



ساده‌سازی و تسریع خطوط لوله داده در برنامه‌های مالی دیجیتال و بیمه

مجری: اهام نوبری

.....





اطلاعات مجری و همکاران

دکتر الهام نوبری استادیار دانشکده علوم دانشگاه علم و فناوری مازندران	مجری طرح
دکتر امیر تیمور پاینده عضو هیات علمی دانشگاه شهید بهشتی	مشاور طرح
دکتر فرزانه خامسیان عضو هیات علمی پژوهشکده بیمه	ناظر علمی طرح
پاییز ۱۴۰۱	تاریخ شروع اجرای طرح
بهار ۱۴۰۲	تاریخ اتمام طرح



بیان مسأله و ضرورت انجام طرح

- فعالیت‌های پیچیده روزانه در شرکتهای بیمه‌ای و مالی مانند محاسبه ریسک، بدهی یا اعتبار
- بزرگی حجم داده‌ها و زمان‌بری تجزیه و تحلیل داده‌ها بیشتر از زمان مجاز برای یک کسب و کار
- نگهداری داده‌های اخیر توسط پایگاه داده از نوع SQL و NoSQL فقط برای مدت کوتاه
- مشکلات استفاده از پایگاه داده‌های تحلیلی کندتر مانند انبارهای داده و نیز دریاچه‌های داده
- لزوم طراحی معماری مناسب پایگاه داده جهت فراهم کردن موارد زیر

صرفه جویی در هزینه فناوری اطلاعات

متناسب بودن هزینه عملیات و پشتیبانی با زمان پردازش
هرچه پردازش کوتاهتر باشد تیم‌های فنی کوچکتر
در پلتفرم‌های ابری هزینه‌ها متناسب با زمان پردازش



بهبود بهره وری

تقسیم کل زمان صرف شده بین زمانی که تحلیل‌گران داده و دانشمندان داده به صورت دستی با داده‌ها کار می‌کنند و زمانی که الگوریتم‌های کامپیوتری پردازش داده را انجام می‌دهند

کاهش خطرات نظارتی و عملیاتی

انجام فرآیندهای دسته‌ای روزانه، هفتگی و ماهانه در بازه زمانی محدود

ارائه قابلیت‌ها و خدمات جدید

استفاده‌ی شرکتهای مالی و بیمه‌ای از توانایی خود مبتنی بر داده‌ها برای پردازش اطلاعات بیشتر در زمان کمتر و در نتیجه ایجاد مزیت رقابتی و ارائه محصول بهتر



اهداف / سوالات / فرضیه های طرح

- خطوط لوله داده چیست؟
- خطوط لوله داده های بانکی و بیمه معمولاً از چه مراحل تشکیل می شوند؟
- با توجه به اهمیت سرعت بخشیدن به خطوط لوله داده، برای تسریع چه کاری می توان انجام داد؟
- در پایگاه داده سنتی ذخیره داده دریافت داده، تبدیل داده ها و تولید خروجی مورد نیاز چگونه انجام می شود؟
- پایگاه داده LeanXcale چطور می تواند بهتر عمل کند؟



سوابق مطالعاتی و پژوهشی مربوطه

LeanXcale با این چشم‌انداز متولد شد تا فناوری پایگاه داده بزرگی در آینده باشد که شکاف بین دنیای داده‌های عملیاتی و دنیای داده‌های تحلیلی را پر کند. LeanXcale ریشه در تحقیقات فنی عمیق در سیستم‌های توزیع شده و مدیریت داده دارد. اساتید آزمایشگاه سیستم‌های توزیع شده تصمیم گرفتند تمام تحقیقات انجام شده برای بیش از ۱۵ سال تحقیق را کنار بگذارند و از ابتدا شروع کنند تا یک مدیر تراکنشی کاملاً متفاوت را طراحی کنند که بتواند بدون محدودیت مقیاس شود. بیش از نه ماه تحقیق برای تولید اولین نسخه از الگوریتم طول کشید، اما نتایج زیبا و ظریف بودند، راه حلی عالی که بزرگترین و مشکل‌سازترین گلوگاه در پایگاه‌های داده را برای چندین دهه حل کرد. سپس تیم تحقیقاتی به دنبال بودجه تحقیق و توسعه در اروپا بود، یک سری کمک هزینه دریافت کرد نمونه اولیه تحقیقاتی در سال ۲۰۱۰ تولید شد و تا سال ۲۰۱۵ تکمیل شد. در سال ۲۰۱۰ یک حق اختراع در مورد چگونگی مقیاس‌پردازی پرس و جو ثبت شد و در سال ۲۰۱۱ یک حق اختراع دیگر ثبت شد.



