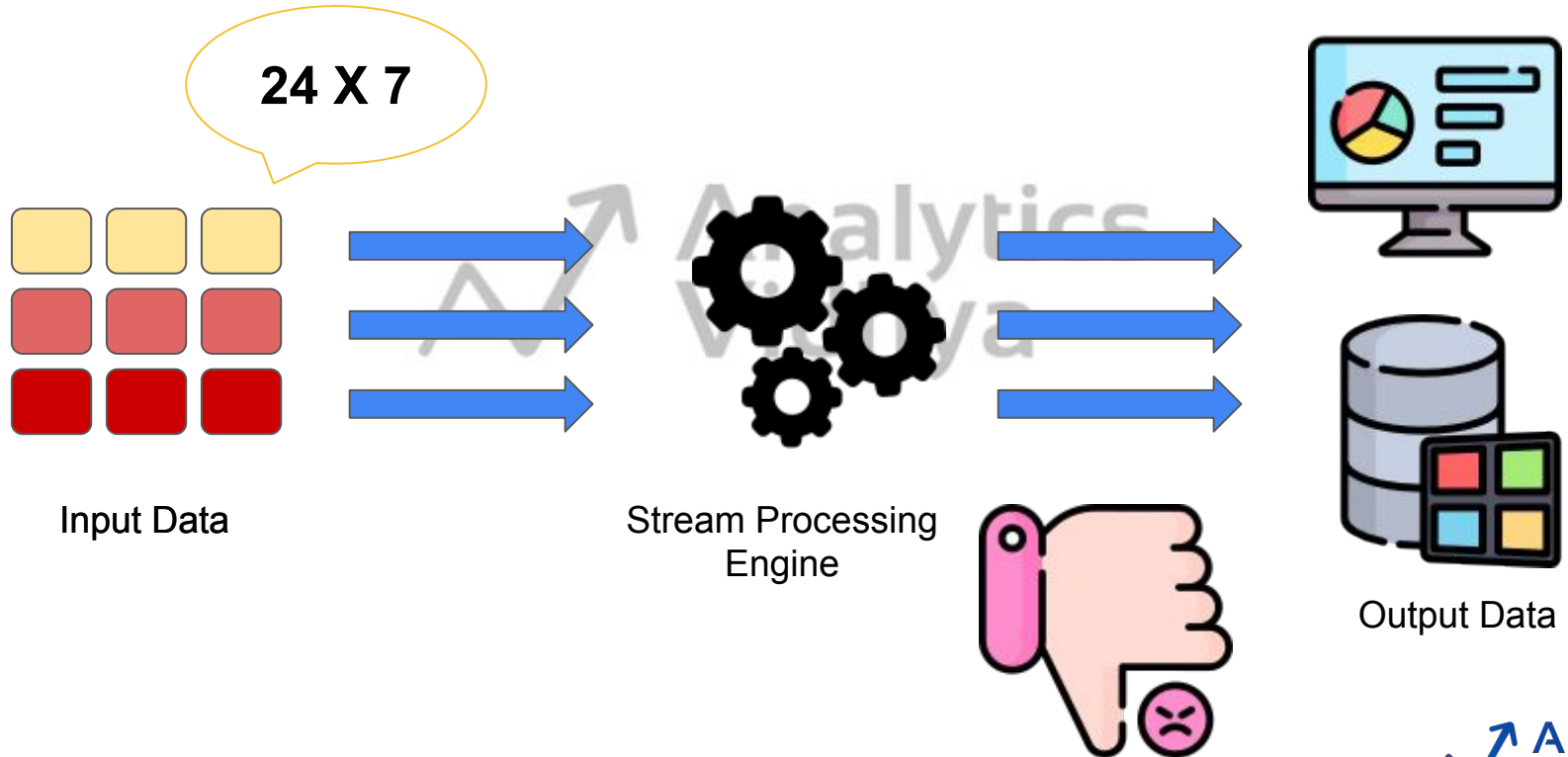


Checkpointing

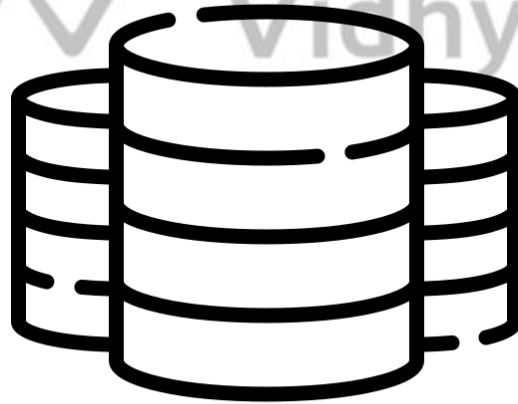


Fault Tolerance

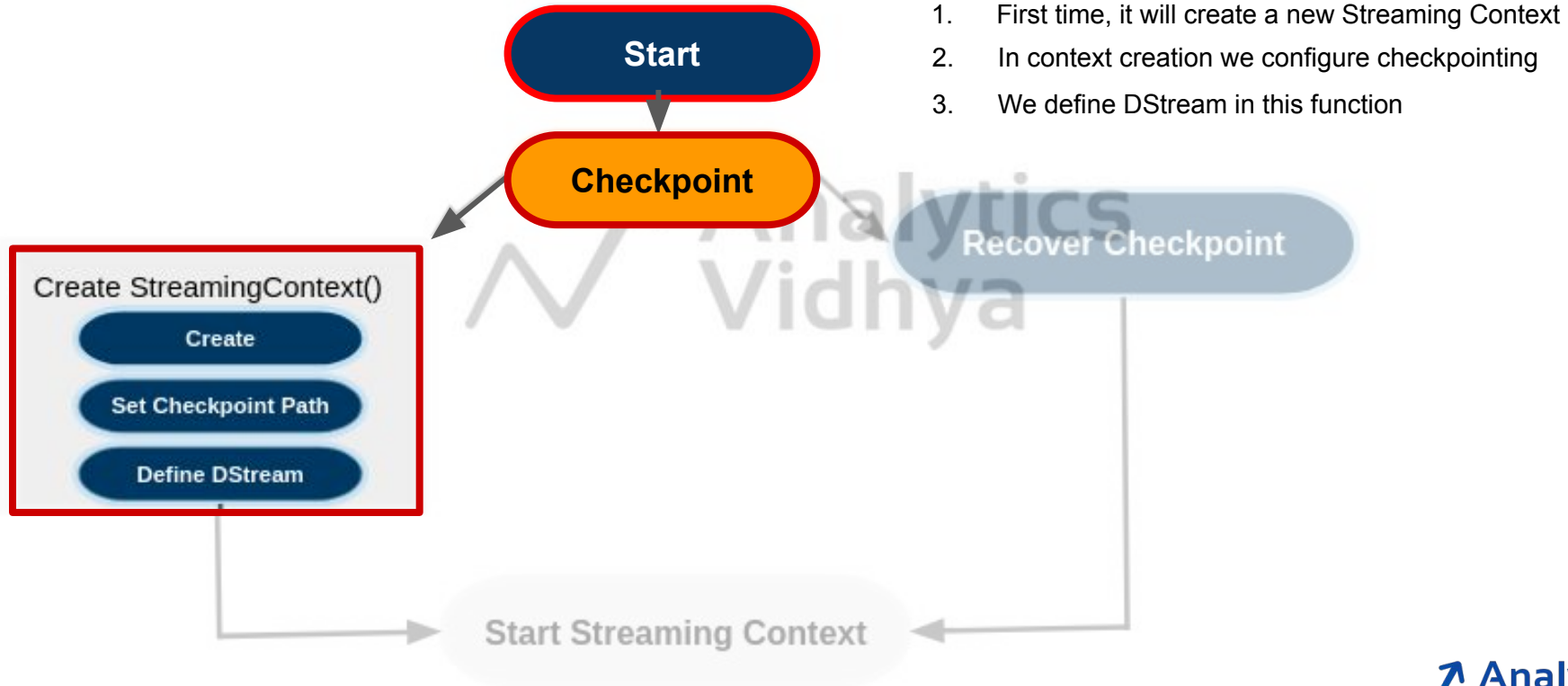


Checkpointing

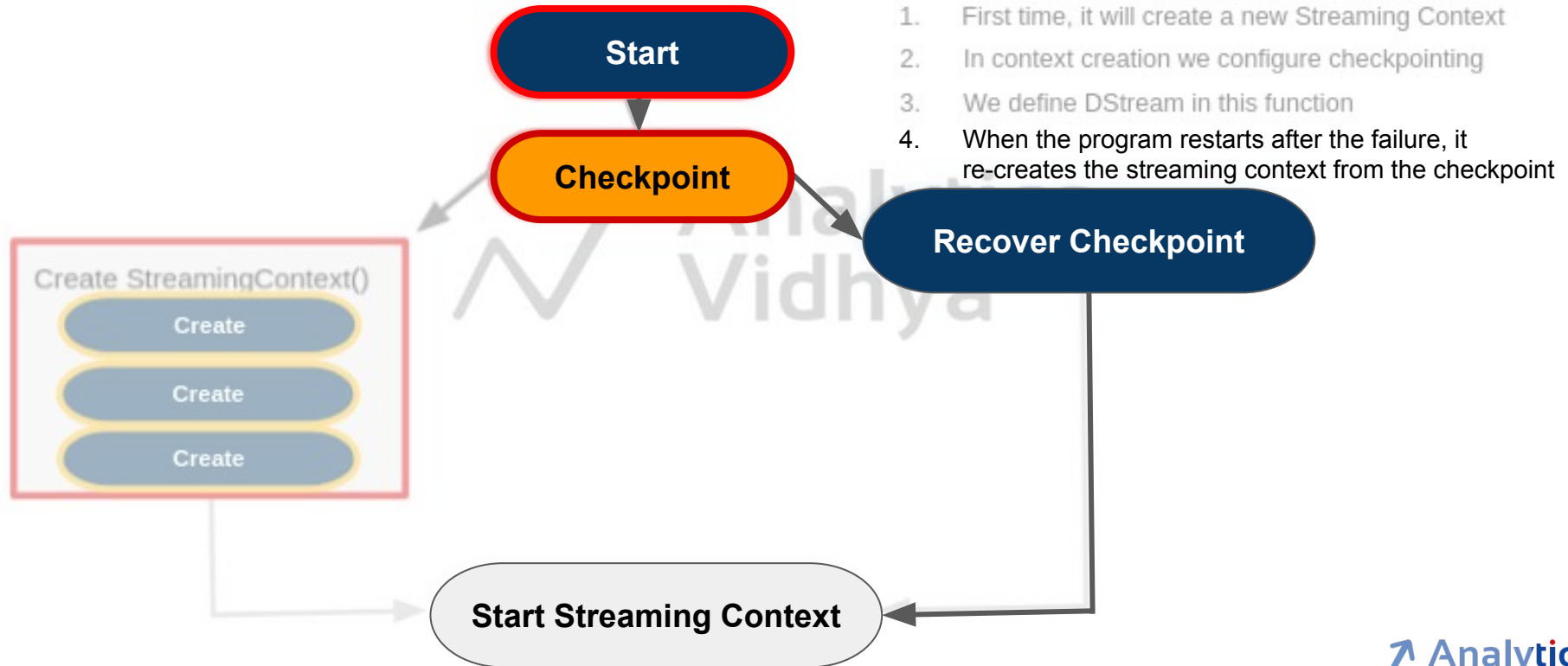
Checkpointing is the process of writing received records at checkpoint intervals to HDFS.



Checkpointing



Checkpointing



