

Seminar Report on “Big Data PySpark”

1. Date of the Seminar/Workshop: 19-11-2020

2. Title of the Seminar/Workshop: “Big Data PySpark”

3. Name of the Speaker/Resource person with Affiliation with the institute/industry:

Mr. Aaron Rebello

Industry:- Capgemini Technology Services India Limited

Designation:-Associate Consultant

Role:-Big Data Engineer

4. Venue of the Seminar/Workshop: Platform used-Google meet

5. Duration of the Seminar: 3 hrs(9.30 am to 12.30 pm)

6. Conducted For: Students of final year computer Engineering

7. Objective of the Seminar/Workshop /Curriculum Gap identified/Other than that:

The objective of the Seminar was basically to ensure that the students should have a clear understanding of big data and how to start projects in this domain as students have to implement mini project based on this subject. And as a speaker is working in industry as a Big Data Engineer so the students should take benefit and a lot of insights could be inferred from his experience in this industry

8. Contents of the Seminar/Workshop:

The 3 hours period was very interactive and aroused key enthusiasm among students. The speaker started with

- Data timeline starting from the 1960s to the present time.
- OLTP vs. OLAP
- OLAP architecture
- ETL tools such as Xplenty, Stitch, Informatica, AWS and Oracle
- Different tools of Big Data like Hadoop, Kafka, Spark, Storm
- The working of MapReduce
- The spark ecosystem
- RDD- the resilient distributed dataset

and finally hands-on/demonstration of basic query structures using cloudera and pyspark

9. Description of the Entire Event

The Computer Department of St. Francis Institute of Technology organized a seminar on “Big Data PySpark” on Thursday 19th November-2020 between 9.30 am to 12 noon , online on Google meet.

There were approximately 80 student participants and 4 faculty participants who attended the seminar. The program was started with welcome speech by Ms. Anuradha Srinivasaraghavan where she introduced the Speaker. Later, Mr Aaron Rebello took over the session.

The speaker started off by introducing himself and his experience in the domain followed by the projects taken up in the organization he works in.

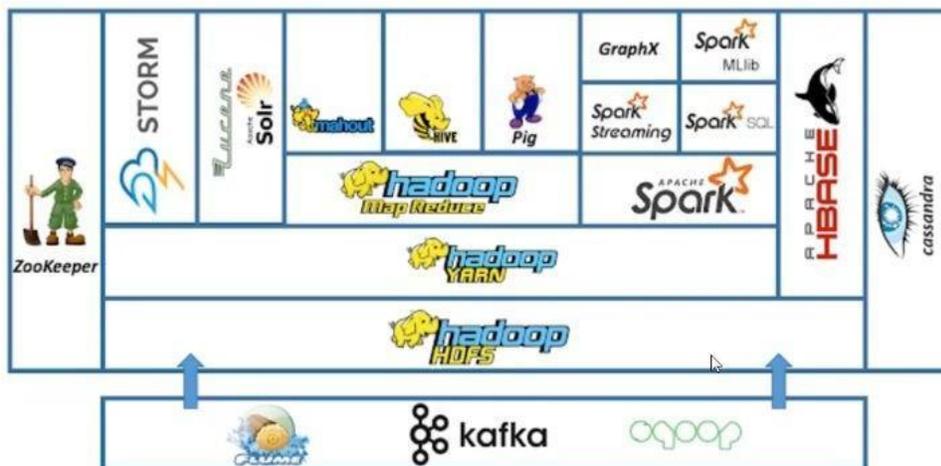
Then he started off with the main topic of the seminar by first talking about the data timeline starting from the 1960s to the present time. He spoke about how in 1960 database was introduced and then in 1970 OLTP came into the picture, followed by OLAP in 1975 and big data in 2005. Then a brief introduction was given about the OLAP architecture and juxtaposed OLAP with OLTP so as to understand the importance of OLAP. Then ETL tools were introduced to us like Xplenty, Stitch, Informatica, AWS and Oracle. He talked about the data distribution pyramid and how there is a very less amount of structured data and a large amount of unstructured data.

