

# The structure spectrum

## Structured

Relational  
Databases

Parquet

Formatted  
Messages

## Semi-structured

HTML

XML

JSON

## Unstructured

Plain text

Generic media



# Tabular data

- Simple data format
- Common variants:
  - Comma Separated Values (CSV),
  - Tab Separated Values (TSV)



# Tabular data

1	Genom	—	PR	PR	—	3	AA	—	—
2	skattereformen	—	NN	NN	—	1	PA	—	—
3	införs	—	VV	VV	—	0	ROOT	—	—
4	individuell	—	AJ	AJ	—	5	AT	—	—
5	beskattning	—	VN	VN	—	3	SS	—	—
6	(	—	IR	IR	—	5	IR	—	—
7	särbeskattning	—	VN	VN	—	5	AN	—	—
8	)	—	IR	IR	—	5	JR	—	—
9	av	—	PR	PR	—	5	ET	—	—
10	arbetsinkomster	—	NN	NN	—	9	PA	—	—
11	.	—	IP	IP	—	3	IP	—	—

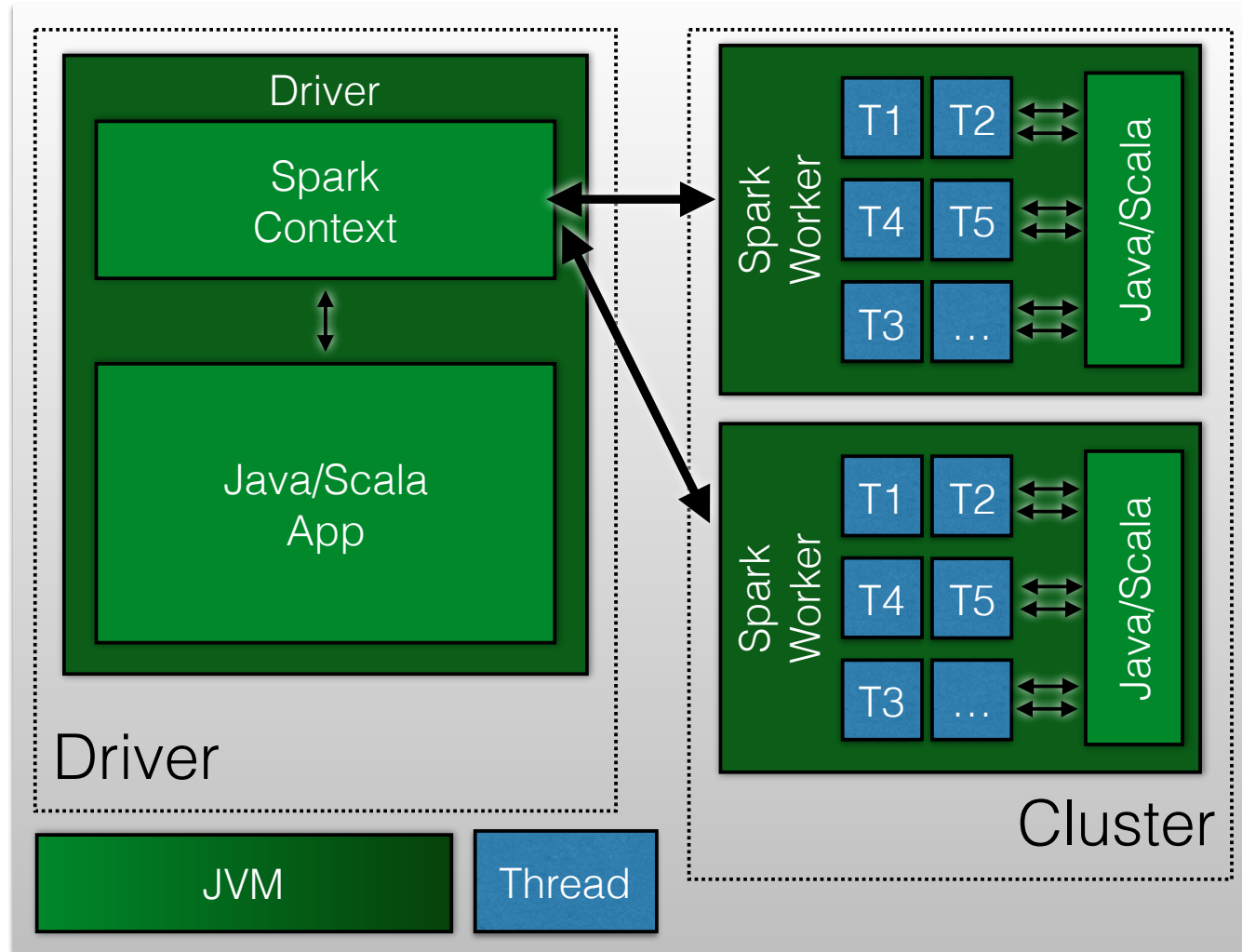


# Spark SQL

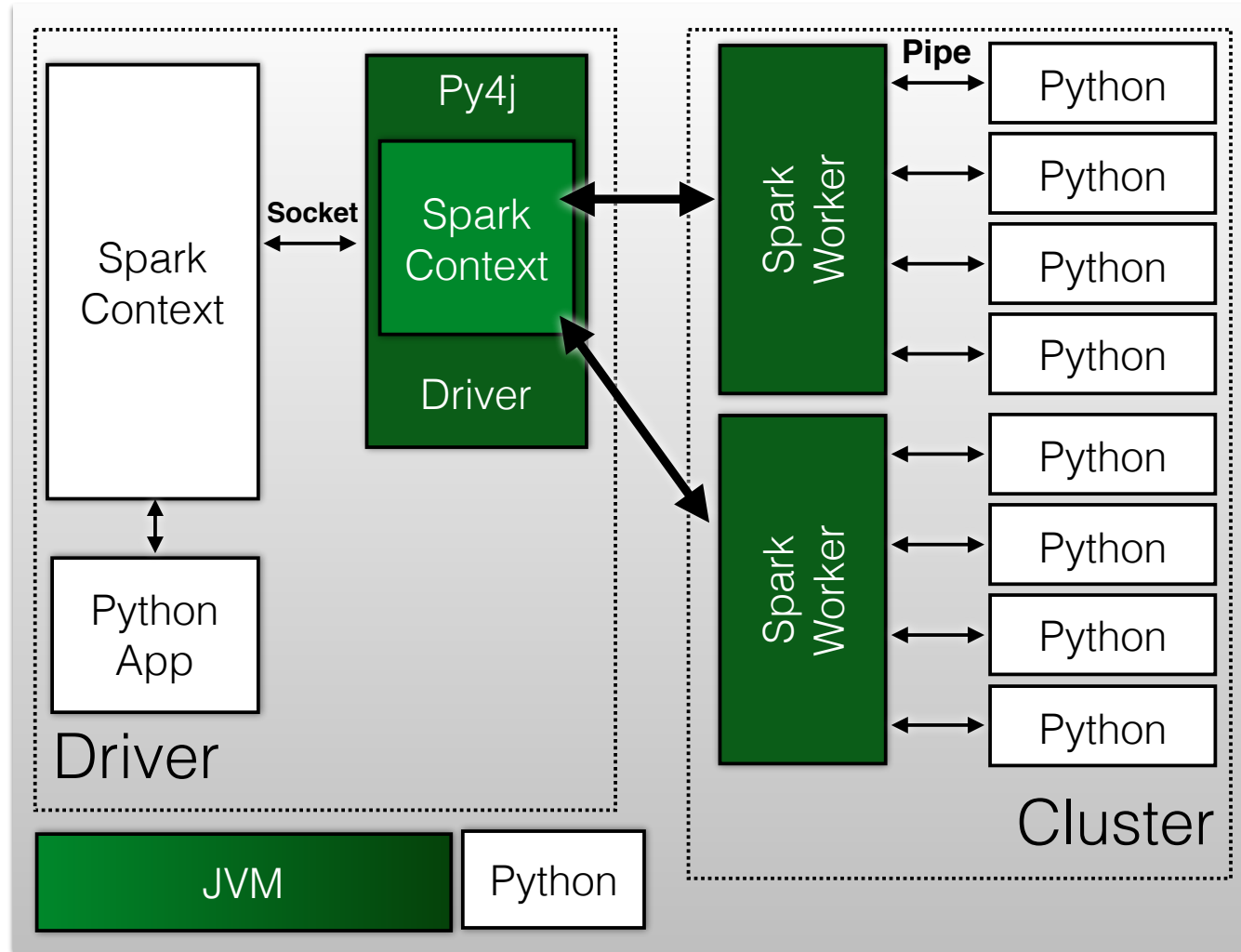
- **Why another API?**
  - Performance



# Spark: Java/Scala

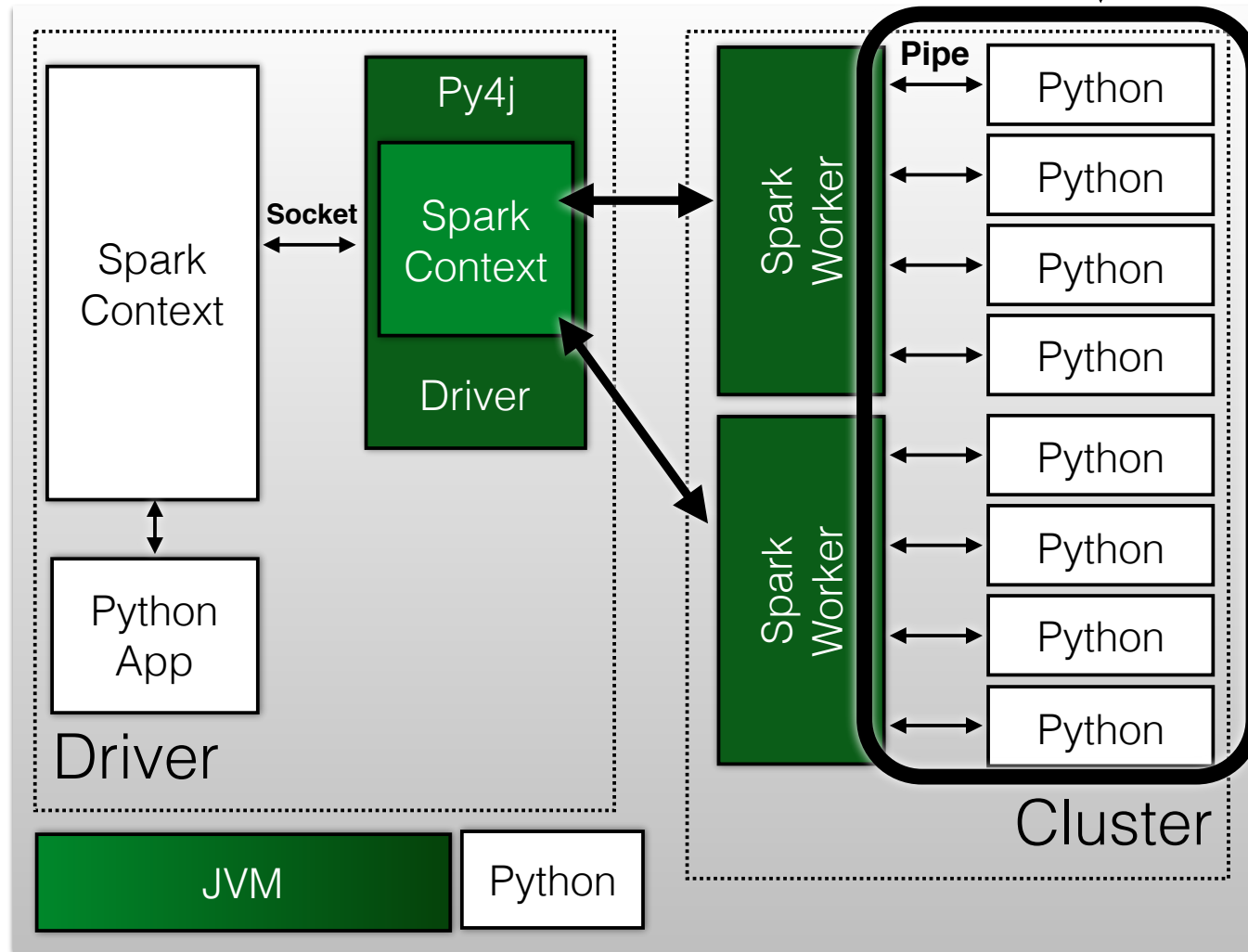


# Pyspark



# Pyspark

## Python and pipe Overhead



# Pyspark: Performance

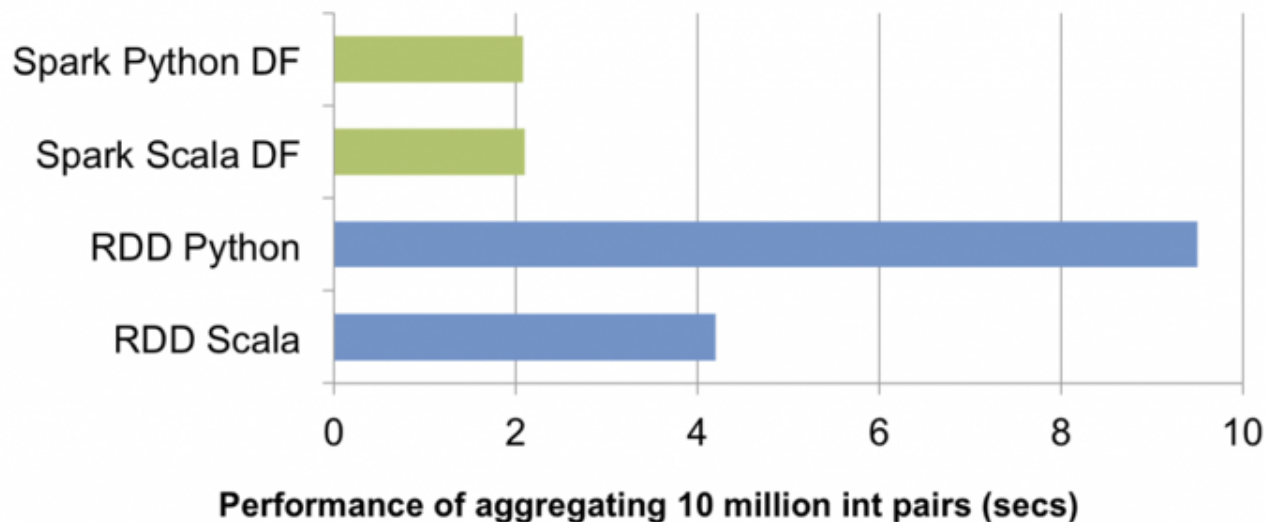
- Serialization/Deserialization overhead
- CPython / PyPy is officially supported.  
CPython still has most support since numpy is not yet supported in PyPy
- Python code performance issues if most time spent computing is not within native libraries
- **How can most of the headache of overhead be reduced?**





