

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه پیام نور تهران  
واحد شهر ری  
سمینار، سمینار تحقیق و تتبع نظری

عنوان:

مدل‌های ایمن و کارآمد برای دریافت و پردازش اطلاعات  
از پایگاه‌های رمزگذاری شده در فضای ابری

استاد راهنما:

جناب آقای دکتر سید علی رضوی ابراهیمی

نگارش:

مجید لطفی

تابستان ۱۴۰۰

## چکیده

امروزه بحث حفظ و امنیت بانک‌های اطلاعاتی با داشتن اطلاعات حساس و استراتژیک در سطح یک شرکت و حتی بعضاً در سطح کشورها، موردبررسی و تحقیق پژوهشگران بخش نرم‌افزار و داده هست.

بخشی از این امنیت جلوگیری از آسیب و افشای اطلاعات از سمت هکرها و گروه‌های خرابکاری بوده و بخش نگران‌کننده دیگر، افشای اطلاعات توسط کاربران و ارائه‌دهندگان سرویس‌های ابری هست زیرا گروه دوم از لایه‌های سخت‌افزاری و نرم‌افزاری رایج برای امنیت، گذر کرده و دسترسی آزاد به منابع اطلاعاتی دارند.

سرویس‌دهندگان ابری و کاربران یک سیستم، به علت داشتن دسترسی آزاد به اطلاعات، براحتی می‌توانند به اطلاعات طبقه‌بندی شده و حساس یک بانک اطلاعاتی دسترسی پیدا نمایند و تبعات سنگینی را برای صاحبان بانک اطلاعاتی ایجاد کنند. بهترین و شاید تنهاترین ایده ای که به ذهن می‌رسد رمزگذاری اطلاعات هست تا این اطلاعات تنها پس از کنترل دسترسی‌ها داخل نرم‌افزار، به نمایش دربیایند.

اما رمزگذاری و به تبع آن رمزگشایی نیاز به محاسبات سنگین و زمان مشغولی سرور دارند و اگر نتوانیم با روش صحیح و مناسب این روش را جلو ببریم، در حجم اطلاعات بالا و پس از مدتی، استفاده از سامانه سخت و حتی غیرمقدور می‌شود.

آشنایی با انواع روش‌های رمزگذاری، میزان بار و هزینه اجرا آن و تعیین روش صحیح برای رمزگذاری، به ما کمک می‌کند به یک راهکار و روش مناسب برای حفظ محرمانگی اطلاعات برسیم.

**کلمات کلیدی:** بانک اطلاعاتی، رمزگذاری، رمزگشایی، سرور ابری، سرویس‌دهنده ابری

## فهرست مطالب

صفحه	عنوان
	فصل اول: مقدمه
۶	مقدمه
۶	۱-۱ تعریف مسئله و سؤال
۷	۲-۱ ضرورت های تحقیق
۷	۳-۱ اهداف
۷	۴-۱ جمع بندی
	فصل دوم: آشنایی با رمزگذاری
۸	مقدمه
۸	۱-۲ روش های رمزگذاری
۹	۲-۲ رمزگذاری متقارن
۹	۳-۲ رمزگذاری نا متقارن
	فصل سوم: رمزگذاری در پایگاه داده
۱۱	مقدمه
۱۱	۱-۳ رویکردهای فعلی رمزگذاری در پایگاه داده
۱۲	۲-۳ معرفی CryptDB
۱۲	۳-۳ معرفی FHOPE
۱۲	۴-۳ معرفی P-McDb
۱۲	۵-۳ معرفی SDB
۱۳	۶-۳ معرفی CASB
۱۳	۷-۳ سایر تحقیقات
	فصل چهارم: استراتژی های رمزگذاری پایگاه داده
۱۴	مقدمه
۱۴	۱-۴ استراتژی رمزگذاری
۱۷	۲-۴ روش های قابل اجرای رمزگذاری

## فصل پنجم: آزمایش و ارزیابی

۲۱	مقدمه
۲۱	۱-۵ آزمایش و ارزیابی
۲۲	۲-۵ ارزیابی
۳۲	۳-۵ بررسی محاسبات Cloud و QM
۳۳	۴-۵ بررسی فضای موردنیاز مدل‌ها

## فصل ششم: بحث و پیشنهادها

۳۴	مقدمه
۳۴	۱-۶ بحث و پیشنهادها
۳۶	۲-۶ نتیجه‌گیری

# فصل اول

## مقدمه

امروزه اکثر شرکت‌ها به موازات اشتراک فایل در شبکه‌های اجتماعی و فضای ابری، به سمت استفاده از فضای ابری جهت نگهداری اطلاعات و دیتای خود رفته‌اند زیرا از سوئی می‌توان با اجاره سرویس‌های آماده، هزینه تهیه سرور را کاهش داد و از طرف دیگر امکان دسترسی هر زمان از دیتاها را فراهم نمود.

به موازات گسترش استفاده از این بستر، آسیب‌ها و مشکلات آن نیز مطرح شده است و باید بتوان این سرویس را در امنیت و با کمترین مشکل اجرا نمود. این مشکلات گاهی مربوط به حملات به سرور و سیستم‌های عامل هست با قصد از کار انداختن سرویس‌ها و گاهی مربوط به نفوذ غیرمجاز به دیتاها هست پس جدا از بحث‌های امنیت شبکه و سیستم‌عامل، خود دیتا نیز حفاظت شود.

حفاظت از داده خارج از روش‌های نرم‌افزاری مانند تعریف نقش‌ها در پایگاه داده و کنترل این نقش‌ها، باید با روش‌هایی مانند رمزگذاری حفاظت شود زیرا سرویس‌دهنده‌های و مالکان سرورهای ابری امکان دسترسی به دیتای مشتریان خود را داشته می‌توانند از این اطلاعات حساس استفاده غیرمجاز نمایند.

رمزگذاری در دو بخش ساختار و دیتا قابل اجرا هست و اجرا هم‌زمان هر دو بخش می‌تواند بالاترین امنیت را ایجاد نماید زیرا با فرض امکان و یا افشای دسترسی به پایگاه داده، داشتن نام جداول و یا ستون‌ها متداول باعث آسیب‌پذیری بیشتر اطلاعات می‌شود برای مثال اگر جدولی با نام Banks داشته باشیم قطعاً بیان‌کننده نام بانک‌ها هست درحالی‌که اگر همین جدول به صورت T۱۲ باشد، امکان افشای اطلاعات کاهش و حتی غیرممکن می‌شود.

## ۱-۱ تعریف مسئله و سؤال

با توجه به آنچه تا اینجا بیان شد، مسئله اصلی ما حفاظت داده‌ها از بالاترین سطح که سرویس‌دهنده ابری هست و دسترسی کامل به اطلاعات دارد تا پایین‌ترین سطح که افراد غیرمجاز و هکرها هست، تعریف می‌شود. این حفاظت به‌گونه‌ای هست که حتی اگر شخص بتواند دیتا را در به‌صورت مستقیم از جداول مشاهده نماید، باز با خطر افشای اطلاعات مواجه نشویم.

سؤالاتی که تا پایان این تحقیق پاسخ داده می‌شود را این‌گونه بیان می‌کنیم.

۱- انواع روش‌های رمزگذاری قابل‌اجرا در پایگاه‌های داده کدام‌اند؟

۲- انواع رویکردهای رمزگذاری در پایگاه‌های داده کدام‌اند؟

۳- بررسی فرآیندها و عملیات پایگاه داده رمزگذاری شده چگونه هست؟

۴- بررسی و مقایسه روش‌ها و رویکردهای رمزگذاری از نظر هزینه اجرا و نگهداری چگونه‌اند؟

## ۲-۱ ضرورت تحقیق

تا اینجا با مسائل و مشکلات یک پایگاه داده از نظر افشای اطلاعات آشنا شدیم و رمزگذاری به‌عنوان راهکار آن معرفی گردید. حال باید این موضوع رو در نظر بگیریم که صرف رمزگذاری و ایجاد لایه‌های رمزگذاری مسئله ما نیست زیرا هرچقدر لایه‌های رمزگذاری اضافه شود، به همان اندازه لایه‌های به دست آوردن دیتاها و هزینه دریافت اطلاعات برای خودمان بالاتر می‌رود زیرا این عبارات رمزگذاری شده باید برای پردازش و استفاده مجدد، رمزگشایی شوند و ما باید سعی کنیم رمزگذاری را بر روی دیتای حساس لحاظ نماییم و نه بر روی تمام پایگاه داده. ایده و نگرش دیگری که در این پایان‌نامه بر آن تأکید شده، کاهش بار محاسبات رمزگذاری و رمزگشایی هست، به این صورت که در هنگام دریافت اطلاعات در یک پرس‌وجو، نیاز به رمزگشایی اطلاعات نامرتب نباشد و سعی شود تا حدودی از میزان اطلاعات واکنشی جهت رمزگشایی کاسته شود.

## ۳-۱ اهداف

با آنچه تا اینجا در خصوص افشای اطلاعات و ضرورت حفظ آن بیان شد، ما در این تحقیق چند هدف را موردبررسی قرار می‌دهیم و سعی می‌کنیم در انتها با جمع‌بندی روش و رویکردها رایج در این خصوص بپردازیم. این اهداف عبارت‌اند از:

۱- آشنایی با روش‌ها رمزگذاری بدون در نظر گرفتن پایگاه داده

۲- بررسی و اجرا روش رمزگذاری مؤثر و رایج در پایگاه داده

۳- آشنایی با مدل‌ها و اصول قابل‌اجرا در رمزگذاری پایگاه داده

۴- مقایسه مدل‌ها و انتخاب مدل مناسب

## ۴-۱ جمع‌بندی

در این تحقیق به‌ضرورت حفظ اطلاعات پایگاه داده از افشای آن و رمزگذاری به‌عنوان بهترین روش جهت انجام آن معرفی می‌گردد. این موضوع چنانچه قبلاً نیز اشاره شد، صرفاً برای جلوگیری از حملات مخرب و غیرمجاز نیست و بررسی‌های این تحقیق در خصوص کاربران و اشخاصی هست که به پایگاه داده، دسترسی مجاز دارند. ابتداها به روش‌های رمزگذاری، اصول و قواعد رایج در آن و سپس به رویکردهای رایج در ایجاد یک پایگاه رمزگذاری شده می‌پردازیم و درنهایت مدل‌ها و روش‌های معرفی‌شده را مقایسه تا بهترین مدل و روش را انتخاب نماییم.

# فصل دوم

## مقدمه

رمزگذاری روند رمز کردن پیام‌ها یا اطلاعات است به گونه‌ای که تنها افراد مجاز قادر به خواندن آن باشند. پیام یا اطلاعات با استفاده از یک الگوریتم، رمزگذاری شده و علائم رمزی به وجود می‌آید که فقط در صورت رمزگشایی قابل خواندن هستند. در رمزگذاری معمولاً یک کلید رمزگذاری شبه تصادفی تولید شده توسط یک الگوریتم، به کار گرفته می‌شود. اگرچه شاید رمزگشایی پیام بدون در اختیار داشتن کلید ممکن باشد، اما در یک رمزگذاری خوب، منابع محاسباتی زیادی برای این کار لازم است. یک گیرنده مجاز به راحتی می‌تواند پیام را با کلید تدارک دیده شده توسط صادرکننده پیام، رمزگشایی کند اما گیرنده غیرمجاز نمی‌تواند. هدف از رمزگذاری اطمینان از این است که فقط کسانی که مجاز به دستیابی اطلاعات هستند، قادر به خواندن آن و استفاده از کلید رمزگذاری باشند.

## ۱-۲ روش‌های رمزگذاری :

رمز عبارت است از تبدیل کاراکتر به کاراکتر یا بیت به بیت بدون آن که به محتویات زبان شناختی (ادبیات) آن پیام توجه شود. رمزگذاری دیجیتال، متن قابل خواندن (که اصطلاحاً به آن plaintext می‌گویند) را می‌گیرد و آن را به متنی غیرقابل تشخیص تبدیل می‌نماید و اصطلاحاً آن را رمزگذاری می‌کند. البته برای انجام چنین کاری، الگوریتم‌های بسیار پیچیده و قدرتمند وجود دارد.

این الگوریتم‌ها از متغیرهایی به نام کلید بهره می‌برند که پیچیدگی رمزگذاری را چندین برابر می‌کنند. این کلیدها به طور تصادفی تولید می‌شوند و منحصر به فرد هستند. یعنی اگر یک هکر بخواهد بانفوذ به یک بانک داده اطلاعات کارت اعتباری کاربران را بدزدد، نه تنها باید از الگوریتم استفاده شده خبر داشته باشد بلکه باید بداند کدام کلید استفاده شده است. چنین کاری به هیچ عنوان کار راحتی نیست و به همین خاطر رمزگذاری‌های دیجیتال در برابر حملات جستجوی فراگیر (brute force) مقاومت بسیار بالایی دارند. دو روش برای رمزگذاری قابل اجرا هست : روش متقارن و روش نامتقارن.



## ۲-۲ رمزگذاری متقارن :

در روش متقارن تنها یک کلید برای رمزگذاری استفاده می‌شود و فرستنده و گیرنده اطلاعات از همان کلید مشترک برای رمزگذاری و رمزگشایی داده استفاده می‌کنند. تمام رمزنگاری‌های کلاسیک از نوع متقارن هستند و تا قبل از دهه ۷۰ تنها نوع رمزنگاری به حساب می‌آمد.

برخی از الگوریتم‌های متداول در رمزگذاری بر مبنای روش متقارن به صورت ذیل هست:

۱- Order-Preserving Encryption: رمزگذاری با حفظ ارزش عبارت رمزگذاری نسبت به اصل عبارت

مثلاً اگر رمزگذاری دو عبارت  $A, B$  بشود  $E_1, E_2$  و  $A > B$  باشد قطعاً  $E_1 > E_2$  و برعکس و این مقایسه بدون رمزگشایی معتبر هست و مناسب برای دیتای عددی و این مزیت خود منجر به ضعف شده و می‌توان توسط استنباط تا حدودی به اطلاعات اصلی رسید.

۲- Deterministic Encryption: یا همان DES الگوریتم رمزگذاری‌ای هست که توسط یک کلید، همواره

یک عبارت با به متن مشخصی رمز می‌نماید و متأسفانه ارزش عبارت اصلی را حفظ نمی‌نماید و در صورت افشای کلید، به راحتی می‌توان سایر عبارات را رمزگشایی کرد! با شکسته شدن الگوریتم DES این استاندارد در سال ۱۹۹۸ تمدید نشد.

