



# CatBoost for Apache Spark Architecture

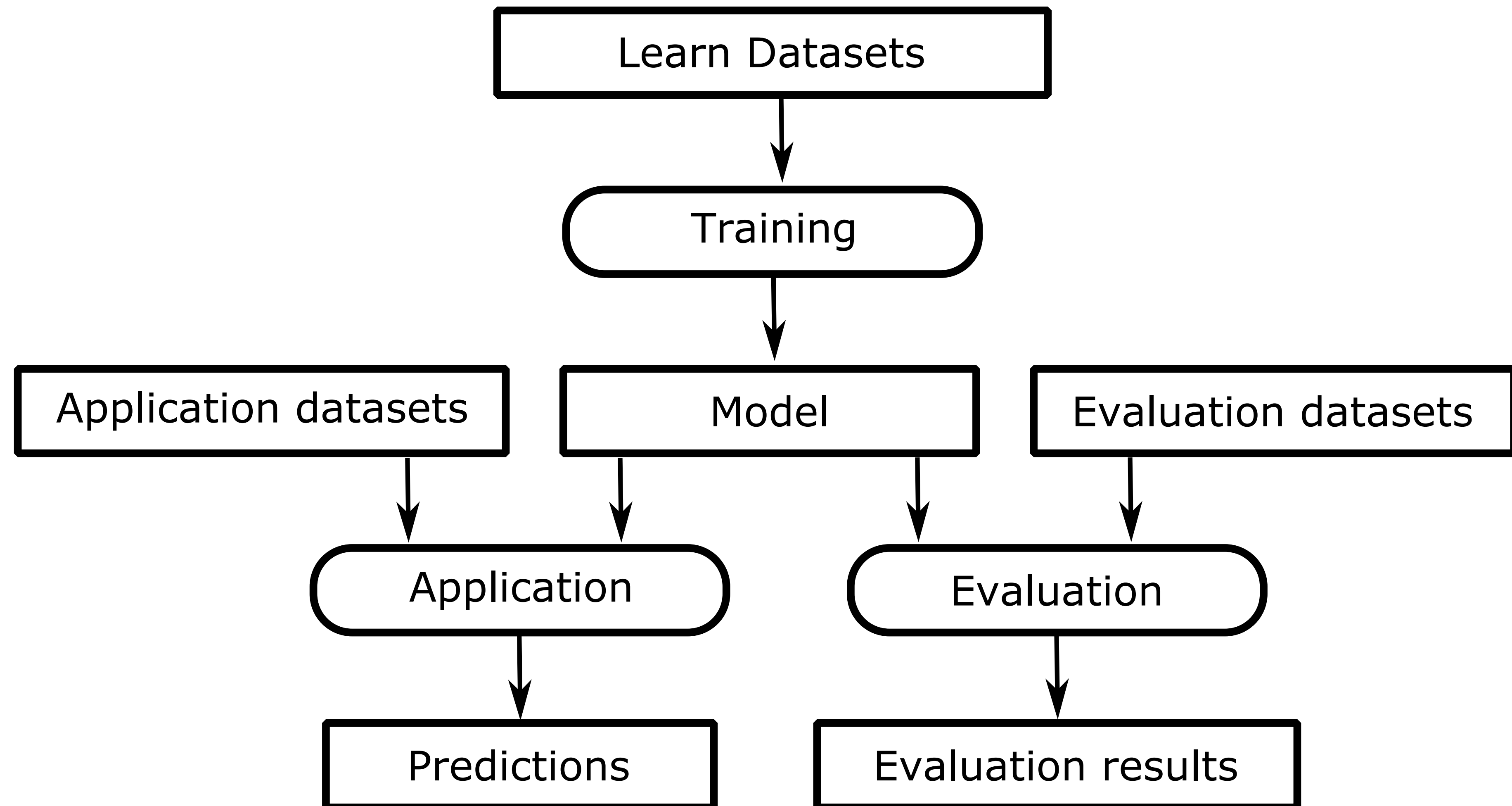
Andrei Khropov,  
CatBoost team @ Yandex

