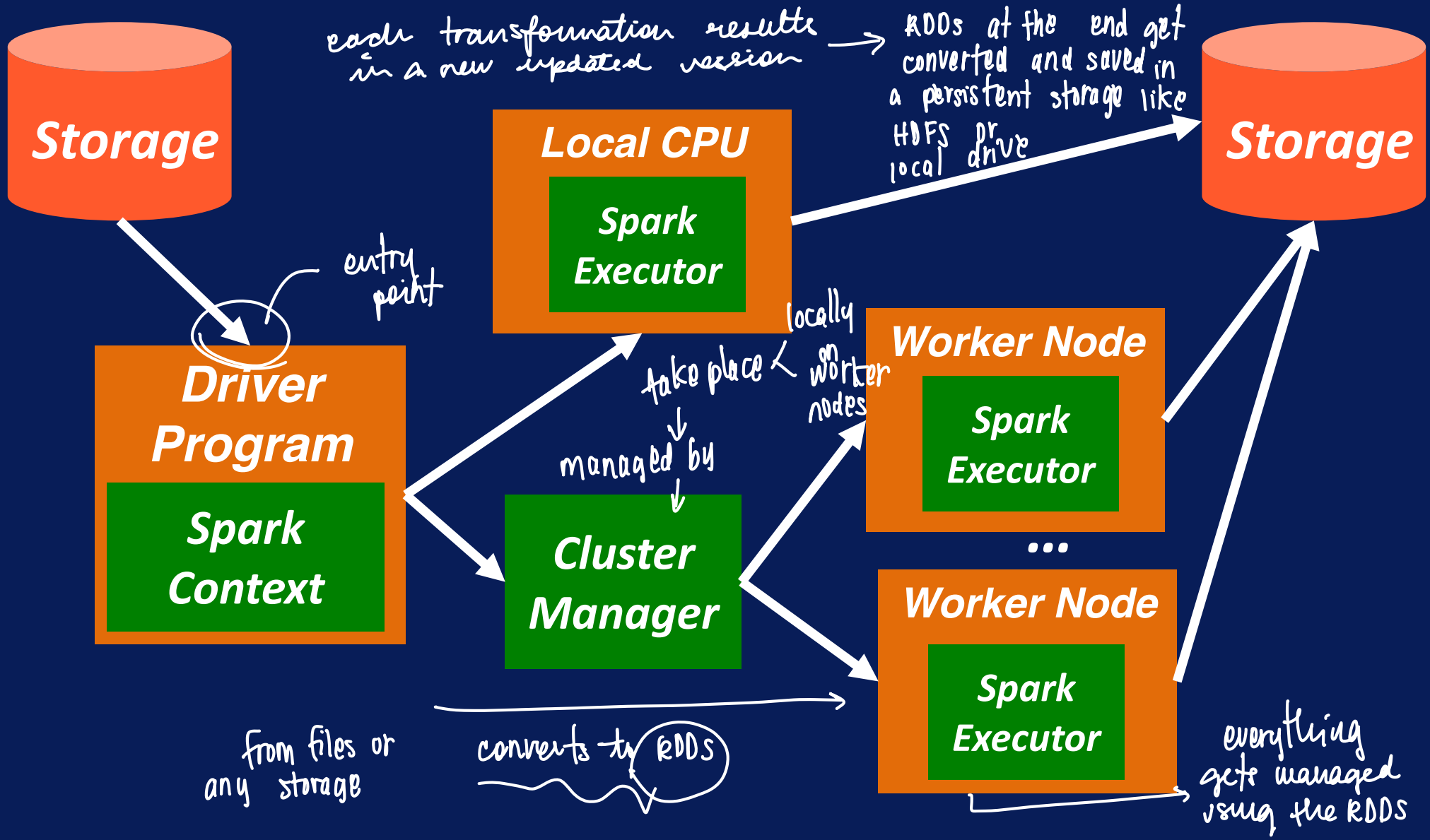


Spark Core: Programming In Spark



After this video you will be able to..

