

# 网络存储系统容错编码技术进展

Erasure Codes for Networked Storage Systems

林 胜 /LIN Sheng

(天津工业大学 信息与工程 学院, 天津 300191)

刘晓光/LIU Xiaoguang

王 刚/WANG Gang

(南开大学 信息技术科学学院, 天津 300071)

(College of Information Technical Science, Nankai University, Tianjin 300071, China)

中图分类号: TP393.03 文 志 : A 文 号: 1009-6868 (2010) 05-0000-00

基 金 : 国家 技 术 发展(“863”) 计划(2008AA01Z401); 国家 学 基 金 (60903028)

摘 要 :

前,专业 大型 存储 均发展为包含多块 大型 列 。 中  
数 不断增加, 失效引 数据丢失 可 性 大。对于 存储 中 分  
失效所引 数据丢失 , 前业 公 好 决方 是使 冗余容 技  
实 容 。在工 实 中, 前广 应 方 大多局 于双容 列 。  
一 加大, 3容 更多容 方 已引 。今后 5  
10年 ,对于3容 或多容 方 将会成为新 。

关 键 词 :

存储 ; 容 ; 列

## Abstract:

Large storage systems are generally large array systems consisted of a great number of hard disks. With the increasing of the number of the hard disks, the chance of data loss also increased. To address this problem, researchers have conducted a lot of research. For the problem of data loss caused by the disk fault in the storage system, the recognized solution is to use coding technique to achieve fault tolerance. Nowadays, only double-erasure array codes are widely used in engineering practice. With the increasing of the system size, triple or more erasure coding schemes will draw more attentions. Experts in this field generally agreed that triple-erasure coding scheme will become the new hot spot in next 5 to 10 years.

## Key words:

storage systems; fault-tolerant coding; array code

## 1 存储容 与 价指

20 年 , 技 发展, 大 存储 发展也十分 。当前, 普 PC 存储器 容 已 到了太 别, 之 20 年前 20 MB 提 了 10 000 倍。了传 动器之外, 新型 固态存储 (SSD) 存储器也已 向市场。尽 单个存储 器 容 发展 , 但是却仍 不上人们对存储容 增 度。 大型 “以 为中心” 向 “以信息处 为中心” 变, 以及信息 式增 , 人们 对 存储 日 提 。 存储 上是将很多 单个存储器件 (下 均以 为例), 接口, 接整合为一个 拟 容 巨大 单一存储器, 即 列。

列中 数 增多, 可 性也 之下 。工业 一 使 平均数据丢失时 (MTTDL) 列 可 性。 单个 平均失效时 为  $MTTF_{disk}$ , 则对于包含  $n$  块 无冗余 列 , 其  $MTTDL$  可 单估 为:  $MTTDL = MTTF_{disk}/n$ 。可 , 当  $n$  大时, 整个 可 性成 例下 。对于 大 是不可接受 。利 冗余数据 提 可 性是公 决 一 好方 。 巧妙 将  $m$  块 准大小 上 数据, 增加 分冗余 信息, 后存放于  $n$  块 上, 使得 : 对于任意  $k$  块 失效, 可以 其他  $n-k$  块 失效 中 数据 恢复, 则 整个 是  $k$  容 , 或  $k$  为 容 数。 分 明<sup>[1]</sup>, 对于  $k$  容 ,  $MTTDL$  可以 似估 为:

$$MTTDL = MTTF_{disk}^{k+1} / (n(n-1) \dots (n-k) MTR^k) \quad (1)$$

因 , 在大 中, 容 数可以 是另一 对 可 性 描 方式。市场 中一  $MTTF_{disk}$  为  $10^5$  左右, 修复时  $MTR$  一 为  $10^1$  左右。 据 (1) 式可以 出, 当 数为  $10^3 \sim 10^4$  时, 一 2 容 或是 3 容 就基 上可以 存储 容