# What is CatBoost?

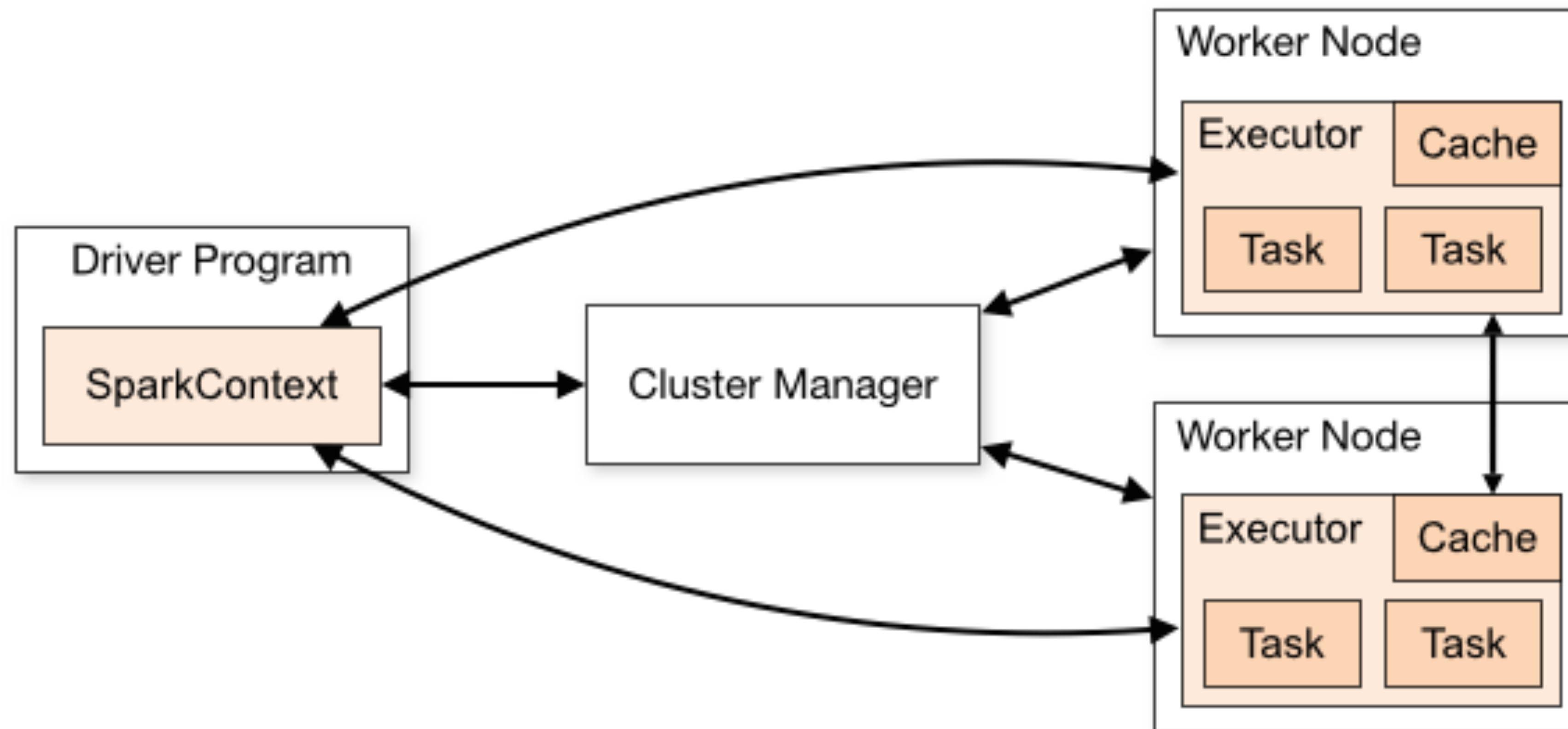


- › CatBoost is a machine learning algorithm that uses gradient boosting on decision trees (GBDT).
- › It is available as an open source library.
- › See “Introducing CatBoost for Apache Spark” presentation for more basic info.

