for 1 petabyte?

```
# vgcreate -V 2.0 -E -S 1p  
Max_VG_size=1p:extent_size=32m
```

The maximum size for a volume group is displayed with *vgdisplay*. It is named “VG Max Size” in the display.

What is the Disadvantage of Over-Provisioning?

The cost of over provisioning has been reduced compared to Version 1.0 volume groups. It costs only disk space. It doesn't cost system memory. When a large Version 2.0 volume group is provisioned, enough disk space is reserved for the metadata to handle any LVM configuration for a volume group of that size.

However, it is a different story for the system memory. When this volume group is created or activated, the memory allocated is based on the extents allocated to logical volumes and not on the maximum size of the volume group. This is different from Version 1.0 volume groups, in that LVM also allocates memory to match the provisioning on disk.

In other words the system memory used by a Version 2.0 volume group depends on how much of the volume group is used. A volume group generously provisioned uses little memory as long as not many extents are allocated to logical volumes. The memory usage grows roughly in proportion to the extents allocated to logical volumes.

As a consequence, you can provision for very large volume groups without wasting system memory.

Metadata Size on Disk versus Maximum Volume Group Size

The following graph illustrates how much space is reserved on disk based on extent size and maximum volume group size.

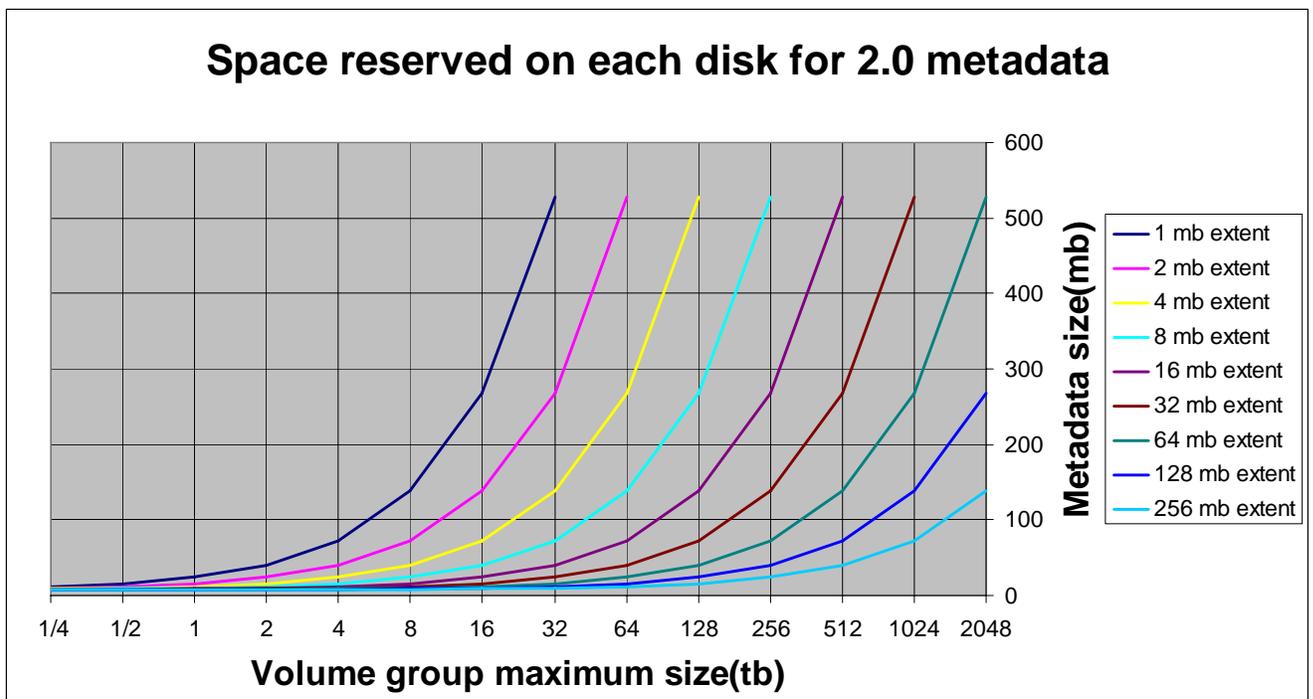


Figure 1

Notes:

- All lines stop at 518 MB on the Y axis because of the limit of 32 million extents per volume group.

What are the Advantages of Over-Provisioning with Version 2.0 Volume Groups?

Over-provisioning does not increase the memory footprint of the volume group. Over-provisioning guarantees you will be able to grow the volume group up to its maximum size without having to create a second bigger volume group and copy the data over.

Version 2.0 Can Leave as Much Space for User Data on Disk as Version 1.0

For an equivalent provisioning (same extent size and same maximum volume group size), the size of the metadata on disk is larger for Version 2.0 compared to Version 1.0. That is because the metadata is provisioned for more logical volumes, physical volumes, more extents per logical volume (or physical volume), and contains additional information such as the volume names.

However, because of the different structures of the Version 2.0 metadata, Version 2.0 may leave as much space for user data as Version 1.0.

Version 2.0 volume groups use the space at the end of the disk that can not be used for user data. Often, the end of the last extent does not coincide with the end of the disk. Version 2.0 volume groups try to use this leftover space between the end of the last extent and the end of the disk to store metadata. As a consequence, even if the on-disk metadata is bigger, the space available on disk for user data may be as large for Version 2.0 compared to Version 1.0.

