

ارایه الگوی انعکاس داده ای نوین بمنظور بهبود قابلیت اطمینان و تحمل پذیری نسبت به خطا

جواد اکبری ترکستانی^۱

دانشکده مهندسی کامپیوتر ، دانشگاه آزاد اسلامی اراک
akbari@jdmarkazi.ac.ir

محمد رضا میبیدی^۲

دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیر کبیر
meybodi@ce.aut.ac.ir

چکیده

این مقاله یک تکنیک ذخیره سازی داده ای نوین را بمنظور پیاده سازی بر روی سیستم های دنباله دیسکی توزیع شده با نیاز امنیتی بالا ارایه مینماید. هدف از طراحی این الگوی داده ای ، افزایش حفاظت داده ای ، ارایه پهنای باند افزوده ، سطح بالاتر اجراء عملیات موازی دستیابی داده ای و بهبود تحمل پذیری نسبت به خطا است . این روش بگونه ای خاص از ترکیب دو تکنیک انعکاس داده ای و ردیفی کردن (در سطح بلاکی) از تکنیکهای RAID داده ها را بر روی دنباله توزیع مینماید. از این روی ، ما مدل پیشنهادیمان را در مقایسه با سایر مدل های مشابه همچون RAID 5 , RAID 1+0 , RAID 1 مورد ارزیابی قرار داده و دریافته ایم که پیاده سازی مدل مزبور موجب افزایش قابلیت اطمینان و سطح بازگرد پذیری سیستم خواهد شد . ویژگی بارز این مدل افزایش تحمل پذیری سیستم بهنگام بروز خطا و کاهش سطح فقدان داده ای و افزایش نرخ احیاء داده ایست .

واژه های کلیدی : افزونگی داده ای ، تحمل پذیری خطا ، انعکاس داده ای ، قابلیت اطمینان

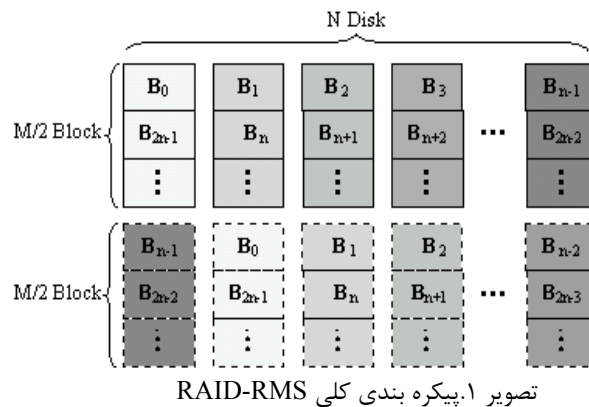
۱- مقدمه

برخی از سیستم ها ، اطلاعات و داده های فوق العاده بحرانی را در بر داشته و بهمین لحاظ در اینگونه از سیستم ها ، حفاظت و امنیت داده ها از اهمیت بسزایی برخوردار است . هر یک از تکنیکهای ذخیره سازی داده ای با توجه به شرایط بحرانی داده درون سیستم و میزان امنیت داده ای مورد نیاز بروشهای مختلفی سطوح گوناگونی از امنیت را برای سیستم فراهم میا ورنند. در تکنیک RAID نیز هر مدل با توجه به ویژگیهای ذاتی خود سطوح گوناگونی از امنیت داده ای را فراهم میاورد . از میان سطوح مختلف تکنیک ، RAID سطح یک دارای افزونگی داده ای نسبتا بالایی بوده اما در مقابل حفاظت ، صحت و امنیت داده ای را در سطح خوبی فراهم مینماید. ایراد اصلی این مدل در سیستم هایی که نیاز چندانی به امنیت داده ای ندارند ، در درجه اول حجم بالای اطلاعات افزوده ای است که به سیستم تحمیل میگردد. اما در سیستم های با نیاز امنیتی بالاتر عدم توازی گرابی عملیاتیهای داده ای سیستم ، عدم استفاده از حداکثر پهنای باند و حتی پایین بودن سطح امنیتی داده ها بهنگام بروز خرابی و احیاء داده ای از معایب مدل فوق میباشد. نا گفته نماند که حفاظت داده ای متناسب با میزان افزونگی هر مدل ، در سایر سطوح نیز فراهم است. برای مثال در سطح پنجم از این تکنیک با توجه به ویژگیهای خاص مدل امنیت داده ای در سطح خوبی برای سیستم فراهم میاید. در RAID 1+0 که از ترکیب دو سطح پایه ای (سطوح یک و دو) ایجاد میگردد تا حدودی سعی گردیده تا با ادغام دو تکنیک Mirroring و Stripping ، سیستمی پیاده سازی گردد که علاوه بر سطح امنیتی بالا ، توازی گرابی و سایر فاکتورهای ذکر شده را نیز بهبود بخشد. این مدل نیز خود دارای کاستیهایی است که در بخشهای بعدی بطور مفصل به بررسی آنها میپردازیم[4][1].

^۱ عضو هیات علمی دانشکده مهندسی کامپیوتر، دانشگاه آزاد اسلامی واحد اراک

^۲ عضو هیات علمی دانشکده مهندسی کامپیوتر و فناوری اطلاعات، دانشگاه صنعتی امیر کبیر

با توجه به ویژگیها، مزایا، معایب و نقاط قوت و ضعف مدل‌های قبلی، ما سعی داریم مدلی ارائه دهیم که علاوه بر فراهم آوردن سطح بالایی از امنیت، صحت و حفاظت داده‌ای بتواند حداکثر توافقی گری، حداکثر پهنای باند موثر بهنگام عملیات‌های دستیابی داده^۳ (عملیات‌های ورودی/خروجی) و احیاء داده‌ای و حداکثر تحمل پذیری خطا را به گونه‌ای برای سیستم فراهم آورد، تا هم برای پیاده‌سازی روی سیستم‌های با نیاز امنیتی بالا مناسب باشد و هم بازدهی سیستم را به میزان خوبی افزایش دهد. در ادامه مقاله، در بخش ۲ به تشریح مدل پیشنهادی خود پرداخته و پس از آن پیاده‌سازی مدل را ارائه نموده ایم. در بخش ۴ مدل ارائه شده را با سایر مدل‌های مشابه از جنبه‌های مختلفی همچون هزینه فضای وزمانی مدل، حداکثر پهنای باند، سطح توافقی سیستم و هزینه احیاء داده‌ای مورد مقایسه قرار داده و در بخش‌های بعدی نیز با شبیه‌سازی مدل و ارائه یک نتیجه‌گیری به مقاله خاتمه داده ایم [3][2].



