

RAID+: Breaking Resource Isolation between Multiple Co-located Volumes for Larger-scope Optimization

Paper #145

1. Proofs on RAID+'s Properties

Below are the major properties associated with MOLS-based normal and interim data layouts. Some of them can be easily derived from the definitions and theorems about MOLS. For the others, we show the proof sketches here.

Property 1. *With normal data layout, any two blocks within a data stripe are placed on separate disk drives.*

Proof. Take two arbitrary blocks within a stripe, which are placed on disks x and y respectively. According RAID+ normal data layout definition, x and y are the elements of two Latin squares L_a and L_b at the same location (i, j) , where $a \neq b$ and $i > 0$.

From Theorem 2, we have $x = f_a(i, j) = a \times i + j$, and $y = f_b(i, j) = b \times i + j$. Given that $a \neq b$ and $i > 0$, we have $x \neq y$. \square

Property 2. *With normal data layout, the n disks involved are assigned equal shares of both data and parity blocks.*

Analysis. The ordinal numbers of disks, on which the k blocks are placed, come from the same locations of the k selected Latin squares. Since the first $k - 1$ blocks within a stripe are data blocks and the last block is parity block, the first $k - 1$ Latin squares are attached with the distribution of data blocks and the k^{th} Latin square is attached with the distribution of parity blocks.

Within any one of the k Latin squares, all elements except those in the first row are used for designing normal data layout. According to Definition 1, each number appears exactly once in each row of a Latin square. It is obvious that each number appears $n - 1$ times in the $n - 1$ rows. According to each of the first $k - 1$ selected Latin squares, RAID+ put $n - 1$ data blocks on each disk. According to the k^{th} selected Latin square, RAID+ put $n - 1$ parity blocks on each disk. Therefore, RAID+ places $(n - 1) \times (k - 1)$ data blocks on each disk, and $n - 1$ parity blocks on each disk. \square

Property 3. *With the n -disk normal layout, all those blocks correlated to blocks on any given drive (i.e., blocks sharing stripes with blocks on this disk) are distributed evenly among the other disks.*

Analysis. For any one drive d ($d \in [0, n - 1]$), all those blocks correlated to this drive must appear within a correlated stripe whose mapping tuple includes the element d . Let us first inspect the set S_i of all the correlated stripes, the i^{th} ($i \in [0, k - 1]$) element of their mapping tuples is d . According to the data distribution principle of normal data layout, the number of the correlated stripes is $n - 1$ and the $n - 1$ d s come from different locations of the i^{th} selected Latin square.

Then, we examine what disks the $n - 1$ blocks, represented by the j^{th} ($j \in [0, k - 1]$ and $j \neq i$) elements of the $n - 1$ mapping tuples in the set S_i , are placed on. According to Theorem 3, the $n - 1$ blocks are placed on all the $n - 1$ drives except Disk d . Further, by traversing i and j in turn, we will easily obtain that for Disk d , all its correlated blocks are distributed evenly among all the other drives. Furthermore, the number of the correlated blocks placed on each other disk is $k \times (k - 1)$. \square

Property 4. *With the $(n - 1)$ -disk interim layout, any two blocks within a data stripe are still to be placed on separate disk drives.*

This property can be obtained easily similar to the analysis of Property 1.

Property 5. *All the $(n - 1)k$ missing blocks on any single failed disk can be redistributed to all the surviving $(n - 1)$ disks evenly, each receiving k additional blocks.*

Analysis. k selected Latin squares are used to construct normal data layout, and the $(k + 1)^{th}$ Latin square is selected from the left $n - k - 1$ orthogonal Latin squares for redistributing all the blocks on any one drive d ($d \in [0, n - 1]$) onto all the other drives. All those blocks on this drive must appear within a stripe whose mapping tuple includes the element d .

Let us first inspect the set S_i of all the stripes, the i^{th} ($i \in [0, k - 1]$) element of their mapping tuples is d . According

to the data distribution principle of normal data layout, the number of the stripes in S_i is $n - 1$ and the $n - 1$ ds come from different rows of the i^{th} selected Latin square. The j^{th} ($j \in [1, n - 1]$) one of the $n - 1$ stripes can be represented by the tuple $(a_0, a_1, \dots, a_i = d, \dots, a_{k-1})$. We assume that θ is in the $(k + 1)^{th}$ Latin square, and has the same location with a_i in the i^{th} Latin square. Then, the data block corresponding to this stripe and on Disk d can be placed on Disk θ . The new tuple for this stripe changes into $(a_0, a_1, \dots, a_{i-1}, \theta, a_{i+1}, \dots, a_{k-1})$. According to Definition 3, all these elements remain different. In other words, no two blocks within a stripe are placed on the same drive.

According to Theorem 3, when we traverse j from 1 to $n - 1$, the corresponding $n - 1$ blocks are placed on all the $n - 1$ drives other than Disk d . When we traverse i from 0 to $k - 1$, each one of the other disks hold k blocks from Disk d . \square

Property 6. In a RAID+ system tolerating dual-disk failures, once two disks fail simultaneously, the ratio, between the numbers of the stripes with two lost blocks and those with one lost block, is $(k - 1) : 2 \times (n - k)$.

