

بررسی پایگاه داده NOSQL

آرزو اکبریان^۱، سیمین امینی^۲

^۱دانشجوی کارشناسی ارشد، دانشگاه آزاد اسلامی واحد شهرکرد، akbarian.arezoo@gmail.com

^۲دانشجوی کارشناسی ارشد، دانشگاه آزاد اسلامی واحد شهرکرد، s_amini144@yahoo.com

چکیده

امروزه با توسعه فناوری‌های مختلف و قابلیت نمونه‌برداری و تولید حجم عظیمی از داده‌ها، امکان ذخیره‌سازی و تحلیل آن‌ها چالشی بزرگ بهشمار می‌آید. ذخیره‌سازی، مدیریت و بازیابی داده‌های گذرا که در بعضی موارد در مقیاس بالایی در برنامه‌های کاربردی امروزی تولید شده‌اند نیز یکی دیگر از چالش‌هایی است که راه حل مدیریت مناسب آن‌ها را، پایگاه‌های داده‌ای NOSQL^۱/راهه کرده‌اند.

این مقاله ویژگیهای اصلی الگوی داده NOSQL را توصیف می‌کند و در مورد معماری جدیدی بنام CDSA^۲ بحث می‌کند که یک معماری پایگاه داده‌ای با حافظه توزیع شده برای محاسبات ابری است تا کارایی داده پرس-وچو و ذخیره حجم داده انبوه در ابر را با استفاده از استراتژی معقول بپسندد. علاوه بر این با اضافه کردن یا حذف هر Node از خوشه پایگاه داده توزیع شده، Node‌های دیگر میتوانند بدون توقف سرویس کار کنند. فضای ذخیره‌سازی را تداوم می‌بخشد، عملکرد را بالا می‌برد و تأخیر دستیابی را کم می‌کند.

کلمات کلیدی: memory database, CDSA, NOSQL

۱. مقدمه

آخرًا NOSQL^۱ به معنای "نه فقط SQL" به دسته بزرگی از پایگاه‌داده‌ها اطلاق می‌شود که خصوصیات پایگاه‌داده‌ای رابطه‌ای را ندارند و برای جستار زدن از زبان توصیفی SQL استفاده نمی‌کنند. از جمله بارزترین ویژگی‌های این دسته از پایگاه‌داده‌ها می‌توان به موارد زیر اشاره نمود:

✓ مدل داده غیر رابطه‌ای^۲:

محدودیت مدل رابطه‌ای در پشتیبانی از ابر داده‌ها^۳ و داده‌هایی با ساختارهای ترکیب شده ساخت یافته، نیمه ساخت یافته و غیر ساخت یافته یکی از دلایل اصلی معرفی NOSQL بود.

¹ NOT ONLY SQL

²Non-relational data model



✓ اهمیت کمتر به سازگاری داده‌ها:

پایگاه‌داده‌های NOSQL از مدل همووندی و تراکنشی ضعیف‌تری نسبت به ACID بهره می‌برند.

✓ شمای منعطفتر:

بیشماقی یا شمای ضعیف در انبار داده‌ها یک پیشرفت در جستارهای تحلیلی حرفه‌های فاقد عمومیت به حساب می‌آمد که پایگاه‌داده‌های NOSQL ای نیز بدان توجه وافری داشته‌اند.

✓ طراحی شده برای محیط‌های توزیع‌شده:

جهش‌های صورت گرفته در توسعه معماری کامپیوتر، پردازش‌های توزیع‌شده و موازی، محاسبات ابری و همچنین نیاز به تکرار و توزیع داده‌ها میان سرویس‌دهنده‌های متعدد؛ نیاز به یک پایگاه‌داده با قابلیت وسعت‌پذیری افقی را بیش از پیش روشن ساخته‌اند. یعنی پایگاهی که بتواند به سادگی و ارزانی با افزودن گره‌های جدید به شبکه‌اش توسعه یابد؛ برخلاف پایگاه‌داده‌های رابطه‌ای که تنها به وسعت پذیری عمودی یعنی ارتقا کارایی یک تک گره با افزودن به منابع آن یا با فناوری‌های مجازی‌سازی اهتمام ورزیده‌اند [1].

(Map Reduce) محیط محاسبات ابری گوگل را تشکیل میدهد، شرکتهای دیگر نیز سیستمهای ابری مشابه مثل EC2 آمازون و ابر آبی IBM را پیاده‌سازی کردند.