Furthermore, he told us about the emerging of big data due to the three papers from Google:

- 1) Google file system (2003)
- 2) MapReduce (2004)
- 3) BigTable (2006)

Then the definition of big data was introduced to us along with its different tools as shown below:

What is Big Data Today



Moving further, he explained the working of MapReduce in detail. MapReduce works in 2 phases. In the first phase, mapping is done where the jobs are split and data is mapped. In the second phase, the reduce task shuffles and consolidates the data. But this process is really slow because data replication takes place. Therefore, Spark came into the picture since it is 10-100 times faster. The spark ecosystem was explained further. It has service APIs like Spark Streaming, Spark SQL and GraphX, computing engine - Spark core, resource manager - yarn and Mesos and distributed file system = HDFS and HBase. Then he told us about RDD which is the resilient distributed dataset and how spark relies on this as it can be stored in memory and also perform various operations. Then he went on to show the code and work in the Java and Python languages.

At the end of the session, the forum was open for doubts where he solved the queries posed by the students. He got us acquainted with the set up of the system so as to be able to work on big data projects. He also told us about the resources available on the internet like the official documentation which can be referred for practical implementation. Due to the session, the students had a clearer understanding of big data and how to start projects in this domain. It was very beneficial and a lot of insights could be inferred from his experience in this industry.

Session Photographs

The screenshot shows a Google Meet session in progress. The browser address bar indicates the meeting URL is `meet.google.com/vmf-kjic-voi`. The meeting interface shows a recording status (REC), the presenter's name (Aaron Rebello), and a list of participants including Vipul Naik, Shelton Jade Pinto, and G. Anuradha.

The main presentation slide is titled "Architecture" and illustrates the flow of data between two environments:

- Operational Environment:** Contains OLTP (On-Line Transaction Processing) at the base, leading to Business / Enterprise Strategy (1), Business Process (2), and Orders (3). Orders are linked to Customer and Products.
- Informational Environment:** Contains OLAP (On-Line Analytical Processing) at the base, leading to Data Warehouse (5) and Data Mart (5).
- Data Flow:** ETL Processes (4) move data from the Operational Environment to the Informational Environment. Data Mining, Analytics, and Decision Making (6) flow from the Informational Environment back to the Operational Environment.

The bottom of the screen shows a taskbar with various application icons and a system tray displaying the time as 9:57 AM on 19/11/2020.

Document1 - Microsoft Word

Meet - vmf-kjic-voi

meet.google.com/vmf-kjic-voi

REC Aaron Rebello is presenting AYUSH NAVGIRI_17... and 66 more

PySparkSession0309 - PowerPoint

OLTP vs. OLAP

ONLINE TRANSACTION PROCESSING	ONLINE ANALYTICAL PROCESSING
Handles recent operational data	Handles all historical data
Size is smaller, typically ranging from 100 Mb to 10 GB	Size is larger, typically ranging from 1 Tb to 100 Pb
Goal is to perform day-to-day operations	Goal is to make decisions from large data sources
Uses simple queries	Uses complex queries
Faster processing speeds	Slower processing speeds
Requires read/write operations	Requires only read operations

Class List (2020-11....csv) Class List (2020-1....html) Show all

9:58 AM 19/11/2020

Document1 - Microsoft Word

Meet - vmf-kjic-voi

meet.google.com/vmf-kjic-voi

REC Aaron Rebello is presenting JENNY D'CRUZ_17... and 66 more

Cluster Computing

Computer Nodes

10100 Mbps Ethernet Switch

Server

Internet

Networked Disk Storage

becmpna ^

Turn on captions

Raise hand

Aaron Rebello is presenting

Class List (2020-11....csv)

Class List (2020-1....html)

