# پروژه پایگاه داده پیشرفته: سامانه پایش شبکه با Apache Cassandra

# مقدمه و اهداف پروژه

هدف این پروژه، طراحی و پیادهسازی یک "سامانه پایش شبکه (Network Monitoring System) " مقیاس پذیر با استفاده از پایگاه داده Apache Cassandra بود .سیستم میبایست قادر به ذخیرهسازی، بازیابی و تحلیل دادههای سری زمانی حجیم باشد که به صورت پیوسته (هر ۵ دقیقه) از صدها منبع شبکه (مانند روترها و فایروالها) جمع آوری می شوند.

## چالش اصلی

پاسخ گویی به دو کوئری کلیدی زیر با کارایی بالا (High Performance) بود:

- **Q1:** نمایش توالی مقادیر یک یا چند متریک خاص برای یک منبع، در بازه زمانی کوتاه (کمتر از یک هفته).
  - **Q2:** استخراج لیست تمام متریکهای مربوط به یک منبع خاص.

#### ابزارهای مورد استفاده

ابزار	کاربرد
Apache Cassandra	پایگاه داده NoSQL مقیاس پذیر برای ذخیرهسازی دادههای سری زمانی.
CQL (Cassandra Query Language)	زبان تعریف جداول و اجرای کوئریها.
Python	زبان اسکریپتنویسی برای فرآیند ETL (تولید، تبدیل و بارگذاری دادهها).
Docker	برای اجرای سریع و ایزوله دیتابیس Cassandra در محیط توسعه.

# مستندات طراحی و منطق انتخاب مدل داده

کلید موفقیت در Cassandra ، استفاده از طراحی مبتنی بر کوئری (Query-Driven Design) است. بر خلاف مدل سازی SQL ، ما ابتدا کوئری های نهایی Q1 و Q2 را شناسایی کرده و سپس جداولی را طراحی می کنیم که مستقیماً به آن کوئری ها پاسخ دهند.

هدف اصلی این است که هر کوئری، پاسخ خود را با خواندن حداکثر یک پارتیشن دریافت کند .برای رسیدن به این هدف، از افزونگی داده (Data Redundancy) استفاده کرده و دو جدول مجزا طراحی نمودیم.

### جدول1: metrics\_by\_resource\_week براى پاسخ به Q1

این جدول برای پاسخ به کوئری تحلیل سری زمانی (Q1) طراحی شده است.

طراحی: کوئری Q1 بازه کمتر از یک هفته را میخواهد ، اما فرضیات پروژه می گوید داده ها برای یک سال نگهداری خواهند شد .اگر ما تمام داده های یک ساله یک متریک را در یک پارتیشن ذخیره می کردیم، آن پارتیشن بسیار حجیم شده و خواندن آن ناکارآمد بود.

راهحل :(Time Bucketing) ما دادهها را در "سطلهای زمانی (Time Buckets) "هفتگی دستهبندی کردیم.

```
PS C:\Users\samin\OneDrive\Desktop> docker run -d --name cassandra-project -p 9042:9042 cassandra:latest
 Unable to find image 'cassandra:latest' locally latest: Pulling from library/cassandra
Digest: sha256:8b55dd41d5d1220e11eb8cf80f26ab655c21f7cf271ca4a7577c1da7d9221624
  Status: Downloaded newer image for cassandra:latest
00510531935ff671e64609298c18158f89ea045aa4f769213724ed54035247a0
  PS C:\Users\samin\OneDrive\Desktop> docker ps
  CONTAINER ID IMAGE
                                                                                                                                            COMMAND
                                                                                                                                                                                                                                                   CREATED
                                                                                                                                                                                                                                                                                                                   STATUS
                                                                                                        NAMES
                                                                                                                                          "docker-entrypoint.s..." 3 minutes ago
  00510531935f cassandra:latest
                                                                                                                                                                                                                                                                                                               Up 3 minutes 7000-7001/tcp, 7199/tcp, 9160/tcp,
00510531935+ cassandra:tatest "docker-entrypoint.s..." 3 minutes ago op 4 
 ... WITH REPLICATION cqlsh> USE network_monitoring; cqlsh:network_monitoring> CREATE TABLE metrics_by_resource_week (
... resource TEXT,
                                                                                                                            year INT,
                                                                                                                            week_of_year INT,
metric_name TEXT,
                                                                                                                            collected_at TIMESTAMP,
                                                                                                           value DOUBLE,
   PRIMARY KEY ((resource, year, week_of_year), metric_name, collected_at)
) WITH CLUSTERING ORDER BY (metric_name ASC, collected_at DESC);
```

#### منطق انتخاب كليدها:

• Partition Key: (resource, year, week\_of\_year)

- این کلید تضمین می کند که تمام متریکها (کمتر از ۱۰) تا ، برای یک منبع خاص، در یک هفته خاص، همگی در یک پارتیشن واحد ذخیره می شوند.
  - بنابراین، وقتی کوئری Q1 مثلاً دادههای هفته ۵۲ سال ۲۰۲۲ برای (router-1) را اجرا
     میکنیم، Cassandra مستقیماً به همان یک پارتیشن مراجعه میکند.
- Clustering Key: (metric\_name ASC, collected\_at DESC)
  - ∞ metric\_name: به ما اجازه می دهد تا به سرعت یک یا چند متریک را فیلتر کنیم.
- ocollected\_at DESC: دادههای درون پارتیشن را بر اساس زمان (جدیدترین اول) مرتب می کند که برای تحلیل سری زمانی ایده آل است.