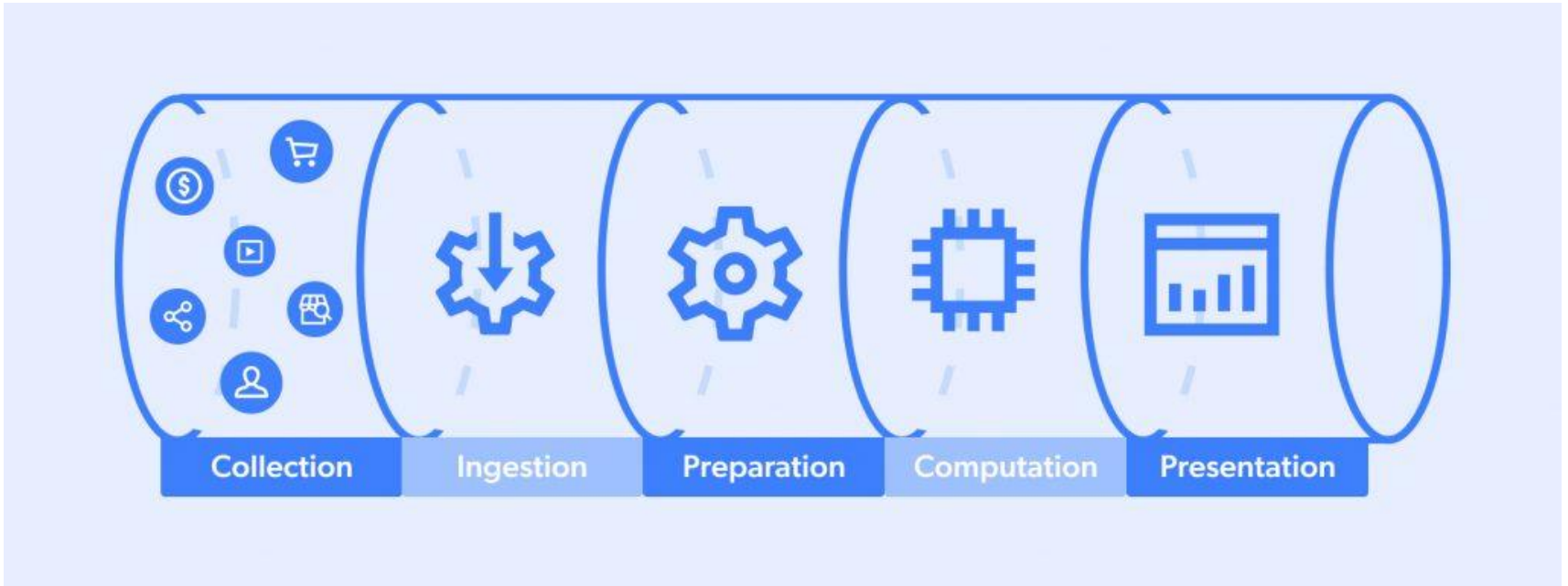
خطوط لوله داده چیست؟

خط لوله داده ها به مجموعه فرآیندهایی گفته می شود که برای جمع آوری، تبدیل، تحلیل و انتقال داده ها به کار می روند تا بتوان در قالبی مناسب از آن ها برای برنامه های کاربردی مختلف استفاده کرد. به طور معمول، برای ساخت این خطوط انتقال داده ها از ابزارها و فرآیندهایی استفاده می شود که در هر مرحله از جمع آوری و پردازش داده ها، داده ها را به صورت خودکار از یک مرحله به مرحله دیگر منتقل می کنند.

در یک Data Pipeline، داده ها معمولاً از منابع مختلف مانند پایگاه های داده، فایل ها، سرویس های وب و سایر منابع جمع آوری می شوند. سپس، این داده ها به صورت خودکار در مراحل مختلفی مانند تبدیل، تحلیل، پاکسازی و ترکیب داده ها پردازش می شوند. در نهایت، داده ها به سیستم مورد نظر انتقال داده می شوند تا برای استفاده در برنامه های کاربردی و سایر سیستم ها موجود قرار گیرند. استفاده از Data Pipeline در سازمان ها، می تواند به بهبود کارایی و دقت در جمع آوری و پردازش داده ها و به کاهش زمان و هزینه مورد نیاز برای پردازش داده ها کمک می کند.



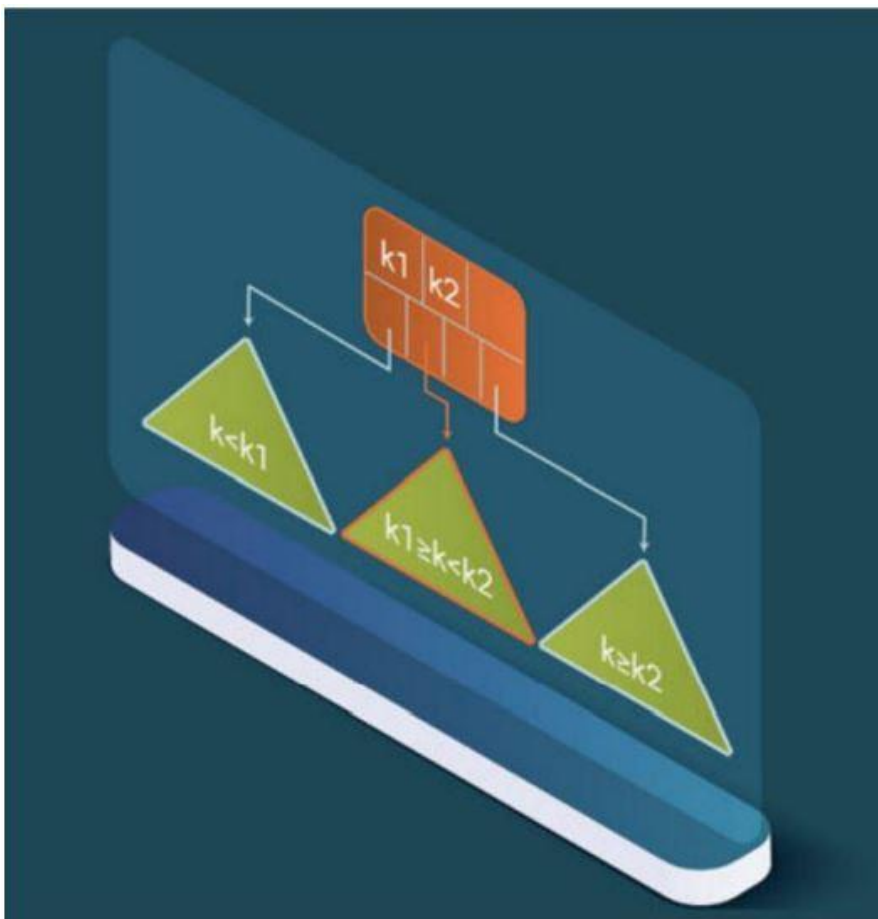
.....



- دریافت داده
- دریافت دادهها بسیار زیاد و صرف زمان قابل توجه
- چالش
- تسريع جذب دادهها بدون به خطر افتادن سرعت پرس و جو برای مرحله بعدی