۲- مدل RAID-RMS

مدل RAID پیشنهادی ما که آنرا RAID-RMS نامیده ایم از ترکیب دو تکنیک ردیفی کردن^۴ داده‌ای در سطح بلاکی و انعکاس داده‌ای^۵ در سطح دیسک‌های سیستم استفاده می‌نماید. اما تفاوت این مدل با سایر مدل‌ها در آنست که این دو تکنیک را بگونه‌ای با یکدیگر در هم می‌آمیزد که موجب افزایش سطح اجراء عملیات موازی داده‌ای، تحمل پذیری خطا و ایجاد پهنای باند بالاتر سیستم گردد. علاوه بر این، در مدل پیشنهادی، الگوی ذخیره‌سازی خاصی پیاده‌سازی می‌شود که قادر است، قابلیت احیاء داده‌ای^۶ را افزایش، داده بگونه‌ای که بهنگام بروز خرابی حداقل داده از دست رود، و انجام عملیات احیاء داده‌ای را تسریع نموده و علاوه بر آن بهنگام بروز خرابی بازدهی عملیات دستیابی داده‌ای در سیستم کمترین افت را داشته باشد [4][3].

همانطور که در تصویر ۱ هم مشاهده می‌گردد، با استفاده از تکنیک انعکاس داده‌ای، بلاک‌های سیستم را به دو دسته بلاک‌های داده‌ای و بلاک‌های پشتیبان^۷ تقسیم بندی نموده ایم و بر خلاف مدل‌های مشابه RAID 1+0، RAID 1 که در آنها بلاک‌های داده‌ای و پشتیبان در دیسک‌های متمایزی قرار دارند، در این مدل داده‌های اصلی و نسخه‌های پشتیبان آنها بر روی دیسک‌های یکسانی توزیع می‌شوند و در واقع ما نیمه بالایی هر دیسک را برای بلاک‌های داده‌ای اصلی و نیمه دیگر از آنها برای بلاک‌های داده‌ای پشتیبان در نظر گرفته ایم. آنچه در این میان از اهمیت بالایی برخوردار است نحوه قرار دادن و توزیع داده‌ها بر روی دنباله دیسکی است، که ما در بخش‌های بعدی برای این نحوه توزیع داده‌ای را در قالب فرمول‌های کلی ارائه مینماییم. لازم بذکر است که در این مدل برای توزیع داده‌ای از تکنیک انتقال براست^۸ بلاک‌های داده‌ای و پشتیبان در هر بلاک و هر دیسک متناسب با سطح بلاک و شماره دیسک استفاده می‌شود [4][3][1]. در این مدل از توزیع داده‌ها، بلاک‌های داده‌ای مختلف و نسخه‌های پشتیبان آنها بر روی دیسک‌های

³ I/O operation

⁴ Data Stripping

⁵ Mirroring

⁶ Data Recovery

⁷ Mirrored Block

⁸ Rotate - Right

مختلف دنباله بطریقی گسترده میشوند که بهنگام بروز خرابی در سطح دیسکی کمترین داده های سیستم مفقود گردد. و علاوه بر آن ، سرعت بازدهی سیستم بهنگام احیاء بلاکهای داده ای از دست رفته افزایش میابد. بر همین اساس برای پیاده سازی بر روی محیطهای ذخیره سازی داده ای که هم نیاز به امنیت وصحت داده ای نسبتا بالایی داشته وهم سرعت دستیابی داده ای وبازدهی سیستم در آنها بالاست مناسب است .

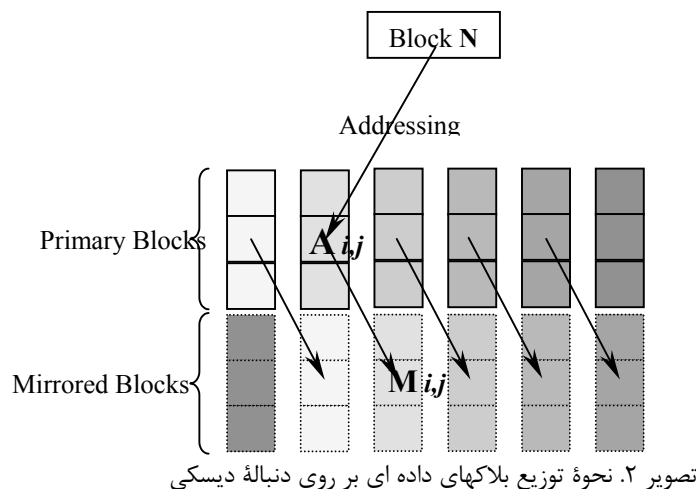
۳- پیاده سازی مدل RAID-RMS^۹

همانگونه که در بخشهای قبلی نیز گفته شد، مدل پیشنهادی از ادغام دو تکنیک انعکاس داده ای (Mirroring) و ردیفی کردن داده ای (Stripping) به همراه تکنیکهایی از انتقال بلاکهای داده ای استفاده مینماید. ما قصد داریم تا در ادامه نحوه توزیع داده ها را در این مدل برای حالات کلی بیان نموده و فرمولهایی را بکمک پارامترهای جدول ۱ برای ایجاد الگوهایی بمنظور برقراری انطباق میان داده های پیوسته فضای حافظه با آدرسهای فضای فیزیکی مجازی روی دنباله دیسکی ارایه میدهیم . برای همین منظور به تجزیه چگونگی آدرس دهی (نگاشت شماره بلاکهای داده ای بر روی آدرسهای دنباله دیسکی) در این مدل پرداخته و آنرا در قالب فرمولهایی بمنظور نگاشت آدرس فضای حافظه تک بعدی (شماره بلاک داده ای) به آدرس فضای دو بعدی دنباله دیسکی ، محاسبه آدرس بلاک داده ای و آدرس نسخه انعکاس یافته آن (بلاک پشتیبان متناسب با آن) از روی بلاک داده ای بیان مینماییم .