۳- Advanced Encryption Standard: یک الگوریتم رمزنگاری بلوک جایگزینی شبکه (SPN) است و

جایگزین DES شد. دقیقاً به این معنی است که الگوریتم یک بلوک از متن ساده را می‌گیرد و مقادیر متناوب جایگزینی را به آن اعمال می‌کند. AES شامل سه بلوک از رمزهای AES-۱۲۸، AES-۱۹۲ و AES-۲۵۶ است. هر رمز، عملیات رمزنگاری و رمزگشایی را بر روی بلوک‌های ۱۲۸ بیتی داده، با استفاده از کلیدهای ۱۲۸، ۱۹۲ و ۲۵۶ بیتی خود انجام می‌دهد و محبوبیت بالایی دارد.

۴- Blowfish: یک رمزنگاری قطعه‌ای سریع محسوب می‌شود که داده‌ها در قطعه‌های ۸ بیتی رمز می‌کند،

به جز هنگامی که کلیدها عوض می‌شوند. هر کلید جدید نیازمند پیش-پردازشی است که معادل رمزگذاری یک فایل متنی تقریباً ۴ کیلوبایتی است، که نسبت به دیگر رمزنگاری‌های قطعه‌ای بسیار کند است.

## ۲-۳ رمزگذاری نامتقارن :

در روش نامتقارن از دو کلید استفاده می‌نماید، یک کلید عمومی و یک کلید خصوصی که کلید عمومی یک عبارت رمزگذاری شده جهت رمزگذاری در سطح کاربر است که خود کلید عمومی توسط کلید خصوصی رمزگذاری شده است. این سطح از رمزگذاری منجر به کاهش سرعت و افزایش هزینه استفاده از سیستم می‌شود.

برخی از الگوریتم‌های متداول در رمزگذاری بر مبنای روش نامتقارن به صورت ذیل هست:

۱- Homomorphic Encryption: نوعی رمزگذاری برای انجام عملیات جمع و ضرب بدون باز کردن عبارت رمزگذاری شده و یا داشتن کلید و مناسب دیتای عددی و دارای قابلیت اجرای دستورات SQL بر روی عبارات رمزگذاری.

۲- Randomized Encryption: این رمزگذاری یک عبارات را هر بار به عبارت رمزگذاری متفاوتی تبدیل می‌نماید و بالاترین سطح امنیت هست اما به علت نوع رمزگذاری قابلیت اجرای دستورات SQL بر روی عبارات رمزگذاری را ندارد.

۳- Onion Encryption: در این رمزگذاری از لایه‌های مختلف با متدهای مختلف رمزگذاری استفاده می‌شود به صورتی که لایه داخلی با متد رمزگذاری امنیت پایین و لایه بیرونی دارای رمزگذاری با امنیت بالا و مشکل آن کند شدن پردازش پرس‌وجو و طولانی شدن عبارات رمزگذاری نسبت به عبارت اصلی هست.

۴- RSA: در این چنین سیستم‌های رمزنگاری، کلید رمزگذاری عمومی است و از کلید رمزگشایی که مخفی است، جداست. هرکسی می‌تواند از این کلید عمومی برای رمزگذاری یک پیام استفاده کند، اما تنها کسی که آن دو عدد اولی که کلید بر اساس آن‌ها ساخته شده را می‌داند، قادر به رمزگشایی پیام است. شکستن رمزگذاری RSA به مسئله‌ی RSA معروف است. تاکنون هیچ روشی برای شکست دادن این سیستم (در صورت استفاده‌ی کلید به اندازه‌ی کافی بزرگ) منتشر نشده است. RSA به صورت نسبی، الگوریتم کندی است و به همین علت، کمتر برای رمزگذاری مستقیم اطلاعات کاربر استفاده می‌شود.

۵- ElGamal: خصوصیت این الگوریتم به گونه‌ای است که در هر مرحله از فرایند رمزنگاری یک کلید تصادفی  $k$  تولید می‌شود مقدار این کلید کاملاً تصادفی بوده و در هر مرحله از اجرای فرایند رمزنگاری کلید تولیدشده با کلید قبلی متفاوت خواهد بود. این خصوصیت موجب می‌شود در دو مرحله متفاوت خروجی متفاوت تولید شود، که این خصوصیت از ویژگی‌های تولید کلید  $K$  در هر مرحله از فرایند رمزنگاری هست. فرایند تولید تصادفی کلید  $K$  موجب می‌شود حدس زدن کلید  $K$  از روی مقدار  $C$  امکان‌پذیر نباشد.

با توجه به تنوع این الگوریتم‌ها ما باید با در نظر گرفتن هزینه اجرا، سرعت و میزان حساسیت اطلاعات، در بخش‌های مختلف پایگاه داده و سیستم نرم‌افزاری، الگوریتم کاربردی‌تر و مناسب‌تر را انتخاب نماییم.

# فصل سوم

## مقدمه

پایگاه داده با توجه به ثبت اطلاعات به صورت متن ساده و داشتن اطلاعات حساس، نیازمند بررسی و مطالعه در جهت حفظ اطلاعات از افشای آن توسط کاربران و صاحبان سرورها دارد. البته نسخه های بالاتر بانک های اطلاعاتی دارای دستورات جهت رمزگذاری و رمزگشایی هست اما ما به دنبال رویکرد و الگوریتم هایی هستیم تا بتوانیم بر روی تمامی پایگاه های داده اجرا نماییم.

امروزه نرم افزارهایی به صورت سرویس دهنده بانک های اطلاعاتی که نقش واسطه بین کاربر و بانک اطلاعاتی دارند، با ایجاد قابلیت مدیریت کدگذاری بانک اطلاعاتی توانسته اند این موضوع را مدیریت کنند اما اگر بخواهیم بانک اطلاعاتی خود را با یک استراتژی قابل تنظیم، از افشای اطلاعات محافظت نماییم، نیازمند آشنایی با رویکردهای فعلی و تحقیقات انجام شده در این زمینه داریم.

## ۳-۱ رویکردهای فعلی رمزگذاری در پایگاه داده :

یکی از روش های رمزگذاری یک پایگاه داده، تبدیل رکوردها به بلوک رمزگذاری هست به این صورت که هر ردیف از اطلاعات با هر تعداد از ستون، تبدیل به یک دیتا رمزگذاری شده در جدول رمزگذاری شود اما این روش امکان اجرای دستورات SQL را بر روی دیتاهای رمزگذاری شده را از بین می برد و حتماً باید دیتا رمزگشایی و دستور اجرا شود و این در اطلاعات با رکوردهای بالا منجر به هزینه و زمان بالا می شود.

برای حل این مشکل محققان روشی را در نظر گرفتند به این صورت که ستونی دیگر به این جداول رمزگذاری شده اضافه گردد که در آن دسته بندی ردیف ها مشخص باشد و در زمان دریافت اطلاعات، با توجه به دسته بندی مورد نظر، ما تعداد رکورد کمتری را دریافت و برای رمزگشایی پردازش نماییم. این دسته بندی ها درجایی به صورت رمزگذاری شده نگهداری می شود.

البته این روش هم امکان افشای این دسته بندی ها و سنگین تر شدن پردازش ها را دارد.

در اینجا بعضی سیستم های رمزگذاری در پایگاه داده معرفی می گردد که دارای امنیت و کارکرد مناسبی می باشند.

### ۳-۲ معرفی CryptDB :

این سیستم به عنوان اولین سیستم عملی برای اجرای دستورات SQL از طریق پایگاه های داده رمزگذاری شده توسعه داد شد که توسط Onion Encryption طراحی گردیده است و هر داده توسط بیش از یک الگوریتم رمزنگاری ، رمزگذاری می شود که در آن رمزهای رمزگذاری بیرونی تولید شده توسط یک الگوریتم رمزگذاری تصادفی انجام می شود. (نمایی از نوع رمزگذاری و پردازش ها در این سیستم)

این نوع رمزگذاری هزینه محاسبات رمزگذاری و رمزگشایی را بالا می برد البته این موضوع توسط MONOMI حل شده است به این صورت که اجرای دستورات به دو بخش تقسیم شده اند : دستوراتی که بر روی اطلاعات رمزگذاری شده قابل اجرا هست و دستوراتی که بر روی اطلاعات رمزگشایی شده در سمت کاربر اجرا می شود. همچنین برای بهبود سرعت رمزگذاری ، از الگوریتم AES-NI بجای AES استفاده شده است که از روش های بهتری برای اجرای رمزگذاری AES استفاده می نماید.

### ۳-۳ معرفی FHOPE :

این سیستم از الگوریتم homomorphic با قابلیت order-preserving برای اجرای دستورات SQL در اطلاعات عددی استفاده می کند و همان طور که قبلاً گفته شد قابلیت order-preserving به ما این امکان را می دهد تا بدون رمزگشایی بتوانیم بر روی داده های رمزگذاری شده کار کنیم. البته این سیستم بر روی تعداد رکورد کم آزمون شده است و باید در دیتاهای بزرگ تر آزمون شود.

### ۳-۴ معرفی P-McDb :

همان طور که گفته شد قابلیت اجرای دستورات بر روی دیتاهای رمزگذاری شده بعضاً منجر به افشای اطلاعات از طریق استنتاج می شود. این سیستم برای جلوگیری از این حمله ، توسط دو سرور که یکی برای ذخیره سازی و جستجوی اطلاعات و دیگری برای تصادفی سازی و تغییر شکل پایگاه داده که این خود منجر به کندی سیستم می شود. این سیستم نیز بجای جستجوی کلی بر روی دیتاهای رمزگذاری شده ، از جستجوی جزئی استفاده می نماید. همچنین در لغو دسترسی یک کاربر به بانک اطلاعاتی برخلاف سایر طرح های SSE (Server Sent Events)، نیاز نیست تمامی اطلاعات مجدد رمزگذاری شود بلکه به تمامی CSP (Content Security Policy) اطلاع داده می شود و دسترسی کاربر بسته می شود و حتی اگر کاربر کلید رمزگشایی و دسترسی به یکی از CSP ها را داشته باشد نیز امکان دسترسی به اطلاعات را نخواهد داشت.

### ۳-۵ معرفی SDB :

یکی از راه های افزایش سرعت و کارایی سیستم در مواجهه با رمزگذاری داده ها ، تقسیم داده به حساس و غیر حساس و رمزگذاری بر روی دیتاهای حساس هست و این سیستم از این روش استفاده می نماید. داده حساس نیز به دو بخش اشتراکی تقسیم می شود، مالک داده ها بخش اول اشتراک و بخش دوم اشتراک را سرویس دهنده ابری نگهداری می کند. در این حالت سرویس دهنده ابری تا زمانی که مالک داده بخش اشتراکی خود را اجرا نماید ، هیچ دسترسی ای به داده ها نخواهد داشت.

### ۳-۶ معرفی CASB :

این سیستم توسط SafeBox تحت عنوان Access Security Broker طراحی شد و به صورت نرم افزار یا سخت افزار داخلی که به عنوان واسطه بین کاربران و ارائه دهندگان خدمات ابری عمل می کند اجرا شده است. این ایده امکان جستجو و اشتراک داده ها و حتی فایل ها با حفظ امنیت آن را فراهم می کند. یکی از امکانات این سیستم جستجوی کلمات کلیدی بر روی اطلاعات رمز گذاری شده را فراهم می کند.

### ۳-۷ سایر تحقیقات :

برخی از محققان از تکنیکی به نام "Bucketization" استفاده کردند که در آن ردیف ها داده به بیش از یک منبع متصل می شود. این فناوری امکان استفاده از پایگاه داده به عنوان یک سرویس دهنده برای اجرای دستورات بر روی داده های رمز گذاری شده را فراهم می کند (Database as service).

برخی از محققان نیز از فضای ابری ترکیبی برای تقسیم داده ها به دو بخش حساس و غیر حساس استفاده نموده که بخش غیر حساس را در فضای عمومی و بخش حساس را فضای خصوصی کاربران قرار می دهند. یکی از اشکالات این روش به این صورت هست که کاربران اکثراً دارای اطلاعات حساس می باشند و از سوی دیگر یکپارچه سازی بین دو بخش مجزا دارای سختی های خود هست.

گروهی دیگر از تکنیکی برای جلوگیری از افشای اطلاعات توسط ارائه دهندگان خدمات ابری غیر قابل اعتماد و مشکوک پیشنهاد کرده اند. این تکنیک بر اساس تقسیم بندی عمودی است، که در آن هر ستون رمز گذاری شده حساس به یک سرور ابر متفاوت برون سپاری می شود. البته این روش دارای تأخیر برای ارتباط بین بخش های مختلف هست خصوصاً اگر دستور SQL پیچیده باشد! واسط این کار یک پراکسی هست که عملیات رمز گذاری و رمز گشایی را انجام می دهد.

با توجه به اینکه تقسیم و برون سپاری داده باعث ایجاد حساسیت ها و نگرانی هایی نسبت به افشای اطلاعات می شود ، گروهی از محققین یک سیستم سخت افزاری / نرم افزاری را برای رفع مشکل نشت محرمانه بودن در پایگاه های داده برون سپاری شده معرفی کرده اند. روش کار به این صورت هست که کاربر کارت هوشمند میانجی را که در کنار آن وصل است نگهداری و کنترل می کند. این کارت هوشمند مسئول رمز گذاری داده ها قبل از قرارداد آن ها در پایگاه داده و رمز گشایی داده ها قبل از ارسال آن ها به کاربر است. عیب عمده این روش این است که کاربر با ظرفیت کارت هوشمند محدود می شود و نمی تواند از فضای ذخیره سازی ارائه شده توسط سرویس های ابری بهره مند شود.

# فصل چهارم

## مقدمه

تا اینجا با انواع رمزگذاری و سیستم‌های رایج رمزگذاری آشنا شدیم و این رمزگذاری را در پایگاه داده مورد بررسی قراردادیم اما در یک بانک اطلاعاتی صرف رمزگذاری اطلاعات تنها می‌تواند در ذخیره اطلاعات مؤثر باشد و ما در یک پایگاه اطلاعاتی عملیات‌های متنوعی مانند دریافت، انتخاب و محاسبه در اطلاعات داریم.

برای ایجاد یک پایگاه داده رمزگذاری شده باید تمامی عملیات‌های قابل اجرا در بانک اطلاعاتی را بررسی نماییم و فرآیندی را پیش‌بینی کنیم تا از ابتدای ذخیره‌سازی، تغییرات و دریافت اطلاعات به مشکل نخوریم. در این فرآیند قصد داریم چنانچه قبلاً گفتیم از یک سیستم واسطه بنام (Query Manager) QM استفاده نماییم که بین کاربران و بانک اطلاعات قرار گرفته و از یک‌سو به کاربران خدمات‌رسانی کند و از سوی دیگر با رمزگذاری اطلاعات از افشای آن‌ها جلوگیری کنیم.