Analysis. Let us inspect a data template in normal data layout. According to Property 3, $k \times (k - 1)$ of the blocks correlated to a given disk are placed on another disk. Therefore, once two disks fail simultaneously, the number of the stripes with two lost blocks is $k \times (k - 1)$. Within a data template, there is $(n - 1) \times k$ blocks on each disk. So, a disk failure results in $(n - 1) \times k$ degraded stripes. Except $k \times (k - 1)$ stripes with two lost blocks, there is $(n - k) \times k$ stripes with one lost block. Considering that two disks fail simultaneously, there is $2 \times (n - k) \times k$ stripes with one lost block in total. As a result, the ratio between them is $(k - 1) : 2 \times (n - k)$. \square

2. Reliability

2.1 Calculation Method of MTTDL

Here we follow the conventional way of calculating the MTTDL (mean time to data loss) of RAID-5, RAID-6, and RAID+ (i.e., RAID-x) using Continuous-time Markov Chains (CTMC) [1–3].

By assuming both disk failures and repairs follow the exponential distribution with rate μ and ν respectively, we first generate the state transition graph and the generator matrix \mathbf{Q} of RAID-x's CTMC, where the state number is denoted as w . In the generator matrix, each element $q_{i,j}$ is defined as the transition speed from state i to state j . Then from the generator matrix, we need to get the embedded Markov chain of RAID-x's CTMC. Note that this embedded Markov chain is a discrete time Markov chain (DTMC). In this embedded Markov chain, the transition probability $p_{i,j}$ from state i to j can be calculated by Equation 1.

$$p_{i,j} = \frac{q_{i,j}}{\sum_{k \neq i} q_{i,k}} \quad (1)$$

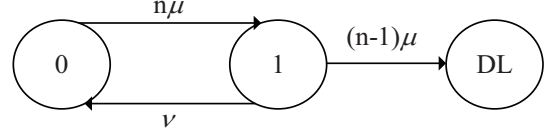


Figure 1. State transition graph of RAID5

As long as we get the state transition probability, the probability transition matrix \mathbf{P} can be formed as $p_{i,j}$ is the value at row i and column j of the matrix. Denote that $\pi_i^e(x)$ is the probability that system will be at state i during step x in embedded Markov chain. These w $\pi_i^e(x)$, $i \in [0, w)$, form the vector $\pi^e(x)$. In matrix \mathbf{P} , the state DL is an absorbing state, so each state in embedded Markov chain has stationary probabilities. Denote the stationary probability vector as π^e and stationary probability for state i as π_i^e .

Take n -disk RAID-5 as an example, we generate the state transition graph, matrix \mathbf{Q}_{RAID-5} , matrix \mathbf{P}_{RAID-5} of RAID-5.

As depicted in Figure 1, the canonical RAID-5's CTMC contains three states, the state 0 is the normal state where n disks are working, the state 1 is the state for one disk failure, and the state DL is the state for data loss. As previous assumption, the disk failures follow the exponential distribution with rate μ . When one of n disk fails, the system will change to state 1. When another one of $n - 1$ disks fails, the system will change to state DL . Thus the state transition from 0 to 1 follows the exponential distribution with rate $n\mu$, while the state transition from 1 to DL follows the exponential distribution with rate $(n - 1)\mu$. As the disk repair makes the system change its state from 1 to 0, the state transition from 1 to 0 follow the exponential distribution with rate ν . Therefore, as shown in Equation 2, the variables in \mathbf{Q}_{RAID-5} can be determined as $q_{0,1} = n\mu$, $q_{1,0} = \nu$, $q_{1,2} = (n - 1)\mu$.

$$\mathbf{Q}_{RAID-5} = \begin{bmatrix} -n\mu & n\mu & 0 \\ \nu & -\nu - (n - 1)\mu & (n - 1)\mu \\ 0 & 0 & 0 \end{bmatrix} \quad (2)$$

As the transition probability of embedded Markov chain can be formulated as Equation 1, we can get the probability transition matrix of embedded Markov chain \mathbf{P}_{RAID-5} in Equation 3, where $p_{i,j}$ is the value at row i and column j of the matrix. So, in Equation 3, we have $r_1 = \frac{\nu}{(n-1)\mu + \nu}$ and $s_1 = 1 - r_1$.

$$\mathbf{P}_{RAID-5} = \begin{bmatrix} 0 & 1 & 0 \\ r_1 & 0 & s_1 \\ 0 & 0 & 1 \end{bmatrix} \quad (3)$$

After the model is determined, we begin to calculate MTTDL of RAID-x. In embedded Markov chain, we denote $E(t_i^c)$ as the expected holding time of state i before the system goes into stationary. As described in equation 4, the

MTTDL is the sum of $E(t_i^c)$, where $i \in [0, w-1]$. Note that the state DL doesn't have holding time as it's the dead state.

$$MTTDL_{RAID-x} = \sum_{i=0}^{w-2} E(t_i^c) \quad (4)$$

Now we need to get the $E(t_i^c)$ of each state except state DL . As the $E(t_i^c)$ is expected holding time before the system goes into stationary, it can be calculated using Equation 5. In Equation 5, $\pi_i^e(x)$ stands for the probability that the system is at state i in step x , and t_i is the average holding time of state i in one step. With distributive law, $E(t_i^c)$ can be regarded as the multiplication of t_i and $\sum_{x=0}^{\infty} \pi_i^e(x)$ and we denote $\sum_{x=0}^{\infty} \pi_i^e(x)$ as σ_i .

