

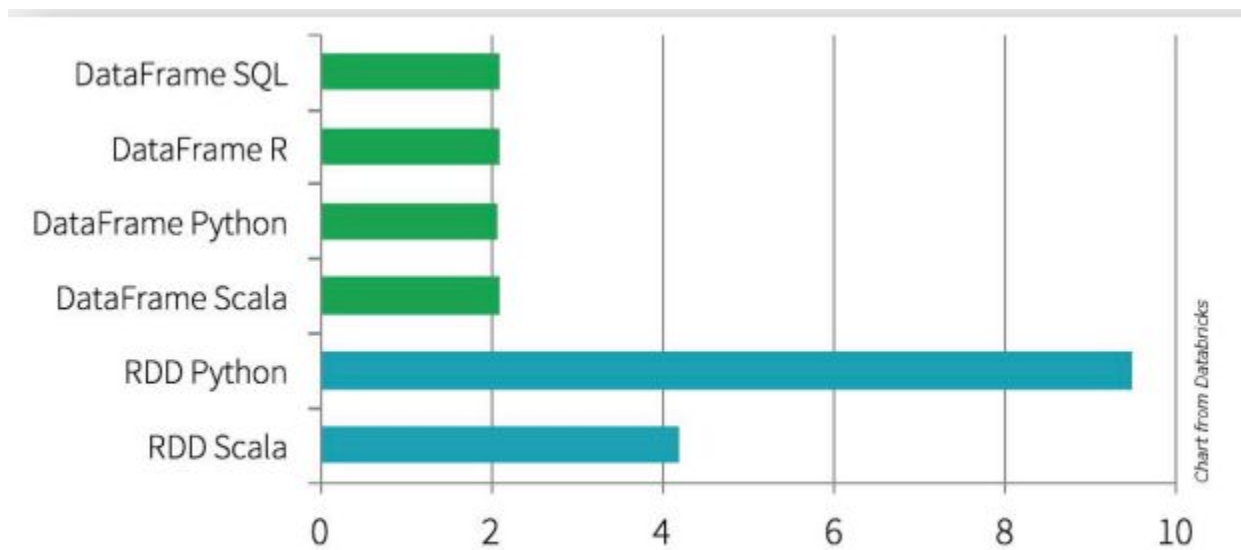
SPARK Interface API

RDD	DataFrame	Dataset
Fault Tolerant	Fault Tolerant	Fault Tolerant
Distributed	Distributed	Distributed
Immutability	Immutability	Immutability
No schema	Schema	Schema
Slow on Non-JVM languages	Faster	Faster
No Execution optimization	optimization Catalyst optimizer	optimization
Low Level	High Level	High Level
No SQL Support	SQL Support	SQL Support
Type Safe	No type Safe	Type Safe
Syntax Error detected at Compile Time	Syntax Error detected at Compile Time	Syntax Error detected at Compile Time
Analysis Error Detected at Compile time	Analysis Error Detected at Run time	Analysis Error Detected at Compile time
JAVA,SCALA, Python,R	JAVA,SCALA, Python,R	JAVA, SCALA
Higher memory is used	Higher memory is used	Low memory is used. Tungsten encoders provide great benefits

We can seamlessly move between DataFrame or Dataset and RDDs by using simple syntax

```
Schema = student (  
    name :String,  
    age :Int,  
    department : String,  
    status : String,  
    remark : String  
)
```