Example

A Version 1.0 volume group (vgtest11) and a Version 2.0 volume group (vgtest12) are created with the same provisioning (64MB extent and maximum volume group size of 128TB). They are created on disks of identical size to simplify the comparison. While the metadata of the Version 2.0 volume group is bigger, the space available for user data in each volume group is the same: 2050 extents (from *pvdisplay* output).

```
# diskinfo /dev/rdisk/disk148
SCSI describe of /dev/rdisk/disk148:
    vendor: COMPAQ
    product id: MSA1000 VOLUME
    type: direct access
    size: 134399790 Kbytes
    bytes per sector: 512

# diskinfo /dev/rdisk/disk149
SCSI describe of /dev/rdisk/disk149:
    vendor: COMPAQ
    product id: MSA1000 VOLUME
    type: direct access
    size: 134399790 Kbytes
    bytes per sector: 512

# vgcreate -V 1.0 -e 16384 -l 255 -p 128 -s 64 vgtest11 /dev/disk/disk148
(16384*128*64MB = 128TB)

# vgcreate -V 2.0 -s 64 -S 128t vgtest12 /dev/disk/disk149

# pvdisplay /dev/disk/disk148
--- Physical volumes ---
PV Name                /dev/disk/disk148
VG Name                /dev/vgtest11
PV Status              available
Allocatable            yes
VGDA                   2
```

Cur LV	0
PE Size (Mbytes)	64
Total PE	2050

pvdisplay /dev/disk/disk149

--- Physical volumes ---

PV Name	/dev/disk/disk149
VG Name	/dev/vgtest12
PV Status	available
Allocatable	yes
VGDA	2
Cur LV	0
PE Size (Mbytes)	64
Total PE	2050

How to Pick an Extent Size

While the extent size has a similar impact on Version 1.0 and 2.0 volume groups, this section discusses only Version 2.0 volume groups unless “Version 1.0” is explicitly mentioned.

The extent size affects the following:

- The maximum volume group size you can select when creating a volume group.
- The amount of disk space reserved for the metadata.
- The memory footprint used by the activated volume group.
- The I/O performance (for the case where the logical volume is not striped).
- The resynchronization time of the mirrors in some particular cases.
- The minimum size of a logical volume.

As shown in Figure 2 “Maximum volume group size versus extent size” on page 19, the maximum size of a volume group dictates a minimum extent size. You may pick this minimum value or a larger value.

A larger extent size uses less disk space for metadata and reduces the volume group memory footprint. To understand the relationship, see Figure 1 “Space reserved on each disk for 2.0 metadata” on page 15 and “Maximal memory footprint, all extents allocated, max number of LVs and PVs” on page 21.

Let’s take an example. You want a volume group that can grow to 128TB. Refer to Figure 2, “Maximum volume group size versus extent size”. The chart shows that the minimum extent size is 4MB. As a consequence, you can pick an extent size between 4MB and 256MB.

For the amount of disk space used for metadata (on each disk), look at Figure 1 “Space reserved on each disk for 2.0 metadata”. The metadata disk space used is between 520MB for 4MB extents and 20MB for 256MB extents approximately.

For the memory footprint, look at Figure 3 “Maximal memory footprint, all extents allocated, max number of LVs and PVs”. The chart shows about 950MB for 4MB extents and 27MB for 256MB extents. However, these figures represent the full usage of 128TB of the volume group. In the following year you expect to use only half the capacity of the volume group (64TB). Consulting the chart using 64TB as the volume group size instead of 128TB you observe about 500MB for 4MB extents and 20MB for 256MB extents. You can examine the chart for other extents sizes between 4MB and 256MB.

Notes:

- This maximum memory footprint chart gives the worst case memory footprint for a given volume group size and extent size. In addition to using all the extents and the maximum number of logical volumes and physical volumes, a specific configuration of the mirrors is needed to reach this worst case. Most likely the actual footprint in a user configuration will be smaller.
- The impact of the number of logical and physical volumes on the memory footprint is negligible.

For the example, you now have a good idea of how much disk and memory overhead you may incur depending on the extent size.

There is an incentive to use large extents because it reduces the LVM overhead in term of disk space and memory footprint. However, before selecting the extent size, you must consider three other factors.

The most important one is the size of logical volumes. The minimum disk space allocation to a logical volume is one extent. As a consequence, it makes sense to select a small extent size for a volume group containing a lot of very small logical volumes. It avoids wasting disk space if a lot of small logical volumes are needed or used.

Another factor to consider is related to user I/O performance. If you use extent-based striping (*lvcreate -D y*), smaller extents give better throughput. However, you can obtain equal or better I/O throughput with striped logical volumes (*lvcreate -i*). As a consequence you can keep large extents and still get excellent I/O throughput by using striped logical volumes.

The last factor is related to resynchronization time. Small extent size may accelerate the resynchronization of a logical volume in the case where an application issues sparse I/Os on the logical volume at the time a physical volume was down. In this particular case the amount of data to resynchronize would be smaller, and as a consequence the resynchronization would be faster. Note that if a resynchronization of a complete mirror of a logical volume is needed, the size of the extent has no impact on the resynchronization time.

Overall, unless you need a lot of very small logical volumes, it is better to pick a large extent size and use striped logical volumes where you need high performance.

The following table illustrates how extent size limits maximum volume group size.

