



وزارت علوم، تحقیقات و فناوری

مرکز منطقه ای اطلاع رسانی علوم و فناوری

گزارش نهایی طرح پژوهشی

بکارگیری سیستمهای NOSQL به عنوان ابزار پردازش کلان داده ها

در جستجو و بازیابی تمام متن

**Using NoSQL Systems as Big Data Processing Tools
for Full-Text Search and Retrieval**

مجری : دکتر بهاره پهلوان زاده

مهر ۱۳۹۸

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِیْمِ

چکیده

امروزه داده ها به عنوان یک دارائی ملی شناخته می شوند. همچنین پردازش، تحلیل و استفاده از داده به عنوان یک عامل کلیدی برای رشد کلیه سازمانها تبدیل شده است و موجب مزیت رقابتی در کسب و کارها، محرک نوآوری، افزایش رقابت و اثرات مثبت اجتماعی خواهد شد. مرکز منطقه ای اطلاع رسانی علوم و فناوری نیز از این قضیه مستثنی نمی باشد و با توجه به افزایش حجم داده ها بصورت تصاعدی (برای مثال در طی سال اخیر رشد ۵۰٪ داده های آن را بر بستر پایگاه داده رابطه ای Microsoft SQL Server شاهد بودیم) و همچنین با لحاظ نمودن ماموریت های این مرکز در ارائه خدمات از طریق پایگاههای داده مقالات تمام متن و موتور جستجوی اختصاصی آن، نیاز به ارائه راهکارهای نوین پردازشی یک ضرورت محسوب میشود.

در این پژوهش با توجه به ورود تکنولوژیهای نوینی همچون پایگاههای داده غیر رابطه ای و سیستمهای NOSQL در عصر داده های حجیم، ضمن کسب دانش کار با سیستمهای نوظهور NOSQL به بررسی بکارگیری و ارزیابی آنها در قیاس با پایگاه داده رابطه ای موجود بعنوان راهکاری احتمالی برای حل مشکلات آتی و پیش رو در اثر رشد تصاعدی داده های مبتنی بر سند مرکز منطقه ای پرداخته شده است.

برای ارزیابی مزایا و معایب پایگاههای داده غیر رابطه ای در قیاس با پایگاههای داده رابطه ای، معیار زمان بازیابی و معیار جدید کیفیت بازیابی تعریف گردید و پایگاه های داده رابطه ای Microsoft SQL و MariaDB و پایگاه های داده غیر رابطه ای MongoDB و Elasticsearch مقایسه گردیدند. نتایج پژوهش برتری چشمگیر Elasticsearch و Microsoft SQL Server نسبت به دو پایگاه داده دیگر از دیدگاه معیارهای زمان بازیابی و کیفیت بازیابی را نشان داد. همچنین مشاهده شد که با افزایش تعداد شاردها در Elastic Search؛ Elasticsearch برتری خود را نسبت به سایر پایگاههای داده افزایش می دهد.

کلید واژه ها: پردازش و تحلیل داده ها، کلان داده ها، پایگاه داده رابطه ای ، پایگاه داده غیر رابطه ای (سیستمهای

(NOSQL

فهرست مطالب

۱	فصل اول : مقدمه	۱
۱-۱	مقدمه	۱
۲-۱	ضرورت انجام طرح پژوهشی و اهداف کاربردی آن	۴
۳-۱	سئوالات پژوهش	۵
۴-۱	روش پژوهش	۶
۵-۱	هم راستایی طرح پژوهشی حاضر با اهداف سند راهبردی مرکز منطقه‌ای	۷
۶-۱	خروجی‌های طرح پژوهشی	۷
۷-۱	استفاده نتایج طرح پژوهشی در سازمان‌ها	۸
۲	فصل دوم : پیشینه پژوهش	۹
۱-۲	مقدمه	۹
۲-۲	پایگاه‌های داده رابطه‌ای	۱۰
۱-۲-۲	مدل رابطه‌ای	۱۰
۲-۲-۲	سیستم‌های مدیریت پایگاه‌های داده رابطه‌ای	۱۱
۳-۲	پایگاه‌های داده غیر رابطه‌ای	۱۲
۱-۳-۲	طبقه بندی پایگاه‌های داده غیر رابطه‌ای NOSQL	۱۳
۲-۳-۲	مروری بر مطالعات پیشین	۱۸
۳-۳-۲	پایگاه داده MongoDB	۲۲
۴-۳-۲	پایگاه داده Elasticsearch	۲۳
۵-۳-۲	مقایسه ویژگی‌های پایگاه‌های داده MongoDB و Elasticsearch	۲۴
۳	فصل سوم: روش شناسی پژوهش	۲۸
	(مدل‌سازی و پیاده‌سازی آزمایشگاهی)	۲۸
۱-۳	مقدمه	۲۸
۲-۳	مشخصات فنی سیستم و نسخه پایگاه‌های داده مورد استفاده	۲۹
۱-۲-۳	مشخصات فنی سیستم	۲۹
۲-۲-۳	نسخه پایگاه‌های داده مورد استفاده	۳۰
۳-۳	کلیات مدل‌سازی و پیاده‌سازی پایگاه‌های داده	۳۱

۳۵ فصل چهارم: یافته‌های پژوهش
۳۵ ۱-۴ مقدمه
 ۲-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای در حالت اجرای YCSB
۳۶ Benchmark
 ۱-۲-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای به ازای تعداد رشته های
۳۸ متفاوت
 ۲-۲-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای با افزایش تعداد
۴۲ عملیات/رکورد-نرخ افزایش زمان کل اجرای بارکاری (Workload)
 ۳-۲-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای با افزایش تعداد
۴۴ عملیات/رکورد-توانش زمانی و تاخیر زمانی
۵۳ ۴-۲-۴ جمع بندی نتایج مربوط به معیار ارزیابی سرویس ابری یاهو!(YCSB)
۵۴ ۳-۴ بررسی پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی تمام متن
 ۱-۳-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی تمام متن
۵۷ بدون خوشه بندی (Single Node) با معیار زمان بازیابی
 ۲-۳-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی تمام متن با
۶۰ قابلیت خوشه بندی (Sharded Cluster)
۶۶ ۳-۳-۴ جمع بندی نتایج جستجوی تمام متن
۶۷ فصل پنجم: بحث و نتیجه گیری
۶۷ ۱-۵ مقدمه
۷۰ ۲-۵ نتیجه گیری
۷۲ ۳-۵ پیشنهادهای آینده

فهرست جداول

- جدول ۱-۲ مقایسه چهار دسته پایگاه داده غیر رابطه ای بر حسب مدل داده.....۱۴
- جدول ۲-۲ مقایسه ویژگیهای کلی پایگاههای داده رابطه ای و غیر رابطه ای کاربردی پژوهش حاضر ۲۵
- جدول ۳-۲ مقایسه مزایا و معایب پایگاههای داده رابطه ای و غیر رابطه ای کاربردی پژوهش حاضر.. ۲۵
- جدول ۱-۳ مشخصات سیستمها ۲۹
- جدول ۲-۳ پایگاههای داده مورد استفاده..... ۳۰

فهرست شکل ها

- شکل ۱-۱ میزان رشد تحقیقات در زمینه پایگاه‌های داده غیررابطه‌ای در مقایسه با پایگاه‌های داده رابطه‌ای..... ۳
- شکل ۱-۳ معماری پیاده سازی شده اجزای تشکیل دهنده بر بستر MongoDB..... ۳۲
- شکل ۲-۳ معماری پیاده سازی شده اجزای تشکیل دهنده بر بستر ElasticSearch..... ۳۴
- شکل ۱-۴ مقایسه توانش زمانی بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای برای WorkloadA با تعداد رکورد/عملیات ۱۰۰۰۰ و تعداد رشته های مختلف 1,2,4,8..... ۴۰
- شکل ۲-۴ مقایسه تاخیر زمانی بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای برای WorkloadA با تعداد رکورد/عملیات ۱۰۰۰۰ و تعداد رشته های مختلف 1,2,4,8..... ۴۰
- شکل ۳-۴ مقایسه توانش زمانی بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای برای WorkloadB با تعداد رکورد/عملیات ۱۰۰۰۰ و تعداد رشته های مختلف 1,2,4,8..... ۴۱
- شکل ۴-۴ مقایسه تاخیر زمانی بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای برای WorkloadB با تعداد رکورد/عملیات ۱۰۰۰۰ و تعداد رشته های مختلف 1,2,4,8..... ۴۱
- شکل ۵-۴ زمان کل اجرای WorkloadA بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای با تعداد رکورد/عملیات ۱۰۰۰۰، ۱۰۰۰ و ۱۰۰۰۰۰ (الف)..... ۴۳
- شکل ۶-۴ زمان کل اجرای WorkloadA بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای با تعداد رکورد/عملیات ۱۰۰۰، ۱۰۰۰۰ و ۱۰۰۰۰۰ (ب)..... ۴۴
- شکل ۷-۴ زمان کل اجرای WorkloadB بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای با تعداد رکورد/عملیات ۱۰۰۰، ۱۰۰۰ و ۱۰۰۰۰۰..... ۴۵
- شکل ۸-۴ مقایسه توانش زمانی بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای برای WorkloadA با تعداد رکورد/عملیات ۱۰۰۰..... ۴۷
- شکل ۹-۴ مقایسه تاخیر زمانی بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای برای WorkloadA با تعداد رکورد/عملیات ۱۰۰۰..... ۴۷
- شکل ۱۰-۴ مقایسه توانش زمانی بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای برای WorkloadB با تعداد رکورد/عملیات ۱۰۰۰..... ۴۸

شکل ۴-۱۱ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB با تعداد رکورد/عملیات

۱۰۰۰.....۴۹

شکل ۴-۱۲ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadA با تعداد رکورد/عملیات

۱۰۰۰۰.....۴۹

شکل ۴-۱۳ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadA با تعداد رکورد/عملیات

۱۰۰۰۰.....۵۰

شکل ۴-۱۴ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB با تعداد رکورد/عملیات

۱۰۰۰۰.....۵۰

شکل ۴-۱۵ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB با تعداد رکورد/عملیات

۱۰۰۰۰.....۵۲

شکل ۴-۱۶ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadA با تعداد رکورد/عملیات

۱۰۰۰۰۰.....۵۲

شکل ۴-۱۷ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadA با تعداد رکورد/عملیات

۱۰۰۰۰۰.....۵۳

شکل ۴-۱۸ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB با تعداد رکورد/عملیات

۱۰۰۰۰۰.....۵۳

شکل ۴-۱۹ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB با تعداد رکورد/عملیات

۱۰۰۰۰۰.....۵۸

شکل ۴-۲۰ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای بدون خوشه بندی با

چهار نوع پرس و جوهای آزمایشی در مقیاس لگاریتمی.....۵۸

شکل ۴-۲۱ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای بدون خوشه بندی با چهار نوع کوئیری.....۵۹

شکل ۴-۲۲ مقایسه تاخیر زمانی پایگاه داده غیر رابطه ای MongoDB با در نظر گرفتن شاردینگ با چهار نوع کوئیری.....۶۱

شکل ۴-۲۳ مقایسه تاخیر زمانی پایگاه داده غیر رابطه ای ElasticSearch با در نظر گرفتن شاردینگ با چهار نوع

کوئیری.....۶۱

شکل ۴-۲۴ مقایسه زمان بازیابی پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی تمام متن با
چهار نوع پرس و جویهای آزمایشی به تفکیک

معماری.....۶۳

شکل ۴-۲۵ مقایسه کیفیت بازیابی داده ها در پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی

تمام متن۶۶

فصل اول : مقدمه

۱-۱ مقدمه

همانطور که می دانیم در عصر داده های حجیم؛ پردازش، تحلیل و استفاده از داده به عنوان یک عامل کلیدی برای رشد کلیه سازمانها تبدیل شده است و موجب مزیت رقابتی در کسب و کارها، محرک نوآوری، افزایش رقابت و اثرات مثبت اجتماعی خواهد شد. بر اساس گزارشهای سالانه سایت^۱ IEEE بعنوان بزرگترین انجمن تخصصی دنیا در پیشرفت تکنولوژی و استاندارد ها برخی از رایجترین و محبوبترین واژه های جستجو شده شامل موارد زیر است:

- کلان داده ها^۲
- اینترنت اشیا^۳
- امنیت سایبری^۴
- رایانش ابری^۵
- نسل جدید شبکه های بی سیم^۶
- شبکه های مشبک توری هوشمند^۷

۱ Institute of Electrical and Electronics Engineers

۲ Big Data

۳ Internet of things

۴ Cyber Security

۵ Cloud Computing

۶ Next Generation Wireless

۷ Smart Grid

از طرفی با توجه به چشم اندازهای آینده و رو به رشد مرکز منطقه ای و متناسب با سند راهبردی و همچنین اساسنامه آن استفاده از شیوه های نوین پردازش برای ارائه برنامه ها و خدمات اطلاع رسانی در منطقه تاکید شده است. با توجه به افزایش روزافزون حجم بزرگ داده ها، گسترش کیفی و کمی سرویسهای متعدد ارائه شده به سایر پژوهشگران در سطح ملی و بعضا بین المللی منطقه؛ حرکت به سمت ارتقا زیرساختهای نرم افزاری در قالب محاسبات توزیع شده و استفاده از روشهای نوین پردازش داده ها اهمیت می یابد.

در این پژوهش هدف اصلی کار بر نمونه چکیده های لاتین از پایگاههای اطلاعاتی مرکز به عنوان داده های مورد پژوهش و بررسی می باشد. ضمن اینکه از داده های در حال حاضر حدود ۶ میلیون رکورد چکیده های لاتین در دست است (که در سال اخیر رشد چشمگیر بالغ بر ۵۰٪ی آن را شاهد بودیم)، که با توجه به خصوصیت جستجوی تمام متنی^۱، هر رکورد به صورت تقریبی حدود ۳۰۰ کلمه است و لذا داده های مورد بررسی بالغ بر یک میلیارد و هشتصد میلیون می باشد که به نوبه خود داده حجیم محسوب می شود.

بنابراین انجام این طرح باعث کسب دانش کار با سیستمها و پایگاههای داده غیررابطه ای NoSQL^۲ و حرکت به سمت خوشه بندی داده ها در مرکز در سالهای آتی می شود. از طرفی استفاده از این دانش نه تنها در حجم بالای داده باعث بهبود سرعت بازیابی می شود بلکه برای ذخیره سازی داده های حجیم در آینده نزدیک برای مرکز ضروری است، چرا که این روزها با وجود چنین تکنولوژیهای هزینه های مقیاس پذیری عمودی^۳ غیر قابل توجیه و بعضا تحمل ناپذیر بوده و باید با

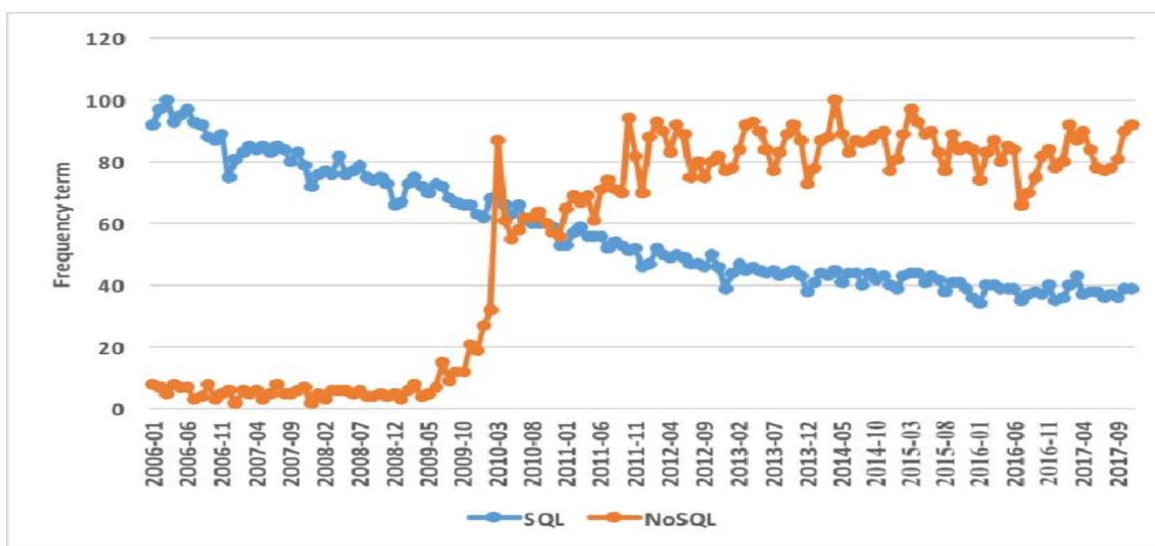
¹ Full-Text Search

² Not Only SQL (Structured Query Language)

³ Vertical Scalability

استفاده از تکنیک های جدید کلاسترینگ در پایگاههای داده غیررابطه ای یا سیستمهای NOSQL، پایگاه داده ها به سمت معماریها و سیاستهای مقیاس پذیری افقی^۱ بروند.

پایگاه های داده غیر رابطه ای NOSQL برای کار با داده های حجیم به صورت توزیع شده طراحی شده اند. می توان با بکارگیری تعداد زیادی گره^۲ و توزیع داده ها و محاسبات بین آنها از ویژگی مقیاس و گسترش پذیری و توسعه منابع در آنها بهره برد. باید در نظر داشت که توزیع خودکار داده و میزان تحمل خطای بالا از ویژگی های بارز پایگاه های داده غیر رابطه ای است. در شکل ۱-۱ نمودار گوگل ترند^۳ رشد استفاده و تحقیقات در زمینه پایگاه های داده غیر رابطه ای NOSQL در قیاس با پایگاه های داده ای رابطه ای را نشان می دهد.



شکل ۱-۱ میزان رشد تحقیقات در زمینه پایگاه های داده غیررابطه ای درمقایسه با پایگاههای داده رابطه ای (Google Trends, 2019)

¹ Horizontal Scalability

² Node

³ Google Trends (<https://trends.google.com/>)

علاوه بر دلایل مطرح شده در بالا و با توجه به اینکه تا به امروز امکان جستجو و بازیابی تمام متن در نرم افزارهای استفاده شده در مرکز منطقه ای به صورت کارا وجود نداشته است، لذا در طرح فعلی برای بهبود کیفیت نتایج جستجو و بازیابی تمام متن به بررسی پایگاه داده های غیر رابطه ای یا همان سیستمهای NOSQL به عنوان روشی نوین برای پردازش داده ها پرداخته می شود. با بررسی انجام شده، بیش از ۲۲۵ نوع سیستم NOSQL وجود دارد که باید با در نظر گرفتن پارامترهایی مانند ساختار ذخیره سازی داده، سرعت خواندن/نوشتن و مکانیزمهای بهینه سازی آن ها متناسب با نوع داده ها و نیازهای مرکز منطقه ای در این طرح پژوهشی محبوب ترینها را منتخب و در قیاس با پایگاه داده رابطه ای موجود در سازمان پرداخت.

۱-۲ ضرورت انجام طرح پژوهشی و اهداف کاربردی آن

نیاز عمده نرم افزارهای نسل جدید کتابخانه ای دیجیتال مبتنی بر اینترنت اشیاء جهت ارائه سرویسهای هوشمندانه و کاربرپسند در کتابخانه های مدرن ، حجم بالای داده ها و رشد سریع داده های تمام متن مرکز منطقه ای در سالیان اخیر، همچنین ضرورت استفاده از تکنیکهای خوشه بندی و سازگاری با ذخیره سازی داده ای غیر ساخت یافته (در پایگاههای غیر رابطه ای) با لحاظ نمودن سرعت پردازشی بهتر به عنوان راهکارهای مقیاس پذیر در راستای مدیریت بهینه منابع سخت افزاری از دیگر مواردی است که انجام پژوهش و تحلیل بر داده های موجود در مرکز منطقه ای را با استفاده از روشهای پردازش نوین ضرورت می بخشد. از اینرو است که در برنامه راهبردی نیز به ضرورت استفاده از پردازشهای نوین بر داده های رو به رشد مرکز تاکید به عمل آمده است؛ لذا این پژوهش شروعی در راستای یکی از اهداف برنامه راهبردی مرکز منطقه ای می باشد.

با توجه به وجود چالشهای متفاوت در حوزه کلان داده ها در سطح جهان و ایران به عنوان

کشوری در حال توسعه و عدم وجود بلوغ کافی در حوزه کلان داده در اکوسیستم جهانی و ازطرفی با توجه به نقشه راه توسعه کلان داده کشور (تدوین شده در دی ماه ۱۳۹۶) در جهت ارائه مدل‌های جدید برای حل مسائل داخلی و آسیب شناسی توسعه فناوری در کشور از طریق شناخت تجربیات و الگوهای موفق و نا موفق به منظور بهره برداری در حوزه کلان داده ها ؛ و نیز با هدف ایجاد ساز و کاری جهت اندازه گیری میزان ارزش استخراج شده از داده های موجود در کشور با روشهای سنتی در مقابل روشهای نوین پیشنهادی؛ بر آن شدیم که با توجه به رشد نسبتا زیاد داده های موجود در پایگاههای داده مرکز منطقه ای، ضمن بررسی بنیادین راهکارهای نوین پردازشی در این حوزه به پیاده سازی و مقایسه مدل‌های انتخابی پرداخته و در صورت کسب نتایج مثبت در راستای سیاستهای آتی مرکز جهت پردازش داده های موجود در پایگاه داده های موجود استفاده نماییم.

مسئله سرمایه گذاری اولیه در این زمینه می تواند گام موثری در راستای شناخت نیازهای واقعی و ایجاد زمینه های پژوهشی و برنامه ریزی برای تولید محصولات در راستای ارائه خدمات بهتر باشد.

۱-۳ سؤالات پژوهش

در انجام این طرح سؤالات پژوهش زیر مطرح می‌باشد:

- آیا با توجه به حجم داده ها و رشد آتی آنها پایگاههای غیر رابطه ای/سیستمهای NOSQL انتخاب مناسبی در راستای برآورده کردن نیازهای مرکز منطقه ای می باشد؟
- کدامیک از پایگاههای غیر رابطه ای/سیستمهای NOSQL بهترین تناسب را با اهداف کاربردی مطرح شده در این طرح دارد؟

۱-۴ روش پژوهش

بصورت اجمالی، مراحل پژوهش و نقشه راه طرح مزبور متناسب با هدفهای ذکر شده و فرضیه های موجود در جهت دانش افزایی مضاعف در این حوزه با توجه به گسترش روز افزون داده ها (برای برنامه های توسعه زیرساختی آتی سازمانی) بشرح زیر است:

- بررسی مدل‌های پایگاه‌های داده غیر رابطه ای / سیستم‌های NoSQL و انتخاب دو مدل برگزیده

- مدل‌سازی سخت‌افزاری مورد نیاز در حالت آزمایشگاهی

- نصب و راه اندازی بستر سخت افزاری جهت راه‌اندازی معماریهای مختلف تک نود و خوشه ای و انجام تنظیمات مورد نیاز

- نصب و راه اندازی بستر نرم افزاری و انجام تنظیمات پایگاه‌های داده غیر رابطه ای/سیستم‌های NoSQL

- بررسی و اجرای داده های معیار^۱ بر روی هر یک از پایگاه‌های داده غیر رابطه ای/سیستم‌های NoSQL و تحلیل و بررسی نتایج

- نمونه‌گیری از پایگاه داده‌های متنی مرکز منطقه‌ای

- درج اطلاعات در مدل‌های پایگاه‌های داده غیر رابطه ای انتخاب شده

- اجرا پرس و جوهای آزمایشی^۲ جهت بازیابی اطلاعات در محیط توزیع شده در پایگاه‌های داده غیر رابطه ای/سیستم‌های NoSQL انتخاب شده

- درج اطلاعات در مدل پایگاه‌های داده رابطه‌ای

¹ Benchmark Data

² Query

- اجرا پرس و جوهای آزمایشی^۱ جهت بازیابی اطلاعات در مدل پایگاه‌های داده رابطه-ای
- بررسی نتایج حاصل شده و نتیجه‌گیری

۵-۱ هم راستایی طرح پژوهشی حاضر با اهداف سند راهبردی مرکز منطقه‌ای

- به طور کلی : متناسب با اساسنامه و سند راهبردی در راستای هدف استفاده از تکنولوژی‌ها و شیوه‌های نوین پردازش برای ارائه برنامه‌ها و خدمات اطلاع رسانی در منطقه
- به طور خاص : متناسب با اهداف برنامه‌های شماره ۳ و ۴ و ۵ جدول ۱ "برنامه یک ساله و پنج ساله نظام جامع رایانه‌ای مرکز" درج شده در فصل هفتم برنامه راهبردی (سند راهبردی مرکز منطقه‌ای اطلاع رسانی علوم و فناوری، ۱۳۹۵)

۶-۱ خروجی‌های طرح پژوهشی

خروجی این طرح پژوهشی متناسب با اهداف و فرضیه‌های مطرح شده در قسمت ۳-۱ می‌باشد. و به عبارت دیگر بطور کلی هدف پژوهش و بررسی سیستمهای NoSQL در پایگاههای داده متنی در مقایسه با پایگاههای داده رابطه‌ای Microsoft SQL Server و Maria DB و نهایتا مقایسه و انتخاب مدل (های) برگزیده مناسب در این حوزه پژوهشی با دیدگاه زمان بازیابی و کیفیت بازیابی برای داده‌های مرکز منطقه‌ای اطلاع رسانی علوم و فناوری می‌باشد. کما اینکه پژوهش در این زمینه منجر به ایجاد دیدگاه تخصصی تر و دانش افزایی دوچندان در حوزه نوین پردازش داده‌های حجیم سازمانی می‌گردد و در صورت استفاده از نتایج آن در راستای سیاستهای آتی در خصوص زیرساختهای نرم

¹ Query

افزایی پایگاههای داده مرکز به عنوان CoreBusiness سازمان و در نتیجه رقابت پذیری با رقبای موثر خواهد بود.