。于增加容 力 加 冗余 多, 外 价也将 。因 在具有 同容 数 前提下, 人们往往 更小 冗余度, 即  $(n-m)/n$  值, 其中  $n$  为 数、 $m$  为存储 户数据 数。 据 Singleton ,  $k$  容 最小冗余度为:  $k/n$ 。到 一最小值 方 做最大最小 可分(MDS) 。前多数存储 中于 不同参数下 MDS 。

了上 指 , 任何 度与效 是 指 。我们 不 如何有效 并 处 多 中 数据 取 ( 是另外一个 大 ), 于冗余 带 外 开 。对于即便是 同 方 , 于 / 不同, 可 效 差异 大。 于在 中, 最 会反映为一些二 制 , 因 常使 总 二 制异或 数 于 外冗余 带

开。对于一个存取存储，小块信息写操作性尤为。中个单元所参与平均异或数可以一指，我们其为更新复度。

合上，存储容可以归为寻对如下指优化方：

- 容性，容数为  $k$ 。
- 有小(或最优)冗余度。
- 有小(或最优)/更新复度。

## 2 性

对于单容，单奇偶使得上3到最优。典RAID1,4,5是使方。对于  $k>1$  情况，决就不是么显易了。从信丰富成中，两有代性方人们挑出，并于决存储容，他们是2制性和RS。

### 2.1 多列

图1所是二列及。二列是奇偶推广，图1很容易出它是双容。

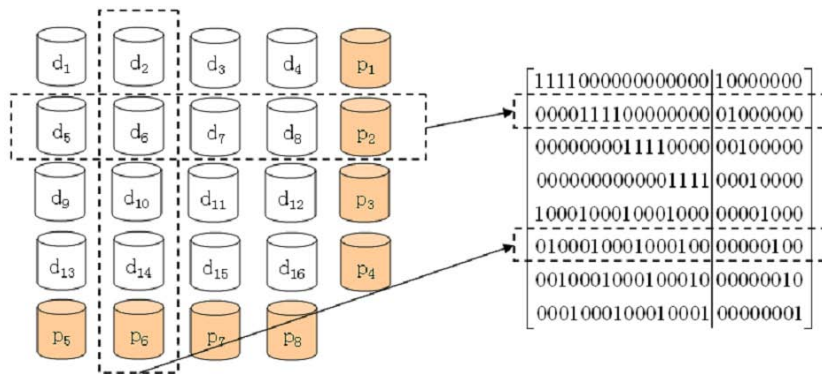


图1二列及二列保持了单容时奇偶最优复度性，但是他冗余度不再是最优了。

二列也很容易推广为  $k$  列，并且容易明  $k$  容性，但是  $k$  增大，冗余大<sup>[2,3]</sup>。

### 2.2 Full

图2所是FULL-2。FULL-2可做是二列推广。

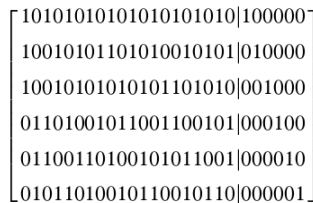


图2 FULL-2

FULL 依保持了最优复度，并且冗余度列好很多。不幸是，当  $k>3$  时，FULL- $k$  不再是  $k$  容<sup>[4]</sup>。

### 2.3 RS

图3所是RS。RS从最佳冗余性出发。到 Singleton RS 人们提出并广应。

$$\begin{pmatrix} 1 & \alpha & \alpha^2 & \dots & \alpha^{q-2} \\ 1 & \alpha^2 & (\alpha^2)^2 & \dots & (\alpha^2)^{q-2} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & (\alpha^{d-1})^2 & (\alpha^{d-1})^2 & \dots & (\alpha^{d-1})^{q-2} \end{pmatrix}$$

图 3 RS

上性变换可以化为，存储亦即显式区分出中哪些单元于存储单元。可以出，中元不再是01，为有域元幂，故使有域。在中，有域元最后是会映射为01单元，时个有域元一会映射为多个01单元，有域也可以分为些01单元复。我们仍以所异或为基单位，则所异或数和巧妙度有关。前好域所异或数大为 $O(n^3)^{[5]}$ ，复度当。RS是MDS，故冗余度是最优。