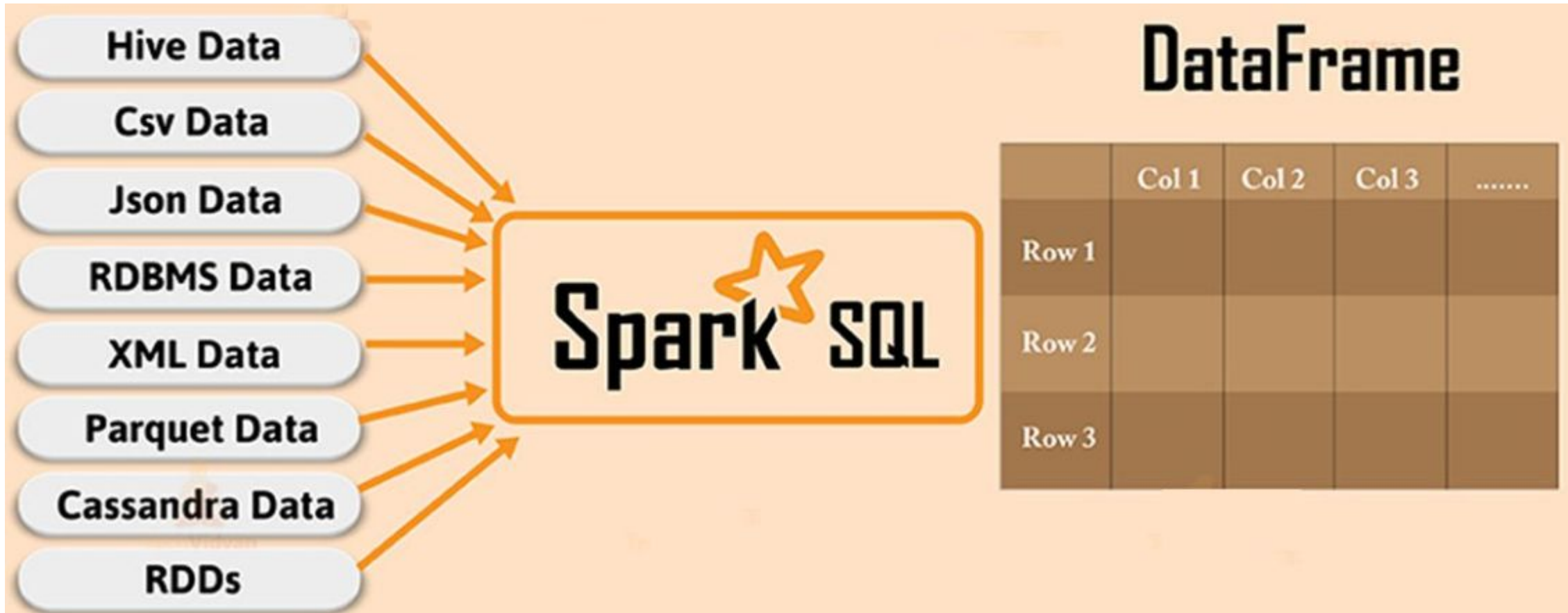
```
Data =  
student(A,17,DS,Pass,Awesome),  
student(B,19,DE,Pass,Excellent),  
student(C,17,AE,Pass,Excellent),  
student(D,18,BD,Pass,Awesome)
```



Time to aggregate 10 million integer pairs (in seconds)

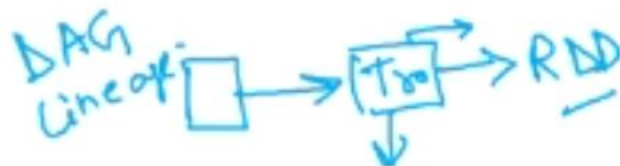
Chart from Databricks

Ways to Create A DATAFRAME



Type Safe

Data Set = DataFrame + RDD



SPARK Interface API

Name	Age
Sh	Int

RDD	DataFrame	Dataset
Fault Tolerant	Fault Tolerant	Fault Tolerant
Distributed	Distributed	Distributed
Immutability	Immutability	Immutability
No schema	Schema	Schema
Slow on Non-JVM languages	Faster	Faster
No Execution optimization	optimization Catalyst optimizer	optimization
Low Level	High Level	High Level
No SQL Support	SQL Support	SQL Support
Type Safe	No type Safe	Type Safe
Syntax Error detected at Compile Time	Syntax Error detected at Compile Time	Syntax Error detected at Compile Time
Analysis Error Detected at Compile time	non-existent element. Analysis Error Detected at Run time	Analysis Error Detected at Compile time
JAVA, SCALA, Python, R	JAVA, SCALA, Python, R	JAVA, SCALA
Higher memory is used	Higher memory is used	Low memory is used. Tungsten encoders provide great benefits

We can seamlessly move between DataFrame or Dataset and RDDs by using simple syntax

Unified DF
SPARK

④

RUN
COMPILE