

## ارائه ی معماری مدیریت داده های کلان در آموزش و پرورش

### جواد افشاری

کارشناس ارشد مهندسی کامپیوتر گرایش نرم افزار، مدرس گروه کامپیوتر دانشگاه علمی کاربردی فرماندهی انتظامی خراسان جنوبی

### چکیده

عبارت کلان داده مدت ها است که برای اشاره به حجم عظیمی از داده ها که توسط سازمان های بزرگی مانند گوگل یا ناسا ذخیره و تحلیل می شوند مورد استفاده قرار می گیرد. اما این اصطلاح بیشتر برای اشاره به مجموعه های بزرگی از داده ها استفاده می شود که به قدری بزرگ و حجیم هستند که با ابزارهای مدیریتی و پایگاه داده های رابطه ای و معمولی قابل مدیریت نیستند. مشکلات اصلی در کار با این نوع داده ها مربوط به توانایی گردآوری، ذخیره سازی، جست و جو، اشتراک گذاری، تحلیل و نمایش آن ها است. داده های ذخیره شده در آموزش و پرورش به قدری حجیم و متنوع هستند که انجام فرآیندهای مختلف روی این پایگاه داده ها با سرعت پایین انجام می گیرد. به نحوی که در پاره ای از مواقع پاسخ به کاربران در زمان مناسب ارائه نمی شود. از طرفی قابلیت های محدود پایگاه داده های رابطه ای در ذخیره سازی داده ها و سرعت رشد داده های آموزش و پرورش، در آینده ی نزدیک با مشکلات عدیده ای در ذخیره سازی و ارائه ی گزارش و تحلیل از این داده ها روبرو خواهد شد.

هدف این پژوهش، ارائه ی یک معماری مدیریت داده های کلان برای آموزش و پرورش به کمک ابزارهای موجود است، به نحوی که زیرساخت نرم افزاری مناسب برای مدیریت و تحلیل داده های کلان روی سخت افزارهای موجود پیاده سازی شود. با کمک این ابزارها، امکان تولید گزارش های درخواستی در بازه ی زمانی مناسب برای کاربران فراهم می شود.

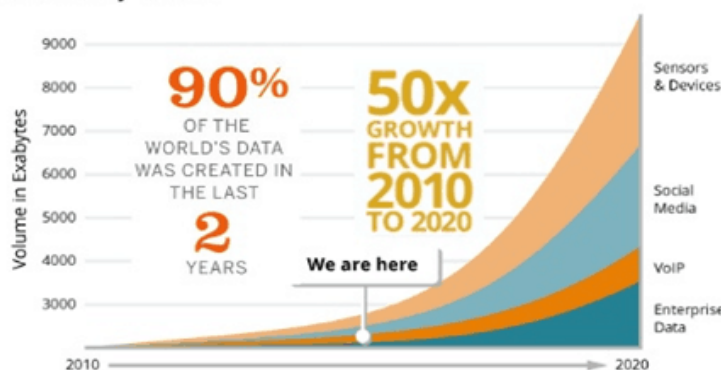
واژگان کلیدی: داده کلان، حجم داده، پایگاه داده ی غیررابطه ای

## مقدمه

امروزه بسیاری از متخصصان و حتی کاربرانی که در حوزه‌های مختلفی از فناوری اطلاعات فعالیت می‌کنند، از روند رشد چشمگیر داده‌ها و چگونگی پردازش آن‌ها شنیده‌اند و حداقل این نمودار رشد داده‌ها را مشاهده کرده‌اند. در مثال‌ها، بیشتر به شبکه‌های اجتماعی پرداخته می‌شود که کاربران آن‌ها در هر سال بیش از ۳ برابر سال قبل از آن، داده تولید می‌کنند. عبارت کلان‌داده مدت‌ها است که برای اشاره به حجم عظیمی از داده‌ها که توسط سازمان‌های بزرگی مانند گوگل یا ناسا ذخیره و تحلیل می‌شوند مورد استفاده قرار می‌گیرد. اما به‌تازگی، این عبارت بیشتر برای اشاره به مجموعه‌های بزرگی از داده‌ها استفاده می‌شود که به‌قدری بزرگ و حجیم هستند که با ابزارهای مدیریتی و پایگاه‌داده‌های رابطه‌ای و معمولی قابل مدیریت نیستند. مشکلات اصلی در کار با این نوع داده‌ها مربوط به برداشت و جمع‌آوری، ذخیره‌سازی، جست‌وجو، اشتراک‌گذاری، تحلیل و نمایش آن‌ها است [۱]. نکته‌ی جالب توجه آن است که ۹۰ درصد داده‌هایی که اکنون در اختیار ما است، تنها در چند سال اخیر تولید شده است. شکل زیر رشد سریع داده‌ها در فاصله سال‌های ۲۰۱۰ تا ۲۰۲۰ را به خوبی نشان می‌دهد:

## BIG IN GROWTH, TOO.

1 exabyte (EB) = 1,000,000,000,000,000 bytes



شکل (۱): حجم داده‌های تولید شده در فاصله‌ی سال‌های ۲۰۱۰ تا ۲۰۲۰ بر حسب اگزابایت [۲]

برخی از مهمترین اقدامات که بر روی کلان‌داده انجام گرفته عبارتند از [۳]:  
ذخیره‌سازی داده‌های حجیم و بازیابی داده‌ها، توزیع داده‌ها، شبکه‌های پرسرعت کامپیوتری، عملکرد محاسباتی بالا، مدیریت موضوع و رویدادها، داده‌کاوی و تجزیه و تحلیل، یادگیری ماشین و تجسم سازی داده‌ها واضح است که تکنیک‌های فوق در گذشته نیز مورد استفاده قرار گرفته‌اند، لیکن اهمیت آن‌ها را در موضوع کلان‌داده‌ها، گردآوری و ذخیره‌سازی داده‌های فراوان و نیز ساخت‌افزارهای مناسب است. دیدگاه صنعتی و علمی کلان‌داده با سه مفهوم مشخص توصیف می‌شود [۳]:

- (۱) داده‌هایی که حجم آن‌ها بسیار زیاد است.
  - (۲) داده‌هایی که در پایگاه‌داده‌های رابطه‌ای منظم طبقه‌بندی نمی‌شوند.
  - (۳) داده‌هایی که با سرعت، ایجاد، دریافت و پردازش می‌شوند.
- آنچه مسلم است، کلان‌داده در طول پنج سال آینده یکی از پنج موضوع حیاتی تحقیقاتی در زمینه‌ی تکنولوژی بوده و موجب تحول در بسیاری از زمینه‌ها از جمله تجارت، تحقیقات علمی و صنعتی، مدیریت اجتماعی، پزشکی و... خواهد گردید [۴].
- مفهوم کلان‌داده دارای سه ویژگی اصلی است [۵]:
- ۱- حجم (Volume): داده‌ها به‌طور میانگین سالیانه دو برابر می‌شوند. یک دلیل ساده برای آن این است که سازمان‌ها به صورت دوره‌ای اطلاعات خود را برای سال‌های متمادی نگه می‌دارند.
  - ۲- سرعت (velocity): مفهوم آن این است که داده‌ها با چه سرعتی تولید می‌شوند و با چه سرعتی برای تحلیل و برآوردن آن‌چه مورد نیاز است پردازش می‌شوند. به‌عنوان مثال، سرعت تولید داده‌ی حسگرهای شناسایی فرکانس رادیویی بسیار بالا

بوده که علاوه بر ذخیره سازی آن‌ها در لحظه، باید مورد تجزیه و تحلیل نیز قرار بگیرند. واضح است سیستم‌های ذخیره سازی و تحلیل اطلاعات رابطه‌ای، به سادگی نمی‌توانند این جریان از اطلاعات را در لحظه مورد بررسی و نمایش قرار دهند. عکس‌العمل سریع در قبال سرعت فوق‌العاده تولید داده‌ها برای سازمان‌ها یک چالش بزرگ محسوب می‌شود.

۳-تنوع (Variety): امروزه اطلاعات در انواع ساختارهای متفاوت از سیستم‌های تولید داده‌ی استاندارد و یا تراکنش‌های پایگاه داده‌های تحلیلی آن‌لاین تولید می‌شوند. به عبارتی منابع تولید داده متفاوت بوده و موجب تولید داده‌ها با ساختارهای متفاوت می‌گردد. از یک‌نظر ساختار داده‌ها را می‌توان در سه دسته‌ی ساخت یافته، نیمه ساختار و بدون ساختار تقسیم‌بندی نمود.

علاوه بر موارد فوق، برخی شرکت‌های بزرگ ویژگی‌های دیگری را نیز برای کلان داده در نظر می‌گیرند:

۴-صحت (Veracity): شرکت IBM معتقد است نباید داده‌هایی که به آن اعتقاد نداریم را به کار ببریم. درواقع نمی‌توان انتظار داشت تمام داده‌ها، به طور کامل صحیح باشند، بلکه کیفیت داده‌ها به میزان زیادی وابسته به منابعی است که از آن‌ها گردآوری می‌شوند [۴].

۵-پیچیدگی (Complexity): کلان داده از منابع مختلف گردآوری شده‌اند و لازم است ارتباط میان آن‌ها به منظور دستیابی به داده‌های ساخت یافته و یا غیرساخت یافته، در نظر گرفته شود.

