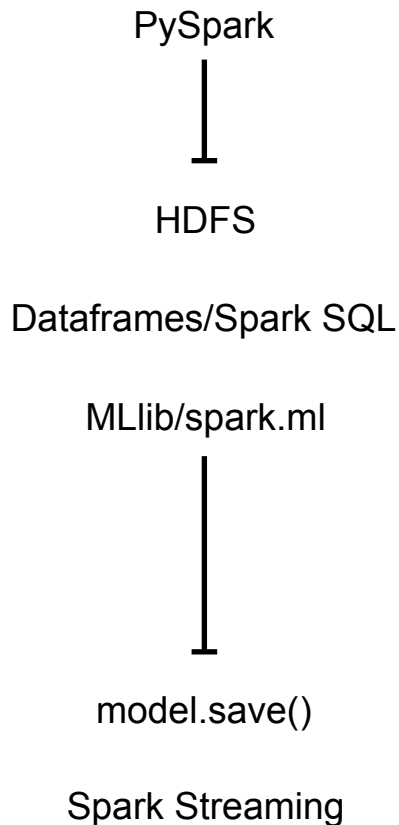


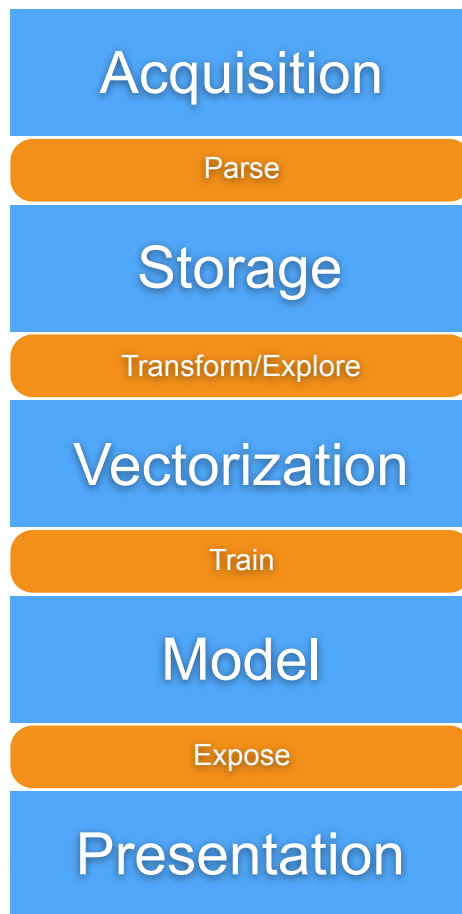
Lesson 5: Advanced Applications

5.1 Machine Learning on Spark: MLlib and `spark.ml`

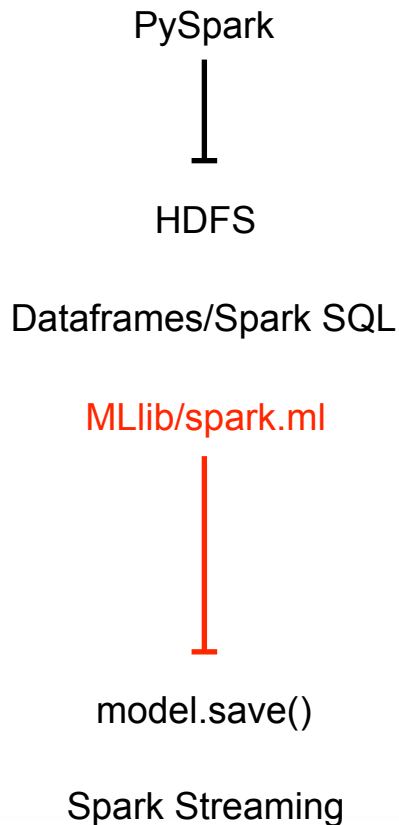
At Scale



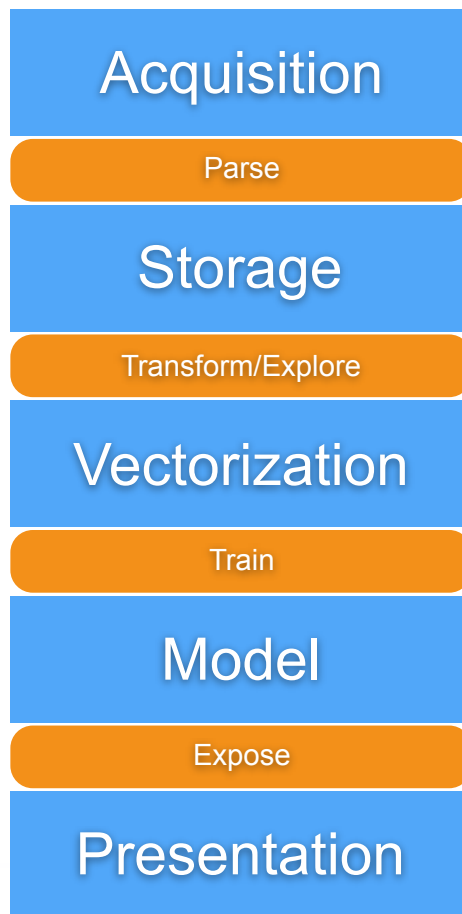
Data Pipeline

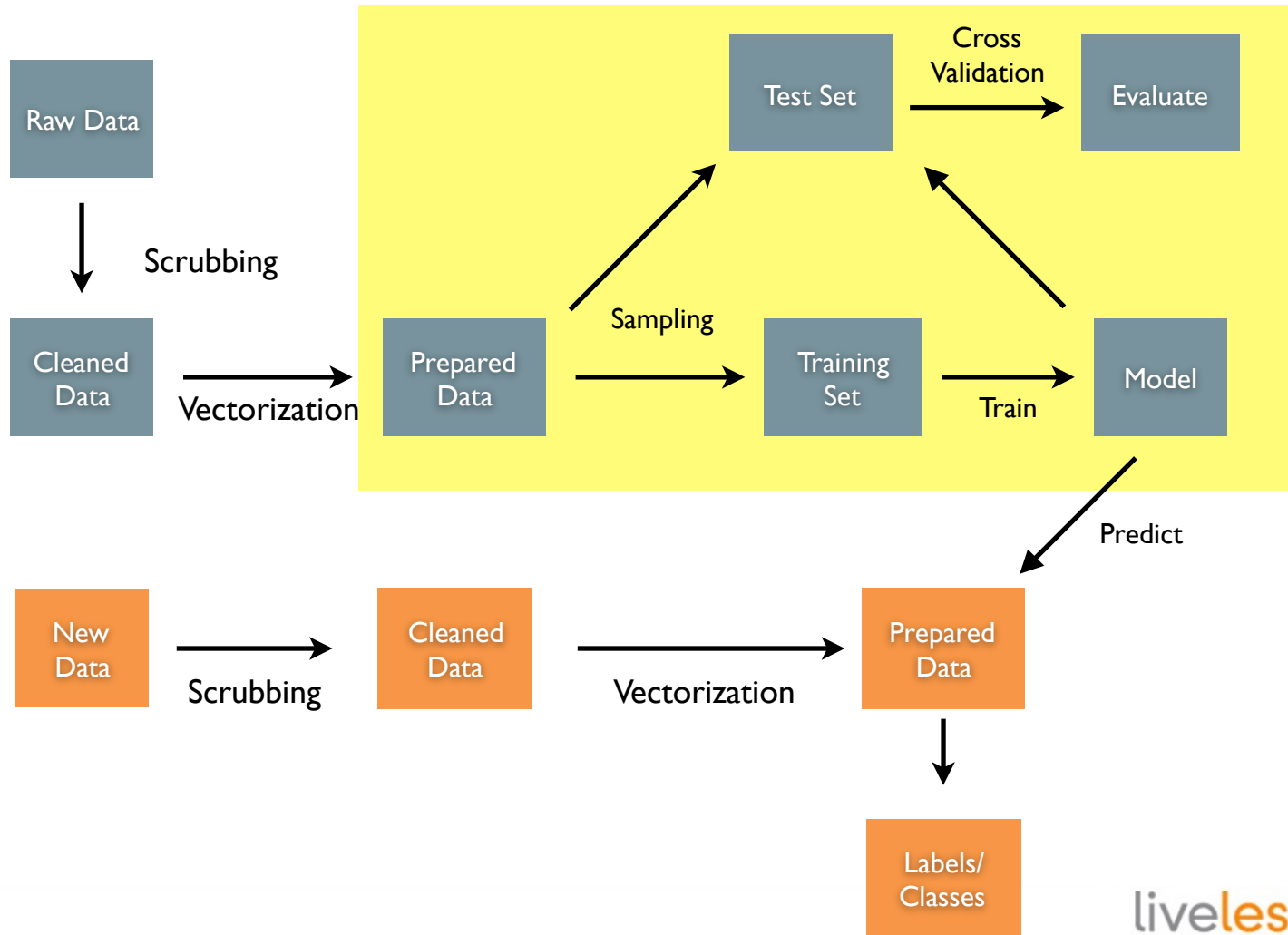


At Scale



Data Pipeline





Unified Platform

Spark SQL

Spark
Streaming

PySpark
SparkR

MLlib
spark.ml

GraphX

Spark Core

Stand Alone Scheduler

YARN

Mesos



Unified Platform

Statistics

Feature
Engineering

Recommendation

Classification/
Regression

Tuning/
Evaluation

spark.ml + MLlib