## جدول ۲: metrics\_by\_resource\_list برای پاسخ به Q2

این جدول صرفاً برای پاسخ سریع به کوئری Q2 (لیست متریکهای یک منبع) طراحی شده است.

#### منطق انتخاب كليدها:

- Partition Key: (resource)
  - ۰ تمام متریکهای مربوط به یک منبع در یک پارتیشن ذخیره میشوند.
- Clustering Key: (metric\_name)
  - این کلید، لیست متریکها را برای آن منبع مرتب و (مهمتر) منحصربهفرد می کند.

## اسکرییت Python برای ورود داده(ETL)

### یادداشتی در مورد بهینهسازی فرآیندETL:

در مستندات پروژه، از ما خواسته شده بود که ابتدا دادهها را در یک فایل CSV ایجاد کرده و سپس آن فایل را با پایتون بخوانیم و در Cassandra درج کنیم.

اولین قدم، تحلیل این نیازمندی بود. بر اساس فرضیات خود پروژه (۲۰۰ منبع، هر ۵ دقیقه، برای ۱ سال)، محاسبات نشان داد که ایجاد این فایل CSV منجر به تولید فایلی با بیش از ۱۰۰ میلیون ردیف داده می شد.

ایجاد یک فایل متنی با این حجم (چندین گیگابایت) و سپس خواندن مجدد کل آن از روی دیسک، یک تنگنای (bottleneck) شدید ۱/۵ ایجاد کرده و فرآیندی بسیار ناکارآمد و زمانبر است.

در نتیجه، برای بهینهسازی این فرآیند، تصمیم گرفتیم که به جای انجام دو فرآیند مجزای نوشتن روی دیسک و خواندن از دیسک، ما اسکریپت generate\_and\_ingest.py را توسعه دادیم.

این اسکریپت هر سه وظیفه ETL (تولید، تبدیل، و بارگذاری) را به صورت یکپارچه و در حافظه انجام میدهد:

- fenerate (تولید): دادهها بر اساس فرضیات پروژه در لحظه تولید میشوند.
- ۲. Transform (تبدیل): ستونهای year و week\_of\_yearدر همان لحظه استخراج میشوند.
  - ۳. Load (بارگذاری): دادهها مستقیماً و به صورت بهینه (با بچهای کوچک) در Cassandra درج میشوند.

این روش، منطق خواسته شده در پروژه ETL با پایتون را در مقیاس بالا به شکلی کارآمد شبیهسازی میکند. فایل کامل این اسکریپت با نام generate\_and\_ingest.py به صورت جداگانه پیوست شده است.

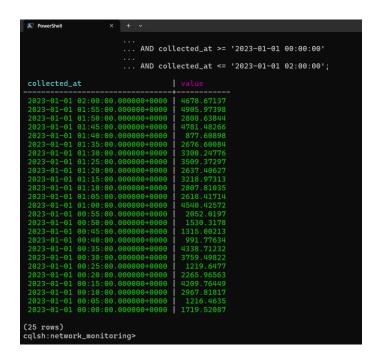
```
3790000 records inserted.
                                               (Last timestamp:
                                               (Last timestamp: 2023-01-07 13:50:00)
(Last timestamp: 2023-01-07 14:15:00)
(Last timestamp: 2023-01-07 14:40:00)
(Last timestamp: 2023-01-07 15:05:00)
(Last timestamp: 2023-01-07 15:30:00)
(Last timestamp: 2023-01-07 15:55:00)
       3800000 records inserted.
       3810000 records inserted.
       3820000 records inserted.
       3830000 records inserted.
       3840000 records inserted.
                                               (Last timestamp: 2023-01-07 16:20:00)
       3850000 records inserted.
                                               (Last timestamp: 2023-01-07 16:20:00)
(Last timestamp: 2023-01-07 16:45:00)
(Last timestamp: 2023-01-07 17:10:00)
(Last timestamp: 2023-01-07 17:35:00)
(Last timestamp: 2023-01-07 18:00:00)
       3860000 records inserted.
       3870000 records
                               inserted.
       3880000 records inserted.
       3890000 records inserted.
       3900000 records
                                               (Last timestamp: 2023-01-07 18:25:00)
                               inserted.
       3910000 records
                               inserted
                                                (Last timestamp:
                                                                         2023-01-07 18:50:00
                                               (Last timestamp: 2023-01-07 19:15:00)
(Last timestamp: 2023-01-07 19:40:00)
(Last timestamp: 2023-01-07 20:05:00)
       3920000 records
                               inserted.
       3930000 records inserted.
       3940000 records inserted.
                                               (Last timestamp: 2023-01-07 20:30:00)
       3950000 records inserted.
       3960000 records inserted.
                                                (Last timestamp:
                                                                         2023-01-07 20:55:00
                                               (Last timestamp: 2023-01-07 21:20:00)
(Last timestamp: 2023-01-07 21:45:00)
(Last timestamp: 2023-01-07 22:10:00)
       3970000 records
                               inserted.
       3980000 records inserted.
       3990000 records inserted.
                                               (Last timestamp: 2023-01-07 22:35:00)
       4000000 records inserted.
     4010000 records inserted. (Last timestamp: 2023-01-07 23:00:00) 4020000 records inserted. (Last timestamp: 2023-01-07 23:25:00) 4030000 records inserted. (Last timestamp: 2023-01-07 23:50:00)
    Data Ingestion Complete
otal records inserted: 4032000
Total execution time: 6223.44 seconds
  Men's Singles

Match results
                                                                               Q Search
```

