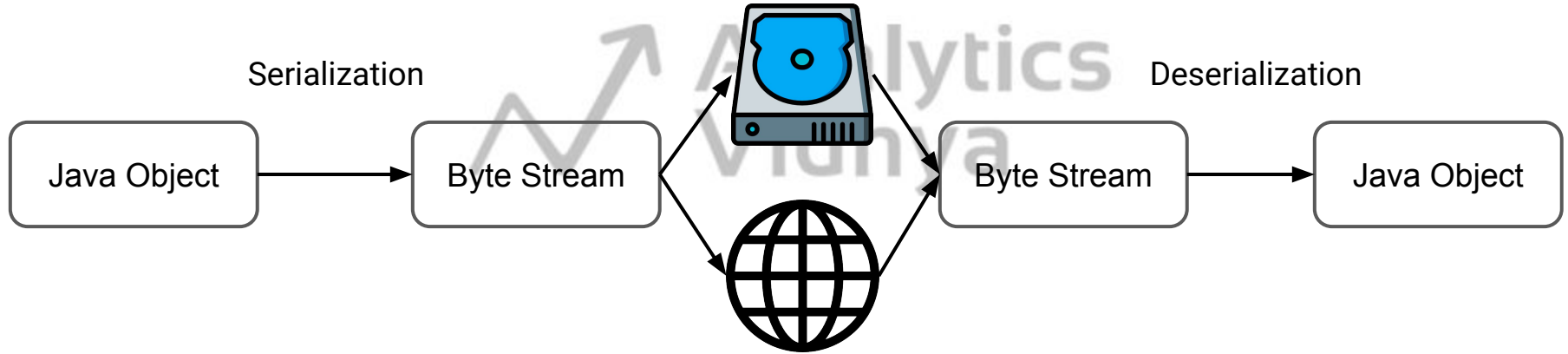


The logo for Analytics Vidhya, featuring a stylized grey arrow pointing upwards and to the right, with the text "Analytics Vidhya" in a grey sans-serif font.