در صورت رویارویی با حجم عظیمی از داده، سیستم پایگاه داده‌ای موجود کارایی را از دست داده و نمی‌تواند از عهده هزینه‌های بالای نرمافزار پایگاه داده برآید و در آینده داده حجمی باید به ابر مهاجرت داده شود، اما ویژگیهای ACID که سیستم پایگاه داده نیاز دارد، وقتی داده توزیع شده را ذخیره می‌کند، منجر به کاهش کارایی خواهد شد. برای اطمینان از دسترسپذیری بالا، قابلیت اطمینان و صرفه‌جویی، انبار داده ابر باید افزونگی داشته باشد تا قابلیت اطمینان را تضمین کند. علاوه بر این سیستم محاسبات ابر باید در یک زمان با نیازهای تعداد زیادی از کاربران مواجه شود و به صورت موازی به کاربران سرویس دهد. بنابراین تکنولوژی انبار داده محاسبات ابری به عملکرد بالا و نرخ انتقال بالا نیاز دارد [2].

۲. پایگاه‌های داده NOSQL

راهکارهای NOSQL، برای مسائلی بسیار فراتر از دنیای سنتی پایگاه‌های داده‌ای به کار می‌روند و عملکردی به شدت بهتر از همتایان سنتی خود ارائه می‌کنند. لازم به تأکید است که گذار به سمت راهکارهای NOSQL، به دلیل مشکلات و محدودیت‌های زبان SQL نبوده است، بلکه به دلیل محدودیت‌های مدل رابطه‌ای پایگاه داده‌ای است.

راه حل‌های NOSQL در بسیاری از شرکت‌هایی که خدمات «وب اجتماعی» ارائه می‌کنند، به کار گرفته شده و به سرعت در حال گسترش است. این امر به دلیل سختی زیاد و محدودیت‌های سیستم‌های SQL کاملاً رابطه‌ای در برآورده کردن نیازهای داده‌ای آن‌ها است. با نگاهی به نیازمندی‌های مقیاس‌پذیری یکی از شبکه‌های اجتماعی به راحتی می‌توانیم به این امر واقف شویم. این نیازمندی‌ها عبارتند از:

³ big data



- ۵۷۰ میلیون مشاهده صفحات در ماه

- آپلود بیش از سه میلیارد عکس در ماه

- پردازش و ارائه بیش از ۱,۲ میلیون عکس در ثانیه

- ارائه ۲۵ میلیون نوع محتوا که با استفاده از ۳۰ هزار سرور انجام می‌پذیرد.

با این نیازمندی‌ها، که به یقین با نیازمندی‌های یک دیپارتمان حسابداری در دهه ۱۹۵۰ تفاوت‌های بسیاری دارد، این شبکه اجتماعی خود را با مجموعه‌ای غنی از ابزارها تطبیق‌داده است که هر کدام یکی از بهترین نمونه‌های پیشرو در حوزه پایگاه‌های داده‌ای Hadoop محسوب می‌شوند:

:Memcached ✓

این شبکه اجتماعی با استفاده از هزاران سرور Memcached، دهها تراپایت داده کش شده گذرا را در هر لحظه پردازش کرده و خدمات مرتبط را به کاربران خود ارائه می‌کند.

که هم اکنون با HBase جایگزین شده است: Cassandra ✓

با استفاده از این پایگاه‌های داده‌ای این شبکه اجتماعی عملیات ذخیره‌سازی گسترده طیف وسیعی از داده‌ها را بدون داشتن هیچ نقطه خط‌دادار یا مشکل‌داری در مجموعه عظیمی از ماشینهای محاسباتی، به بهترین نحو به اجرا در می‌آورد.

:Hive و Hadoop ✓

با استفاده از این ابزارهای پیشرفته، این شبکه اجتماعی تحلیل داده‌های عظیم و تحلیلهای بازاری و تبلیغاتی را با کارایی بالایی به انجام می‌رساند.

داده‌ای مانند داده‌های هوشنگی، فعالیت‌های آنلاین کاربران یا تحلیل‌های اقتصادی در قالب پایگاه‌های داده‌ای سنتی کارایی چندانی نخواهد داشت و در ذخیره‌سازی‌های بدون قالب و توزیع شده‌ای مانند هادوپ به بهترین روش کار خواهد کرد [3].