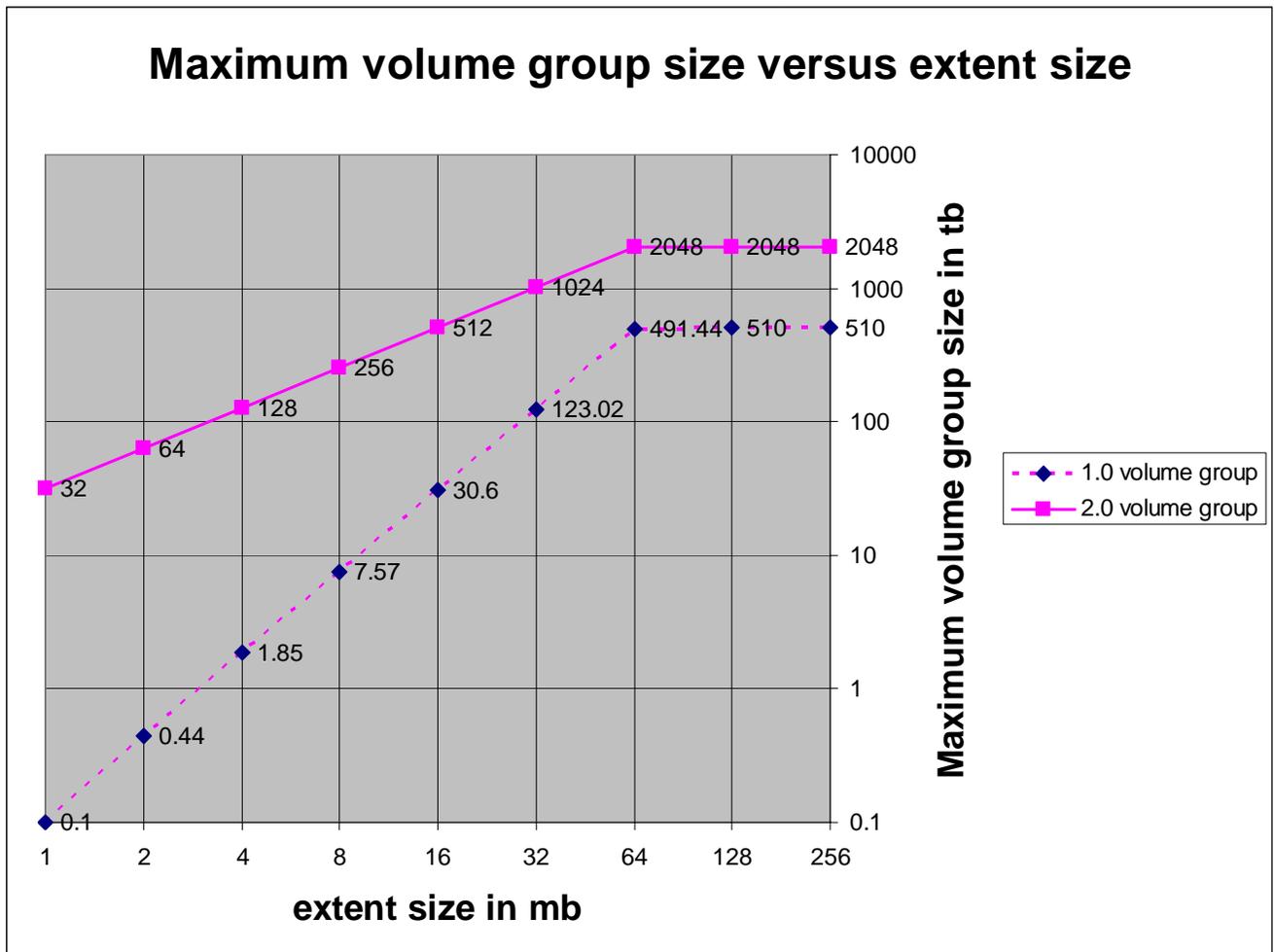


Figure 2

Notes:

- The Version 2.0 graph stops at 2048 TB because of the limit of 32 million extents per volume group.
- The Version 1.0 graph stops at 510 TB because the metadata must fit in one extent.

Memory Footprint

The memory footprints given in this chapter are estimates.

The memory footprint of a Version 2.0 volume group provisioned for a given extent size and maximum volume group size is somewhere between the minimum and the maximum memory footprint as shown in Figures 3 and 4.

The Version 2.0 volume group memory footprint starts very small when the volume group is created (independent of the provisioned maximum size of the volume group) and grows roughly proportional to the number of extents allocated to logical volumes.

Let's take the example of a Version 2.0 volume group provisioned for 128TB with an extent size of 64MB. Its minimum footprint is 2MB and its maximum memory footprint is 73MB. The memory footprint of this volume group is 2MB when the volume group is created and may eventually go all the way up to 73MB.

Minimum Memory Footprint

The graph below shows the memory footprint of Version 1.0 and 2.0 volume groups when activated and with no extents yet allocated to logical volumes. This is called the "minimum footprint". That would be the memory footprint just after a `vgcreate` for example.