$$E(t_i^c) = \sum_{x=0}^{\infty} t_i \pi_i^e(x) = t_i \sum_{x=0}^{\infty} \pi_i^e(x) = t_i \sigma_i \quad (5)$$

Here, t_i for state i can be calculated by Equation 6.

$$t_i = E(H_i) = \frac{1}{\sum_{j \neq i} q_{i,j}} \quad (6)$$

So far, in order to calculate $MTTDL$, we only need to get the value of σ_i . First, we define the vector σ which have w elements, where the first $w-1$ elements are $\sigma_0, \sigma_1, \dots, \sigma_{w-2}$, while the last element is set to 0. Then we get the $\pi^{e0}(x)$, π^{e0} by setting the last element of $\pi^e(x)$, π^e to 0. By setting all the elements in column $w-1$ of P to zero, we get a new matrix P^0 . As we have properties $\pi^e(x+1) = \pi^e(x)P$ and $\pi^e = \pi^e P$ in DTMC, π^{e0} , $\pi^{e0}(x)$, P^0 also have the similar properties, we list these properties in equation 7. That's because state DL is the absorbing state.

$$\begin{aligned} \pi^{e0}(x+1) &= \pi^{e0}(x)P^0, \\ \pi^{e0} &= \pi^{e0}P^0 \end{aligned} \quad (7)$$

From the definition of σ_i , we have $\sigma = \sum_{x=0}^{\infty} \pi^{e0}(x)$. Thanks to absorbing state DL , all the elements in π^{e0} are 0, therefore, we have $\sigma = \sigma + \pi^{e0}$. Using this and the

properties in equation 7, we have equation 8 which gives us a solution to get the value of σ .

$$\begin{aligned} \sum_{x=0}^{\infty} \pi^{e0}(x) &= \sum_{x=0}^{\infty} \pi^{e0}(x) + \pi^{e0} \\ &= \pi^{e0}(0) + \left(\sum_{x=1}^{\infty} \pi^{e0}(x) \right) + \pi^{e0} \\ &= \pi^{e0}(0) + \left(\sum_{x=1}^{\infty} \pi^{e0}(x-1) \right) P^0 \\ &= \left(\sum_{x=0}^{\infty} \pi^{e0}(x) \right) P^0 + \pi^{e0}(0) \\ &= \left(\sum_{x=0}^{\infty} \pi^{e0}(x) \right) P^0 + \pi^e(0), \\ \sigma &= \sigma P^0 + \pi^e(0) \end{aligned} \quad (8)$$

We rewrite equation 8 into element-wise simultaneous equation 9 which has $w-1$ equations.

$$\sigma_i = \sum_{j \neq i} \sigma_j p_{j,i} + \pi_i^e(0), i, j \in [0, w-1] \quad (9)$$

The initial probability vector of each state in RAID-x's embedded Markov chain $\pi^e(0)$ can be formulated as $[1, 0, 0, \dots]$, therefore, simultaneous equation 9 has unique solution.

Let's continue the example of RAID-5. First, we can get the average holding time t_i of each state except DL in RAID-5's Markov chain by equation 6. Therefore, the average holding time of state 0 is $t_0 = \frac{1}{n\mu}$ and the average holding time of state 1 is $t_1 = \frac{1}{(n-1)\mu + \nu}$. As the DL state is the dead state we assume its holding time is 0.

The σ_0 and σ_1 of RAID-5 can be calculated by simultaneous equation 10, which is derived from Equation 9.

$$\begin{cases} \sigma_0 &= r_1 \sigma_1 + 1, \\ \sigma_1 &= \sigma_0 \end{cases} \quad (10)$$

The following equations give the rest inference of getting the MTTDL of RAID-5.

$$\begin{aligned} \sigma_0 &= 1 + r_1 \sigma_1 = 1 + \frac{\nu}{(n-1)\mu + \nu} \sigma \\ \sigma_1 &= \sigma_0 \\ \sigma_0 &= \sigma_1 = \frac{\nu + (n-1)\mu}{(n-1)\mu} \\ MTTDL_{RAID-5} &= \sigma_0 t_0 + \sigma_1 t_1 \\ &= \frac{(n-1)\mu + \nu}{(n-1)\mu} \left(\frac{1}{n\mu} + \frac{1}{(n-1)\mu + \nu} \right) \\ &= \frac{(2n-1)\mu + \nu}{n(n-1)\mu^2} \\ &\approx \frac{\nu}{n(n-1)\mu^2} \end{aligned}$$

