

1. Explain what is in-Memory computation in details?

**IN-MEMORY:** IN-MEMORY is nothing but storing the data in RAM

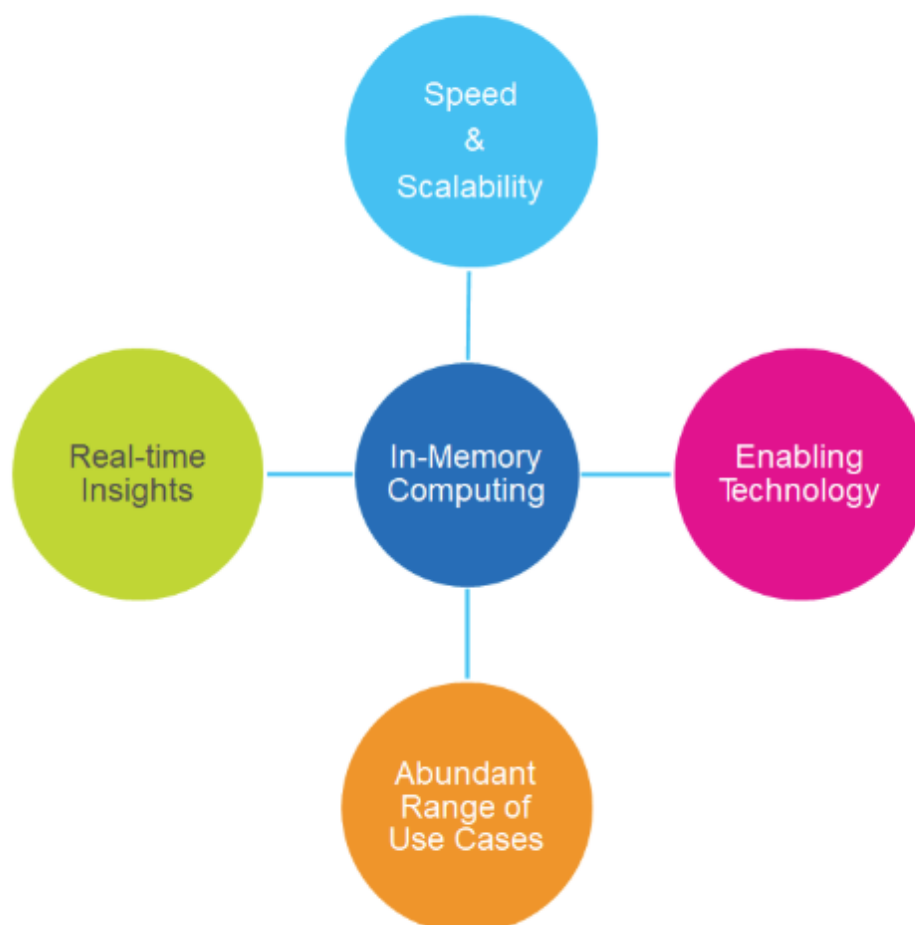
Where apache spark is nothing but open source cluster computing frame work

In-memory computation works by eliminating all slow data accesses and relying exclusively on data stored in RAM.

In-memory computation is often done in the technology known as **in-memory data grids** (IMDG).

In-Memory Computing provides super-fast performance (thousands of times faster) and scale of never-ending quantities of data, and simplifies access to increasing numbers of data sources

By storing data in RAM and processing it in parallel, it supplies real-time insights that enable businesses to deliver immediate actions and responses.