جدول ۱. جدول نمادها

نماد	توضیحات
N	شماره بلاک داده ای
n	تعداد دیسکهای موجود در دنباله
m	تعداد بلاکهای موجود در هر دیسک
i	موقعیت افقی بلاک داده ای درون دنباله
j	موقعیت عمودی بلاک داده ای درون دنباله

همانطور که در تصویر ۲ نیز مشاهده میگردد بلاکهای داده ای را با استفاده از تکنیک ردیفی کردن داده ای در سطح بلاکی بر روی دیسکهای مختلف ومتوالی دنباله توزیع مینماییم که در هر سطح بلاکی از هر دیسک از تکنیک انتقال داده ای متناسب با همان سطح استفاده میگردد . از سوی دیگر نیز برای هر داده ای که بر روی یک بلاک داده ای از دنباله دیسکی کپی میگردد، نسخه ای بر روی یکی از بلاکهای پشتیبان که آدرس آن با توجه به فرمولهای زیرین بدست می آید قرار داده می شود و در واقع بدین نحو از تکنیک انعکاس داده ای نیز در این مدل استفاده میشود. پس از پیاده سازی الگوهای سخت افزاری ساختار مدل پیشنهادی ، بایستی یک نگاشت بمنظور برقراری ارتباط میان شماره بلاکهای داده ای و فضای آدرس دهی مجازی دنباله دیسکی تعریف کرد، تا بهنگام قراردادن بلاکهای داده ای بر روی آدرس دو بعدی دنباله دیسکی ، شماره بلاک متناظر با هر زوج (i, j) را بتوان بر اساس این نگاشت محاسبه کرد . تعداد بلاکهای موجود در هر دیسک را m و تعداد دیسکهای دنباله را n در نظر گرفته ایم. ما در این مرحله سه نوع نگاشت برای ذخیره سازی داده ها روی دنباله دیسکی محاسبه مینماییم که عبارتند از :



تصویر ۲. نحوه توزیع بلاکهای داده ای بر روی دنباله دیسکی

^۹Rotating – Mirroring - Striping

۱. محاسبه نگاشت آدرس دنباله دیسکی به شماره بلاک داده ای $N_{A(i,j)}$
۲. محاسبه نگاشت شماره بلاک داده ای به آدرس دنباله دیسکی بلاک اصلی $A_{(i,j)}$
۳. محاسبه نگاشت شماره بلاک داده ای به آدرس دنباله دیسکی بلاک پشتیبان $M_{(i,j)}$

ابتدا نگاشتی را محاسبه مینماییم که آدرس بلاک داده ای روی دنباله دیسکی را بکمک فرمولی به شماره بلاک ذخیره شونده در آن آدرس تبدیل نماید .

$$N_{A(i,j)} = \begin{cases} (i \times n) + (j \bmod n) - i & ; i \leq j \\ 0 \leq i \leq \frac{m}{2} - 1, 0 \leq j \leq n - 1 & \text{فرمول ۱} \\ (i + 1) \times n + j - i & ; i > j \end{cases}$$

در فرمول شماره ۱ ابتدا شماره نسخه اصلی بلاک داده ای مورد نظر بر اساس پارامترهای مختلف تعیین و سپس بکمک فرمولهای دیگری آدرس وموقعیت بلاک داده ای پشتیبان محاسبه میگردد بر همین اساس دامنه تغییر پارامتر i از ۰ تا $m/2$ می باشد چرا که نیمه بالایی بلاکهای درون هر دیسک را برای نگهداری بلاکهای داده ای اصلی در نظر گرفته ایم. بجهت آنکه برای توزیع داده ها در سطح دنباله دیسکی پس از بکارگیری تکنیک انعکاس وردیفی کردن داده ای ، داده ها را در سطح بلاکی تحت یک چرخش براست قرار میدهیم. بهمین لحاظ بلاکهای داده ای که آدرس افقی آنها از آدرس عمودیشان بیشتر است به ابتدای آن سطح منتقل می شوند و در فرمول ، ما آنها را با شرط $i > j$ از سایر بلاکها متمایز مینماییم .

بکمک فرمول ۱ می توان مشخص نمود که در هر موقعیت افقی وعمودی از دنباله دیسکی ، کدامین بلاک داده ای میبایستی ذخیره گردد. اما بهنگام انتقال داده ها بر روی دنباله دیسکی آنچه از اهمیت بالاتری برخوردار است محاسبه یک نگاشت معکوس برای فرمول فوق میباشد که بر اساس آن بتوان موقعیت افقی وعمودی هر بلاک داده ای را بهنگام ذخیره سازی آن بر روی دنباله مشخص نموده و داده مورد نظر را در آن بلاک قرار داد. که بسادگی وبکمک فرمول ۲ آدرس افقی (شماره بلاک درون دیسک) و آدرس عمودی (شماره دیسک درون دنباله) محاسبه میگردد.

$$A_{(i,j)} = \begin{cases} i = \left\lfloor \frac{N}{n} \right\rfloor \\ 0 \leq i \leq \frac{m}{2} - 1, 0 \leq j \leq n - 1 & \text{فرمول ۲} \\ j = \left[\left\lfloor \frac{N}{n} \right\rfloor + (N \bmod n) \right] \bmod n \end{cases}$$

در این روش بمنظور افزایش سطح حفاظت وامنیت داده ای از تکنیک انعکاس داده ای (کپی های دوگانه داده ای) استفاده میگردد. یعنی بهنگام نوشتن هر بلاک داده ای یک کپی از آن نیز بعنوان نسخه پشتیبان بر روی دیسک نگهداری میگردد تا بهنگام بروز خطا در یکی از بلاکهای مورد نظر داده مزبور از روی نسخه دیگر واکشی شود در این مدل پیشنهادی بلاکهای داده ای پشتیبان بر روی بلاکهای نیمه پایینی هر دیسک ($m/2 \leq i \leq m - 1$) و با استفاده از روابط حاکم در فرمول ۳ قرار داده می شوند.