Types of Data to Checkpoint

Metadata Checkpointing

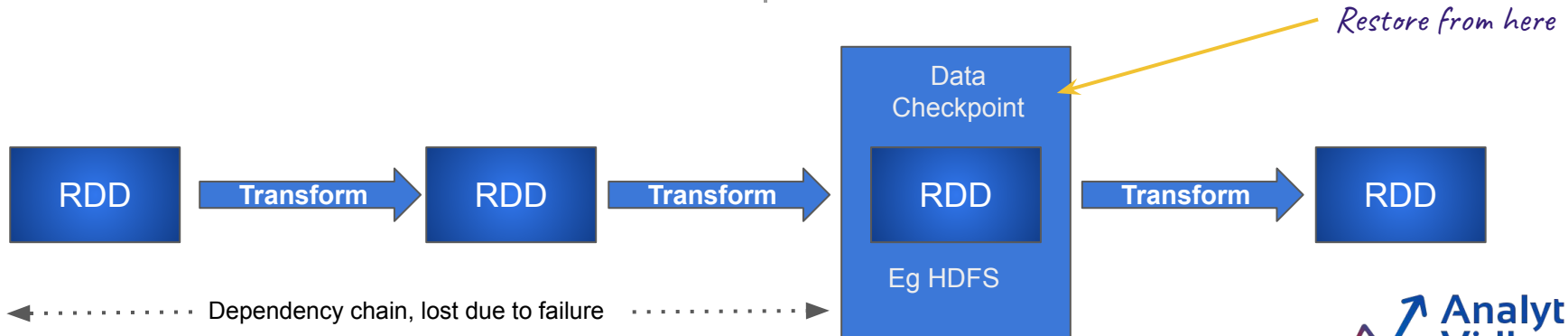


1. Configurations
2. DStream Operations
3. Incomplete batches

Data Checkpointing



1. When upcoming RDD's depend on RDD's of previous batches



Types of Checkpointing

- **Reliable Checkpointing**

- Actual RDD is stored in a reliable distributed system
- Call `SparkContext.setCheckpointDir(directory: String)`

- **Local Checkpointing**

- RDD is persisted to local storage in the executor

Caching/Persisting

persist()

- Cache computed RDD and its **lineage** in **memory**
- **Next operation** on the same dataset will be faster
- DStreams from **window-based operations** are automatically persisted

cache()