# Storage Levels in Apache Spark

# Serialization & Deserialization



# 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

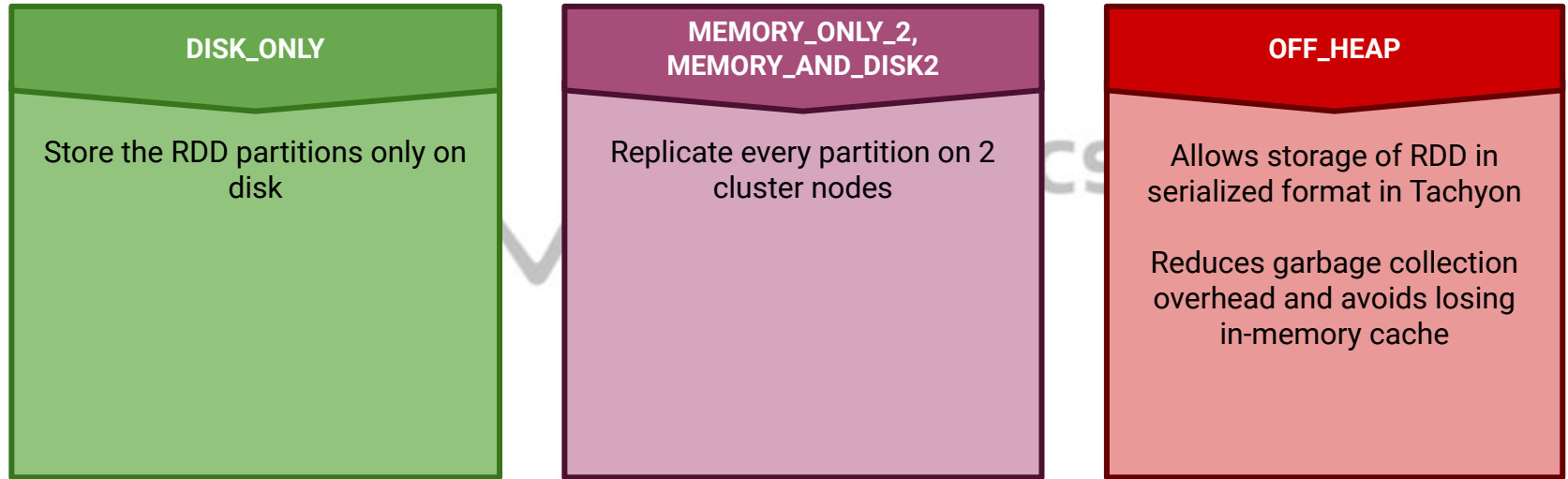
Stores RDD as serialized Java objects

Enables better space efficiency

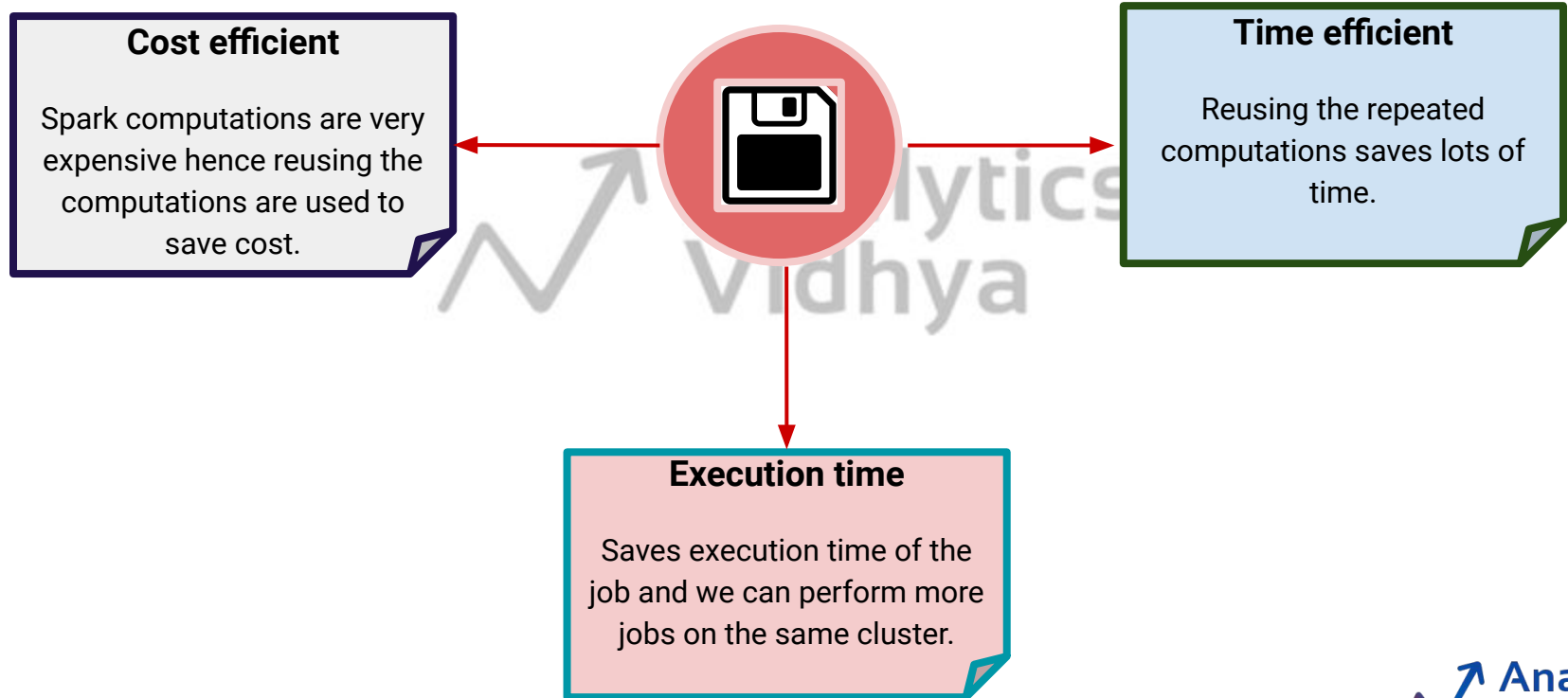
## MEMORY\_AND\_DISK\_SER

Similar to **MEMORY\_ONLY\_SER**, but spills partitions not fitted in memory to disk

# Storage Levels in Apache Spark



# Features of RDD Persistence



# Storage Levels in PySpark

*In Python, stored objects will always be serialized with the [Pickle](#) library, so it does not matter whether you choose a serialized level.*

