

3 耐故障 RAID の実装

中村 祐司† 上原 稔†

東洋大学総合情報学部総合情報学科†

1 はじめに

今日、ブログなどの CGM の普及、写真ないし動画投稿サイトの流行などにより、かつてないほどの情報がネットワーク上に氾濫している。Google はこれらの情報をインデックス化しようと努力しており、ある程度成功しているように見えるが、まだまだ未分類の情報が実世界に放置されている。最近では、日々の生活を記録する人々が増えている。これはライフログと呼ばれる。ライフログで動画が使われるようになると情報の量はいっそう増大する。仮想世界は実世界の情報をすべて取り込もうとしている。仮想世界が実世界を飲み込もうとするとき、それらの情報を保存するストレージが必要になる。今日では、大量生産される安価な PC を用いて大規模なストレージを実現可能となった。また、HDD 技術の進歩によりストレージ容量も増えている。

大規模ストレージでは多量のディスクを運用する必要がある。しかし、ストレージのディスク数が増加すると信頼性が減少するという問題が生じる。一般にストレージの信頼性を増すには RAID を用いる。RAID の中でも RAID6 は 2 つのパリティで 2 台までの故障に耐える。しかし、大規模ストレージでは 2 耐故障の RAID6 でも不十分である。現在、3 台以上の故障に耐える技法には RAID2 と階層型 RAID の HiRAID(全 RAID による多階層 RAID システム)があるが、RAID2 はシステムが複雑になるため実用化がされていない。HiRAID では 2 階層の RAID5 である RAID55 が 3 耐故障であるが、HiRAID には問題点として容量効率の低さがある。

本研究では同耐故障で階層型 RAID よりも容量効率に優れる RAID である 2 進 RAID および 2 進 RAID を拡張した NaryRAID の実装を目的とする。

第 2 章では信頼性を高めるストレージ技術と仮想的な大規模ストレージを構築する VLSB について、第 3 章では NaryRAID を実装する方法について提案し、第 4 章では今後の課題について述べる。

2 関連研究

2.1 RAID [1][2]

RAID(レイド)とは Redundant Array of Inexpensive(Independent) Disks の略で、データを複数のディスクに分散して記録して、高速化や信頼性の向上を得るための技術であり、ディスクへのデータ配置や、データの冗長化の方法により 0 から 6 までの 7 つのレベルで定義されて、信頼性を増す効

果があるものには 1 耐故障性のある RAID1,3,4,5 および 2 耐故障性のある RAID6 がある。

2.2 階層型 RAID

階層型 RAID は RAID で構築したディスク(論理ドライブ)を 1 つのディスクとして見立てそのディスクでさらに RAID を構築するシステムのことである。構築した階層型 RAID の耐故障性は各層の RAID のデータを修復不可能な台数の積から 1 台引いたものとなるが、容量効率は各層の容量効率の積となる。

HiRAID の代表例として 16 台のディスクからなる RAID55 を図 1 に示す。

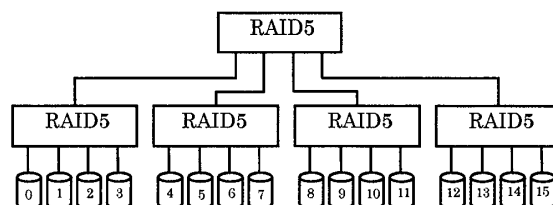


図 1 RAID55

図 1 の 16 台からなる RAID55 の容量効率は $3/4 * 3/4 = 9/16 = 56\%$ で、3 耐故障である。

2.3 2 進 RAID[3]

同ドライブ内で 2 進数の各桁に対応するグループを構成し RAID4 形式でパリティをとる論理ディスクを構築する。2 進 RAID では 2^n 台のデータに対して 2^n 台のパリティを用いるが、3 耐故障であるためには $n \geq 3$ である必要がある。例として $n=3$ の 2 進 RAID の構成を図 2 に示す。8 台のデータディスク $d_0 \sim d_7$ に対して 6 台のパリティディスク $p_0 \sim p_5$ がある。

d	0	1	2	3	4	5	6	7	0	1
2^0	0	1	0	1	0	1	0	1	p_0	p_1
2^1	0	0	1	1	0	0	1	1	p_2	p_3
2^2	0	0	0	0	1	1	1	1	p_4	p_5

図 2 2 進 RAID($n=3$)

ここで、 p_0, p_1 はディスク番号の 0 桁がそれぞれ 0,1 のディスクからなるグループ、 p_2, p_3 はディスク番号の 1 桁がそれぞれ 0,1 のディスクからなるグループ、 p_4, p_5 はディスク番号の 2 桁がそれぞれ 0,1 のディスクからなるグループである。

2.4 NaryRAID

2 進 RAID では 2 進数の各桁に対応するグループを構成していたものを N 進数で構成できるよう

An Implementation of 3 Fault Tolerant RAID

Nakamura Yuji, Minoru Uehara

† Dept. of Information Sciences & Arts, Toyo Univ.