### 3 列

上几各有优，么是否存在对于多指同时最优 $k$ 容方呢？文[5]提出EVENODD，一大只使异或列提出并广。多列或FULL二进制性块只取一个单元。列则在块上取多个单元，一交叉。同二进制性一，只使二进制异或，但冗余度却可以与RS同。

#### 3.1 EVENODD

EVENODD想很单，块中取干单元，排成方，后将些单元分成不同，另外加两块于存储单元。所有均使单二进制奇偶。

##### 1 平与对

Disk1	Disk2	Disk3	Disk4	Disk5	Disk6	Disk7
D15	D16	D17	D18	D19	P1	P5
D26	D27	D28	D29	D25	P2	P6
D37	D38	D39	D35	D36	P3	P7
D48	D49	D45	D46	D47	P4	P8
0	0	0	0	0	0	P9

平与对如1所。1中D代户数据单元，P代冗余单元。可以出，Disk 1~5存储户数据单元；Disk6、7存储冗余单元。Disk6各单元为户数据各平和，Disk7各单元为户数据对和。

存储户数据数为 $p$ （如上例中 $p=5$ ），则包含 $p+2$ 块，前 $p+1$ 块中最后一个单元为拟0元，故实包含 $p-1$ 个单元，最后一块包含 $p$ 个单元。可以明，当 $p$ 为数时是双容。

单可时冗余度为 $(2p-1)/((p+2)(p-1)+1)$ 。于最后多出一个单元，所以冗余度大于最优 $2/(p+2)$ 。为了到最优值，文[5]中使如下技巧：将多出单元（即对交和）叠加到其他单元上，MDS EVENODD如2所。

##### 2 MDS EVENODD

Disk1	Disk2	Disk3	Disk4	Disk5	Disk6	Disk7
-------	-------	-------	-------	-------	-------	-------

D15	D16	D17	D18	D19	P1	P5+ P9
D26	D27	D28	D29	D25	P2	P6+ P9
D37	D38	D39	D35	D36	P3	P7+ P9
D48	D49	D45	D46	D47	P4	P8+ P9

2 也可 为如 3 所 。

### 3 MDS EVENODD

Disk1	Disk2	Disk3	Disk4	Disk5	Disk6	Disk7
D15	D16	D17	D18	D15678	P1	P5
D26	D27	D28	D25678	D25	P2	P6
D37	D38	D35678	D35	D36	P3	P7
D48	D45678	D45	D46	D47	P4	P8

也就是 当 一 对 和为 1 时, 其他各对 为奇 ; 当 一 对 和为 0 时, 其他各对 为偶 。 就是它 命名为 EVENODD 原因。

### 3.2 RDP

从 2 可以 出, 为了得到冗余最优, EVENODD 对 上 单元 更新复 度很 。 更新 些单元 数据时 同时更新其他  $p$  个 单元, 对于双容 , 最优值为 2。文 [6]中 RDP 将 些单元 更新复 度均 到 个单元, 从 有效地 了小写操作中更新性 不均 。包含 平 对 如 4 所 。

4 包含 平 对

Disk1	Disk2	Disk3	Disk4	Disk5	Disk6
D159	D169	D179	D189	P1	P5
D256	D256	D258	D259	P2	P6
D367	D368	D369	D356	P3	P7
D478	D479	D457	D467	P4	P8
0	0	0	0	0	P9

与 EVENODD 不同处在于, 做对 时也包含了 平 单元 一列 (因 , 数据单元 也 EVENODD 少了一列)。

同 , RDP 最后一个 多出一个单元, 使得整个 不是 MDS 。但 RDP 优势在于, 单地将多出 单元删去, 仍 为双容 。即得到如 5 所 列。

### 5 RDP MDS

Disk1	Disk2	Disk3	Disk4	Disk5	Disk6
D15	D16	D17	D18	P1	P5
D256	D256	D258	D25	P2	P6
D367	D368	D36	D356	P3	P7
D478	D47	D457	D467	P4	P8