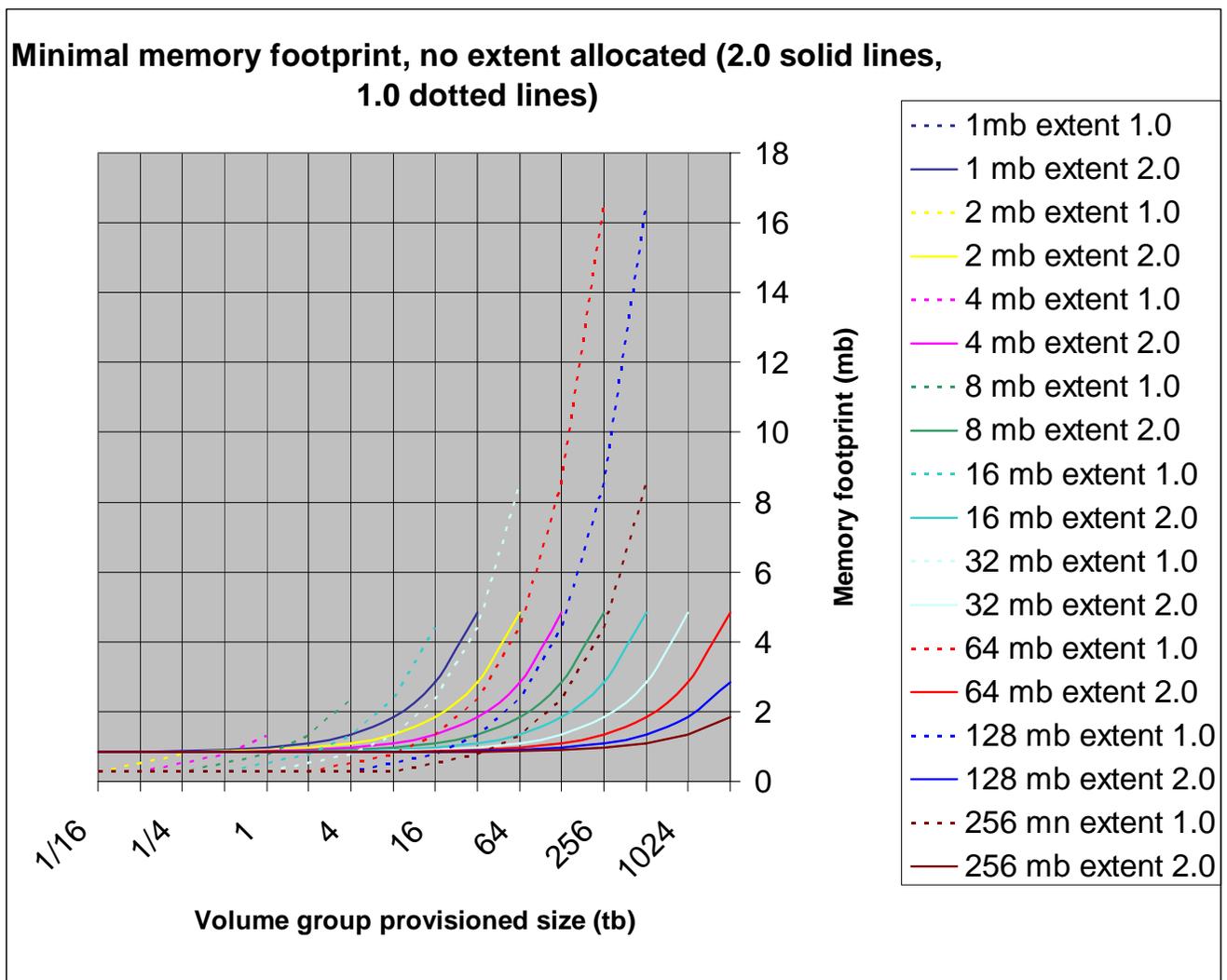


Figure 3

A Version 2.0 volume group allows the user to provision a very large size, and yet the minimal memory footprint is small. For example the minimal memory footprint of a Version 2.0 volume group provisioned for 2PB with 64MB extents is around 5MB while the minimal memory footprint of a Version 1.0 volume group with the same extent size provisioned for only 256TB (8 times smaller) is around 17MB.

The minimum footprint of a Version 2.0 volume group whatever its provisioned size doesn't go above 5MB.

Maximum Memory Footprint

The graph below shows the memory footprint of Version 1.0 and 2.0 volume groups activated, fully populated, and with the worst case volume group configuration. The worst case volume group configuration in terms of memory footprint is when all the extents of the volume group are allocated to logical volumes, the maximum number of logical volumes and physical volumes are used and everything is mirrored. Most of the time users will not reach the maximum footprint event even if all the extents are allocated to logical volumes.

Note that this comparison is not one to one, between Version 1.0 and 2.0 because the volume groups in this chapter charts are assumed to contain the maximum number of logical and physical volumes. And the maxima for Version 2.0 are higher. In other words the memory footprint of a Version 2.0 volume group includes more volumes than Version 1.0. For an exact comparison, see "Comparison 1.0 versus 2.0 for a same number of logical and physical volumes". However the number of logical or physical volumes does not make much difference in the memory footprint.