に拡張したもので、構成に必要なディスクの台数のパターンを多くすることができる。必要となるディスクの台数は N 進数では N^n 台のデータに対して Nn 台のパリティを用いる。3 耐故障であるためには 2 進 RAID と同様 $n >= 3$ である必要がある。

2.5 VLSD [4]

VLSD (Virtual Large-Scale Disk) とは PC 教室などの教育環境の数百台の PC で安価で高信頼・大規模なファイルストレージを構築することを実現するためのツールキットである。

VLSD には仮想ディスクに関する `vlsd`、仮想ディスクの入出力に関する `vlsd.io`、ディスクサーバに関する `vlsd.server`、仮想ディスクのテストを行うための `test.vlsd`、ディスクサーバのテストを行うための `test.vlsd.server` などがパッケージされている。

3 提案

本研究ではこの VLSD ツールを用いて NaryRAID を実装する。NaryRAID は RAID4 のパリティをとるグループで構成するため、VLSD の RAID4 クラスを N 進数の各桁に対応するグループで構成するように改良を行った。以下に今回の実装に関するプログラムの説明を行う。

クラス `DiskArray` は複数ディスクからなるディスク・ラッパーの基底クラスである。`DiskArray` を作成し、標準では RAID1 を作成する。外部で引数 `disks` (要素ディスクの配列) を変更しても内部に影響しない。

クラス `RAID` は RAID の基底クラスで簡単な RAID1 の実装を `DiskArray` から引き継いでいる。引数の `disks` から初期故障のディスクを間引いて保持する。

クラス `RAID4` は RAID のサブクラスで RAID4 の実装 (ブロック単位ストライピングとパリティ専用ディスクの作成) を行う。

4 評価および検討

ここでは、JUnit のテストによって NaryRAID の 3 耐故障性能および既存の RAID との読み書き性能を比較、評価する。

まず NaryRAID にさまざまな組み合わせのディスクの故障を与えて 3 耐故障性を実現できているかを確認した。表 1 は NaryRAID に様々な故障を与えた結果で、故障ディスクに書かれた数字がデータディスクの故障の番号、`p` つきの数字がパリティディスクの故障の番号、「*」印は無故障を表している。いずれの場合も問題なく読み書きを行うことができ 3 耐故障性を実現しているといえる。

次に同台数の既存 RAID と NaryRAID との性能の比較を行った。表 2 は NaryRAID ($N=2, n=3$) とその必要台数である 93 台の RAID4,5,6 との読み書き性能の比較である。実験の結果、総ディスク台数が増えたとわずかながら他の RAID より書き込み時間が長くなったが、十分実用可能な範囲であると考え

られる。

表 1 NaryRAID の故障テスト

進数[N]	回数[n]	台数	size	故障ディスク		
2	3	14	1MB	*	*	*
2	3	14	1MB	0	2	p1
2	3	14	1MB	6	p0	p3
2	3	14	1MB	p0	p2	p4
2	3	14	1MB	5	6	4
3	3	36	1MB	*	*	*
3	3	36	1MB	8	23	p7
3	3	36	1MB	0	8	23

表 2 NaryRAID と既存 RAID の読み書き時間

	N	n	Small Read[s]	Small Write[s]	Large Read[s]	Large Write[s]
NaryRAID	3	4	0.437	4.813	2.718	43.469
RAID4			0.437	2.438	2.734	38.172
RAID5			0.453	3.016	2.750	37.984
RAID6			0.416	4.375	2.687	40.610

5 現状

現段階で、VLSD での NaryRAID の実装は完了し、NaryRAID の性能評価として 3 耐故障性の確認と既存の RAID である RAID4、5、6 との読み書き性能の比較を行った。

しかし、階層型 RAID との読み書き性能における比較評価は階層型 RAID のテストがうまくいかなかったため行うことができなかった。

6 今後の課題とまとめ

今回実装した NaryRAID を用いることで容量効率のよい 3 耐故障を実現できる RAID を構築することができる。また NaryRAID は同じディスク台数の RAID4、RAID5、RAID6 とほぼ変わらない時間で読み書きが行えるので十分実用性のある RAID であると考えられる。

今後の課題としては更なる高速化や信頼性の向上などが考えられる。

参考文献

- [1] ストレージ ETERNUS ストレージの基礎用語：富士通
<http://storage-system.fujitsu.com/jp/term/beginner/>
- [2] RAID Level
http://homepage3.nifty.com/sekira/raid_level/raid_level.html
- [3] 上原 稔 "2 進 RAID: もう一つの 3 耐故障 RAID", 信学技報 FIIS-08-228, (2008.3.7)
- [4] チャイエリアント, 上原稔, 森秀樹 "PC 教室のための仮想的な大規模ストレージの構築", マルチメディア、分散、協調とモバイル(DICOMO 2007) シンポジウム論文集, pp.617-622, (2007.7.4-6)