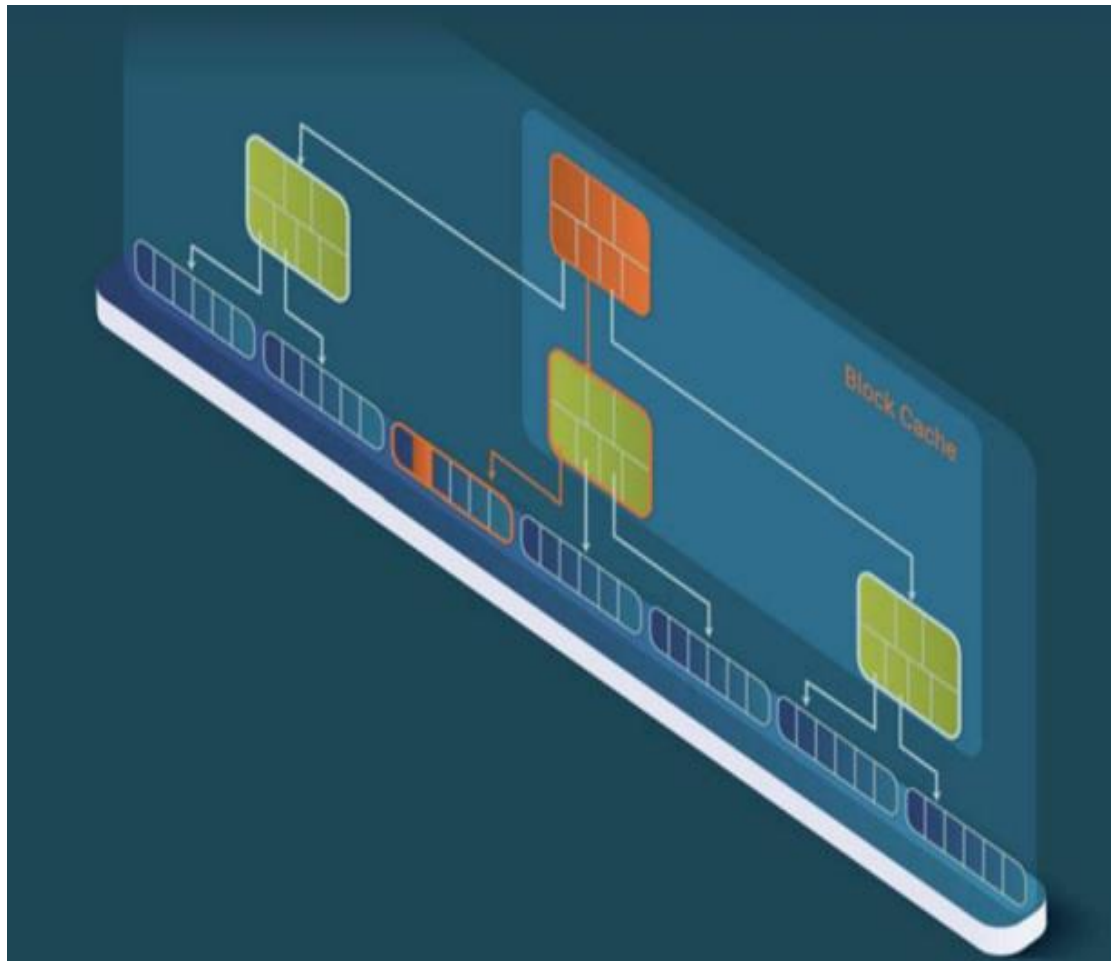


راه حل؟؟؟؟



راه حل سنتی جذب داده در SQL

استفاده از ساختار داده پایدار به نام درخت
B+ یافتن داده‌ها با پیچیدگی لگاریتمی



حافظه cash

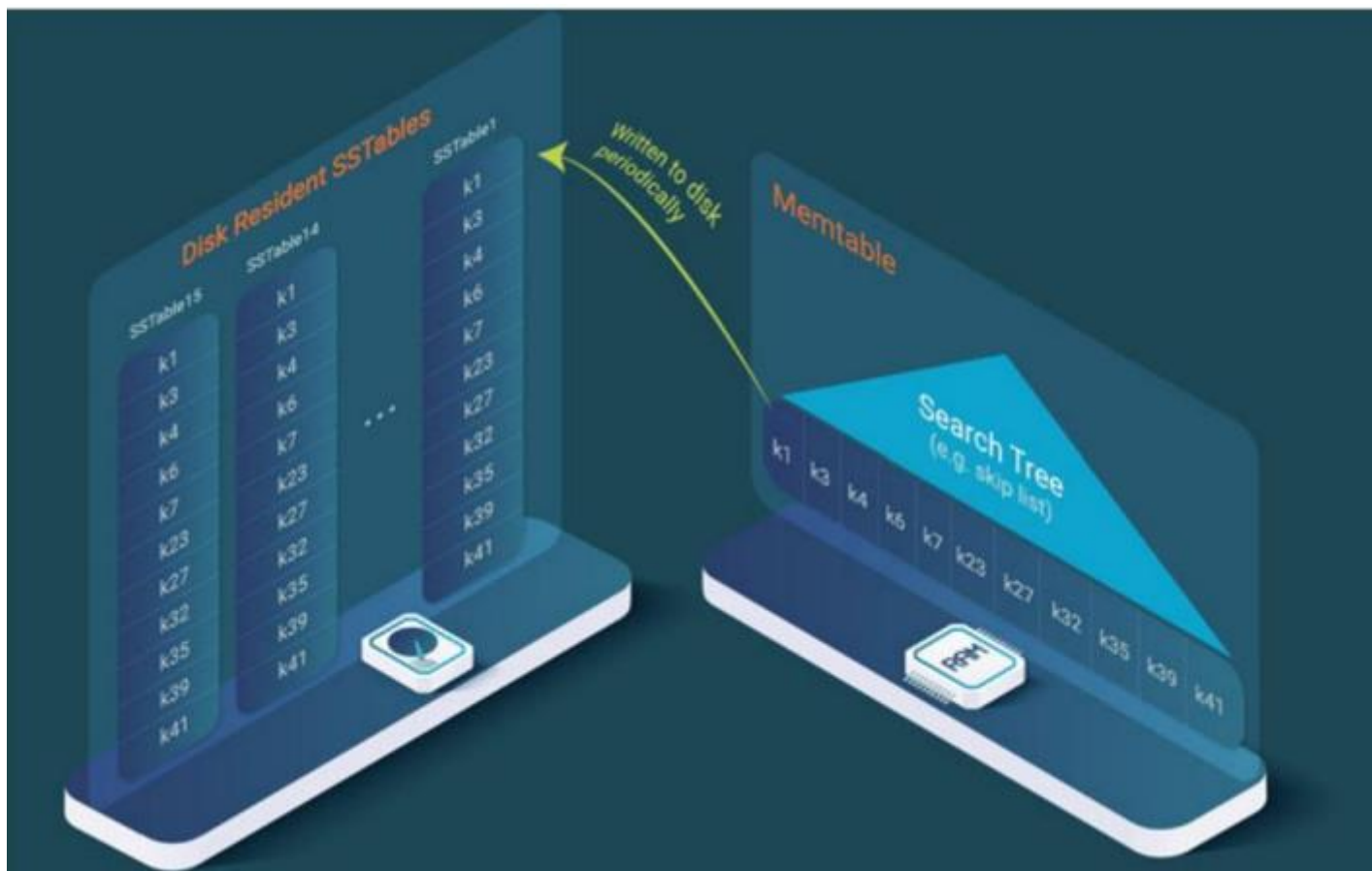
با افزایش اندازه
جدول، دریافت داده‌ها
گران‌تر و گران‌تر می
شود زیرا سطوح
بیشتری در درخت
+ B به وجود می‌آید و
هر ردیفهزینه بیشتری
برای خواندن و خارج
کردن از حافظه پنهان
دارد



.....

سریعترین پایگاه‌های داده **NoSQL** برای دریافت داده‌ها، از ذخیره‌سازی داده کلید-مقدار و پایگاه‌های داده ستونی استفاده می‌کنند. راه حل‌های مبتنی بر هش حتی از پرس و جوهای محدوده ساده پشتیبانی نمی‌کنند.

راه حل‌های دیگر مبتنی بر رویکرد جداول مرتب‌سازی شده رشته‌ای (**SSTables**) است که از آنچه در **Google Big Table** انجام می‌شود، تقلید می‌کنند.





مرحله دوم خطوط لوله داده در بانکداری و بیمه

فرآیند تجزیه و تحلیل و غنی‌سازی داده‌هاست
تجمع اطلاعات با استفاده از پرس‌وجوهای SQL

انبارهای داده در پاسخ به این نوع پرس و جوها خوب هستند اما از به روزرسانی‌ها پشتیبانی نمی‌کنند و فقط پرس و جوهای خواندنی را پشتیبانی می‌کنند

پایگاه های داده NoSQL

عدم پشتیبانی از داده‌های انبوه
عدم توانایی پس‌وجو برای راه‌حل‌های مبتنی بر هش
گرانی راه‌حل‌های مبتنی بر SSTable
فقدان ویژگی‌های ACID



ویژگی های ACID

ویژگی های **ACID** در پایگاه داده به مجموعه ای از اصول اشاره دارد که پردازش قابل اعتماد و سازگار تراکنش های پایگاه داده را تضمین می کند

Atomicity: این ویژگی تضمین می کند که یک تراکنش به عنوان یک کار واحد و غیرقابل تقسیم تلقی می شود.