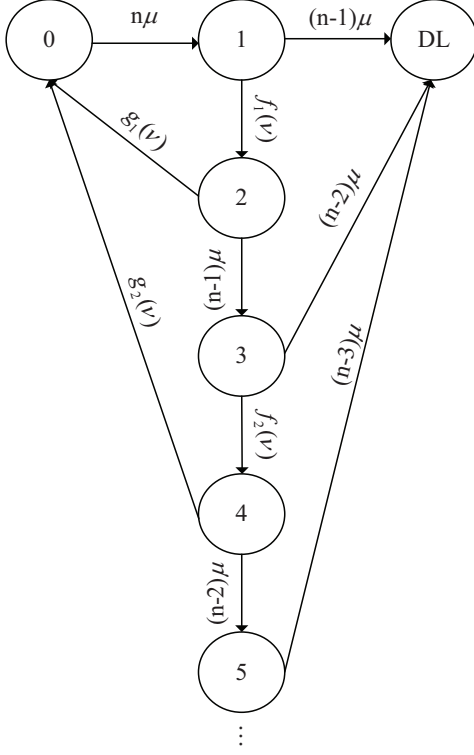


Figure 2. State transition graph of RAID+

At the end of this subsection, we conclude above method as a recipe-like style. By following these steps, the MTDL of RAID-x can be calculated.

1. Get the state transition graph and CTMC of RAID-x.
2. Get the generator matrix Q .
3. Get the average holding time of each state t except state DL .
4. Get the transition probability matrix of embedded Markov chain P .
5. Get the σ by simultaneous equation $\sigma_i = \sum_{j \neq i} \sigma_j p_{j,i} + \pi_i^e(0)$, $i, j \in [0, w-1)$, where $\pi_i^e(0)$ is the initial probability of state i , and w is the number of states in Markov chain.
6. Finally, the MTDL is calculated by $MTDL = \sum_{i=0}^{w-2} t_i \sigma_i$.

2.2 RAID-5 organized RAID+'s Reliability and its comparison to RAID-50

From the lesson learned in subsection 2.1, we calculate the MTDL of RAID-5 organized RAID+. Figure 2 gives us a part of state transition graph of RAID-5 organized RAID+. Suppose we have n disks and organized them into RAID+, in figure 2 there are 7 states where the state 0 stands for normal state that there are total n disks servicing for RAID+. State 1 is the state in which one disk fails and there are $(n-1)$ disks still online. As analyzed before, the $q_{0,1} = n\mu$. State 2 is a

State	t_i	r_i
0	$\frac{1}{n\mu}$	$\frac{1}{n\mu}$
1	$\frac{1}{f_1(\nu) + (n-1)\mu}$	$\frac{f_1(\nu)}{f_1(\nu) + (n-1)\mu}$
2	$\frac{1}{g_1(\nu) + (n-1)\mu}$	$\frac{g_1(\nu)}{g_1(\nu) + (n-1)\mu}$
3	$\frac{1}{f_2(\nu) + (n-2)\mu}$	$\frac{f_2(\nu)}{f_2(\nu) + (n-2)\mu}$
4	$\frac{1}{g_2(\nu) + (n-2)\mu}$	$\frac{g_2(\nu)}{g_2(\nu) + (n-2)\mu}$
5	$\frac{1}{f_3(\nu) + (n-3)\mu}$	$\frac{f_3(\nu)}{f_3(\nu) + (n-3)\mu}$
\vdots	\vdots	\vdots
$2i-1$	$\frac{1}{f_i(\nu) + (n-i)\mu}$	$\frac{f_i(\nu)}{f_i(\nu) + (n-i)\mu}$
$2i$	$\frac{1}{g_i(\nu) + (n-i)\mu}$	$\frac{g_i(\nu)}{g_i(\nu) + (n-i)\mu}$
\vdots	\vdots	\vdots

Table 1. The average holding time t and the recover transition probability r of the states in RAID+ (except state DL)

recovered state where parities are moved into the rest $n-1$ disks, therefore, there are $n-1$ disks online while RAID+ is still tolerant for one disk failure. The average repair time in RAID+ a function of conventional RAID-5's, therefore we use $f_1(\nu)$ to represent the disk repair rate with $n-1$ disks working online. Then disk adjunction is needed to maintain the disk numbers in RAID+ and we denote the adjunction rate for i disks as $g_i(\nu)$. In figure 2, we can see $q_{2,0} = g_1(\nu)$ as there is only one disk needed to maintain n -disk-RAID+. If another disk fails, State 2 will turn into state 3, then if recovered then the system would be in state 4 otherwise if another disk fails the data loss happens. The graph is only drawn with 5 abnormal state but in real application the abnormal state can be extended to tolerate more disk failures if the parities are stored in online disks.

After forming the model of RAID+, we need to calculate the t , P , and σ to get the MTDL. The average holding time t of the states in RAID+ (except state DL) and the probability transition matrix of embedded Markov chain P are given in table 1 and equation 3. Note that the r_i s and the s_i s are the recover transition probability and failure transition probability which indicates the transition probability of a safer state and a more dangerous state for each state. The values of r_i s are given in table 1 and the values of s_i s can be inferred with equation $r_i + s_i = 1$.

$$P_{RAID-5+} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & r_1 & 0 & 0 & 0 & 0 & \cdots & s_1 \\ r_2 & 0 & 0 & s_2 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 0 & r_3 & 0 & 0 & \cdots & s_3 \\ r_4 & 0 & 0 & 0 & 0 & s_4 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots & s_5 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \quad (11)$$