Show all

Page: 6

10:05 AM 19/11/2020

Microsoft Word - aaronseminar

Meet - vmf-kjic-voi

meet.google.com/vmf-kjic-voi

REC Aaron Rebello is presenting JESDIN RAPHAEL_17... and 70 more

Spark

What We'd Like

The diagram illustrates two data processing workflows. The top workflow shows a sequential process: 'Input' (cylinder) feeds into 'iter. 1' (blue box), which feeds into a server rack icon, which then feeds into 'iter. 2' (blue box), followed by another server rack icon, and so on. The bottom workflow shows a parallel process: 'Input' (cylinder) feeds into 'one-time processing' (blue box), which feeds into 'Distributed memory' (server rack icon). From 'Distributed memory', three arrows point to 'query 1', 'query 2', and 'query 3' (blue boxes), each followed by a green box representing output. An orange banner at the bottom states '10-100x faster than network and disk'.

10-100x faster than network and disk

Class List (2020-11....csv) Class List (2020-1....html) Show all

Page: 1

10:40 AM 19/11/2020

Meet - vmf-kjic-voi

meet.google.com/vmf-kjic-voi

REC Aaron Rebello is presenting

Spark Architecture

The diagram illustrates the Spark Architecture. On the left, a box labeled 'Driver program' contains 'SparkContext'. In the center, a box labeled 'Cluster Manager' is connected to the 'SparkContext' by a bidirectional arrow. On the right, there are two 'Worker Node' boxes. Each 'Worker Node' contains an 'Executor' and a 'Cache' at the top, and two 'Task' boxes at the bottom. Arrows point from the 'Cluster Manager' to each 'Worker Node'. Additionally, curved arrows point from each 'Worker Node' back to the 'SparkContext' in the 'Driver program'.

Class List (2020-11....csv) Class List (2020-1....html) Show all

10:41 AM 19/11/2020

Meet - vmf-kjic-voi

meet.google.com/vmf-kjic-voi

REC Aaron Rebello is presenting SURAJ MAURYA_18... and 65 more

SPARK ECOSYSTEM

Service APIs: Spark Streaming, GraphX, MLLib, SparkSQL, Pig, Hive, Search

Computing Engine: Spark Core, Spark Core Or Hadoop MapReduce

Resource Manager: Yarn, Mesos

Distribute File system: HDFS, HBase

Participants: You, Aaron Rebello, G. ANURADHA, KEVIN DSOUZA_172059

Class List (2020-11....csv) Class List (2020-1....html)