چنانچه بخواهیم تمامی اطلاعات اعم از حساس و غیر حساس را رمزگذاری کنیم، ناچار هزینه بالای ذخیره و خواندن اطلاعات را قبول کردیم اما اگر تنها اطلاعات حساس را رمزگذاری کنیم کمی از بار هزینه کم می‌شود. در این رویکرد یا اطلاعات رمزگذاری شده را در یک ستون نگهداری می‌کنیم و یا در جدولی مجزا.

## ۴-۱ استراتژی رمزگذاری :

برای تأمین حریم خصوصی و سطح بالاتری از امنیت، ما از AES-CBC برای رمزگذاری داده‌های حساس استفاده می‌نماییم فقط مدیر پرس‌وجو (QM) کلیدهای مخفی (SK) را نگهداری می‌کند. ما برای دستیابی به سطح بالاتری از امنیت و پردازش سریع‌تر رمزنگاری، از الگوریتم متقارن به جای الگوریتم نامتقارن استفاده کردیم.

مدیر پرس‌وجو (QM) به عنوان یک سرور قابل اعتماد در فضای خصوصی یک سازمان یا شرکت قرار می‌گیرد و به عنوان یک واسطه بین کاربران و سرویس ابری عمل می‌کند و وظیفه پردازش درخواست‌ها (Query) و رمزگذاری (BV (bit vector) برای هر جدول را دارد.

عنصر اصلی همه سیستم‌های پیشنهادی PT (Partitioning Tree) است که در آن پرس‌وجو به‌طور مناسب برای اجرا توسط سرور ابری بازنویسی می‌شود. مقادیر موجود در هر ستون یک جدول به چند پارتیشن تقسیم می‌شوند که در آن‌ها هر پارتیشن شامل مجموعه‌ای از مقادیر است سپس ، QM بر اساس این مشخصات PT را می‌سازد . مثال در جدول زیر ، بر اساس حرف اول Name می‌تواند چندین Range تعریف نمود مثال A-I و J-P و برای هر کدام از این محدوده‌ها یک شناسه اختصاص دهیم و مجموع این محدوده‌ها می‌تواند یک PT را تشکیل دهد .

Table 3.1 The Students Table

ID	Name	SSN	VisaType	Department
01	Alice	12701	J-1	MATH
02	Ryan	25678	F-2	BUSN

بر اساس کاراکتر اول Name	
ID=1	A-I
ID=2	J-P
ID=3	Q-Z

Alice : ۱۲۱۱۱
Ryan : ۳۳۱۲

ما در راه‌حل‌های خود نگران مصرف حافظه نیستیم زیرا اطلاعات حساس در هر ردیف به بیت‌ها رمزگذاری می‌شوند (یعنی کوچک‌ترین واحد محاسبه) و حافظه زیادی هم برای جستجو نیاز ندارد . برای آشنایی بهتر می‌خواهیم عملیات اصلی و جبر رابطه‌ای یک پایگاه داده رابطه‌ای را توسط QM توضیح دهیم.

#### ۴-۱-۱ عملیات اصلی پایگاه داده :

- دستور **Insert** یک عملیات مستقیم هست. QM دستور درج را دریافت می‌کند ، BV جدیدی برای رکورد تازه‌وارد شده ایجاد می‌کند ، آن را به ماتریس بیتی مربوطه (BVM - Bit Vectors as a Matrix) اضافه می‌کند و سپس داده حساس را رمزگذاری می‌کند و به فضای ابری می‌فرستد.
- دستور **Select** یک عبارت اساسی در تمام برنامه‌های پایگاه داده است. در QM ، فرآیند اجرای دستور این‌گونه است که بررسی کنید کدام یک از موارد زیر قابل اجرا است و سپس آن را اجرا کنید .
- الف) هیچ یک از ستون‌های درخواستی در عملیات Select دارای اطلاعات حساس نیست ، پس رمزگذاری‌ای در ذخیره انجام نشده است و مستقیماً اطلاعات از پایگاه داده دریافت و نمایش داده می‌شود.
- ب) تمامی ستون‌ها دارای اطلاعات حساس هست پس این اطلاعات به‌صورت رمزگذاری شده ذخیره شده‌اند. در این حالت QM برای هر کدام از ستون‌هایی که در دستور وجود دارد ، موقعیت BV ها از BVM دریافت می‌نماید.
- سپس عملیات منطقی AND / OR را بر اساس شرایط موجود در فهرست‌های برگشتی برای هر ستون انجام می‌دهد . به‌عنوان مثال ، اگر شرط پرس‌وجو "F-۲" AND visa type = 'Mark' باشد ، QM هر BV را پیدا می‌کند که بیتی را نشان می‌دهد که مقادیر "QZ" را به گره نشان می‌دهد زیرا Mark زیرشاخه این گره در PT هست و تمامی BV ها به‌دست آمده در یک آرایه ثبت می‌شود. سپس عملیات منطقی AND را بین دو لیست انجام می‌دهد و برای دستور را مطابق رکوردهای یافت شده اصلاح می‌نماید.
- دستور **Update** نیز یکی از عملیات QM هست. برای این کار ابتدا QM توسط فرآیند توضیح داده شده Select ، رکورد را دریافت و رمزگشایی می‌نماید، سپس اطلاعات جدید را در رکورد یافت شده اصلاح می‌نماید و سپس آن را به سرور ابری ارسال می‌نماید.
- دستور **Delete** نیز مانند Update عمل می‌کند با این تفاوت که پس از یافتن رکورد موردنظر دیگر نیاز به رمزگشایی نیست

در سه دستور **Update** و **Delete** و **Insert** تمامی تغییرات بر روی **BVM** نیز اعمال می‌گردد تا در مراجعات بعدی معتبر باشد.

- در دستور **Alter** نیز در صورت حذف یک ستون، تمامی **BV** ها از **BVM** حذف می‌شود و سپس دستور حذف ستون ارسال می‌شود و چنانچه ما بخواهیم ستونی را اضافه نماییم و تعیین کنیم دارای اطلاعات حساس هست، تمامی رکوردهای مورد نیاز **BV** در **BVM** ساخته و آماده می‌شود و اگر ستون دارای اطلاعات حساس نباشد، عملیات بدون پردازش خاصی انجام می‌شود.

#### ۴-۱-۲ عملیات جبر رابطه‌ای :

- در دستور **Join** شاید بهترین و تنهاترین راه حل، دریافت کل اطلاعات رمزگذاری شده جدول و انجام این عملیات بر روی اطلاعات رمزگشایی شده باشد اما در حجم بالای اطلاعات، این راهکار مناسب نیست برای همین باید تا جایی که امکان دارد، ستون‌های مورد پیوند **Join** بدون رمزگذاری در پایگاه ابری ثبت شود مگر آنکه ستون‌های پیوند دارای اطلاعات حساس باشد.

زمانی که پیوند ما با اطلاعات حساس روبرو باشد، ما عملکردی مانند **Select** را خواهیم داشت یعنی **QM** ابتدا به ازای جداول پیوند فهرستی از **PT** را ساخته که شامل محدوده مقادیر مربوط به ستون‌های پیوند هست مثل تصویر روبرو و

Table A	Table B
Partitions of salary	Partitions of salary
10,000 to 20,000	10,000 to 25,000
20,001 to 30,000	25,001 to 40,000
30,001 to 40,000	

سپس ردیف‌هایی که دارای شرط پیوند هست را تعیین می‌کند و سپس دستور را بازنویسی کرده و ارسال می‌نماید. درواقع این گونه لیست سازی منجر می‌شود که ما بتوانیم اطلاعات کمتری جهت رمزگشایی دریافت و محاسبات کمتری انجام دهیم.

- در دستور **Union** هدف یکی کردن خروجی دو یا چند دستور هست که در آن باید مشخصات خروجی‌ها باید یکی باشد. این دستور مقادیر تکراری را نیز از خروجی نهایی حذف می‌نماید. البته باید جداول و خروجی‌هایی که در این عملیات شرکت می‌کنند حتماً با یک کلید یکسان رمزگذاری شده باشد. **QM** ابتدا ردیف‌های شرکت کننده در این عملیات را رمزگشایی می‌کند و آن‌ها را به یک **LinkedHashSet** (اگر نمی‌خواهید نظم درج را حفظ کنید اما می‌خواهید اشیاء منحصر به فرد را ذخیره کنید از **HashSet** استفاده می‌کنید اما اگر نظم درج مهم باشد از **LinkedHashSet** استفاده می‌شود) منتقل می‌کند.

- دستور **Intersection** فرآیندی است برای یافتن خروجی مشترک از دو یا چند خروجی و این کار بدون رمزگشایی اطلاعات، کار سختی هست. حال اگر حجم اطلاعات بالا باشد، بار سنگینی برای به دست آوردن اطلاعات ایجاد می‌شود و باید سعی کنیم این کار را بر روی سرور انجام دهیم چون امکان پردازش از سمت کاربر نیست. در این دستور هم مانند **Union** باید خروجی‌ها دارای مشخصات یکسان باشد و چون دنبال اعضای مشترک هستیم، می‌توانیم از فرآیند دستور **Join** استفاده کنیم و خروجی حاصل شده را برای کاربر ارسال نماییم.



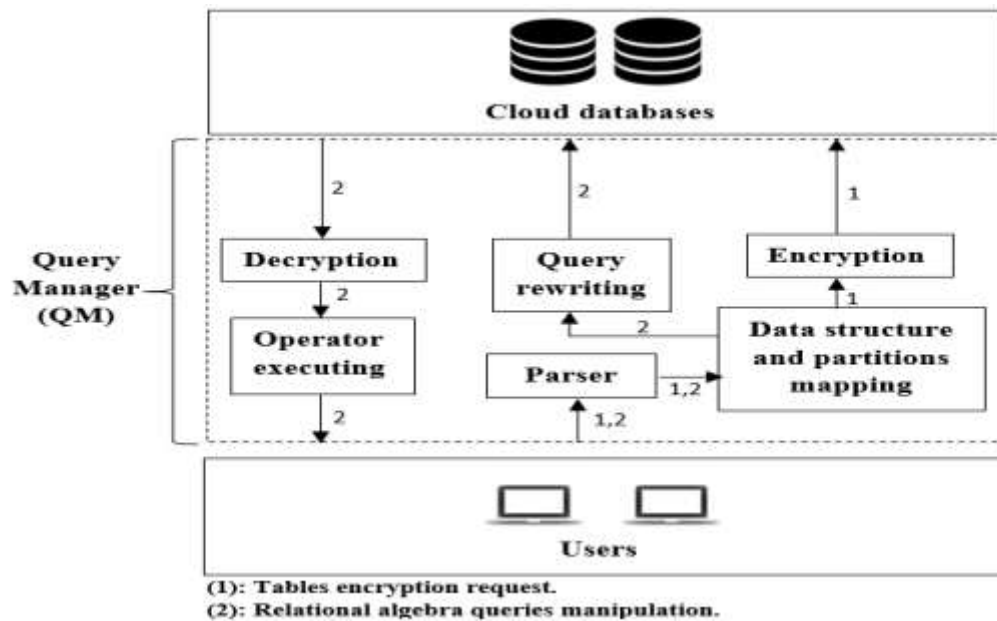
- عملگر **Difference** فرآیندی برعکس Intersection برای به دست آوردن خروجی خای غیرمشترک هست. برای انجام این عملیات از طریق Intersection ابتدا ردیف‌های مشترک را به دست آورده و آن را از خروجی اصلی خارج می‌کنیم تا خروجی به دست بیاید.
- عملگر **Duplication Removal** یا همان distinct فرآیندی برای حذف مقادیر تکراری از خروجی‌ها هست. در فضای ابری این عملیات غیرممکن است زیرا ما از الگوریتم رمزگذاری غیرقطعی (AES-CBC) استفاده می‌کنیم. بنابراین، ما قبل از ارسال خروجی به کاربر، این عمل را در QM انجام می‌دهیم. در این مورد نیز مانند Union می‌توانیم از یک LinkedHashSet با قابلیت حذف مقادیر تکراری استفاده کنیم.
- عملگرهای **Aggregation and Sort** شامل max, min, and count هستند. ما در اینجا ستون‌ها با مقادیر دامنه‌ای و غیر دامنه‌ای داریم که در ستون‌های غیر دامنه‌ای می‌توانیم پردازش را در QM و در کاربری انجام دهیم. ما از محاسبات بر روی اطلاعات رمزگشایی شده به خاطر حجم محاسبات، اجتناب می‌کنیم. چنانچه قبلاً در دستور Select ما ابتدا اطلاعات مربوطه را همراه شرط‌هایی که لحاظ شده بود را تو BV و BVM یافت می‌کردیم، این بار هم مرحله اول به این صورت هست. در مورد دستور count عملیات راحت هست چون ما ردیف‌ها را پیدا کردیم و کافی ست تعداد آن را به دست بیاوریم اما در خصوص Sum مجبور به رمزگشایی هستیم و در دستور Avg باید هم Sum و هم Count را انجام دهیم.
- عملگر **Project** یک عملیات بر اساس ستون‌ها هست که QM تمام ردیف‌ها مطابق ستون‌های درخواستی انتخاب می‌کند و چون انتخاب ستون هست ما نیازی به PT نداریم اما به‌هرحال نیاز به رمزگشایی هست و البته نیازی به حفظ اطلاعات تکراری نداریم.

## ۴-۲ روش‌های قابل اجرای رمزگذاری :

ما برای ثبت اطلاعات رمزگذاری شده، دو روش را موردبررسی قرار می‌دهیم. یک روشی که در آن اطلاعات رمزگذاری شده به‌صورت یک ستون در همان جدول نگهداری می‌شود و روش دوم، روشی که اطلاعات رمزگذاری شده را در جدول مجزا نگهداری کنیم.

### ۴-۲-۱ استفاده Bit Vectors به‌صورت یک ستون (BVSAC):

ما مدلی را بررسی می‌کنیم که Bit Vectors را به‌عنوان یک ستون اضافی در جدول رمزگذاری شده اصلی ذخیره می‌کند و این مدل را برای اجرای دستورات جبر رابطه‌ای مختلف بر روی داده‌های رمزگذاری شده طراحی کردیم. علاوه بر این، ما محاسبه را به دو طرف تقسیم کردیم: سمت سرویس‌گیرنده و سرویس‌دهنده ابر (CP – Cloud Process) که در آن ما بیشترین محاسبه را با بازنویسی دستورات که CP را قادر به اجرای آن‌ها می‌کند، به سایت CP منتقل می‌کنیم. این مدل در تصویر زیر نمایش داده شده است و این مدل برای تمام عملگرهایی که اینجا تحلیل شدند اجرا می‌گردد. در این مدل ما از رمزگذاری BV که تا اینجا توضیح داده شد استفاده کرده‌ایم و در آن به‌جای آنکه محاسبات را مداوم توسط QM در سمت کاربر انجام دهیم، آن را به CP منتقل می‌کنیم.



