# Big Data Assignment-4

**Hruday Vishal Kanna Anand**

**1006874517**

**Part-a**

**1.**

**Resources created in resource group**



**2.**

**The input data is copied into the blob storage**

**The pipeline created in azure data factory is executed manually-** the name of the pipeline in the first picture and the pipeline executed in the monitor view is different as they were done at different times and I had deleted all my resources in between.



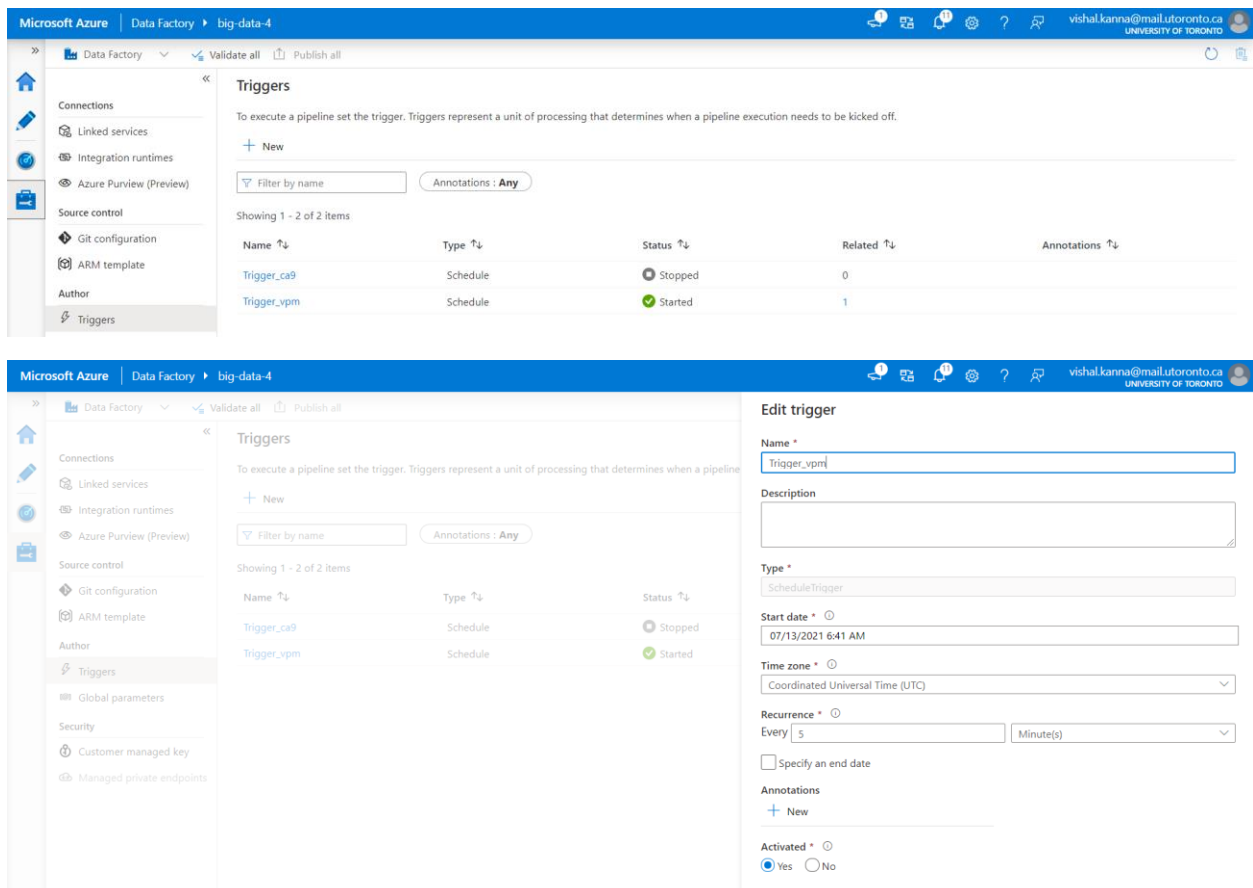**The data is copied into SQL DB using azure data factory pipeline**

**3.**

**The schedule trigger is created**





**The trigger is executed 5 times**



There are mainly 4 types of triggers in data factory-

1. Manual trigger
2. Schedule trigger
3. Tumbling window trigger

4. Event based trigger

Manual trigger is a trigger that activated by the user manually by either using the UI or a Powershell script.