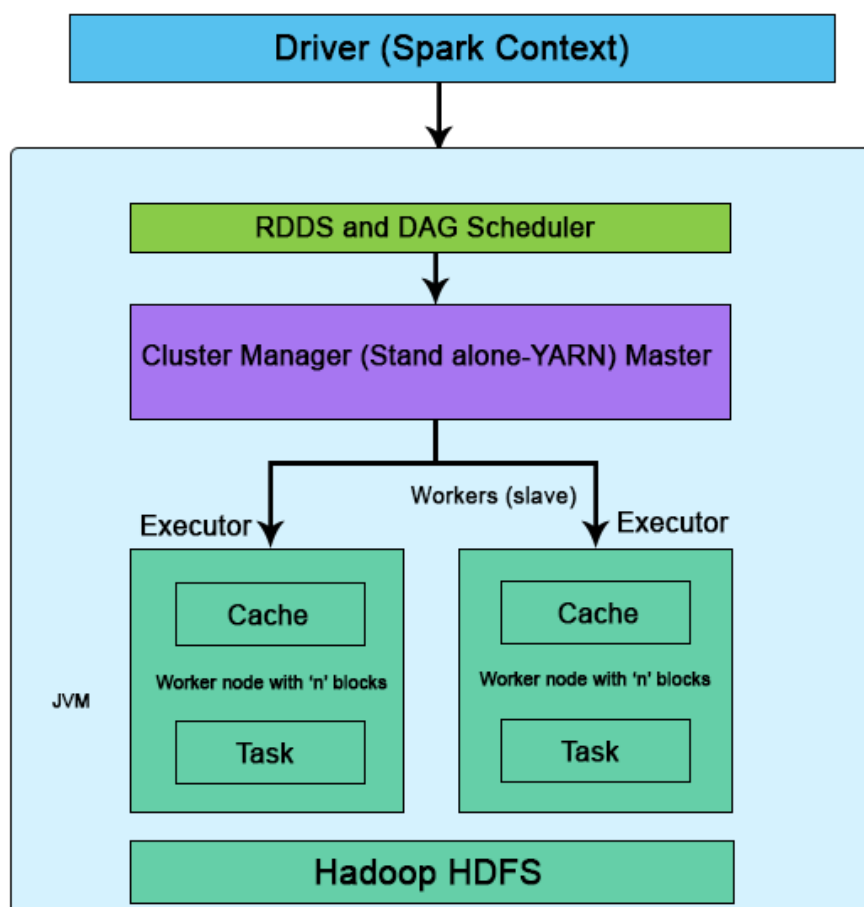


2. Explain advantages of Spark framework ?

### Benefits of Apache Spark:

1. Speed
2. Ease of Use
3. Advanced Analytics
4. Dynamic in Nature
5. Multilingual
6. Apache Spark is powerful
7. Increased access to Big data
8. Demand for Spark Developers
9. Open-source community

3. Explain components of Spark with block diagram ?



4. Explain benefits of in-Memory computation ?

- Investment banking
- Insurance claim processing & modeling
- Real-time ad platforms
- Real-time sentiment analysis
- Merchant platform for online games
- Hyper-local advertising
- Geospatial/GIS processing
- Medical imaging processing
- Natural language processing & cognitive computing
- Real-time machine learning
- Complex event processing of streaming sensor data

5. Explain major difference between Hadoop & Spark ?

Hadoop	Spark
Hadoop is an open source framework which uses a MapReduce algorithm	Spark is lightning fast cluster computing technology, which extends the MapReduce model to efficiently use with more type of computations
Hadoop's MapReduce model reads and writes from a disk, thus slow down the processing speed	Spark reduces the number of read/write cycles to disk and store intermediate data in-memory, hence faster-processing speed.
Hadoop is designed to handle batch processing efficiently	Spark is designed to handle real-time data efficiently.
With Hadoop MapReduce, a developer can only process data in batch mode only	Spark can process real-time data, from real time events like twitter, facebook
Hadoop is a high latency computing framework, which does not have an interactive mode	Spark is a low latency computing and can process data interactively.
Hadoop is a cheaper option available while comparing it in terms of cost	Spark requires a lot of RAM to run in-memory, thus increasing the cluster and hence cost.