Then we get the value of σ according to the simultaneous equation $\sigma_i = \sum_{j \neq i} \sigma_j p_{i,j} + \pi_i^e(0)$. As we put the calculated figures into simultaneous equation, we get simultaneous equation 12.

$$\begin{aligned}\sigma_0 &= 1 + \sigma_2 r_2 + \sigma_4 r_4 + \dots \\ \sigma_1 &= \sigma_0 \\ \sigma_2 &= \sigma_1 r_1 \\ \sigma_3 &= \sigma_2 s_2 \\ \sigma_4 &= \sigma_3 r_3 \\ \sigma_5 &= \sigma_4 s_4 \\ &\vdots\end{aligned}\tag{12}$$

From the simultaneous equation 12, we can infer that $\sigma_0 = \frac{1}{1 - \sum_{i=1}^n \prod_{j=1}^{2i-1} (I_{\text{even}}(j)r_i + (1 - I_{\text{even}}(j))s_i)}$, in which $I_{\text{even}}(x)$ stands for the indicator function for even number, if the x is even, $I_{\text{even}}(x) = 1$ and $I_{\text{even}}(x) = 0$ if not. The rest elements in σ can be inferred by σ_0 .

Since the result of σ is nearly unwritable in the paper, a simplified model of RAID+ with some approximation loss of MTTDL will be adopted. This simplified model is based on the simplification of the equation $MTTDL = \sum_{i=0}^{w-1} t_i \sigma_i$ and equation $\sigma_0 = \frac{1}{1 - \sum_{i=1}^n \prod_{j=1}^i (I_{\text{even}}(j)r_i + (1 - I_{\text{even}}(j))s_i)}$. In the first equation, we just remain the first three entries while in the second equation, we remain the first entry of $\sum_{i=1}^n \prod_{j=1}^i (I_{\text{even}}(j)r_i + (1 - I_{\text{even}}(j))s_i)$ in denominator. Note that all the two simplifications will lead to a smaller σ and MTTDL but the model can't express the property that RAID+ could tolerate more failures after repairing.

In this paragraph, we will finish all the preparation for comparing the MTTDL of RAID-50 and RAID-5 organized RAID+, including the configurations, the analysis of the relationship between $f_i(\nu)$, $g_i(\nu)$, ν . Assuming that there are n disks, RAID-50 organized k disks into RAID-5 and puts n/k RAID-5 into one RAID-0. In RAID+, the stripe size is set to k . Suppose we have S bytes contents in each disk, then if one disk fails, RAID-5 will repair with the average time of S/w_{disk} in which w_{disk} stands for the write bandwidth of one disk. In subsection 2.1, we assume that the disk repairs follow the exponential distribution with rate ν , therefore, the average repair time can be formulated by $1/\nu$. From above two analyses, we can get the equation $\nu = w_{\text{disk}}/S$. For RAID+, as the parity per disk is $S \frac{k}{n-i}$ when there are $n-i$ disks online, we can get $f_i(\nu) = \frac{(n-i)}{nk} \nu$. When adding the disks on RAID+, each disk is written S bytes with the speed w_{disk} thus $g_i(\nu) = \nu$.

Based on the preparation in the last paragraph, we can get the simplified MTTDL of RAID+ is equation 13. As released in subsection 2.1, the MTTDL of RAID-5 is $\frac{\nu}{k(k-1)\mu^2}$. When n/k RAID-5s organized into RAID0, the MTTDL of RAID-50 becomes $\frac{\nu}{n(k-1)\mu^2}$. Therefore, the ratio of MTTDL for two RAIDs $\frac{MTTDL_{RAID-5+}}{MTTDL_{RAID-50}} = \frac{k-1}{k}$. Considering the approxi-

mation when counting the estimated MTTDL of RAID+, the ratio is slightly larger than $\frac{k-1}{k}$.

$$\begin{aligned}MTTDL_{RAID-5+} &> \frac{1}{1-r_1} \frac{1}{n\mu} + \frac{1}{1-r_1} \frac{1}{(n-1)\mu + f_1(\nu)} \\ &= \frac{\nu}{nk\mu^2}\end{aligned}\tag{13}$$

2.3 RAID6 organized RAID+'s Reliability and its comparison to RAID-60

Following the conclusion in subsection 2.1, we calculated the MTTDL of RAID6 step by step:

1. Get the state transition graph and CTMC of RAID6, the state transition graph is shown in figure 3.
2. Get the generator matrix Q , the generator matrix Q_{RAID6} is given in equation 14.
3. Get the average holding time of each state t , the result is given in equation 15.
4. Get the transition probability matrix of embedded Markov chain P , the result is given in equation 16.
5. Get the σ by simultaneous equation $\sigma_i = \sum_{j \neq i} \sigma_j p_{i,j} + \pi_i^e(0)$, where $\pi_i^e(0)$ is the initial probability of state i , the result is given in equation 17. The $p_{i,j}$ refers to the i, j -th element of P_{RAID6} .
6. The MTTDL can be calculated by $MTTDL = \sum_{i=0}^{w-1} t_i \sigma_i$, where w is the number of states in Markov chain. The inference is given in equation 18.