Schedule trigger runs a pipeline using generic clock timing such as hours, minutes, weeks, etc. It can be a simple timing such as weekly or more complex such as every Monday at 2:00 P.M. There can also be many schedule triggers for a single pipeline.

Tumbling window trigger is similar to a schedule trigger as it also runs the pipeline at specific time intervals but with the inclusion of state. So, in case of failure it can automatically try again as it stores past states making this more reliable. There can only be one tumbling window trigger for a pipeline.

Event based trigger runs a pipeline in response to an event. There are 2 types of event triggers-

1. storage based- triggered when a change occurs in a storage account
2. custom event


**4.**

They would first need to create ADLS gen 2 accounts in both Canada central and west Europe regions. This can be done very similar to creating storage accounts. We will specify the region as either Canada central or west Europe and we will enable hierarchical namespace for both the accounts. This will make them ADLS gen 2 accounts.

We will then have to create data factory account to facilitate the data transfer from one ADLS to another. The data factory can communicate with the ADLS's with the help of a linked service that must be created to both the data lakes.

Now we must create a pipeline for transferring files (objects) from one ADLS to another ADLS. Let's say we transfer from Canada central to west Europe. this pipeline will be run with a storage event trigger; hence it will be triggered when a new file is added to the Canada central ADLS. This trigger will have parameters of the file that caused trigger with the help of @triggerbody().folderpath and @triggerbody().filename. Using these parameters, we run the first activity of our pipeline which is validation. This is done to check if west Europe ADLS has the folder and file path that triggered the pipeline. If the folder and file is found in west Europe ADLS then the pipeline proceeds to do nothing (wait activity) as we don't want replicate files to be created by multiple copies. If it is not found, then a copy data activity is done where data is copied from Canada central ADLS using the trigger properties above to the same folder and file as it was in Canada central to west Europe. We can see this method prevents unnecessary copies of files if the file already exists in the other ADLS.

As we want this to be bi directional we set up a similar pipeline but this time it will be from west Europe ADLS to Canada central ADLS following the same steps as above only switching the data flow path(west Europe to Canada central) and the validation will be done on Canada central and the storage event trigger will be based on west Europe ADLS.

**Part-b**

**1.**

**The output is attached in a csv file**



**2.**

**3.**

**The output is attached in a csv file**

| Query editor (preview) ···                                                     ✕

Feedback

Query 1 ✕   Query 2 ✕   Query 3 ✕   Query 4 ✕   Query 5 ✕   Query 6 ✕   Query 7 ✕   **Query 8 ✕**

▷ Run   ☐ Cancel query   ↓ Save query   ↓ Export data as  ∨   ▦ Show only Editor

```
1  select * from newgender.jobs where occupation='Bus drivers'
2  order by year;
```

**4.**

) | Query editor (preview) ···                                                     ✕

Feedback

Query 1 ✕   Query 2 ✕   **Query 3 ✕**   Query 4 ✕   Query 5 ✕   Query 6 ✕

▷ Run   ☐ Cancel query   ↓ Save query   ↓ Export data as  ∨   ▦ Show only Editor

```
1  select year ,SUM( CAST(workers_female as int)) as "female workers"
2  from newgender.jobs where major_category='Management, Business, and Financial' group by year
3  order by year;
```

**Results**   Messages

| year | female workers |
| --- | --- |
| 2013 | 7748347 |
| 2014 | 8061480 |
| 2015 | 8381812 |
| 2016 | 8617853 |

✓ Query succeeded | 0s

**5.**

| Query editor (preview)    ⋯                                                                                    ✕

Feedback

Query 1 ✕   Query 2 ✕   Query 3 ✕   Query 4 ✕   Query 5 ✕   Query 6 ✕   Query 7 ✕   Query 8 ✕   **Query 9 ✕**

▷ Run    ☐ Cancel query    ↓ Save query    ↓ Export data as  ⌄    ▦ Show only Editor

```
1    select SUM( CAST(total_earnings_male as int)) as "male total earnings in 2015"
2    from newgender.jobs where major_category='Service' and year='2015'
```

**Results**    Messages

🔍 Search to filter items...