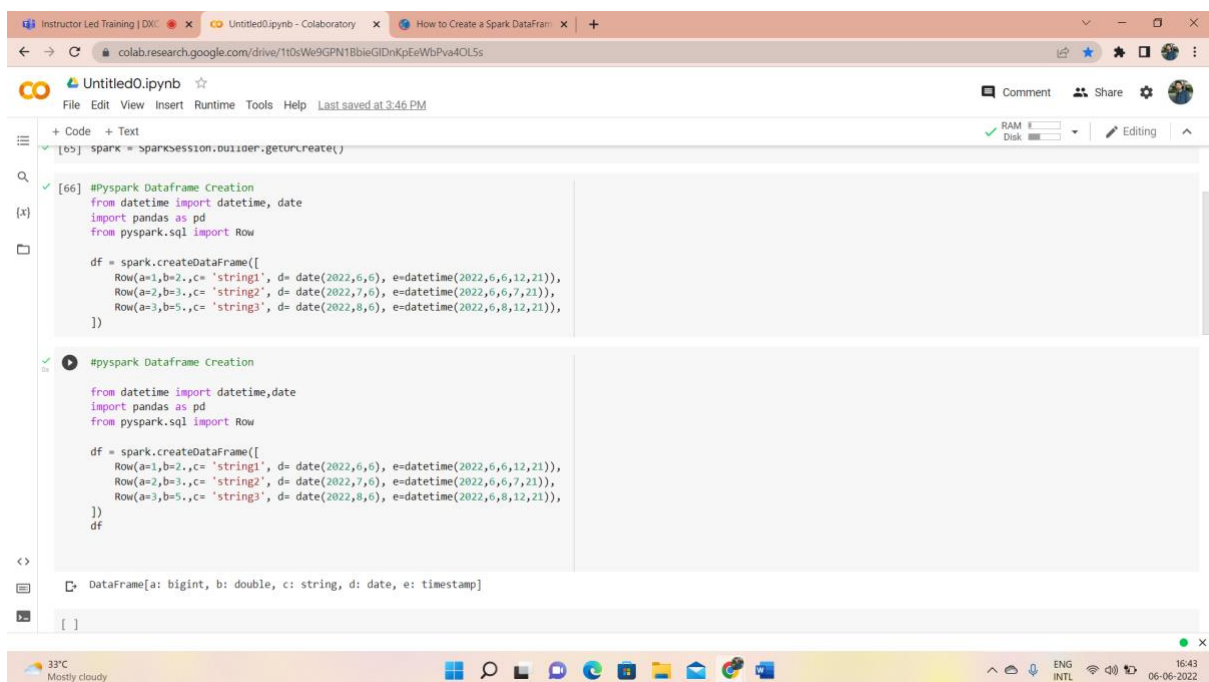
## 6. Explain features of Spark?

- Lightning-fast processing speed.
- Ease of use.
- It offers support for sophisticated analytics.
- Real-time stream processing.
- It is flexible.
- Active and expanding community.
- Spark for Machine Learning.
- Spark for Fog Computing.

## 7. Write a Py-Spark program to create Dataframe from RDD & explain with screenshots & steps ?

```
from pyspark.sql import SparkSession
spark = SparkSession.builder.getOrCreate()
```

```
df = spark.createDataFrame(data)
type(df)
```



The screenshot shows a Google Colaboratory notebook with the following code and output:

```
[66]: #Pyspark Dataframe Creation
from datetime import datetime, date
import pandas as pd
from pyspark.sql import Row

df = spark.createDataFrame([
    Row(a=1,b=2,c= 'string1', d= date(2022,6,6), e=datetime(2022,6,6,12,21)),
    Row(a=2,b=3,c= 'string2', d= date(2022,7,6), e=datetime(2022,6,6,7,21)),
    Row(a=3,b=5,c= 'string3', d= date(2022,8,6), e=datetime(2022,6,8,12,21)),
])

#pyspark Dataframe Creation
from datetime import datetime,date
import pandas as pd
from pyspark.sql import Row

df = spark.createDataFrame([
    Row(a=1,b=2,c= 'string1', d= date(2022,6,6), e=datetime(2022,6,6,12,21)),
    Row(a=2,b=3,c= 'string2', d= date(2022,7,6), e=datetime(2022,6,6,7,21)),
    Row(a=3,b=5,c= 'string3', d= date(2022,8,6), e=datetime(2022,6,8,12,21)),
])
df
```