## اجرای کوئریها و تحلیل نتایج

پس از درج موفقیتآمیز بیش از ۴ میلیون رکورد داده، کوئریهای کلیدی پروژه برای اعتبارسنجی مدل اجرا شدند.

الف) پاسخ به :Q1 تحلیل سری زمانی (یک متریک)



نتیجه: کوئری با موفقیت اجرا شد و ۲۵ ردیف داده را برگرداند. این نشان میدهد که مدل

"Time Bucketing" به درستی کار می کند. (دلیل استفاده از 2022 year=2022 و week=52 برای تاریخ الله استفاده از ISO بیروی اسکریپت پایتون از استاندارد تقویم ISO است.

### ب)یاسخ به :Q1 تحلیل سری زمانی (چند متریک)

```
PowerShell
(25 rows)
cqlsh:network_monitoring> SELECT metric_name, collected_at, value
                              ... FROM metrics_by_resource_week
... WHERE resource = 'router-1'
                              ... WHERE resource = 'router-1'
... AND year = 2022
... AND week_of_year = 52
... AND metric_name IN ('cpu_usage', 'memory_usage')
... AND collected_at >= '2023-01-01 00:00:00'
                              ... AND collected_at <= '2023-01-01 02:00:00';
 metric_name | collected_at
     cpu_usage
     cpu_usage
     cpu_usage
     cpu_usage
                                                                    4781.48266
                      2023-01-01 01:40:00.000000+0000
                                                                    877.60898
2676.60084
     cpu_usage
                      2023-01-01 01:35:00.000000+0000
     cpu_usage
                      2023-01-01 01:30:00.000000+0000
                                                                    3300.24776
     cpu_usage
     cpu_usage
                      2023-01-01 01:20:00.000000+0000
     cpu usage
```

نتیجه : کوئری با موفقیت اجرا شد و ۵۰ ردیف داده (۲۵ ردیف برای هر متریک) را برگرداند. این ثابت می کند. که مدل ما قادر است "یک یا چند متریک" را به صورت کارآمد و تنها با خواندن یک پارتیشن واکشی کند.

#### پاسخ به :Q2 استخراج لیست متریکها

```
cqlsh:network_monitoring> SELECT metric_name
                                 ... FROM metrics_by_resource_list
... WHERE resource = 'router-1';
 metric_name
      cpu_usage
    latency_ms
 memory_usage
packets_in
packets_out
(5 rows)
cqlsh:network_monitoring> SELECT * FROM metrics_by_resource_list LIMIT 10;
                    | metric_name
 firewall-31
firewall-31
firewall-31
                      active_connections
                                    bytes_in
bytes_out
 firewall-31
firewall-31
firewall-31
switch-104
switch-104
                           cpu_usage
dropped_packets
                              cpu_usage
packets_total
                              port_1_errors
port_2_errors
   switch-104
   switch-104
switch-104
                                   uptime_sec
(10 rows)
cqlsh:network_monitoring> |
```

نتیجه : کوئری بلافاصله اجرا شد و لیست ۵ متریک مربوط به router-1 را برگرداند. این ثابت می کند که جدول دوم (جدول افزونه) به درستی به کوئری Q2 پاسخ می دهد.

## نتيجهگيري

این پروژه با موفقیت تمام اهداف خود را محقق کرد. با استفاده از اصول Query-Driven Design و Pucketing در Apache Cassandra ، یک مدل داده بسیار کارآمد طراحی شد که می تواند میلیون ها رکورد داده سری زمانی را مدیریت کند و در عین حال، کوئری های تحلیلی کلیدی را با کارایی بالا (خواندن از یک یارتیشن) یاسخ دهد.