پس از بکارگیری فرمولهای ارایه شده بمنظور تعیین موقعیت بلاکهای داده ای ، اگر به نحوه توزیع داده ای اصلی وپشتیبان نگاه نماییم بلاکهای داده ای پشتیبان متناظر با هر بلاک داده ای در نیمه پایینی دیسک بعدی قرار میگیرد، که به این ترتیب در اثر از کار افتادن یک دیسک منفرد ، دنباله دیسکی باز هم می تواند بمنظور دستیابی داده ها با حداکثر پهنای باند فراهم شده توسط سایر دیسکهای باقیمانده در دنباله بکار خود ادامه دهد. ویا آنکه بهنگام بروز یک خطا در سطح دیسک ، داده های از دست رفته با سرعت بالاتری قابل احیاء بوده و حتی اگر دو دیسک متوالی از دنباله بطور همزمان دچار خرابی شوند، نسبت به سایر مدلهای دیگر کمترین داده موجود درون سیستم از دست خواهد رفت .

$$M_{(i,j)} = \begin{cases} i' = \left\lfloor \frac{N}{n} \right\rfloor + \frac{m}{2} \\ j' = \left[\left\lfloor \frac{(N+1)}{n} \right\rfloor + ((N+1) \bmod n) \right] \bmod n \end{cases} \quad \frac{m}{2} \leq i' \leq m-1, \quad 0 \leq j' \leq n-1 \quad \text{فرمول ۳}$$

۴- ارزیابی مدل پیشنهادی

همانگونه که در بخشهای قبلی نیز به آن اشاره گردید مدل پیشنهادی ما از تلفیق تکنیک انعکاس داده ای با مدل ردیفی کردن داده ای در سطح بلاکی استفاده مینماید و در واقع مدل مزبور برای پیاده سازی در سیستم هایی مناسب است که اطلاعات و داده های آنها نیاز به سطح امنیتی بالایی دارند. بنابراین بمنظور ارزیابی این مدل آنرا با سایر تکنیکهای ذخیره سازی داده ای مطرح در سیستم های امن مقایسه مینماییم. تکنیک سطح یک (RAID 1) بدلیل امنیت بالا و تشابه در استفاده از تکنیک Mirroring و سطح پنج از این تکنیک (RAID 5) بجهت بکارگیری تکنیک Striping و تکنیک ترکیبی RAID 1+0 سه تکنیکی هستند که ما مدل پیشنهادیمان را با آنها مقایسه نموده و در نهایت به بررسی نتایج بدست آمده و ارزیابی نهایی مدل خواهم پرداخت. در این بخش ما سعی داریم مدل پیشنهادی را با سایر مدل‌های مشابه قبلی از جنبه های گوناگونی همچون هزینه فضای افزوده مورد نیاز بمنظور تأمین امنیت و حفاظت داده ای^{۱۰} مدل، هزینه عملیات ورودی/ خروجی مورد نیاز برای دستیابی به داده ها (خواندن و نوشتن داده)، امکان کنترل همزمانی^{۱۱}، تحمل پذیری مدل نسبت به خطا^{۱۲} و هزینه احیاء خرابی دیسک و داده ارزیابی نموده و در هر مورد مزایا و معایب مدل مورد بررسی را ارایه کرده و آنها را با همدیگر مقایسه نموده و در صورت امکان راه حل هایی را هم برای مرتفع ساختن ناکارآمدیهای مدل در مواردی خاص ارایه میدهیم.

۴-۱- حداکثر پهنای باند^{۱۳}

همانگونه که در بخشهای قبلی نیز بدان اشاره گردید یکی از نتایج ارزشمندی که بواسطه پیاده سازی تکنیک RAID فراهم میگردد، پهنای باند بیشتر برای سیستم است. در این تکنیک سعی میگردد با بکارگیری روشهای مختلف توزیع داده ای، امکان انجام عملیات ورودی/ خروجی موازی را افزایش داده و بدین ترتیب به حداکثر پهنای باند برای انتقال داده ها در سیستم دست یافت. در این بخش ما تکنیکهای مطرح شده را مورد بررسی قرار داده و حداکثر پهنای باند مؤثری را که هر یک از این مدلها برای سیستم فراهم میاورند محاسبه نموده و در نهایت آنها را با یکدیگر مقایسه مینماییم. برای این منظور ابتدا برخی از نمادهای بکار رفته در مقایسات را در قالب جدول ۲ ارایه میدهیم.

جدول ۲. جدول نمادها

نماد	توضیحات
N	تعداد دیسک در یک سیستم RAID توزیع شده
B	حداکثر پهنای باند هر دیسک
S	اندازه هر بلاک از دیسک
R	متوسط زمان خواندن هر بلاک از دیسک
W	متوسط زمان نوشتن هر بلاک از دیسک
M	تعداد بلاکهای موجود در هر قابل

پهنای باند فراهم آورده شده در هر یک از مدل‌های مورد بررسی برای حالات مختلف و بازاء عملیتهای مختلف دیسکی مقادیر متفاوتی را شامل میگردد، از این روی ما در تحلیل پهنای باند مؤثر سیستم عملیتهای خواندن و نوشتن داده ای را بطور جداگانه

¹⁰Data Protection

¹¹Concurrency Control

¹²Fault tolerance

¹³Maximum Bandwidth

مورد بررسی قرار داده ایم و حتی در برخی موارد که عملیات دیسکی (ورودی/ خروجی) برای فایلها و مقادیر داده ای طولانی و کوتاه نتایج متفاوتی را در بر دارد. آنها را از یکدیگر متمایز نموده و تحلیل کرده ایم [4][3].

مفاهیم Small Read / Write و یا Large Read / Write که در جدول ۳ مشاهده می شود بنابر نیاز تقسیم بندی عملیات دیسکی که در بالا گفته شد صورت گرفته است. منظور از خواندن یا نوشتن کوتاه آن دسته از عملیاتیهای ورودی / خروجی داده ای هستند که تنها بر روی تمام یا قسمتی از یک بلاک داده ای تاثیر می گذارند و در واقع عملیاتیهای موازی ورودی / خروجی که از این نوع هستند می توانند روی بلاکهای متفاوت از یک Stripe منفرد اعمال شده و یا بر بلاکهای متفاوت یا یکسان از Stripe های متفاوت تاثیر گذارند. که هر یک از انواع فوق نیز در هر یک از تکنیکهای مورد بحث ما نتایج متفاوتی را در پی دارند. با این همه دامنه عملیاتی آنها از یک بلاک فراتر نمی رود. اما منظور از خواندن و نوشتن های طولانی آن دسته از عملیات دیسکی را شامل می شود که در آنها داده ها بر روی چندین بلاک متوالی از یک یا چند Stripe متوالی قرار داده شده و یا از آنها خوانده می شوند [12][11].

