



Reliability of systems and components

GNS Systems GmbH
Am Gaußberg 2
38114 Braunschweig / Germany

Tue 19 Jul 2011, Revision 0.1

Please send comments about this document to: norman.krebs@gns-systems.de

Contents

- 1 Business continuity 3**
- 2 Down times 3**
- 3 Availability 4**
- 4 Reliability of RAID 5**
 - 4.1 Annualized Failure Rate 5
 - 4.2 RAID 0 - Stripe 6
 - 4.3 RAID 1 / RAID 5 6
 - 4.4 RAID 6 7
 - 4.5 RAID 01 7
 - 4.6 RAID 10 7
 - 4.7 RAID 01 vs. RAID 10 8
 - 4.8 Estimation of Failure 9
 - 4.9 Bit Error Rate 10
- 5 RTO and RPO 11**
 - 5.1 Sample Recovery Levels 11
- References 12**

The reliability theory and practice consists of several objectives. This document mentions a few of them to show some key words and calculations quickly.

1 Business continuity

Business continuity consists of:

High availability - keep the IT infrastructure running in case of minor outages

Continuous operations - find a way to accomplish maintenance and backups without or with less disturbance of the running processes

Disaster Recovery - keep the IT infrastructure running or get it running again quickly after major outages

2 Down times

availability	down time/year
98%	7.3 days
99%	3.65 days
99.5%	1.83 days
99.9%	8.76 hours
99.99%	52.6 minutes
99.999%	5.26 minutes
99.9999%	31.5 seconds

3 Availability

short	unit	description
⚡		failure
👍		recovery
$MTTR$	[h]	mean time to return
$MTTF$	[h]	mean time to failure
$MTBF$	[h]	mean time between failure
A	probability	availability
A_p	probability	availability - parallel components
A_s	probability	availability - serial components

Availability as a function of times:

$$A = \frac{MTTF}{MTTF + MTTR} \quad (1)$$

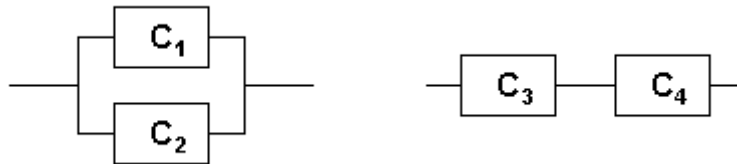


Figure 1: parallel (C_1/C_2) and serial (C_3/C_4) components

The components C_1 and C_2 with the availabilities A_1 and A_2 are installed parallelly and have the combined availability A_p :

$$A_p = 1 - ((1 - A_1) \cdot (1 - A_2)) \quad (2)$$

The components C_3 and C_4 with the availabilities A_3 and A_4 are installed serially and have the combined availability A_s :

$$A_s = A_3 A_4 \quad (3)$$

4 Reliability of RAID

This section contains the equations to calculate the MTTF-rates for several RAID levels. This is needed to calculate the probability of failing within a special time interval.

short	unit	description
$MTTF$	[h]	mean time to failure; this value is given by the hardware vendor. He gains it from a mix of knowledge and experience and the ruthless lies of the marketing. It means the failure of the whole disk and is now within the range of 500000 hours.
$MTTR$	[h]	mean time to return; it consists of the whole time needed to rebuild a proper array $MTTR = MFDT + MHRT + MTDR + MTRA$
$MFDT$	[h]	mean failure detection time; the time needed to detect a failure
$MHRT$	[h]	mean human reaction time; the time needed to bring a skilled person to handle the problem
$MTDR$	[h]	mean time to replace the disk
$MTRA$	[h]	mean time to rebuild the array
$A(x)$		probability of failure within x years
N		number of drives in the array
$PoH(x)$	[h]	power on hours per x years
AFR	$\frac{\%}{year}$	annualized failure rate ($[\%] = [\frac{1}{100}]$:keep this in mind)
K_1		correction factor for the MTTF of aged disks - first failure
K_2		correction factor for the MTTF of aged disks - second failure

4.1 Annualized Failure Rate

The **AFR** is the probability of failure for a component. It is the most important value for the user.