۶-ارزش (Value): باید در نظر داشت، مهم‌ترین ویژگی کلان داده، ارزش آن است. گرچه دسترسی به حجم بالای داده‌ها خوب است، اما تنها زمانی این دسترسی مفید است که داده‌ها، دارای ارزش واقعی باشند. در این صورت، زیرساخت تکنولوژی برای ذخیره سازی و سرمایه‌گذاری تجاری کلان داده، بسیار ارزشمند خواهد بود [۴].

برخی دیگر از موارد کاربردی کلان داده‌ها را می‌توان در قالب زیر برشمرد:

۱-بهداشت و درمان: تحلیل کلان داده‌ها می‌تواند در صنعت بهداشت و درمان در قالب ارائه‌ی خدمات بهتر به عموم مردم کمک کند. شخصی سازی درمان می‌تواند منتج به افزایش سلامت جامعه و کاهش هزینه‌های دولت در بخش بهداشت و درمان شود.

۲-آموزش: کلان داده‌ها در صنعت آموزش به شخصی سازی فرایند یادگیری کمک می‌کند. این شخصی سازی به نوبه‌ی خود می‌تواند باعث شکوفایی استعدادهای دانش‌آموزان و دانشجویان شود و پویایی محیط یادگیری را افزایش دهد.

۳-تولید: در صنعت تولید استفاده از کلان داده‌ها می‌تواند به تولید طبق نیازهای مشتری کمک کند، زمان تولید محصول را کاهش دهد، همین‌طور با استفاده از شبیه سازی و بهینه سازی به کمک کلان داده‌ها می‌توان خط تولید را به صورت بهینه طراحی کرد و بسیاری از عیوب خط تولید کالاها را پیش از شروع به کار خط تولید شناسایی کرد.

۴-دولت: دولت می‌تواند از کلان داده‌ها برای ایجاد شفافیت، خدمت رسانی بهتر به مردم، استفاده‌ی بهینه از منابع محدود و تخصیص بودجه به فعالیتهای موجود استفاده کند. همین‌طور می‌تواند برای کمک و اطلاع رسانی در زمان بحران به مردم مبارزه با فقر و جرم و جنایت کلان داده‌ها را بکار بگیرد.

۵-علوم اجتماعی: در مطالعات علوم اجتماعی کلان داده‌ها ابزاری برای بررسی پیچیدگی رفتار فردی و اجتماعی انسان‌ها بوده و دریچه‌ای جدید برای طرح سوال‌های جالب‌تر و الگوهای ناشناخته باشد.

۶-ورزش: در علوم ورزشی از کلان داده‌ها برای افزایش کارایی ورزشکاران در تمرین و مسابقه، پیشگیری از بروز مصدومیت و یافتن بهترین راهبرد برای مسابقات پیش‌رو استفاده می‌شود [۵].

در واقع رشد فزاینده‌ی کاربردهای کلان داده در زمینه‌های مختلف، از توانایی‌های نرم‌افزاری به کاررفته با هدف دریافت، مدیریت و پردازش این حجم داده‌ها بسیار فراتر است. مهم‌ترین چالش در کاربرد کلان-داده کاوش در حجم انبوهی از داده‌ها و استخراج داده‌های مفید و قابل استفاده برای کاربردهای آینده است. از سویی این حجم بی‌سابقه از داده‌ها نیازمند سیستم‌های تحلیلی بسیار قدرتمندی است [۵].

رویکرد استفاده از صدها ابزار ذخیره سازی داده‌های کلان گرچه امری بی‌فایده تلقی می‌گردد، اما در کنار آن امکان ذخیره سازی صدها مجموعه داده با یک ترابایت داده فراهم می‌شود. بدین ترتیب می‌توان انتظار داشت امکان دسترسی مشترک به داده‌ها برای کاربران علی‌رغم سرعت پایین در پاسخ به عملیات تحلیلی فراهم می‌شود. با این حال این روش دو چالش اساسی دارد:

- ۱- اولین چالش زمانی است که سخت افزار دچار نقص شود. در این حالت امکان رخ دادن نقص سخت-افزار همزمان با افزایش تعداد آن ها به تناسب بیشتر می شود.
- ۲- دومین چالش زمانی است که بسیاری از عملیات های تحلیل، نیازمند گردآوری داده ها به روش های مختلف است. به عنوان مثال، داده های خوانده شده ی یک دیسک نیازمند ترکیب داده های ۹۹ دیسک دیگر باشد [۵].

## روش تحقیق

مفاهیم کلیدی مورد نیاز برای معماری پیشنهادی شامل پایگاه داده ی رابطه ای، مقیاس پذیری، دریاچه داده، هدوپ، ایمپالا می باشد. این مفاهیم را به طور خلاصه بررسی می نماییم.

### ۱. پایگاه داده ی رابطه ای

مفهوم سیستم های مدیریت پایگاه داده (DBMS) از قریب چهار دهه ی گذشته توسعه یافت. سیستم های مدیریت پایگاه داده ی رابطه ای، شامل نرم افزارهایی هستند که برای مدیریت مجموعه ای از داده های مرتبط و همه ی متعلقات آن ها به کار می روند. براین اساس، سیستم های مدیریت پایگاه داده یا از نوع متمرکز هستند که پس از نصب بر روی یک سکو، به صورت فیزیکی توسط کامپیوترهای متعدد که از طریق شبکه توزیع شده اند، پردازش می شوند و یا به صورت توزیع شده، که در واقع شامل چندین پایگاه داده ی توزیع شده بر روی یک شبکه است که مانند یک سیستم واحد عمل می کنند. از رایج ترین آن ها، به سیستم مدیریت پایگاه داده ی رابطه ای (RDBMS) می توان اشاره کرد که در انواع گوناگونی از پایگاه داده ها مانند ۲DB، SQLSERVER، Oracle و ... توسعه یافته است. پایگاه داده های رابطه ای، شامل جداولی هستند که داده ها را در قالب ماتریسی از ردیف ها که هر کدام شامل ستون هایی با یک کلید اصلی یا کلید خارجی می باشند، ذخیره می کنند. تکنولوژی پایگاه داده های رابطه ای، بسیار سازگار و ثابت است که داده ها را در قالب جداول کاملاً ساخت یافته ذخیره می نمایند. گرچه این نوع پایگاه داده ها برای انجام پردازش تراکنش آن لاین یا به اختصار OLTP مناسبند، در عین حال قادر به انتقال پاسخ سریع در صورت افزودن پردازنده یا حافظه بیشتر نمی باشند [۵].

### ۲. مقیاس پذیری

یک سامانه که کارایی آن پس از افزودن سخت افزار جدید به تناسب ارتقا یافته است، یک سیستم مقیاس پذیر نامیده می شود. این اصطلاح برای الگوریتم های مناسب برای اجرا در موقعیت های بزرگتر نیز به کار می رود. با دو رویکرد می توان یک سامانه را پس از افزودن منابع سخت افزاری مقیاس پذیر نمود:

اولین رویکرد، زمانی است که یک سیستم به صورت عمودی مقیاس پذیر می گردد. در این حالت منابع به یک گره تنها در یک سامانه افزوده می شود که در حالت رایج، شامل افزودن یک پردازنده یا حافظه به یک رایانه است. مقیاس پذیری عمودی سیستم های موجود، آن ها را در استفاده از تکنولوژی مجازی سازی توانمند می سازد و منابع بیشتری برای استفاده ی سیستم های عامل میزبانی شده و برنامه های کاربردی اشتراک گذاری شده، در اختیار قرار می دهد.

دومین رویکرد، زمانی است که یک سیستم به صورت افقی مقیاس پذیر می شود. در این حالت گره های جدیدی به سامانه افزوده می شود. افزودن کامپیوتر جدید به یک نرم افزار کاربردی توزیع شده، افزودن یک وب سرور به یک سامانه با تعداد ۳ عدد وب سرور مثال هایی از این نوع مقیاس پذیری است [۶].