۷-۱ استفاده نتایج طرح پژوهشی در سازمانها

هدف اولیه و اصلی مرکز منطقه ای (بصورت نمونه و مطالعه موردی)^۱ است و جامعه آکادمیک نیز از نتایج این بررسی با توجه به نوین و عملیاتی بودن کار میتوانند استفاده ببرند. همانطور که در ضرورت پژوهش عنوان شد، علاوه بر مرکز منطقه ای، با توجه به نوین بودن پژوهش طرح حاضر هر سازمان دیگری همچون آن که دارای حجم زیادی پایگاه داده متنی می باشد در سطح ملی و بین المللی از نتایج این طرح استفاده نماید. از آن جمله می توان به ISC یا ایران داک با ماموریت های مشابه با مرکز منطقه ای اشاره نمود.

¹ Case Study

فصل دوم : پیشینه پژوهش

۱-۲ مقدمه

پایگاه‌های داده رابطه ای از سال ۱۹۷۰ تا به امروز در حال توسعه اند (Chen, J. K. et al., 2019). استفاده و مدیریت این پایگاه‌های داده با استفاده از سیستم های مدیریت پایگاه داده رابطه ای^۱ به سادگی امکان پذیر است و همین امر منجر به محبوبیت این پایگاه‌های داده شده است. با این وجود، با توجه به استفاده روزافزون از فناوری هایی که با تولید اطلاعات در ابعاد و انواع مختلف همراه هستند، سازمانها نیاز به ذخیره سازی حجم بالاتری از اطلاعات نسبت به گذشته دارند. همچنین نیاز است که بازیابی این اطلاعات در سریع ترین زمان ممکن انجام شود.

پردازش این حجم بالای اطلاعات نیازمند سرعت کافی، ساختار انعطاف پذیر و پایگاه‌های داده توزیع شده می‌باشد و پایگاه‌های داده غیر رابطه ای یا همان سیستمهای NOSQL برای پاسخگویی به این نیازها طراحی شده اند.

در نقطه مقابل پایگاه‌های داده رابطه ای، برای بازیابی اطلاعات پیچیده در بسیاری از مواقع نیاز به انجام SQL Join بین تعداد دو یا بیشتر جدول دارد که ممکن است باعث افت عملکرد عملیات شود. همچنین عدم سازگاری این نوع پایگاه‌های داده با انبوه اطلاعات بدون ساختار مشخص و نداشتن قابلیت توزیع پذیری از دیگر نقاط منفی آن ها نسبت به پایگاه‌های داده غیر رابطه ای

¹ RDBMS (Relational Database Management System)

می‌باشد (Chen, J. K. et al., 2019; Li, Y, et al., 2013). در ادامه به توضیح مفصل تری از مفاهیم ارائه شده در این قسمت می‌پردازیم.

۲-۲ پایگاه‌های داده رابطه‌ای

پایگاه داده رابطه‌ای نوعی از پایگاه داده است که داده‌های مرتبط به هم را ذخیره کرده و دسترسی به آن‌ها را فراهم می‌آورد. به عبارتی این نوع پایگاه‌های داده مبتنی بر مدل رابطه‌ای^۱ می‌باشند که یک راه شهودی و قابل درک برای نشان دادن اطلاعات در جداول است. در این نوع پایگاه داده هر سطر یک جدول، یک رکورد با یک کلید شناسایی منحصر به فرد است و هر ستون نمایانگر یک ویژگی در مورد رکورد های آن جدول است (Oracle-RDBMS, 2017).

۲-۲-۱ مدل رابطه‌ای

در سال‌های اولیه توسعه پایگاه‌های داده، هر برنامه کاربردی^۲ داده‌های خود را در ساختار منحصر فرد خود ذخیره می‌کرد و همین امر باعث می‌شد توسعه این برنامه‌ها و برنامه‌های دیگر با استفاده از این داده‌ها، نیازمند مطلع شدن از جزئیات آن ساختار داده خاص باشد. اینگونه ساختارهای داده ناکارآمد بودند و بهینه‌سازی و ارتقای آن‌ها دشوار بود. برای حل این مشکل مدل پایگاه داده رابطه‌ای طراحی و به عنوان جایگزینی برای ساختارهای داده دلخواه پیشین مطرح شد. مدل رابطه‌ای روشی استاندارد برای نمایش داده‌ها و جستجو در آن‌ها فراهم می‌آورد که برای هر برنامه‌ای قابل استفاده می‌باشد. از همان ابتدا، توسعه دهندگان دریافتند که قدرت اصلی این

¹ Relational

² Application

مدل، در استفاده از مفهوم جدول است که روشی بصری، کارآمد و تا حدودی انعطاف پذیر برای ذخیره و دسترسی به اطلاعات ساختاریافته است.

با گذشت زمان، قدرت دیگری از مدل رابطه‌ای پدیدار شد و آن استفاده از زبان جستجوی ساختاریافته SQL برای ارتباط با پایگاه‌های داده رابطه‌ای بود. این زبان بر اساس جبر رابطه‌ای، یک زبان تقریباً ریاضی در اختیار کاربران قرار می‌دهد و به بهبود عملکرد پرس و جوهای آزمایشی بر روی پایگاه‌داده کمک شایانی می‌کند.

۲-۲-۲ سیستم‌های مدیریت پایگاه‌های داده رابطه‌ای

یک سیستم مدیریت پایگاه داده رابطه‌ای، مجموعه‌ای از برنامه‌ها و قابلیت‌هایی است که تیم‌های IT و سایر کاربران را قادر می‌سازد به ایجاد یک پایگاه‌داده رابطه‌ای، مدیریت و ارتباط با آن پردازند. بیشتر سیستم‌های مدیریت پایگاه‌های داده رابطه‌ای های تجاری برای دسترسی به پایگاه‌داده زبان SQL^۱ استفاده می‌کنند (Gudivada, V. N. et. al, 2014) از محبوب‌ترین این سیستم‌ها، سیستم‌های مدیریت پایگاه داده Microsoft SQL Server و MySQL می‌باشند که هر دو به زبان‌های C و C++ توسعه داده شده‌اند و اولین نسخه آن‌ها به ترتیب در سال‌های ۱۹۸۸ و ۱۹۹۵ عرضه شد (DB-Engines Ranking, 2017). پس از انتقال مالکیت کامل پروژه MySQL از شرکت MySQL به شرکت Oracle در بین سال‌های ۲۰۰۸ و ۲۰۰۹، توسعه دهندگان اصلی MySQL، پایگاه داده

¹ Structured Language Query

MariaDB را به عنوان یک شاخه و انشعاب^۱ از MySQL معرفی کردند و تا به امروز این پایگاه داده نوپا در بسیاری از شرکت ها جایگزین MySQL شده است (Kenler, E., 2015).

۲-۳ پایگاه‌های داده غیر رابطه‌ای

از اینکتومی^۲ به عنوان اولین موتور جستجوی واقعی تا گوگل رهبر فعلی موتورهای جستجو، پژوهشگران و توسعه دهندگان به خوبی محدودیت‌های پایگاه‌های داده رابطه‌ای را دریافتند. از جمله این محدودیت‌ها که بیشتر در مقیاس بزرگ اطلاعات خود را نشان می دهند، مقیاس پذیری، موازی سازی عملیات و هزینه است (Vaish, G. 2013). پژوهش‌هایی از این دست و موفقیت اولین پایگاه‌های داده غیررابطه‌ای جدول بزرگ گوگل^۳ و پایگاه داده داینامو آمازون^۴ (Amazon DynamoDB, 2017) و ، آغازگر توسعه پایگاه‌های داده غیر رابطه‌ای متعددی بود که بصورت متن باز عرضه شدند (Li, Y. et al., 2013). این پایگاه‌های داده دارای ویژگی‌های اصلی زیر می باشند که باعث افزایش روزافزون محبوبیت آن‌ها شده است (Vaish, G. 2013).

- ذخیره و بازیابی داده‌های بدون ساختار در قالب فرمت‌های مختلف
- افزایش سرعت توسعه برنامه‌های مختلف چرا که برخلاف پایگاه‌داده‌های رابطه‌ای دسترسی به داده‌های مرتبط نیازی به استفاده از پرس و جوهای آزمایشی پیچیده‌ی SQL ندارد.
- افزایش سرعت دسترسی به داده‌های مختلف

¹ Fork

² Inktomi

³ Google's BigTable

⁴ Amazon's Dynamo

- قابلیت مقیاس پذیری افقی به معنای اینکه بتوانیم به منظور بهبود عملکرد، بجای اینکه سخت افزار یک سیستم واحد را ارتقاء دهیم (مقیاس پذیری عمودی)، از سیستم های متعددی در قالب یک سیستم توزیع شده بهره مند شویم.

۲-۳-۱ طبقه بندی پایگاه های داده غیر رابطه ای NOSQL

پایگاه های داده غیر رابطه ای NOSQL بر اساس ساختار ذخیره سازی داده ها بصورت زیر به چهار دسته کلی قابل طبقه بندی می باشد (Han, J., et al., 2011; Abramova, V. et al., 2014; Patil, M. M., 2017) که در ادامه به شرح هریک از آنها بصورت اجمالی خواهیم پرداخت. جدول ۱-۲ نیز به مقایسه اجمالی ۴ دسته پایگاه داده غیر رابطه ای می پردازد.

- پایگاه داده کلید-مقدار^۱
- پایگاه داده مبتنی بر ستون^۲
- پایگاه داده سند گرا^۳
- پایگاه داده مبتنی بر گراف^۴

¹ Key-Value Database

² Column oriented

³ Document Based

⁴ Graph Based

جدول ۲-۱ مقایسه چهار دسته پایگاه داده غیر رابطه ای بر حسب مدل داده

مدل داده	کار آبی	مقیاس پذیری	انعطاف پذیری	پیچیدگی	عملکرد
کلید-مقدار	بالا	بالا	بالا	-	متغیر
مبتنی بر ستون	بالا	بالا	متوسط	پایین	حداقل
سند گرا	بالا	متغیر	بالا	پایین	متغیر
مبتنی بر گراف	متغیر	متغیر	بالا	بالا	متغیر

۲-۳-۱-۱ پایگاه داده کلید-مقدار

پایگاه داده کلید-مقدار ساده ترین مدل پایگاه داده غیر رابطه ای از لحاظ پیاده سازی می باشد. این مدل از یک جدول هاش^۱ استفاده می کند و به هر رکورد یک کلید منحصر به فرد و یک اشاره گر^۲ به داده اصلی اختصاص می یابد. به عبارتی دیگر، این پایگاه داده ها را می توان به عنوان یک جدول هاش توزیع شده^۳ یا دیکشنری بزرگ در نظر گرفت که هر داده (مقدار) با یک کلید منحصر به فرد آدرس دهی می شود. داده ها از یکدیگر مستقل بوده و ارتباطات بین آن ها در سطح برنامه کاربردی مدیریت می شوند. درج داده در آن ها به راحتی انجام می شود. هیچ رابطه ای در این

¹ Hash Table

² Pointer

³ Distributed Hash table

ساختار وجود ندارد و برای ذخیره سازی سریع داده‌های پایه استفاده می‌شود (مانند ذخیره سازی نام کاربری و رمز عبور، اطلاعات سژن/نشست^۱، سویچینگ^۲ و . . .). گروه‌بندی کلیدها و درج مقادیر با ساختارهای تودرتو از جمله روش‌های ذخیره سازی داده‌های پیچیده در این نوع پایگاه داده می‌باشد (R. Hecht, S. Jablonski, 2011)

مهمترین ویژگی عملیاتی این نوع از پایگاه داده‌ها سرعت پاسخ‌دهی برخط^۳ و مدیریت حجم بالایی از داده با قابلیت مقیاس‌پذیری خوب بین گره‌های یک خوشه می‌باشد. شعار بانک اطلاعاتی کلید-مقدار آئرواسپاک^۴ در این گروه قابل تامل است که ادعا می‌کنند ۹۹ درصد از درخواست‌ها را در کمتر از ۱ میلی‌ثانیه پاسخ می‌دهد. مثال‌های این نوع پایگاه‌داده غیر رابطه‌ای پایگاه داده غیررابطه ای اوراکل^۵ و پایگاه داده ساده آمازون^۶ (Amazon SimpleDB, 2017) و ردیس^۷ (Macedo, T., & Oliveira, F. (2011). و لدمورت^۸ (Voldemort Project, 2017)، ریاک^۹ (Riak, 2017)، پایگاه داده آمازون داینامو (DynamoDB, 2017) می‌باشند.

۲-۱-۳-۲ پایگاه داده مبتنی بر ستون

پایگاه داده مبتنی بر ستون همانطور که از نامش پیداست، داده‌ها را در قالب ستون، به‌جای سطر، ذخیره می‌کند. این مدل به منظور پردازش حجم بالای اطلاعات توزیع شده روی سرورهای متعدد طراحی شده است. در این مدل ذخیره داده روی یک نقشه مرتب شده توزیع شده^{۱۰} و

¹ Session information

² Switching

³ Real-Time

⁴ AeroSpike

⁵ Oracle NoSQL Database

⁶ Amazon simpleDB

⁷ Redis

⁸ Voldemort

⁹ RIAK

¹⁰ Distributed Sorted-map

چندبعدی صورت می گیرد. از نظر مفهومی، یک پایگاه داده مبتنی بر ستون مانند یک سیستم های مدیریت پایگاه های داده رابطه ای با شاخص روی هر ستون بدون سربار مرتبط است. این پایگاه ها، برای اولین بار با ارائه جدول بزرگ گوگل (Chang, F., 2008) و معرفی "سیستم ذخیره سازی توزیع شده برای مقیاس پذیری داده های غیر ساخت یافته در اندازه بسیار بزرگ" آغاز به کار کردند و بعد از آن نسخه های مختلفی از پایگاه داده های مبتنی بر ستون ایجاد شدند. ستون ها در دسته های بزرگتری به نام "خانواده ستون ها"^۱ دسته بندی می شوند. انعطاف پذیری بالایی در افزودن سطرها و ستون ها وجود دارد اما خانواده ستون ها باید از پیش تعریف شده باشد. زمانی که یک مقدار تغییر می یابد، مقدار جدید به عنوان نسخه متفاوتی از مقدار قبلی با یک مهر زمان ذخیره می گردد. به بیان دیگر مفهوم به روزرسانی به طور مؤثر وجود ندارد. پایگاه داده مبتنی بر ستون، عملیات تجمیع مانند محاسبه ماکزیمم، مینیمم، میانگین و جمع را روی پایگاه داده وسیع با بهترین راندمان انجام می دهند. پایگاه های داده غیر رابطه ای اچ-بیس^۲ (HBase, 2017)، جدول هایپر^۳ (Hypertable, 2017) و کاساندر^۴ (Carpenter, J., and Hewitt, E., 2020) مثال هایی از این مدل هستند. پایگاه های داده اچ-بیس و جدول هایپر بر مبنای جدول بزرگ بنا نهاده شده اند اما پایگاه داده کاساندر ساختار متفاوتی داشته و در آن مفهومی بنام "فوق ستون"^۵ که شامل چندین خانواده ستون می باشد، وجود دارد؛ بنابراین برای ذخیره سازی داده های پیچیده و معنایی مناسبتر می باشد. یکی از تفاوت های پایگاه داده های رابطه ای در مقایسه با پایگاه داده های مبتنی بر ستون نحوه ذخیره سازی مقدار خالی نول^۶ می باشد. در پایگاه داده های رابطه ای باید به ازای هر صفت در یک سطر در صورت نبود

¹ Column family

² HBase

³ Hypertable

⁴ Cassandra

⁵ Super column

⁶ Null

مقدار واقعی، مقدار خالی نول ذخیره شوند در صورتیکه در پایگاه داده‌های مبتنی بر ستون فقط مقادیر موجود کلید (صفت) مقدار ذخیره می‌شود (Hecht, R., S. Jablonski, 2011).

۲-۳-۱-۳ پایگاه داده مبتنی بر گراف

انواع پایگاه‌های ذکر شده تا کنون برای ذخیره سازی توزیع شده حجم زیادی از داده‌های نرمال نشده تعریف شده‌اند و برقراری ارتباط بین داده‌ها باید در سطح برنامه کاربردی صورت گیرد. برخلاف این پایگاه‌ها پایگاه‌های مبتنی بر گراف برای ذخیره‌سازی ارتباطات بین داده‌ها ایجاد شده‌اند و از گره‌ها، ارتباطات بین آن‌ها و ویژگی‌های هر گره ساخته شده‌اند. بعبارتی دیگر، در این مدل پایگاه داده، داده‌ها و ارتباطات بین آن‌ها بصورت شبکه‌ای از راس‌ها و یال‌ها مدل می‌شود و فرآیند مدل سازی و مشخص کردن روابط بین داده‌ها بسیار ساده است. به عنوان مثال برای پیدا کردن دوستان یک فرد در شبکه‌های اجتماعی از الگوریتم‌های کوتاهترین مسیر در گراف برای یافتن نتایج استفاده می‌شود در صورتی که در پایگاه داده‌های رابطه‌ای هزینه بیشتری صرف پیدا کردن چنین رابطه‌هایی می‌شود (Hecht, R., S. Jablonski, 2011). از میان ساختارهای موجود پایگاه‌های داده غیر رابطه‌ای مدل‌های مبتنی بر گراف بیشترین کارایی را در داده‌هایی مانند گراف‌های شبکه‌های اجتماعی دارند. مثال‌هایی از این نوع پایگاه‌های داده نئوفورجی^۱ و زرافه (جراف)^۲ است.

۲-۳-۱-۴ پایگاه داده مبتنی بر سند

برخلاف پایگاه‌های داده رابطه‌ای که در آن‌ها داده‌ها در سطر و ستون ذخیره می‌شود، این مدل از پایگاه‌های داده غیررابطه‌ای داده‌ها را در سندهایی با فرمت‌هایی همچون JSON، XML و BSON

¹ Neo4j

² Giraph

ذخیره می کنند. این مدل برای ذخیره داده‌های بدون ساختار بسیار مناسب است چرا که هر سند می تواند از تعداد متفاوتی فیلد^۱ تشکیل شده باشد و لزومی ندارد که سند های مشابهی که در یک دسته بندی قرار می گیرند حتما از ساختار مشابهی پیروی کنند.

از طرفی پایگاه داده سند گرا، قابلیت‌های بیشتری نسبت به پایگاه داده کلید-مقدار فراهم میکند. واحد ذخیره در این پایگاه یک سند است. سند یک شی است که مجموعه دلخواهی از ویژگیها را شامل می شود که به طور مثال میتواند در قالب^۲ جی-سون ارائه شود. هر سند یک کلید منحصر به فرد دارد که توسط آن شناسایی می شود. ذخیره در قالب سند از جستجو بر اساس فیلدهای یک سند و شاخص گذاری پشتیبانی میکند. به طور معمول شمای داده از پیش تعریف شده ندارند. ذخیره سازی سند جدید با صفات جدید براحتی قابل انجام است. یک امتیاز این سیستمها نسبت به سیستمهای کلید-مقدار، پشتیبانی از قابلیت مجموعه‌ای از پرس وجوها روی یک سند با ایجاد چندین محدودیت روی فیلدهاست. این سیستمها می توانند پرس وجوهای تجمعی را اجرا کنند، نتایج را مرتب کنند و از ایندکس گذاری روی فیلدهای یک سند پشتیبانی کنند (Hecht, R., S. Jablonski, 2011). نمونه‌هایی از این مدل پایگاه‌های داده غیر رابطه ای مبتنی بر سند (Hoberman, S. 2014) Elastic, MongoDB (Divya, M. S., Goyal, S. K. 2013) search می باشند.

۲-۳-۲ مروری بر مطالعات پیشین

امروزه پایگاه‌داده‌های NoSQL در شرکت‌های بزرگ دنیا جهت نگهداری و پردازش داده‌ها به طور گسترده مورد استفاده قرار می گیرند. البته به این معنی نیست که در این شرکت‌ها از پایگاه-داده‌های رابطه‌ای استفاده نمی شود بلکه NoSQLها برای کاربردهای خاص جستجو و پردازش و

¹ Field

² Jason

ذخیره‌سازی در کنار پایگاه‌داده‌های رابطه‌ای قرار می‌گیرند. به عنوان مثال پایگاه‌داده MongoDB در شرکت‌های تجاری مانند اتو^۱، ای بی^۲، گپ^۳، ساکس فیفت اونیو^۴، گیلت گروپ^۵ و ... مورد استفاده قرار گرفته است (MongoDB, 2019) و یا پایگاه‌داده کاساندر در شرکت‌های تجاری مانند کانستانت کانتکت^۶، نت فلیکس^۷، اینستاگرام^۸، هولو^۹، گو ددی^{۱۰}، گیت هاب^{۱۱}، کام کست^{۱۲}، سرن^{۱۳}، د ودر چنل^{۱۴}، ردیت، اینتوییت^{۱۵}، و ... استفاده می‌شود (Han, J., et al., 2011). شایان ذکر است که شبکه اجتماعی تلگرام برای جستجو در هشتگ‌های خود از پایگاه‌داده کاساندر استفاده می‌کند. این مثال‌ها نمونه‌هایی از شرکت‌های بزرگی است که پایگاه‌داده‌های NoSQL را بخدمت گرفته‌اند.

در ادامه به برخی از تحقیقاتی که به بررسی چالش‌ها و کاربردهای کلان‌داده‌ها و پردازش

آن‌ها پرداخته‌اند اشاره می‌کنیم.

کوپر^{۱۶} و همکارانش (۲۰۱۰) با ارائه چارچوب وای-سی-اس-بی^{۱۷} (معیار ارزیابی سرویس ابری یا هوا) راهکاری را برای ارزیابی سرویس تحت ابر ارائه دادند. آن‌ها توسط چارچوب پیشنهادی خود داده‌های

¹ Otto

² eBay

³ Gap

⁴ Saks Fifth Avenue

⁵ Gilt Group

⁶ Constant contact

⁷ Netflix

⁸ Instagram

⁹ Hulu

¹⁰ GoDaddy

¹¹ GitHub

¹² Comcast

¹³ Cern

¹⁴ The Weather Channel

¹⁵ Intuit

¹⁶ Cooper

¹⁷ YCSB (Yahoo! Cloud Service Benchmark)

تولید کرده و پایگاه داده‌های اچ-بیس^۱، کاساندر^۲، پی-نات یا هو^۳، مای اس-کسو-ال ساده شارد^۴ را مقایسه نمودند.