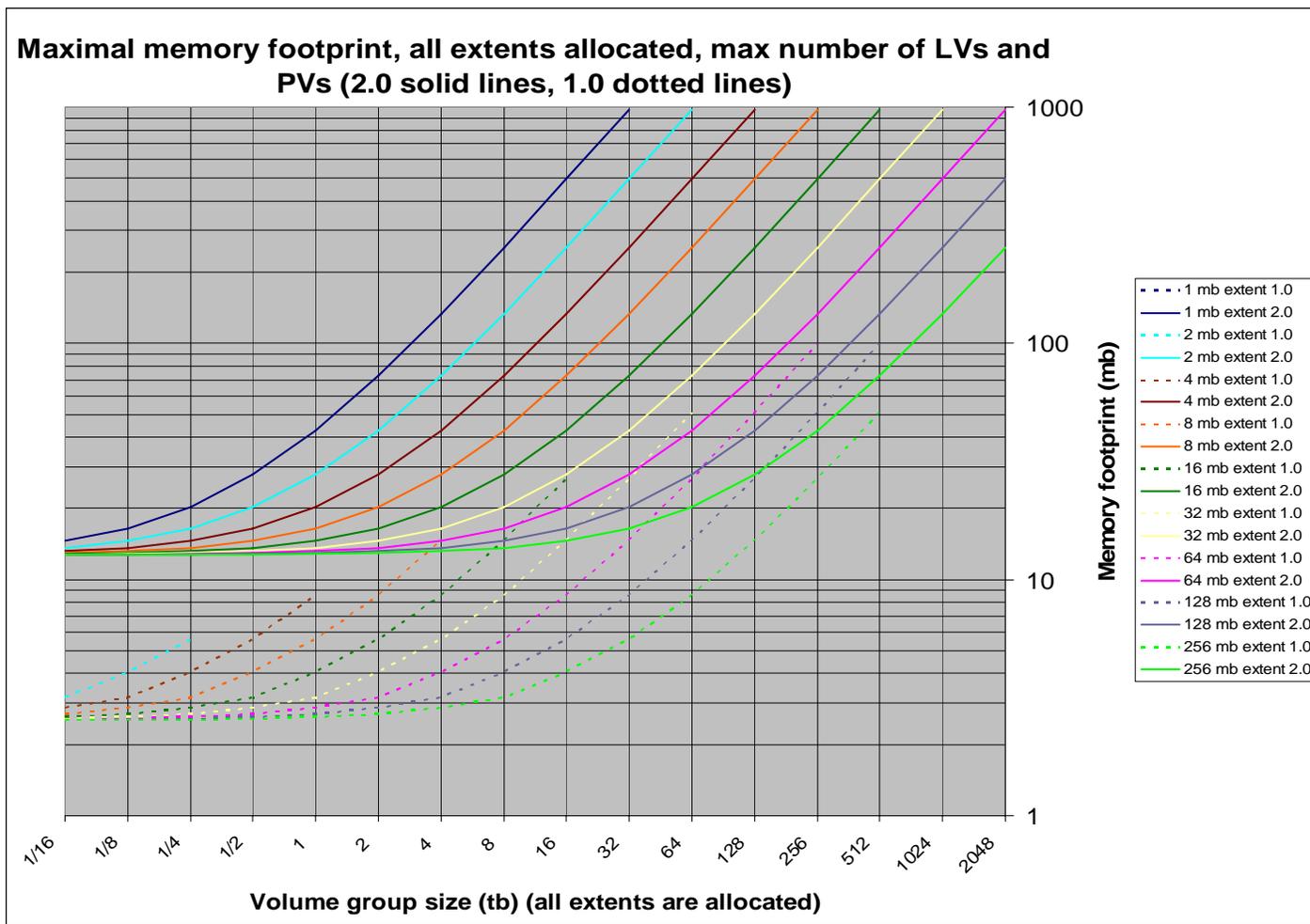


Figure 4

Notes:

- For a given extent size, the Version 2.0 line goes further right compared to the Version 1.0 line because a Version 2.0 volume group can be bigger.
- The Version 2.0 lines don't go further up and right because of the limit of 32 million extents per volume group.

For the same volume group size and same extent size, a Version 2.0 volume group uses about 10MB more than Version 1.0 because of two reasons:

1. A Version 2.0 volume group at the maximum contains 511 logical volumes and 511 physical volumes while a Version 1.0 volume group is limited to only 255 logical volumes and 255 physical volumes.
2. This 10MB increase enables the removal of the Version 1.0 limitation (logical volume size, number of extents per physical volume or logical volumes, and so forth) and enables new features such as storing the logical volume name in the metadata.

Comparison of Memory Footprints for the Same Number of Logical and Physical Volumes

This comparison uses 255 logical volumes and 10 or 100 physical volumes per volume group. The first chart uses an extent size of 8MB while the second uses 64MB.