در این مدل جدول رمزگذاری شده دارای یک ستون است که حاوی اطلاعات BV هست و چون اطلاعات به صورت بیت هست نیاز به نگرانی بابت حجم آن نیست اما چون نهایت سایز ذخیره برای هر ویژگی ۶۴Bit هست نیاز به بیش از یک ویژگی داریم. حتی اگر یک مهاجم بتواند از طریق استنتاج BV به اطلاعات دست پیدا کند، چون ما نام جداول و ستون‌ها را توسط الگوریتم رمزگذاری قطعی (AES) (که در آن متن‌های رمز برای هر متن ساده همیشه یکسان است) رمزگذاری می‌کنیم، این احتمال ضعیف‌تر و مشکل مرتفع می‌شود.

به عنوان مثال، ستون "Student-Rank" شامل محدوده‌ای از مقادیر مانند junior، grad، senior و غیره است. یک مهاجم ممکن است بتواند مقادیر احتمالی را استنباط کند، حتی وقتی آن‌ها به طور تصادفی با AES-CBC رمزگذاری شده‌اند.

اشتراک بیش از یک متن رمزگذاری با کلید یکسان باعث نگرانی نمی‌شود زیرا هرکدام نمی‌تواند متن اصلی را به دست آورد مگر اینکه کلید رمزگذاری را به دست آورد و برای حل این نگرانی، کلید به سرور ابری منتقل نمی‌شود.

برای مثال می‌خواهیم یک جدول را با کلید رمزگذاری کنیم. ابتدا QM نام جدول و ستون‌های حساس را توسط AES-DET رمزگذاری می‌کنیم به این صورت که جدول دیگری در فضای ابری ساخته می‌شود و با آنچه تا اینجا درباره PT و الگوریتم‌های رمزگذاری گفته شد، QM ردیف‌ها را به صورت BV رمزگذاری می‌کنید و جدول رمزگذاری شده دارای یک ستون Index است که در آن هر Index بخشی از BV هر ردیف است.

این فرآیند برای جدول List در شکل زیر توضیح داده شده است به این صورت که مشاهده می‌کنید ستون‌های Name - Rank - Visa Type و Department در جدول رمزگذاری شده توسط رمزگذاری AES-DET و یک کلید مشترک به صورت نام غیر آشنا درآمده است.

ID	Name	Rank	Visa type	Department
110	Alice	freshman	F1	Computer science
111	Sara	senior	J1	Computer engineering
112	John	junior	None	Information system
113	Ryan	Sophomore	J2	Math

Index	ID	ct <sub>n</sub>	ct <sub>r</sub>	ct <sub>vt</sub>	ct <sub>d</sub>	Reference
01	110	* & ^	* ^ %	*/d	^ % ^ H	10001000010000100000
02	111	% ^ &	/+ \$	& ^ /	& & % \$	00010100000100010000
03	112	)(#	% \$ / *	+ - * &	) * # R	01000010000001001000
04	113	\$ # !	! @ ~ K	*/f	@ \$ % *	00100001000010000010

زمانی که دستوری برای خواندن اطلاعات ارسال می‌شود QM ابتدا ستون‌های استفاده‌شده در دستور به ستون‌های رمزگذاری شده ترجمه می‌کند مانند `Select Name From List Where [Visa Type]=?` تبدیل به `Select Ctn From Cta where Ctv=?`

حال با پردازش بخش `Where` شروع به دست آوردن ردیف‌ها با توجه `PT` می‌نماید و ردیف‌های مربوطه یافت می‌شود و دستور بازنویسی شده را با توجه به `ID` ها می‌سازد. البته چنانچه قبلاً گفته‌شده `PT` شامل محدوده مقادیر هست و `Query` بازنویسی شده نیز محدوده‌ای از اطلاعات را می‌آورد و نه ردیف‌های قطعی! مثلاً اگر شرط `[Visa Type]='F۱'` باشد با توجه به حرف اول `F` و خاصیت `Bit` بودن ستون محدود ای به این صورت در نظر گرفته می‌شود `Reference&۵۲۴۲۸۸>۰` سپس با توجه به حرف دوم که `۱` هست یک شرط ترکیبی ساخته می‌شود و درنهایت ممکن است دستوری مانند زیر ساخته شود `Reference&۶۵۵۳۶>۰ And Reference&۵۲۴۲۸۸>۰`

درنهایت ردیف‌های به‌دست‌آمده توسط `QM` رمزگشایی می‌شود و مجدد شرط بر روی آن اعمال می‌گردد و اطلاعات به خروجی ارسال می‌شود.

عملیات‌هایی مانند `Insert` و `Update` و `Delete` نیز چنانچه قبلاً گفته شد اعمال می‌گردد با این تفاوت که این بار از `BVM` استفاده نمی‌کنیم و این اطلاعات را در ستون `Reference` داریم.

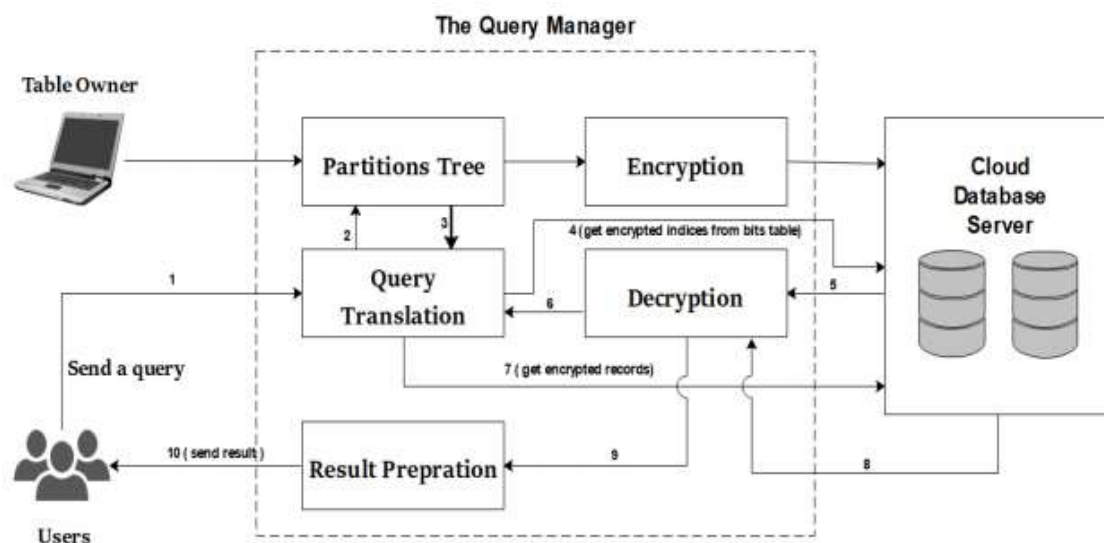
#### ۴-۲-۲ استفاده Bit Vectors به صورت یک جدول مستقل (BVSIT):

در این روش می‌خواهیم `BV` را در یک جدول مستقل در سرور ابری ذخیره می‌کنیم و بین جدول رمزگذاری شده و جدول بیت‌ها یک ستون `Index` اضافه می‌کنیم که در آن مقدار `Index` هر رکورد در جدول رمزگذاری شده برابر با `Index` رکورد مربوطه در جدول بیت‌ها است.

برای اطمینان از اینکه در جدول بیت‌ها هیچ بیت داده‌ای توسط هکر از طریق استنتاج نشت نمی‌کند، ما `Index` موجود در جدول بیت‌ها را رمزگذاری می‌کنیم و سعی می‌کنیم حتی از لحاظ قرارگیری موقعیت رکوردها نیز یکسان نباشد.

برای امنیت در `BVSIT`، داده‌ها در تمام ستون‌های حساس با الگوریتم متقارن (`AES-CBC`) که رمزگذاری تصادفی است رمزگذاری می‌شوند.

فرآیند کلی این روش در شکل زیر نمایش داده‌شده است.



در این عملیات برای پردازش دستورات ، QM دو پرس و جو صادر می کند. پرس و جو اول ، شاخص های رمزگذاری شده تمام ردیف های دارای شرایط درخواستی را برای پرس و جو انتخاب شده از جدول بیت ها بازیابی می کند. پرس و جو دوم ، ردیف های رمزگذاری شده را مطابق جدول بیت ها انتخاب می کند و اطلاعات رمزگذاری شده را دریافت و رمزگشایی می کند و نتیجه به خروجی ارسال می شود. نحوه ساخته شدن جدول بیت ها و جدول رمزگذاری شده در تصویر زیر نمایش داده شده است. اطلاعات قرار گرفته در جدول بیت ها همان PT بر اساس محدوده های کاراکترها هست که قبلاً توضیح داده شده است

> جدول بیت ها

Index	1	2	3	4	5	6	7	8	9	10	11	12	13
CT(3)	0	1	0	1	0	0	1	0	0	1	0	0	0
CT(4)	0	1	0	1	0	1	0	0	0	0	0	1	0
CT(1)	1	0	0	1	0	0	1	0	0	0	1	0	0
CT(2)	0	0	1	0	1	0	0	1	0	0	0	1	0

ID	Name	SSN	Rank	Salary
01	Alice	12701	Secretary	30,000
02	Ryan	25678	Admin	60,000
03	Mark	46932	Secretary	29,000
04	John	42213	Manager	55,000

> جدول رمزگذاری شده

Index	ID	Name	SSN	Rank	Salary
1	^%#	)@{(	\$@^#	@^#S	FS^@
2	(#*	^#*	)@H2	+ @	%@*&
3	(#(	(#&	\$@%	@^#(	@%TW
4	)^@	)^#H	!@M	*@^#	+_)#FQ

عملیاتی هایی مانند Insert و Update و Delete از همان روش های گفته شده استفاده می کند و با این تفاوت که برای هر کدام از عملیات ها دو دستور ارسال می شود ، یکی برای جدول رمزگذاری شده و دیگری جدول بیت ها.

# فصل پنجم

## مقدمه

هر ایده یا استراتژی در بخش نرم‌افزار باید در مرحله تست و آزمایش قرار بگیرد تا بتوان نسبت به اجرا و منابع موردنیاز، بررسی و نیازسنجی شود زیرا استفاده از ایده نامناسب و پرهزینه منجر به شکست سیستم می‌شود. با آنچه از رمزگذاری و پایگاه داده گفته‌ایم می‌توان مدل‌های معتبر و رایج را بررسی و مقایسه نمود. در آزمایش بین مدل‌ها و تئوری‌های اجرای پایگاه داده رمزگذاری شده، پارامترهای متعددی مدنظر قرار می‌گیرد مانند: حجم اطلاعات، زمان اجرا، هزینه‌های سخت‌افزاری مانند حافظه و امنیت مدل‌ها. مدل‌های پایگاه داده رمزگذاری شده با توجه به ایده و استراتژی که در آن قرار دارد در عملیات‌های مختلف دارای پاسخ و ارزیابی متفاوتی هست مانند آنکه مدلی در هنگام ذخیره اطلاعات به علت طراحی آن باعث هزینه بالا می‌شود و یا مدلی در ذخیره‌سازی اطلاعات دارای نتایج بهتر ولی در بازخوانی اطلاعات دچار کندی یا هزینه بالای اجرا می‌شود. البته عملیات‌های پایگاه داده تنها ذخیره و خواندن اطلاعات نیست و عملیات‌های محاسباتی نیز موردبررسی قرار می‌گیرد پس آن چیزی که به‌عنوان مدل مناسب انتخاب می‌گردد، مدلی هست که دارای میانگین بهتری نسبت به سایر مدل‌ها باشد و درنهایت باید مدل انتخابی را در اجرای فضای ابری موردبررسی قرارداد.

## ۵-۱ آزمایش و ارزیابی :

تا اینجا، ما به‌طور مفصل مدل‌های پیشنهادی را توضیح دادیم و برای هر مدل، الگوریتم‌های توسعه‌یافته را برای اجرای عملگرهای مختلف جبر رابطه‌ای ارائه دادیم. همان‌طور که قبلاً گفته شد، این پایان‌نامه باهدف پیاده‌سازی، ارزیابی و مقایسه عملکرد نمونه‌های پیشنهادی برای انواع مختلف عبارات با استفاده از معیارهای مختلف (به‌عنوان مثال، تأخیر اجرای نمایش داده‌ها، نیازهای فضا در QM و سرور ابری، و درصد سربار محاسبات در QM) ارائه دادیم. حال برای مقایسه مدل‌ها و سیستم‌های ارائه‌شده نیاز به آزمایش آن‌ها داریم. برای انجام همه آزمایش‌ها برای همه سیستم‌ها BVM، BVSAC، BVSIT، OBT، CBF، CryptDB از یک سیستم با ۶ گیگابایت رم، ۱ ترابایت HDD و پردازنده Core i۵ با ۲٫۸ گیگاهرتز استفاده کردیم. برای اجرای توابع

QM در مدل‌های پیشنهادی، از Java برای استفاده کردیم. سرور MySQL بر روی دستگاه کاربر استفاده شد و ما از Java Database Connectivity (JDBC) به عنوان اتصال دهنده به MySQL استفاده کردیم و تمام آزمایش‌ها بر روی دستگاه محلی انجام شده است تا تأخیر ارتباط از گزارش‌ها حذف شود.

در تمامی مدل‌های خود، از نسخه تصادفی AES-CBC برای رمزگذاری داده‌های حساس استفاده کردیم و برای رمزگذاری و رمزگشایی از کتابخانه "javax.crypto" استفاده نمودیم.

ما OBT و CBF را اجرا کردیم زیرا پیاده‌سازی آن‌ها برخلاف CryptDB که برای استفاده عمومی در GitHub در دسترس است به صورت آنلاین در دسترس نیست و برای ذخیره BV آن‌ها را به صورت محلی در QM در مدل اول ذخیره کردیم و آن‌ها را برای استفاده در آینده در یک فایل متنی نوشتیم و در سایر مدل‌ها آن‌ها را در داده ابری بارگذاری می‌کنیم. برای شروع آزمایش ما به طور تصادفی چهار جدول ایجاد کردیم و برای هر ستون فهرستی از مقادیر تعریف کردیم و به یک برنامه java اجازه می‌دهیم تا با انتخاب مقادیر از فهرست‌ها، ردیف‌هایی را بسازد. اندازه جداول (تعداد رکوردها) ۱۰ هزار – ۲۰ هزار – ۵۰ هزار و ۱۰ هزار ردیف بود.

ما در کل ۲۴ ستون برای همه جداول داشتیم و همه آن‌ها را به استثنای ستون ID به عنوان اطلاعات حساس در نظر گرفتیم. در مطالعه خود، اگرچه می‌توانستیم از مدل‌های پیشنهادی با جداول کوچک استفاده کنیم، اما ما بر روی جدول‌های بزرگ تمرکز کردیم.