غزال و همکارانش (۲۰۱۳) سیستمی بنام بیگ-بنچ^۵ را برای تولید معیارسنجی^۶ های داده‌های کلان ارائه نمودند که از نظر تنوع و حجم داده، ساختار داده (ساخت یافته، نیمه ه ساخت یافته و بدون ساختار) ویژگی‌های کلان داده‌ها را شبیه‌سازی کند.

آبرامووا و همکارانش (۲۰۱۴) با ایجاد انواع مختلفی از بارکاری‌ها، پنج پایگاه داده غیر رابطه‌ای اچ-بیس، اورینتد دی-بی، ردیس، مانگو دی-بی^۷ را از نظر سرعت خواندن و نوشتن و به روزرسانی با یکدیگر مقایسه کرده‌اند. ردیس و کاساندر در تمام آزمایشات نسبت به سایر پایگاه داده‌ها کارایی بهتری داشتند.

ماوریدیس^۸ و همکارانش (۲۰۱۷) از ابزارهای پردازش کلان داده برای آنالیز فایل لاگ وب سرور^۹ استفاده نمودند. آن‌ها با راه‌اندازی یک کلاستر ۶ گره‌ای کارایی امکانات هادوپ (Hadoop, 2017) و اسپارک (Chambers, B., & Zaharia, M. 2018) ذخیره‌سازی و پردازش کلان داده را با هم مقایسه کردند.

دده^{۱۰} و همکارانش (۲۰۱۳) به ارزیابی الگوریتم مپ ردیوس^{۱۱} در هادوپ و MongoDB پرداختند. آن‌ها برای ارزیابی خود از مجموعه داده ۳۰۰ گیگا بایتی سنسوس یو-اس^{۱۲} استفاده کردند. در تمام

¹ Hbase

² Cassandra

³ Yahoo!'s PNUTS

⁴ Simple sharded MySQL

⁵ BigBench

⁶ Benchmark

⁷ MangoDB

⁸ Mavridis

⁹ Web server log file Analysis

¹⁰ Dede

¹¹ MapReduce

¹² U.S. Census dataset

بارهای کاری تست شده هادوپ کارایی بهتری از خود نشان داده است.

یکی از کاربردهای سیستم‌های کلان داده تجزیه و تحلیل موقعیت حرکتی افراد و تجزیه و تحلیل رفتار فرد یا بیمار برای حفظ سلامت و نظارت بر نحوه درست حرکت آنهاست. دولو و همکارانش (1395) با ارائه یک نرم‌افزار موبایل موقعیت مکان افراد را استخراج و از ابزارهای کلان داده برای پردازش آن استفاده نمودند. آنها از محیط هادوپ^۱ و ابزار هایو^۲ و ابزار هایو برای استخراج الگوهای حرکتی در داده‌ها و کمک به سلامت افراد بهره بردند.

نوریان و همکارانش (۱۳۹۲) از تفسیر کلان-داده‌ها برای اعتبارسنجی معیارهای تجربی مربوط به توسعه مبتنی بر حمل و نقل عمومی و برنامه‌ریزی حمل و نقل شهری در خصوص سفرسازی استفاده نمودند. بدین منظور الگوی ترافیکی کاربری آموزش عالی با روش ارزیابی همبستگی تحلیل می‌گردد. با به کارگیری نگرش شهرسازی در پالایش کلان-داده‌های ترافیکی موجود و در دسترس برای دانشگاه ایالتی کالیفرنیا در شهر لوس آنجلس، الگوی ترافیکی سفرهای تحصیلی استخراج شده است. نصری فلاح (۱۳۹۵) به بررسی کاربردها و چالش‌های کلان‌داده‌ها پرداخته است. مدبرفرزام و همکارانش (۱۳۹۴) ضمن بررسی چالش‌های کلان‌داده‌ها ویژگی‌ها و انواع پایگاه‌داده‌های NoSQL را بررسی کرده‌اند.

در نهایت با بررسی پژوهش‌های پیشین در این حوزه و در نظر گرفتن نوع داده‌های در دست در مرکز منطقه ای و هدف این پژوهش که بررسی عملیات جستجو و بازیابی تمام متن بر روی آنهاست؛ پایگاه‌های داده غیررابطه ای /سیستم‌های NOSQL محبوب و منتخب MongoDB و

¹ Hadoop

² Hive

Elasticsearch که در دسته مبتنی بر سند قرار می گیرند برای ارزیابی درقیاس با پایگاههای داده رابطه ای لحاظ شده اند؛ بنابراین در ادامه به شرح مختصری در خصوص آنها می پردازیم.

۲-۳-۳ پایگاه داده MongoDB

MongoDB یک پایگاه داده غیر رابطه ای مبتنی بر سند است که رکورد های داده را تحت قالب BSON ذخیره می کند. مفاهیم جدول^۱، ردیف^۲ و ستون^۳ که در پایگاه های داده رابطه ای مطرح هستند در این پایگاه داده به ترتیب با مفاهیم جدیدی به نام های^۴؛ کلکسیون^۴، سند^۵ و فیلد^۶ جایگزین می شوند. همچنین برخلاف پایگاه های داده رابطه ای که برای بازیابی اطلاعات مرتبط که در جدول های متفاوتی قرار گرفته اند، احتیاج به SQL Join می باشد، در MongoDB به این علت که اطلاعات مرتبط به هم بصورت سندهای تعبیه شده^۷ در یک جا قرار می گیرند، مفهوم Join وجود ندارد. راهکار این پایگاه داده برای ارائه قابلیت مقیاس پذیری افقی، خرد کردن به روش شاردینگ^۸ است. در صورتی که این قابلیت بر روی یک پایگاه داده فعال شود، کلکسیونهای آن پایگاه داده می توانند بین نمونه های^۹ متفاوتی از MongoDB که در یک^{۱۰} خوشه خرد شده قرار دارند، توزیع شوند. در صورت فعال سازی قابلیت شاردینگ بر روی یک کلکسیون ، سندهای آن کلکسیون بر روی نمونه

¹ Table

² Row

³ Column

⁴ Collection

⁵ Document

⁶ Field

⁷ Embedded Document

⁸ Sharding

⁹ Instances

¹⁰ Sharded Cluster

های مجزا توزیع خواهند شد. وظیفه توزیع متعادل این سندها بر روی نمونه های مختلف، بر عهده بالانس کننده^۱ می باشد. (ElasticSearch, 2019; Dixit, B.,2016; MongoDB, 2019)

۲-۳-۴ پایگاه داده Elasticsearch

Elasticsearch بیش از آنکه به عنوان یک پایگاه داده شناخته شود، به عنوان یک موتور جستجوی-تمام متن^۲ و توزیع شده مطرح است که بر پایه کتابخانه لوسنس^۳ و به زبان جاوا^۴ توسعه داده شده است. پس از انتشار نسخه اولیه Elasticsearch در سال ۲۰۱۰، این موتور جستجو به طور گسترده در سازمان های بزرگی همچون ناسا^۵، ویکی پدیا^۶ و گیت هاب^۷ N به کار گرفته شده است. آخرین نسخه های انتشار یافته از Elasticsearch، به هدف اینکه علاوه بر موتور جستجو بودن به عنوان پایگاه داده نیز مورد اعتماد کاربران قرار گیرد، بیشتر بر ویژگی انعطاف پذیری^۸ متمرکز بوده است. Elasticsearch بر پایه معماری رست^۹ بنا شده و برقراری ارتباط با این پایگاه داده برای انجام عملیات مختلف و تغییر تنظیمات، توسط درخواست های HTTP صورت می گیرد. در این پایگاه داده، اصطلاحات کلکسیون و پایگاه داده که در MongoDB معرفی شدند، به ترتیب با مفاهیم مپینگ^{۱۰} و ایندکس^{۱۱} ایندکس جایگزین می شوند. سندها نیز در Elasticsearch با فرمت جی-سون^{۱۲} ذخیره می شوند. Elasticsearch نیز به منظور فراهم کردن قابلیت مقیاس پذیری افقی مفهوم شاردرینگ و

¹ Balancer

² Full-Text Search

³ Lucene

⁴ Java

⁵ Nasa

⁶ WikiPedia

⁷ Github

⁸ Resiliency

⁹ REST

¹⁰ Mapping

¹¹ Index

¹² JSON

خرد کردن را معرفی می کند. در یک Elasticsearch کلاستر این قابلیت وجود دارد که ایندکس ایندکس های متفاوت با تعداد شارد خرد شده متفاوتی تعریف شوند. به عبارتی کاربر این اختیار را دارد که هر ایندکس ایندکس را به تعداد قطعات متفاوتی بشکند و بر روی نود های کلاستر توزیع کند (Chambers, B., & Zaharia, M. 2018).

۲-۳-۵ مقایسه ویژگی های پایگاه های داده MongoDB و Elasticsearch

با توجه به اینکه پایگاه های داده Elasticsearch و MongoDB هر دو جزء دسته بندی پایگاه های داده غیر رابطه ای مبتنی بر سند قرار می گیرند، نقاط مشترک متعددی بین آن ها وجود دارد. از جمله دارا بودن ویژگی هایی همچون پذیرفتن داده های بدون ساختار بودن^۱، ذخیره و بازیابی داده ها بصورت توزیع شده^۲ (و دسترسی بالا^۳ قابل ذکر است).

با این وجود تفاوت های اساسی بین آن ها وجود دارد که باعث شده هر کدام در جایگاه خاص خود مورد استفاده قرار گیرد. از جمله این تفاوت ها می توان به موارد زیر اشاره کرد (MongoDb Manual Document, 2019; Github, 2019; Elasticsearch Resiliency Status, 2019)

- MongoDB به عنوان یک پایگاه داده همه منظوره^۴ شناخته می شود و در نقطه مقابل Elasticsearch برای هدف خاص تری طراحی شده است و آن قابلیت ایندکسینگ و جستجوی قدرتمند و توزیع شده بین سند ها در ابعاد بزرگ است.

- با توجه به همین تفاوت ماهیت این دو پایگاه داده دور از انتظار نیست که Elasticsearch به اندازه MongoDB، به عنوان یک ابزار ذخیره سازی داده ها، انتظارات را برآورده نکند. همینطور قابلیت

¹ Scheme Free

² Sharding

³ High Availability

⁴ Genral Purpose

های ایندکسینگ و جستجوی تمام متن قدرتمند Elasticsearch را نمی توان از MongoDB انتظار داشت.

- علی رغم تلاش های انجام شده، بزرگترین ضعف Elasticsearch همچنان عدم قابلیت اعتماد بودن آن در نگهداری و حفظ داده ها به عنوان یک سند است. این مشکل معمولا زمانی نمود پیدا می کند که کلاستر تغییر وضعیت داده باشد و نتیجتا باعث می شود که Elasticsearch به عنوان پایگاه داده اصلی مورد اطمینان نباشد.

- Elasticsearch از واسط ای-پی-آی رست^۲ و پروتکل قدرتمند HTTP بهره می گیرد و این مورد از نقاط قوت آن نسبت به MongoDB به شمار می آید.

در پایان این فصل به مقایسه ای بین پایگاههای داده رابطه ای و غیر رابطه ای کاربردی در پژوهش حاضر در قالب جدول های ۲-۳ و ۲-۲ می پردازیم.

جدول ۲-۲ مقایسه ویژگیهای کلی پایگاههای داده رابطه ای و غیر رابطه ای مورد بررسی پژوهش حاضر

مدل داده ها	کار آیی	مقیاس پذیری	قابلیت انعطاف پذیری	پیچیدگی	عملکرد
پایگاه داده رابطه ای	متغیر	متغیر	پایین	متوسط	متغیر
پایگاه داده غیر رابطه ای	بالا	بالا	بالا	متغیر (پایین)	متغیر

¹ Reliable

² API REST

جدول ۲-۳ مقایسه مزایا و معایب پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای کاربردی در پژوهش حاضر
(Chen, J. K. et al., 2019; DB-Engines Ranking, 2017; Oracle-RDBMS, 2017; Kenler, E., 2015; Li, Y, et al., 2013)

موارد مقایسه	پایگاه داده رابطه‌ای / Relational Data Base		پایگاه داده غیررابطه‌ای (NoSQL) / Non-Relational Data Base	
	Microsoft SQL Server	Maria DB	MangoDB	Elastic Search
شناسنامه	<ul style="list-style-type: none"> توسعه یافته : Microsoft Corporation انتشار اولیه : ۲۴ آوریل ۱۹۸۹ سیستم عامل: Linux, Microsoft Windows Server, Microsoft Windows زبان توسعه : C, C++ 	<ul style="list-style-type: none"> توسعه یافته : MariaDB Corporation AB, MariaDB Foundation انتشار اولیه : ۲۹ اکتبر ۲۰۰۹ سیستم عامل: Windows, Solaris, FreeBSD, Linux, X 10.7 OS, Vista زبان توسعه: Windows, macOS 	<ul style="list-style-type: none"> توسعه یافته : MongoDB Inc. انتشار اولیه : ۱۱ فوریه ۲۰۰۹ سیستم عامل: Solaris, Windows, FreeBSD, Vista, Linux زبان توسعه : جاوا اسکریپت، پایتون، Go، C++ 	<ul style="list-style-type: none"> توسعه یافته : Elastic NV انتشار اولیه : ۸ فوریه ۲۰۱۰ سیستم عامل: Windows Vista and later, Linux, OS X 10.7 and later, Solaris, FreeBSD زبان توسعه : جاوا

کاربرد (use case)	<ul style="list-style-type: none"> • اگر انواع داده های مختلفی دارید که در قسمت های مختلف برنامه تان از آن ها استفاده می کنید بهتر است از SQL ها استفاده کنید. • اگر داشتن یک ساختار ثابت برایتان مهم است و احتمال تغییر داده هایتان وجود ندارد می توانید از SQL ها استفاده کنید. 	<ul style="list-style-type: none"> • اگر انواع داده های مختلفی دارید که در قسمت های مختلف برنامه تان از آن ها استفاده می کنید بهتر است از SQL ها استفاده کنید . • اگر داشتن یک ساختار ثابت برایتان مهم است و احتمال تغییر داده هایتان وجود ندارد می توانید از SQL ها استفاده کنید. 	<ul style="list-style-type: none"> • اگر انتظار دارید داده هایتان در آینده تغییر کند یا ساختار داده ها در آینده مشخص نیست . • اگر برنامه ی شما درخواست های دریافت خوانش (Read) زیادی دارد اما داده ها را زیاد تغییر نمی دهید. • اگر بعدا به مقیاس دهی افقی نیاز پیدا خواهید کرد باید NoSQL -Elastic Search را انتخاب کنید. 	<ul style="list-style-type: none"> • اگر انتظار دارید داده هایتان در آینده تغییر کند یا ساختار داده ها در آینده مشخص نیست. • اگر برنامه ی شما درخواست های دریافت خوانش (Read) زیادی دارد اما داده ها را زیاد تغییر نمی دهید. • اگر بعدا به مقیاس دهی افقی نیاز پیدا خواهید کرد باید NoSQL -MongoDB را انتخاب کنید.
----------------------	---	--	---	---

ریسک (Risk)	<ul style="list-style-type: none"> • ریسک بالای در معرض قرار گرفتن حملات SQL-INJECTION • عدم مقیاس پذیری افقی و در نتیجه ساپورت حجم محدودی داده 	<ul style="list-style-type: none"> • ریسک بالای در معرض قرار گرفتن حملات SQL-INJECTION • عدم مقیاس پذیری افقی و در نتیجه ساپورت حجم محدودی داده 	<ul style="list-style-type: none"> • بدلیل ساختار اولیه طراحی با ریسک کمتری از بابت حملات پایگاه داده ای مواجه است • انعطاف پذیری بیش از حد ممکن است باعث تنبل شدن نیروی سازمانی و پیاده نکردن ساختاری مشخص برای پایگاه داده تان شود. 	<ul style="list-style-type: none"> • بدلیل ساختار اولیه طراحی با ریسک کمتری از بابت حملات پایگاه داده ای مواجه است • انعطاف پذیری بیش از حد ممکن است باعث تنبل شدن نیروی سازمانی و پیاده نکردن ساختاری مشخص برای پایگاه داده تان شود.
------------------------	---	---	---	---

مزایای کلی (Pros)	<ul style="list-style-type: none"> • قدرت اصلی این مدل، در استفاده از مفهوم جدول است که روشی بصری، کارآمد و تا حدودی انعطاف پذیر برای ذخیره و دسترسی به اطلاعات ساختاریافته است. • از محبوب ترین و کاربرپسندترین پایگاه‌های داده رابطه‌ای است. 	<ul style="list-style-type: none"> • از محبوب ترین و کاربرپسندترین پایگاه‌های داده رابطه‌ای است که نسبت به MySQL برتری‌هایی دارد. • قدرت اصلی این مدل، در استفاده از مفهوم جدول است که روشی بصری، کارآمد و تا حدودی انعطاف پذیر برای ذخیره و دسترسی به اطلاعات ساختاریافته است. • مهاجرت راحت از سایر سیستم های پایگاه داده به MariaDB از دیگر مزایای آن است. • منبع باز و رایگان 	<ul style="list-style-type: none"> • یک پایگاه داده همه منظوره (General Purpose) با خصوصیات بهینه برای داده های بدون ساختار ذخیره و بازیابی داده‌های بدون ساختار در قالب فرمت های مختلف (Schema Free) • ذخیره و بازیابی داده ها بصورت توزیع شده • انعطاف پذیری و سازگاری بیشتر • دسترسی بالا (High Availability) • عدم نیاز به استفاده از پرس و جوهای آزمایشی پیچیده 	<ul style="list-style-type: none"> • بهره بردن از واسط API REST و پروتکل قدرتمند HTTP • قابلیت های ایندکسینگ و جستجوی تمام متن قدرتمند • ذخیره و بازیابی داده‌های بدون ساختار در قالب فرمت های مختلف (Schema Free) • ذخیره و بازیابی داده ها بصورت توزیع شده • انعطاف پذیری و سازگاری بیشتر • دسترسی بالا (High Availability) • عدم نیاز به استفاده از پرس و جوهای آزمایشی پیچیده • افزایش سرعت توسعه برنامه‌های مختلف • افزایش سرعت دسترسی به داده‌های مختلف
------------------------------	--	---	---	--

			<ul style="list-style-type: none"> • افزایش سرعت توسعه برنامه‌های مختلف • افزایش سرعت دسترسی به داده‌های مختلف • قابلیت مقیاس پذیری عمودی و افقی بدون حد نصاب برای داده‌های عظیم • قابلیت مقیاس پذیری افقی در قالب یک سیستم توزیع شده • منبع باز و رایگان 	<ul style="list-style-type: none"> • قابلیت مقیاس پذیری عمودی و افقی بدون حد نصاب برای داده‌های عظیم • قابلیت مقیاس پذیری افقی در قالب یک سیستم توزیع شده • منبع باز و رایگان
--	--	--	--	--

<p style="text-align: center;">معایب کلی (Cons)</p>	<p>بطور کلی محدودیتهایی از نظر مقیاس پذیری، موازی سازی عملیات و هزینه به شرح زیر دارد:</p> <ul style="list-style-type: none"> • بازیابی اطلاعات پیچیده که منجر به افت عملکرد عملیات می‌شود. • انعطاف پذیری کم به دلیل ساختار از پیش تعیین شده. • بسیار سخت و یا غیر ممکن بودن تغییر ساختار جداول در آینده • روابط ممکن است باعث ایجاد پرس و جوهای آزمایشی بسیار پیچیده ی شوند. 	<p>بطور کلی محدودیتهایی از نظر مقیاس پذیری، موازی سازی عملیات و هزینه به شرح زیر دارد:</p> <ul style="list-style-type: none"> • بازیابی اطلاعات پیچیده که منجر به افت عملکرد عملیات می‌شود. • انعطاف پذیری کم به دلیل ساختار از پیش تعیین شده. • بسیار سخت و یا غیر ممکن بودن تغییر ساختار جداول در آینده • روابط ممکن است باعث ایجاد پرس و جوهای آزمایشی بسیار پیچیده ی شوند. 	<p>ضعف نسبی در ابزار ایندکسینگ و جستجوی تمام متن داده‌ها</p> <ul style="list-style-type: none"> • انعطاف پذیری بیش از حد ممکن است باعث پیاده نکردن ساختاری مشخص و مدون برای پایگاه داده سازمانی شود. • داده های تکراری باعث می شوند که برای انجام عملیاتی خاص مانند به روز رسانی مجبور شوید چندین دستور به روز رسانی را برای چندین مقدار مختلف اجرا کنید 	<ul style="list-style-type: none"> • بزرگترین ضعف Elasticsearch عدم reliable بودن آن در نگهداری و حفظ داده‌ها به عنوان یک پایگاه داده است. • انعطاف پذیری بیش از حد ممکن است باعث پیاده نکردن ساختاری مشخص و مدون برای پایگاه داده سازمانی شود • داده های تکراری باعث می شوند که برای انجام عملیاتی خاص مانند به روز رسانی مجبور شوید چندین دستور به روز رسانی را برای چندین مقدار مختلف اجرا کنید.
---	--	--	--	--

	<ul style="list-style-type: none"> • مقیاس دهی افقی در کار سختی است و در اکثر مواقع مقیاس دهی فقط عمودی است بنابراین اگر شرکت شما بسیار بزرگ شود به حد نصاب خاصی خواهید رسید • هزینه لیسانس بالا. 	<ul style="list-style-type: none"> • مقیاس دهی افقی کار سختی است و در اکثر مواقع مقیاس دهی فقط عمودی است بنابراین اگر شرکت شما بسیار بزرگ شود به حد نصاب خاصی خواهید رسید. • هیچگونه استاندارد زبانی پرس و جویهای آزمایشی ندارد. 		
--	---	--	--	--

سایر نکات	<ul style="list-style-type: none"> • از قدیمیترین و شناخته شده ترین پایگاه داده های رابطه ای با محبوبیت و قابلیت کاربری خوب با استفاده از رابط کاربری آسان • در این نوع پایگاه داده هر سطر یک جدول، یک رکورد با یک کلید شناسایی منحصر به فرد است و هر ستون نمایانگر یک ویژگی در مورد رکورد های آن جدول است • پایگاه داده ای است که داده های مرتبط به هم را ذخیره کرده و دسترسی به آن ها را فراهم می آورد. به 	<ul style="list-style-type: none"> • پس از انتقال مالکیت کامل پروژه MySQL از شرکت MySQL AB به شرکت Oracle در بین سال های ۲۰۰۸ و ۲۰۰۹، توسعه دهندگان اصلی MySQL، پایگاه داده MariaDB را به عنوان یک fork از MySQL معرفی کردند و تا به امروز این پایگاه داده نوپا در بسیاری از شرکت ها جایگزین MySQL شده است • MariaDB در مقایسه با MySQL از موتورهای بیشتری پشتیبانی می کند. • MariaDB با توجه به منبع باز بودن در مقایسه با 	<ul style="list-style-type: none"> • MongoDB یک پایگاه داده همه منظوره (General Purpose) با خصوصیات بهینه برای داده های بدون ساختار است. • به اندازه Elasticsearch، به عنوان یک ابزار قدرتمند ایندکسینگ و جستجوی تمام متن داده ها، انتظارات را برآورده نمی کند 	<ul style="list-style-type: none"> • Elasticsearch طراحی شده است برای اهداف خاص از جمله قابلیت ایندکسینگ و جستجوی قدرتمند و توزیع شده بین Document ها در ابعاد بزرگ • Elasticsearch به اندازه MongoDB، به عنوان یک ابزار ذخیره سازی داده ها، انتظارات را برآورده نمی کند • Document ها نیز در Elasticsearch با فرمت JSON ذخیره می شوند. • به هدف اینکه علاوه بر موتور جستجو بودن به عنوان پایگاه داده نیز مورد اعتماد کاربران قرار گیرد، بیشتر بر ویژگی resiliency متمرکز بوده است.
------------------	---	--	--	---

<p>عبارتی این نوع پایگاه‌های داده مبتنی بر مدل Relational می‌باشند که یک راه شهودی و قابل درک برای نشان دادن اطلاعات در جداول است.</p> <p>• در این نوع پایگاه داده هر سطر یک جدول، یک رکورد با یک کلید شناسایی منحصر به فرد است و هر ستون نمایانگر یک ویژگی در مورد رکورد های آن جدول است</p>	<p>MySQL به طور مداوم در حال توسعه است.</p> <p>• MariaDB همچنین یک پایگاه داده خوشه ای برای استفاده تجاری ارائه می دهد ، که همچنین امکان تکثیر چند اسناد را فراهم می کند.</p> <p>• MariaDB برای کارایی بالا بهینه شده است و برای مقادیر زیادی مجموعه داده بسیار قدرتمندتر از MySQL است.</p>		<p>• Elasticsearch نیز به منظور فراهم کردن قابلیت Horizontal Scalability مفهوم شاردینگ را معرفی می کند</p>
---	---	--	--

فصل سوم: روش شناسی پژوهش

(مدل سازی و پیاده سازی آزمایشگاهی)

۳-۱ مقدمه

با در نظر داشتن هدف اصلی پژوهش حاضر مبنی بر بررسی عملکرد جستجوی تمام متن بر پایگاه های داده رابطه ای در قیاس با پایگاه داده های غیر رابطه ای مبتنی بر متن و متناسب با پژوهشهای انجام شده در خصوص پایگاههای داده غیر رابطه ای (MongoDB, 2019; Cooper B.f.) (et al., 2010) ، دو مدل پایگاه های داده غیر رابطه ای MongoDB و Elasticsearch به عنوان سیستمهای NOSQL محبوب جهت مدل سازی و پیاده سازی آزمایشگاهی برگزیده شدند. به روش مشابه Microsoft SQL Server و MariaDB نیز بعنوان پایگاه داده های رابطه ای منتخب برگزیده شدند. سپس سایر مراحل پیاده سازی و مدل سازی سخت افزاری و نرم افزاری مورد نیاز در حالت آزمایشگاهی (مشروح در ادامه فصل) صورت گرفت تا بتوان به ارزیابی مد نظر در دو فاز زیر پردازیم :

- بررسی عملکرد بارهای کاری ۱۰۰٪ خوانش و ۱۰۰٪ اسکن برای پایگاه های داده مورد بحث به کمک ارزیابی وای-سی-اس-بی^۱

¹ YCSB Benchmark

• بررسی عملکرد پرس و جوهای آزمایشی جستجوی تمام متن برای پایگاههای داده مورد

بحث و همچنین در نظر گرفتن تاثیر قابلیت مقیاس پذیری افقی بر عملکرد این پرس و

جوهای آزمایشی در پایگاههای داده غیر رابطه‌ای MongoDB و Elasticsearch

۲-۳ مشخصات فنی سیستم و نسخه پایگاههای داده مورد استفاده

۱-۲-۳ مشخصات فنی سیستم

برای انجام ارزیابیها از هفت سیستم با مشخصات فنی مندرج در جدول زیر استفاده شده است.