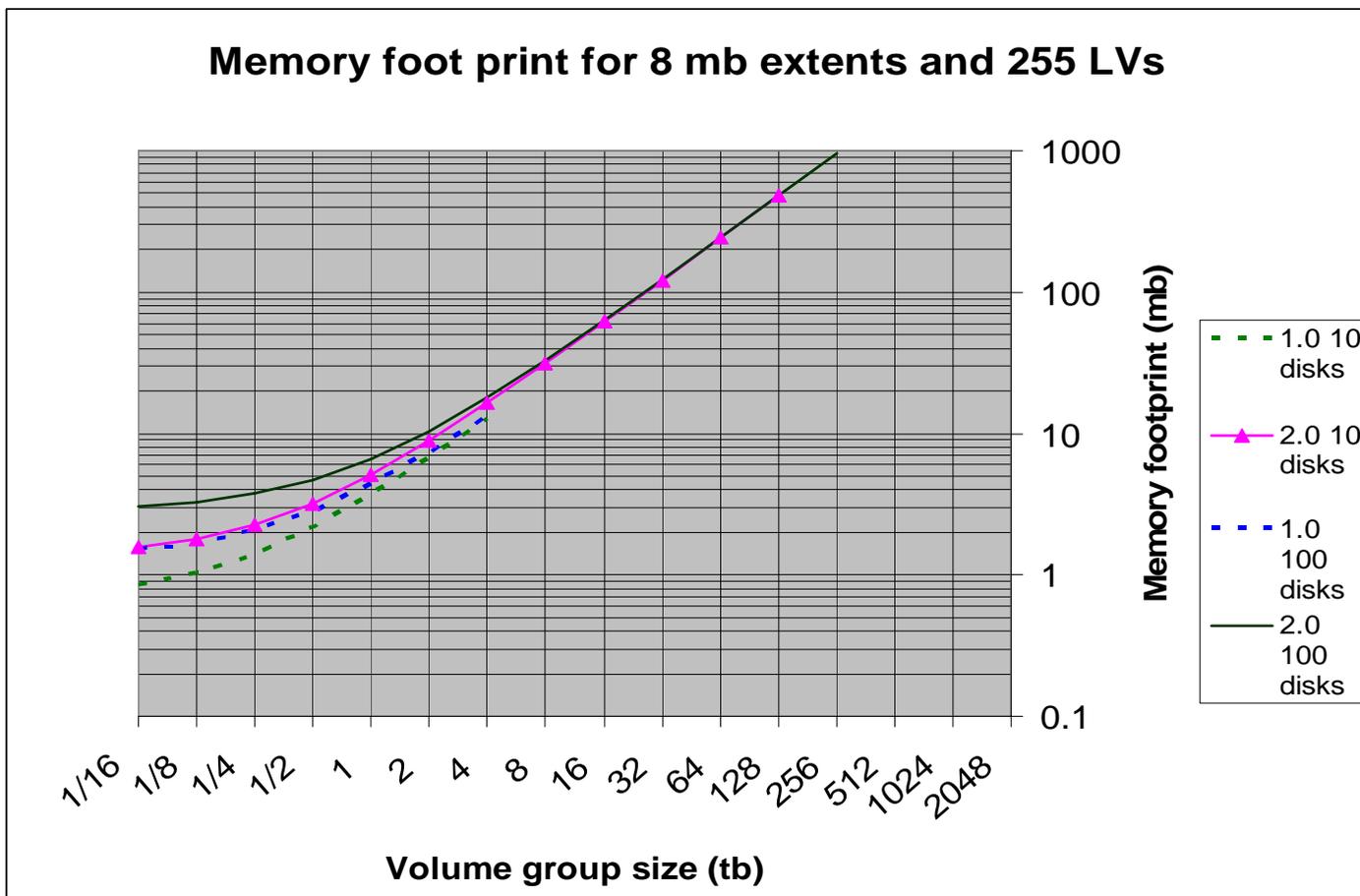


Figure 5

Memory foot print for 64 mb extents and 255 LVs

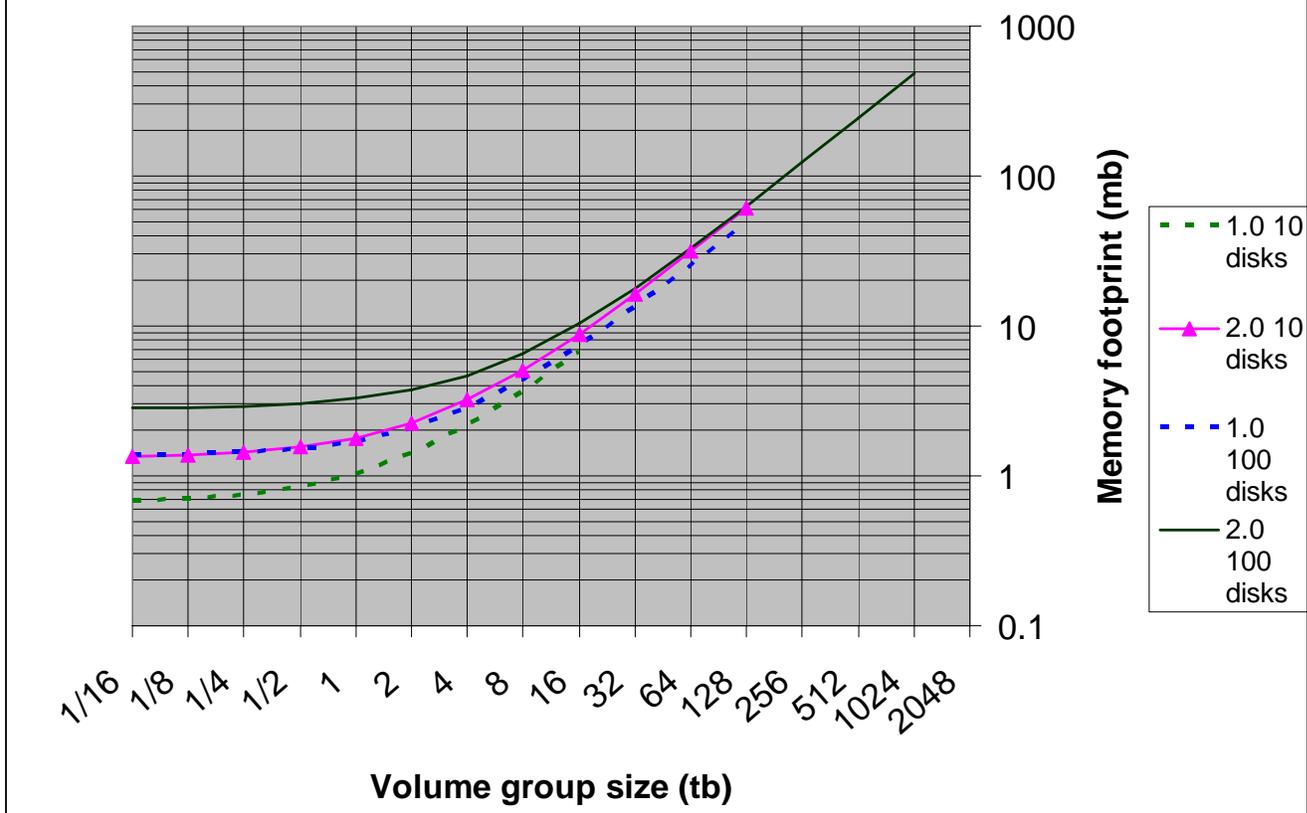


Figure 6

How to Configure the New Limits

Number of Logical Volumes or Physical Volumes in a Volume Group

There is nothing to do to on a Version 2.0 volume group (at creation time or later) to reach the new limits for the number of physical volumes (511) or logical volumes (511). As indicated earlier, Version 2.0 volume groups are provisioned in a way that LVM can always add physical volumes or logical volumes until the maximum capacity of the volume group is reached. Just keep adding physical volumes (*vgextend*) and logical volumes (*lvcreate*) as needed.