با توجه به موارد ذکر شده در بالا، می‌توان معماری داده‌ای جدید و کارا را کلید رشد و توسعه سریع این شبکه اجتماعی دانست که به عنوان دلیل اصلی مقیاسپذیری خوب آن نیز به شمار می‌آید. عاملی که زمینه رشد و توسعه شرکت‌های بزرگ دیگری مانند یاهو، Foursquare و Twitter را نیز به ارمغان آورده است. با این‌که این‌گونه شرکت‌ها در زمینه استفاده از این فناوری پیشگام هستند، اما هسته اصلی فناوری NOSQL به کار گرفته شده در بسیاری از کاربردهای موجود به صورت کلی در دسترس همگان قرار دارد که در بیشتر موارد به صورت Open Source توسعه داده شده‌اند. به همین دلیل، طیف وسیعی از توسعه‌دهندگان در برنامه‌های کاربردی و تجاری خود در حال آزمایش و تطبیق با این فناوری نوپا هستند و به زودی شاهد موج عظیمی از به‌کارگیری چنین فناوری‌هایی در گوش و کنار دنیای نرم - افزارها خواهیم بود [1].

این مقاله ویژگی‌های اصلی الگوی داده NOSQL را توصیف می‌کند و کارایی بالای معماری انبار داده برای محاسبات ابری ارائه میدهد تا با استفاده از استراتژی عقلانی، کارایی ذخیره داده توزیع شده را بهبود ببخشد [2].



۳. کارهای مرتبط

۱.۳. محاسبات ابری^۴

در یک سیستم^۵ محاسبات ابری، سرور میتواند سرور فیزیکی یا مجازی باشد. ابرهای محاسبه‌ی پیشرفته معمولاً شامل تعداد زیادی از منابع محاسباتی هستند، مثل شبکه‌های حوزه ذخیره‌سازی، تجهیزات شبکه، فایروال و سایر تجهیزات امن.

هر کاربر میتواند از طریق تجهیزات مناسب و یک مرورگر اینترنت استاندارد به برنامه‌های محاسبات ابری دسترسی پیدا کنند.

با وجود مفاهیم متفاوت از محاسبات ابری، اساساً سه ویژگی اصلی زیر را دارا هستند: پیدایش بر روی خوشه‌های سرور کم هزینه در مقیاس بزرگ، برنامه‌های توسعه داده شده به منظور همکاری با سرویسهای زیربنایی برای به حداقل رساندن استفاده از منابع، تکثیر بر روی چندین سرور کم هزینه برای تضمین دسترسی‌پذیری بالا.

قابلیت انبار داده در محاسبات ابری این است که: داده به صورت اتوماتیک روی Node‌های ذخیره‌سازی متفاوت در خوشه توزیع می‌شود، هر Node ذخیره‌سازی داده فقط بخشی از داده را نگه میدارد، در همان زمان کاربر میتواند داده را بر روی Node‌های متفاوت دیگری ذخیره کند تا مطمئن شود که یک نقطه خرابی باعث از دست داده نخواهد شد [2].

۲.۳. پایگاه داده ابر

با توسعه سرویس توزیع شده با مقیاس بزرگ و انبار توزیع شده که محاسبات ابری نیاز دارند، پایگاه داده رابطه‌ای سنتی با چالشهای زیادی روبرو است. پایگاه داده NOSQL نمونه قید پایگاه داده رابطه‌ای سنتی را از بین برد. از لحاظ ذخیره‌سازی، NOSQL یک پایگاه داده رابطه‌ای نیست، بلکه یک پایگاه داده hash با فرمت کلیدمقدار^۶ است. به خاطر متروک شدن زبان پرسچویی قوی SQL، سازگاری تراکنش و قیدهایی در نمونه پایگاه داده سنتی، پایگاه داده NOSQL چالشهایی که پایگاه داده رابطه‌ای سنتی با آن روبرو بود را حل کرد.

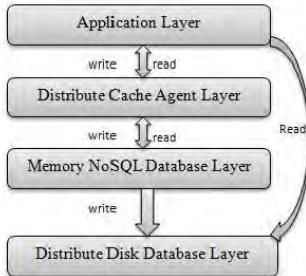
۴. معماری CDS

طبق عملکرد اجزای متفاوت، سیستم به سه لایه منطقی تقسیم می‌شود: لایه Memory Cache، لایه Data Cache و لایه Disk پایگاه داده [2].

⁴ Cloud Computing

⁵ platform

⁶Key-Value



شکل ۱: معماری CDS

^۷DCL . لایه ۱.۴

در حال حاضر روش گسترش متمنکز، برای بررسی مشکل از حافظه اصلی نهان^۸ استفاده می‌شود. گسترش توزیع شده پایگاه داده، وقتی که داده در هر یک از Node‌ها بروز رسانی شود منجر به تغییر حافظه اصلی نهان خواهد شد. ساختار حافظه اصلی نهان ذخیره شده باید مشابه ساختار جدول روی دیسک باشد و داده نمیتواند افزونگی داشته باشد. اگر چه استفاده منطقی از این روش میتواند کارایی را در پرسوچوهای پایگاه داده اثبات کند، اما بعضی از درخواستهای معمولی به موقع پاسخ داده نمی‌شود.