جدول ۱-۳ مشخصات سیستمها

نام سیستم	IP آدرس
سیستم ۱	172.16.13.26
سیستم ۲	172.16.13.25
سیستم ۳	172.16.13.24
سیستم ۴	172.16.13.23
سیستم ۵	172.16.13.22
سیستم ۶	172.16.13.21
سیستم ۷	172.16.13.20

مشخصات سخت افزاری سیستم ۱ به شرح زیر می باشد:

- Operating System: Linux CentOS 7
- Architecture: x86_64
- Processor: Intel(R) Core(TM) i7-6700K CPU @ 4.00GHz
- Memory: 16GiB DIMM DDR4 Synchronous 2400 MHz (0.4 ns)
- HDD: 2TB WDC WD20EZRZ-00Z

در قسمت ارزیابی جستجو و بازیابی تمام متن، برای بررسی تاثیرشاردینگ یامقیاس پذیری افقی، ۶ سیستم دیگر (سیستم ۲ الی ۶) با ویژگی های مشترک سخت افزاری، برای میزبانی شاردها استفاده شده است.

- Operating System: Linux CentOS 7
- Architecture: x86_64
- Processor: Intel(R) Core(TM) i7 CPU 920 @ 2.67GHz
- Memory: 11GiB DIMM DRAM EDO
- HDD: 429GB Virtual disk

۳-۲-۲ نسخه پایگاههای داده مورد استفاده

نام پایگاه داده و نسخه های مورد استفاده هر کدام در این پژوهش در جدول زیر به اختصار آمده است.

جدول ۳-۲ پایگاههای داده مورد استفاده

نسخه	نام پایگاه داده
15.0.1800.32	Microsoft SQL Server 2019 (CTP3.2)
5.5.1	Elasticsearch
10.3.17	MariaDB Server
4.0.9	MongoDB

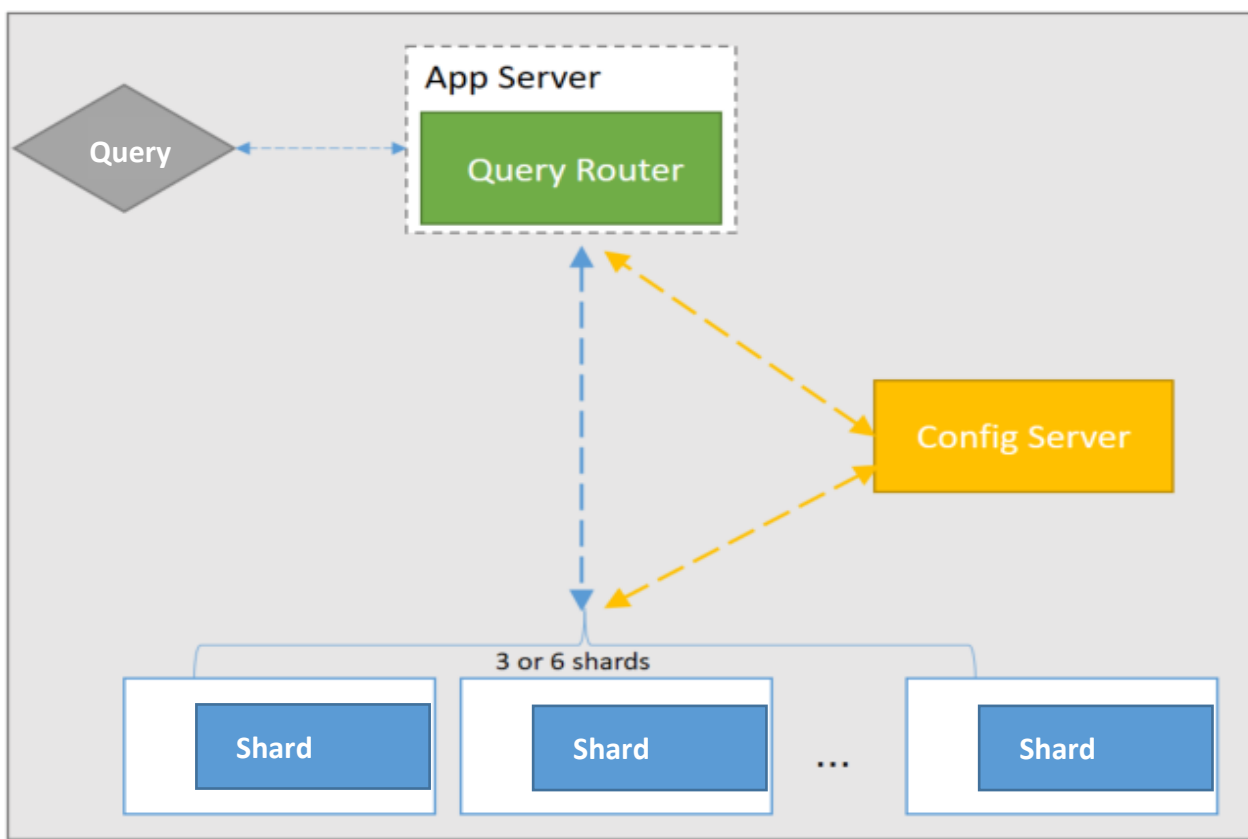
۳-۳ کلیات مدل سازی و پیاده سازی پایگاه‌های داده

در فاز نخست که در قسمت ۴-۲ به شرح جزئیات آن پرداخته می شود، چهار پایگاه داده Microsoft SQL Server، MariaDB، MongoDB و Elasticsearch به صورت تک نمونه^۱ روی سیستم 1 (که مشخصات فنی آن در قسمت ۳-۳ آمده است) پیاده سازی شده و همچنین اجرای امعیار ارزیابی وای-سی-اس-بی^۲ بر روی این پایگاه‌های داده انجام می شود.

در فاز بعدی که در قسمت ۴-۳ به شرح جزئیات آن پرداخته می شود به بررسی مقایسه ای بین پایگاه‌های داده رابطه ای و غیر رابطه ای در جستجوی تمام متن در دو حالت بدون خوشه بندی (۴-۳-۱) و با خوشه بندی (۴-۳-۲) می پردازیم. در حالت خوشه بندی، کلاسترهای MongoDB و Elasticsearch نیز در حالت^۳ شارد و ۶ شارد^۴ پیاده سازی و عملکرد پرس و جوهای آزمایشی جستجوی تمام متن بر روی آن ها نیز بررسی می شود.

یک کلاستر در MongoDB و یا اصطلاحاً کلاستر شارد و خرد شده^۴ از سه جزء روتر پرسش و پاسخ^۵ (برای دریافت درخواست هایکلاینت)، کانفیگ سرور^۶ و شارد/خرد شده^۷ تشکیل شده است. در شکل ۳-۱ معماری اجزای تشکیل دهنده MongoDB و ارتباط بین آنها نمایش داده شده است .

¹ Single Instance
² YCSB Benchmark
³ Shard
⁴ Sharded Cluster
⁵ Query Router
⁶ Config Server
⁷ Shard



شکل ۱-۳ معماری پیاده سازی شده اجزای تشکیل دهنده بر بستر MongoDB

در این پژوهش برای پیاده سازی یک کلاستر شارد و خرد شده^۱ از سیستم ۱ به عنوان میزبان بخش های سرور کانفیگ^۲ و روتر پرسش و پاسخ^۳ و از سایر سیستم ها برای میزبانی شاردها استفاده شده است (در حالت ۶ شارد از هر شش سیستم و در حالت ۳ شارد فقط از سه سیستم استفاده شده است). جزئیات پیاده سازی یک کلاستر شارد خرد شده در (MongoDB Manual Document, 2019) که قسمتی از مستندات رسمی MongoDB است، آمده است.

¹ Sharded Cluster

² Config Server

³ Query Server

یک کلاستر Elasticsearch مجموعه ای از نودهاست که هر کدام می توانند نقش هایی از قبیل مستر^۱، داده^۲، کوردیناتینگ^۳ و یا ترکیبی از آن ها را داشته باشند (ElasticSearch, 2019). در این پژوهش، سیستم 1 به عنوان نود کوردیناتینگ^۴ تنظیم شده که صرفاً درخواست های کلاینت را دریافت کند. سیستم ۷ هم به عنوان نود داده^۵ (برای میزبانیشارد) و هم به عنوان مستر نود^۶ (برای مدیریت کلاستر) استفاده شده و ۵ سیستم دیگر صرفاً به عنوان نود داده تنظیم شده است. جزئیات پیکرندی نودها در (ElasticSearch Manual Document1, 2019 ; ElasticSearch Manual) که قسمتی از مستندات رسمی Elasticsearch می باشد آمده است. در شکل ۳-۲ معماری اجزای تشکیل دهنده Elasticsearch و ارتباط بین آنها نمایش داده شده است.

¹ Master

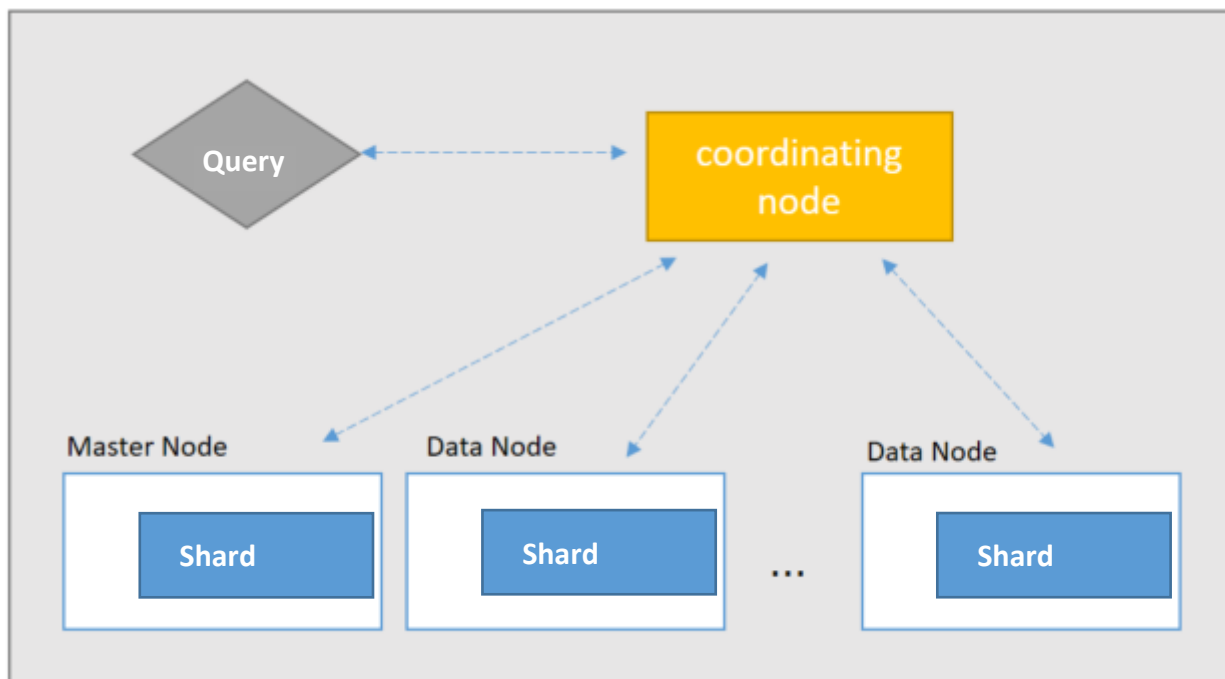
² Data

³ Coordinating

⁴ Coordinating Node

⁵ Node Data

⁶ Master Node



شکل ۲-۳ معماری پیاده سازی شده اجزای تشکیل دهنده بر بستر Elasticsearch

فصل چهارم: یافته‌های پژوهش

۴-۱ مقدمه

در این فصل به بررسی و ارزیابی نتایج و یافته‌های حاصل از سناریوهای مختلفی که با هدف بررسی مقایسه‌ای عملکرد کلی و جستجوی تمام متن پایگاه‌های داده رابطه‌ای SQL Server و MariaDB و پایگاه‌های داده غیر رابطه‌ای MongoDB و Elasticsearch در دایرۀ دایرۀ کاملاً مجزا پرداخته می‌شود. در بخش اول آزمایشات، ابزار مورد استفاده معیار ارزیابی سرویس ابری یا هو^۱ و یا به اختصار وای-سی-اس-بی^۲ است که در بسیاری از تحقیقات آکادمیک در زمینه مقایسه پایگاه‌های داده مختلف به عنوان معیار ارزیابی^۳ استاندارد مورد استفاده قرار می‌گیرد. در بخش دوم نیز به بررسی عملکرد جستجوی تمام متن برای پایگاه‌های داده SQL Server، MariaDB، MongoDB و Elasticsearch با استفاده از چهار نوع پرس و جوهای آزمایشی^۴ پرداخته خواهد شد. جزئیات بیشتر نتایج این آزمایشات با انواع سناریوهای طراحی شده در قسمت‌های ۴-۲ و ۴-۳ بتفصیل قابل مطالعه و ارزیابی است.

¹ Yahoo! Cloud Service Benchmark

² YCSB

³ Benchmark

⁴ Query

۲-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای در حالت اجرای YCSB Benchmark

در این قسمت به بررسی عملکرد کلی پایگاه های داده رابطه ای SQL Server و MariaDB و پایگاه های داده غیر رابطه ای MongoDB و Elasticsearch و مقایسه آنها با یکدیگر می پردازیم. ابزار مورد استفاده در این قسمت معیار ارزیابی سرویس ابری یاهو!^۱ و یا به اختصار وای-سی-اس-بی است که در بسیاری از تحقیقات آکادمیک در زمینه مقایسه پایگاههای داده مختلف به عنوان معیار ارزیابی استاندارد مورد استفاده قرار می گیرد. هدف اصلی پروژههای وای-سی-اس-بی ایجاد یک چارچوب و مجموعه ای از بارهای کاری متداول برای ارزیابی عملکرد سرویس های پایگاه داده مختلف است. این پروژه از دو قسمت اصلی تشکیل شده است (Yahoo! Cloud Serving Benchmark (YCSB), 2019):

- کلاینت وای-سی-اس-بی که بارهای کاری^۲ قابل توسعه را تولید می کند.
- بارهای کاری هسته مرکزی^۳ که مجموعه ای بارهای کاری اصلی هستند که توسط تولید کننده^۴ اجرا می شوند.

YCSB که به زبان جاوا توسعه داده شده است برای پایگاه های داده متعددی اینترفیس ارائه کرده و همچنین با توجه به متن باز بودن این پروژه برای سایر پایگاه داده ها نیز قابل استفاده و توسعه توسط افراد مختلف است.

¹ Yahoo! Cloud Service Benchmark

² Workload

³ Core Workloads

⁴ Generator

در این پژوهش با توجه به اهداف پژوهش، دو بارکاری^۱ مورد استفاده قرار گرفته است. بارکاری تحت عنوان WorkloadA که تماماً شامل عملیات خوانش^۲ می‌شود (معادل همان WorkloadC که بطور پیش فرض در YCSB تعریف شده است) و WorkloadB که تماماً شامل عملیات اسکن^۳ می‌باشد. اجرای هر بارکاری شامل دو مرحله است. در مرحله اول با مشخص کردن تعداد رکورد های مورد نظر، بایستی رکورد ها در پایگاه داده بارگذاری^۴ شوند و در مرحله بعدی با مشخص کردن تعداد عملیات مورد نظر، بارکاری بر روی پایگاه داده اجرا^۵ شود. همچنین ساختار هر رکورد بدین صورت است که هر کدام شامل ۱۰ فیلد می باشند و اندازه هر فیلد ۱۰۰ بایت است.

در این پژوهش در هر مرحله بارگذاری و اجرای بار کاری، تعداد رکورد ها و تعداد عملیات برابر در نظر گرفته شده و به ازای عدد های هزار، ده هزار و صد هزار نتایج جمع آوری شده است. هنگام اجرای هر بارکاری به منظور بهره وری بیشتر از منابع سیستم می توان از قابلیت چند رشته ای^۶ استفاده کرد و تعداد رشته های^۷ مورد نظر را تعیین کرد. در این پژوهش به ازای تعداد رشته های ۱، ۲، ۴ و ۸ نتایج جمع آوری و مورد بحث واقع شده است. پارامتر های در نظر گرفته شده برای مقایسه توانش زمانی^۸ و میانگین تاخیر زمانی^۹ است. همچنین اجرای هر بارکاری با مشخصات خاص، برای ۵ مرتبه تکرار شده و نمودار ها با توجه به میانگین این نتایج رسم شده اند.

¹ Workload

² Read Operation

³ Scan

⁴ Load

⁵ Run

⁶ Multi-threading

⁷ Thread

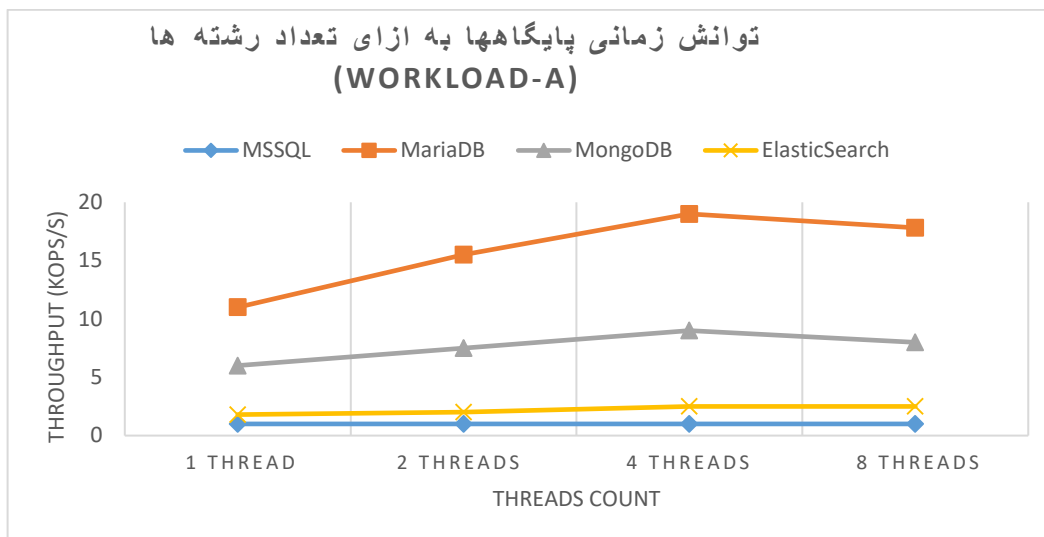
^۸ Throughput: توانش زمانی یا تعداد عملیات انجام شده در واحد زمان
^۹ Latency: در ترجمه لفظ latency زمان بازیابی داده، زمان پاسخگویی به یک درخواست و تاخیر زمانی الفظی است که در طول متن با توجه به ماهیت و بافت متن تخصصی با آن برخورد خواهیم داشت.

۴-۲-۱ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای به ازای تعداد رشته های^۱ متفاوت

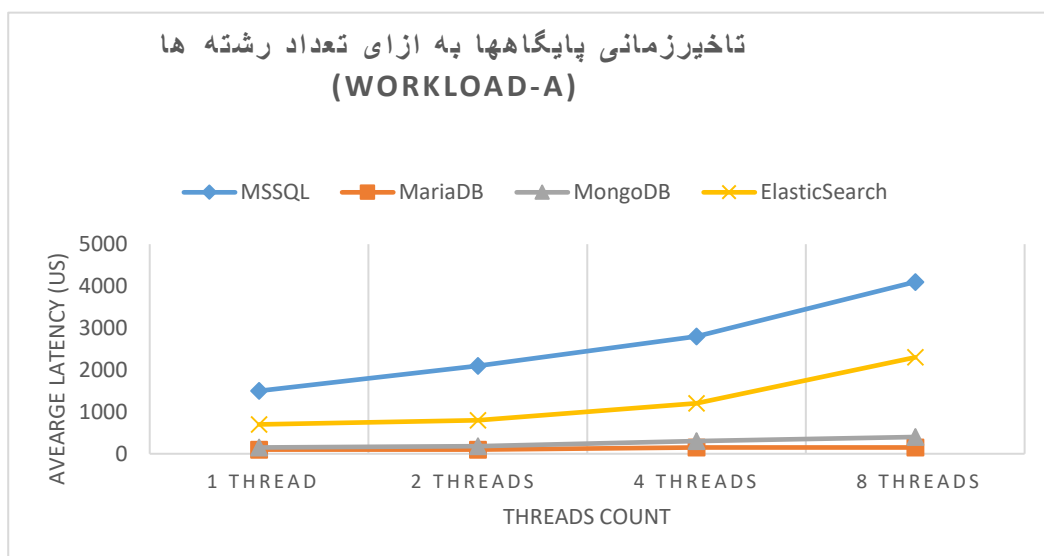
در این بخش عملکرد پایگاههای داده رابطه ای و غیر رابطه ای برای تعداد رشته های مختلف و به ازای تعداد رکورد/عملیات ده هزار بررسی شد و تعداد رشته ای که به ازای آن رفتار مناسب تری مشاهده می شود را تشخیص می دهیم.