Maximum Size of the Volume Group

The size limit is displayed in the chart of the section "[How to Pick an Extent Size](#)". You set the maximum size of the volume group when you create it (*vgcreate* option *-S*). If the volume group is not provisioned sufficiently, and you wish to go beyond the volume group maximum size that was provisioned, with March 2008 HP-UX 11i v3 release, you must copy the data to a different volume group that is larger.

For example if you provision a volume group for 32TB with an extent size of 4MB (*vgcreate -v 2.0 -s 4 -S 32t myvg /dev/disk/disk149*), you can not grow the volume group up to the limit allowed by Version 2.0, that is 128 MB for 4 MB extent. You can grow the volume group only up to 32 TB.

Size of Logical Volume or Number of Mirrors

To reach the new limits just increase the values passed in to *lvcreate* or *lvextend* or *fsweb*.

Migration Between Version 1.0 and Version 2.0

Why Would You Migrate?

You might want to migrate from Version 1.0 to Version 2.0 because your Version 1.0 volume group is too small. You can use *vgmodify* to extend the limits of a Version 1.0 volume group but that may not work or may not be enough.

Migrating from Version 2.0 to Version 1.0 is impossible as soon as the Version 1.0 limits are exceeded. For example a Version 2.0 volume group containing 256 logical volumes (one more than the Version 1.0 limit) cannot be migrated to Version 1.0.

How to Migrate

At this time, there is no in-place migration available. You must migrate by copying data (with *dd* for example) to another volume group created with the desired version.

You face the same situation when you can't increase the size of a Version 1.0 volume group and *vgmodify* cannot do it either. You must copy over the data to a new volume group.

To speed up the copy you can do several *dd* commands in parallel on different volume groups or different parts of volume groups.

When to Migrate

Because the migration is done by copy, a good time to migrate is when you replace your storage.

High Availability

With Version 2.0 volume groups, you now have the ability to provide additional copies of user data, thus increasing the availability of your data. Version 1.0 volume groups are limited to two mirrors (three replicas of the data). Version 2.0 volume groups can have up to five mirrors (six replicas of the data).

When Does It Make Sense to Use More Than Two Mirrors?

Using more than three replicas of the data makes sense as part of a setup for disaster recovery. Below are two examples of a logical volume whose data is replicated six times. The first example is a two site setup (Figure 7). In this case, if one site goes down the data stays highly available on the surviving site because it is still replicated three times on that site.

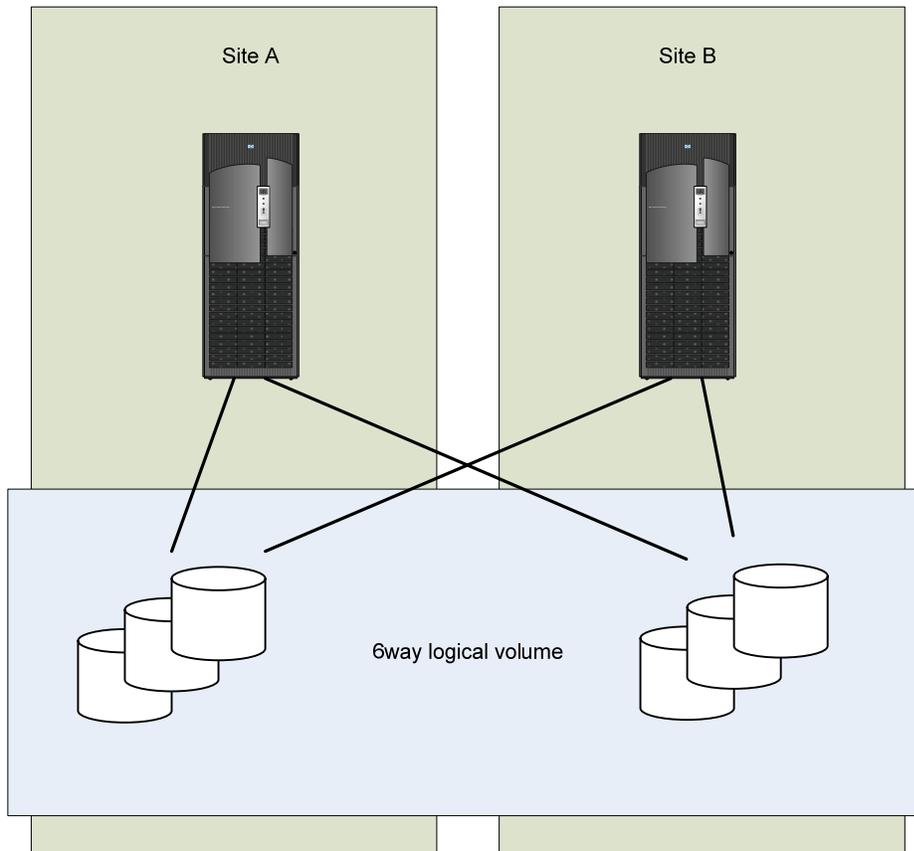


Figure 7

The advantage of having three mirrors at one site is if a backup strategy is deployed where a mirror logical volume is split off (*lvsplit*) to perform backups. In this example, either site can employ a backup strategy that would not jeopardize data availability during a site outage.

A second example is illustrated in Figure 8 where a three node cluster is deployed. In this example, each node could be geographically located together with two mirrors of data. This provides data redundancy at each site in the event of lost connectivity between the sites.