- Shorthand for **persist** at the default storage level: **MEMORY_ONLY**

Storage Levels in Apache Spark

MEMORY_ONLY

Allows storage of RDD as deserialized Java objects

Recomputes any RDDs not fitted in memory

MEMORY_AND_DISK

Allows storage of RDD as deserialized objects

Also stores RDDs on disk

MEMORY_ONLY_SER

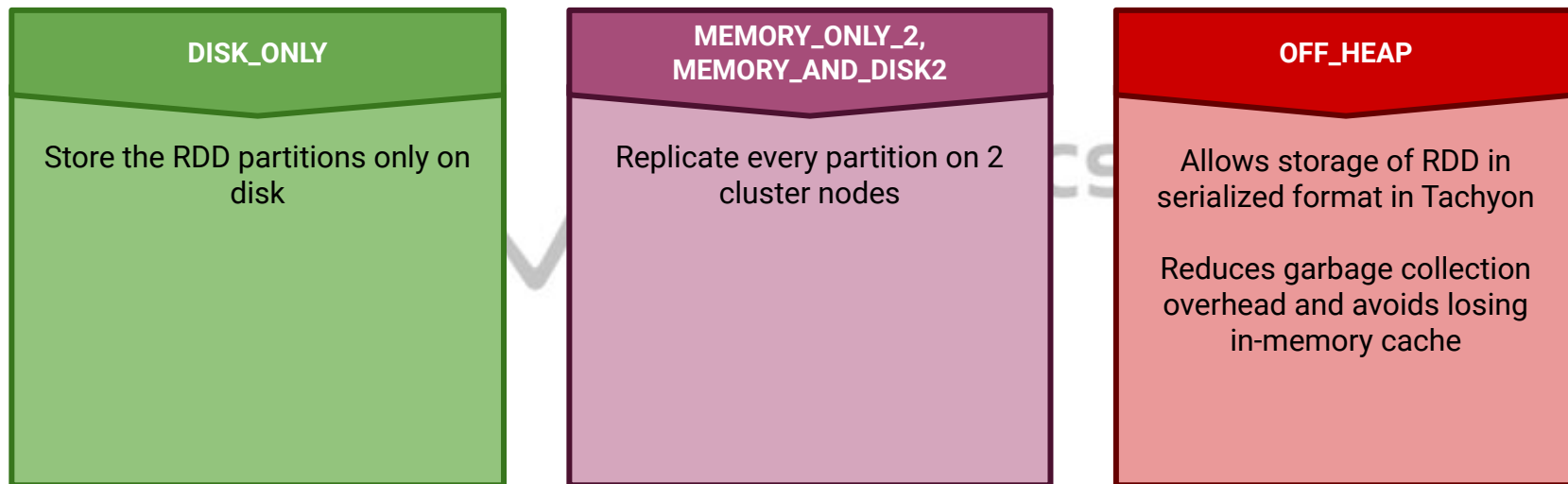
Stores RDD as serialized Java objects

Enables better space efficiency

MEMORY_AND_DISK_SER

Similar to **MEMORY_ONLY_SER**, but spills partitions not fitted in memory to disk

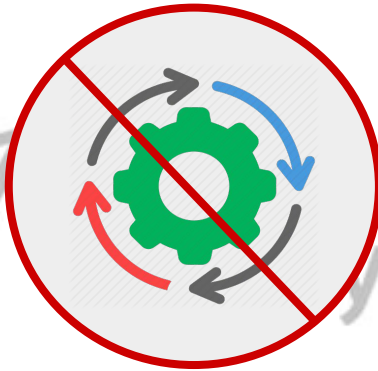
Storage Levels in Apache Spark



Benefits of Caching/Persistence



Saves the data by default in memory



Avoids re-computation of the whole lineage



Improves Performance

When to Enable Checkpointing

Checkpointing must be enabled for applications with any of the following requirements:

- Usage of **stateful transformations**
- Recovering from **failures of the driver** running the application

How to Configure Checkpointing

```
# Function to create and setup a new StreamingContext
def functionToCreateContext():
    sc = SparkContext(...) # new context
    ssc = StreamingContext(...)
    lines = ssc.socketTextStream(...) "
    ...
    ssc.checkpoint(checkpointDirectory) # set checkpoint directory
    return ssc

# Get StreamingContext from checkpoint data or create a new one
context = StreamingContext.getOrCreate(checkpointDirectory, functionToCreateContext)

# Do additional setup on context that needs to be done,
# irrespective of whether it is being started or restarted context. ...

# Start the context
context.start()
context.awaitTermination()
```

Checkpointing can be enabled by setting a directory in a fault-tolerant, reliable file system (e.g., HDFS, S3, etc.) to which the checkpoint information will be saved



Thank You Analytics
Vidhya