شکلهای ۴-۱ و ۴-۲ مربوط به نتایج حاصل از اجرای WorkloadA بر روی تمامی پایگاه های داده رابطه ای و غیر رابطه ای مورد بحث، برای تعداد رکورد/عملیات ده هزار و به ازای تعداد رشته های متفاوت می باشند. تعداد رشته مناسب تعدادی است که به ازای آن توانش زمانی افزایش یابد و حتی الامکان میانگین تاخیر زمانی افزایش نیابد. همانطور که در شکلهای ۴-۱ و ۴-۲ قابل ملاحظه است در مجموع می توان گفت که به ازای افزایش تعداد رشته تا ۴، شاهد افزایش توانش زمانی بطور قابل ملاحظه ای بوده ایم و به نسبت تعداد رشته ۸، تاخیر زمانی نیز افزایش چشمگیری نداشته است و همانطور که شکلهای ۴-۳ و ۴-۴ نشان می دهند، این موضوع برای WorkloadB نیز صادق است. لازم به ذکر است با توجه به نتایج حاصل، در ادامه نتایج آزمایشات نیز که با بررسی جزئی تری همراه هستند به ازای تعداد ۴ رشته که در کل رفتار مناسب تری نشان داده، از شکلهای ۴-۵ به بعد (تا پایان این گزارش پژوهشی) نمایش داده شده اند.

¹ Thread

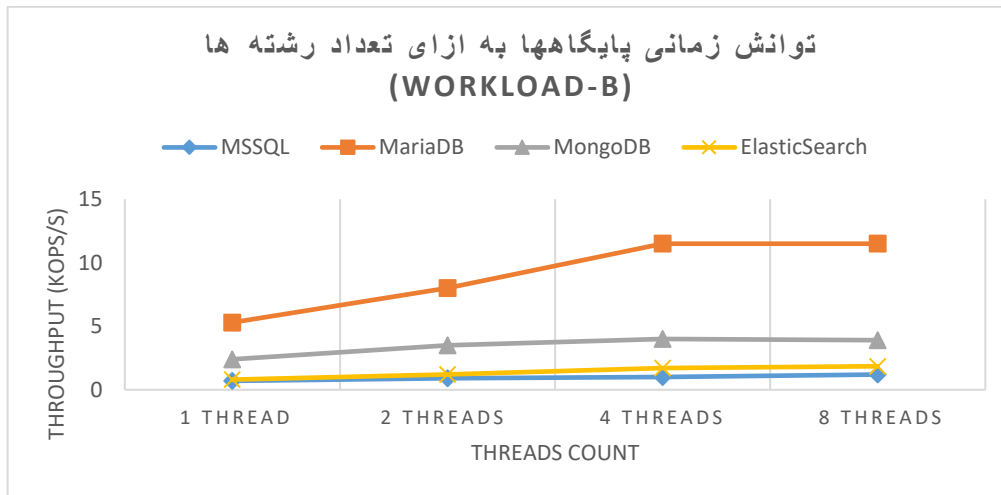


شکل ۴-۱ مقایسه توانش زمانی^۱ بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadA با تعداد رکورد/عملیات ۱۰۰۰۰ و تعداد رشته های مختلف 1,2,4,8

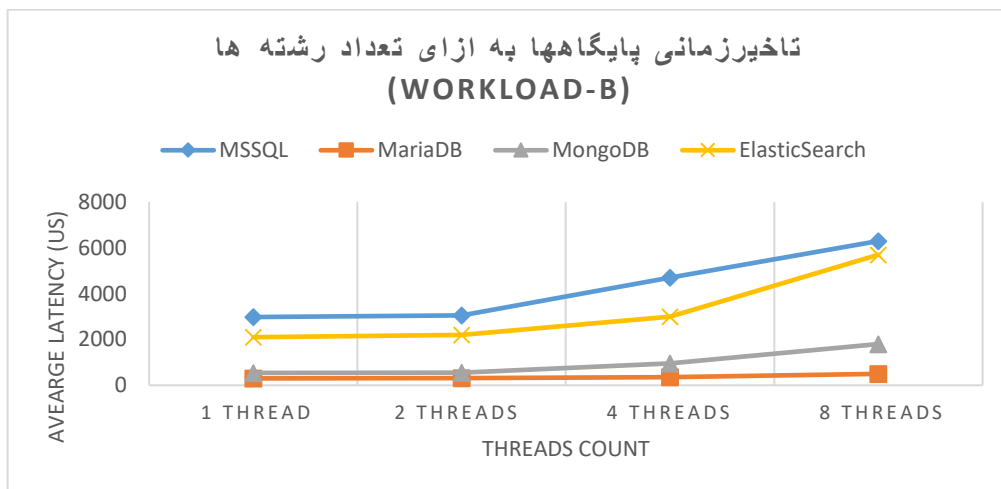


شکل ۴-۲ مقایسه تاخیر زمانی^۲ بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadA با تعداد رکورد/عملیات ۱۰۰۰۰ و تعداد رشته های مختلف 1,2,4,8

^۱Throughput: توانش زمانی یا تعداد عملیات انجام شده در واحد زمان (KOps/Sec: Kilo-Operations/ Second)



شکل ۳-۴ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB با تعداد رکورد/عملیات ۱۰۰۰۰ و تعداد رشته های مختلف 1,2,4,8



شکل ۴-۴ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB با تعداد رکورد/عملیات ۱۰۰۰۰ و تعداد رشته های مختلف 1,2,4,8

از دیدگاه دیگر با بررسی عملکرد کلی سیستمها مشخص است که توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای هر دو WorkloadA و WorkloadB با تعداد رکورد/عملیات ده هزار و تعداد رشته های مختلف 1,2,4,8 به اختصار به شرح زیر است.

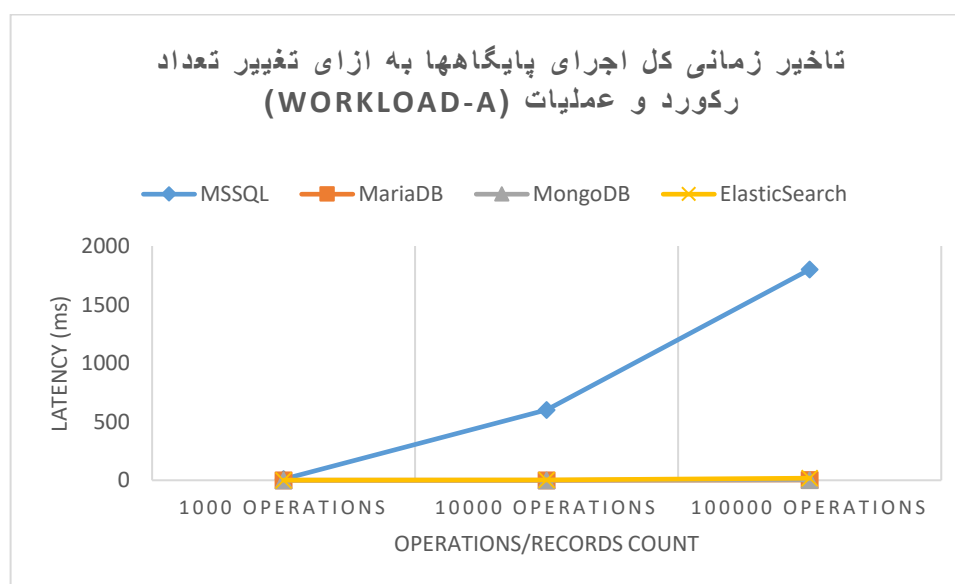
Throughput: MariaDB > MongoDB > ElasticSearch > MSSQL

Average Latency: MariaDB < MongoDB < ElasticSearch < MSSQL

به بیان بهتر توانش زمانی برای MariaDB در کلیه رشته های مختلف در بهترین وضعیت، MongoDB در درجه دوم، ElasticSearch در درجه سوم و MSSQL در درجه آخر قرار دارد. و مسلماً میزان تاخیر زمان بازیابی در MariaDB در بهترین وضعیت با کمترین میزان تاخیر در بازیابی است، MongoDB در رده دوم و ElasticSearch در درجه بعدی و MSSQL در بدترین حالت و با بیشترین میزان تاخیر در بازیابی در آخرین رده قرار می گیرد.

۲-۲-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای با افزایش تعداد عملیات/رکورد-نرخ افزایش زمان کل اجرای بارکاری

در این بخش عملکرد پایگاههای داده رابطه ای و غیر رابطه ای به ازای تعداد رکورد/عملیات ۱۰۰۰، ۱۰۰۰۰ و ۱۰۰۰۰۰ با در نظر گرفتن نرخ افزایش زمان کل اجرای بارکاری مورد بررسی قرار می گیرد. با توجه به رابطه معکوس بین توانش زمانی و تاخیر زمانی از هم اکنون به بعد در نتایج آزمایشات گزارش شده صرفا به نمایش نمودار تاخیر زمانی که در واقع یک پارامتر بسیار مهم در بررسی عملکرد پایگاههای داده (با توجه به سرعت بازیابی اطلاعات^۱ در حین جستجو) خواهیم پرداخت.



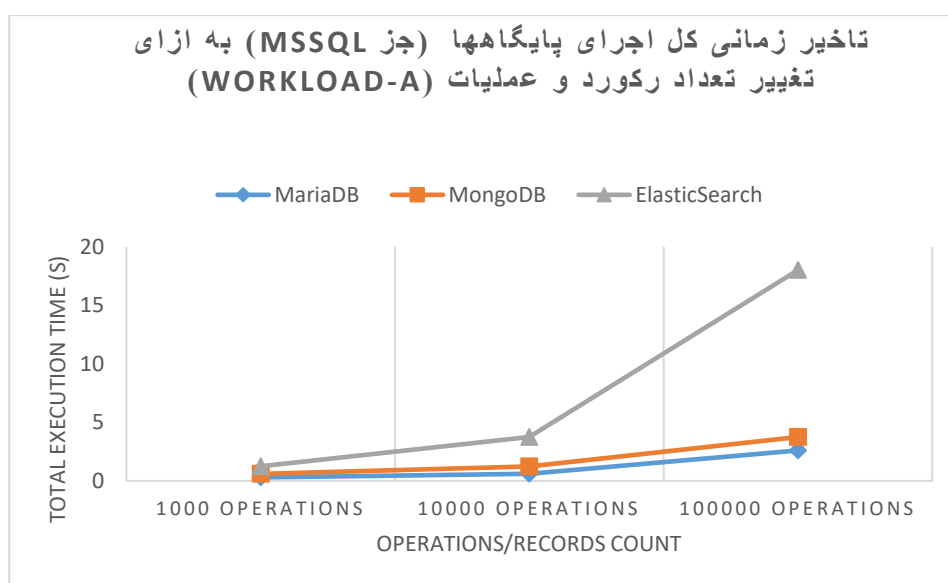
شکل ۴-۵ زمان کل اجرای^۲ WorkloadA بر روی پایگاههای داده رابطه ای و غیر رابطه

ای با تعداد رکورد/عملیات ۱۰۰۰، ۱۰۰۰۰ و ۱۰۰۰۰۰ (الف)

¹ Information Retrieval

² Total Execution Time

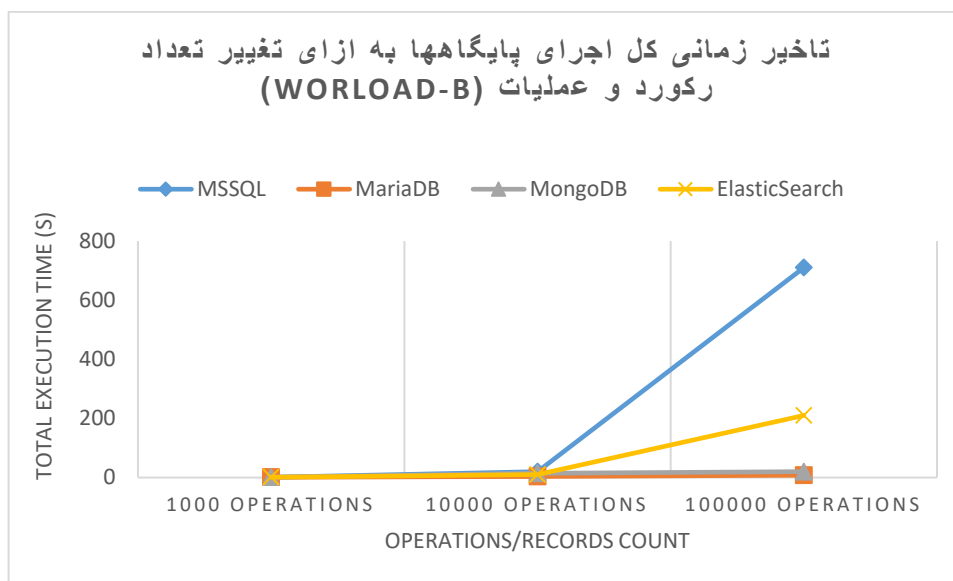
شکل ۴-۵ رفتار پایگاه های داده مورد بحث را در برابر افزایش تعداد عملیات رکورد برای WorkloadA (خوانش ۱۰۰٪) نشان می دهد. نکته ی قابل توجه در این نمودار نرخ افزایش شدید زمان کل اجرای عملیات با افزایش ده هزار به صد هزار عملیات/رکورد برای SQL Server را نشان می دهد، تا حدی که بررسی رفتار سایر پایگاه های داده و مقایسه آن ها با یکدیگر را دشوار ساخته است. به همین علت در شکل ۴-۶ برای پرداختن به این مقایسه SQL Server را حذف کرده است.



شکل ۴-۶ تأخیر زمانی کل اجرای WorkloadA بر روی پایگاههای داده رابطه ای و غیر رابطه ای با تعداد رکورد/عملیات ۱۰۰۰، ۱۰۰۰۰، و ۱۰۰۰۰۰ (ب)

همانگونه که در شکل ۴-۶ قابل مشاهده است، نرخ افزایش زمان اجرای بارکاری برای پایگاه های داده MariaDB و MongoDB تقریباً مشابه است اما برای Elasticsearch این نرخ به طرز قابل توجهی بیشتر است. در مجموع می توان گفت در برابر افزایش تعداد رکورد/عملیات برای WorkloadA (۱۰۰٪ خوانش)، بدترین رفتار مربوط به پایگاه داده SQL Server و در درجه بعدی

مربوط به Elasticsearch می‌باشد و پایگاه‌های داده MongoDB و MariaDB به نسبت این دو رفتار با ثبات تری از خود نشان می‌دهند. همین موضوع همچنان که در شکل ۴-۷ قابل مشاهده است برای WorkloadB (۱۰۰٪ اسکن) نیز صدق می‌کند.



شکل ۴-۷ زمان کل اجرای WorkloadB بر روی پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای با تعداد رکورد/عملیات ۱۰۰۰، ۱۰۰۰۰، و ۱۰۰۰۰۰

۴-۲-۳ بررسی عملکرد پایگاه‌های داده رابطه‌ای و غیر رابطه‌ای با افزایش تعداد عملیات/رکورد - توانش زمانی و تاخیر زمانی

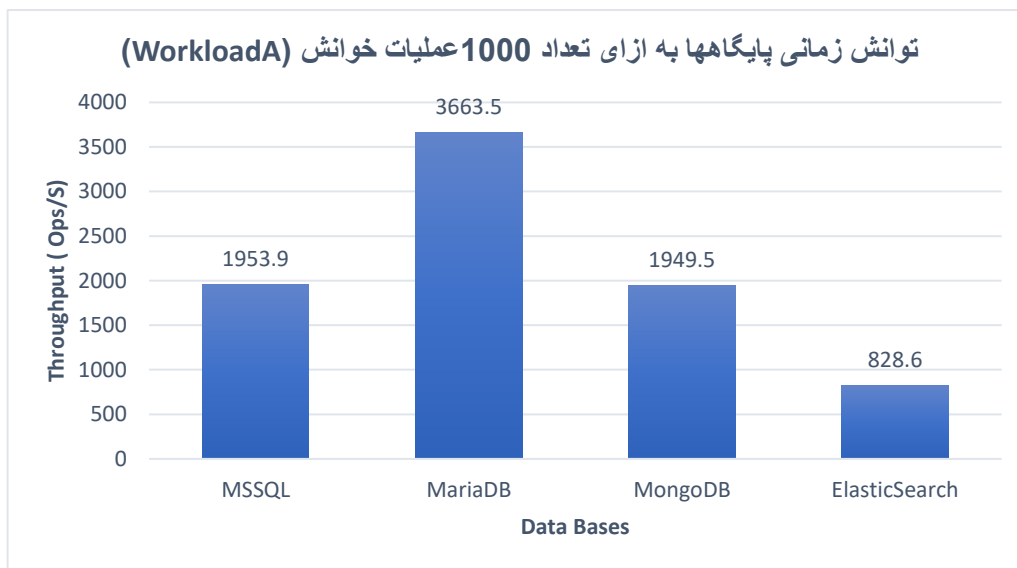
در ادامه به صورت جزئی‌تر رفتار پایگاه‌های داده مورد بحث را به ازای تعداد رکورد/عملیات ۱۰۰۰، ۱۰۰۰۰، و ۱۰۰۰۰۰ در قیاس با یکدیگر با در نظر گرفتن توانش زمانی و تاخیر زمانی می‌سنجیم:

همانطور که دو شکل ۴-۸ و ۴-۹ نتایج حاصل از WorkloadA را نشان می دهند بهترین رفتار مربوط به پایگاه داده MariaDB است که بالاترین توانش زمانی و کمترین تاخیر زمانی را ثبت کرده است و در مقابل آن پایگاه داده Elasticsearch قرار می گیرد که کمترین توانش زمانی و بالاترین تاخیر زمانی را ثبت کرده است. آمار ثبت شده توسط پایگاه های داده MongoDB و SQL Server نزدیک به هم می باشد با این تفاوت که MongoDB آمار بهتری را برای Latency ثبت کرده است. (خصوصاً برای ۹۹ امین صدک تاخیر^۱ که مربوط به ۱ درصد درخواست هایی است که بدترین تاخیر زمانی را تجربه کردند و این نشان از ثبات بیشتر MongoDB دارد).

دو شکل ۴-۱۰ و ۴-۱۱ که نتایج حاصل از WorkloadB را نشان می دهند نیز کلیات تحلیل مربوط به WorkloadA صادق است با این تفاوت که با مقایسه دقیق تر عملکرد پایگاه های داده MongoDB و SQL Server می توان دریافت که برتری MongoDB نسبت به SQL Server در WorkloadB (۱۰۰٪ اسکن) مشهود تر است.

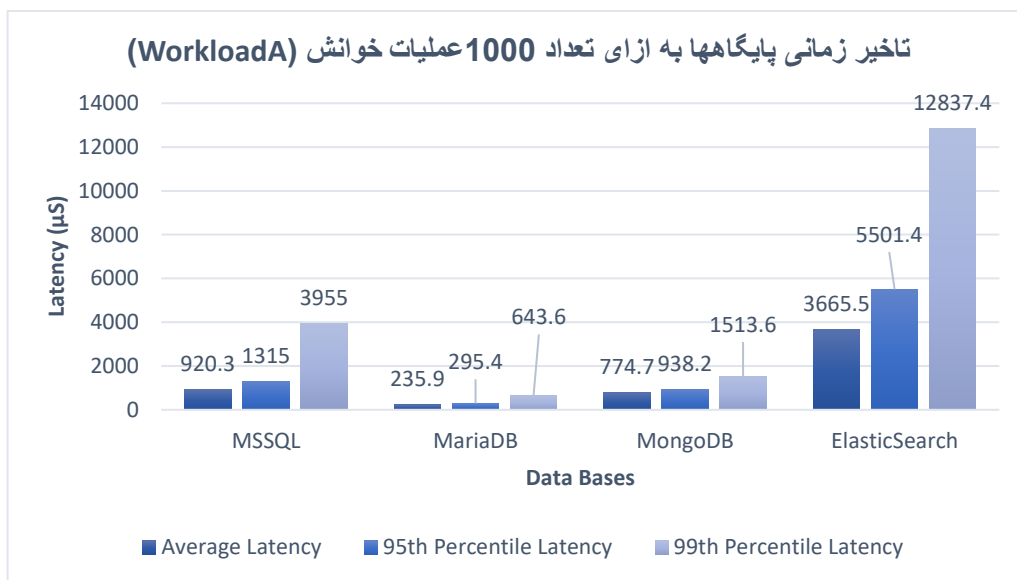
^۱ 99th Percentile Latency

تعداد عملیات/رکورد ۱۰۰۰



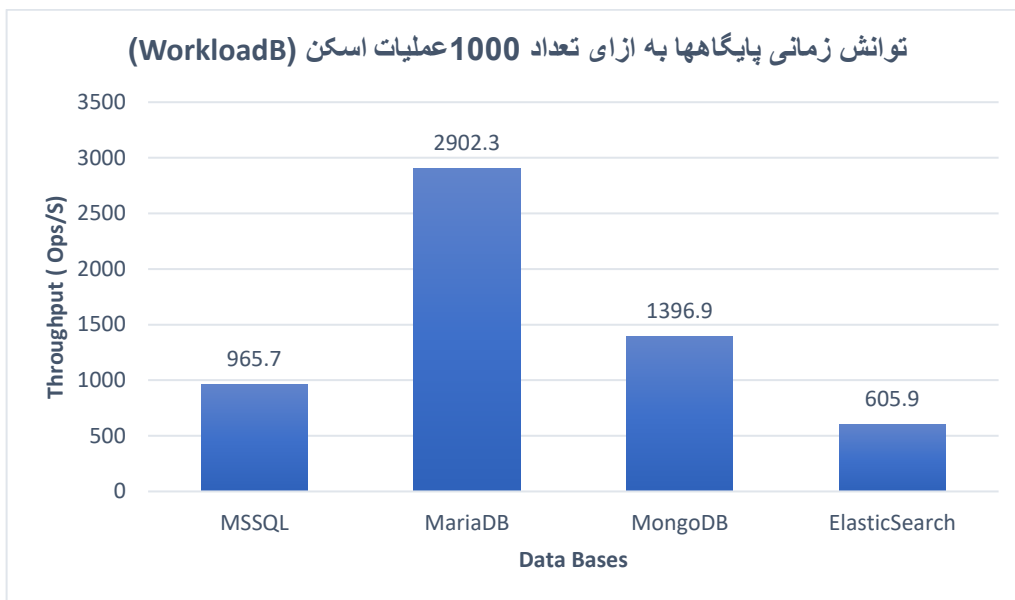
شکل ۴-۸ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadA

با تعداد رکورد/عملیات ۱۰۰۰



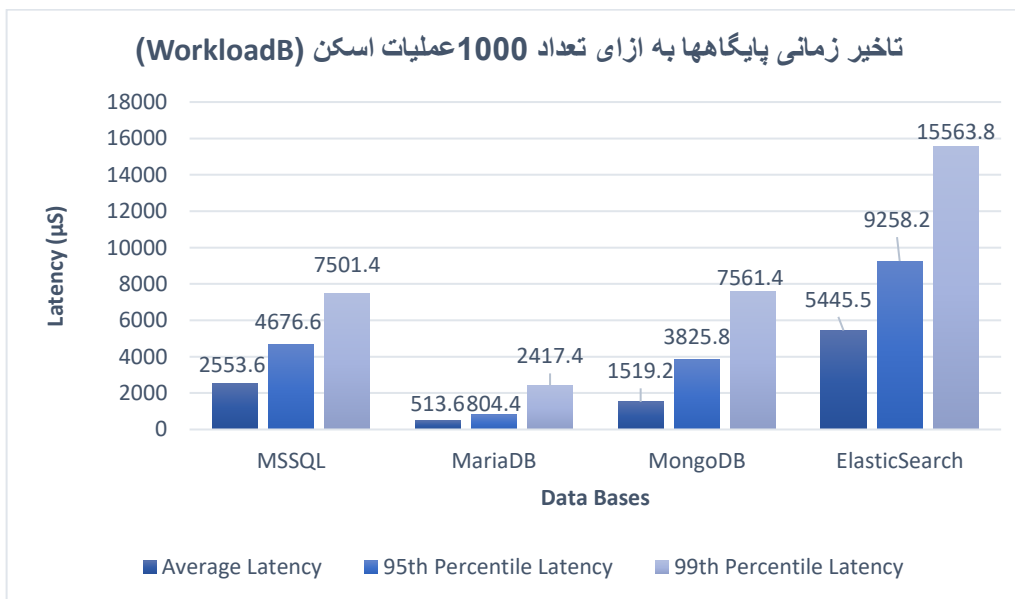
شکل ۴-۹ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadA

با تعداد رکورد/عملیات ۱۰۰۰



شکل ۴-۱۰ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای

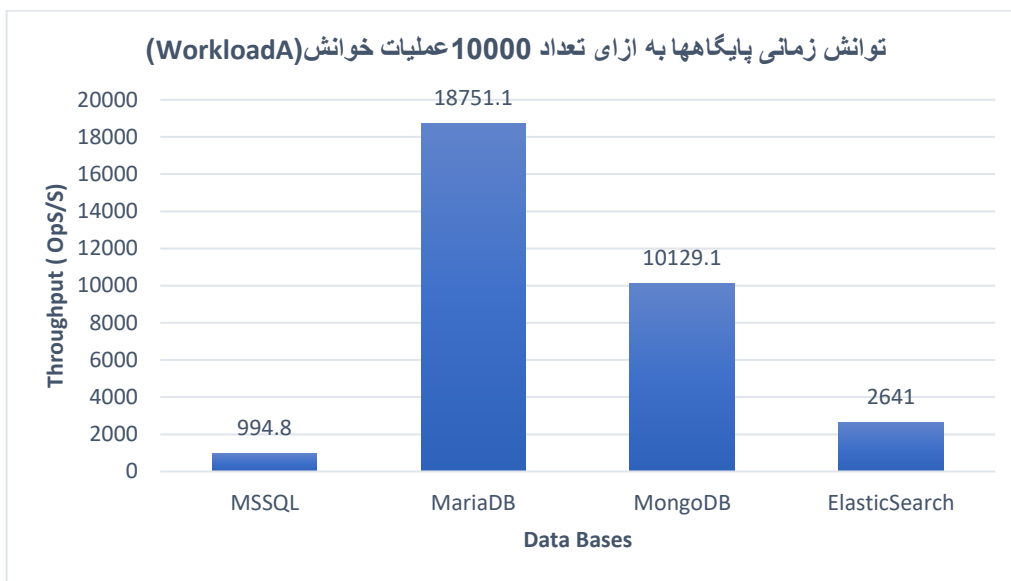
WorkloadB با تعداد رکورد/عملیات ۱۰۰۰



شکل ۴-۱۱ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB

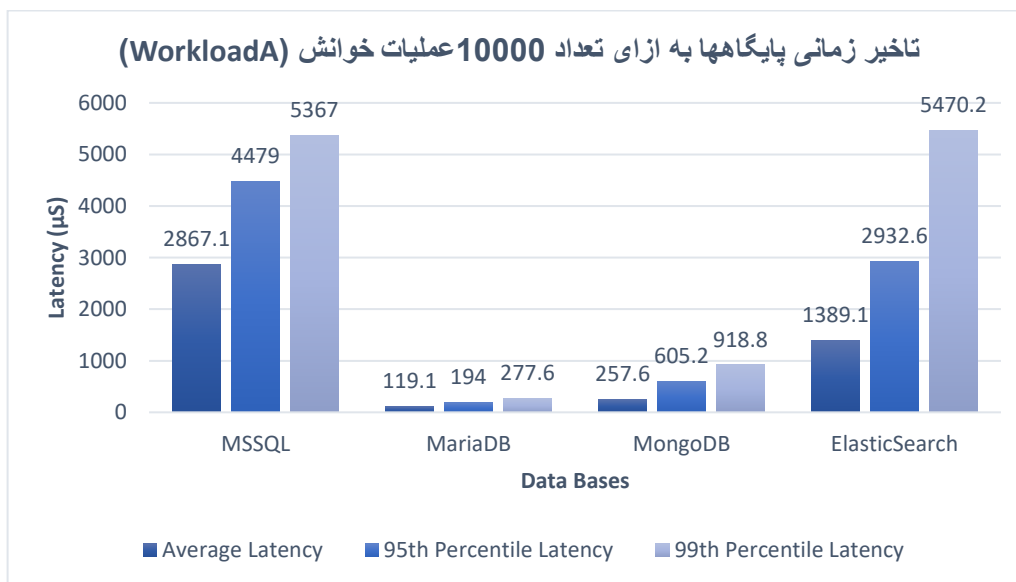
با تعداد رکورد/عملیات ۱۰۰۰

تعداد عملیات/رکورد ۱۰۰۰۰



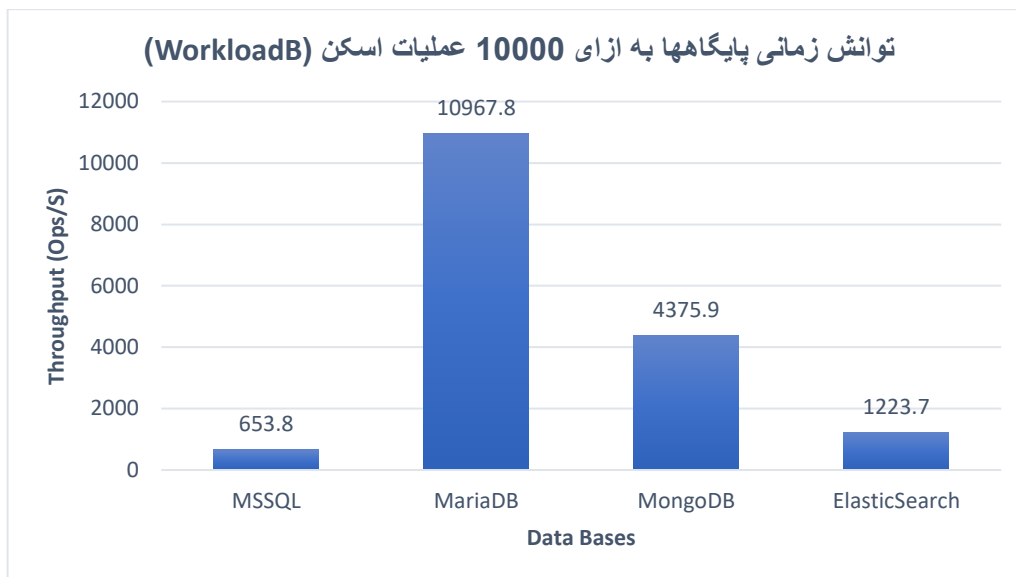
شکل ۴-۱۲ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای

WorkloadA با تعداد رکورد/عملیات ۱۰۰۰۰



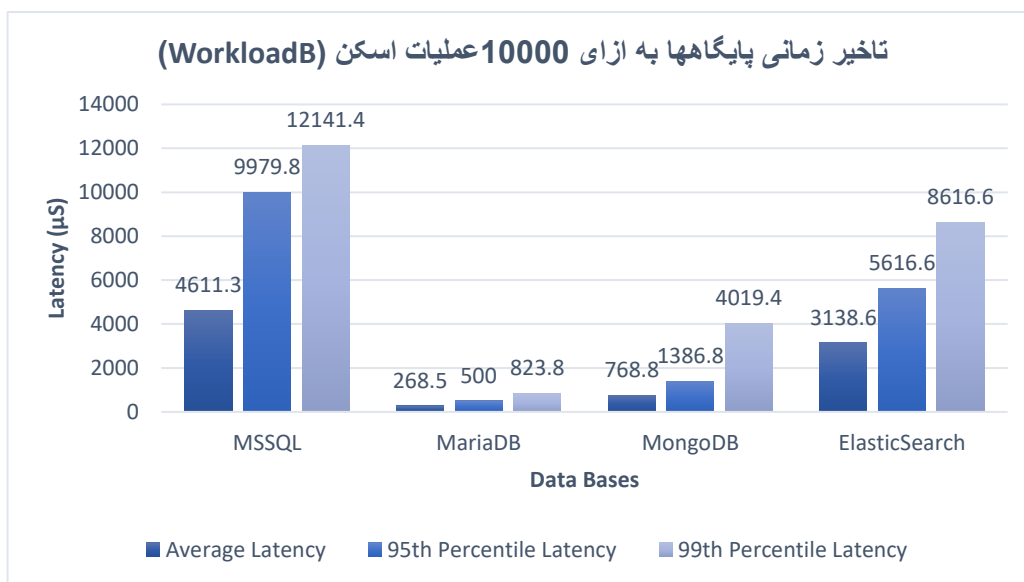
شکل ۴-۱۳ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadA

با تعداد رکورد/عملیات ۱۰۰۰۰



شکل ۴-۱۴ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای

WorkloadB با تعداد رکورد/عملیات ۱۰۰۰۰



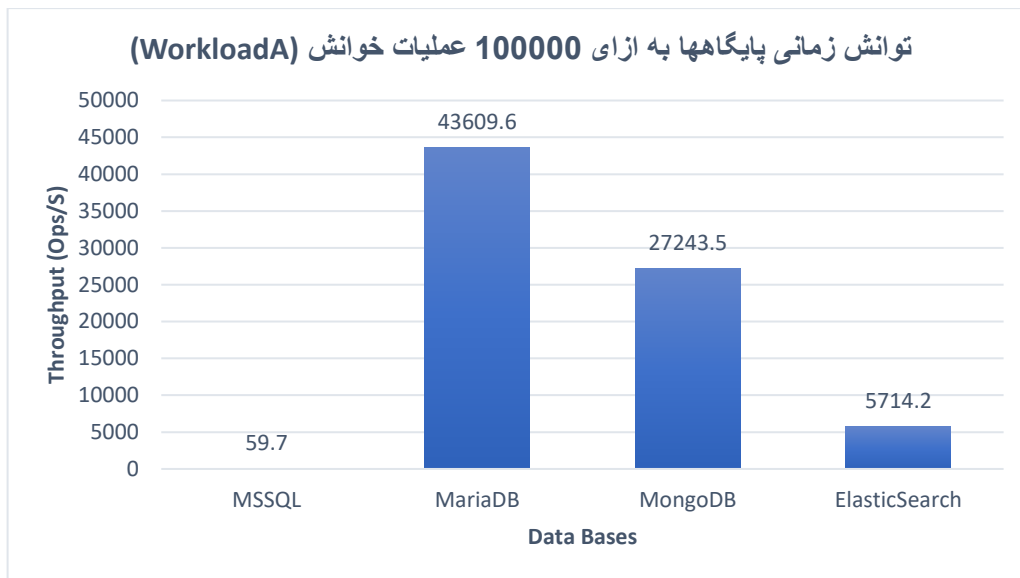
شکل ۴-۱۵ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای

WorkloadB با تعداد رکورد/عملیات ۱۰۰۰۰

با مقایسه شکل های ۱۲-۴ و ۱۳-۴ با شکل های ۸-۴ و ۹-۴ می توان دریافت که با افزایش تعداد عملیات/رکورد ها از ۱۰۰۰ به ۱۰۰۰۰، SQL Server افت شدیدی داشته و به نسبت آن Elasticsearch عملکرد بهتری از خود نشان می دهد و برخلاف SQL Server، توانش زمانی آن افزایش و تاخیر زمانی آن کاهش یافته است. بنابراین با افزایش تعداد درخواست ها به ده هزار در این قسمت، بدترین رفتار مربوط به SQL Server است. بهترین عملکرد کماکان مربوط به پایگاه داده MariaDB است و پس از آن پایگاه داده MongoDB قرار می گیرد. با مقایسه شکل های ۴-۱۴ و ۴-۱۵ که (مربوط به WorkloadB) با شکل های ۴-۱۰ و ۴-۱۱ نیز همین نتایج حاصل می شود.

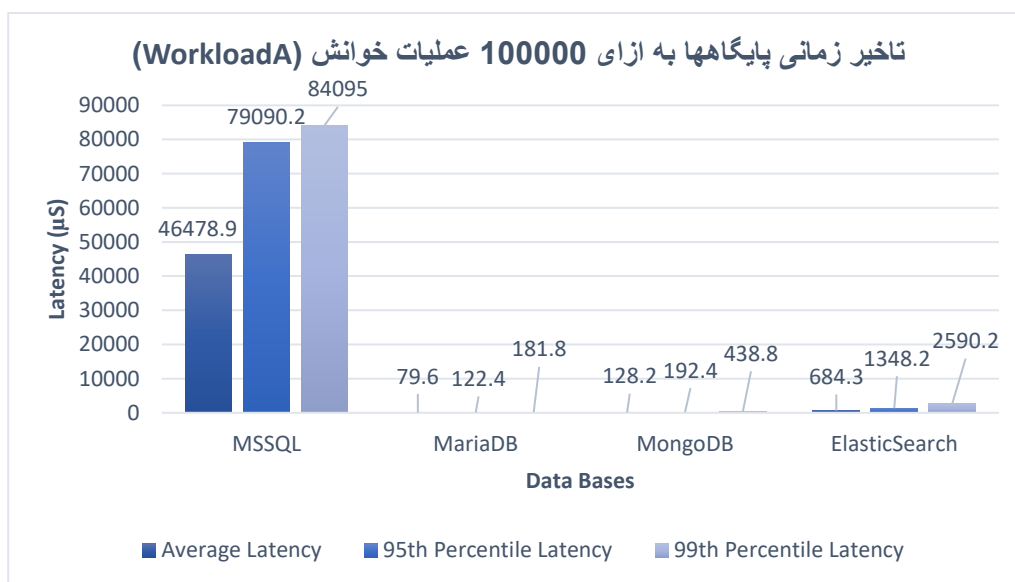
شکل های ۴-۱۶ تا ۴-۱۹ نشان می دهند که با افزایش تعداد عملیات/رکورد ها به ۱۰۰۰۰۰ روند تغییراتی که در قسمت قبلی (مربوط به تعداد عملیات/رکورد ها به ۱۰۰۰۰) توضیح داده شد در این قسمت نیز تکرار می شود و SQL Server به روند نزولی عملکرد خود را با شدت بیشتری ادامه می دهد.

تعداد عملیات/رکورد ۱۰۰۰۰۰



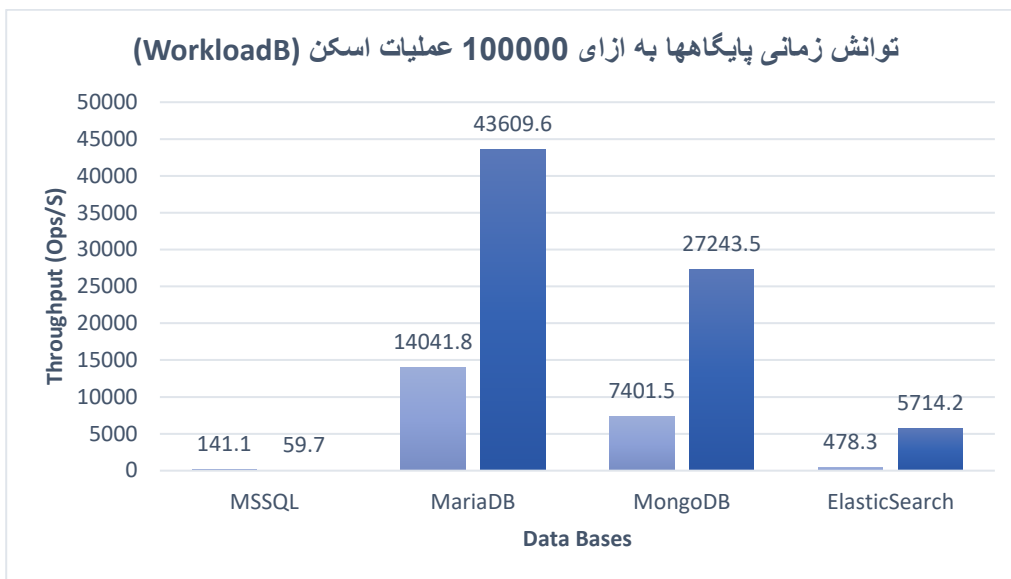
شکل ۴-۱۶ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای

WorkloadA با تعداد رکورد/عملیات ۱۰۰۰۰۰



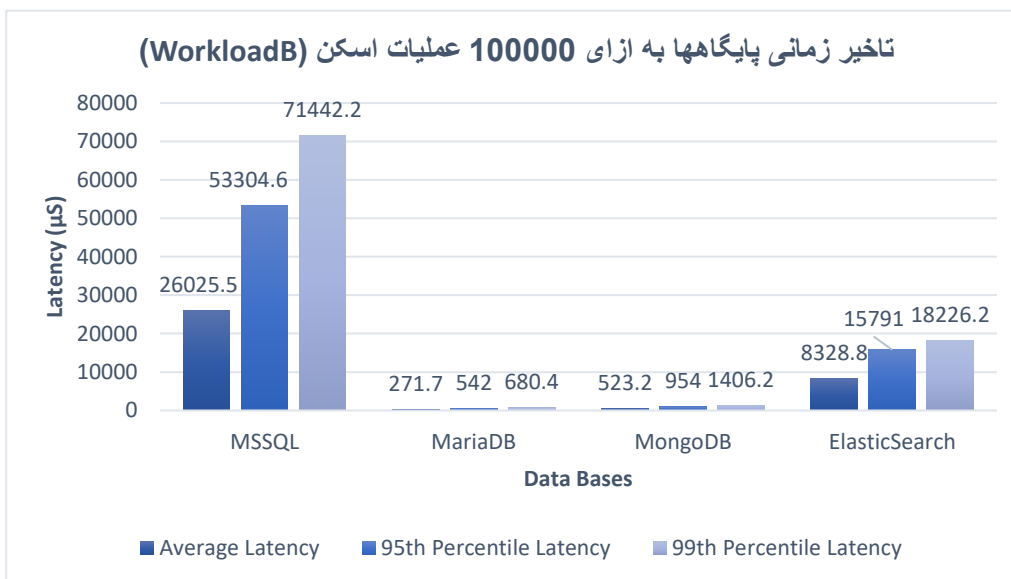
شکل ۴-۱۷ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای

WorkloadA با تعداد رکورد/عملیات ۱۰۰۰۰۰



شکل ۴-۱۸ مقایسه توانش زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB

با تعداد رکورد/عملیات ۱۰۰۰۰۰



شکل ۴-۱۹ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای برای WorkloadB

با تعداد رکورد/عملیات ۱۰۰۰۰۰

۴-۲-۴ جمع بندی نتایج مربوط به معیار ارزیابی سرویس ابری یا هو! (YCSB)

در مقیاسی که پایگاه های داده Elasticsearch و MongoDB، MariaDB، SQL Server را بوسیله YCSB و برای بارهای کاری ۱۰۰٪ خوانش و ۱۰۰٪ اسکن مورد ارزیابی قرار دادیم، حکایت از برتری چشمگیر MongoDB و MariaDB نسبت به دو پایگاه داده دیگر داشت. همچنین مشاهده شد که با افزایش تعداد عملیات/رکورد ها SQL Server برتری خود را نسبت به Elasticsearch به طور چشمگیری از دست می دهد و در بین ۴ پایگاه داده مورد بحث ضعیف ترین عملکرد را به خود اختصاص می دهد.

۴-۳ بررسی پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی تمام متن

در این قسمت به بررسی عملکرد جستجوی تمام متن برای پایگاههای داده SQL Server، MariaDB، MongoDB و Elasticsearch می پردازیم. برای ارزیابی این عملکرد چهار نوع پرس و جوهای آزمایشی در نظر گرفته شده است که به شرح زیر می باشند:

- Single: به جستجوی یک کلمه مشخص در پایگاه داده می پردازد.
- OR: به جستجوی سه کلمه در پایگاه داده می پردازد بدین صورت که هر رکورد که حداقل یک کلمه از سه کلمه را در متن خود داشته باشد جزء نتایج پذیرفته شود.
- AND: به جستجوی سه کلمه در پایگاه داده می پردازد بدین صورت که هر رکورد پذیرفته شده بایستی هر سه کلمه را در متن خود داشته باشد.
- Exact Phrase: به جستجوی یک عبارت مشخص (که متشکل از سه کلمه می باشد) در پایگاه داده می پردازد.

برای ساخت یک دیتاست با عنوان NoSQL-Dataset، از روش نمونه گیری تصادفی از ۶ میلیون چکیده لاتین موجود در مخزن داده های مرکز منطقه ای استفاده شده است. شایان ذکر است که در سال اخیر رشد چشمگیر بالغ بر ۵۰٪ی در تعداد رکوردهای لاتین در دو پایگاه مقالات انگلیسی^۱ و مقالات انگلیسی مهندسی^۲ شاهد بودیم. با توجه به خصوصیت جستجوی تمام متن، هر رکورد به صورت تقریبی حدود ۳۰۰ کلمه است و داده های مورد بررسی بالغ بر یک میلیارد و هشتصد میلیون می باشد که داده حجیم محسوب می شود.

¹ English Articles

² Engineering English Articles

به منظور بررسی عملکرد جستجوی تمام متن در انواع پرس و جوهای آزمایشی^۱ فوق؛ برای پرس و جوی آزمایشی single هزار کلمه، برای کوئیری OR هزار گروه سه کلمه ای، برای پرس و جوهای آزمایشی AND هزار گروه سه کلمه ای دیگر و برای پرس و جوهای آزمایشی Exact Phrase هزار عبارت معنادار سه کلمه ای به صورت تصادفی از دیتاست NoSQL-Dataset استخراج شده و بر روی پایگاههای داده MariaDB، SQL Server، MongoDB و Elasticsearch اجرا شده است.

دو معیار برای مقایسه پایگاه های داده انتخاب شده اند معیار اول زمان بازیابی است. این زمان عبارت است از مجموع مدت زمانی که صرف یافتن تمامی رکوردهایی که دارای ویژگی مورد نظر باشند، مرتب کردن آنها بر اساس امتیاز^۲ (relevance score)، ذخیره بیست رکورد اول (بالاترین امتیاز) در یک فایل و محاسبه تعداد کل این رکوردهاست. متوسط زمان بازیابی سه بار اجرای هر یک از پرس و جوهای آزمایشی (۳۰۰۰ پرس و جوهای آزمایشی) به عنوان نتیجه نهایی زمان بازیابی در گزارشات آمده است.

معیار دوم کیفیت بازیابی نتایج است که به مقایسه بهترین نتیجه بازیابی شده در ۴ نوع پرس و جوهای آزمایشی می پردازد. به این منظور روابط پیشنهادی QRRT در این پژوهش در فرمولهای ۴-۱ الی ۴-۴ ارائه شده است که قابلیت مقایسه کیفیت نتایج بازیابی شده بین پایگاه های مختلف با توجه به نوع پرس و جوهای آزمایشی را می دهد که در قسمت ۳،۲،۳،۴ بتفصیل به آن پرداخته شده است.

¹ Query

² Relevance Score

تمام آزمایشات برای پایگاههای داده غیر رابطه ای که قابلیت کلاستر /خوشه بندی^۱ را دارند در سه حالت تک نود^۲، کلاستر سه نود شاردرده^۳ و کلاستر شش نود شاردرده^۴ اجرا و تاثیر قابلیت خوشه بندی که از مزایای پایگاههای داده غیر رابطه ای محسوب می شود گزارش شده است.

گفتنی است پیش از اجرای پرس و جوهای آزمایشی لازم است تا فرآیند ایندکسینگ تمام متن بر روی مجموعه رکوردها در پایگاهها انجام شده باشد. برای اینکه ویژگیهای این فرآیند حتی الامکان در چهار پایگاه داده مورد بحث یکسان باشد، ویژگیهای افزوده ای همچون ریشه یابی^۵ و حذف کلمه های ایستا^۶ از فرآیند نمایه سازی ایندکسینگ، که برای پایگاه دادههای مختلف به صورت های مختلفی اجرا می شود، از تنظیمات پیش فرض پایگاههای داده مورد بحث، حذف شده است.