### ۳. دریاچه داده

دریاچه داده یک راه حل برای مدیریت داده های سازمانی است که مقدار زیادی از داده ها در یک مخزن بزرگ برای تجزیه و تحلیل بیشتر ذخیره می شوند [۷]. در واقع دریاچه داده یک مخزن بسیار بزرگ از داده های خام در همان قالب اصلی خود است که داده ها را برای مدت زمان طولانی و تا زمان نیاز در خود نگهداری می کند. رویکرد دریاچه داده در برخی جنبه ها دارای تفاوت هایی با سایر موارد مشابه است :

- داده‌ها کاملاً در دریاچه‌داده قرار می‌گیرند و از هیچ داده‌ای صرف‌نظر نمی‌شود.
  - داده‌های پایین‌ترین سطوح (مثلاً توضیحات یک فرد در یک مقاله یا یک سایت) بدون تغییر یا با حداقل تغییرات به دریاچه‌داده منتقل می‌شوند. این رویکرد ذخیره‌سازی داده‌ها که در آن داده، بدون توجه به ساختار و منبع ذخیره می‌شود را به اصطلاح خواندن با ساختار می‌نامند. داده‌ها به صورت خام در فرمت‌های اولیه نگهداری می‌شوند و ابزارهایی در اختیار کاربران قرار داده می‌شود که بتوانند داده‌ها را با همین شکل مورد ارزیابی و تحلیل قرار دهند. گاهی این ارزیابی مشخص می‌کند چه بخشی از داده‌ها ارزش لازم برای تحلیل‌های کامل‌تر را دارند و ممکن است برای این بخش از داده‌ها، ساختار مورد نیاز را طراحی کنند.
- از چالش‌های اصلی دریاچه‌داده، درک صحیح اهداف تحلیل‌های مورد نیاز و محدودیت‌های پیاده‌سازی به سبب در اختیار نبودن ابزارهای تکنولوژی را می‌توان برشمرد [۸]. رویکرد دریاچه‌داده در تضاد با رویکرد رابطه‌ای انبارداده است. در انبارداده، داده‌ها پیش از بارگذاری در یک شمای از پیش تعیین شده مورد تجزیه و تحلیل قرار می‌گیرند که باعث صرف زمان زیاد قبل از استفاده مفید داده‌ها و کاهش کارایی سیستم‌های مربوط می‌شود [۹]. از دیدگاه علوم داده نگهداری همه‌ی داده‌ها به صورت کامل مفیدتر و کاربردی‌تر است. زیرا مشخص نیست کدام بخش از داده‌ها برای تحلیل‌های داده‌ای پیش‌رو دارای ارزش خواهند بود [۱۰]. انبارداده مبتنی بر عملیات استخراج، انتقال و بارگذاری داده‌ها یا به اختصار ETL می‌باشد که دیگر برای سرعت رشد تولید داده‌های امروزی مناسب نمی‌باشد [۱۱]. دریاچه‌داده بر خلاف رویکرد انبارداده، داده‌ها در فرم خام و اصلی خود ذخیره می‌شوند و لذا نیاز به فرآیند ETL نبوده و داده‌ها بلافاصله در اختیار کاربران قرار دارد. تفاوت بین دریاچه‌داده و انبارداده را می‌توان در حوزه‌های جدول زیر با یکدیگر مقایسه کرد [۱۲]:

جدول (۱): تفاوت دو رویکرد انبارداده و دریاچه‌داده در تحلیل داده‌ها

حوزه مقایسه	انبار داده	دریاچه داده
نوع داده‌ها	پردازش شده، ساختار یافته	ساختار یافته، نیمه ساختار یافته، غیر ساختار یافته، خام
نحوه پردازش	ساختار داده‌ها به طور دقیق تعریف شده و به صورت از پیش تعیین شده پردازش می‌شوند. (داده‌ها به صورت ساختار یافته و با هدف خاصی ذخیره می‌شوند).	ساختار داده‌ها به طور دقیق تعریف نشده و به صورت خام ذخیره می‌شوند. (داده‌ها به صورت خام ذخیره می‌شوند و بعداً پردازش می‌شوند).
هزینه ذخیره‌سازی	هزینه ذخیره‌سازی بالا برای حجم بسیار زیاد داده‌ها	هزینه ذخیره‌سازی پایین برای حجم بسیار زیاد داده‌ها
تغییر پذیری	تغییرات داده کم و به صورت مشخص	تغییرات داده زیاد و به صورت نامشخص
امنیت	بالا	در حال رشد
کاربران	تحلیل‌گران تجاری	دانشمندان داده و تحلیل‌گران فنی
روش تحلیل داده‌ها	تحلیل داده به صورت ساختار یافته و با هدف خاصی انجام می‌شود.	تحلیل داده به صورت خام و با هدف‌های مختلف انجام می‌شود.

#### ۴.۴. هادوپ

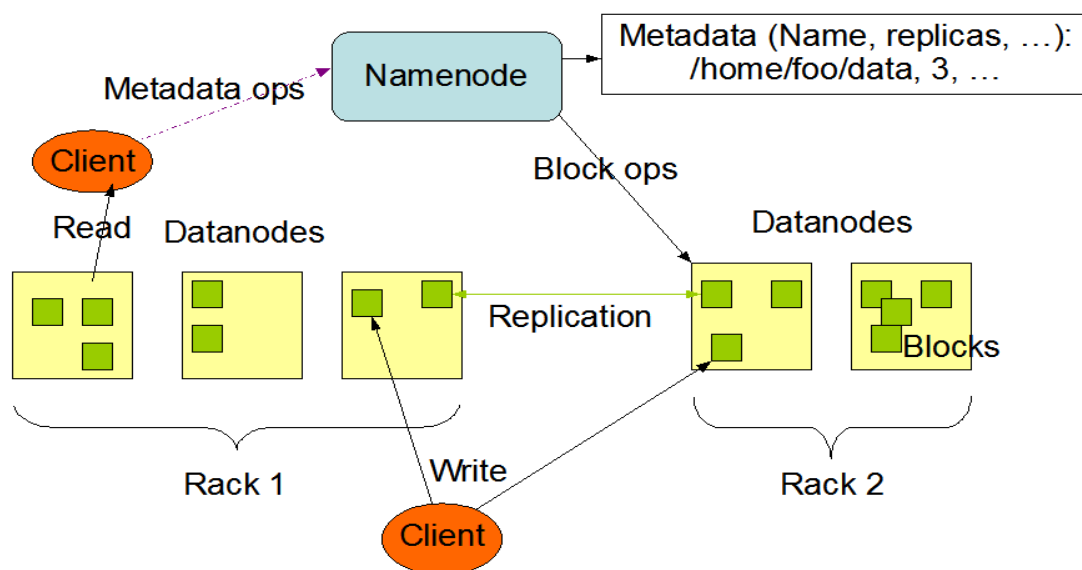
ساخت یک موتور جستجوی اینترنت ابتدا یک هدف جاه‌طلبانه ارزیابی می‌شد. انتقال حجم بالایی از صفحات وب برای عملیات کاوش یک چالش بزرگ محسوب می‌شد که توانایی مقیاس‌پذیری برای مواجهه با میلیون‌ها صفحه‌ی وب ذخیره شده را نداشت. با توسعه ابزار جدیدی به نام GFS امکان مواجهه با حجم بسیار زیادی از صفحات وب جمع‌آوری شده توسط ابزارهای جستجو فراهم گردید. این ابزار امکان مدیریت گره‌های ذخیره‌سازی را نیز داشت. همزمان توسعه ابزار متن‌باز جدید NDFS در

سال ۲۰۰۴ آغاز گردید و همان زمان، گوگل ابزار متن باز MapReduce را معرفی کرد. در سال ۲۰۰۸ ابزار هادوپ معرفی شد که به صورت متن باز امکان ذخیره سازی و اجرای برنامه های کاربردی روی تعداد زیادی از سخت افزارهای همگام را فراهم می نمود [۱۴][۳]. برخلاف پایگاه داده های رابطه ای، هادوپ توانایی کار تحلیلی روی حجم زیادی از داده ها را دارد [۱۴]. هادوپ در واقع یک چارچوب متن باز پردازش، ذخیره سازی و تحلیل حجم زیادی از داده های توزیع شده و بدون ساختار است. آغاز کار هادوپ، به شرکت های بزرگ جستجوی اینترنتی بر می گردد که دنبال ابزارها و مدل های پردازشی جدید برای استفاده در شاخص گذاری صفحات اینترنتی و جستجو روی آن ها بوده اند [۱۴]. مهمترین عامل گرایش سازمان ها در بکارگیری هادوپ، توانایی بالای آن در ذخیره سازی و پردازش حجم زیادی از انواع مختلف داده است.

اجزای اصلی سازنده هادوپ عبارتند از [۱۳]:

۱- Hadoop Common: این بسته کتابخانه ها و ابزارهای کاربردی برای راه اندازی سایر اجزای هادوپ است و به عنوان چارچوب اصلی ارائه دهنده خدمات و فرآیندهای اساسی مانند انتزاع از سیستم عامل و سیستم فایل عمل می کند.

۲- Hadoop Distributed File System (HDFS): این ابزار حجم بسیار زیادی از داده ها در ماشین های مختلف که معمولا شامل صدها یا هزاران گره متصل به هم می باشد را ذخیره می کند. دو نسخه از آن ها در یک گروه و نسخه ی دیگر در گروهی دیگر قرار می گیرد. معماری HDFS شامل دسته بندی هایی است که از طریق یک ابزار نرم افزاری به نام NameNode امکان مدیریت و رصد سیستم فایل های این دسته بندی ها و سازوکار دسترسی کاربران به آن ها را فراهم می کند و روی یک ماشین نصب می شود، در دسترس قرار می گیرد. روی سایر ماشینها یک نمونه از ابزار DataNode برای مدیریت ذخیره سازی دسته ها نصب می شود. این سیستم فایل دارای قابلیت هایی از قبیل ذخیره سازی توزیع شده، قابل حمل بودن، قابلیت اعتماد (حفظ و نگهداری از داده ها) می باشد.