## ۵-۲ ارزیابی :

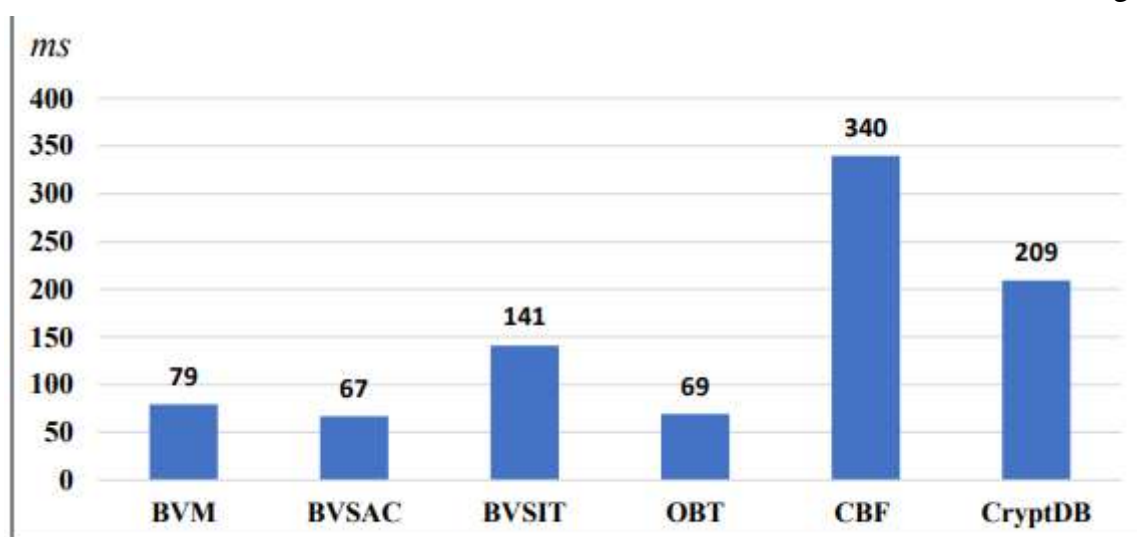
ارزیابی شامل موارد زیر هست :

- ۱) آزمایش و ارزیابی هزینه اجرای عملیات پایه پایگاه داده (Insert – Update – Deleted - Insert)
  - ۲) آزمایش و ارزیابی هزینه اجرای عملیات محاسباتی (Sum – Min – Max - Avg)
  - ۳) آزمایش و ارزیابی هزینه‌های اجرای عبارات جبر رابطه‌ای (Join – Union - Intersection)
  - ۴) شناسایی نیازهای فضایی برای هر سیستم
- در مدل‌های پیشنهادی، مرحله رمزگذاری پایگاه داده اصلی شامل تجزیه رکوردها، تولید شاخص‌ها، ساخت BV، رمزگذاری داده‌های حساس و قراردادن داده‌های رمزگذاری شده در جدول رمزگذاری شده در سرور ابری است. در جدول زیر ما زمان صرف شده توسط هر سیستم برای رمزگذاری هر جدول را مقایسه کرده‌ایم و همان‌طور که دیده می‌شود، مدل BVSAC کمترین تأخیر در رمزگذاری و درج را در بین همه سیستم‌ها تجربه کرده است زیرا QM نیازی به ذخیره و مدیریت BV‌ها به صورت محلی در QM مانند BVM ندارد.

**The delay of the original database encryption comparison among all systems in minutes**

N.R	BVM	BVSAC	BVSIT	CryptDB	CBF	OBT
10k	12	11	25	44	55	13
20k	26	23	48	66	112	24
50k	65	54	115	158	291	55
100k	147	112	223	308	578	113

پس از BVSAC مدل OBT است که هر ردیف را به صورت یک بلوک رمزگذاری می کند و آن را در ستونی قرار می دهد و در روی دیگر، مدل CBF بالاترین تأخیر فرایند رمزگذاری را دارد زیرا به ازای هر ستون باید جداگانه درج انجام شود و رکوردها را به جداول مختلف در سرورهای مختلف ابری قرار می دهد. مقایسه مدل ها برای عملیات Insert در جدول زیر نمایش داده شده است



**The delay comparison of insert statements for all systems, in millisecond**

مدل کند بعدی مدل CryptDB هست که به دلیل محاسبات سنگین رمزگذاری Onion هست.

## ۵-۲-۱ آزمایش شماره یک :

از آنجاکه پراستفاده ترین عملیات ، عمل انتخاب (Select) هست ، ما بیشترین تمرکز خود را بر روی این عملیات می گذاریم. در این آزمایش ، ما درصد متوسط ردیف های رمزگذاری شده واكشی شده از جداول رمزگذاری شده برای مدل های پیشنهادی را محاسبه کردیم .ما همچنین مطالعه کردیم که چگونه تعداد بندهای موجود در شرایط پرس و جو می تواند به کاهش دامنه موارد واكشی کمک کند.



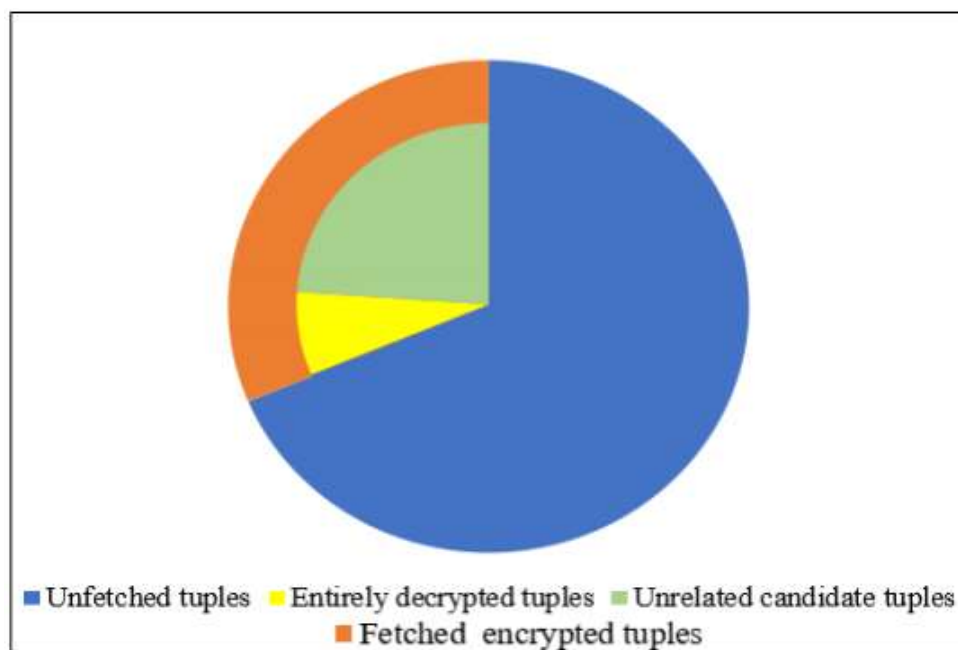
**The percentage of retrieved encrypted tuples for all proposed models**

چنانچه در نمودار می بینید ، بدون استفاده از هیچ یک از نمونه های اولیه پیشنهادی برای مدیریت پایگاه داده رمزگذاری شده (یعنی بدون Index) باید کل جدول رمزگذاری شده برون سپاری شده را بازیابی کنیم.

#### ۵-۲-۲ آزمایش شماره دو :

چنانچه تا اینجا گفته شد هدف ما این است که برای جلوگیری از محاسبات سربار ، که باعث تسریع در زمان اجرای درخواست می شود ، محاسبات رمزنگاری را در QM کاهش دهیم.

در تصویر زیر می بینیم که از کل مجموعه رمزگذاری شده واکنشی شده ، میانگین ردیف های کاملاً رمزگشایی شده کمتر از ۲۴٪ ردیف های واکنشی شده است . بر این اساس ، مدل های ما کارآمد هستند زیرا ما نه تنها دامنه داده های برون سپاری شده بازیابی شده را محدود می کنیم ، بلکه در صورت امکان محاسبات غیرضروری (محاسبات رمزنگاری) را نیز حذف می کنیم.



**Average percent of entirely decrypted tuples**



### ۵-۲-۳ آزمایش شماره سه :

در این مطالعه ، ما بر روی عبارات *Select* متمرکز شدیم ، زیرا به طور گسترده در سیستم های پایگاه داده استفاده می شود. بنابراین ، ما عبارات *Select* را با تعداد بندهای مختلف و منفرد ، مانند انتخاب با نام ستون و انتخاب همه ستون ها ، با مدل های معرفی شده تا اینجا CryptDB ، OBT و CBF آزمایش و ارزیابی کردیم.

جدول زیر کل زمان اجرا را برای همه سیستم ها هنگام اجرای عبارات *Select* برای بازیابی مقادیر با خاصیت ستون حساس ارائه می دهد و مدت زمانی که اندازه گیری کردیم زمان تجزیه پرس و جو بود تا نتیجه پرس و جو نهایی در میلی ثانیه (ms) است.

	1 clause						2 clauses						3 clauses					
	BVM	BVSAC	BVSIT	OBT	CBF	CryptDB	BVM	BVSAC	BVSIT	OBT	CBF	CryptDB	BVM	BVSAC	BVSIT	OBT	CBF	CryptDB
10k	242	246	229	200	231	870	136	153	133	98	238	1329	115	132	112	77	324	1558
20k	326	382	352	275	338	1633	259	263	215	179	483	2016	213	224	192	142	542	12626
50k	624	510	783	481	547	10764	479	367	614	289	746	18969	440	303	545	221	893	21634
100k	1107	918	1240	791	970	20135	840	702	1128	562	1478	37078	835	610	1021	493	1604	59910

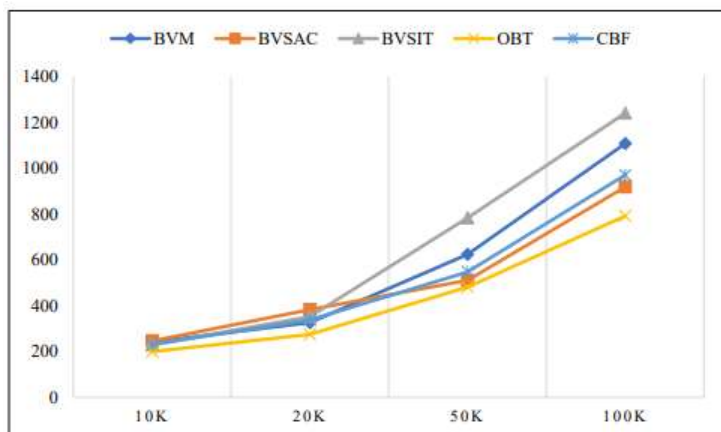
چنانچه مشاهده می شود مدل BVSAC به طور متوسط کوتاه ترین زمان اجرا را از بین تمام مدل های پیشنهادی برای همه موارد دارد. این نتیجه به این معنی است که BVSAC از نظر زمان اجرا ، با پایگاه داده های بزرگ تر عملکرد بهتری دارد و از زمان شروع بیشتری نسبت به سایر مدل های پیشنهادی برخوردار است ، به همین دلیل سرعت آن نسبت به برخی دیگر از سیستم ها برای پردازش پایگاه داده با ۲۰ هزار و ۱۰ هزار ردیف کمتر است. عامل دیگری که سرعت پردازش پرس و جو را افزایش می دهد این است که ما از عملیات bitwise ارائه شده توسط MySQL در فضای ابری بهره مند می شویم زیرا بیت ها را به عنوان یک ستون مستقل همراه با هر رکورد رمزگذاری شده ذخیره می کنیم.

برای پایگاه داده های کوچک تر (۱۰ هزار یا کمتر )، مدل BVSIT سریع ترین زمان اجرا را دارد زیرا عملیات bitwise به صورت ستون پایه انجام می شود ، که امکان جستجوی سریع در جدول بیت ها را فراهم می کند . با این حال ، وقتی پایگاه های داده رشد می کنند ، با رشد مقدار شاخص های رمزگشایی در QM ، عملکرد BVSIT کاهش می یابد.

مدل BVM در جایگاه دوم برای پایگاه داده با بیش از ۵۰ هزار ردیف است و اصلی ترین عاملی که عملکرد BVM را تحت تأثیر قرار می دهد زمان بارگیری BV ها از دیسک سخت به حافظه اصلی و سپس جستجوی آن هاست.

از آنجاکه مدل CryptDB دارای تایم خیلی بالاتری نسبت به سایر مدل ها هست، ما آن را حذف و میانگین سایر مدل ها در هر سه حالت را در آورده ایم.

نمودار روبرو حاصل مقایسه مدل‌ها هست.



(b) Average total processing time for all models, excluding CryptDB.

با توجه به نمودار و جدول بالا مشاهده می‌شود که مدل OBT دارای سرعت بهتری نسبت به سایر مدل‌ها هست زیرا مقدار داده‌های رمزگشایی کمتر از پیشنهاد ما بود (همه مقادیر به‌عنوان یک بلوک ذخیره می‌شوند که منجر به رمزگشایی کمتر برای هر سطر می‌شود) البته این زمانی است که دستور Select با تعداد ستون بیشتر انجام شود زیرا در دریافت تک‌ستونی، سایر مدل‌ها از سرعت بهتری برخوردار هستند.

## ۵-۲-۴ توان عملیاتی :

توان عملیاتی به‌عنوان مقدار داده منتقل‌شده در یک‌زمان معین تعریف می‌شود. با اندازه‌گیری توان عملیاتی، می‌توان فهمید که کدام سیستم با افزایش داده‌های درخواستی، پاسخ بیشتری به دستورات کاربر می‌دهد، زیرا کاربر نهایی کسی است که تحت تأثیر کاهش سرعت سیستم قرار خواهد گرفت.