¹ Sharded Cluster

² Single Node

³ 3-Sharded Cluster

⁴ 6-Sharded Cluster

⁵ Stemming

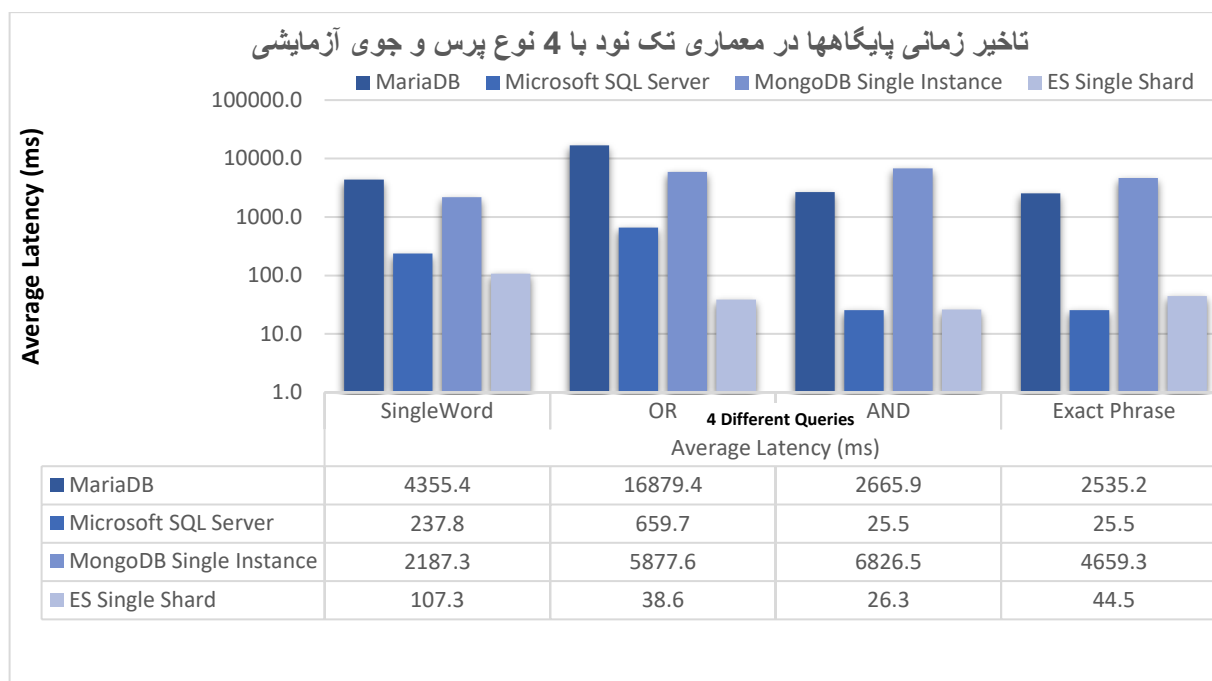
⁶ Stopword

۳-۴-۱ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی تمام متن بدون خوشه بندی (Single Node) با معیار زمان بازیابی

شکل ۴-۲۰ و ۴-۲۱ به مقایسه عملکرد جستجو و بازیابی تمام متن در پایگاههای داده رابطه

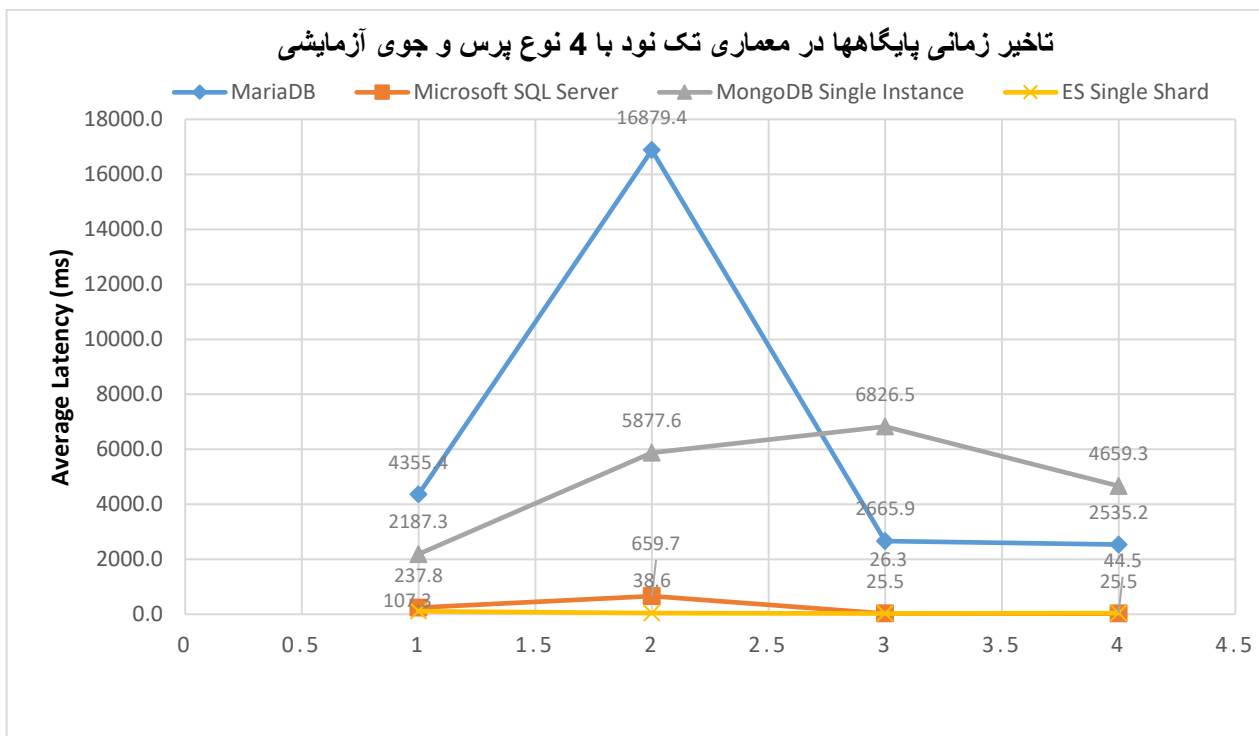
ای و غیر رابطه ای در حالت ساده یک تک نود (single node) یا به اصطلاح بدون خوشه بندی بر

روی چهار نوع پرس و جوهای آزمایشی می پردازد.



شکل ۴-۲۰ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای بدون

خوشه بندی با چهار نوع پرس و جوهای آزمایشی - در مقیاس لگاریتمی



شکل ۴-۲۱ مقایسه تاخیر زمانی بر روی پایگاههای داده رابطه ای و غیر رابطه ای بدون خوشه بندی با چهار نوع پرس و جوهای آزمایشی (نمایش مقایسه ای دیگری از نوسانات داده های ارائه شده در شکل ۴-۲۰ در حالت غیر لگاریتمی)

شکل ۴-۲۰ و ۴-۲۱ به مقایسه عملکرد جستجو و بازیابی تمام متن در پایگاههای داده رابطه ای و غیر رابطه ای در حالت ساده یک تک نود (single node) یا به اصطلاح بدون خوشه بندی بر روی چهار نوع پرس و جوهای آزمایشی می پردازد.

در ابتدا نکته قابل توجه این شکل، طولانی تر بودن اجرای پرس و جوهای آزمایشی نوع OR است. این موضوع نشان می دهد که هر چه تعداد رکورد های نتیجه پرس و جوهای آزمایشی بیشتر باشد، سربار محاسبه امتیاز آن ها و مرتب کردن آن ها بر اساس این امتیاز بیشتر خواهد بود و همین امر باعث می شود زمان اجرای پرس و جوهای آزمایشی نوع OR در این پایگاههای داده به مراتب بیشتر از سایر پرس و جوهای آزمایشی باشد.

مقایسه عملکرد جستجو و بازیابی تمام متن در پایگاه‌های داده رابطه‌ای Microsoft SQL Server و MariaDB بدون خوشه بندی حاکی از آن است که SQL Server در هر چهار نوع پرس و جوهای آزمایشی با برتری بسیار قابل چشمگیری قادر به ثبت میانگین تاخیر کمتری نسبت به MariaDB بوده است و بعبارت دیگر از نظر سرعت عملکردی کارایی بهتری را از خود نشان داده است.

به همین نحو مقایسه عملکرد جستجو و بازیابی تمام متن در پایگاه‌های داده غیررابطه‌ای ElasticSearch و MongoDB بدون خوشه بندی حاکی از آن است که ElasticSearch در هر چهار نوع پرس و جوهای آزمایشی با برتری بسیار قابل چشمگیری قادر به ثبت میانگین تاخیر کمتری نسبت به MongoDB بوده است و بعبارت دیگر از نظر سرعت عملکردی کارایی بهتری را از خود نشان داده است.

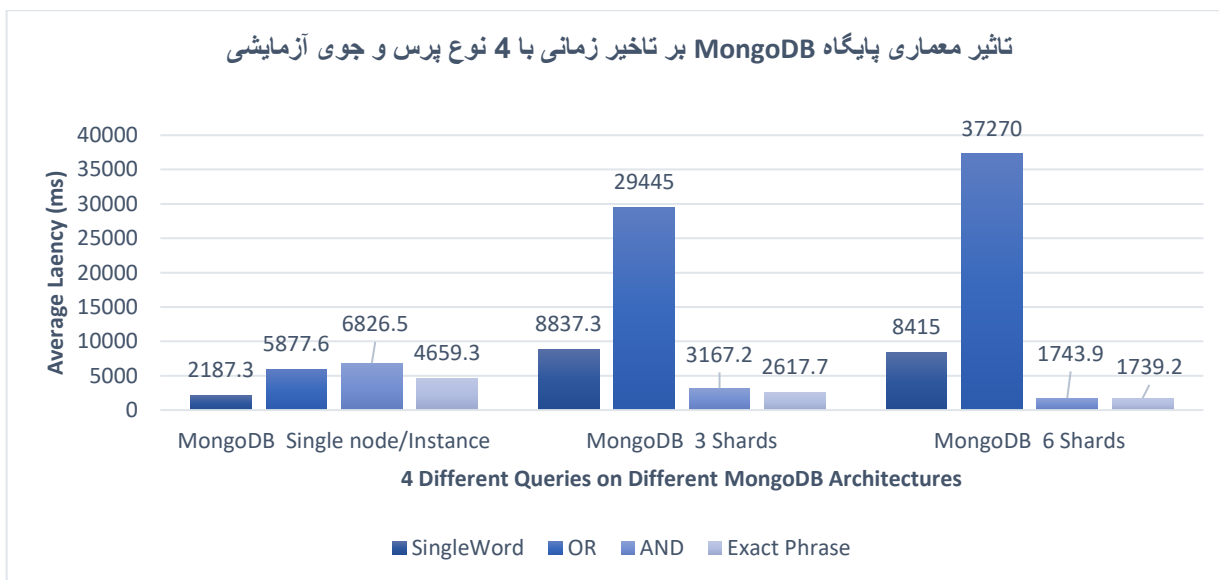
همانطور که در شکل‌های ۴-۲۰ و ۴-۲۱ نمایش داده شده است، بطور کلی در حالتی که معماری بدون خوشه بندی مطرح است، در اکثر مواقع در هر چهار نوع پرس و جوهای آزمایشی برتری عملکرد بسیار قابل توجه ElasticSearch و Microsoft SQL Server نسبت به دو نمونه دیگر پایگاه داده MariaDB و MongoDB مشهود است.

۲-۳-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی تمام متن با قابلیت خوشه بندی (Sharded Cluster)

در این قسمت به بررسی دقیقتر تاثیر اضافه نمودن شاردینگ بر دو پایگاه داده غیررابطه ای MongoDB و ElasticSearch که از قابلیت‌های مهم آنهاست می پردازیم. بدین منظور دو حالت 3 Shard و 6 Shard در نظر گرفته می شود و در قیاس با حالت تک نود^۱ نتایج بررسی ها برای هر ۴ نوع پرس و جوهای آزمایشی اعلام میشود.

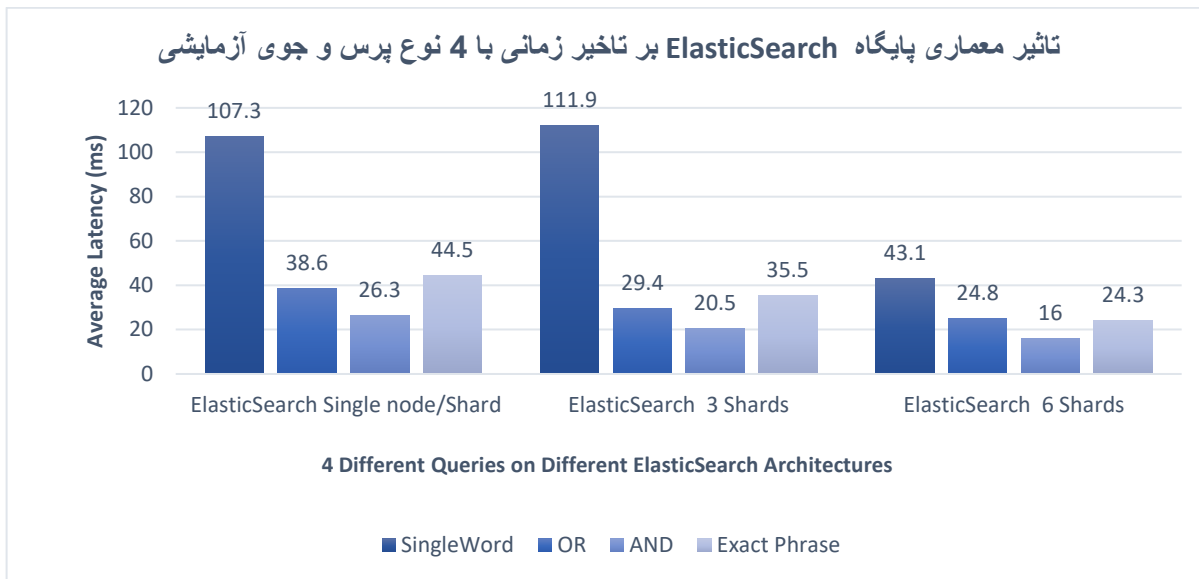
۱-۲-۳-۴ بررسی عملکرد پایگاههای داده غیر رابطه ای با در نظر گرفتن قابلیت خوشه بندی (Sharded Cluster) به تفکیک معماری با معیار زمان بازیابی

در این قسمت به بررسی تاثیر شاردینگ در عملکرد جستجوی تمام متن Full-Text Search برای پایگاه داده غیررابطه ای MongoDB و ElasticSearch می پردازیم.



شکل ۲۲-۴ مقایسه تاخیر زمانی پایگاه داده غیر رابطه ای MongoDB با در نظر گرفتن شاردینگ با چهار نوع پرس و جوهای آزمایشی

¹ Single Node/ Single Instance /Single Shard



شکل ۴-۲۳ مقایسه تاخیر زمانی پایگاه داده غیر رابطه ای ElasticSearch با در نظر گرفتن

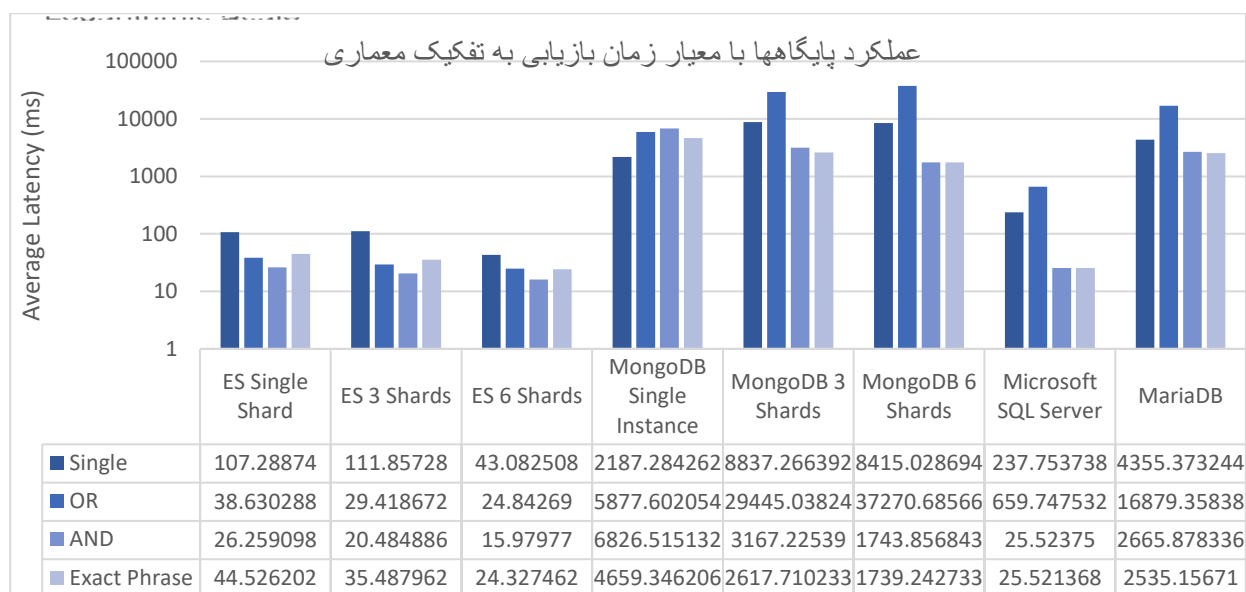
شاردینگ با چهار نوع پرس و جوی آزمایشی

شکل ۴-۲۲ به بررسی تاثیر افزایش تعداد Shard ها بر عملکرد جستجو و بازیابی تمام متن پایگاه داده MongoDB می پردازد. همانطور که قابل ملاحظه است با افزایش تعداد Shard ها عملکرد پرس و جوهای آزمایشی AND و Exact Phrase بهبود پیدا می کند و در نقطه مقابل عملکرد پرس و جوهای آزمایشی Single و OR از نظر زمان بازیابی افت می کند.

اما با توجه به شکل ۴-۲۳ می توان دریافت که این افزایش تعداد Shard بطور کلی برای هر چهار نوع پرس و جوهای آزمایشی عملکرد پایگاه داده Elasticsearch را بهبود می دهد. بطور خلاصه می توان نتیجه گرفت افزایش تعداد shardها در MongoDB تاثیر منفی در کارایی و عملکرد خواهد گذاشت و این در حالیست که برای ElasticSearch افزایش shardها منجر به بهبود وضعیت عملکرد و کارایی خواهد شد.

۲-۲-۳-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای در جستجو و بازیابی تمام متن به تفکیک معماری با معیار زمان بازیابی

در این قسمت به نمایش نمودارهای مقایسه ای کلی همه حالات معماری (پیاده سازی شده در آزمایشات) در یک نگاه می پردازیم تا بتوان به جمع بندی بهتری در خصوص انتخاب وضعیت بهینه پایگاه داده برگزیده دست یافت.



شکل ۴-۲۴ مقایسه زمان بازیابی پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی

تمام متن با چهار نوع پرس و جوهای آزمایشی به تفکیک معماری (لگاریتمی)

شکلهای ۴-۲۴ (که برای مشاهده بهتر به دو حالت ساده و لگاریتمی نیز نمایش داده شده است)، جمع بندی کلیه آزمایشات برای معیار زمان بازیابی است. مشاهده می شود SQL Server و ElasticSearch عملکرد بهتری را در زمان بازیابی جستجوی تمام متن در داده

های مرکز منطقه ای داشته اند. جمع بندی وضعیت پایگاههای داده در تمام حالات (انواع پرس و جوهای آزمایشی و شیوه های مختلف خوشه بندی¹ Sharded-Cluster) از نظر زمان بازیابی به شرح زیر قابل خلاصه شدن است.

Latency: ElasticSearch < SQL < MongoDB ~ MariaDB

در واقع همانطور که مشخص است، معیار زمان بازیابی در ElasticSearch دارای بهترین وضعیت (سریعترین زمان بازیابی با کمترین تاخیر) و سپس در رده دوم SQL و در نهایت MongoDB و MariaDB در جایگاه مشابهی در بدترین وضعیت با بیشترین زمان بازیابی (کندترین زمان بازیابی با بیشترین تاخیر) در جستجوی تمام متن قرار دارند

۳-۲-۳-۴ بررسی عملکرد پایگاههای داده رابطه ای و غیر رابطه ای در جستجوی تمام متن به تفکیک معماری با معیار کیفیت بازیابی داده

یکی از اهداف این پژوهش بررسی کیفیت داده های بازیابی شده می باشد؛ بدین منظور معیار کیفیت بازیابی داده ها QRRT² که تلفیقی از زمان بازیابی و کیفیت بازیابی داده است را به صورت زیر تعریف نمودیم؛ تا بتوان به جمع بندی بهتری در خصوص انتخاب کیفی پایگاه داده برگزیده دست یافت.

فرمول ۴-۱ $QRRT_{and, or} = QRR_{and, or} / (\text{زمان بازیابی})$

فرمول ۴-۲ $QRRT_{single, exact} = QRR_{single, exact} / (\text{زمان بازیابی})$

¹ Sharded Cluster

² Quality of Retrieved Result per Time

فرمول ۴-۳ (تعداد کل کلمات سند^۳ بازیابی شده) / $QRR^{1}_{or, and} = NAWC^2$

فرمول ۴-۴ (تعداد کل) / (تعداد کلمات سند شده در مدرک) $QRR_{single, exact}$

(کلمات سند بازیابی شده)

NAWC میانگین نرمال شده تعداد کلمات بازیابی شده در نتیجه پرس و جوهای آزمایشی AND یا OR می باشد. برای نرمال نمودن تعداد کلمات بر بیشینه تعداد بازیابی ها تقسیم میشود. مثال زیر چگونگی محاسبه NAWC را نشان می دهد:

Query = *receive* OR *proposed* OR *source*

Number of each word in the best result:

	<i>receive</i>	<i>proposed</i>	<i>source</i>
Count	2	5	4

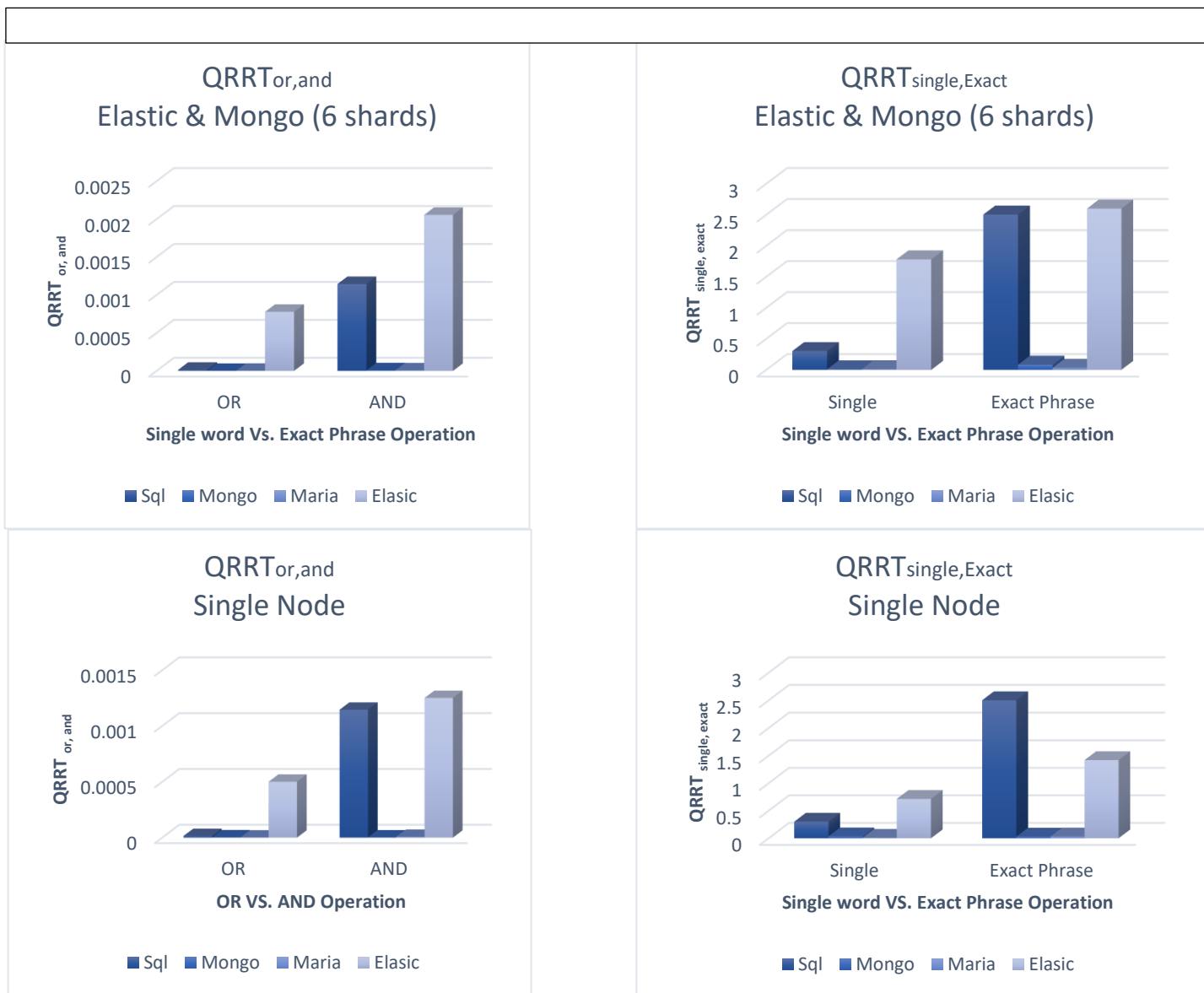
$$NAWC = \left(\frac{2+4+5}{5} \right) / 3$$

افزایش مقادیر عددی برای QRRT نشان دهنده عملکرد بهتر پایگاههای داده مورد بحث می باشد. در ادامه به نمایش نمودارهای مقایسه ای با معیار کیفیت بازیابی داده ها میپردازیم تا بتوان به جمع بندی بهتری در خصوص انتخاب وضعیت بهینه پایگاه داده برگزیده دست یافت.