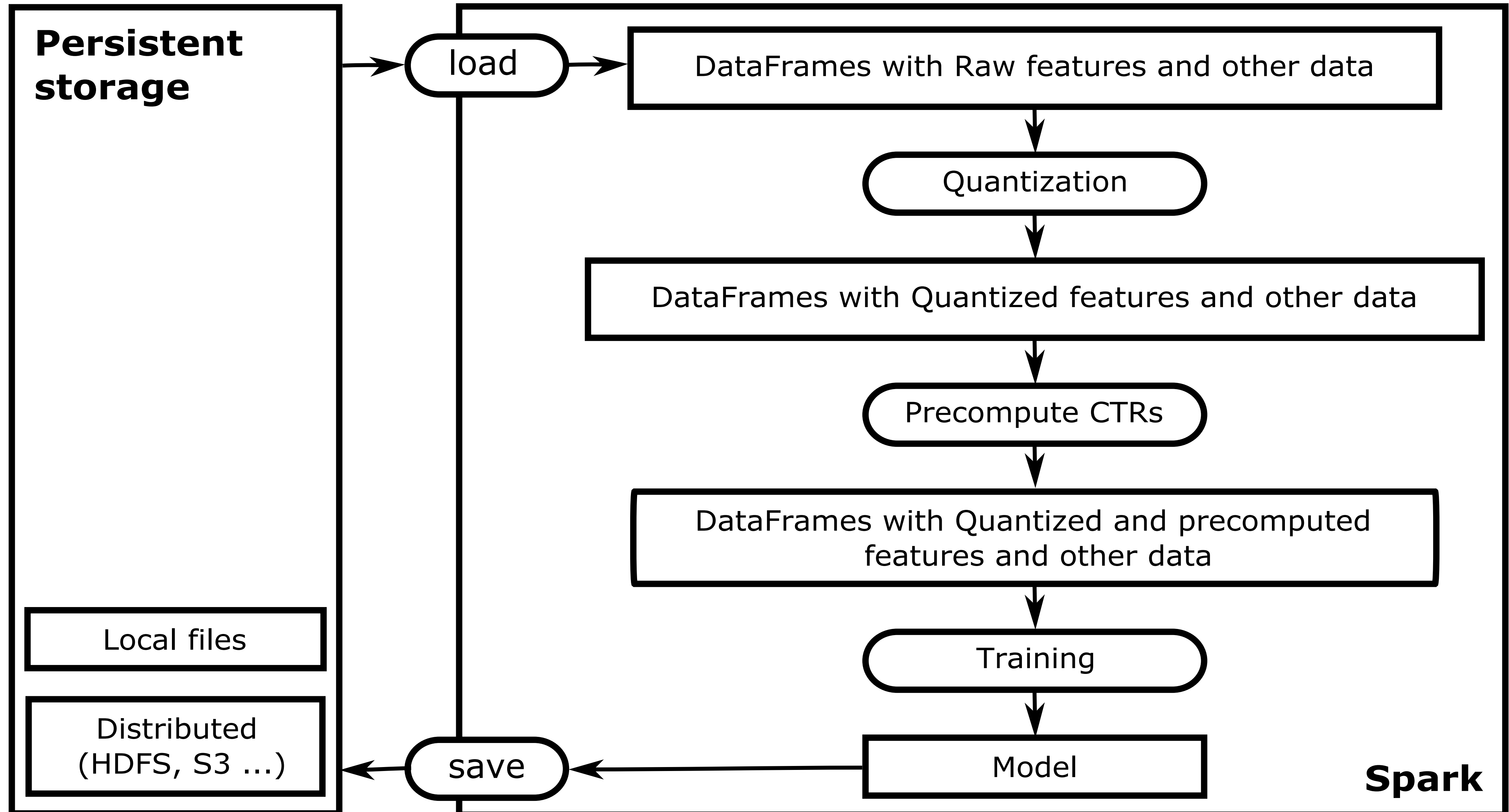
# Generic Machine Learning Workflow