$$\begin{aligned}Q_{RAID6} &= \begin{bmatrix} -n\mu & n\mu & 0 & 0 \\ \nu & -\nu - (n-1)\mu & (n-1)\mu & 0 \\ 0 & \nu & -\nu - (n-2)\mu & (n-2)\mu \\ 0 & 0 & 0 & 0 \end{bmatrix}\end{aligned}\tag{14}$$

$$\begin{aligned}t_{RAID6} &= \begin{bmatrix} \frac{1}{n\mu} & \frac{1}{\nu + (n-1)\mu} & \frac{1}{\nu + (n-2)\mu} \end{bmatrix}\end{aligned}\tag{15}$$

$$\begin{aligned}P_{RAID6} &= \begin{bmatrix} 0 & 1 & 0 & 0 \\ \frac{\nu}{\nu + (n-1)\mu} & 0 & \frac{(n-1)\mu}{\nu + (n-1)\mu} & 0 \\ 0 & \frac{\nu}{\nu + (n-2)\mu} & 0 & \frac{(n-2)\mu}{\nu + (n-2)\mu} \\ 0 & 0 & 0 & 1 \end{bmatrix}\end{aligned}\tag{16}$$

$$\begin{aligned}\sigma_{RAID6} &= \begin{bmatrix} 1 + \frac{p_{1,0}}{(1-p_{1,0})(1-p_{2,1})} & \frac{1}{(1-p_{1,0})(1-p_{2,1})} & \frac{1-p_{1,0}}{(1-p_{1,0})(1-p_{2,1})} \end{bmatrix}\end{aligned}\tag{17}$$

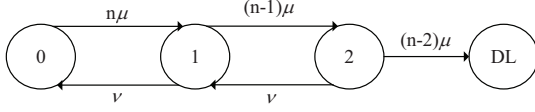


Figure 3. State transition graph of RAID6

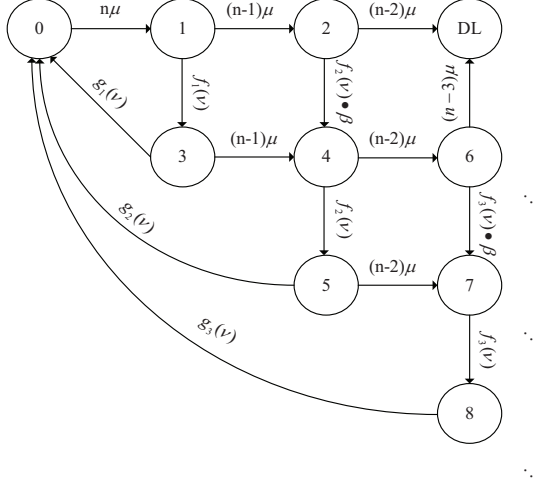


Figure 4. State transition graph of RAID6 based RAID+

$MTTDL_{RAID6}$

$$\begin{aligned}
&= \sum_{i=0}^2 \sigma_i t_i \\
&= \left(1 + \frac{(\mu(n-1) + \nu)(\mu(n-2) + \nu)}{\mu^2(n-1)(n-2)}\right) \frac{1}{n\mu} \\
&+ \left(\frac{(\mu(n-1) + \nu)(\mu(n-2) + \nu)}{\mu^2(n-1)(n-2)}\right) \frac{1}{(n-1)\mu + \nu} \\
&+ \frac{\mu(n-2) + \nu}{\mu(n-2)} \frac{1}{(n-2)\mu + \nu} \\
&= \frac{2(n-1)(n-2)\mu^2}{\mu^3 n(n-1)(n-2)} \\
&+ \frac{\mu^2 n(2n-3) + \mu n \nu + \nu^2 + \mu(2n-3)\nu}{\mu^3 n(n-1)(n-2)} \\
&\approx \frac{\nu^2}{\mu^3 n(n-1)(n-2)}
\end{aligned} \tag{18}$$

The MTTDL of RAID6 organized RAID+ is even more complicated than RAID-5 organized RAID+. The complexities which RAID+ brings to us are follows. First, the RAID+ will recover its fault-tolerant ability before the disk added on, which means the states of the same-level-tolerant are not the same state. This property of model is similar to RAID-5 organized RAID+. Second, in RAID6 organized RAID+, when two disks fail, the first level recovery will have less data to transfer, we denote that the ratio of original data transfer and optimized data transfer as β . with the property

i	j	$q_{i,j}$	$p_{i,j}$
0	1	$n\mu$	1
0	0	$-n\mu$	0
1	2	$(n-1)\mu$	s_1
1	3	$f_1(\nu)$	r_1
1	1	$-(n-1)\mu - f_1(\nu)$	0
2	9	$(n-2)\mu$	s_2
2	4	$f_2(\nu)\beta$	r_2
2	2	$-(n-2)\mu - f_2(\nu)\beta$	0
3	0	$g_1(\nu)$	r_3
3	4	$(n-1)\mu$	s_3
3	3	$-(n-1)\mu - g_1(\nu)$	0
4	5	$f_2(\nu)$	r_4
4	6	$(n-2)\mu$	s_4
4	4	$-(n-2)\mu - f_2(\nu)$	0
5	0	$g_2(\nu)$	r_5
5	7	$(n-2)\mu$	s_5
5	5	$-(n-2)\mu - g_2(\nu)$	0
6	7	$f_3(\nu)\beta$	r_6
6	9	$(n-3)\mu$	s_6
6	6	$-(n-3)\mu - f_3(\nu)\beta$	0
7	8	$f_3(\nu)$	1
7	7	$-f_3(\nu)$	0
8	0	$g_1(\nu)$	1
8	8	$-g_1(\nu)$	0
9	9	0	1