¹ Quality of Retrieved Result

² Normalized Average of Word Count

³ Document



شکل ۴-۲۵ مقایسه کیفیت بازیابی داده ها (QRRT)^۱ در پایگاههای داده رابطه ای و غیر

رابطه ای در جستجوی تمام متن

QRRT: Quality of Retrieved Result (برای درک مفهوم " معیار کیفیت بازیابی داده" و سایر موارد مندرج در نمودار بالا به شرح توضیحات فرمولهای آن در بخش ۳-۲-۳-۴ مراجعه شود)

همانطور که در بخش های شکل ۴-۲۵ نمایش داده شده است، جمع بندی وضعیت پایگاههای داده در تمام حالات (انواع پرس و جوهای آزمایشی و شیوه های مختلف خوشه بندی Sharded-Cluster) از نظر معیار کیفیت بازیابی (QRRT) به شرح زیر قابل خلاصه شدن است.

QRRT: ElasticSearch > SQL > MongoDB ~ MariaDB

در واقع همانطور که مشخص است، معیار کیفیت بازیابی در ElasticSearch دارای بهترین وضعیت و سپس در رده دوم SQL و در رده بعدی MongoDB قرار دارد و نهایتا MariaDB در بدترین وضعیت از نظر معیار کیفیت بازیابی قرار دارد.

۴-۳-۳ جمع بندی نتایج جستجوی تمام متن

در مقیاسی که پایگاه های داده SQL Server، MariaDB، MongoDB و Elasticsearch را بوسیله جستجوی تمام متن مورد ارزیابی قرار دادیم، حکایت از برتری چشمگیر Elasticsearch و Microsoft SQL Server نسبت به دو پایگاه داده دیگر داشت. همچنین مشاهده شد که با افزایش تعداد Shardها در Elastic Search؛ Elasticsearch برتری خود را نسبت به Microsoft SQL Server افزایش می دهد.

فصل پنجم: بحث و نتیجه گیری

۵-۱ مقدمه

همانطور که گفته شد، امروزه داده ها به عنوان یک دارائی ملی شناخته می شوند . همچنین پردازش، تحلیل و استفاده از داده به عنوان یک عامل کلیدی برای رشد کلیه سازمانها تبدیل شده است و موجب مزیت رقابتی در کسب و کارها، محرک نوآوری، افزایش رقابت و اثرات مثبت اجتماعی خواهد شد. مرکز منطقه ای نیز از این قضیه مستثنی نمی باشد و با توجه به افزایش حجم داده ها بصورت تصاعدی (برای مثال در طی سال اخیر رشد ۵۰٪ داده های پایگاه داده آن را شاهد بر بستر پایگاه داده رابطه ای Microsoft SQL Server را شاهد بودیم) و با لحاظ نمودن ماموریت‌های این مرکز در ارائه خدمات از طریق پایگاه‌های داده مقالات تمام متن و موتور جستجوی اختصاصی آن ، نیاز به ارائه راهکارهای نوین پردازشی یک ضرورت محسوب میشود. به بیان بهتر با توجه به چشم اندازهای آینده و رو به رشد مرکز منطقه ای و متناسب با سند راهبردی و همچنین اساسنامه آن استفاده از شیوه های نوین پردازش برای ارائه برنامه ها و خدمات اطلاع رسانی در منطقه تاکید شده است. با توجه به افزایش روزافزون حجم بزرگ داده ها، گسترش کیفی و کمی سرویسهای متعدد ارائه شده به سایر پژوهشگران در سطح ملی و بعضا بین المللی منطقه؛ حرکت به سمت ارتقا زیرساختهای نرم افزاری در قالب محاسبات توزیع شده و استفاده از روشهای نوین پردازش داده ها اهمیت می یابد.

لذا در این پژوهش بر آن شدیم تا برای کسب دانش کار با پایگاههای غیر رابطه ای NOSQL

، حرکت به سمت خوشه بندی (که از نقاط قوت پایگاههای داده غیر رابطه ای است)، و حل مشکلات

احتمالی آتی و پیش رو؛ با بکارگیری سیستمهای NOSQL به عنوان ابزار پردازش کلان داده ها در جستجو و بازیابی تمام متن، به مطالعه مقایسه ای وضعیت عملکردی و کارایی آنها در قیاس با پایگاه داده رابطه ای Microsoft SQL Server مرکز منطقه ای بپردازیم.

فرضیه و سوالاتی که در این پژوهش مدنظر واقع شد و در طی این پژوهش با انجام

آزمایشات در دو بخش اصلی توانستیم پاسخهای آن را دریابیم به شرح زیر است:

• آیا با توجه به حجم داده ها و رشد آتی آنها پایگاههای غیر رابطه ای/سیستمهای NOSQL

انتخاب مناسبی در راستای برآورده کردن نیازهای مرکز منطقه ای می باشد؟

مسئله در آینده نزدیک در مرکز منطقه ای با حجم بی سابقه داده ها و نیازمندی های مقیاس -

پذیری و کارایی برای این حجم از داده ها، این عرف را دستخوش تغییر می کند که پایگاه

داده های رابطه ای تنها رویکرد مبتنی بر مدیریت داده ها هستند. بطور عمومی مبتنی بر

تجارب دیگر پژوهشها و سازمانها با داده های حجیم، پایگاه داده های رابطه ای دارای مشکلاتی

در مقیاس پذیری موثر، ذخیره سازی حجم بسیار بالا داده، پاسخ گویی کویری های پیچیده

در حجم بسیار بالای داده و ذخیره سازی داده های غیرساخت یافته می باشند.

لذا با توجه به کارایی مشهود در نتایج آزمایشات در فصل ۴ میتوان مدعی شد که پایگاههای

داده غیر رابطه ای در حجم بالای داده باعث بهبود زمان و کیفیت بازیابی، توزیع خودکار

داده و میزان تحمل خطای بالا برای ذخیره سازی داده های حجیم میشوند. لذا در آینده

نزدیک استفاده از آنها برای مرکز ضروری است، چرا که این روزها با وجود چنین

تکنولوژیهای هزینه های مقیاس پذیری عمودی غیر قابل توجه و بعضا تحمل ناپذیر بوده

و باید با استفاده از تکنیک های جدید خوشه بندی در NOSQL پایگاههای داده غیررابطه

ای، پایگاه داده ها به سمت معماریها و سیاستهای مقیاس پذیری افقی بروند.

• کدامیک از پایگاههای غیر رابطه ای/سیستمهای NOSQL بهترین تناسب را با اهداف

کاربردی مطرح شده در این طرح دارد؟

همانطور که واضح است با گسترش روز افزون حجم داده ها و لذا نیاز به پایگاههای داده غیر رابطه ای؛ دانشگاهها و شرکتهای تجاری زیادی بصورت رایگان و غیر رایگان در حال توسعه بسترهای نرم افزاری برای پایگاههای داده غیر رابطه ای می باشند و تعداد زیادی از پلتفرمها، ابزارها و نرم افزارهای این حوزه در حال توسعه و تولید است. اما با توجه به نوع داده های حداکثری موجود و ذخیره شده در پایگاههای داده مرکز منطقه ای بهترین و محبوب ترین پایگاه داده غیر رابطه ای، پایگاه داده مبتنی بر سند (می باشد که قابلیت ذخیره سازی داده ها با فرمت هایی همچون JSON، BSON و XML را فراهم می آورند. این مدل برای ذخیره داده های بدون ساختار بسیار مناسب است چرا که هرسند می تواند از تعداد متفاوتی فیلد تشکیل شده باشد و لزومی ندارد که سند های مشابهی که در یک دسته بندی قرار می گیرند حتما از ساختار مشابهی پیروی کنند. این سیستمهای غیر رابطه ای قادرند پرس وجوهای تجمعی را اجرا کنند، نتایج را مرتب کنند و از ایندکس گذاری روی فیلدهای یک سند پشتیبانی کنند. نمونه هایی از این مدل پایگاههای داده غیر رابطه ای مبتنی بر سند Elastic, search MongoDB می باشند. که در پژوهش حاضر به صورت جامع مورد بحث و بررسی واقع شدند و نتایج آن در فصل ۴ ارائه گشت.

۵-۲ نتیجه گیری

در عصر اینترنت تعداد اسناد تمام متن روز به روز در حال افزایش است. در چنین شرایطی استخراج اطلاعات مورد نیاز از این اسناد با سرعت و کیفیت قابل قبول از اهمیت خاصی برخوردار است. راه حل پیاده سازی شده برای حل این مشکل در پایگاه‌های داده مختلف جستجو و بازیابی تمام متن می‌باشد.

هدف از این پژوهش بررسی عملکرد جستجو و بازیابی تمام متن برای پایگاه‌های داده Microsoft SQL Server و MariaDB در قیاس با MongoDB و Elasticsearch (از محبوب ترین پایگاه‌های داده غیر رابطه‌ای مبتنی بر سند جهانی) می‌باشد. با توجه به هدف پژوهش در گام اول به ارزیابی عملکرد پایگاه‌های داده معرفی شده به ازای عملیات اسکن و خوانش و به کمک معیار ارزیابی وای-سی-اس بی^۱ پرداخته و در گام بعدی بطور دقیق تر عملکرد عملیات جستجو و بازیابی تمام متن در این پایگاه‌های داده بررسی شد. همچنین در این گام، در مورد تاثیر شاردینگ - که مرتبط با مهم ترین مزیت پایگاه‌های داده غیر رابطه‌ای نسبت به پایگاه‌های داده رابطه‌ای یعنی مقیاس پذیری افقی می‌باشد - بر عملکرد جستجو و بازیابی تمام متن MongoDB و Elasticsearch از دیدگاه معیار زمان بازیابی و کیفیت بازیابی بحث شد.

سایر نتایج حاصل شده به صورت جامع با جزییات در فصل ۴ شرح داده شد. در اینجا صرفاً به اختصار به نتایج کلی و نهایی کار پس از انجام پژوهش اشاره می‌نماییم. در آزمایشات مربوط به جستجو و بازیابی تمام متن که عمده تمرکز این پژوهش محسوب میشود؛ در مقیاسی که پایگاه‌های داده SQL Server، MariaDB، MongoDB و Elasticsearch را بوسیله جستجوی تمام متن مورد

¹ YCSB Benchamrk

ارزیابی قرار دادیم، برتری چشمگیر ElasticSearch و Microsoft SQL Server نسبت به دو پایگاه داده دیگر مشهود بود. همچنین مشاهده شد که با افزایش تعداد Shardها در Elastic Search؛ Elasticsearch برتری خود را نسبت به Microsoft SQL Server افزایش می دهد.

بطور خلاصه میتوان اینگونه بیان نمود که؛ سیستمهای مربوط به NOSQL/پایگاههای داده غیررابطه ای MongoDB و رابطه ای MariaDB که عملکرد خوبی را در وای-سی-اس-بی ثبت کرده بودند، به عنوان پایگاه داده صرفاً برای ذخیره سازی و پایداری داده ها مناسب هستند. این دو پایگاه داده در قسمت ارزیابی جستجوی تمام متن که عمده تمرکز این پژوهش محسوب می شود به مراتب عملکرد ضعیف تری نسبت به رقبای خود یعنی پایگاههای داده Microsoft SQL Server و Elasticsearch نشان دادند و بطور کلی در این قسمت بهترین عملکرد به Elasticsearch تعلق گرفت. همچنین استفاده از قابلیت شاردینگ در پایگاه داده غیر رابطه ای بخصوص در حالت ElasticSearch باعث بهبود قابل توجه عملکرد آن در اجرای جستجوی تمام متن گشت. بنحوی که جمع بندی وضعیت پایگاههای داده در تمام حالات (انواع پرس و جوهای آزمایشی و شیوه های مختلف خوشه بندی¹ از نظر معیار کیفیت بازیابی (QRRT) به شرح زیر قابل خلاصه شدن است.

QRRT: ElasticSearch > SQL > MongoDB ~ MariaDB

در واقع همانطور که مشخص است، معیار کیفیت بازیابی در ElasticSearch دارای بهترین وضعیت و سپس در رده دوم SQL و در رده بعدی MongoDB قرار دارد و نهایتاً MariaDB در بدترین وضعیت از نظر معیار کیفیت بازیابی قرار دارد.

¹ Sharded-Cluster

با در نظر گرفتن نتایج این پژوهش با مطالعه انجام شده بر روی داده های مرکز منطقه ای و با تکیه بر ویژگی مقیاس پذیری افقی پایگاه های داده غیر رابطه ای، در زمان حاضر این پژوهش (با توجه به سرعت رشد تکنولوژی در حوزه سیستم های NOSQL) در صورتیکه هدف حرکت به سمت استفاده از پایگاه های داده غیر رابطه ای باشد، بهترین تصمیم می تواند استفاده از پایگاه داده غیر رابطه ای قدرتمند همچون MongoDB به عنوان پایگاه داده اصلی و استفاده از Elasticsearch به عنوان ابزاری برای نمایه سازی¹ و جستجوی تمام متن بر روی پایگاه داده اصلی باشد. در صورتیکه همچنان نظر بر مقیاس پذیری عمودی و استفاده از پایگاه های داده رابطه ای در مرکز منطقه ای باشد، راهکار ترکیبی² استفاده از SQL به عنوان پایگاه داده و استفاده از موتور جستجوی قدرتمند ElasticSearch بر روی آن می باشد.

۳-۵ پیشنهاد های آینده

با در نظر گرفتن نتایج حاصل شده پژوهش، در پایان به عنوان پیشنهاد های آینده این طرح پژوهشی میتوان بر روی موارد زیر پژوهش های بیشتری را با توجه به زیرساخت های فعلی و متناسب با ساختار و ماهیت داده های مرکز منطقه ای اطلاع رسانی علوم و فناوری لحاظ نمود:

- بررسی عمیق و جزئی بر نحوه عملکرد الگوریتمها و پرس و جوهای آزمایشی مختلف این دو پایگاه داده غیر رابطه ای MongoDB و Elasticsearch که خود از چالش های مهم پایگاه های داده غیر رابطه ای / سیستم های NOSQL در قیاس با پایگاه های داده رابطه ای محسوب میشود.

¹ Indexing

² Hybrid

- مطالعه و پژوهش بر چگونگی استفاده همزمان بصورت ترکیبی از هر دو بستر نرم افزاری پایگاههای داده رابطه ای و غیر رابطه ای و بررسی نتایج حاصل از آنها با هدف استفاده حداکثری از ویژگیهای هر دو (از جمله پایگاه داده رابطه ای SQL و پایگاه داده غیر رابطه ای یا به اصطلاح سیستم NOSQL - Elastic Search)
- مطالعه بر روی معماریهای متنوع احتمالی دیگر جهت برآورد هزینه سودمندی استفاده از مقیاس پذیری افقی در صورت فراهم آمدن زیرساختهای لازم

منابع و ماخذ

منابع و ماخذ فارسی

دولو، صدیقه و کوروش کیانی (۱۳۹۵). تحلیل تحرک فیزیکی افراد با استفاده از داده های کلان جمع آوری شده توسط برنامه کاربردی نصب شده روی موبایل افراد. اولین کنفرانس بین المللی دستاوردهای نوین پژوهشی در مهندسی برق و کامپیوتر ۱۳۹۵، تهران، کنفرانسیون بین المللی مخترعان جهان (IFIA)، دانشگاه جامع علمی کاربردی. https://www.civilica.com/Paper-CBCONF01-CBCONF01_1075.html

شعله ارسطوپور (۱۳۹۵). سند راهبردی مرکز منطقه ای اطلاع رسانی علوم و فناوری. شیراز: انتشارات مرکز منطقه ای اطلاع رسانی علوم و فناوری.

نوریان، فرشاد و سمانه حجازی (۱۳۹۲). کاربرد کلان-داده ها در نقد توسعه مبتنی بر حمل و نقل عمومی. فصلنامه مطالعات شهری دانشگاه کردستان دوره ۸، شماره ۸: صفحه ۸۳-۹۱.

مدیرفرزام، نیلوفر و احمد فراهی (۱۳۹۴). مدیریت داده های حجیم با استفاده از سیستم های NoSQL، کنفرانس بین المللی پژوهش های کاربردی در فناوری اطلاعات ۱۳۹۴، کامپیوتر و مخابرات، تربت حیدریه، شرکت مخابرات خراسان رضوی، https://www.civilica.com/Paper-ITCC01-ITCC01_124.html

نصری فلاح، پروین (۱۳۹۵). کاربرد و چالش های کلان داده ها، اولین کنفرانس بین المللی چشم انداز های نو در مهندسی برق و کامپیوتر ۱۳۹۵، تهران، کنفرانسیون بین المللی مخترعان جهان (IFIA)، دانشگاه جامع علمی کاربردی، https://www.civilica.com/Paper-NPECE01-NPECE01_388.html

منابع و ماخذ لاتین

- Abramova, V., & Bernardino, J. (2013). NoSQL databases: MongoDB vs Cassandra. In Proceedings of the international C* conference on computer science and software engineering (pp. 14-22).
- Abramova, V., J. Bernardino, P. Furtado, Which nosql database? a performance overview, Open J. Databases. 1 (2014) 17-24.

- Amazon DynamoDB: Fast and flexible NoSQL database service for any scale.(2017). Available from: <https://aws.amazon.com/dynamodb/>. Accessed September 20, 2017.
- Amazon SimpleDB (2017). Available from: <https://aws.amazon.com/simpledb> (2017). Accessed September 2017.
- Carpenter, J., & Hewitt, E. (2020). Cassandra: the definitive guide: distributed data at web scale. O'Reilly Media.
- Chambers, B., & Zaharia, M. (2018). Spark: The definitive guide: Big data processing made simple. "O'Reilly Media, Inc."
- Chang, F., Dean, J., Ghemawat, S., Hsieh, W. C., Wallach, D. A., Burrows, M., ... & Gruber, R. E. (2008). Bigtable: A distributed storage system for structured data. ACM Transactions on Computer Systems (TOCS), 26(2), 1-26.
- Dede, E., Govindaraju, M., Gunter, D., Canon, R. S., & Ramakrishnan, L. (2013). Performance evaluation of a mongodb and hadoop platform for scientific data analysis. In Proceedings of the 4th ACM workshop on Scientific cloud computing (pp. 13-20).
- Chen, J. K., & Lee, W. Z. (2019). An Introduction of NoSQL Databases Based on Their Categories and Application Industries. Algorithms, 12(5), 106.
- Cooper B.F., A. Silberstein, E. Tam, R. Ramakrishnan, R. Sears, Benchmarking cloud serving systems with YCSB, Proc. 1st ACM Symp. Cloud Comput. - SoCC '10. (2010) 143–154. doi:10.1145/1807128.1807152.
- DB-Engines Ranking (2017). Available from: <https://db-engines.com/en/ranking>. Accessed September 20, 2017.
- Divya, M. S., & Goyal, S. K. (2013). Elasticsearch: An advanced and quick search technique to handle voluminous data. Compusoft, 2(6), 171.
- Dixit, B. (2016). Elasticsearch Essentials. Packt Publishing Ltd. Elasticsearch. Available from: <https://www.elastic.co/guide/en/elasticsearch/reference/5.5/modules-cluster.html>. Accessed August 2, 2019.
- Elasticsearch Resiliency Status (2019). Available from: <https://www.elastic.co/guide/en/elasticsearch/resiliency/current/index.html> Accessed August 2, 2019.
- ElasticSearch Manual Document1 (2019). Available from: <https://www.elastic.co/guide/en/elasticsearch/reference/current/modules-node.html> . Accessed August 2, 2019.
- ElasticSearch Manual Document2 (2019). Available from: Document1<https://www.elastic.co/guide/en/elasticsearch/reference/5.5/important-settings.html>. Accessed August 2, 2019.

- Ghazal, A., T. Rabl, M. Hu, F. Raab, M. Poess, A. Crolotte, H.-A. Jacobsen, BigBench: towards an industry standard benchmark for big data analytics, in: Proc. 2013 ACM SIGMOD Int. Conf. Manag. Data, 2013: pp. 1197–1208.
- Github. Available from: <https://github.com/elastic/elasticsearch/issues/10933>. Accessed August 2, 2019.
- Google Trends. “Using NOSQL versus SQL”, Available from: <https://trends.google.com/>. Accessed August 2, 2019.
- Gudivada, V. N., Rao, D., & Raghavan, V. V. (2014). NoSQL systems for big data management. In 2014 IEEE World congress on services (pp. 190-197). IEEE.
- Hadoop (2019). Available from : <http://hadoop.apache.org/>. Accessed September 20, 2019.
- Han, J., Haihong, E., Le, G., & Du, J. (2011). Survey on NoSQL database. In 2011 6th international conference on pervasive computing and applications (pp. 363-366). IEEE.
- HBase (2017). Available from: <https://hbase.apache.org/>. Accessed September 20, 2017.
- Cloud Serv. Comput., IEEE, 2011: pp. 336–341. doi:10.1109/CSC.2011.6138544.
- Hoberman, S. (2014). Data Modeling for MongoDB: Building Well-Designed and Supportable MongoDB Databases. Technics Publications.
- Hypertable. Available from: www.hypertable.com. Accessed September 20, 2017.
- Kenler, E., & Razzoli, F. (2015). MariaDB Essentials. Packt Publishing Ltd.
- Li, Y., & Manoharan, S. (2013). A performance comparison of SQL and NoSQL databases. In 2013 IEEE Pacific Rim Conference on Communications, Computers and Signal Processing (PACRIM) (pp. 15-19). IEEE.
- Macedo, T., & Oliveira, F. (2011). Redis Cookbook: Practical Techniques for Fast Data Manipulation. “O'Reilly Media, Inc.”
- Mavridis, I., H. Karatza, Performance evaluation of cloud-based log file analysis with Apache Hadoop and Apache Spark, J. Syst. Softw. 125 (2017) 133–151.
- MongoDB (2019). Available from: <https://www.mongodb.com/industries/retail>. Accessed August 2, 2019.
- MongoDb Manual Document (2019). Available from: <https://docs.mongodb.com/manual/tutorial/deploy-shard-cluster/>. Accessed August 2, 2019.

Oracle-RDBMS (2017). Available from: <https://www.oracle.com/database/what-is-a-relational-database>. Accessed September 20, 2017.

Patil, M. M., Hanni, A., Tejeshwar, C. H., & Patil, P. (2017). A qualitative analysis of the performance of MongoDB vs MySQL database based on insertion and retrieval operations using a web/android application to explore load balancing—Sharding in MongoDB and its advantages. In 2017 International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC) (pp. 325-330). IEEE.

Riak (2017). Available from: <http://basho.com/products/#riak>. Accessed September 20, 2017.

Redis (2017). Available from: <https://redis.io/>. Accessed September 20, 2017.

Vaish, G. (2013). Getting started with NoSQL. Packt Publishing Ltd.

Voldmort Project (2017). Available from: <http://project-voldemort.com>. Accessed August 2, 2019.

Yahoo! Cloud Serving Benchmark (YCSB) (2019). Available from: <https://github.com/brianfrankcooper/YCSB/wiki>. Accessed August 2, 2019.