$$AFR_{disk} = \frac{100 \cdot PoH(1)_{disk}}{MTTF_{disk}} \quad (4)$$

$$MTTF_{disk} = \frac{100 \cdot PoH(1)_{disk}}{AFR_{disk}} \quad (5)$$

4.2 RAID 0 - Stripe

A RAID 0 is defined as a set of at least two striped disks. There is no redundancy and no safety but its the fastest.

$$MTTF_{R0} = \frac{MTTF_{disk}}{N} \quad (6)$$

$$\begin{aligned} A(1)_{R0} &= \frac{PoH(1)_{R0}}{MTTF_{R0}} \\ &= \frac{PoH(1)_{R0} \cdot N}{MTTF_{disk}} \end{aligned} \quad (7)$$

4.3 RAID 1 / RAID 5

RAID 1 is defined as a set of two mirrored drives. RAID 5 is defined as a set of N disks with parity blocks on all disks. The volume size is $N - 1$ disks. One RAID stripe consists of N blocks distributed over the disks. One block contains the parity. RAID 5 writes slowly. The change of a block requires:

- $N - 2$ reads to get the informations to recalculate the parity
- 2 writes: the data block itself + the parity block

The disks in a RAID are usually bought from one production charge and may have the similar failures. They have got the same influences of transportation, temperature, falling down etc. Thus the probability is high, that after the first failure a second one happens. To represent this fact the correction factors can reduce the remaining MTTF: p.e.: $K_1 = 0.1$ and $K_2 = 0.01$.

$$\begin{aligned} MTTF_{R15} &= \underbrace{\frac{MTTF_{disk}}{N}}_{1st\ failure} \cdot \underbrace{\frac{MTTF_{disk} \cdot K_1}{(N-1) MTTR}}_{2nd\ failure} \\ &= \frac{MTTF_{disk}^2 \cdot K_1}{N \cdot (N-1) MTTR} \end{aligned} \quad (8)$$

$$A(1)_{R15} = \frac{PoH(1)_{R15}}{MTTF_{R15}} \quad (9)$$

4.4 RAID 6

RAID 6 is defined as a set of N disks. The redundant information or parity is stored on all disks. The volume size is $N - 2$ disks. The change of a block requires:

- $N - 3$ reads to get the informations to recalculate the parity
- 3 writes: the data block itself + two parity blocks

$$\begin{aligned}
 MTTF_{R6} &= \underbrace{\frac{MTTF_{disk}}{N}}_{1st\ failure} \cdot \underbrace{\frac{MTTF_{disk} \cdot K_1}{(N-1) MTTR}}_{2nd\ failure} \cdot \underbrace{\frac{MTTF_{disk} \cdot K_2}{(N-2) MTTR}}_{3rd\ failure} \\
 &= \frac{MTTF_{disk}^3 \cdot K_1 \cdot K_2}{N(N-1)(N-2) MTTR^2}
 \end{aligned} \tag{10}$$

$$A(1)_{R6} = \frac{PoH(1)_{R6}}{MTTF_{R6}} \tag{11}$$

4.5 RAID 01

The RAID 01 is defined as a mirror of striped disks.

$$\begin{aligned}
 MTTF_{R01} &= \underbrace{\frac{2MTTF_{disk}}{N}}_{1st\ failure} \cdot \underbrace{\frac{2MTTF_{disk} \cdot K_1}{N \cdot MTTR}}_{2nd\ failure} \\
 &= \frac{4MTTF_{disk}^2 \cdot K_1}{N^2 \cdot MTTR}
 \end{aligned} \tag{12}$$

4.6 RAID 10

The RAID 10 is defined as a stripe of mirrored disks.

$$\begin{aligned}
 MTTF_{R10} &= \underbrace{\frac{MTTF_{disk}}{2}}_{1st\ failure} \cdot \underbrace{\frac{MTTF_{disk} \cdot K_1}{MTTR}}_{2nd\ failure} \cdot \frac{2}{N} \\
 &= \frac{MTTF_{disk}^2 \cdot K_1}{N \cdot MTTR}
 \end{aligned} \tag{13}$$

4.7 RAID 01 vs. RAID 10

$$MTTF_{R10} \sim \frac{1}{N} \quad (14)$$

$$MTTF_{R01} \sim \frac{4}{N^2} \quad (15)$$

The table shows, depending on the number of disks, the increasing probability of failure of a RAID 01 against a RAID 10.