شکل (۲): معماری سیستم فایل توزیع شده ی هادوپ

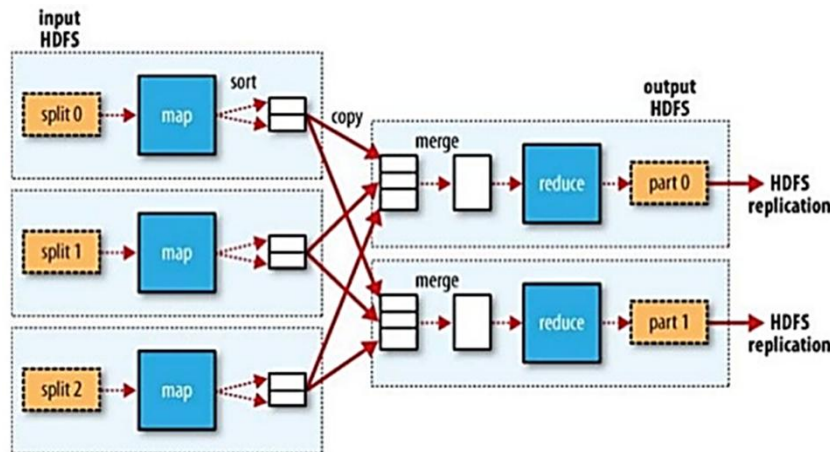
۳- MapReduce: در واقع یک مدل برنامه نویسی است که توسط شرکت گوگل برای پردازش و تولید مجموعه داده های بزرگ روی دسته هایی از کامپیوترها با استفاده از روش تقسیم و غلبه برای درهم شکستن مسائل داده های حجیم به بخشهای کاری کوچک و پردازش موازی آن ها معرفی شده است. در واقع MapReduce روی دسته های بزرگی از کامپیوترهای همگام اجرا می شود و بالقوه مقیاس پذیر است. چارچوب MapReduce از دو بخش اصلی تشکیل شده است:

- **تابع Map**: این تابع اجازه می دهد نقاط مختلف دسته ی توزیع، فعالیت خود را توزیع نمایند. در این مرحله داده ی گره اصلی به تعدادی زیر مسئله ی کوچکتر خرد شده و به دنبال آن یک گره کارگر تعدادی زیر مجموعه از



مسئله‌های کوچکتر را زیر نظر مدیر درخواست پردازش کرده و نتایج را در سیستم‌فایل محلی ذخیره میکند، به نحوی که یک تابع کاهنده قادر به دسترسی به آن باشد.

- تابع Reduce: این تابع برای کاهش شکل نهایی نتایج دسته‌ها به یک خروجی طراحی شده و داده‌های ورودی از مراحل نگاشت را تحلیل و ادغام می‌کند. چندین وظیفه‌ی کاهش برای موازی‌سازی اجتماع می‌تواند همزمان اجرا شده و این وظایف روی گره‌های کارگر تحت نظر مدیر درخواست اجرا می‌شوند.



شکل (۳): معماری مدل MapReduce در هدوپ

## ۵. کلودر ایمیپالا :

کلودر ایمیپالا نرم‌افزار متن‌بازی است که شرکت کلودرا برای پردازش موازی سنگین ارائه نموده است. ایمیپالا فناوری پایگاه‌داده‌ی مقیاس‌وسیع را به هدوپ آورده و کاربران را قادر به ارائه‌ی درخواست و دریافت پاسخ در کمترین زمان ممکن می‌نماید و دسته‌ای از نرم‌افزارهای هدوپ را برای طیف گسترده‌ای از تحلیل‌گران پایگاه‌داده‌ها، کاربران و توسعه‌دهندگان در اختیار قرار می‌دهد. این پروژه در اکتبر ۲۰۱۲ به صورت عمومی معرفی شد. این محصول حاوی یک موتور SQL برای ذخیره‌سازی داده‌ها در کلاسترهای رایانه‌ای مبتنی بر هدوپ است. موتور پردازش‌های موازی ایمیپالا، درخواست‌های SQL هدوپ را به اندازه‌ی کافی برای تحلیل‌گرانی که با SQL آشنایی دارند و کاربران ابزارهای هوش کسب‌وکار، آسان نموده و برای انجام اکتشافات تعاملی و آزمایشات مختلف به اندازه‌ی کافی سریع است. این درخواست‌ها، در قالب SQL روی داده‌های مستقر شده در HDFS و HBase اجرا می‌شود. در هنگام درخواست، داده‌ها روی HDFS جابجا یا منتقل نمی‌شوند. ایمیپالا با هدوپ تجمیع شده تا از فایل‌ها، ابرداده‌ها، حفاظت و مدیریت منابع به صورت مشترک بهره برد [۱۵].

ایمیپالا همچنین امکان کار با چارچوب‌های MapReduce، Hive و Pig را فراهم می‌کند و در واقع کلودر ایمیپالا برای ترکیب اعتمادپذیری و مقیاس‌پذیری که از هدوپ مورد انتظار است با عملکرد و پشتیبانی از SQL که توسط پایگاه‌داده‌های با موتور پردازش موازی تجاری پیشنهاد شده طراحی گردیده است. در حال حاضر کلودر ایمیپالا درخواست‌ها را با سرعت بین ۱۰ تا ۱۰۰ برابر سریعتر از راه‌حل‌های موجود هدوپ اجرا می‌کند که قابل مقایسه با پایگاه‌داده‌هایی با قابلیت پردازش موازی بوده و اجازه‌ی تحلیل‌های اکتشافی و تعاملی روی کلان‌داده‌ها را می‌دهد [۱۶]. کلودر ایمیپالا دارای یک لایه‌ی موثر ورودی، خروجی برای داده‌های ذخیره شده روی HDFS است که به‌ازای هر دیسک روی گره‌ها یک نخ ورودی و خروجی تولید می‌نماید [۱۲]. این ابزار جایگزین چارچوب‌های کاری روی MapReduce مانند Hive قرار نمی‌گیرد، بلکه درخواست‌های سریع و تعاملی به سبک SQL را روی داده‌های ذخیره شده‌ی HDFS فراهم می‌آورد. همچنین برای فراهم نمودن یک سکوی ذخیره‌سازی یکپارچه، ایمیپالا از فراداده‌ها و دستورات مشابه SQL استفاده می‌کند.

## پژوهش‌های مرتبط

سازمان‌های بزرگ به‌طور ویژه دنبال معماری و یا سیستم‌های مقیاس‌پذیری هستند که توانایی مواجهه با داده‌ها را داشته باشند. چالش‌های مدیریتی مرتبط با ماهیت انفجاری داده‌ها به‌ندرت در گذشته وجود داشته است. در حال حاضر این چالش‌ها با دسترس‌بودن جریان داده‌ها در سرعت بالا از منابع مختلف ایجاد شده است [۱۷]. تجزیه و تحلیل داده‌های کلان به بهبود رویکرد سازمان‌ها کمک شایسته‌ای می‌کند. اما لازم است داده‌ها از منابع جداگانه برای تجزیه و تحلیل و ارزیابی‌های مورد نیاز به پایگاه و مخزن اصلی آن به‌صورت مداوم و مستمر در جریان باشند. وجود یک دریاچه‌داده برای ذخیره‌سازی مقادیر زیاد داده‌ها به‌صورتی که همه‌ی آن‌ها در یک منبع اصلی گردآوری شوند و تجزیه و تحلیل نهایی روی آن‌ها صورت گیرد یک راه‌کار مناسب بدون نیاز به ساختار خاص امنیتی و صلاحیت‌سنجی می‌باشد [۱۸]. دریاچه‌داده در واقع یک انبار عظیم از داده‌های ذخیره‌سازی شده مورد نیاز در فرم و قالب اصلی خود است که با رویکرد انبارداده رابطه‌ای به‌گونه‌ای در تضاد است. به‌این صورت که در دریاچه‌داده به فرضیات و طراحی‌های بیشتری در مورد نحوه به‌کارگیری و استفاده از داده‌ها نیاز است در حالی که در انبارداده فرآیند پیش‌پردازش تعیین‌شده قبل از بارگذاری داده‌ها روی آن انجام می‌گیرد. در این صورت تلاش و وقت فراوانی را قبل از استفاده مفید از داده‌ها برای برنامه‌های کاربردی همراه دارد [۱۹]. دریاچه‌داده جایگزین و یا تکمیل‌کننده‌ی انبارداده نیست. در هسته‌ی دریاچه‌داده مجموعه‌ای از مولفه‌ها و ابزارها مانند پایگاه‌داده‌های رابطه‌ای، مراکز داده عملیاتی و خوشه‌های سیستم‌فایل توزیع‌شده وجود دارند [۲۰]. در سال ۲۰۰۳ شرکت Google معماری سیستم فایل توزیع‌شده‌ی خود را با نام GFS منتشر کرد. این معماری مسئله‌ی کمبود فضای فایل‌های حجیم تولیدشده توسط موتورهای جستجو را مدیریت و برطرف می‌نمود. در سال ۲۰۰۴ آن‌ها تصمیم گرفتند یک نسخه‌ی پیاده‌سازی شده‌ی متن‌باز از آن معماری را ایجاد نموده و آن‌را سیستم‌فایل توزیع‌شده Nutch یا به اصطلاح NDFS نامیدند. در سال ۲۰۰۴ شرکت Google مقاله‌ای را با عنوان MapReduce ارائه نمود. خیلی زود در سال ۲۰۰۵ توسعه‌دهندگان Nutch شروع به پیاده‌سازی یک نسخه از آن نمودند و طولی نکشید در اواسط همان سال تمامی الگوریتم‌های Nutch برای استفاده از MapReduce و NDFS تغییر ساختار دادند. مدل برنامه‌نویسی نگاشت‌کاهش برای پردازش مجموعه‌های عظیمی از داده‌های رایانه‌ها (گره‌ها) که روی موضوعی خاص فعالیت می‌کنند طراحی شده است. این مجموعه روی هم رفته در صورتی که از سخت‌افزار یکسان بهره‌برند به‌عنوان خوشه شناخته می‌شوند. پردازش محاسباتی روی داده‌های ذخیره شده درون سامانه‌ی فایل بدون ساختار یا بر روی پایگاه‌داده‌ی ساختاریافته قابل اجرا است. این مدل در واقع شامل دو گام اساسی است:

- گام نگاشت: گره اصلی ورودی را به قطعات کوچک‌تر تقسیم می‌کند. در این حالت یک مسئله‌ی بزرگ به چند مسئله‌ی کوچک تقسیم می‌شود. پس از آن این مسائل کوچک بین گره‌های کارگر توزیع می‌شوند. یک گره کارگر نیز ممکن است این عملیات را به‌نوبه‌ی خود تکرار نماید، که ایجادکننده‌ی ساختار درختی و چند مرحله‌ای است. هر گره کارگر زیرمسئله‌ی خود را حل نموده و نتیجه را به گره اصلی خود برمی‌گرداند.
- گام کاهش: در مرحله بعد گره اصلی جواب زیرمسائل را از گره‌های کارگر گرفته و خروجی را می‌سازد تا خروجی، که حل مسئله ورودی است، را ایجاد نماید.

در اوایل ۲۰۰۸ هدوپ به‌عنوان سریع‌ترین سیستم مرتب‌سازی یک ترابایت داده رکورددار شد. با استفاده از یک کلاستر ۹۱۰ تایی هدوپ یک ترابایت داده را در ۲۰۹ ثانیه مرتب‌سازی نمود. در اواخر همان سال Google ادعا کرد می‌تواند همان حجم داده را در ۶۸ ثانیه مرتب‌سازی نماید. در اواسط ۲۰۰۹ اعلان شد تیمی از Yahoo! توانسته همان حجم داده را در ۶۲ ثانیه مرتب‌سازی نماید.

## سامانه‌ی ثبت داده‌های محصلین آموزش و پرورش

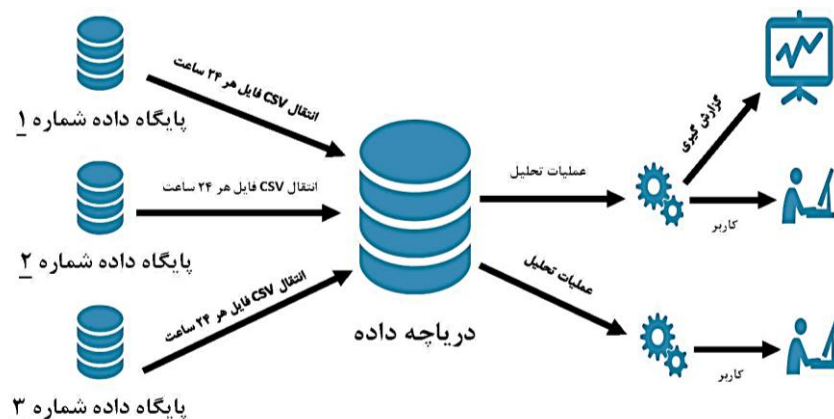
سامانه‌ی ثبت داده‌های محصلین آموزش و پرورش به‌منظور ثبت و نگهداری داده‌های تحصیلی دانش‌آموزان طراحی و پیاده‌سازی شده است. این سامانه در سال ۸۹ برای استفاده در کلیه‌ی مراکز آموزشی کشور راه‌اندازی و مورد بهره‌برداری قرار گرفت. در طراحی این سامانه از پایگاه‌داده‌ی رابطه‌ای ۲۰۰۵ SQLSERVER استفاده شده و کلیه‌ی داده‌های مرتبط با وضعیت تحصیلی



افراد برای مدت زمان نامحدود در این پایگاه داده ذخیره سازی و مورد استفاده قرار می گیرد. اطلاعات فردی، مشخصات والدین، نمرات، کارنامه های سالیانه، وضعیت انضباطی و فارغ التحصیلی مقاطع تحصیلی مختلف، داده های مربوط به محل سکونت و تماس، مشخصات فردی دبیران و آموزشگاه ... در این پایگاه داده ذخیره می شوند و در تمامی سال ها برای افراد قابل ارائه و استناد می باشند. در سامانه ی ثبت اطلاعات تحصیلی دانش آموزان از پایگاه داده ی رابطه ای SQLSERVER ۲۰۰۵ استفاده شده که به منظور کاهش سربار و افزونگی می بایست نرمال سازی، اصول مربوط به تراکنش ها، سازگاری داده ها، استقلال فرآیندها و جامعیت پایگاه داده را رعایت کرد. باین حساب در پایگاه داده ی رابطه ای فوق ۶۱ جدول مورد نیاز است تا اطلاعات مربوط به وضعیت تحصیلی دانش آموزان در ادوار مختلف را ذخیره سازی نماید. این تعداد جدول نرمال سازی شده فضای ذخیره سازی را افزایش می دهند و از طرفی لزوم وجود ارتباط زیاد بین آن ها برای استخراج اطلاعات کامل مورد درخواست از پایگاه داده سبب کاهش سرعت گزارش گیری و طولانی شدن انجام درخواست ها روی این نوع پایگاه داده ها می شود. از طرفی به علت محلی بودن پایگاه داده های فوق ارائه ی گزارش های جامع و کامل به منظور تحلیل دقیق تر از وضعیت موجود و اتخاذ رویکردهای درست و منطقی کار دشواری است.

### کلیات رویکرد پیشنهادی

با توسعه ی ابزارهایی مانند هدوپ که شرایط مطلوب برای افزایش بهره وری، سرعت و کارایی عملیات روی داده ها را فراهم نموده و نیز لزوم ارائه ی گزارش های تحلیلی جامع و کامل در زمان مطلوب برای آموزش و پرورش، در ابتدا می توان روی یک سخت افزار مناسب هدوپ و سرویس های لازم آن را فعال کرده و کلاودرایمپالا را برای ایجاد جداول و پایگاه داده ی دلخواه و گردآوری داده ها و اجرای درخواست ها فعال نمود. با توجه به قابلیت های ایمپالا که امکان بارگذاری از پایگاه داده ی رابطه ای را فراهم می کند، می توان شمای لازم از پایگاه داده را متناسب با پایگاه داده ی رابطه ای ایجاد نمود و سپس بارگذاری آن ها را از پایگاه داده ی رابطه ای به پایگاه داده ی جدید انجام داد. سپس انجام درخواست ها و تحلیل ها در زمان مطلوب تری انجام می شود. لذا با توجه به استفاده از ساختار فایل هدوپ برای ذخیره سازی داده ها به هیچ سخت افزار اضافی برای ذخیره سازی و یا عملیات پیش پردازش نیاز نمی باشد. در این رویکرد به سخت افزار جدید و قدرتمند نیازی نبوده و در هر گام امکان اجرای درخواست های جامع و کامل بر مبنای داده های ذخیره شده در زمان مطلوب تر نسبت به پایگاه داده های رابطه ای فراهم است. از آنجا که در گام اول داده ها و جداول لازم برای انجام درخواست های تعیین شده در یک پایگاه داده ی رابطه ای ذخیره شده اند، لذا باید تمامی این داده ها یا بخشی از آن ها که مورد نیاز تحلیل گران و کاربران نهایی است به دریاچه داده در یک پایگاه داده ی غیررابطه ای منتقل و درخواست ها روی آن ها اجرا شود. در این صورت اجرای درخواست ها روی داده های موجود در دریاچه داده ی پایگاه داده ی غیررابطه ای به مراتب سریعتر از پایگاه داده ی رابطه ای خواهد بود. معماری ارائه شده در قالب کلی به صورت زیر است :



#### شکل (۴): معماری پیشنهادی

### ابزارهای پیاده سازی

در این معماری پیکربندی و ایجاد دریاچه داده برای ذخیره سازی داده های مورد نیاز در نظر گرفته شده است. در اولین قدم باید پایگاه داده ای منبع که حاوی داده های خام هستند، برای استخراج و انتقال داده ها به دریاچه داده مورد تجزیه و تحلیل قرار گیرند. داده های پایگاه داده ای منبع با تبدیل به فرمت استاندارد CSV قابلیت بارگذاری و انتقال به پایگاه داده ای غیررابطه ای را دارند. در واقع تمامی این مراحل را به شکل ساده تر شامل عملیات استخراج، تغییر شکل و بارگذاری می توان در نظر گرفت. داده های اصلی ما در این تحقیق روی نرم افزار مدیریت پایگاه داده ای SQLSERVER قرار گرفته است و از سرویس های این نرم افزار برای تغییر شکل داده ها به فرم استاندارد CSV استفاده می شود.