- Use two methods to create RDDs in Spark
- Explain what immutable means
- Interpret a Spark program as a pipeline of transformations and actions
- List the steps to create a Spark program



Creating RDDs

in Driver Program

you can just read files →

```
In [1]: lines = sc.textFile("hdfs:/user/cloudera/words.txt")
```

```
lines = sc.parallelize(["big", "data"])
```

provide existing collection (like list) to be turned into a distributed collection

```
numbers = sc.parallelize(range(10), 3)
```

integer RDD + provide number of partitions

creates 3 partitions

* spark decides how to assign partitions to assign partitions to DLR executors and worker nodes

```
numbers.collect()
```

Parallelize range output into 3 partitions

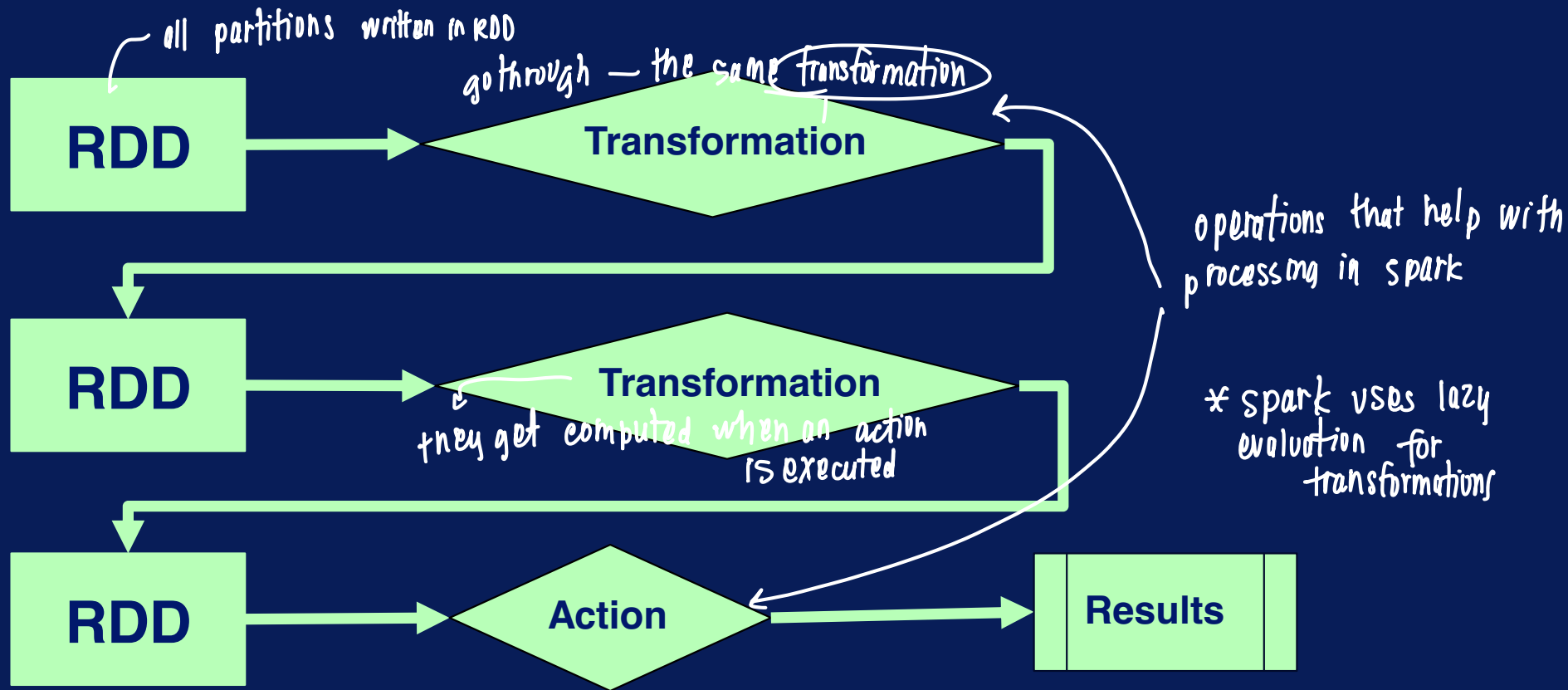
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]