برای اندازه‌گیری توان عملیاتی، مجموعه‌ای از پرس‌وجوها را برای بازیابی ۲۵٪ و سپس ۵۰٪ رکوردها از جدول دارای ۱۰۰۰۰۰ رکورد اجرا کردیم. سپس، ما مقدار داده‌های ساده (داده‌های سوابق پس از رمزگشایی) را اندازه‌گیری می‌کنیم و

$$\text{Throughput} = \frac{\sum_{k=0}^n \text{unencrypted record's size (in bytes)}}{\text{Total time taken to deliver the data (MS)}} \quad \text{آن‌ها را بر زمان موردنیاز هر سیستم تقسیم می‌کنیم تا}$$

داده‌های موردنیاز را به کاربر تحویل دهیم

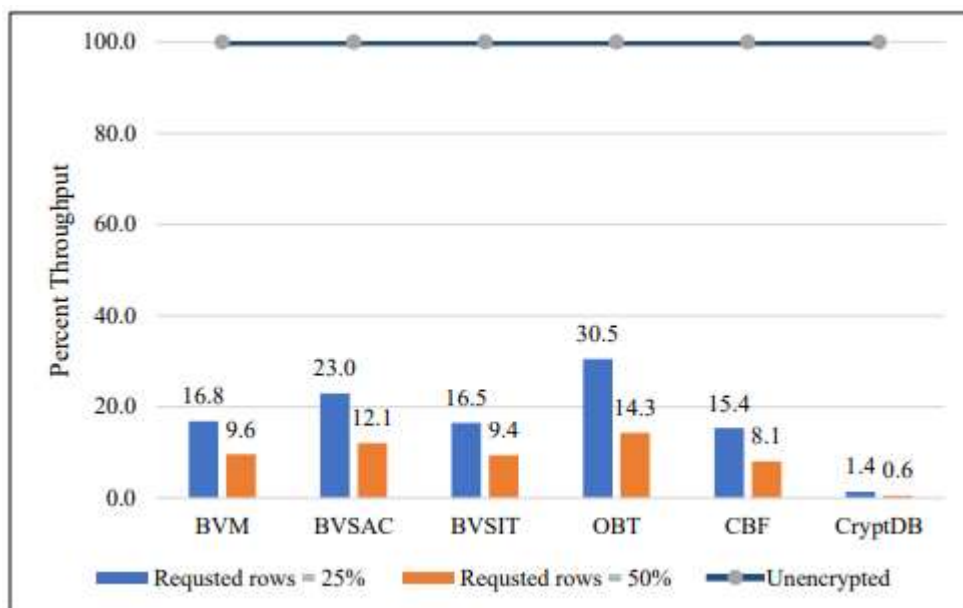
جدول زیر این محاسبات را برای تمامی مدل‌ها نمایش می‌دهد

	kilobytes per second (kB/s)					
	BVM	BVSAC	BVSIT	OBT	CBF	CryptDB
Requested rows = 25%	794	1083	778	1440	727	66
Requested rows = 50%	971	1221	947	1447	816	57

**The amount of plain data in kB that each system can deliver to the user per second.**

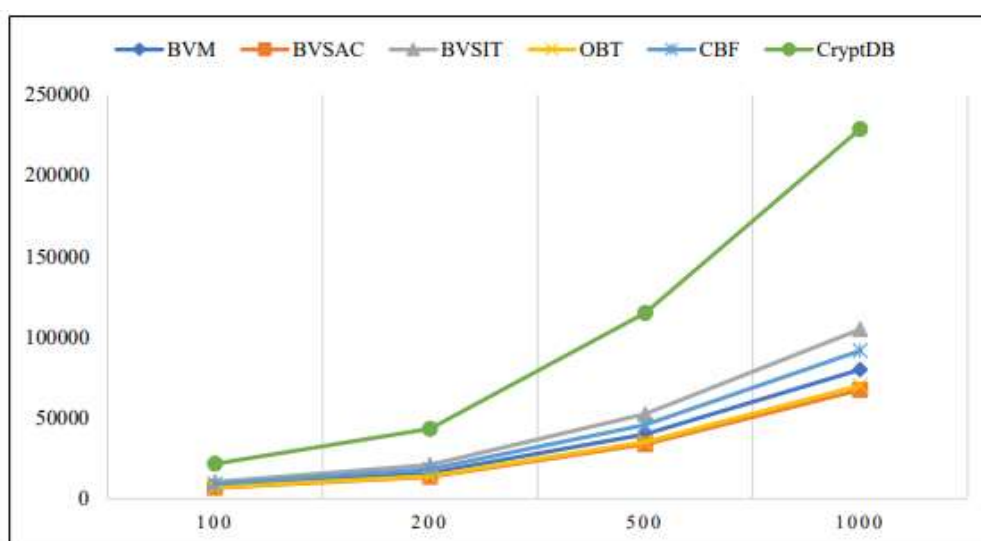
چنانچه مشاهده می‌شود مدل BVSAC توان عملیاتی بالاتری را در بین سیستم‌های پیشنهادی دارا هست که آن را بهترین گزینه برای کاربران نهایی می‌کند که به دنبال یک سیستم پاسخگویی سریع‌تر هستند. همچنین مدل OBT

بالاترین توان عملیاتی را دارد، به این دلیل که برای رمزگذاری داده‌ها به مقدار کمتری از بایت در میان همه سیستم‌ها نیاز دارد. برای درک بهتر مقادیر را به صورت نمودار درصد درآوردیم.



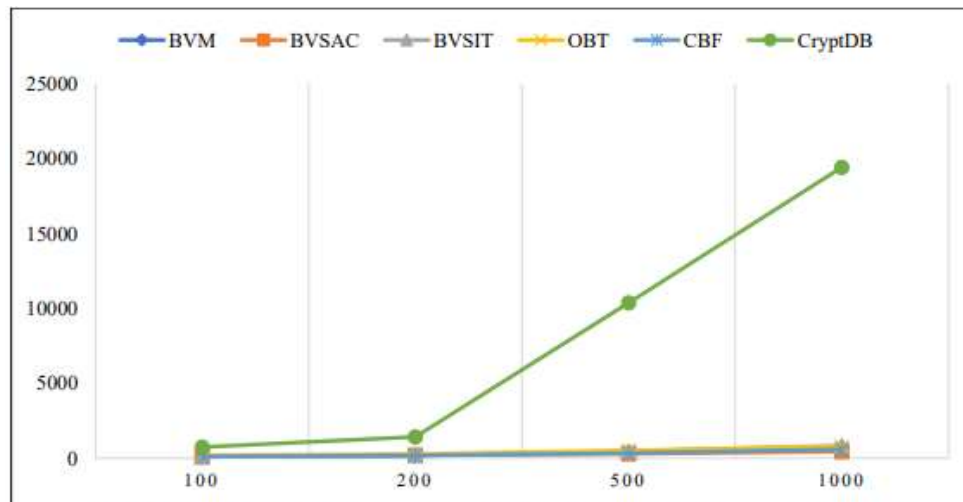
#### ۵-۲-۵ آزمایش شماره چهار (آزمایش دستورات Update و Delete):

ما این آزمایش را با مدل‌های مختلف انجام و نمودار زیر متوسط زمان هزینه شده توسط هر سیستم برای اجرای عبارات *Update* را نشان می‌دهد. در این آزمایش، عبارات *Update* را فقط با یک فیلد اجرا کردیم. همان‌طور که در شکل دیده می‌شود، هزینه زمان به‌روزرسانی در همه سیستم‌ها زیاد است و نتیجه به‌روزرسانی یک فیلد رمزگذاری شده در سیستم‌های پایگاه داده است. محور X نشان‌دهنده تعداد ردیفی است که از دستورات *Update* تأثیر می‌پذیرد، به این معنی که داده‌های موردنیاز را انتخاب کرده و سپس در حالت درج ۱۰۰ ردیف، در حالت دوم ۲۰۰ ردیف و غیره صادر می‌کند.



Comparison of the average delay of update statements for all models to update a number of existing tuples (100, 200, 500, and 1,000 tuples).

در مدل‌های پیشنهادی، BVSIT کندترین سیستم پس از CryptDB است زیرا *Update* باید هم در جدول رمزگذاری شده اصلی و هم در جدول بیت‌ها انجام شود و دومین مدل کند است. علاوه بر این، BVM کمی تأخیر بالاتر از BVSAC را تجربه کرد اما هنوز هم سریع‌تر از CBF عمل می‌کند، که سومین مدل کندترین در این مقایسه است. در CryptDB، هزینه به‌روزرسانی بالاترین هست. به‌طور خلاصه، می‌توان گفت BVSAC کارآمدترین مدل برای اجرای *Update* است. آزمایش را با انجام عملیات Delete انجام دادیم و نمودار زیر زمانی است که هر مدل برای حذف تعداد متفاوت (۱۰۰، ۲۰۰، ۵۰۰ و ۱۰۰۰ ردیف) صرف می‌کند، درحالی‌که شرط حذف فقط یک‌بند دارد.



Comparison of the average delay of delete statements for all models to delete a different number of tuples (100, 200, 500, and 1,000 tuples).

فرآیند حذف، ردیف‌های موردنیاز را انتخاب کرده و سپس آن‌ها را حذف می‌کند. حذف در همه مدل‌های پیشنهادی کارآمد است زیرا حذف پس از اجرای عملیات *Select* برای بازیابی ردیف‌های موردنیاز انجام می‌شود. به‌جای ارسال یک پرس‌وجو برای حذف هر رکورد، ما Index رکورد را به دست می‌آوریم و سپس یک دستور صادر می‌کنیم تا هر رکورد از شاخص آن را در دستور حذف، یعنی حذف از TABLE\_NAME که در آن index است حذف کند. CryptDB کندترین سیستم برای اجرای دستورات حذف به همان دلایلی است که قبلاً ذکر کردیم.

#### ۵-۲-۶ آزمایش شماره پنج (دستورات محاسباتی):

ما این آزمایش را با دو بخش تقسیم و در بخش اول متوسط زمان هزینه شده توسط هر سیستم برای اجرای عملیات *Count*, *Average*, *Sum* در مدل‌های مختلف را به دست آورده‌ایم. در بخش دوم سایر دستورات مانند Max و Min مورد آزمایش قرار گرفته‌اند.

تمامی آزمایش‌ها برای محدوده مشخصی از اطلاعات هست و چیزی حدود ۶٪ از کل اطلاعات.

	SUM						AVERAGE						COUNT					
	DVM	BVSAC	BVSIT	OBT	CBF	Cryptdb	DVM	BVSAC	BVSIT	OBT	CBF	Cryptdb	DVM	BVSAC	BVSIT	OBT	CBF	Cryptdb
10k	110	121	107	206	143	878	111	121	107	206	147	878	104	115	101	200	131	746
20k	196	212	232	284	212	1718	198	214	234	286	208	1718	189	205	225	277	193	1417
50k	440	304	473	525	361	10762	441	305	474	527	363	10762	431	295	464	516	328	10332
100k	733	460	850	807	645	20185	734	462	853	806	647	20185	722	449	839	796	599	19325

**Delay comparison in milliseconds of executing SUM, AVERAGE, AND COUNT functions. (only one predicate in the queries)**

در سیستم OBT، طبق ساختار، اجباراً کل ردیف‌های رمزگذاری شده واکشی شد، که نتیجه آن محاسبه هزینه‌های سنگین‌تر در QM است. در CBF و CryptDB، الگوریتم رمزنگاری مورد استفاده برای رمزگذاری قسمت عددی برای پشتیبانی از عملیات جمع‌آوری روی مقادیر رمزگذاری شده، رمزگذاری homomorphic است که مدل‌های آن بر ضرب مدو لا برای تولید مجموع رمزگذاری شده استفاده می‌کند، که نیاز به محاسبه سنگین‌تر دارد. جدول زیر آزمایش را در هر سیستم برای اجرای عملیات *Min*، *Max* نشان می‌دهد.

	Average delay of MAX and MIN					
	BVM	BVSAC	BVSIT	OBT	CBF	CryptDB
10k	112	123	109	208	147	878
20k	201	217	237	289	208	1714
50k	447	311	480	532	363	10762
100k	741	468	858	815	647	20185

**Delay comparison in milliseconds of executing MAX and MIN functions. (only one predicate in the queries)**

با بررسی این دو جدول در آزمایش‌ها عبارات محاسباتی، می‌توانیم مدل که دارای تأخیر کمتر هست را به دست آوریم که آن را می‌توان برای پایگاه داده کوچک در نظر گرفت، اما برای پایگاه‌های داده بزرگ‌تر، مدل BVSAC برنده است، درحالی‌که BVSIT محاسبات سنگین‌تری را تجربه می‌کند، همان‌طور که گفته شد در آزمایش قبلی CryptDB به دلیل عملیات رمزگشایی فشرده، کندترین مدل هست.

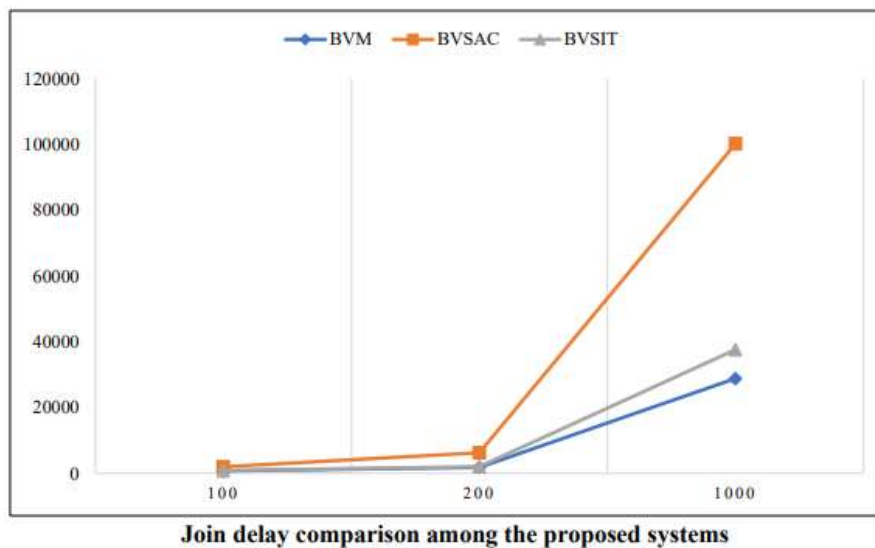
## ۵-۲-۷ آزمایش شماره شش (دستورات Join, Union):

در این آزمایش قصد داریم هزینه زمانی دستورات *Join*، *Union* را مقایسه کنیم و چنانچه قبلاً گفتیم در مدل CBF و CryptDB این دو دستور پشتیبانی نمی‌شود، لذا از آزمایش‌ها حذف می‌شود. در این آزمایش از دو جدول که توسط ستون رمزگذاری شده ID باهم Join شده‌اند و با آزمایش‌ها ۵۰۰ و ۱۰۰۰ و ۱۰۰۰۰ ردیف استفاده کرده‌ایم که به ترتیب دارای ۱۰۰ و ۲۰۰ و ۱۰۰۰ رکورد در شرط Join استفاده شده هست.

