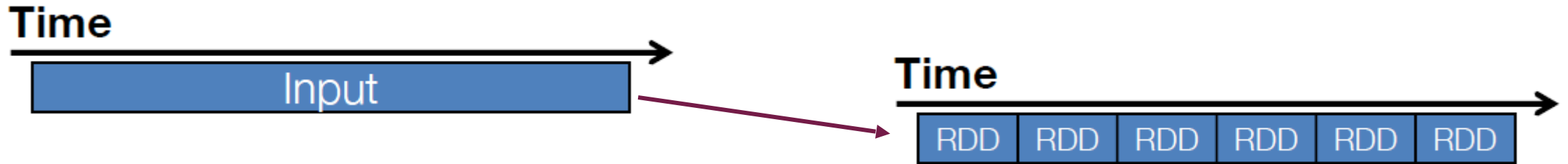


Example: Spark Streaming



Represents streams as a series of RDDs over time (typically sub second intervals, but it is configurable)

```
val spammers = sc.sequenceFile("hdfs://spammers.seq")
sc.twitterStream(...)
  .filter(t => t.text.contains("Santa Clara University"))
  .transform(tweets => tweets.map(t => (t.user, t)).join(spammers))
  .print()
```

Summary

Spark is a powerful “manager” for big data computing.

It centers on a job scheduler for Hadoop (MapReduce) that is smart about where to run each task: co-locate task with data.

The data objects are “RDDs”: a kind of recipe for generating a file from an underlying data collection. RDD caching allows Spark to run mostly from memory-mapped data, for speed.

Next: Resource Scheduling