从 5 可以 出, 所有数据单元 更新 为 2 或 3, 分布 EVENODD 均匀, 不会 产 方式带 外 , 但 平均更新复 度是 同 。

### 3.3 Liberation

从前 几 可以 出, 所使 方 是 平 加其他一 共同 成双容 。不同之处就在于“另一 ” 不同 择。如将另一 上 元 作一个 01 向 , 块数据 上 单元对 些 元 影响可 一个 01 。如 5 中 1 列 4 个数据单元对 Disk7 中 各 元 影响可 为如图 4 所 。

$$\begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

图 4 列

在下, 前所更新复度就对应中 1 个数。于是 一个双容列 就变为: 寻找 若干个, 使得其中 1 个数尽少, 并且任意 2 个之和为。

在  $p$  为 数时, 文 [7] 中 Liberation 使得  $p \times p$  1 数不  $p+1$ , 其  $p$  个可 单地描 为: 各对 加一个 外单元。  $k$  个 外 1 单元 位 可描 为  $(k(p-1)/2 \text{ Mod } p, 1+k(p-1)/2 \text{ Mod } p)$ 。得到 如 6 所。

6 Liberation

Disk1	Disk2	Disk3	Disk4	Disk5	Disk6	Disk7	Disk8	Disk9
D1a	D1g	<b>D1fg</b>	D1e	D1d	D1c	D1b	P1	Pa
D2b	D2a	D2g	D2f	<b>D2ef</b>	D2d	D2c	P2	Pb
D3c	D3b	D3a	D3g	D3f	D3e	<b>D3de</b>	P3	Pc
D4d	<b>D4cd</b>	D4b	D4a	D4g	D4f	D4e	P4	Pd
D5e	D5d	D5c	<b>D5bc</b>	D5a	D5g	D5f	P5	Pe
D6f	D6e	D6d	D6c	D6b	<b>D6ab</b>	D6g	P6	Pf
D7g	D7f	D7e	D7d	D7c	D7b	D7a	P7	Pg

### 3.4 PDHLatin

前 些 为 MDS 充 件均为: 与 数 关 (RDP 为  $p+1$ , 其他为  $p+2$ )。它们 双容 方 均为 据一个已 单元, 后 关 与失效单元形成 式关 依 恢复所有单元。使人们 到其容 力 是任意两列 可以形成 关。

。文 [8] 中利 拉丁方 了 PDHLatin, 使得 不再必 关 一个 数。所 拉丁方是指  $n \times n$  方 中填入  $n$  个不同 号, 使得 列 号 不 复。显 拉丁方 两列 成一个  $n$  元 换, 所 密尔 拉丁方是指拉丁方 任何两列 成 换为 单。图 5 为一个 9 密尔 拉丁方。

$$\begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ 2 & 4 & 8 & 9 & 3 & 5 & 1 & 7 & 6 \\ 3 & 1 & 9 & 2 & 8 & 7 & 5 & 6 & 4 \\ 4 & 5 & 2 & 3 & 1 & 8 & 6 & 9 & 7 \\ 5 & 7 & 4 & 1 & 6 & 9 & 8 & 3 & 2 \\ 6 & 9 & 5 & 8 & 7 & 4 & 2 & 1 & 3 \\ 7 & 8 & 6 & 5 & 9 & 2 & 3 & 4 & 1 \\ 8 & 6 & 1 & 7 & 4 & 3 & 9 & 2 & 5 \\ 9 & 3 & 7 & 6 & 2 & 1 & 4 & 5 & 8 \end{bmatrix}$$

图 5 密尔 拉丁方

从一个 定 密尔 拉丁方, 我们可以 与 EVENODD 似 方, 只不 各单元对于 二 关 不再依单元所在对 位 决定, 是 据拉丁方 应位 号决定。据图 5, 得到 7 所 PDHLatin。

7 PDHLatin

Disk1	Disk2	Disk3	Disk4	Disk5	Disk6	Disk7	Disk8	Disk9	Disk10	Disk11
Da1	Da2	Da3	Da4	Da5	Da6	Da7	Da8	Da9	Pa	P1
Db2	Db4	Db8	Db9	Db3	Db5	Db1	Db7	Db6	Pb	P2