مشخصات دو جدول در تصویر زیر آماده است

students table.			international students		
Name of the Attribute	Data type	Is sensitive?	Name of the Attribute	Data type	Is sensitive?
ID	int	Yes	ID	int	Yes
Name	varchar	Yes	Name	varchar	Yes
SSN	int	Yes	Address	int	Yes
Visa_Type	varchar	Yes	Citizenship	varchar	Yes
Salary	int	Yes			
Department	varchar	Yes			

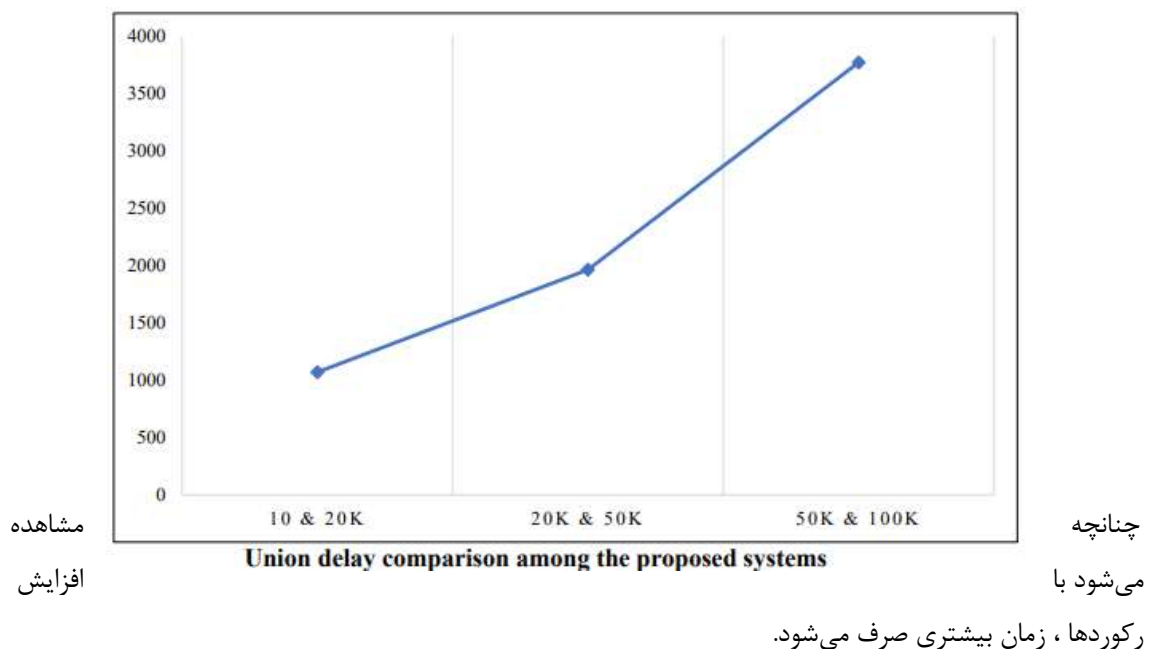
نمودار مقایسه مدل‌ها در عملیات Join به شکل زیر آماده است.



این نمودار نشان می‌دهد که تأخیر BVSAC بالاترین و پس از آن BVSIT است و BVM دارای تأخیر کمتری است. تأخیر زیاد در دستورات Join در همه مدل‌های پیشنهادی به دلیل ماهیت محاسبه Join و مقدار داده‌های رمزگشایی شده است. سپس، پس از اجرای شرایط Join (در QM)، تاپلهای رمزگشایی شده به داخل LinkedHashSet برده شدند تا تکراری را حذف کند.

در پایان، اگرچه تأخیرهای بالایی را تجربه کردیم، اما با استفاده از سیستم‌های پیشنهادی بیش از پایگاه داده‌هایی که با الگوریتم رمزگذاری تصادفی رمزگذاری شده بودند، مانند AES-CBC، دستورات join را اجرا کردیم که این تأخیر برای امنیت بهتر هست.

نمودار هزینه زمان در تمامی مدل‌ها در دستور Union به شکل زیر آماده است و از میانگین کلیه مدل‌ها استفاده شده است زیرا در آزمایش‌ها متوجه شدیم که زمان مدل‌های خیلی به هم نزدیک هست. سعی شده است دستور Union به صورتی انتخاب شود که بدون تغییر توسط QM به سرور ارسال شود. تأخیری که اندازه‌گیری کردیم از زمانی است که QM پرس‌وجو را قطع می‌کند تا نتیجه نهایی پرس‌وجو تشکیل شود. برای حذف موارد تکراری، ما ردیف‌های رمزگشایی شده را از هر دو جدول به یک LinkedHashSet هدایت کردیم.



#### ۵-۲-۸ آزمایش شماره هفت (بازنویسی دستورات توسط QM) :

چنانچه قبلاً گفتیم در رمزگذاری ما از QM استفاده می کنیم تا دستورات را بازنویسی کند تا به سرعت بهتر و تأخیر کمتری در پردازش و نمایش اطلاعات برسیم. این آزمایش میانگین زمان صرف شده توسط هر مدل برای بازنویسی یک پرس و جو را نشان می دهد.

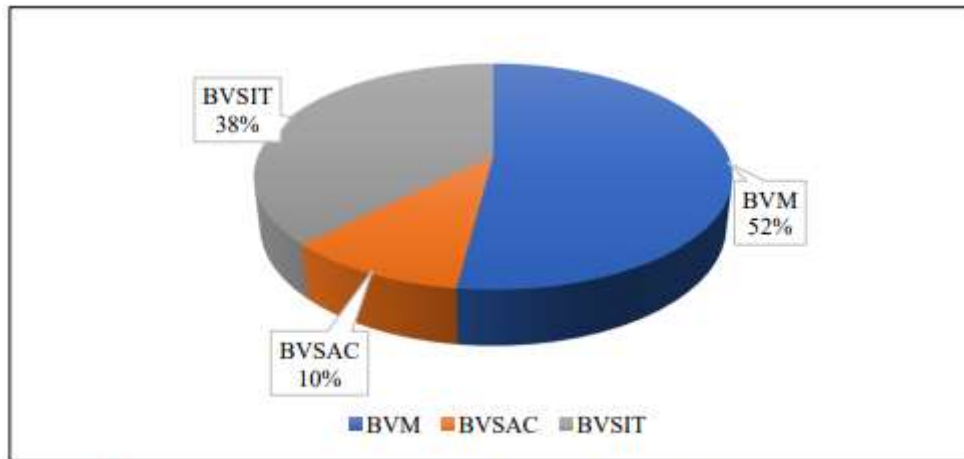
بازنویسی دستورات شامل چند مرحله هست که در جدول زیر برای مدل های انتخابی بررسی شده است

	Query parsing	PT search	BVs store and search (at QM)	Indices decryption
BVM	✓	✓	✓	x
BVSAC	✓	✓	x	x
BVSIT	✓	✓	x	✓

برای محاسبه میانگین ، مجموعه ای از Query را برای موارد مختلف *Select* با یک ، دو و سه بند اجرا کردیم و جدول زیر حاصل شد.

ما روی عبارات *Select* متمرکز شدیم زیرا اینها بخشی از دستورات به روزرسانی و حذف هستند ، بنابراین بررسی عبارت *Select* برای اندازه گیری هزینه بازنویسی دستورات در هر مدل کافی است همچنین عملیات رمزگذاری را از مقادیر آزمایش حذف کردیم.



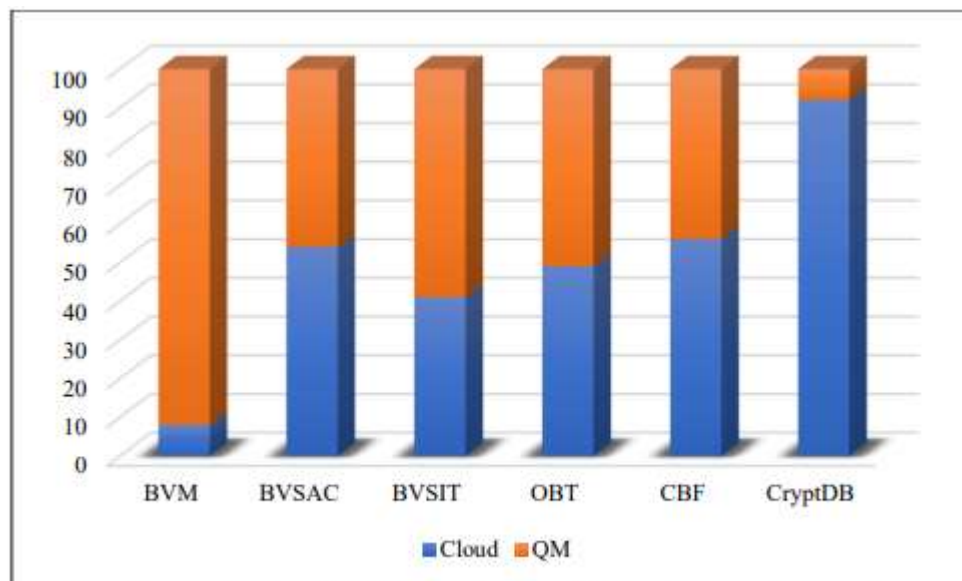


**The average query translation cost at the QM for each model**

طبق این نمودار ، BVM بیشترین تأخیر را دارد (حدود ۵۲٪) زیرا بالاترین درصد تولید و پردازش BV را ارائه می‌دهد ، پس‌از آن BVSIT (حدود ۳۸٪) ، و BVSAC تنها ۱۰٪ از هزینه کل را گرفته است.

### ۳-۵ بررسی محاسبات Cloud و QM :

در این بررسی سعی شده است حجم محاسبات در سرور ابری و QM محاسبه و مقایسه شود . نمودار زیر این مقدار را به‌صورت میانگین درصد زمان محاسبه نشان می‌دهد.



**The percentage of the computation at the QM versus at the cloud server**

چنانچه مشاهده می‌شود که CryptDB بیش از ۹۰٪ تأخیر محاسبات پردازش پرس‌وجو را به سرور منتقل کند ، و کمتر از ۱۰٪ تأخیر محاسبه را به QM واگذار می‌کند (QM معادل پروکسی در CryptDB و CBF است) .

همچنین می‌توانیم ببینیم که BVM بیشترین هزینه محاسبات را در QM دارد زیرا علاوه بر عملیات رمزنگاری ، BV ها به‌صورت محلی در QM پردازش می‌شوند. مدل BVSIT به دلیل رمزگشایی شاخص‌های رمزگذاری شده BV پس از واکنشی آن‌ها از سرور حدود ۶۰٪ محاسبات در QM را دارد. مدل OBT حدود ۵۳٪ محاسبه در QM را دارد زیرا QM به‌طور کامل رمزگشایی می‌کند ، حتی برای اگر دستور خواندن تمامی ستون‌ها نباشد.



از طرف دیگر ، BVSAC و CBF درصد نزدیک به هم دارند و حدود ۵۰٪ محاسبات به سرور منتقل شد زیرا در BVSAC همان طور که قبلاً توضیح داده شد ، BV ها در سرور جستجو می شوند و ما از پردازش آن ها در QM جلوگیری می کنیم. بعلاوه ، در CBF ، ترتیب مجموعه های مختلف شاخص ها و بازنویسی پرس و جو پس از انجام عملیات پیوندی بین مجموعه شاخص ها ، از علل اصلی به دست آوردن درصد محاسبات بالاتر در QM است.

## ۵-۴ بررسی فضای مورد نیاز مدل ها :

در این بخش ، ما هر سیستم را از نظر فضای مورد نیاز در سرور ابری و QM ارزیابی می کنیم. ستون شاخص یک عدد صحیح است که به ۴ بایت نیاز دارد ، در حالی که بقیه ستون ها (ستون هایی که رمزگذاری می شوند) به ۳۳ بایت نیاز دارند (۳۲ بایت از آنجا که کلید مخفی ۳۲ بایت است و ۱ بایت برای ضبط طول نوع داده) علاوه بر این ، ما فضای مورد نیاز برای ذخیره BV ها را در سرور برای هر مدل ارائه می دهیم.

**Storage requirement per record at the server (cloud) of the encrypted student table**

Name of the Attribute	Data type	BVM (bytes)	BVSAC (bytes)	BVSIT (bytes)	OBT	CBF	CryptDB
index	int	4	4	4	65 + 24	28	0
ID	varbinary	17	17	17		284	291
Name	varbinary	33	33	33		66	50
SSN	varbinary	17	17	17		284	291
Visa_Type	varbinary	17	17	17		34	50
Salary	varbinary	17	17	17		284	291
Department	varbinary	33	33	33		66	50
Bit vectors storages		0	4	37*	0	0	0
<b>Total Space required</b>		<b>202</b>	<b>206</b>	<b>239</b>	<b>89</b>	<b>1046</b>	<b>1023</b>

\*37 = 20 bytes for storing bits as individual columns (we have 20 bits that require 20 columns of bit datatype) + 17 bytes for the encrypted index column

همان طور که در دیده می شود ، OBT کمترین نیاز ذخیره سازی را در سمت سرور دارد زیرا کل دیتا به یک بلوک تبدیل می شود و پس از آن BVM و سپس BVSAC قرار دارد. BVSIT به فضای بیشتری نسبت به BVSAC نیاز دارد زیرا ما BV ها را به عنوان یک جدول مستقل و نه ستونی در همان جدول رمزگذاری شده مانند BVSAC ذخیره می کنیم. CBF برای هر ضبط حدود ۱۰۴۶ بایت نیاز دارد که بالاترین میزان فضای مورد نیاز در بین سیستم ها است. CryptDB دومین بالاترین میزان نیاز به فضا را کسب می کند ، حدود ۱۰۲۳ بایت / رکورد.

# فصل ششم

## مقدمه

بحث در خصوص بانک‌های اطلاعاتی و دغدغه‌های افشای اطلاعات همواره پژوهشگران و صاحبان داده را به دنبال بررسی و انتخاب راهکارهای امن و مناسب کشانده و در این مسیر ایده‌های متفاوتی مطرح می‌شود چنانچه تا اینجا مدل‌ها و استراتژی‌های مختلف عنوان گردید و مورد بررسی قرار گرفت.

مدل‌ها و ایده‌هایی که در این پایان‌نامه مطرح گردید را می‌توان به‌عنوان راهکارهای اولیه دانست زیرا دارای نقص و کمبودهایی هست، علی‌الخصوص که بانک‌های اطلاعاتی مانند SQL با فراهم آوردن گزینه رمزگذاری و رمزگشایی به‌عنوان دو ابزار درونی، این دغدغه را کاهش داده است.

در این تحقیق، تمامی تمرکز بر یک سیستم به‌عنوان واسطه بنام QM بوده تا بتواند با ترجمه دستورات و درخواست‌ها به بانک اطلاعاتی و نگهداری اطلاعات به‌صورت رمزگذاری شده، از افشای آن توسط کاربران غیرمجاز و یا صاحبان سرورها جلوگیری نماید اما آنچه به‌عنوان راهکار در این QM در نظر گرفته‌شده دارای اشکالات و نقص‌های فراوانی هست و در بعضی بخش‌ها حتی به بن‌بست خواهیم رسید.