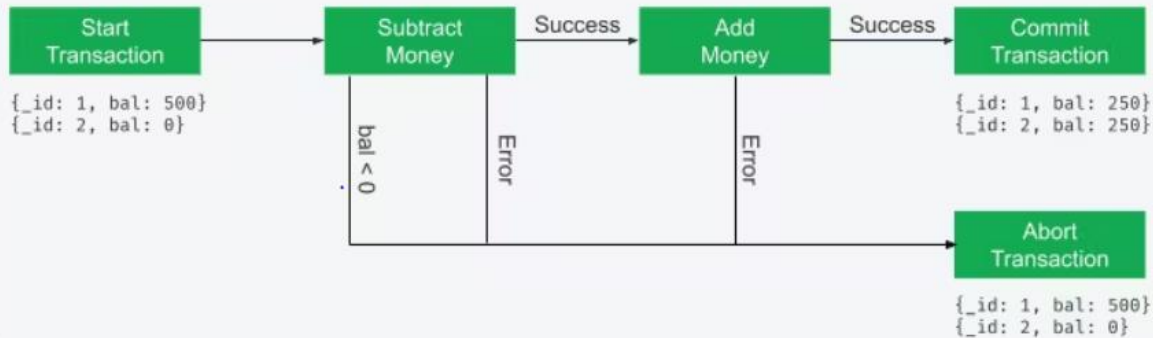
Consistency: این ویژگی تضمین می کند که یک تراکنش پایگاه داده را از یک وضعیت معتبر به حالت دیگر می آورد.

Isolation: این ویژگی تضمین می کند که تراکنش ها جدا از یکدیگر انجام می شوند، حتی اگر به طور همزمان انجام شوند. این از تداخل بین تراکنشها جلوگیری می کند.

Durability: این ویژگی تضمین می کند که به محض انجام یک تراکنش، تغییرات آن دائمی است و حتی در صورت خرابی سیستم قابل بازگشت نیست.



Transaction 1: Transfer \$250 from Account 1 to Account 2



Power Failure

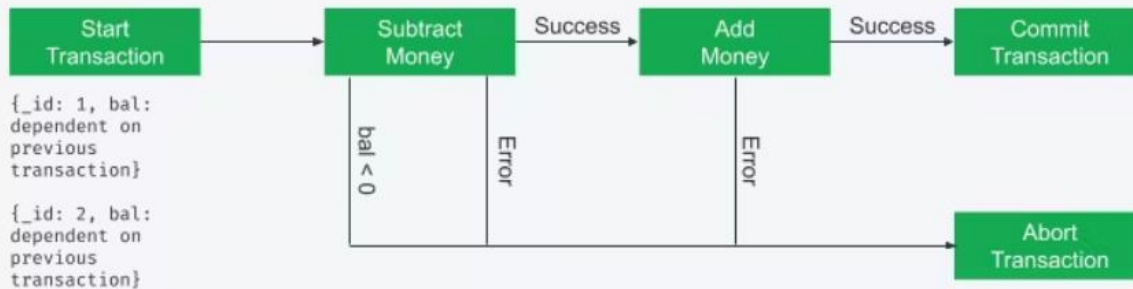
System Online

{_id: 1, bal: 250
if the transaction
committed or 500 if
the transaction was
aborted}

{_id: 2, bal: 250
if the transaction
committed or 0 if
if the transaction
aborted}

Waits for Transaction 1 to Finish

Transaction 2: Transfer \$100 from Account 1 to Account 2





ایجاد خروجی مورد نیاز

مرحله سوم خطوط لوله داده

تولید گزارش‌ها یا فایل‌های خروجی

پایگاه‌های داده **SQL** عملیاتی سنتی مناسب نیستند زیرا برای بسیاری از پرس و جوهای کوتاه به جای پرس و جوهای تحلیلی بزرگ بهینه‌سازی شده‌اند

پایگاه داده‌های **NoSQL** برای این فرآیند بد هستند، زیرا معمولاً از بیان لازم برای پرس و جو و استخراج داده‌ها پشتیبانی نمی‌کند





چگونه LeanXcale خطوط لوله داده را ساده و تسریع می کند

فراتر از قابلیت‌های معمولی یک پایگاه داده SQL، LeanXcale دارای مجموعه‌ای از ویژگی‌های منحصر به فرد است که آن را برای ساده‌سازی و تسریع خطوط لوله داده بهینه می‌کند. این ویژگی‌ها دو اثر دارند:

به لطف یک موتور ذخیره کلید-مقدار کارآمد و پارتیشن‌بندی دوبعدی سرعت درج داده‌ها افزایش می‌یابد

قابلیت موازی‌سازی فرآیند در چندین رشته اجرایی و در عین حال اجتناب از قفل شدن پایگاه داده به دلیل تجمع آنلاین و مقیاس‌پذیری افقی خطی





نرخ‌های درج بالا

Leanxcale رویکردی مشابه درختان LSM دارد.

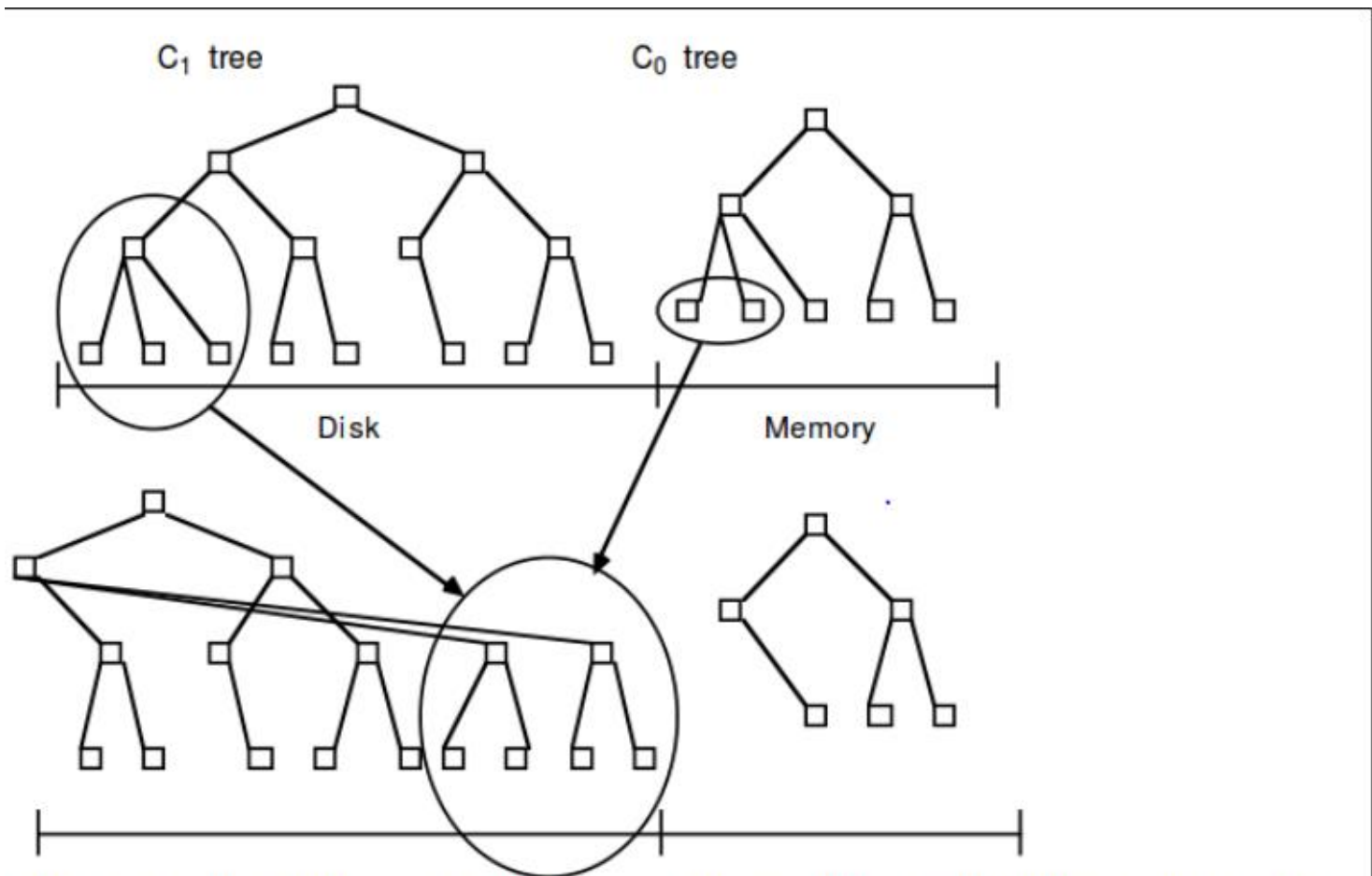
درخت LSM (Merge-Structured-Log) یک ساختار داده‌ای است که در علوم کامپیوتر