همانطور که گفته شد پایگاه داده ای غیررابطه ای که به صورت توزیع شده راه اندازی شده است، باید سه ویژگی سازگاری، در دسترس بودن و تحمل تقسیم بندی را دارا باشد. برای مقیاس پذیری افقی باید تحمل تقسیم بندی آن قوی بوده و از ویژگی های سازگاری و در دسترس بودن تنها یکی را انتخاب نمود. زیرا:

- اگر سیستمی تحمل تقسیم بندی را دارا نباشد دیگر آن را به عنوان یک سیستم توزیع شده و مقیاس پذیر افقی نمی توان در نظر گرفت.

- این ویژگی ها در تضاد با یکدیگرند و در یک زمان بیشتر از دو مورد آن ها امکان کنار هم بودن را ندارند. لذا با توصیه ای متخصصان تحمل تقسیم بندی انتخاب قطعی بوده و انتخاب گزینه ی بعدی از بین سازگاری و در دسترس بودن کامل به ماهیت نرم افزار و اولویت های سازمانی بستگی دارد.

وضعیت انتخاب یکی از دو گزینه ی سازگاری و یا در دسترس بودن به شرح زیر است :

- انتخاب گزینه ی در دسترس بودن و تحمل تقسیم بندی (AP): انتخاب این دو گزینه به معنی این است که هر درخواستی حتما پاسخی دریافت خواهد کرد و این پاسخ تا حد ممکن شامل جدیدترین اطلاعات خواهد بود. اما شاید این اطلاعات قدیمی باشند. از طرفی شاید نوشتن اطلاعات جدید در لحظه روی گره ها ممکن نباشد. اما سیستم همچنان به فعالیت خود ادامه داده و در نهایت به ثبات اطلاعات خواهد رسید. در نهایت انتخاب این دو گزینه بیانگر این است برای سیستم در دسترس بودن و سرعت بیشتر آن اولویت بالاتری نسبت به ثبات و دوام اطلاعات دارد.
- انتخاب گزینه ی سازگاری و تحمل تقسیم بندی (CP): انتخاب این دو گزینه به معنی این است که هر درخواستی برای دریافت اطلاعات به طور قطعی و دقیق باید آخرین نسخه ی موجود را به عنوان پاسخ ارسال نماید. از طرفی این آخرین نسخه ی موجود از هر اطلاعاتی روی سیستم، به طور کامل مشخص و واضح است. یعنی اگر دو درخواست یکسان از دو مکان مختلف به سیستم ارسال شود، سیستم جواب یکسان و مشخصی را به هر دو درخواست خواهد داد. نکته ی دیگر این است که افزایش سازگاری زمان رکود سیستم را افزایش می دهد. هر میزان که سازگاری سیستم افزایش داده شود سیستم دچار رکود شده و سرعت پاسخگویی کاهش می یابد.

از آن جا که برای ارائه ی تحلیل های دقیق و جامع از وضعیت آموزشی و تحصیلی، همواره باید آخرین و بروزترین اطلاعات و داده های موجود در دسترس باشند، لذا برای انتخاب پایگاه داده ی غیررابطه ای مناسب ویژگی های سازگاری و تحمل تقسیم بندی (CP) را می بایست در نظر داشت. از این رو درمیان پایگاه داده های غیررابطه ای که از این ویژگی برخوردار هستند با در نظر گرفتن رده بندی شکل زیر پایگاه داده ی غیررابطه ای HBase برای پیاده سازی انتخاب شد.

☐ include secondary database models

11 systems in ranking, September 2020

Rank			DBMS	Database Model	Score		
Sep 2020	Aug 2020	Sep 2019			Sep 2020	Aug 2020	Sep 2019
1.	1.	1.	Cassandra +	Wide column	119.18	-0.66	-4.22
2.	2.	2.	HBase +	Wide column	48.35	-0.76	-7.37
3.	3.	3.	Microsoft Azure Cosmos DB +	Multi-model i	31.67	+0.94	+0.80
4.	4.	4.	Datastax Enterprise +	Wide column, Multi-model i	8.51	-0.13	-0.12
5.	5.	5.	Microsoft Azure Table Storage	Wide column	5.77	+0.19	+1.61
6.	6.	6.	Accumulo	Wide column	4.11	+0.09	+0.07
7.	7.	7.	Google Cloud Bigtable	Wide column	3.67	+0.29	+1.82
8.	8.	8.	ScyllaDB +	Multi-model i	2.59	+0.10	+0.96
9.	9.	9.	MapR-DB	Multi-model i	0.79	+0.05	+0.17
10.	10.	↑ 11.	Alibaba Cloud Table Store	Wide column	0.43	+0.05	+0.20
11.	11.		Elassandra	Wide column, Multi-model i	0.39	+0.06	

شکل (۵): معماری پیشنهادی

### نحوه ی پیاده سازی

برای ارزیابی معماری پیشنهادی پایگاه داده های آموزش و پرورش استفاده شد. این داده ها در پایگاه داده ی رابطه ای SQLSERVER ذخیره شده است. در پایگاه داده ی غیررابطه ای میان جداول رابطه وجود ندارد. لذا همه ی اطلاعات جداول مرتبط در پایگاه داده ی رابطه ای باید در یک جدول پایگاه داده ی غیررابطه ای ذخیره سازی شوند. از این رو جدول پایگاه داده ی غیررابطه ای دارای ستون های زیر خواهد بود :

جدول (۲): ستونهای سند پایگاه داده ی غیررابطه ای

نام فیلد	نام فیلد	طول فیلد
کد دانش آموز	StudentCode	۱۲
نام دانش آموز	FirstName	۱۰۰
نام خانوادگی دانش آموز	LastName	۱۰۰
کد ملی	NationalCode	۱۰

نام پدر	FatherName	۱۰۰
کد مدرسه	SchoolCode	۱۰
سال تحصیلی	TimeYear	۱۲
کد درس	CourseCode	۵
نام درس	CourseName	۱۰۰
نمره	Val	۵

این داده‌ها در پایگاه داده‌ای رابطه‌ای در جداول جداگانه‌ای نرمال ذخیره‌سازی شده‌اند. با الحاق جداول به شکل مناسب در محیط SQLSERVER اطلاعات مورد نیاز از پایگاه داده استخراج شده و در قالب CSV ذخیره می‌شوند. پس از راه اندازی و پیکربندی پایگاه داده‌ی غیررابطه‌ای و ایجاد بانک اطلاعاتی و جداول مورد نیاز، داده‌های فوق به این پایگاه بارگذاری شده و درخواست‌ها روی آن‌ها اجرا شد. برای اجرای درخواست‌ها و مقایسه زمان اجرا در این تحقیق از داده‌های آزمایشی استفاده شد.

### محیط اجرا

در این تحقیق برای پیاده‌سازی و ارزیابی معماری پیشنهادی، پایگاه داده‌ی غیررابطه‌ای HBase استفاده شد که از سیستم فایل هدوپ (HDFS) برای ذخیره‌سازی داده‌ها استفاده می‌کند. برای شبیه‌سازی هدوپ و راه اندازی پایگاه داده‌ی HBase از محیط مجازی کلاودرا نسخه ۵.۵ استفاده شد.



شکل (۶): ماشین مجازی کلاودرا

درخواست‌هایی که برای مقایسه‌ی زمان اجرا در پایگاه داده‌ی رابطه‌ای SQLSERVER و پایگاه داده‌ی غیررابطه‌ای HBase در نظر گرفته شده‌اند عبارتند از :

- محاسبه‌ی میانگین نمرات یک درس دانش‌آموزان به تفکیک سال تحصیلی
- محاسبه‌ی میانگین نمرات یک درس دانش‌آموزان به تفکیک آموزشگاه
- محاسبه‌ی میانگین نمرات یک درس دانش‌آموزان به تفکیک سال تحصیلی و آموزشگاه
- محاسبه‌ی میانگین نمرات یک درس دانش‌آموزان تاکنون
- محاسبه‌ی معدل دانش‌آموزان به تفکیک سال تحصیلی
- محاسبه‌ی معدل دانش‌آموزان به تفکیک آموزشگاه محل تحصیل



- محاسبه‌ی معدل دانش‌آموزان به تفکیک سال تحصیلی و آموزشگاه
- محاسبه‌ی معدل دانش‌آموزان تاکنون

داده‌های آزمایشی این تحقیق روی پایگاه داده‌ی رابطه‌ای SQLSERVER در جداول مجزا به شکل نرمال ذخیره سازی شده و حاوی یکصد هزار رکورد از اطلاعات دانش‌آموزان می باشد. در ابتدای راه اندازی ماشین مجازی، حافظه‌ی رم اختصاص یافته به آن در این پژوهش ۴۰۹۶ مگابایت است. از طرفی برای ایجاد شرایط یکسان مقایسه، پایگاه داده‌ی رابطه‌ای SQLSERVER نیز روی یک ماشین مجازی با تنظیمات مشابه نصب و راه اندازی شده و داده‌ها به آن منتقل شدند. کد محاسبه‌ی زمان اجرای درخواست‌های پایگاه داده‌ی رابطه‌ای در شکل زیر نشان داده شده است :

```

DECLARE @Time1 DATETIME
DECLARE @Time2 DATETIME
SET @Time1 = GETDATE()
SET @Time2 = GETDATE()
SELECT [Firstname]
      ,[LastName]
      ,[NationalCode]
FROM [TestDb].[dbo].[Students]

SET @Time2 = GETDATE()
SELECT DATEDIFF(MILLISECOND,@Time1,@Time2) AS Elapsed_MS
PRINT DATEDIFF(MILLISECOND,@Time1,@Time2);