Dc3	Dc1	Dc9	Dc2	Dc8	Dc7	Dc5	Dc6	Dc4	Pc	P3
Dd4	Dd5	Dd2	Dd3	Dd1	Dd8	Dd6	Dd9	Dd7	Pd	P4
De5	De7	De4	De1	De6	De9	De8	De3	De2	Pe	P5
Df6	Df9	Df5	Df8	Df7	Df4	Df2	Df1	Df3	Pf	P6
Dg7	Dg8	Dg6	Dg5	Dg9	Dg2	Dg3	Dg4	Dg1	Pg	P7
Dh8	Dh6	Dh1	Dh7	Dh4	Dh3	Dh9	Dh2	Dh5	Ph	P8
0	0	0	0	0	0	0	0	0	0	P9

### 3.5 X

上介几方到了冗余最优，但在更新复度方均于最优值，么是否可以到两同时最优呢？文 [9]提出 X 是一双容。X 想也很单，仍是在列中主对和对两，但是巧妙地单元分布到各个中（不是像其他方中，单元分出，存放于），使得同时到了两方指同时最优。为了双容，X 也列中包含列数（或）为数。为数  $p$  X 中，一列包含  $p-2$  个户数据单元，2 个冗余单元。

### 3.6 B

是否存在与 X 同性其他方呢？显将两个 X 列叠，仍保持最优冗余与最优更新复度。

得到新，在数不变情况下，块关单元数加倍。在实中，为了化实，我们实上块关单元数尽量少。对于  $n$  块，在保持最优冗余与最优更新复度条件下，块最少多少个单元关呢？文 [10] 提出 B 在双容情况下，决了一。

将同于图中完全图完一因子分。并据图已有，出一各方性均到最优。从一个完全图  $K_6$  一完 1 因子分方，我们可以如 8 所双容。

### 8 B

Disk1	Disk2	Disk3	Disk4	Disk5
D23	D34	D14	D12	D13
P1	P2	P3	P4	D24

，块包含多 1 个单元，并且只有一块不包含单元。它将  $n$  个号所有 2 元分划为各列，并且双容，因在保持了最优冗余度与更新复度前提下，到最。因也做最最低密度列。

### 3.7 T

对于 3 容最最低密度列之双容复很多，文 [11]最先出了一，并利助明了些参数下，3、4 容最最低密度列 MDS 性。我们在文 [12]中同并利合 NRB（乎可分不完全区）出了合，同时也出了明代数明。

T 从形式上与 B 同，块包含多 1 个单元，并且只有一块不包含单元。文 [12]明了对于任意容最最低密度列均性。对于普参数 T，或任意容最最低密度列，仍是困。

### 3.8 Weaver

前将优化冗余最为一，同时兼/更新复度。但在一些中，如冗余当损失可换更好性或更易于，则也是可择。文 [13]从优先/更新复度度，提出了易于 Weaver。

B、T 也可以出，在保持更新复杂度最优前提下，单元分布在各中容易。为了化，文 [13] 择具有循对性列，也就是：(1) 所有数据单元参与数为常数；(2) 所有包含单元数为常数；(3) 如  $i$  上数据单元  $j$  参与  $k$  上单元  $p$  所代，则必有对于任何  $0 \leq x < n$   $i+x \bmod n$  块上数据单元  $j$  参与  $k+x \bmod n$  单元  $p$  所代。为了更容易地得到  $k$  容，文 [13] 放宽了冗余，只块中，冗余单元不少于户数据单元情况。Weaver 最好冗余只有 50%。

#### 4

列尽有很多性优势，但在前存储中，是 RS 及层叠 RAID (如 RAID1+0) 使得多。为其原因主为以下几个方：

先是实上单性因：RS 已是工业技，无件有成实方，层叠 RAID 原十分单。所以两实施最简单易。与之对，列多、原复，实施一定投入。前存储处于发展，什么是“最好”尚不形成定，因就前，最单就是最好。其，受到前大分应存储影响：尽将多个单个件合成一个一拟件会有好处，但也会有应。如对 10 000 块是合成 1 个好呢？是成 10 个包含 1 000 块小好呢？据判断。一小一些会更容易、护。前只有少应对 1 000 块容数据并处，因将分为多个小是有。