برای ذخیره و بازیابی موثر داده‌ها استفاده می‌شود. این یک ساختار داده مبتنی بر دیسک است که

برای بارهای کاری فشرده، مانند مواردی که در پایگاه داده‌ها و سیستم‌های فایل یافت می‌شود، بهینه

شده است.





2.2. Conceptual picture of rolling merge steps, with result written back to disk





ویژگی‌های منحصر به فرد LeanXcale

پارتیشن بندی دو بعدی

leanXcale به طور موثر با داده‌های تاریخی طولانی سروکار دارد.

درج داده در جدول اغلب به لحاظ زمانی دارای خصیصه موقتی در اطراف ستون زمان یا ستون خودافزایشگر است .

قطعات داده را به طور خودکار تقسیم می‌کند، بنابراین داده‌ها خیلی بزرگ نمی‌شوند.





ID	User Name	Profile Data	Connections	Content authored
1				
2				
3				
4				

Vertical partition 1

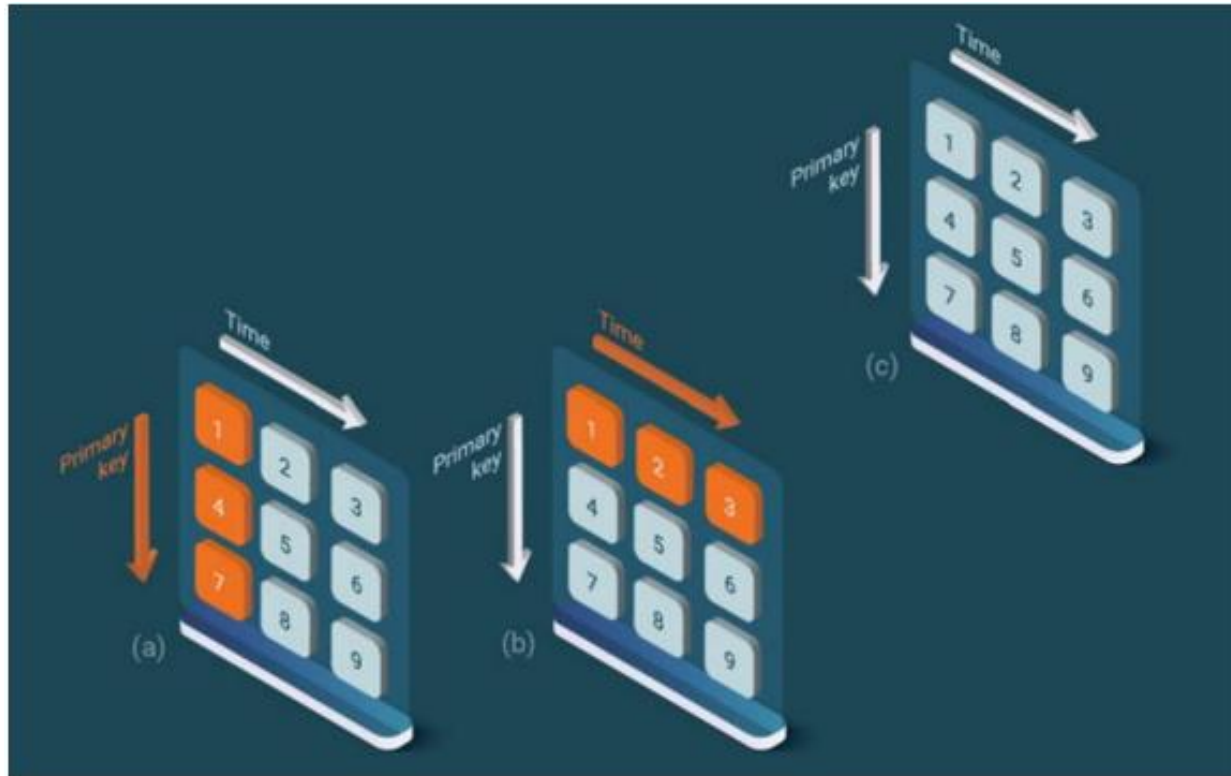
ID	User Name	Profile Data
1		
2		
3		
4		

Vertical partition 2

ID	Connections	Content authored
1		
2		
3		
4		

Vertical partitioning







انبوه‌های آنلاین

روش کنترل همزمانی جدید به نام کنترل همزمانی چند نسخه‌ای معنایی

امکان تجمیع در یک ردیف در زمان واقعی با سطوح همزمانی بالا (مثلاً هزاران در ثانیه) بدون هیچ گونه اختلافی و با سازگاری کامل **ACID**



مقیاس پذیری

مقیاس پذیری افقی را با مقیاس بندی سه لایه پایگاه داده:
موتور ذخیره سازی
مدیر تراکنش
موتور پرس و جو
لایه ذخیره سازی یک فناوری اختصاصی به نام **KiVi**
است که به شکل کلید-مقدار رابطه ای است.