همانطور که در جدول ۳ هم مشاهده می شود، تمامی سطوح RAID مورد بررسی برای عملیات خواندن کوتاه پهنای باند $N * B$ را فراهم میاورند چرا که تعداد N عملیات داده ای Small Read قادرند بطور همزمان روی بلاکهای N دیسک مختلف انجام گیرند. برای عملیات Large Read نیز در RAID 1 بدلیل استفاده از تکنیک Mirroring داده ها و کپی آنها روی بلاکهای متوالی دیسک قرار داده میشوند و این شرایطی را فراهم میاورد تا هر پردازنده که عملیات Large Read را انجام میدهد بتواند بلاکهای داده ای متوالی را بطور یک در میان از روی دیسک اصلی و پشتیبان بطور موازی و همزمان بخواند و بدین ترتیب بهنگام اجراء $N / 2$ پردازنده موازی روی دیسکهای مختلف از پهنای باند نیمی از دیسکهای دنباله N دیسکی استفاده کامل میگردد. اما در سایر موارد بدلیل استفاده از تکنیک Stripping به مدل امکان انجام عملیات داده ای موازی به مقدار قابل توجهی افزایش میابد.

در RAID 1 انجام هر عملیات Large Write بدلیل آنکه داده ها ی منطقاً متوالی روی بلاکهای فیزیکی متوالی هر دیسک گسترده شده اند تنها دو دیسک را درگیر کرده و پهنای باند $2 * B$ را مورد استفاده قرار میدهد. اما برای اجراء عملیات Large Write پردازنده های موازی که هر یک بر روی یک دیسک از دنباله اثر میکنند پهنای باند سیستم تا $(N/2)*B$ نیز افزایش میابد. بغیر از مدل RAID 5 که بدلیل Stripping در سطح بلاکی و استفاده از یک بلاک توازن در مقابل $N - 1$ دیسک برای نگهداری اطلاعات داده ای، حداکثر پهنای باند $(N - 1) * B$ میباشد سایر مدلها پهنای باندی همچون RAID 1 خواهند داشت [5][2][1].

۴-۲- هزینه زمانی عملیات ورودی / خروجی^{۱۴}

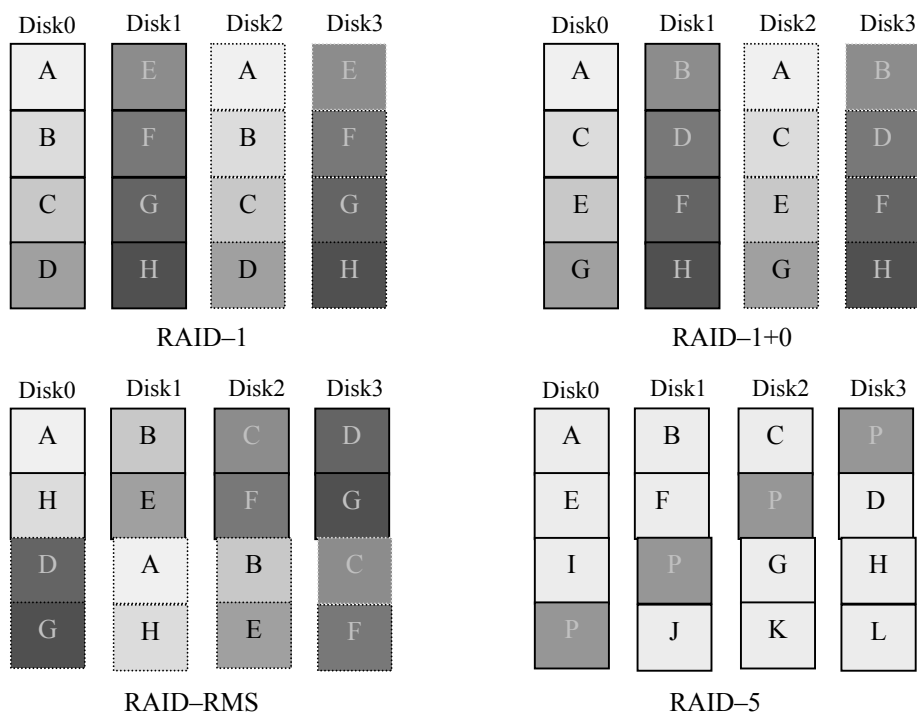
در این مرحله به ارزیابی هزینه های زمانی اجراء عملیات ورودی / خروجی روی دنباله دیسکی پرداخته و همچون مراحل قبلی هزینه مدلهای گوناگون را بازنه حالات مختلف عملیات دیسکی مورد بررسی قرار میدهیم. تمامی سطوح مورد بررسی برای انجام عملیات Small Read بایستی هزینه زمانی R را پردازند. اما در صورتیکه بخواهیم یک فایل یا یک دنباله داده ای بطول M را بر روی دنباله دیسکی توزیع شده درج نماییم در مواردی که بلاکهای داده ای بر روی دنباله دیسکی Stripe شده باشند بسادگی از فرمول زیر برای تعیین هزینه زمانی مدل برای عملیات Large Read می توان استفاده کرد.

$$(L_{Data} * T_R) / (N_{Total Disk} - N_{Check Disk})$$

$N_{Total Disks}$: تعداد کل دیسکهای دنباله
 $N_{Check Disks}$: تعداد دیسکهای حاوی اطلاعات افزوده
 L_{Data} : طول فایل یا دنباله داده ای بر حسب بلاک
 T_R : هزینه زمانی متوسط عملیات خواندن یک بلاک

در این میان تنها مورد استثناء RAID سطح یک است [8] که با استفاده از تکنیک Mirroring بلاکهای داده ای منطقاً متوالی را بر روی بلاکهای متوالی فیزیکی بصورت دوگانه قرار می دهد. در این روش امکان ذخیره سازی موازی بلاکهای متوالی بلحاظ ساختار ذخیره سازی موجود فراهم نمیشود و تمامی یک داده M بلاکی بطور پی در پی در یک دیسک ذخیره و نگهداری می شود و بهمین لحاظ هزینه زمانی عملیات متوالی خواندن یک داده M بلاکی بصورت $M * R$ محاسبه میگردد.