三，价低且发展：使得人们可以“奢侈”地使存储，因大型存储建前处于“旷”。对于易实施性、易护性、易扩展性，当前冗余并不是主决定因。但是，单容日和、对性、容、不断变化，存储应也会不断发展。明天存储将会具备什么性形式，我们不断地探。

#### 5 参文

- [1] HARTLINE J R, RAO K K. Notes on Reliability Models for Non-MDS Erasure Codes [R]. Research Report. RJ10391(A0610-035). San Jose, CA, USA: IBM. 2006.
- [2] 新，国。——原与方 [M]。安：安子技大学出，2001.
- [3] . 存储容及列 [D]。天：南开大学，2010.
- [4] HELLERSTEIN L, GIBSON G A, KARP R M, et al. Coding Techniques for Handling Failures in Large Disk Arrays [J]. Algorithmic, 1994, 12(3/4):182-208.
- [5] BLAUM M, BRADY J, BRUCK J, et al. EVENODD: An Optimal Scheme for Tolerating Double Disk Failures in RAID Architectures [J]. ACM SIGARCH Computer Architecture News, 1994, 22(1):245-254.
- [6] CORBETT P, ENGLISH B, GOEL A. Row-diagonal Parity for Double Disk Failure Correction [C]//Proceedings of the 3rd USENIX Conference on File and Storage Technologies (FAST'04), Mar 31-Apr 2, 2004, San Francisco, CA, USA. Berkeley, CA, USA: USENIX Association, 2004: 14p.
- [7] PLANK J S. The RAID-6 Liberation Codes [C]//Proceedings of the 6th USENIX Conference on File and Storage Technologies (FAST'08), Feb 26-29, 2008, San Jose, CA, USA. Berkeley, CA, USA: USENIX Association, 2008:97-110.
- [8] WANG Gang, LIN Sheng, LIU Xiaoguang, et al. Combinatorial Constructions of



Multi-erasure-correcting Codes with Independent Parity Symbols for Storage Systems [C]//Proceedings of the 13th IEEE Pacific Rim International Symposium on Dependable Computing(PRDC'07), Dec 17-19, 2007, Melbourne, Australia. Piscataway, NJ, USA: IEEE, 2007:61-68.

[9] XU Lihao, BRUCK J. X-code: MDS Array Codes with Optimal Encoding [J]. IEEE Transactions on Information Theory, 1999,45(1):272-276.

[10] XU Lihao, BOHOSSIAN V, BRUCK J, et al. Low Density MDS Codes and Factors of Complete Graphs [J]. IEEE Transactions on Information Theory, 1999,45(6): 1817-1826.

[11] LOUIDOR E, ROTH R M. Lowest-density MDS Codes over Extension Alphabets [C]//Proceedings of the IEEE International Symposium on Information Theory (ISIT'03), Jun 29-Jul 4,2003, Yokohama, Japan. Piscataway, NJ, USA: IEEE, 2003:58.

[12] LIN Sheng, WANG Gang, STONES D S, et al. T-code: 3-erasure Longest Lowest-density MDS Codes [J]. IEEE Journal on Selected Areas in Communications, 2010, 28(2):289-296.

[13] HAFNER J L. WEAVER Codes: Highly Fault Tolerant Erasure Codes for Storage Systems [C]//Proceedings of the 4th USENIX Conference on File and Storage Technologies(FAST'05), Dec 13-16,2005, San Francisco, CA,USA. Berkeley, CA, USA: USENIX Association, 2005.

收 日 : 2010-07-07

作 介

，南开大学 专业博士 业，天 工大学副教授， 方向为存储 、 合

。

刘晓光，南开大学 专业博士 业，南开大学信息技 学学 副教授， 方向为 性 、 信息存储技 。

刚，南开大学 专业博士 业，南开大学信息技 学学 教授， 方向为 信 息存储技 、 并 。