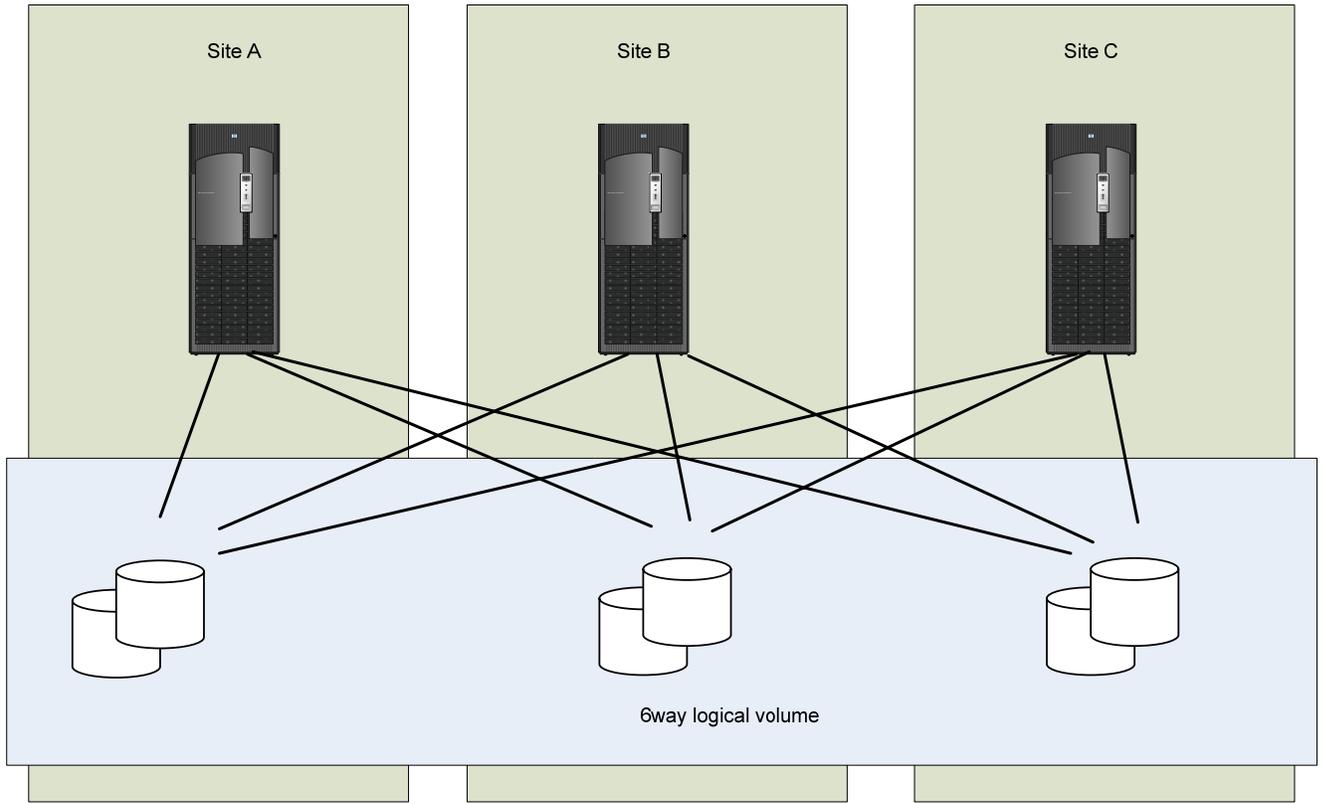


Figure 8

Glossary

DSF

Device Special File. A file associated with an I/O device. DSFs are read and written the same as ordinary files, but requests to read or write are sent to the associated device.

LUN

A SCSI logical unit. This refers to an end storage device such as a disk, tape, floppy, or CD. This is the logical unit itself and does not represent the path to the logical unit.

Metadata

The on-disk structures that LVM uses to manage a volume group. This space is not available for application data.

KB

A kilobyte unit of information equal to 2^{10} or 1024 bytes

MB

A megabyte unit of information equal to 2^{20} or 1,048,576 bytes

GB

A gigabyte unit of information equal to 2^{30} or 1,073,741,824 bytes

TB

A terabyte unit of information equal to 2^{40} or 1,099,511,627,776 bytes

PB

A petabyte unit of information equal to 2^{50} or 1,125,899,906,842,624 bytes

For More Information

To learn more about some of the LVM features, see the following document on HP documentation website:

<http://docs.hp.com> (Use search with the given name of the whitepaper)

<http://www.docs.hp.com/en/oshpux11iv3#LVM%20Volume%20Manager>

- LVM Online Disk Replacement (LVM OLR)
- LVM Volume Group Dynamic LUN expansion (DLE)/vgmodify
- LVM Volume Group Quiesce/Resume
- HP-UX LVM Performance Assessment (The whitepaper will be available soon)
- SLVM Single-Node Online Reconfiguration (SLVM SNOR)
- HP-UX System Administrator's Guide Logical Volume Management
- When Good Disks Go Bad: Dealing with Disk Failures under LVM

Call to Action

HP welcomes your input. Please give us comments about this whitepaper, or suggestions through our technical documentation feedback website: <http://docs.hp.com/en/feedback.html>.

© 2008 Hewlett-Packard Development Company, L.P. The information contained herein is subject to change without notice. The only warranties for HP products and services are set forth in the express warranty statements accompanying such products and services. Nothing herein should be construed as constituting an additional warranty. HP shall not be liable for technical or editorial errors or omissions contained herein.