¹⁴Input / Output Cost



تصویر ۳. ساختار توزیع داده ای در مدل پیشنهادی وسایر مدل‌های مشابه

۳-۴ - هزینه عملیات احیا، داده ای^{۱۵}

همانگونه که در بخشهای ابتدایی نیز بدان اشاره گردید، مزیت‌های مدل پیشنهادی را نسبت به سایر مدل‌های مشابه می توان در مواردی همچون ، تحمل پذیری بالاتر نسبت به خطا ، صحت و حفاظت داده ای بیشتر ، توازی گرابی بالاتر به همراه حداکثر پهنای باند موثر بهنگام انجام عملیات‌های دستیابی ویا احیاء داده ای خلاصه کرد. در این بخش ما قصد داریم مدل‌های مورد نظر را در موارد فوق بررسی ویا یکدیگر مقایسه نماییم [16][14][13].

به لحاظ صحت و حفاظت داده ای ، مدل پیشنهادی ما از هر سه مدل دیگر محیط امن تری را برای داده های درون سیستم فراهم میاورد . چرا که در دو مدل RAID 1 و RAID 1+0 بهنگام بروز خرابی در بیش از یک دیسک ، بدلیل آنکه تمامی بلاکهای داده ای پشتیبان دیسک مزبور دقیقاً بر روی یک دیسک دیگر متمرکز شده اند بنا بر این ، در یک سیستم دنباله دیسکی با N دیسک ، با احتمال $1/(N-1)$ داده های از دست رفته قابل بازگرداندن نمیباشد. این در حالی است که در این دسته از سیستم ها در برخی از حالات خاص نیز حتی تا $N/2$ خطای دیسکی نیز قابل پوشش بوده و داده های از دست رفته قابل احیاء میباشد. در مدل RAID 5 نیز بهنگام بروز خرابی در یک دیسک، تشخیص خطا ویا تعیین موقعیت خطاهای بروز کرده تا حدود زیادی دشوار و گاه در مواردی همچون تعداد خطاهای زوج غیر ممکن است و در واقع این مدل نسبت به دو مدل قبلی نیز امنیت کمتری دارد . اما در مدل پیشنهادی بجهت استفاده از نوعی روش انتقال بلاکی خاص که با دو تکنیک انعکاس داده ای وریفی کردن داده ای همراه شده ، داده ها بگونه ای بر روی دنباله دیسکی توزیع می شوند که بهنگام بروز هر نوع خطای داده ای و دیسکی در سطح دنباله کمترین داده از دست خواهد رفت . نکته جالب در این روش آنست که بدلیل توزیع مورب نیمه های دیسکی اگر خطاهای بروز کرده فقط بر روی دیسکهای با شماره زوج ویا فقط بر روی دیسکهای با شماره فرد گسترده شده باشند، هر نوع وهر تعداد از خرابیها و خطاهای دیسکی (حتی تا $N/2$ خطا نیز) قابل بازگرداندن بوده و موجب از دست رفتن هیچ داده ای نمی شود. در این روش نرخ احتمال خطای غیر قابل برگشت کاهش میابد و با احتمال $1/((N-1) * (N-2))$ داده های از دست رفته یک دیسک قابل احیاء، نمیباشد (به [4] مراجعه گردد) .

¹⁵Recovery Cost

در یکی از بخشهای قبلی به ارزیابی پهنای باند سیستم پرداخته و حداکثر پهنای باند موثر را درحالی که سیستم عملکردی عادی دارد(حالتی که خطا در سیستم بروز نکرده) بزاء عملیتهای داده ای کوتاه و طولانی برای مدل‌های مختلف محاسبه نموده ایم . اما آنچه را که ما در این قسمت مورد توجه قرار میدهیم، حداکثر پهنای باندی است که سیستم در هنگام بروز خطا بطور مؤثر برای پردازشها فراهم میاورد. و در خلال این بررسی به دو سؤال زیرپاسخ میدهیم.

۱. هنگام بروز خرابی در یک یا عده ای از دیسکهای دنباله، مدل بکار رفته قادر است تا چه مقدار از پهنای باند سایر دیسکهای دنباله را برای انجام عملیات دیسکی مورد استفاده قرار دهد؟

۲. هنگام بروز خرابی چه مقدار از پهنای باند مؤثر سیستم برای انجام عملیات احیاء دسترس پذیر است؟

در مدل RAID 1، در شرایط عادی، هر پردازش تنها قادر است برای انجام عملیات Large Read تنها از پهنای باند دودیسک استفاده نماید که این در مقابل مدل پیشنهادی ما بسیار ضعیف است اما در چنین شرایطی RAID 1+0 پهنای باندی مشابه مدل پیشنهادی ما دارد. هنگام بروز خطای منفرد پهنای باند مؤثری که برای انجام عملیات Small Read مورد استفاده قرار میگیرد در تمامی مدلها مشابه یکدیگر و برابر است با:

$$(N-1) * B$$

برای عملیات خواندن فایلها بزرگ و یا داده های طولانی پهنای باند مدل پیشنهادی ما نتیجه جالب توجهی را ارائه کرده و پهنای باندی حتی بیش از مدل RAID 5 را فراهم میاورد. نتایج بررسی و مقایسه مدل‌های مختلف در جدول ۳ آورده شده است. پس از بروز هر خرابی در یک سیستم دنباله دیسکی بلافاصله بایستی داده هایی که بر روی دیسک مزبور قرار داشته اند احیاء و تا هنگام تصحیح و یا جایگزینی دیسک خراب بر روی یک دیسک یدکی نگهداری شود. در خلال عملیات احیاء داده ای، انجام عملیات عادی سیستم متوقف میگردد. بنابر این هر مدلی که مدت زمان کوتاه تری را صرف انجام عملیات احیاء نماید مناسب تر بوده و بازدهی بالاتری را برای سیستم فراهم میاورد. در تکنیک RAID 1+0، RAID 1، هنگام عملیات احیاء دیسکی تنها از پهنای باند یک دیسک استفاده میگردد و عملیات احیاء بصورت سریال انجام میگردد [9][7][6]. اما در RAID-RMS بدلیل نحوه توزیع داده ها بر روی دنباله دیسکی برای عملیات احیاء N-1 دیسک باقیمانده قادرند بطور موازی و همزمان داده های دیسک خراب را احیاء نمایند [18][17][10].