Table 2. Q_{RAID6+} and P_{RAID6+}

6 in the main paper, we can get that $\beta = \frac{2(n-k)+k-2}{k-2}$. Figure 4 illustrates the state transition graph of RAID6 based RAID+, in this figure, we listed part of the states in RAID+ and here we'll explain them respectively. The state 0 is the normal state in which there are n disks in RAID+. In state 1, one disk fails and there are $n-1$ disks online. In state 2, two disks fails and there are $n-2$ disks online. In state 3, the disk failure in state 1 is recovered and there are $n-1$ disks online. In state 4, only one disk failure is recovered and there are $n-2$ disks online. In state 5, all the disk fails are recovered and there are $n-2$ disks online. In state 6, there are two disk failures and only $n-3$ disks online. In state 7, one disk failure is recovered in state 6 and only $n-3$ disks online. In state 8, disk failures are all recovered in state 6, 7 but only $n-3$ disks online. In state DL , more than two disks fail.

For the same simplification solution with RAID-5 organized RAID+ (with the MTTDL loss for RAID+), we only use those 10 states to estimate the MTTDL of RAID+. Follow the conventional way, we first get the Q_{RAID6+} , for better fitting the size of paper, we use sparse matrix representation COO to present it. The Q_{RAID6+} is shown in table 2.

From the generator matrix Q_{RAID6+} , we can get the t_{RAID6+} and P_{RAID6+} . The results of them are listed in

i	t_i	r_i
1	$\frac{1}{n\mu}$	0
2	$\frac{1}{(n-1)\mu + \frac{n-1}{k}\nu}$	$\frac{\frac{n-1}{k}\nu}{(n-1)\mu + \frac{n-1}{k}\nu}$
3	$\frac{1}{(n-2)\mu + \frac{n-2}{k}\frac{2(n-k)+k-2}{k-2}\nu}$	$\frac{\frac{n-2}{k}\frac{2(n-k)+k-2}{k-2}\nu}{(n-2)\mu + \frac{n-2}{k}\frac{2(n-k)+k-2}{k-2}\nu}$
4	$\frac{1}{(n-1)\mu + \nu}$	$\frac{\nu}{(n-1)\mu + \nu}$
5	$\frac{1}{(n-2)\mu + \frac{n-3}{k}\nu}$	$\frac{\frac{n-3}{k}\nu}{(n-2)\mu + \frac{n-3}{k}\nu}$
6	$\frac{1}{(n-2)\mu + \nu}$	$\frac{\nu}{(n-2)\mu + \nu}$
7	$\frac{1}{(n-3)\mu + \frac{n-3}{k}\frac{2(n-k)+k-2}{k-2}\nu}$	$\frac{\frac{n-3}{k}\frac{2(n-k)+k-2}{k-2}\nu}{(n-3)\mu + \frac{n-3}{k}\frac{2(n-k)+k-2}{k-2}\nu}$
8	$\frac{1}{\frac{n-3}{k}\nu}$	1
9	$\frac{1}{\nu}$	1

Table 3. t_{RAID6+} and r_{RAID6+}

table 2 and table 3. Note that the r_i in P_{RAID6+} is given in 3 and the $s_i = 1 - r_i$.

After P_{RAID6+} is calculated, the simultaneous equation 20 can be listed. The simultaneous equation yields equation 19.

$$\begin{aligned} \sigma_1 &= 1 \\ &/((1 - r_4r_2 - r_5r_6s_4r_2 \\ &- r_6r_5r_3s_2 - r_7s_5r_3r_2 \\ &- r_7s_5s_4r_2 - s_6r_5s_4r_2 - s_6r_5r_3s_2)) \end{aligned} \quad (19)$$

$$\begin{aligned} \sigma_1 &= 1 + \sigma_9 + r_6\sigma_6 + r_4\sigma_4 \\ \sigma_2 &= \sigma_1 \\ \sigma_3 &= s_2\sigma_2 \\ \sigma_4 &= r_2\sigma_2 \\ \sigma_5 &= r_3\sigma_3 + s_4\sigma_4 \\ \sigma_6 &= r_5\sigma_5 \\ \sigma_7 &= s_5\sigma_5 \\ \sigma_8 &= r_7\sigma_7 + s_6\sigma_6 \\ \sigma_9 &= \sigma_8 \end{aligned} \quad (20)$$

Different from RAID-5 organized RAID+, with simplification model the result of MTTDL of RAID6 is still too complicated to directly write here. Thus further approximation is adopted. We throw away all the states that have k failure tolerance and disks less than $n + k - 2$ and just approximate the model as the disk adjunction is immediate. This approximation is not promise to have a bigger or smaller MTTDL thus it's only a reference method. We take the result of $f_i(x)$ and $g_i(x)$ in subsection 2.2 and set $n = 56$ and $k = 7$. The RAID-60's MTTDL is $\frac{\nu^2}{\mu^3n(k-1)(k-2)}$ and the estimated MTTDL of RAID-6 organized RAID+ is in 21.