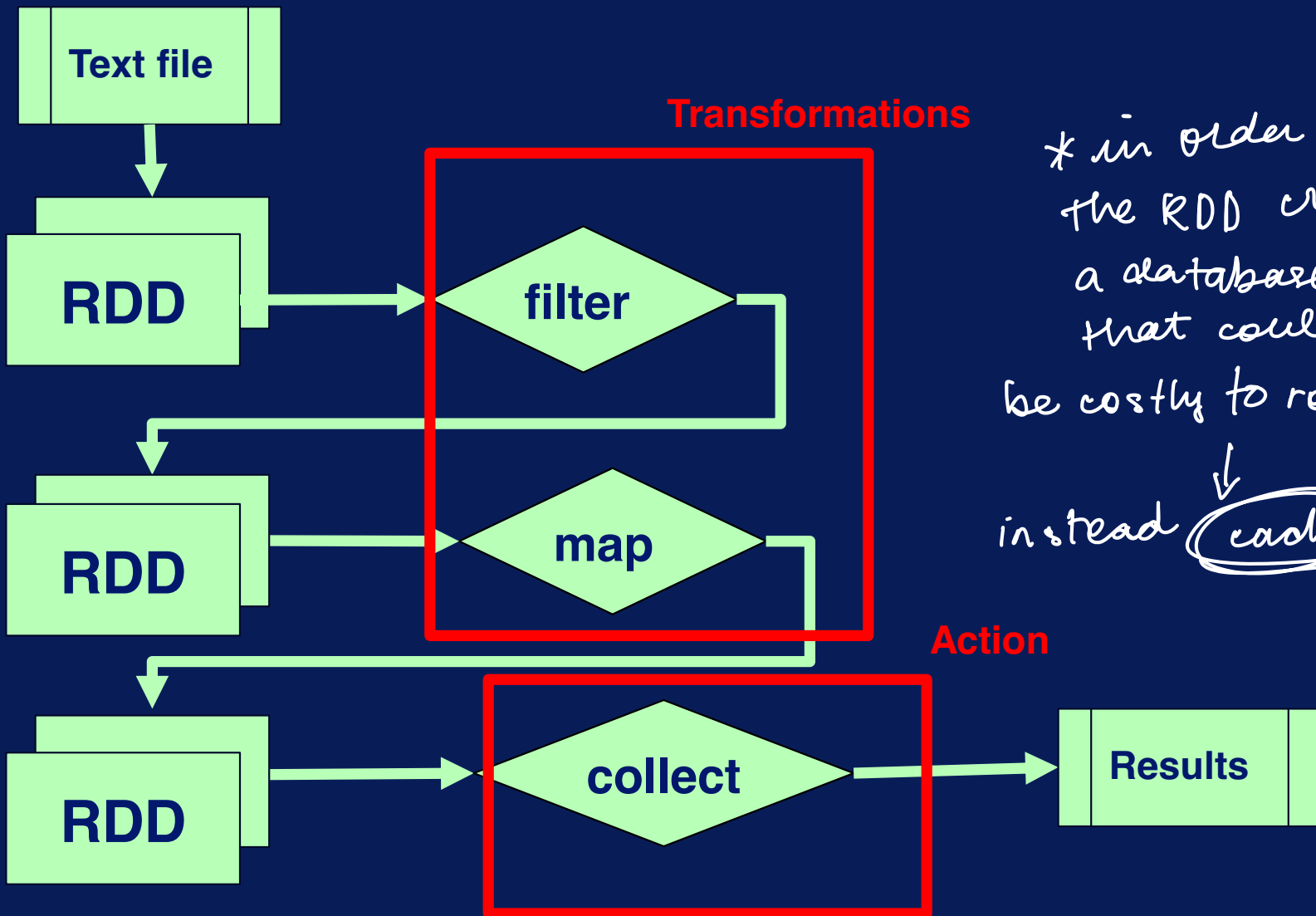
[0, 1, 2], [3, 4, 5], [6, 7, 8, 9]

[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]

can then be gathered into a single partition on the driver using the collect transform