جدول ۳. جدول مقایسه هزینه های مدل پیشنهادی و سایر مدل‌های مشابه

	RAID Level	Small Read	Large Read	Small Write	Large Write
Max.I/O Bandwidth	RAID 1	$N * B$	$(N/2) * B$ or $2 * B$	$(N/2) * B$	$(N/2) * B$ or $2 * B$
	RAID 1+0	$N * B$	$N * B$	$(N/2) * B$	$(N/2) * B$
	RAID 5	$N * B$	$(N-1) * B$	$(N/2) * B$	$(N-1) * B$
Read / write Time (parallel)	RAID 1	R	$M * R$	W	$M * W$
	RAID 1+0	R	$M / N * R$	W	$M * W / (N/2)$
	RAID 5	R	$M * R / (N-1)$	$R + W$	$M * (R + W) / (N-1)$
	RAID-RMS	R	$M * R / N$	W	$M * W / (N/2)$
Recovery Max.I/O Bandwidth	RAID 1	$(N-1) * B$	B		
	RAID 1+0	$(N-1) * B$	$(N/2) * B$		
	RAID 5	$(N-1) * B$	$(N-2) * B$		
	RAID-RMS	$(N-1) * B$	$(N-1) * B$		

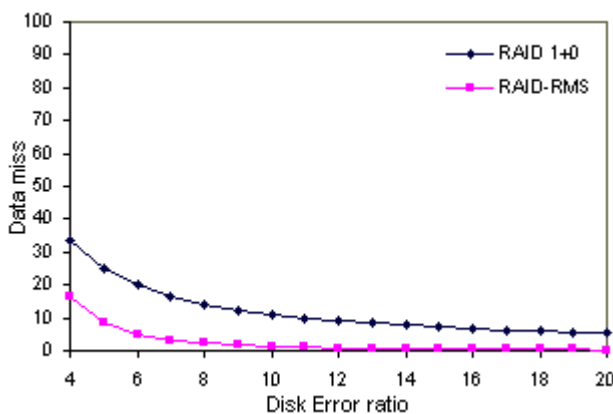
۵- شبیه سازی مدل پیشنهادی

در این فرایند شبیه سازی بمنظور ارزیابی دقیقتر نتایج و داده های حاصل از بررسیهای تحلیلی انجام شده بر روی مدل پیشنهادی و سایر مدل‌های مشابه، و مطالعه نحوه عملکرد و رفتار واقعی مدل‌های تحت بررسی در مقابل تغییرات پارامترهای مختلف، مدل پیشنهادی خود را به همراه تکنیک RAID 1+0، RAID 1 در یک محیط برنامه سازی سطح بالا شبیه سازی نموده ایم. بدلیل تشابه RAID 1، RAID 1+0 و آنکه در پاره ای از موارد RAID 1+0 به لحاظ احیاء، داده ای نتایج بهتری را ارائه میدهد، در شبیه سازی از این مدل استفاده گردیده است. مدل‌های تحت بررسی را در معرض گونه های مختلف خطا های داده ای در سطح بلاکی و دیسکی قرار داده، رفتارها، نحوه عملکرد، میزان حفاظت و امنیت داده ای و تحمل پذیری مدلها نسبت به گونه های مختلف خطا بررسی و نتایج حاصل، در قالب نمودارهایی ارائه گردیده است.

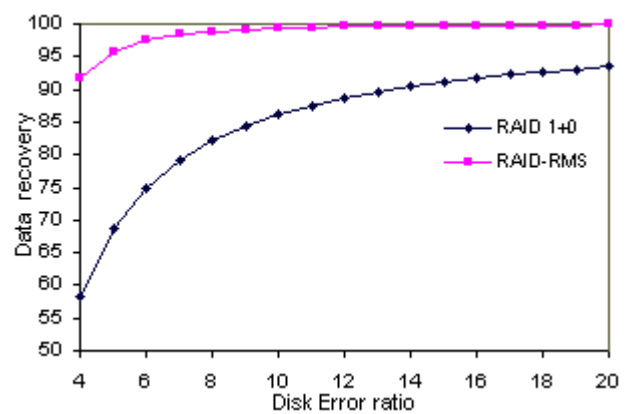
در مرحله نخست هر دو مدل را در معرض گونه های مختلفی از خطا ها قرار داده و پاسخ سیستم را برای دامنه ای از خطا های بروز کرده در سیستم مشخص نموده و بمنظور افزایش سطح دقت در تعیین نتایج، انجام آزمون را برای هر مدل تا ۱۰۰۰۰ بار تکرار نموده ایم. در حقیقت در این مرحله از آزمون قصد داریم مشخص نماییم که هر مدل قادر است تا چند درصد از خطا های بروز کرده را بطور کامل احیا، نماید.

همانگونه که در نمودار ۱ هم مشاهده می شود، مدل پیشنهادی در شرایط یکسان (بازا، تعداد خطا و دیسک) نرخ احیا، داده ای بسیار بالاتری را داراست. هر دو مدل در مقابل تغییرات افزایشی تعداد دیسکهای درون دنباله تغییرات افزایشی مشابهی را از خود نشان میدهند. با توجه به نمودار ۱ میتوان دریافت که مدل RAID1+0 بازا، تعداد دیسکهای بیشتر درون دنباله نرخ احیا، بسیار خوبی را فراهم میآورد. اما مدل پیشنهادی بازا، تعداد دیسکهای کمتر نیز مؤفق عمل مینماید.

در مرحله بعد رفتار ونحوه عملکرد هر دو مدل را در مقابل نرخ تغییرات تعداد مشابه خطا و تعداد دیسکهای درون دنباله آزموده و نتایج حاصله را در نمودار ۲ ارائه کرده ایم. با کمی دقت میتوان گفت که در مدل پیشنهادی در بدترین حالت وبا بیشترین نسبت تعداد خطا به دیسک درون دنباله کمتر از ۲۰٪ داده از دست خواهند رفت. اما در هر دو مدل با افزایش نسبت تعداد دیسک در مقابل تعداد خطا های یکسان، تغییرات کاهشی مشابهی را از خود نشان میدهند ونرخ داده های از دست رفته بمقدار قابل توجهی کاهش خواهد داشت. با بررسی نتایج حاصل از شبیه سازی مدل پیشنهادی و انطباق آنها با نتایج تحلی که در مراحل قبلی در قالب فرمولهایی ارائه شده می توان گفت مدل پیشنهادی در شرایط یکسان ودر مقایسه با سایر مدلهای مشابه نرخ احیا، داده ای بسیار بالاتری را فراهم میآورد.