البته ایده و استراتژی کاهش حجم اطلاعات در فراخوانی و انتخاب اطلاعات را می‌توان ایده مناسبی دانست و در راهکارهای بعدی استفاده نمود. در ادامه راهکارها و ایده بهتری را جهت تحقیقات آتی اعلام می‌داریم.

## ۶-۱ بحث و پیشنهادها:

عواملی که باعث افزایش کارایی سیستم می‌شوند:

تعداد پارتیشن‌های یک ستون در کارایی مدل‌های پیشنهادی ما نقش بسزایی دارد. هرچه تعداد پارتیشن‌ها بیشتر باشد، رکوردهای کمتری از ابر گرفته می‌شود زیرا کارایی مدل‌ها تحت تأثیر تعداد رکوردهای رمزگذاری شده واکشی شده از ابر است.

از آنجاکه ما ردیف‌ها را بر اساس تقسیم ویژگی‌ها به بیت کد می‌کنیم، جستجوی بیت‌ها سریع است و ما نگران طولانی بودن PT نیستیم، زیرا این امر به کاهش تعداد ردیف‌های رمزگذاری شده بازایی شده از ابر کمک می‌کند. البته داشتن تعداد پارتیشن متعادل خود نیز چالشی هست که طراح و مالک داده‌ها با آن روبرو هست.

## ۶-۱-۱ بحث :

همان‌طور که در آزمایش‌های انجام‌شده مشاهده شد، ما هر مدل پیشنهادی را ارزیابی کردیم. سپس، ما مدل‌ها را از نظر سرعت پردازش پرس‌وجو، محاسبات سربار، نیازهای ذخیره‌سازی در QM و سرور ابری مقایسه کردیم. از نظر سرعت، مهم‌ترین دلیل تأخیر، مقدار داده‌هایی است که باید رمزگشایی شود. با این حال، تمام نمونه‌های پیشنهادی در نظر دارند مقدار داده‌های رمزگذاری شده بازیابی شده از ابر را به کمتر از ۳۵ درصد کاهش دهند و همان‌طور که دیده شد علت ایجاد تغییر در سرعت پردازش پرسش به نحوه و محل ذخیره‌سازی BV مرتبط است.

تأخیر BVM عمدتاً تحت تأثیر انتقال BV از هارد دیسک سخت به حافظه است، در حالی که در BVSIT، ما با ذخیره‌سازی BV ها به‌عنوان یک جدول رابطه‌ای مستقل در ابر، از این موضوع جلوگیری می‌کنیم. در BVSAC نیز ما BV ها را به‌عنوان یک ستون اضافی در همان جدول رمزگذاری شده در ابر ذخیره کردیم.

برای سیستم‌های تک کاربر که توسط فضای ابری محدود شده‌اند، BVM مدل بهینه است زیرا برای مدیریت و پردازش جداول رابطه‌ای رمزگذاری شده به حداقل فضای اضافی نیاز دارد. با این حال، اگر بیش از یک کاربر منبع محاسبات مشترک را به اشتراک بگذارد، محاسبه سربار محاسبه و میزان مصرف حافظه در QM دچار مشکلاتی خواهد بود، زیرا برای اجرای یک پرس‌وجو، کل مجموعه BV های جدول باید در حافظه اصلی بارگیری شوند. بنابراین، BVM برای یک سیستم چندکاربره توصیه نمی‌شود، و مدل مطلوب BVSAC است که از روش BVSIT استفاده می‌نماید. همچنین BVSAC کمترین زمان اجرا، دومین کمترین میزان فضای مورد نیاز در ابر و کمترین سیستم مورد نیاز برای ذخیره‌سازی را در QM دارا هست، به این معنی که برای سیستم‌های تک کاربر و چندکاربره مناسب است.

از نظر امنیت، همه مدل‌های پیشنهادی در بالاترین سطح امنیتی بودند. در CBF، اگر ارائه‌دهندگان سرویس ابر تباری نکنند، امنیت بالا بود. به هر حال تضمین شده نیست و رتبه سطح امنیت پایینی دارد.

در CryptDB، نویسندگان پیشنهاد نکردند لایه‌های رمزگذاری را عقب بگذاریم که باعث کاهش سطح امنیت شود. مدل OBT در صورت استفاده از الگوریتم رمزگذاری تصادفی، ایمن بود و گزینه سطح امنیت کاهش پیدا می‌کند. برای جداول با ستون‌های حساس زیاد، مانند هزاران، عملکرد OBT به خطر می‌افتد زیرا داده‌های تمام زمینه‌های هر سطر در یک بلوک به هم پیوسته‌اند.

در CBF، هزاران سرویس ابری برای ذخیره ستون‌های توسعه یافته مورد نیاز بود، که این سیستم را غیرعملی می‌کند. در مقابل، سیستم‌های پیشنهادی چنین مشکلی ندارد زیرا رمزگذاری مبتنی بر سلول را انجام دادیم. ما تمام داده‌ها را در همان سرور ابری ذخیره کردیم که باعث تسریع در زمان پردازش جستجو برای پایگاه‌های مقیاس‌پذیر می‌شود.

## ۶-۱-۲ پیشنهادها :

آنچه در این پایان‌نامه بررسی و مورد بحث قرار گرفت راهکارها و مدل‌های ایجاد و ارتباط با یک بانک اطلاعاتی رمزگذاری شده، بوده و به‌عنوان راهکار کلی برای تمامی انواع پایگاه اطلاعاتی رابطه‌ای قابل انجام هست. امروزه شرکت‌های ارائه‌دهنده پایگاه‌های داده رابطه‌ای مانند SQL با ایجاد دستورات و ابزارهای آسان و قابل فهم‌تر، قدم بزرگی در ایجاد بانک‌های اطلاعاتی رمزگذاری شده برداشته‌اند به صورتی که دیگر نیاز به یک نرم‌افزار واسط نباشد.

آن چیزی که در این پایان‌نامه ایدئال بود، ایده کاهش تعداد ردیف پردازش رمزگذاری هست زیرا حتی با داشتن ابزارهایی مانند آنچه بیان شد، باز باید تمامی ردیف‌ها مورد پردازش رمزگشایی قرار گیرد، پس ایده کاهش تعداد ردیف جهت پردازش نهایی به‌عنوان ایده مناسبی می‌تواند در کنار این ابزارها و امکانات منجر به ایجاد یک بانک و سرویس اطلاعاتی امن گردد. این نکته را باید در نظر بگیریم که عملیات‌های پایگاه داده خیلی گسترده‌تر و متنوع‌تر از آن چیزی است که در این پایان‌نامه مورد بررسی قرار گرفت و در مدل‌های مختلف اجرا گردید مثلاً عملیات‌هایی مانند Like که قابلیت جستجو در محتوای فیلدهای کاراکتری هست در این پایان‌نامه قرار نگرفته و عملاً برای اجرای آن نیاز به رمزگشایی در تمامی اطلاعات هست که ناگزیر منجر به هزینه بالای اجرا می‌گردد.

همچنین ایده نرم‌افزار واسط به‌عنوان یک ایده مناسب برای کار با بانک اطلاعاتی رمزگذاری شده را می‌توان به‌صورت یک حافظه موقت رمزگشایی‌شده طراحی و اجرا نمود چنانچه این نرم‌افزار هر زمان که اطلاعاتی از جدولی درخواست شود آن را به‌طور کامل رمزگشایی و در یک حافظه موقت قرار دهد و تمامی تغییرات هم‌زمان بر حافظه موقت و بانک اطلاعاتی لحاظ گردد. در این حالت بار هزینه تنها یک‌بار برای مدت اجرای حافظه موقت اتفاق می‌افتد و در تمامی ارجاعات، این نرم‌افزار با اطلاعات رمزگشایی‌شده کار خواهد کرد و دسترسی به آن نیز پس از گذراندن مجوزهای لازم هست.

## ۶-۲ نتیجه‌گیری :

رایانش ابری یک محیط محاسباتی جذاب برای انواع کاربران و شرکت‌ها است. اما، نقض حریم خصوصی، نه‌تنها توسط مهاجمان مخرب بلکه توسط ارائه‌دهندگان کنجکاو، نقطه‌ضعف این نوع خدمات است، زیرا کاربران کنترل دسترسی بر داده‌های خارج از منابع را از دست می‌دهند. راه‌حل‌های زیادی برای این مشکل وجود دارد و رمزگذاری داده‌ها یک راه‌حل مؤثر است. با این حال، اجرای نمایش داده‌های SQL بر روی داده‌های رمزگذاری شده چالش‌برانگیز است، به‌خصوص اگر از الگوریتم رمزگذاری تصادفی مانند AES-CBC برای رمزگذاری استفاده شود. در این تحقیق، ما ابتدا QM را معرفی می‌کنیم، یک سرور قابل اعتماد، که به‌عنوان واسطه بین سرور ابری و کاربر (ها) کار می‌کند و همه فرایندهای رمزنگاری را انجام می‌دهد. علاوه بر این، ما یک تکنیک جدید برای نمایه‌سازی بر اساس پارتیشن‌های از پیش تعریف‌شده برای هر ویژگی حساس طراحی می‌کنیم، و سپس هر ردیف از داده را به‌صورت بیتی رمزگذاری می‌کنیم.

از بیت‌ها برای بازیابی ردیف‌های انتخابی برای یک پرس‌وجو خاص استفاده می‌شود که دامنه تاپل‌های رمزگذاری شده بازیابی شده را به حداقل می‌رساند. بر اساس این طرح رمزگذاری، ما سه سیستم ایمن مختلف را پیشنهاد داده‌ایم که هر یک از آن‌ها برای ذخیره و نگهداری داده‌های شاخص (به‌عنوان مثال، BV) یا به‌صورت محلی در QM یا با انتقال BV ها به فضای ابری، از روش متفاوتی استفاده می‌کنند.

برای هر نمونه اولیه پیشنهادی، ما الگوریتم‌های مختلفی را برای انجام پرس‌وجو از عملگرهای مختلف جبر رابطه‌ای SQL طراحی می‌کنیم و آن را در برابر سناریوهای حمله، مانند حملات استنتاج، مقاوم می‌کنیم. ما مدل‌های خود را با پیاده‌سازی آن‌ها و مقایسه عملکرد آن‌ها با سیستم‌های پیشرفته مانند CryptDB، آزمایش می‌کنیم.

ما آن‌ها را از نظر نیاز به زمان اجرا و فضای مورد نیاز ارزیابی می‌کنیم. درمی‌یابیم که سیستم‌های پیشنهادی در مقایسه با اکثر سیستم‌های رقیب، به زمان اجرا و فضای کمتری نیاز دارند.

# منابع و مراجع

- ۱) R. A. Popa, C. M. S. Redfield, N. Zeldovich, and H. Balakrishnan, "CryptDB: Processing queries on an encrypted database." Commun. ACM, vol. ۵۵, no. ۹, pp. ۱۰۳-۱۱۱, ۲۰۱۲.
- ۲) Y. D. Jang and J. H. Kim, "A comparison of the query execution algorithms in secure database system," Int. J. Electr. Comput. Eng., vol. ۶, no. ۱, pp. ۳۳۷-۳۴۳, ۲۰۱۶.
- ۳) W. Wang, Y. Hu, L. Chen, X. Huang, and B. Sunar, "Exploring the feasibility of fully homomorphic encryption," IEEE Trans. Comput., vol. ۶۴, no. ۳, pp. ۶۹۸-۷۰۶, ۲۰۱۳.
- ۴) A. Alsirhani, P. Bodorik, and S. Sampalli, "Improving database security in cloud computing by fragmentation of data," in ۲۰۱۷ International Conference on Computer and Applications (ICCA), ۲۰۱۷, pp. ۴۳-۴۹.
- ۵) H. Hacigümüş, B. Iyer, C. Li, and S. Mehrotra, "Executing SQL over encrypted data in the database-service-provider model," in Proceedings of the ۲۰۰۲ ACM SIGMOD international conference on Management of data, ۲۰۰۲, pp. ۲۱۶-۲۲۷.
- ۶) H. Hacigümüş, B. Iyer, C. Li, and S. Mehrotra, "Executing SQL over encrypted data in the database-service-provider model," in Proceedings of the ۲۰۰۲ ACM SIGMOD international conference on Management of data, ۲۰۰۲, pp. ۲۱۶-۲۲۷.
- ۷) A. Boicea, F. Radulescu, C.-O. Truica, and C. Costea, "Database encryption using asymmetric keys: a case study," in ۲۰۱۷ ۲۱st International Conference on Control Systems and Computer Science (CSCS), ۲۰۱۷, pp. ۳۱۷-۳۲۳.
- ۸) O. M. Ben Omran and B. Panda, "A new technique to partition and manage data security in cloud databases," in The ۹th International Conference for Internet Technology and Secured Transactions (ICITST-۲۰۱۴), ۲۰۱۴, pp. ۱۹۱-۱۹۶.
- ۹) S. Cui, M. R. Asghar, S. D. Galbraith, and G. Russello, "P-McDb: Privacy-preserving search using multi-cloud encrypted databases," in ۲۰۱۷ IEEE ۱۰th International Conference on Cloud Computing (CLOUD), ۲۰۱۷, pp. ۳۳۴-۳۴۱.
- ۱۰) O. M. Omran, "Data Partitioning Methods to Process Queries on Encrypted Databases on the Cloud," ۲۰۱۶.
- ۱۱) V. H. Hacigumus, B. R. Iyer, and S. Mehrotra, "Query optimization in encrypted database systems." Google Patents, Mar-۲۰۱۰.
- ۱۲) V. H. Hacigumus, B. R. Iyer, and S. Mehrotra, "Query optimization in encrypted database systems." Google Patents, Mar-۲۰۱۰.
- ۱۳) "Package javax.crypto." [Online]. Available: [https://docs.oracle.com/javase/۷/docs/api/javax/crypto/package-summary.html#package\\_description](https://docs.oracle.com/javase/۷/docs/api/javax/crypto/package-summary.html#package_description)

## **ABSTRACT**

Nowadays, the issue of maintaining and securing databases with sensitive and strategic information at the level of a company and sometimes even at the level of countries, is studied and researched by software and data researchers. Part of this security is the prevention of damage and disclosure of information by hackers and sabotage groups, and another worrying part is the disclosure of information by users and providers of cloud services because the second group of hardware and software layers common for security , Have and have free access to information resources.

Cloud servers and users of a system, due to their free access to information, can easily access the classified and sensitive information of a database and have serious consequences for the owners of the database. The best and perhaps the only idea that comes to mind is to encrypt information so that it can only be displayed after controlling access within the software.

But encryption and consequently decryption require heavy calculations and server busy time, and if we can not advance this method with the right and appropriate method, in the high volume of information and after a while, the use of the system becomes difficult and even impossible.

Familiarity with the types of encryption methods, the amount of load and the cost of its implementation and determining the correct method for encryption, helps us to reach a suitable solution and method to maintain the confidentiality of information.