<i>MLlib</i>	<code>spark.ml</code>
<i>original primary ML API</i>	<i>Higher level API to construct ML workflows/ pipelines</i>
<i>Library</i>	<i>Framework</i>
<i>Statistics, Classification/Regression, Clustering, Feature Extraction, etc.</i>	<i>Dataset, Transformer, Estimator, Pipeline, Parameter</i>



MLlib

- Data types
- Basic Statistics
- Classification and Regression
- Collaborative Filtering
- Clustering
- Dimensionality Reduction
- Feature Extraction
- Optimization
- Model Export

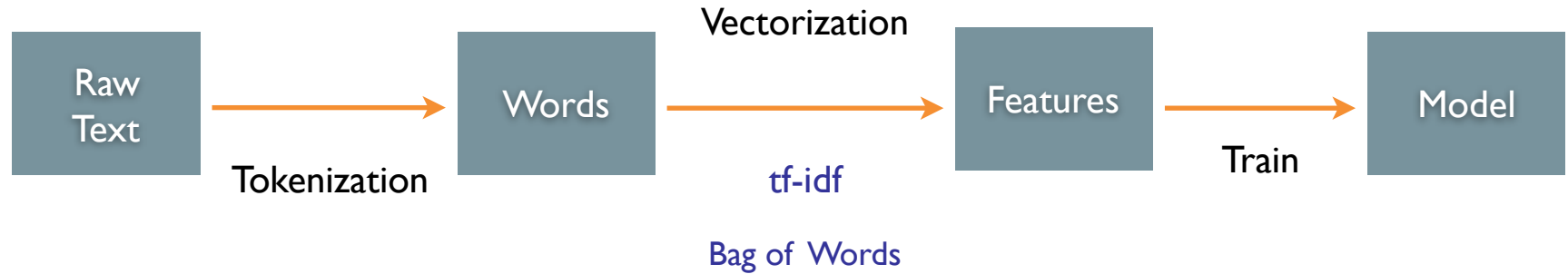


MLlib

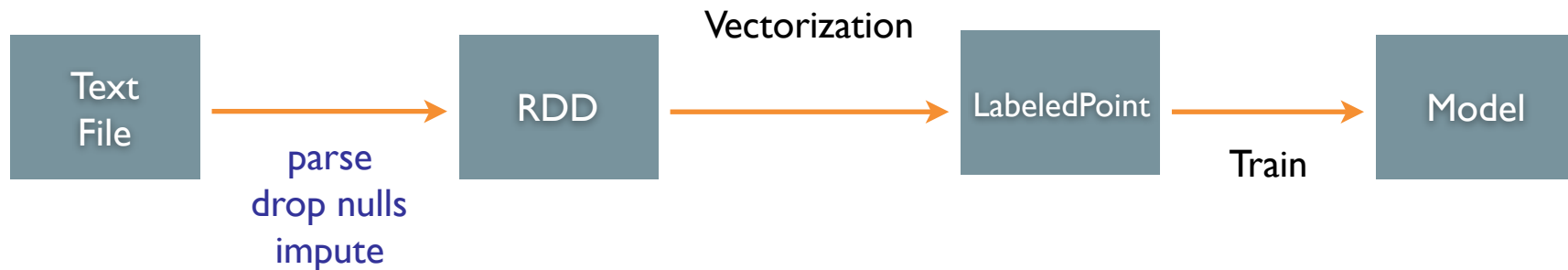
- Data types
- Basic Statistics
- Classification and Regression
- Collaborative Filtering
- Clustering
- Dimensionality Reduction
- Feature Extraction
- Optimization
- Model Export



Machine Learning Pipeline



Machine Learning Pipeline



2009 KDD Cup: Customer Relation Prediction

- French telecom company Orange
- Optimize customer relationship:
 - **Churn**: switching providers
 - **Appetency**: buying new products and services
 - **Up-selling**: buying upgrades/add-ons for current sale



Challenges

- “Very large dataset”
- Heterogeneous (continuous and discrete) noisy data
- Imbalanced classes
- (no context on data — anonymized)