نمودار ۲. مقایسه میزان فقدان داده



نمودار ۱. مقایسه نرخ احیاء داده ای

۶- نتیجه گیری

با توجه به مطالعات و بررسیهای انجام شده و بر اساس نتایج تحلیلی بدست آمده در بخشهای قبل، میتوان گفت که سایر مدلهای مشابه همچون RAID 1, RAID 1+0 با آنکه هزینه هایی مشابه با مدل پیشنهادی را به سیستم تحمیل مینمایند، اما قادر به فراهم آوردن پهنای باند مناسبی بمنظور افزایش سطح دسترسی پذیری خصوصا بهنگام عملیات احیا، داده ای نبوده و یا بهنگام بروز خطا در سطح دنباله دیسکی، بخوبی نمی توانند صحت وامنیت داده ای را درون سیستم حفظ نمایند. مزیت عمده مدل پیشنهادی ما آنست که بدون هیچگونه هزینه افزوده و تنها بکمک الگوریتمهای خاص توزیع داده ای در سطح دنباله، موجب افزایش سطح اجراء عملیات داده ای موازی، بهبود تحمل پذیری نسبت به خطا و ایجاد پهنای باند بالاتر در سیستم میگردد. علاوه بر این، مدل RAID-RMS، الگوی توزیع داده ای خاصی را پیاده سازی می کند که قادر است، نرخ احیاء داده ای^{۱۶} را بگونه ای افزایش دهد که، بهنگام بروز هر خرابی در سطح دنباله دیسکی، احتمال از دست رفتن یک داده بمقدار قابل توجهی کاهش پیدا

¹⁶Data Recovery

کند . حفظ پهنای باند مؤثر در سطح سیستم بهنگام بروز خطا ، یکی دیگر از ویژگیهای بارز این مدل بشمار می رود که موجب افزایش سطح دسترس پذیری داده ای و در نتیجه تسریع انجام عملیات احیاء داده ای میگردد.

مراجع

- [1] P. M. Chen, E. K. Lee, G. A. Gibson, R. H. Katz and D. A. Patterson; "RAID: High-Performance, Reliable Secondary Storage", ACM Computing Surveys, Vol.26, No.2, June 1994, pp.145-185.
- [2] G. Gibson, D. Nagle, K. Amiri, F. Chang, H. Gobioff, E. Riedel, D. Rochberg and J. Zelenka, "A Cost-effective, High-bandwidth Storage Architecture", Proc. of the 8th Conf. on Architectural Support for Programming Languages and Operating Systems, 1998, pp.97-106.
- [3] J. Akbari and A. T. Haghghat . "A New Redundancy Algorithm For Distributed Environment", Operating system & Security Conference-OSSC 2003, Sharif university of technology, 24 -25 Dec. pp.70-81.
- [4] T. Cortes, "Software RAID and Parallel Filesystems", in High Performance Cluster Computing--Architectures and Systems, Rajkumar Buyya (ed.), Prentice Hall PTR, 1999, pp.463-496.
- [5] Michael Stonebraker ,Gerhard A. Schloss. "Distributed Raid – A New Multiple Copy Algorithm" , University of California, Berkeley, CA 94720, 2000
- [6] T. Anderson, M. Dahlin, D. Patterson, and R. Wang. "Serverless Network FileSystems", ACM Trans. on Computer Systems, Jan. 1996, pp.41-79.
- [7] S. Asami, N. Talagala, and D. A. Patterson, "Designing a self-maintaining storage system", Proceedings of 16th IEEE Symposium on Mass Storage Systems, March 1999, pp. 222-233.
- [8] L. F. Cabrera, and D. E. Long, "Using Distributed Disk Striping to Provide High I/O Data Rates", Proceedings of USENIX Computing Systems, Fall 1991, pp.405-433.
- [9] P. Cao, S. B. Lim, S. Venkataraman, and J. Wilkes, "The TickerTAIP Parallel RAID Architecture", ACM Trans. on Computer System, Vol.12, No.3, August 1994, pp.236-269.
- [10] P. F. Corbett, D. G. Feitelson, J.-P. Prost, and S. J. Baylor. "Parallel Access to Files in the Vesta File System". Proceedings of Supercomputing'93, 1993.
- [11] T. H. Cormen and D. Kotz, "Integrating Theory and Practice in Parallel File Systems", Proceedings of DAGS '93 Symposium, June 1993, pp. 64-74.
- [12] I. Foster, D. Kohr, Jr., R. Krishnaiyer, and J. Mogill, "Remote I/O: Fast Access to Distant Storage". Proc. of the Fifth Workshop on I/O in Parallel and Distributed Systems, November 1997, pp.14-25.
- [13] M. Harry, J. M. del Rosario, and A. Choudhary, "VIP-FS: a Virtual, Parallel File System for High Performance Parallel and Distributed Computing", Proceedings of the 9th International Parallel Processing Symposium (IPPS'95), April 1995, pp. 159-164.
- [14] R. S. Ho, K. Hwang, and H. Jin, "Design and Analysis of Clusters with Single I/O Space", Proceedings of 20th International Conference on Distributed Computing Systems (ICDCS 2000), April 2000, Taiwan, pp.120-127.
- [15] J. H. Howard, M. L. Kazar, S. G. Menees, D. A. Nichols, M. Satyanarayanan, R. N. Sidebotham, and M. J. West, "Scale and Performance in a Distributed File System". ACM Trans. on Computer System, Vol.6, No.1, pp.51-81, February 1988.
- [16] H. I. Hsiao and D. DeWitt, "Chained Declustering: A New Availability Strategy for Multiprocessor Database Machines", Proc. of 6 th International Data Engineering Conf., 1990, pp.456-465.
- [17] Y. Hu, Q. Yang and T. Nightingale. "RAPID-Cache --- A Reliable and Inexpensive Write Cache for Disk I/O Systems", Proceedings of the 5th International Symposium on High Performance Computer Architecture (HPCA-5), Orlando, Florida, Jan. 1999, pp. 204 – 213.
- [18] J. Huber, C. L. Elford, D. A. Reed, A. A. Chien, and D. S. Blumenthal, "PPFS: A High Performance Portable Parallel File System", Proceedings of the 9th ACM International Conference on Supercomputing, Barcelona, July 1995, pp.385-394.