Page: 1

10:48 AM 19/11/2020

Microsoft Word window: aaronseminar - Microsoft Word

Browser window: Meet - vmf-kjic-voi

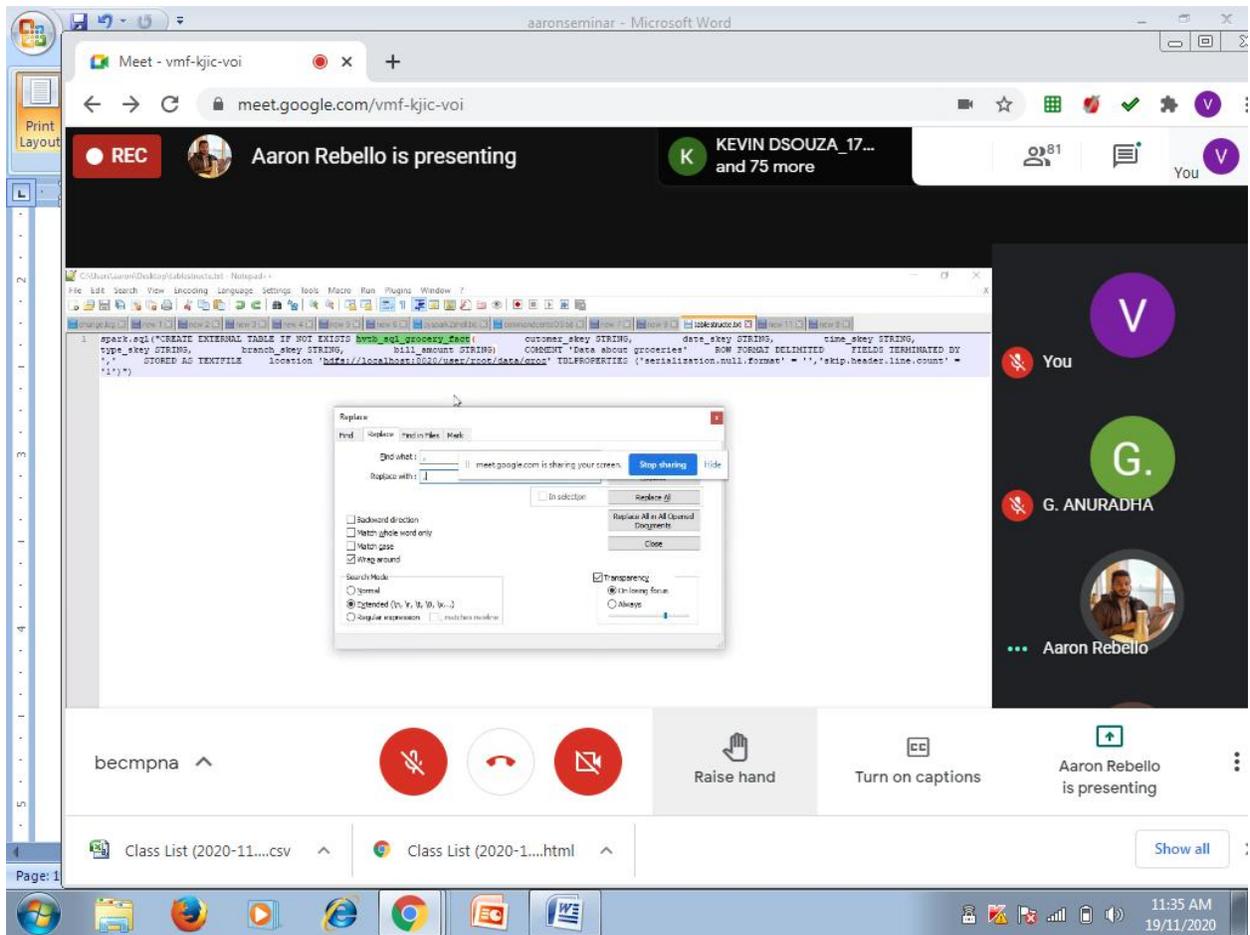
URL: meet.google.com/vmf-kjic-voi

Meeting controls: REC, Aaron Rebello is presenting, SHIVANI BISHT_17... and 75 more

Terminal window (root@quickstart:usr/src):

```
[root@quickstart src]#
[root@quickstart src]#
[root@quickstart src]# java -version
java version "1.8.0_172"
Java(TM) SE Runtime Environment (build 1.8.0_172-b11)
Java HotSpot(TM) 64-Bit Server VM (build 25.172-b11, mixed mode)
[root@quickstart src]# python -version
Python 2.7.17
[root@quickstart src]# python -V
Python 2.7.17
[root@quickstart src]# java -version
java version "1.8.0_172"
Java(TM) SE Runtime Environment (build 1.8.0_172-b11)
Java HotSpot(TM) 64-Bit Server VM (build 25.172-b11, mixed mode)
[root@quickstart src]# pyspark
WARNING: User-defined SPARK_HOME (/opt/cloudera/parcels/SPARK2-2.4.0-cloudera2-1.cdh5.13.3.p0.1041012/lib/spark2/) overrides detected (/opt/cloudera/parcels/SPARK2/lib/s
park2/).
WARNING: Running pyspark from user defined location.
Python 2.7.17 [(default, Nov 29 2019, 13:22:37)
[GCC 4.4.7 20120313 (Red Hat 4.4.7-23)] on linux2
Type "help", "copyright", "credits" or "license()" for more information.
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
█
```

Taskbar: Class List (2020-11....csv), Class List (2020-1....html), 11:31 AM 19/11/2020



Ms. Varsha N and Ms. Snehal K.
Seminar Incharge

Dr. Kavita Sonawane
HOD, CMPN