N	$A_{R10} \sim$	$A_{R01} \sim$
4	4	4
6	6	9
8	8	16
10	10	25
12	12	36
20	20	100

⇒ Avoid RAID 01!

4.8 Estimation of Failure

short	unit	description
r^*	[pcs]	number of failing drives
N	[pcs]	number of drives
r	[%]	relative number of failing drives
$f(t)$	[pcs/h]	relative number of failing drives per time interval
λ	[h ⁻¹]	failure rate
$F(t)$	probability	probability of failure
$R(t)$	probability	probability of survival

While the mechanism of failing can be described by a exponential or Weibull distribution, we can use the following expressions:

The relative number of failing drives is:

$$r = \frac{r^*}{N} \quad (16)$$

The number of failing drives per time interval, or the density function is:

$$f(t) = \frac{dr}{dt} \quad (17)$$

The probability of failure is:

$$F(t) = \int f(t)dt \quad (18)$$

The probability of survival is:

$$R(t) = 1 - F(t) \quad (19)$$

The failure rate λ is:

$$\lambda = \frac{f(t)}{R(t)} \quad (20)$$

The percental rate of failing disks within a time interval is:

$$R(t) = e^{-(\frac{t}{MTTF})^b} \quad (21)$$

$$R(MTTF) = e^{-(\frac{MTTF}{MTTF})} = e^{-1} = \frac{1}{e} = 0.368 \quad (22)$$

This tells us that $\approx 36.8\%$ of the drives will survive the MTTF period.

4.9 Bit Error Rate

short	unit	description
N	[pcs]	number of disks
M_{disk}	[bit]	amount of bits in one disk
S_{disk}	[byte]	hard disk size
M_{RAID}	[bit]	amount of bits in the RAID
BER	[bit]	probability of one uncorrectable bit read error. It is in the range of 10^{-14} to 10^{-15} - One bit error every 10-100TB read.

The BER is value given by the vendor for one harddrive. It contains all sources of bit failures like the disks themselves, the electronic and mechanic components of the harddrive and the assumed influences from the outer space around the harddrive.

The number of bits in the RAID:

$$M_{RAID} = S_{disk} \cdot 8 \cdot N \quad (23)$$

If one drive fails in a RAID, all data have to be read to restore parity. The probability to read all bits in the RAID without error during a restore after losing one drive:

$$A_{rRX} = (1 - BER)^{M_{disk} \cdot N} \quad (24)$$

The probability to lose bits:

$$A_{be} = 1 - (1 - BER)^{M_{disk} \cdot N} \quad (25)$$


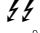

The probability to lose bits depends on the amount of bits read and written, not on time. So we cannot calculate the probability of bit loss within a time interval without knowing the amount of data shoveled around in this time.

Example:

$$\begin{aligned} BER &= 10^{-14} \\ M_{RAID} &= 10^{14} \\ \rightarrow A_{rRX} &\approx 0.36 \\ \rightarrow A_{be} &\approx 0.64 \end{aligned}$$

5 RTO and RPO



-  last backup
-  failure
-  recovery
- RTO* Recovery Time Objective
- RPO* Recovery Point Objective

RPO maximum tolerable data loss (time since last backup)
RTO time in minutes needed from failure to recover

Both objectives have to be considered separately.

5.1 Sample Recovery Levels

	<i>RTO</i> [h]	<i>RPO</i> [h]
Immediate Recovery	≤ 4	≈ 0
Rapid Recovery	≤ 48	< 24
Medium Recovery	≤ 120	< 24
Best-Effort Recovery	> 120	< 48

References

- [1] D.A.PATTERSON,G.GIBSON,R.KATZ *A case for redundant arrays of inexpensive disks (RAID)* UC Berkeley, Dept. of Electrical Engineering and Computer Sciences, 1988
- [2] SCOTT SPEAKS *Reliability and MTBF Overview* VICOR Reliability Engineering, 2004
- [3] HOLGER WILKER *Weibull-Statistik in der Praxis* Lauffen, 2004