```

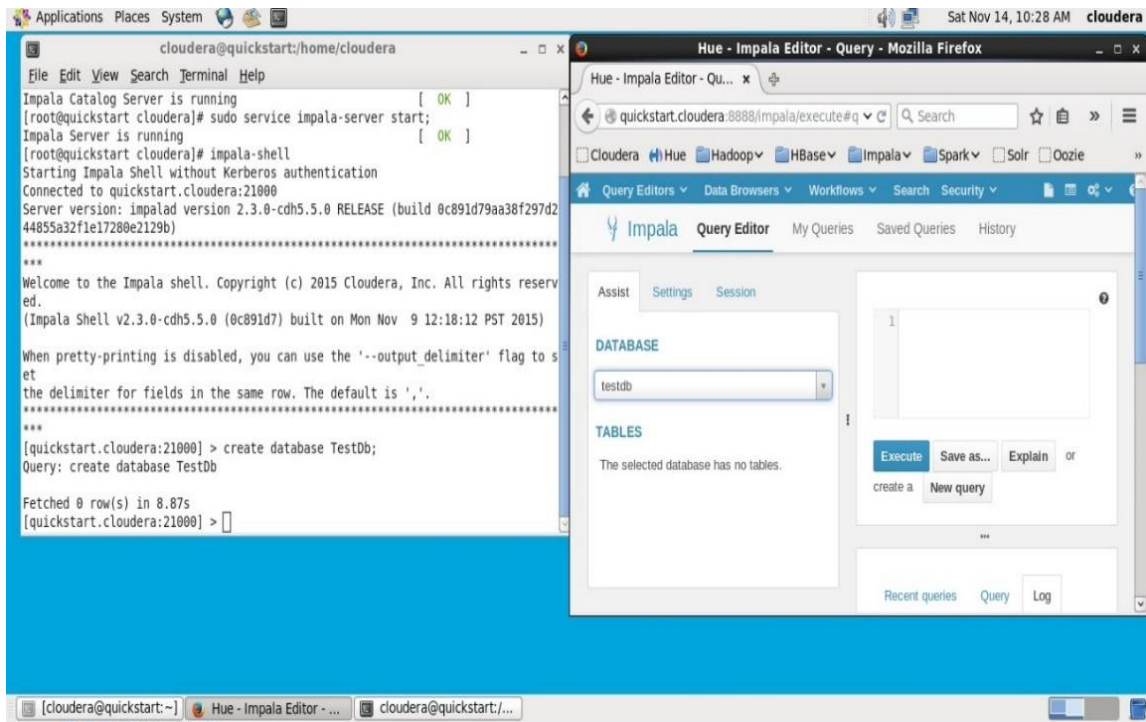
Firstname	LastName	NationalCode
رضا	قاسمی پور	640752773
رضا	حسن آهانی	640752986
مهدی	مؤنن	640753966
روح الله	عسکری	640754450
عباس	جضری	640755861
محمد رضا	احمدی	640756070
سینجه‌مندی	حسینی	640773390
مهدی	رضائی	640820592
...	...	...

Elapsed_MS
780

Query executed successfully. DESKTOP-5KC084F (15.0 RTM) | sa (53) | TestDb | 00:00:00 | 100,001 rows

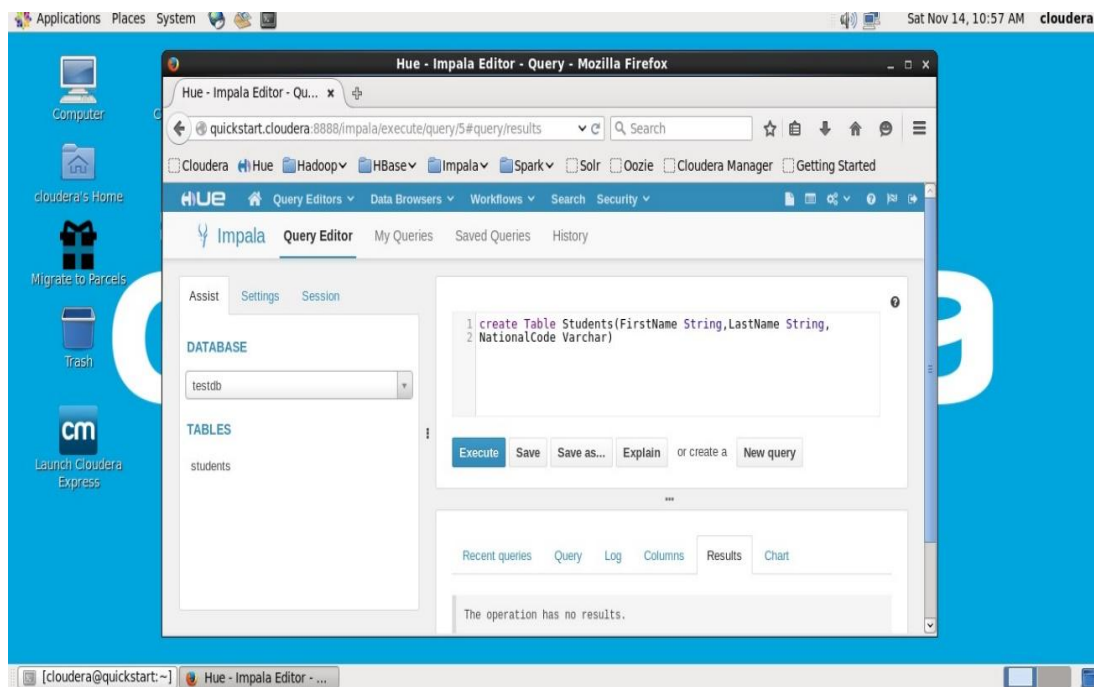
شکل (۷): اجرای درخواست‌ها در محیط پایگاه داده‌ی رابطه‌ای

راه اندازی سرویس‌های ایمپالا، پایگاه داده و جداول در پوسته‌ی آن یا مرورگر Hue به صورت زیر است:



شکل (۸): راه اندازی پایگاه داده در ایمپالا

اکنون می توان داده ها را از پایگاه داده ی رابطه ای در قالب فایل CSV استخراج و در جداول پایگاه داده ی غیررابطه ای با اجرای دستور انتقال، بارگذاری نمود.



شکل (۹): ایجاد جدول در پایگاه داده ی ایمپالا



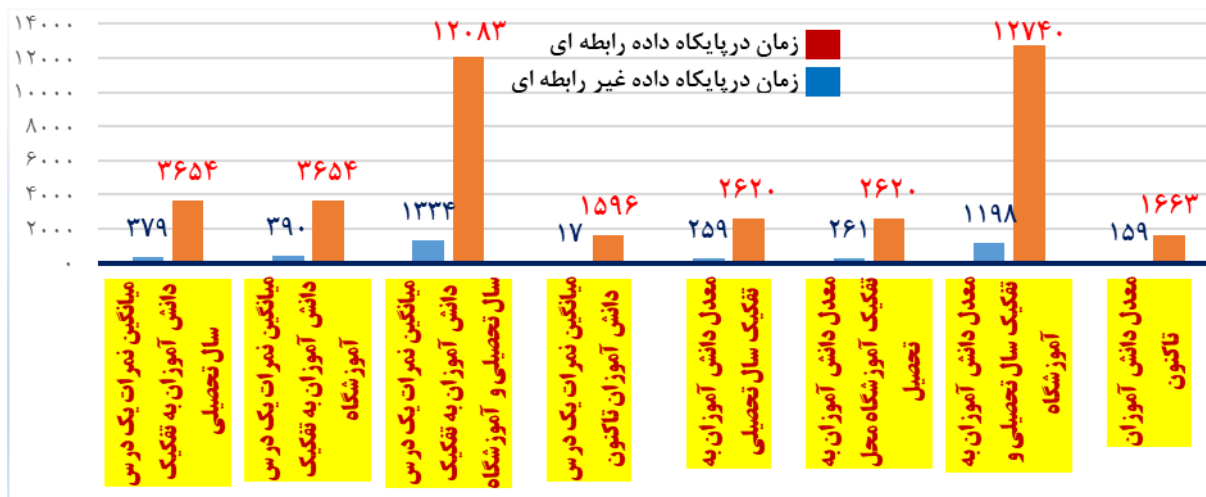
## ارزیابی

برای ارزیابی معماری پیشنهادی، تعدادی از درخواست‌های پرکاربرد برای نمونه انتخاب شده و روی پایگاه‌داده‌ی رابطه‌ای اجرا شد و سپس همان درخواست‌ها روی پایگاه‌داده‌ی غیررابطه‌ای نیز اجرا گردید و زمان اجرای آنها با یکدیگر مقایسه شد. زمان اجرای این درخواست‌ها در جدول زیر با یکدیگر مقایسه شده است :