The output of the second code block is:

```
DataFrame[a: bigint, b: double, c: string, d: date, e: timestamp]
```

```
Instructor Led Training | DXC x Untitled0.ipynb - Colaboratory x How to Create a Spark DataFrame x
colab.research.google.com/drive/110sWe9GPN1BbieGIDnKpEeWbPva4OL5s

Untitled0.ipynb
File Edit View Insert Runtime Tools Help Last saved at 3:46 PM
+ Code + Text
#Creating data frame from Explicit Schema
df = spark.createDataFrame([
    (1,2,'string1', date(2022,6,6), datetime(2022,6,6,12,21)),
    (2,3,'string2', date(2022,7,6), datetime(2022,6,6,7,21)),
    (3,4,'string3', date(2022,8,6), datetime(2022,6,8,12,21)),
],
    schema = 'a long, b double, c string, d date, e timestamp')

DataFrame[a: bigint, b: double, c: string, d: date, e: timestamp]

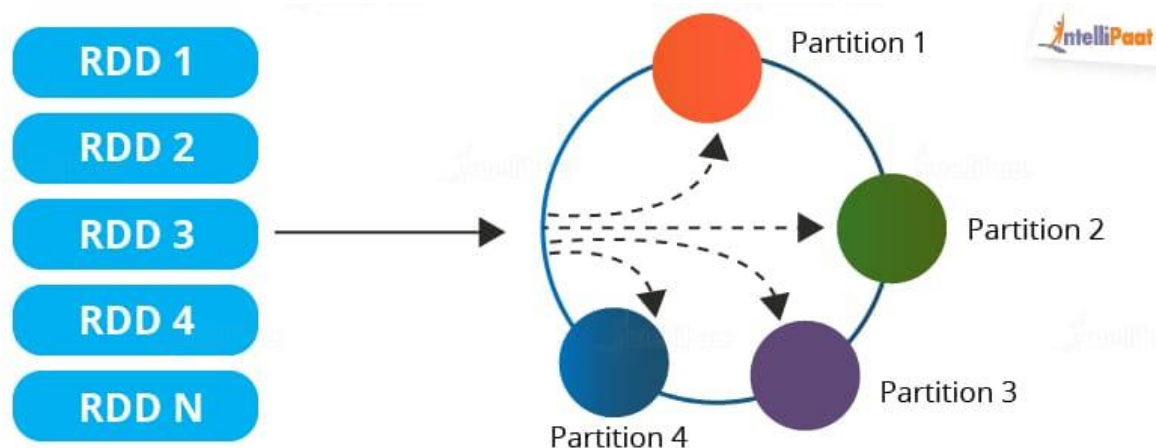
[71] #Create pyspark dataframe from pandas DF
pandas_df = pd.DataFrame({
    'a': [1,2,3],
    'b': [2,3,4],
    'c': ['string1', 'string2', 'string3'],
    'd': [date(2022,6,6),datetime(2022,7,6),datetime(2022,8,6)],
    'e': [datetime(2022,6,6,12,30),datetime(2022,6,7,12,30),datetime(2022,6,8,12,30)]
})

df = spark.createDataFrame(pandas_df)

DataFrame[a: bigint, b: bigint, c: string, d: timestamp, e: timestamp]
```

8. Explain what is RDD & why it is needed ?

---→RDD (Resilient Distributed Dataset)



1. You want low-level transformation and actions and control on your dataset;
2. Your data is unstructured, such as media streams or streams of text;
3. You want to manipulate your data with functional programming constructs than domain specific expressions;
4. You don't care about imposing a schema, such as columnar format while processing or accessing data attributes by name or column; and

5. You can forgo some optimization and performance benefits available with DataFrames and Datasets for structured and semi-structured data.