**male total earnings in 2015**

2502426

**6.**

| Query editor (preview)    ⋯                                                                                    ✕

Feedback

Query 1 ✕   Query 2 ✕   Query 3 ✕   Query 4 ✕   Query 5 ✕   Query 6 ✕   Query 7 ✕   Query 8 ✕   **Query 9 ✕**

▷ Run    ☐ Cancel query    ↓ Save query    ↓ Export data as  ⌄    ▦ Show only Editor

```
1    select SUM( CAST(workers_female as int)) as "total female workers in 2015"
2    from newgender.jobs where minor_category='Management' and year='2015'
```

**Results**    Messages

🔍 Search to filter items...

**total female workers in 2015**

5166720

**7.**

) | Query editor (preview) ···

♡ Feedback

Query 1 ✕    Query 2 ✕    Query 3 ✕    Query 4 ✕    Query 5 ✕    Query 6 ✕    Query 7 ✕    Query 8 ✕    Query 9 ✕    Query 10 ✕

▷ Run    ☐ Cancel query    ↓ Save query    ↓ Export data as ⌄    ▦ Show only Editor

```
1    create view compare_3 as
2    select year, SUM( CAST(total_earnings_male as int)) as "total earning male",
3    SUM( CAST(total_earnings_female as int)) as "total earning female"
4    from newgender.jobs
5    group by year;
6
```

) | Query editor (preview) ···

⟲ Feedback

Query 1 ✕    Query 2 ✕    Query 3 ✕    Query 4 ✕    Query 5 ✕    Query 6 ✕    Query 7 ✕    Query 8 ✕    Query 9 ✕    Query 10 ✕

▷ Run    ☐ Cancel query    ↓ Save query    ↓ Export data as ⌄    ▦ Show only Editor

```
1    select * from compare_3 order by year;
```

**Results**    Messages

| year | total earning male | total earning female |
|------|--------------------|-----------------------|
| 2013 | 27050782 | 22054404 |
| 2014 | 27470450 | 22491208 |
| 2015 | 27754851 | 22768521 |
| 2016 | 28463638 | 23075602 |

**8.**

uery editor (preview)   ...                                          ✕

◡ Feedback

Query 1 ✕    **Query 2** ✕

▷ Run   ☐ Cancel query   ↓ Save query   ↓ Export data as ∨   ▦ Show only Editor

```sql
1   select sum(cast(total_earnings_female as int))
2   as "total earnings female in 2016" from newgender.jobs
3   where occupation like '%engineers%' and year='2016';
```

**Results**   Messages

🔍 Search to filter items...

total earnings female in 2016

1618757

---

**9.**

eview)   ...                                                         ✕

Feedback

**Query 1** ✕    Query 2 ✕

▷ Run   ☐ Cancel query   ↓ Save query   ↓ Export data as ∨   ▦ Show only Editor

```sql
1   create view compare_gender as
2   select year, sum(cast(full_time_male as float)*workers_male*0.01) as "total full time male",
3   sum(cast(part_time_male as float)*workers_male*0.01) as "total part time male",
4   sum(cast(full_time_female as float)*workers_female*0.01) as "total full time female",
5   sum(cast(part_time_female as float)*workers_female*0.01) as "total part time female"
6   from newgender.jobs
7   group by year;
```

erver/sql-db)

**ver/sql-db) | Query editor (preview)** ...

×

Query 1 ×      **Query 2** ×

▷ Run      ☐ Cancel query      ↓ Save query      ↓ Export data as ∨      ⊞ Show only Editor

```
1    select * from compare_gender order by year;
```

**Results**      Messages

⌕ Search to filter items...

| year | total full time male | total part time male | total full time female | total part time female |
|------|----------------------|----------------------|------------------------|------------------------|
| 2013 | 48827487.577 | 7360645.423 | 31568143.22 | 11091509.78 |
| 2014 | 50330271.951 | 7321815.049 | 32313480.43 | 11235684.57 |
| 2015 | 51720573 | 7321177 | 33414427.86 | 11257267.14 |
| 2016 | 52526792.592 | 7435299.408 | 34274127.486 | 11363858.514 |

⊘ Query succeeded | 0s