Ultra-Scalable
Transactions

Transaction Mng

OLAP & OLTP
SQL Queries

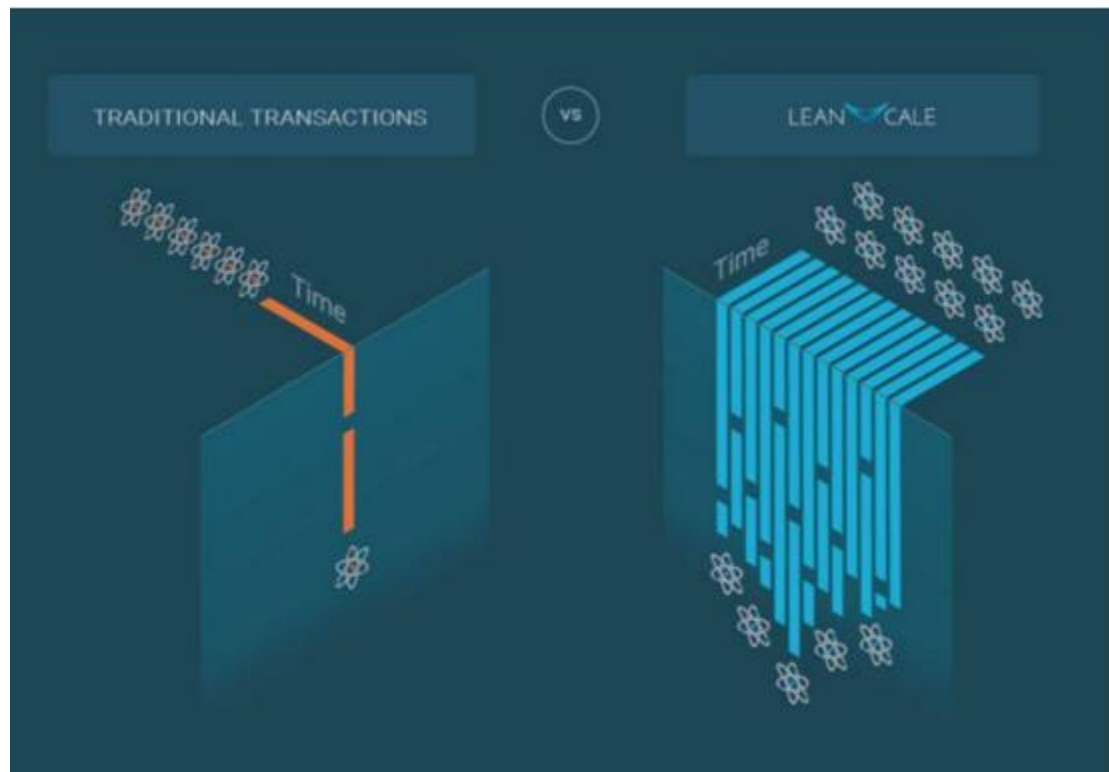
Query Engine

Key-Value
Data Store

KiVi Storage



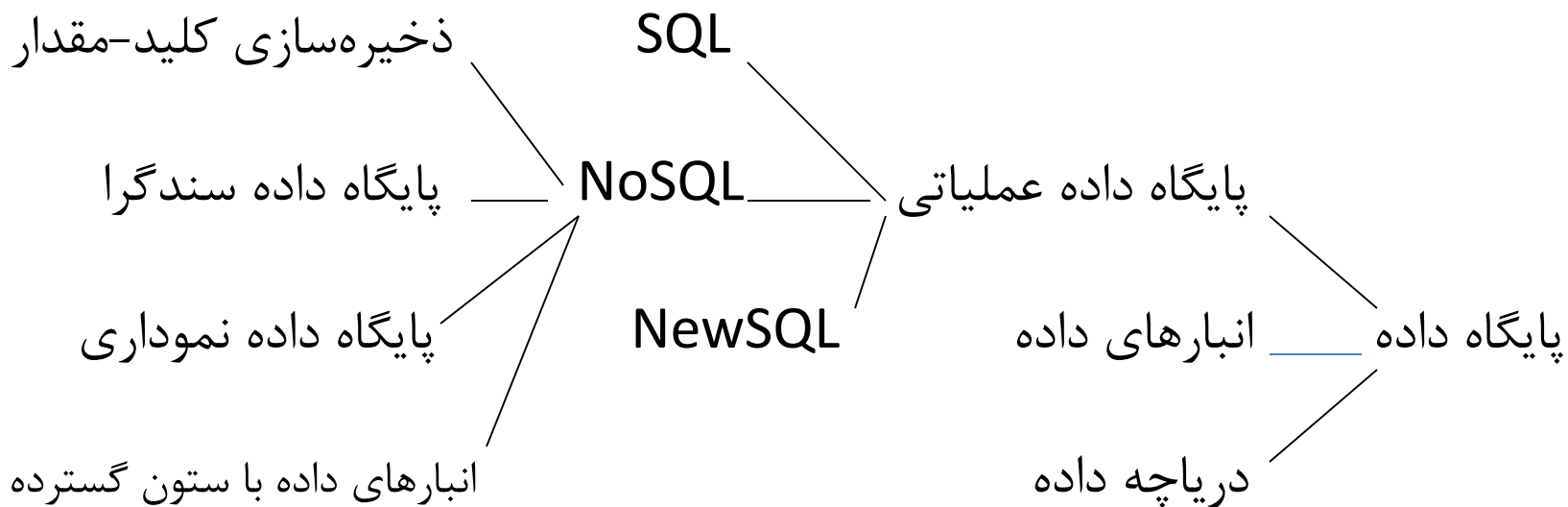
اساساً همه تراکنشها را به صورت موازی و بدون هیچ هماهنگی انجام می دهد





لوله گذاری

به دلیل استفاده از پایگاه داده‌های مختلف، داده‌ها باید از یک پایگاه داده به پایگاه داده دیگر منتقل شوند، این عمل لوله‌گذاری داده نامیده می‌شود.





پایگاه‌های داده عملیاتی را میتوان بیشتر به سه دسته کلی تقسیم کرد:

پایگاه‌های داده SQL سنتی

محدودیت اصلی مقیاس‌پذیری آنهاست: آنها یا مقیاس‌بندی نمی‌شوند یا فقط به صورت لگاریتمی مقیاس‌بندی می‌شوند، به این معنی که هزینه آنها به طور تصاعدی با مقیاس حجم کاری که قرار است پردازش شود، افزایش می‌یابد.

پایگاه‌های داده NoSQL

میتوانند تضمین‌های سازگاری **ACID** را فراهم کنند. از سوی دیگر، اکثر آنها می‌توانند مقیاس شوند، اگرچه همه انواع آن این توانایی را ندارند. برخی از آنها می‌توانند مقیاس شوند اما نه به صورت خطی یا نه با تعداد گره‌های زیاد.





