RDD — (Resilient Distributed Dataset).

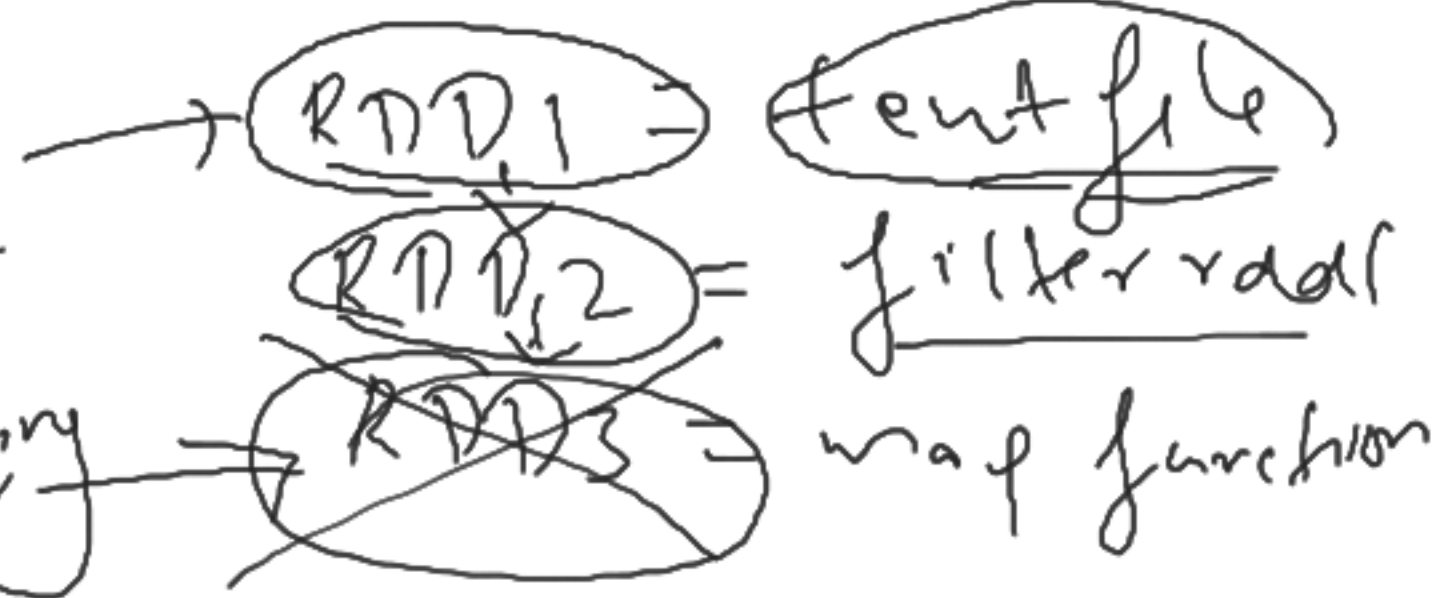Fault tolerance.    spread across    collections of
                    cluster          data

Immutable →

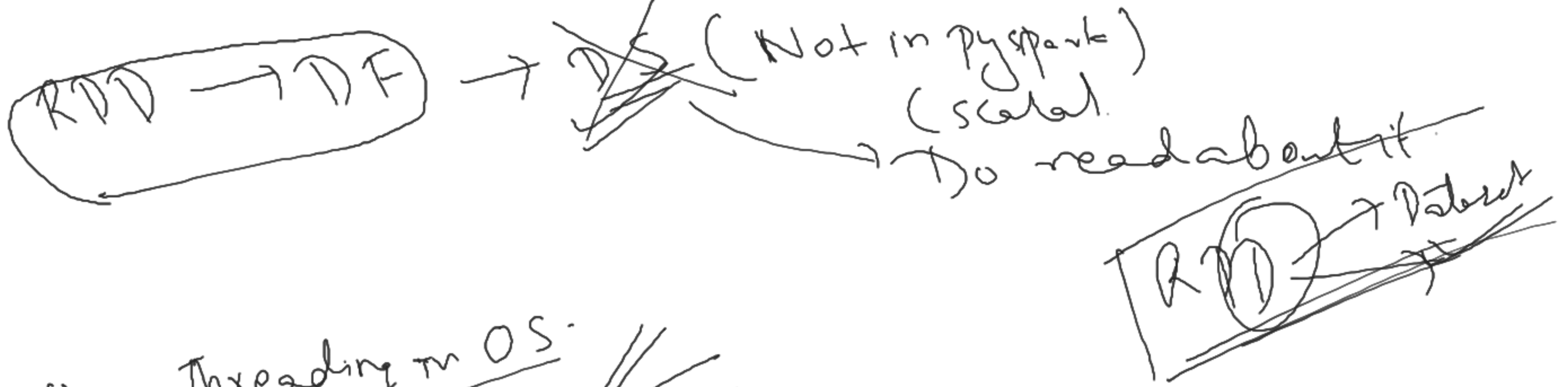recompute →    RDD-lineage graph (DAG). →    RDD 1 = text file

                                             RDD,2 = filter rdd

Apache Spark is an (In-memory data processing) →    RDD 3    map function

framework.                    RAM.

① Spark Context — (Entry point to the spark functionality).

Narrow — No shuffle operations.

Wide — there will be shuffle operations.

Cad dyst optimister and tungsten execution engine

RDD ⟶ DF ⟶ DS (Not in pyspark)
(Scala.
⟶ Do read about it

RDD ⟶ Dataset

Hyper Threading in O.S.

1) Quad core – (4 cores) ⟶ Program running to behave like hardware.

octacore ( )

smartphones – octacore – (Hyper Threading)

② octacore – by hyperthreading

coalesce

4 cores

(1)

(3) → every time
particle

2 cores.

repartition

update

Coalesce —

minimum shuffling ( It doesn't matter whether a are

increasing | decreasing

(2)

(1)

repartition -1 maximum shuffling → Consequences into your consideration

(2) → (8)