# Pyspark SQL

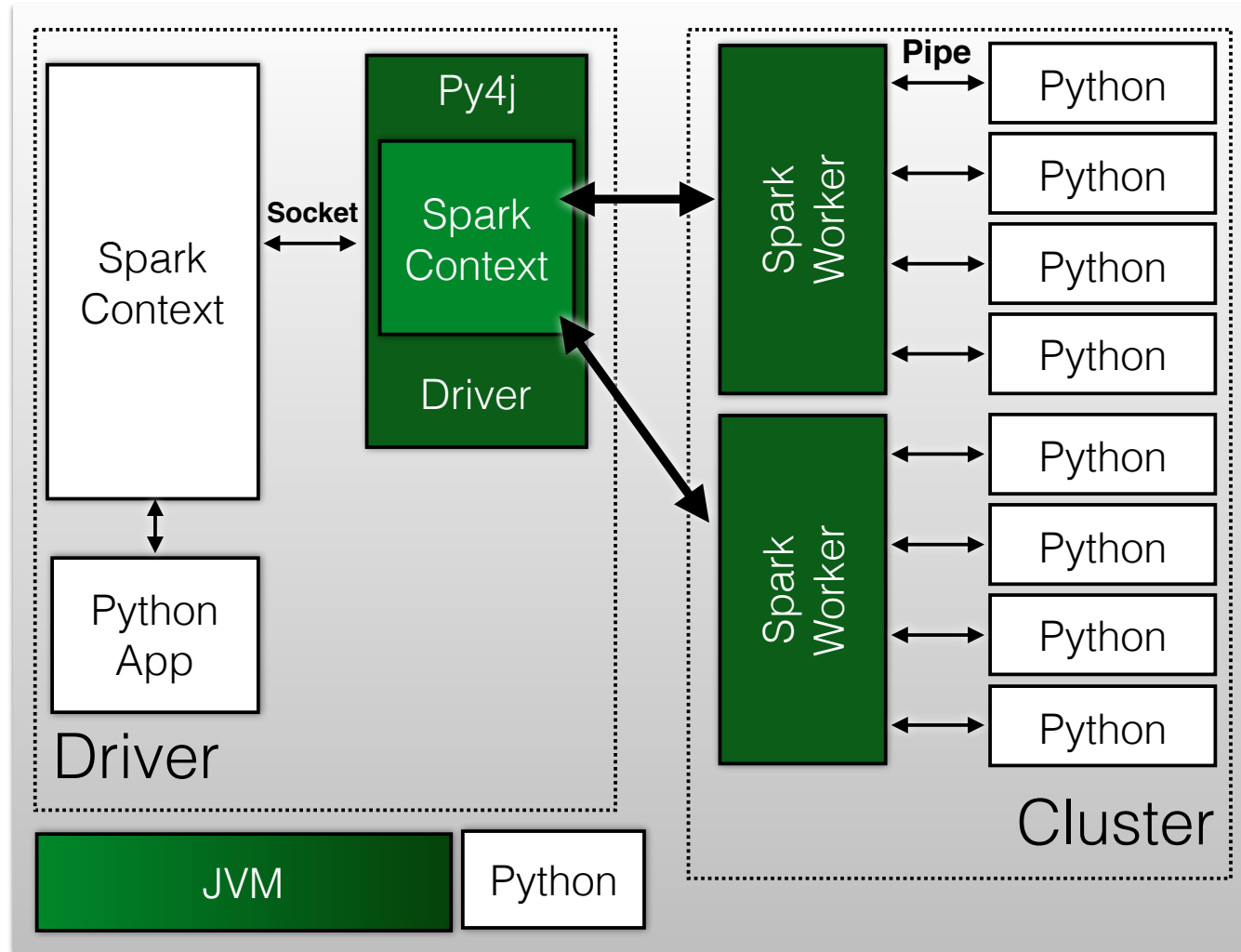
- Optimized data processing using optimized data structures and spark code, no extra overhead.
- New in Spark 1.3.0: Dataframes (SQL)