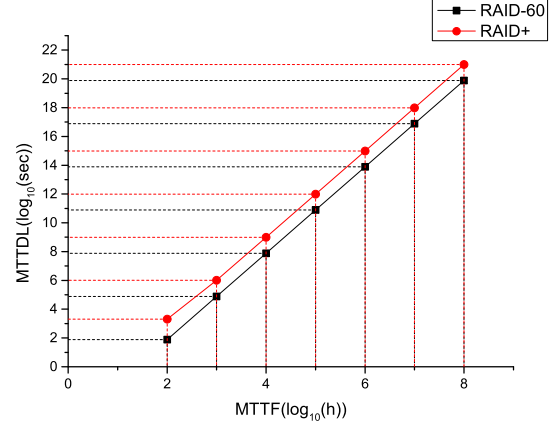


Figure 5. Comparison for RAID-60 and RAID+'s MTTDL

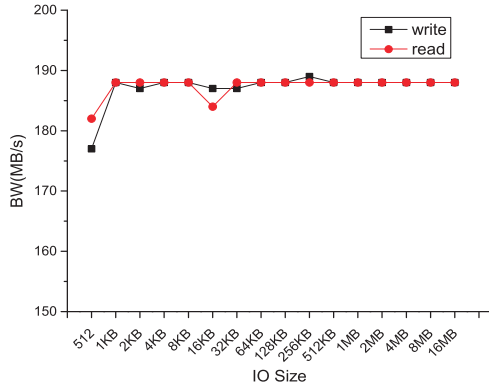
Thus the ratio of two MTTDL is $\frac{MTTDL_{estimated,RAID-6+}}{MTTDL_{RAID-60}} = \frac{(2(n-k)+k-2)(k-1)}{k^2} \approx 12.61$.

$$\begin{aligned} &MTTDL_{estimated,RAID-6+} \\ &= \frac{\nu_1\nu^2}{\mu^3n(n-1)(n-2)} \\ &= \frac{\frac{2(n-k)+k-2}{k-2} \frac{n-2}{k} \nu \frac{n-1}{k} \nu}{\mu^3n(n-1)(n-2)} \\ &= \frac{(2(n-k)+k-2)\nu^2}{\mu^3nk^2(k-2)} \end{aligned} \quad (21)$$

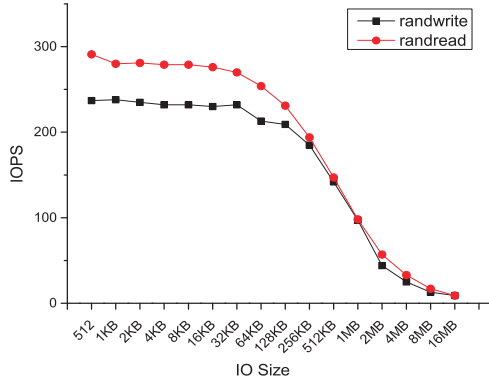
Last paragraph only gives an estimated ratio between two RAIDs. Here we use MATLAB to calculate the MTTDL for both RAIDs with different MTTF. The rest parameters are fixed as $n = 56$, $k = 7$, $\nu = \frac{1}{10000}$. We set MTTF to $3600 \times [100, 1000, 10000, 100000, 1000000, 10000000, 100000000]$ seconds and plot the MTTDL of both RAIDs. The result is shown in 5, in this figure, the smallest ratio of two RAIDs is $(\frac{MTTDL_{RAID-6+}}{MTTDL_{RAID-60}})_{min} = 12.6$. From this conclusion, we can infer that the approximation in the last paragraph is rather accurate and the reason comes that the probabilities of transmitting into most of the states are low.

3. Disk Performance

Figure 6 gives the synthetic workload performance on single disks using fio as these disks are used in previous experiments. In sequential access evaluation, the read and write performance are rather stable with different IO sizes. Besides IO size 512 Bytes, both the sequential read and write bandwidth of one disk is approximately stabilized at 188MB/s. In random access evaluation, the IOPS of random read and write are stable at around 280 and 230 respectively when IO size is less than 64KB. Then in the IO size ranged from 64KB to 2MB, a sharp reduction of IOPS occurs, however, this reduction effect of IOPS is weakened when the IO size is larger than 2MB.



(a) Sequential performance



(b) Random performance

Figure 6. Disk performance under synthetic workloads

References

- [1] GREENAN, K. M. *Reliability and power-efficiency in erasure-coded storage systems*. University of California, Santa Cruz, 2009.
- [2] GREENAN, K. M., PLANK, J. S., WYLIE, J. J., ET AL. Mean time to meaningless: Mttld, markov models, and storage system reliability. In *HotStorage* (2010).
- [3] XIN, Q., MILLER, E. L., SCHWARZ, T., LONG, D. D., BRANDT, S. A., AND LITWIN, W. Reliability mechanisms for very large storage systems. In *Mass Storage Systems and Technologies, 2003.(MSST 2003). Proceedings. 20th IEEE/11th NASA Goddard Conference on* (2003), IEEE, pp. 146–156.