جدول (۳): مقایسه‌ی زمان اجرا

عنوان گزارش		زمان (میلی ثانیه)	
پایگاه داده‌ی رابطه‌ای	پایگاه داده‌ی غیررابطه‌ای		
۳۶۵۴	۳۷۹	محاسبه‌ی میانگین نمرات یک درس دانش‌آموزان به تفکیک سال تحصیلی	
۳۶۵۴	۳۹۰	محاسبه‌ی میانگین نمرات یک درس دانش‌آموزان به تفکیک آموزشگاه	
۱۲۰۸۳	۱۳۳۴	محاسبه‌ی میانگین نمرات یک درس دانش‌آموزان به تفکیک سال تحصیلی و آموزشگاه	
۱۵۹۶	۱۷	محاسبه‌ی میانگین نمرات یک درس دانش‌آموزان تاکنون	
۲۶۲۰	۲۵۹	محاسبه‌ی معدل دانش‌آموزان به تفکیک سال تحصیلی	
۲۶۲۰	۲۶۱	محاسبه‌ی معدل دانش‌آموزان به تفکیک آموزشگاه محل تحصیل	
۱۲۷۴۰	۱۱۹۸	محاسبه‌ی معدل دانش‌آموزان به تفکیک سال تحصیلی و آموزشگاه	
۱۶۶۳	۱۵۹	محاسبه‌ی معدل دانش‌آموزان تاکنون	

نمودار مقایسه‌ی زمان اجرا در پایگاه‌داده‌ی رابطه‌ای و غیررابطه‌ای این تحقیق در زیر ارائه شده است :



نمودار (۱): نمودار مقایسه‌ی زمان درخواست‌ها در پایگاه‌داده‌ی رابطه‌ای و غیررابطه‌ای

## نتیجه‌گیری

همانطور که از نتایج اجرای درخواست‌ها روی پایگاه‌داده‌ی رابطه‌ای و غیررابطه‌ای و مقایسه‌ی زمان اجرا در آنها استنباط می‌شود، اجرای درخواست‌ها روی پایگاه‌داده‌ی غیررابطه‌ای سریعتر از پایگاه‌داده‌ی رابطه‌ای انجام می‌شود. همچنین به دلیل اجرای درخواست‌های همزمان روی گره‌های مختلف در فرآیند نگاشت-کاهش، زمان اجرا به نسبت یک سرور با حجم داده‌ی بالا

بسیار کمتر بوده و به دلیل ویژگی مقیاس پذیری پایگاه داده های غیررابطه ای نگرانی در زمینه ی افزایش حجم داده ها و کاهش سرعت اجرای درخواست ها آن گونه که در پایگاه داده های رابطه ای مدنظر است وجود نخواهد داشت. راه اندازی پایگاه داده ی غیررابطه ای به دلیل سرعت بالای اجرای درخواست ها و مقیاس پذیری که امکان ذخیره سازی حجم بیشتری از داده ها را فراهم می کند، در مقایسه با تجهیز و ارتقای سرورهای پایگاه داده های رابطه ای به مراتب کاربردی تر و در بازه زمانی طولانی تر مقرون به صرفه تر خواهد بود. گرچه درخواست های اجرا شده در این پژوهش شامل داده های عددی و اجرای عملیات محاسباتی بود، لیکن در مورد داده های آماری نیز تا حدودی این جنبه از تفاوت قابل توجه است. همان گونه در فصل های قبل این پژوهش مشاهده گردید، پایگاه داده های غیررابطه ای برای پردازش حجم بالای داده ها و ذخیره سازی اقسام مختلف داده توسعه یافته اند. به کارگیری این پایگاه داده ها برای ایجاد منبع عظیم داده که از آن به عنوان دریاچه داده یاد می شود علاوه بر افزایش سرعت در پردازش و اجرای درخواست ها، به دلیل خاصیت مقیاس پذیری می تواند به نسبت پایگاه داده های رابطه ای حجم بیشتری از داده ها را در خود ذخیره سازی نمایند. در آموزش و پرورش حجم بالای داده ها و سرعت بالای تولید داده، نیاز به دسترسی به داده های بیشتر به منظور تحلیل جامع و برنامه ریزی های مدیریتی می توان از پایگاه داده ی غیررابطه ای برای ذخیره سازی حجم زیادی از داده ها و اجرای درخواست ها در زمان مطلوب تر استفاده نمود. وجود معماری دریاچه داده افق های تازه ای را در مورد مباحث مربوط به ذخیره سازی، امنیت و پردازش موازی، تکنیک های استخراج و بروزرسانی داده ها و همچنین تکنیک ها و ابزارهای داده کاوی نوین فراهم می نماید. از ایده های پژوهشی بیشتر در همین زمینه می توان به موارد زیر اشاره نمود :

- استفاده از پایگاه داده های غیررابطه ای مقیاس پذیر برای دسترسی و ذخیره سازی همه ی داده های آموزش و پرورش و اجرای درخواست ها در زمان مطلوب تر
- به کارگیری هوش مصنوعی در فرآیند داده کاوی و تحلیل و ارزیابی داده ها
- استفاده از ابزارهای پردازش ابری
- استفاده از ابزارها و سیستم های امنیت و کنترل جامعیت داده ها و حفظ پیوستگی و صحت داده ها- ترمیم و حذف داده های تخریب شده

## منابع

- [۱] A. Mohamed, M.K. Najafabadi, B.W. Yap, E.A. Kamaru-Zaman, and R. Maskat, "The state of the art and taxonomy of big data analytics: view from new big data framework". AI Review, ۲۰۱۹, ۱-۴۹.
- [۲] <https://avora.com/blog/rise-of-the-data-warehouse/>
- [۳] S. Singh and N. Singh, "Big Data Analytics", ۲۰۱۲ International Conference on Communication, Information & Computing Technology Mumbai India, IEEE, October ۲۰۱۱
- [۴] B. Gerhardt, K. Griffin and R. Klemann, "Unlocking Value in the Fragmented World of Big Data Analytics", Cisco Internet Business Solutions Group, June ۲۰۱۲,
- [۵] C. Tankard, "Big Data Security", Network Security Newsletter, Elsevier, ISSN ۱۳۵۳-۴۸۵۸, July ۲۰۱۲
- [۶] <http://inside-bigdata.com/۲۰۱۳/۰۹/۱۲/beyond-volume-variety-velocity-issue-big-data-veracity>
- [۷] H. Fang, "Managing data lakes in big data era: What's a data lake and why has it become popular in data management ecosystem," in ۲۰۱۵ IEEE International Conference on Cyber Technology in Automation, Control, and Intelligent Systems (CYBER), ۲۰۱۵: IEEE, pp. ۸۲۰-۸۲۴.
- [۸] P. Kortvelyesi, "Top ۷ challenges of building a data lake," ed, ۲۰۱۷.
- [۹] R. Hai, S. Geisler, and C. Quix, "Constance: An intelligent data lake system," in Proceedings of the ۲۰۱۶ International Conference on

Management of Data, ۲۰۱۶: ACM, pp. ۲۰۹۷-۲۱۰۰.

[۱۰] IBM, "IBM Industry Model support for a data lake architecture," ۲۰۱۶.

[۱۱] O. Mendelevitch, C. Stella, and D. Eadline, Practical Data Science with Hadoop and Spark: Designing and Building Effective Analytics at Scale.

Addison-Wesley Professional, ۲۰۱۶.

[۱۲] <https://www.grazitti.com/blog/data-lake-vs-data-warehouse-whichone-should-you-go-for/>

[۱۳] Ibrahim Abaker Targio Hashem, Ibrar Yaqoob, Nor Badrul Anuar, Salimah Mokhtar, Abdullah Gani, Samee Ullah Khan, The rise of "big data" on cloud computing: Review and open research issues, Information Systems, Volume ۴۷, January ۲۰۱۵, Pages ۹۸-۱۱۵, ISSN ۰۳۰۶-۴۳۷۹

[۱۴] C. Eaton, D. Deroos, T. Deutsch, G. Lapis and P.C Zikopoulos, Understanding Big Data Analytics for Enterprise Class Hadoop and Streaming Data, Mc Graw-Hill Companies ۹۷۸-۰-۰۷-۱۷۹۰۵۳-۶, ۲۰۱۲

[۱۵] S. Madden, "From Databases to Big Data", IEEE Internet Computing, June ۲۰۱۲, v.۱۶, pp. ۴-۶

[۱۶] Ishwarappa, Anuradha J, 'A Brief Introduction on Big Data Vs Characteristics and Hadoop Technology', Procedia Computer Science ۴۸(۲۰۱۵)۳۱۹-۳۲۴

[۱۷] Ruoran Liu, Haruna Isah, Farhana Zulkernine, A Big Data Lake for Multilevel Streaming Analytics

[۱۸] H. Fang, "Managing data lakes in big data era: What's a data lake and why has it become popular in data management ecosystem," in ۲۰۱۵ IEEE

International Conference on Cyber Technology in Automation, Control, and Intelligent Systems (CYBER), ۲۰۱۵: IEEE, pp. ۸۲۰-۸۲۴.

[۱۹] R. Hai, S. Geisler, and C. Quix, "Constance: An intelligent data lake system," in Proceedings of the ۲۰۱۶ International Conference on Management of Data, ۲۰۱۶: ACM, pp. ۲۰۹۷-۲۱۰۰.

[۲۰] IBM, "IBM Industry Model support for a data lake architecture," ۲۰۱۶