DISK\_ONLY

DISK\_ONLY\_2

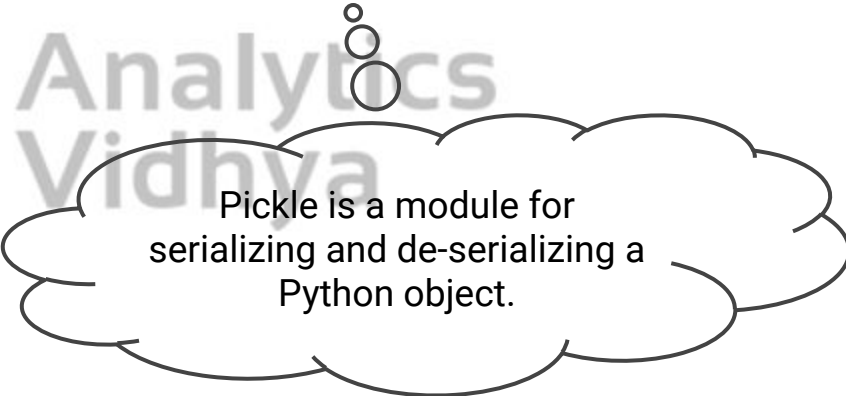
MEMORY\_AND\_DISK

MEMORY\_AND\_DISK\_2

MEMORY\_ONLY

MEMORY\_ONLY\_2

OFF\_HEAP



Pickle is a module for serializing and de-serializing a Python object.

# Which Storage Level to Choose?

Spark's storage levels are meant to provide different trade-offs between memory usage and CPU efficiency.

If your RDDs fit comfortably with the default storage level, leave them that way.

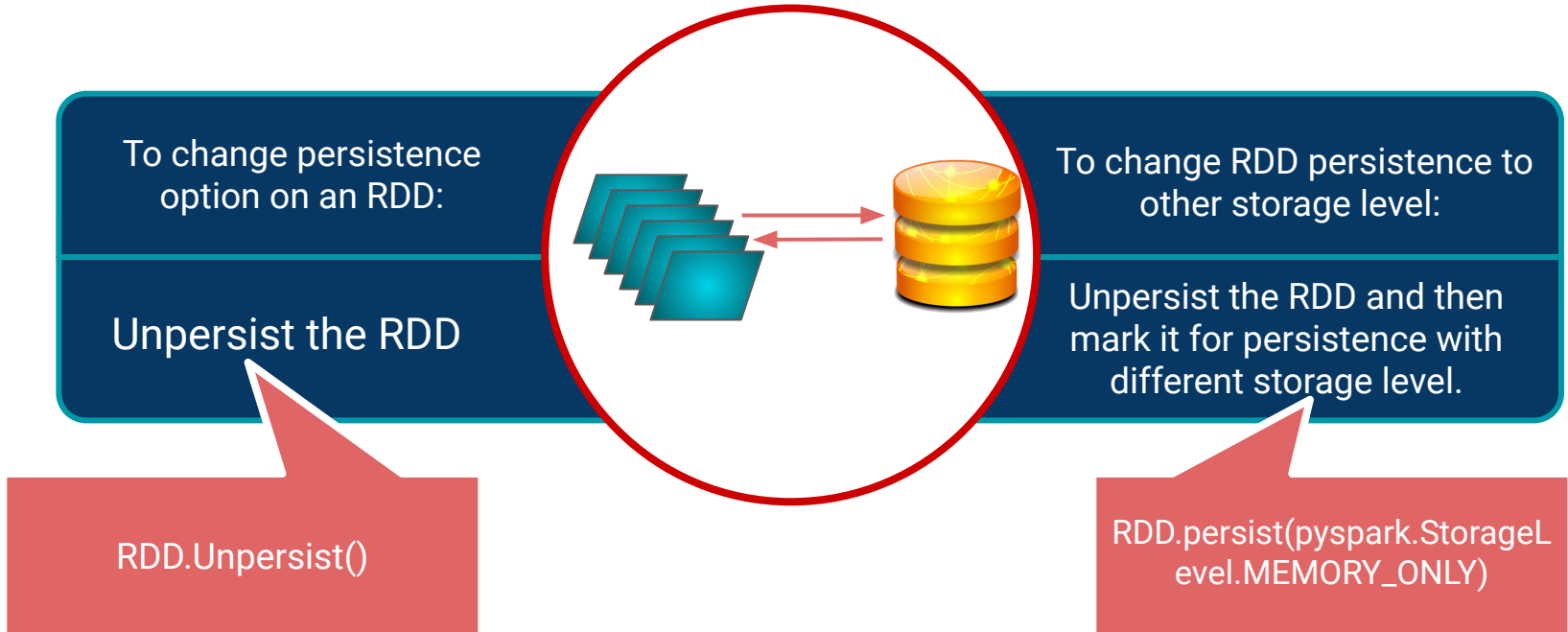
If not, try using `MEMORY_ONLY_SER` and selecting a fast serialization library to make the objects much more space-efficient, but still reasonably fast to access.

Use the replicated storage levels  
If you want fast fault recovery.

Don't spill to disk unless the functions that computed your datasets are expensive.



# Changing Persistence Options







Thank You!!