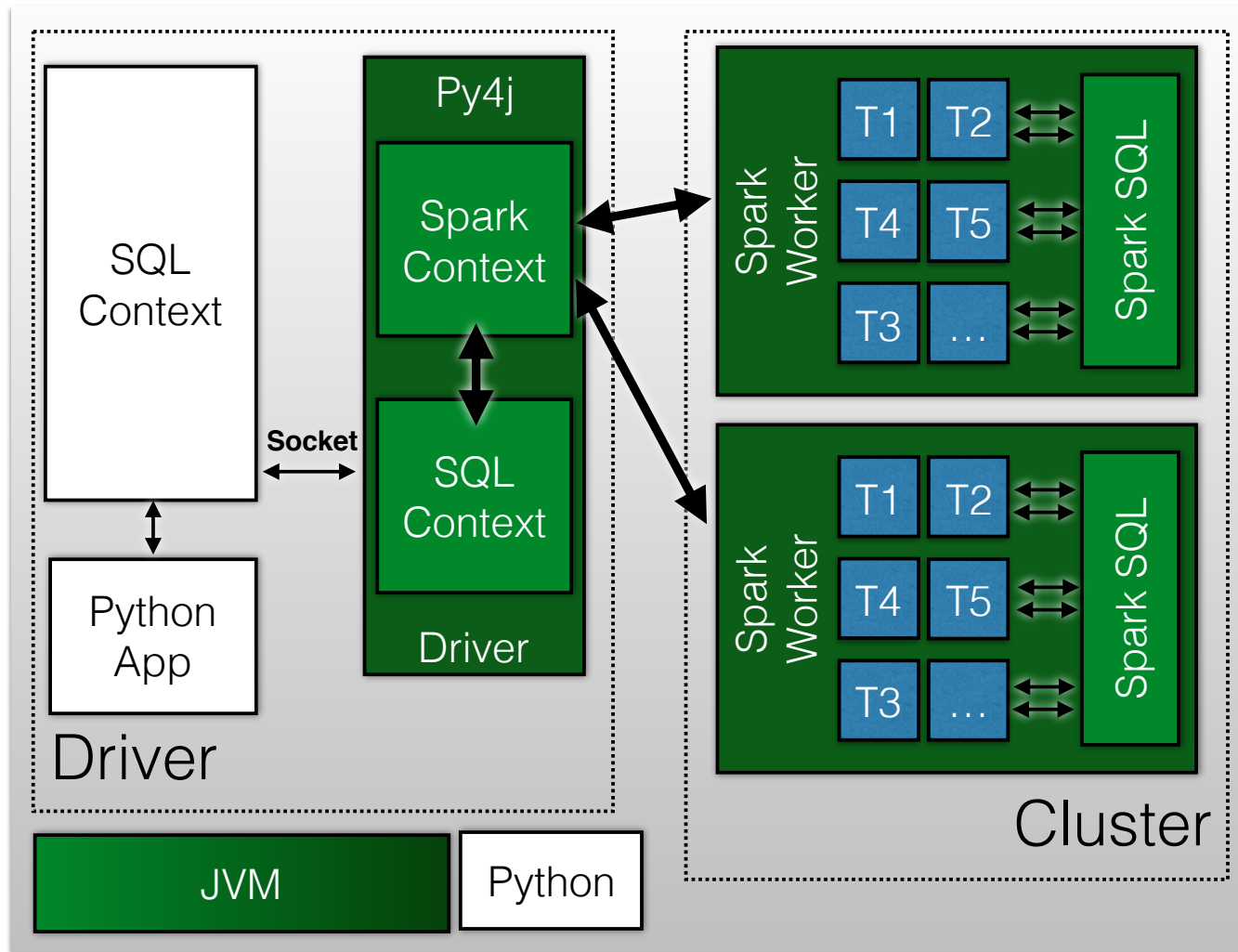
Source: <https://databricks.com/blog/2015/02/17/introducing-dataframes-in-spark-for-large-scale-data-science.html>



# Pyspark



# Pyspark SQL



# Pyspark

**No overhead  
+ Native Spark SQL  
Code only.**