NoSQL

ذخیره سازی داده‌های کلید-مقدار

ذخیره سازی داده‌های کلید-مقدار بدون طرحواره هستند
به همین دلیل آنها قابلیت های بسیار کمی برای پرس وجو ارائه می دهند

پایگاه‌های داده سند‌گرا

از داده‌های نیمه ساختاریافته نوشته شده به زبان محبوبی مانند JSON یا XML پشتیبانی می کنند. قابلیت اصلی آنها این است که میتوانند داده‌ها را در یکی از این زبانها به طور کارآمد ذخیره کنند و پرس‌وجوهایی را برای این داده‌ها به روشی مؤثر انجام دهند





NoSQL

پایگاه داده‌های نموداری

پایگاه‌های داده گراف در ذخیره‌سازی و پرس‌وجو داده‌های مدل‌سازی شده به صورت گراف تخصص دارند. داده‌های نموداری که در قالب رابطه‌ای نمایش داده می‌شوند برای پرس‌وجو بسیار گران می‌شوند.

انبارهای داده با ستون گسترده

نسبت به ذخیره‌سازی داده‌های کلیدی، قابلیت‌های بیشتری را ارائه می‌کنند. آنها معمولاً پارتیشن بندی محدود را انجام می‌دهند، بنابراین از پرس‌وجوهای محدود پشتیبانی می‌کنند.





NewSQL

NewSQL یک کلاس از سیستم های مدیریت پایگاه داده رابطه ای است که به دنبال ارائه مقیاس پذیری سیستم های NoSQL برای بارهای کاری پردازش تراکنش آنلاین (OLTP) و در عین حال حفظ ضمانت های ACID یک سیستم پایگاه داده سنتی است.

بسیاری از سیستم های سازمانی که داده های با مشخصات بالا را مدیریت می کنند (به عنوان مثال، سیستم های پردازش مالی و سفارش) برای پایگاه های اطلاعاتی رابطه ای معمولی بسیار بزرگ هستند، اما دارای الزامات تراکنش و سازگاری هستند که برای سیستم های NoSQL عملی نیستند تنها گزینه هایی که قبلاً برای این سازمان ها در دسترس بود، خرید رایانه های قوی تر یا توسعه میان افزار سفارشی بود که درخواست ها را روی DBMS معمولی توزیع می کرد. هر دو رویکرد دارای هزینه های زیرساختی بالا و/یا هزینه های توسعه هستند. سیستم های NewSQL سعی می کنند تضادها را با هم آشتی دهند.



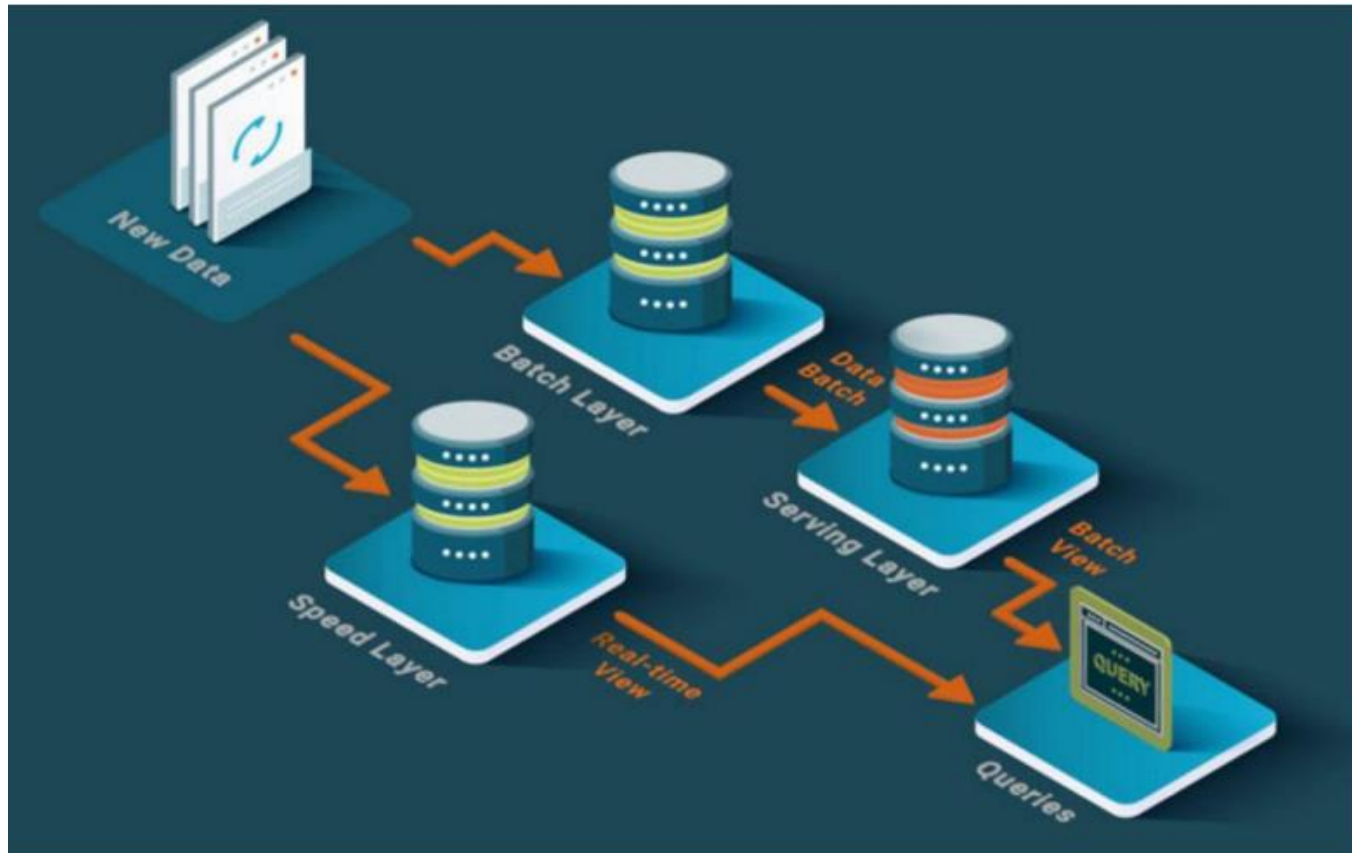


معماری لاندا

لایه دسته‌ای

لایه سرعت

لایه سرویس





فراتر از معماری لاندرا

- معماری لاندرا به سادگی با پایگاه داده **LeanXcale** جایگزین می‌شود که تمام قابلیت‌های معماری لاندرا را فراهم می‌کند و هیچ گونه پیچیدگی و هزینه توسعه و نگهداری ندارد.
- ناکارآمدی در جذب: همانطور که پایگاه داده رشد می‌کند هر درج نیاز به خواندن یک گره برگ از درخت + **B** دارد که ابتدا باید یک گره را از حافظه پنهان خارج کرده و روی دیسک بنویسد. این بدان معناست که هر نوشتن حداقل به دو **IO** نیاز دارد.
- **LeanXcale** این مشکل را با استفاده از نوع جدیدی از درختان **LSM** به لطف ترکیب قابلیت‌های **SQL** و **NoSQL** با ارائه کارایی ذخیره‌سازی داده‌های کلید-مقدار در مصرف حل می‌کند.