برای حل این مشکل از متد Optimal Prefetch استفاده می‌شود که یک متد پیشواکشی را ارائه میدهد، اما الگوی داده باید از قبل وارد شود. برای داشتن الگوی داده ناشناخته از پایگاه داده ابر، ایده‌های اصلی از متد OP چنین است: بروز رسانی حافظه اصلی نهان توزیع شده قبل از اولین جستار زدن^۹ کاربر، پیشینی جستار زدنهایی که ممکن است اجرا شود، جستار زدنها در زمان مناسب برای فعالسازی حافظه اصلی نهان و ذخیره کردن اطلاعات مربوط در هر حافظه اصلی نهان توزیع شده. این متد مشکل اولین بروخورد با حافظه اصلی نهان و کاهش زمان انتظار کاربر را در زمان اجرای اولین جستار زدن حل خواهد کرد، بنابراین حافظه اصلی نهان کارکرد بهتری میتواند داشته باشد.

OP یک متد برای ضبط تغییرات داده است و مبنی بر الگوریتم پیشینی است. OP در لایه حافظه اصلی نهان استفاده می‌شود تا سرعت پاسخ به جستار زدن را بالا ببرد. سابقه جستار زدن را ضبط می‌کند، اسکریپت جستار زدن مناسبی ایجاد می‌کند، جستار زدن را با مقادیر پارامترهای متفاوتی به عنوان یک الگو، بصورت آمار Sort شدهای از پارامترهای جستار زدن احتمالی قبلی از تمام جستار زدنهای خلاصه می‌کند. اشیاء OP الگویی با احتمالات بالا و پارامترهایش است [2].

^{۱۰}MDL . لایه ۲.۴

قبل از ماندگاری داده، داده ابر ابتدا در Memory ذخیره می‌شود، چرا که کارایی Memory بیشتر از حافظه Disk است.

⁷ Data Cache Layer

⁸ Cache

⁹ Query

¹⁰ Memory Database Layer



- با ترکیب نیازمندیهای فضای ذخیره‌سازی داده برای بالا بردن کارایی داده که قبلاً تحلیل شد و فضای ذخیره سازی این‌به در محیط محاسبات ابری، NOSQL به عنوان تنها راه حل به شمار می‌رود.

بنابراین این دو مورد قابل ملاحظه را می‌توان از یک ایده استنباط کرد: حافظه پایگاه داده Memory، NOSQL را بعنوان یک Disk استفاده می‌کند. عملیات خواندن و نوشتن مستقیماً با پایگاه داده Disk انجام نمی‌گیرد بلکه با پایگاه داده‌های Memory صورت می‌گیرد و باعث می‌شود از تأخیر زمانی و مشکلات آن که در NOSQL ساده وجود داشت جلوگیری شود.

- در سیستم پایگاه داده تراکنشی، در شروع کار پایگاه داده Memory برای سازگاری با پایگاه داده Disk نیاز می‌شود. بنابراین پایگاه داده Memory نیاز دارد که تعریف جدول یکسان با پایگاه داده Library را داشته باشد، در اولین شروع بکار لازم است همه داده‌ها از جدول پایگاه داده Disk به پایگاه داده Memory بارگذاری شود [2].

۱۱. DDL لایه

پایگاه داده Key-Value اساساً ویژگیهای از کارایی بالای خواندن و نوشتن هم زمان است. پایگاه داده حافظه Value همه پایگاه داده را در Memory لود می‌کند و داده را به طور منظم از پایگاه داده Memory با اعمال آسنکرون روی Disk مینویسد. علاوه بر این با تنظیم زمان انقضای کلید-مقدار می‌توان داده تغییر داده شده در پایگاه داده Memory را گرفت. از آنجا که عملیات کلأ در Memory اجرا می‌شود، کارایی پایگاه داده کلید-مقدار خیلی بالا است و بیش از میلیونها خواندن و نوشتن در هر ثانیه می‌تواند انجام شود.

- تغییرات داده در پایگاه داده حافظه باید در پایگاه داده Disk کپی شود. پس کپی داده از Disk به می‌تواند به عنوان فرایند نوشتن ناهمگام اصلی دیده شود، به وضوح نوشتن ناهمزمان سریعتر می‌شود [2].