Processing RDDs





Transformations

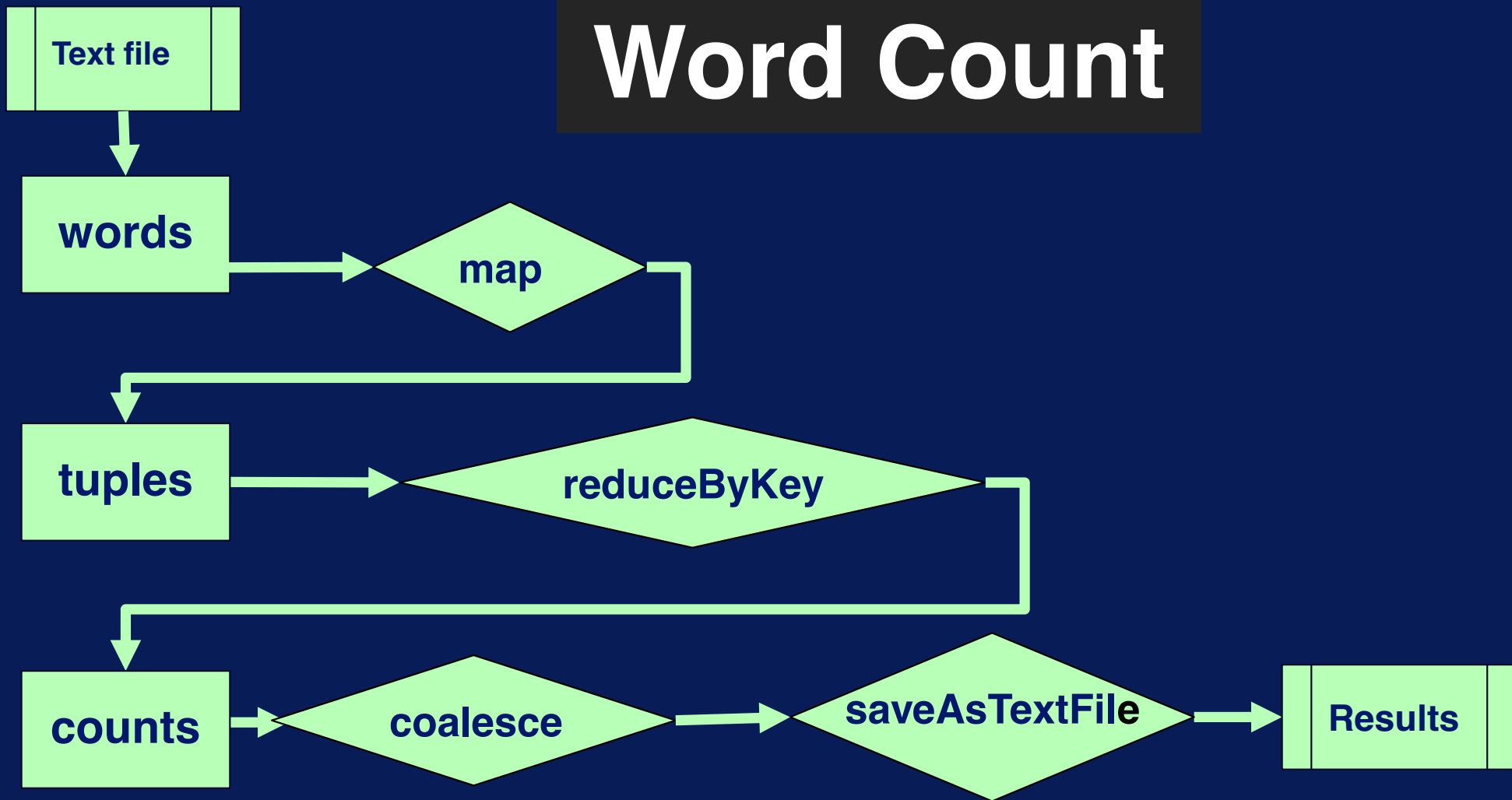
* in order to reuse the RDD created from a database query that could otherwise be costly to re-execute

instead cache these RDDs

Action

⚠ caution
- it can consume too much memory and generate bottleneck

Word Count



Programming in Spark

Create RDDs

from external collections
or from local collections



Apply transformations

like filter
map
reduceBykey

they get lazily evaluated until
action is performed



Perform actions

both for local
parallel computation
to generate results