- سهولت در پرس و جو : معماری لاندا نیاز به توسعه برنامه نویسی هر پرس و جو با سه پایه کد مختلف برای هر یک از سه لایه دارد. در **LeanXcale**، پرس و جوها به سادگی در **SQL** نوشته می شوند. پرس و جوهای **SQL** بر خلاف پرس و جوهای برنامه های در معماری لاندا که نیاز به بهینه سازی دستی در سه پایه کد مختلف برای هر یک از لایه ها دارند، به طور خودکار بهینه می شوند.





تقسیم داده های تاریخی کنونی





برای حل این مشکل، در LeanXcale، یک الگوی جدید، به نام ذخیره‌سازی داده‌ها در زمان واقعی استفاده می‌شود. این الگو با یک نوآوری که در LeanXcale معرفی می‌شود، حل خواهد شد، یعنی توانایی تقسیم پرس‌وجوهای تحلیلی بر روی LeanXcale و یک انبار داده خارجی. اساساً، قطعات قدیمی‌تر داده را به صورت دوره‌ای در انبار داده کپی می‌کند.





Data marts

مشکل رایج:

اشباع انبار داده

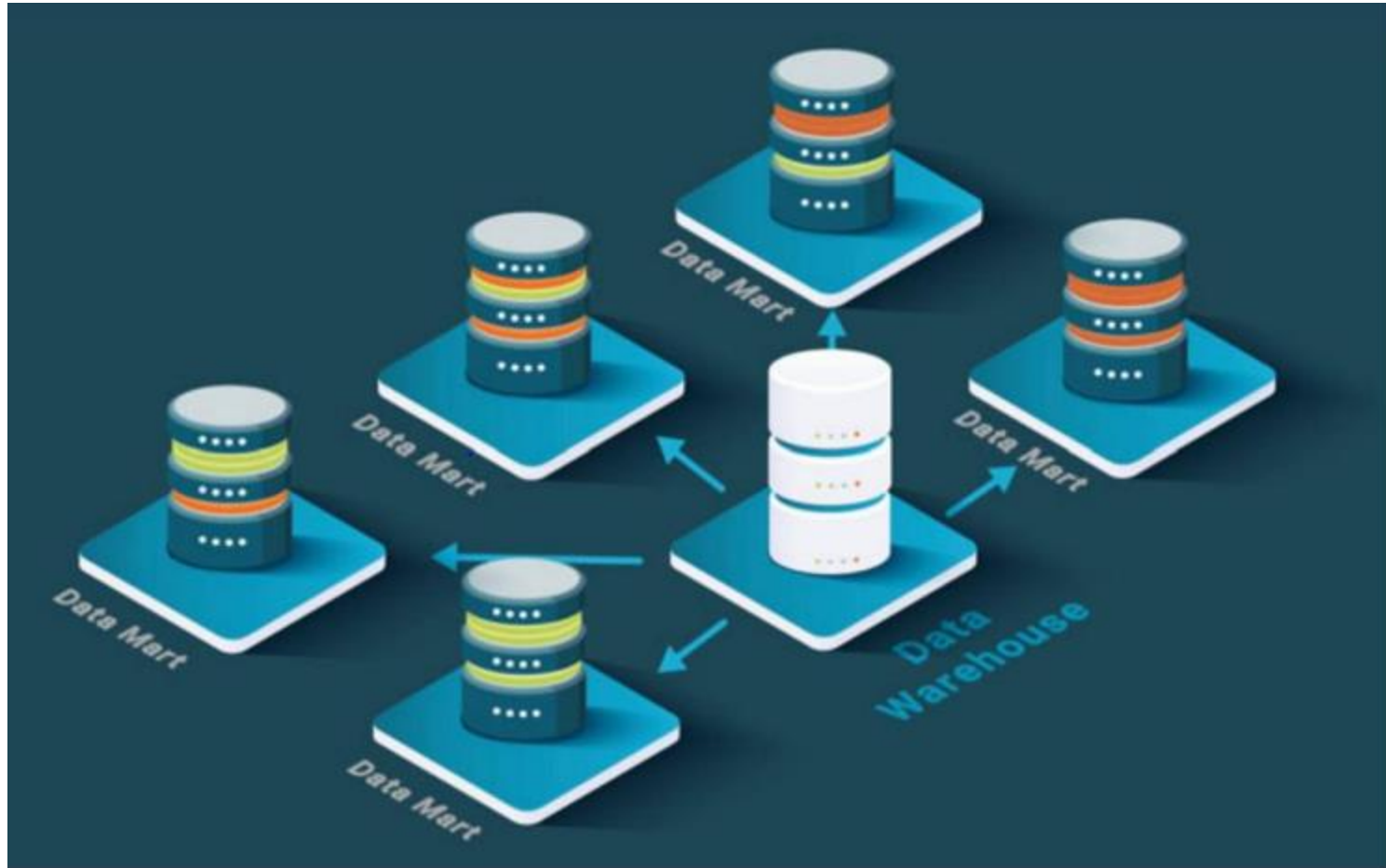
راه حل اول:

ایجاد نماهای کوچک از داده‌های موجود در انبار داده و ذخیره آنها را در پایگاه‌های داده مستقل





Data marts





Snapshotting

با Snapshot گرفتن شما می‌توانید وضعیت پایگاه داده را در لحظه فعلی ثبت و داخل فایل نگهداری کنید تا بتوان در زمان مورد نیاز از امکاناتش استفاده نمایید. شاید بتوان گفت به نوعی بکاپ‌گیری است و می‌توان این فایل را **Restore** کرد.

یک عکس فوری از صورت دوره‌های داده‌ها با تازگی یکسان یا بالاتر از داده‌های **mart** ذخیره می‌شود. مزیت آن این است که کپی به جای اینکه از انبار داده بیاید، مستقیماً از پایگاه داده عملیاتی می‌آید، در نتیجه داده‌های جدیدتری به دست می‌آید.

به لطف کارایی و سرعت دریافت داده‌های **LeanXcale**، می‌توان عکسبرداری فوری را به جای هفتگی روزانه انجام داد زیرا فرآیندهای بارگذاری که روزها طول میکشد به چند دقیقه کاهش می‌یابد.





جمع‌بندی

LeanXcale، منجر به معماری‌های ساده‌تر با تازگی داده‌ها یا حتی داده‌های بلادرنگ می‌شود. این معماری‌ها به دلیل سادگی، فرآیند توسعه را سرعت می‌بخشند و زمان بین جمع‌آوری نیازمندی‌ها و شروع ایده تا نرم‌افزارهای آماده تولید را کاهش می‌دهند نگهداری آنها مقرون به صرفه تر است زیرا متخصصان کمتر، سرورهای مختلف و مجوزهای پایگاه داده مورد نیاز هستند. علاوه بر این، قابلیت‌های پردازشی آنها راه چابک‌تری برای سرعت بخشیدن به فرآیندهای تجاری و در عین حال کاهش ریسک عملیاتی فراهم می‌کند. در نهایت، این معماری‌ها واقعاً می‌توانند از ایجاد جریان‌های درآمدی جدید با ایجاد راه‌های جدید برای برآوردن نیازهای مشتری، با استفاده از نفت امروزی یعنی داده پشتیبانی کنند.





با سپاس از شما