۱۲. CDSA لایه ذخیره‌سازی توزیع شده در

پایگاه داده توزیع شده در محاسبات ابری یک سیستم سرویس پایگاه داده‌ای است که از پایگاه داده توزیع شده در Node‌های مختلف تشکیل یافته است و بر طبق این معماری توزیع شده مقیاسپذیری قابل انعطاف آنلاین فراهم می‌شود. به عنوان مثال، Node‌های داده بیشتری می‌توانند بدون توقف سرویس اضافه یا حذف شوند.

با جداسازی اعمال نوشتن و خواندن، تقسیم عمودی و افقی، داده این‌به می‌تواند درون خوشه پایگاه داده حافظه NOSQL ابر ذخیره شود و Proxy توزیع شده خوشه پایگاه داده مسئول کنترل‌ها مثل مسیریابی داده است.

وقتی داده از پایگاه داده Disk به پایگاه داده Memory کپی می‌شود، از متد پارتیشن‌بندی ترکیبی استفاده خواهد کرد. پارتیشنی که توسط یک گره NOSQL در لایه Memory پایگاه داده نهایی می‌شود توسط یک گره Memory و یک Node NOSQL نهایی می‌شود. طرح پایگاه داده از دو پایگاه داده شکل گرفته است [2].

۱۳. کارایی CDSA

^{۱۱} Disk Database Layer

^{۱۲} Cloud Data Storage Architecture



MDL از CDSA با تجمعیع حافظه‌های اصلی سرورها بوجود آمده است. CDSA میتواند حافظه‌های ذخیره‌سازی در دسترس و با دوام با توان عملیاتی صد برابر از استفاده تنها از پایگاه داده Disk و با تأخیر دسترسی کمتر را فراهم کند.

ترکیب تأخیر پایین و مقیاس بزرگ، نسل جدیدی از برنامه‌های داده‌ای قدرتمند را ایجاد خواهد کرد.

این بخش روی تست الگوی دسترسی CDSA تمرکز دارد. از بین عملیات اصلی insert, read, update, delete هم عملکرد و هم متوسط زمان پاسخ در یک زمان تست شد تا صحت و کارایی CDSA بررسی شود. محیط آزمایش به صورت زیر است:

برای تست کارایی و زمان پاسخ CDSA، ابتدا nmoll به عنوان ابزاری برای اندازه‌گیری کارایی برنامه آزمایشی و نوشتن در یک فایل بکار رفت و سپس nmoll به عنوان تحلیلگر برای آنالیز فایل داده استفاده شد.

در نهایت نتایج آزمایش CPU و I/O دیسک برای ایجاد یک گزارش گرافیکی انتخاب میشود. برنامه آزمایشی اعمال insert, read, update, delete پانصد میلیون رکورد را به ترتیب انجام داد و زمان مصرفی ثبت شد. اول برنامه آزمایشی پایگاه داده را باز کرد، رکوردها از یک پایگاه داده Memory به Disk به متناظر خوانده شد و سپس insert شد. اعمال read, update, delete نهایتاً داده را از Memory به فایل Disk همگامسازی کرد. نتایج مقالات بررسی شده نشان میدهد که میانگین زمان پاسخ در حد مایکرو ثانیه است و از آنجایی که از لحاظ منطق بسیار ساده است به وضوح CDSA انتخاب خوبی برای ابر میباشد.

۶. نتیجه

این مقاله در ابتدا پیش زمینه NOSQL را توصیف کرد، سپس با معرفی پایگاه داده Memory و NOSQL برای اهداف آدرسدهی محل ذخیره‌سازی داده حجمی و برای مسائل دستیابی همزمان بالا در محیط‌های محاسبات ابری، بر طبق ویژگیهای انبارهای ذخیره‌سازی توزیع شده در ابر، یک معماری از پایگاه داده حافظه توزیع شده NOSQL برای ابر پیشنهاد داده است. داده آماری نشان میدهد این معماری خیلی خوب انجام پذیرفته و همچنین فواید زیادی به همراه داشته است. البته محدودیتی در محدودیت منابع داده‌ای که توسط کاربر تعریف میشود وجود دارد، مثل انعطاف در کنترلهای سرویس ابر و جستجو در منابع داده‌ی مختلف. برنامه‌ریزیها برای دنبال کردن این مشکلات در آینده است.

مراجع

- [1] Han, Jing ,E, Haihon, Le , Guan, Du ,Jian " Survey on NoSQLDatabase",IEEE,978-1-4577-0208,363-366 ,2011.
- [2] Han ,Jing, Song,Meina, Song, Junde," A Novel Solution of Distributed Memory NoSQL Database for Cloud Computing", IEEE, 10.1109/ICIS.2011.61 ,351-355,2011.
- [3] Farmad,A."NO SQL". Fourth International Conference on Information Technology Management.