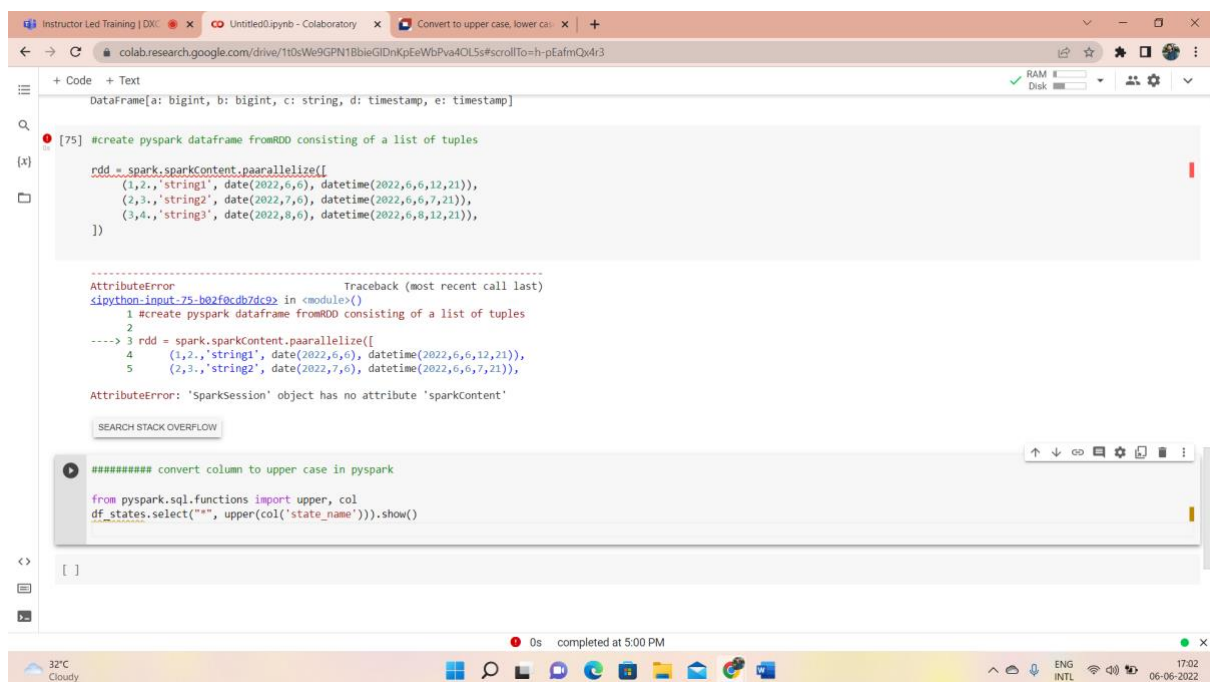
9. Write a Py-Spark program to make the column in Upper case & explain with screenshots & steps ?

- Convert column to upper case in pyspark – upper() function
- Convert column to lower case in pyspark – lower() function
- Convert column to title case or proper case in pyspark – initcap() function

We will be using dataframe **df\_states**

# convert column to upper case in pyspark

```
from pyspark.sql.functions import upper, col
df_states.select("*", upper(col('state_name'))).show()
```



The screenshot shows a Google Colaboratory notebook interface. At the top, there are three tabs: 'Instructor Led Training | DX...', 'Untitled0.ipynb - Colaboratory', and 'Convert to upper case, lower ca...'. The main code area contains a cell with a red error icon. The code in the cell is:

```
[75] #create pyspark dataframe fromRDD consisting of a list of tuples

rdd = spark.sparkContent.parallelize([
    (1,2,'string1', date(2022,6,6), datetime(2022,6,6,12,21)),
    (2,3,'string2', date(2022,7,6), datetime(2022,6,6,7,21)),
    (3,4,'string3', date(2022,8,6), datetime(2022,6,8,12,21)),
])
```

Below the code, a traceback is shown for an `AttributeError`. The error message is: `AttributeError: 'SparkSession' object has no attribute 'sparkContent'`. The traceback indicates the error occurred in the code cell at line 3, column 10.

Below the error, there is a section titled 'SEARCH STACK OVERFLOW'. Below that, there is a code cell with a green icon, containing the following code:

```
##### convert column to upper case in pyspark

from pyspark.sql.functions import upper, col
df_states.select("*", upper(col('state_name'))).show()
```

At the bottom of the notebook, there is a status bar showing '0s completed at 5:00 PM' and a system tray with weather information (32°C Cloudy) and the date/time (17:02 06-06-2022).