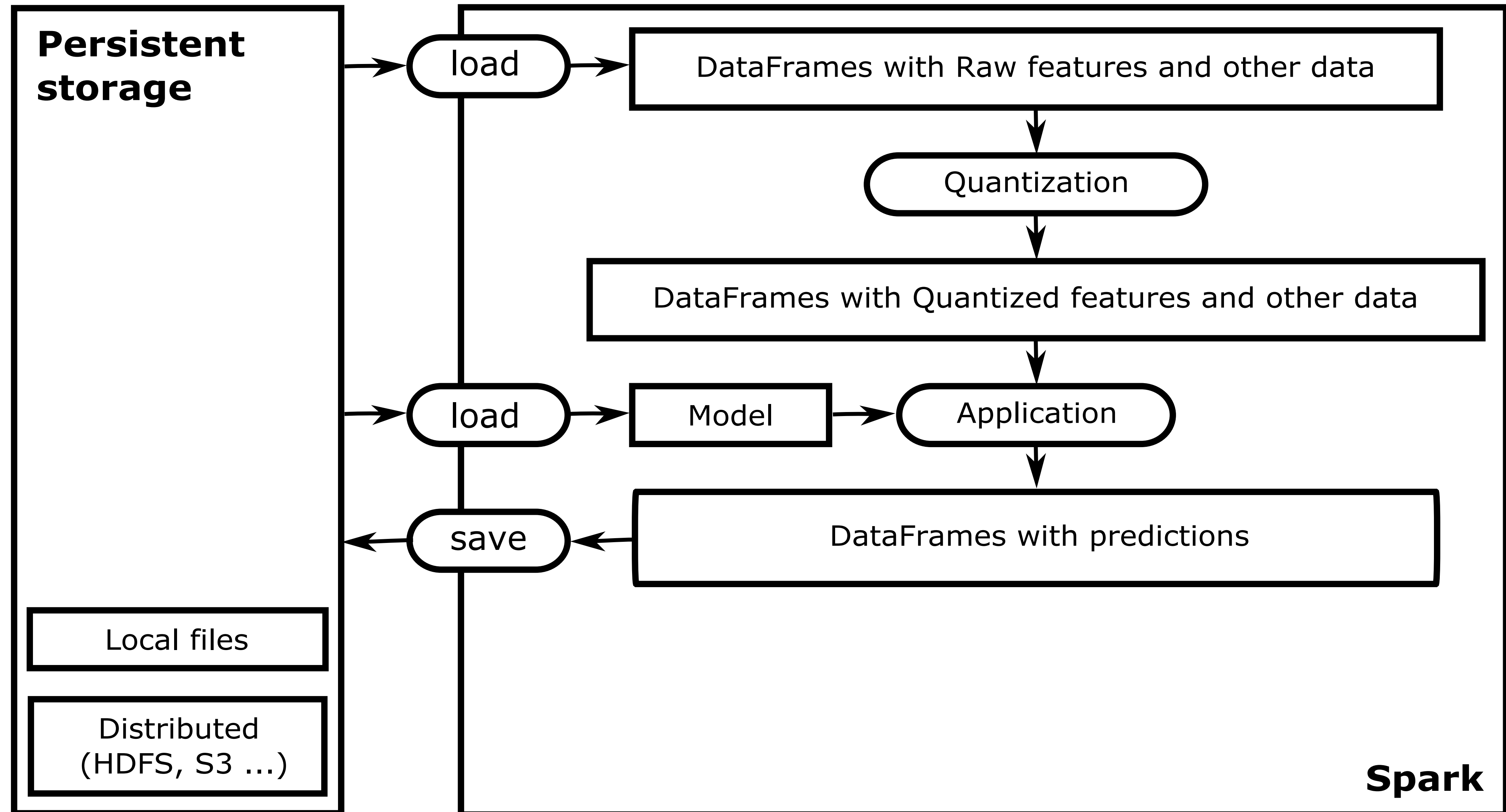
# Spark Cluster Overview



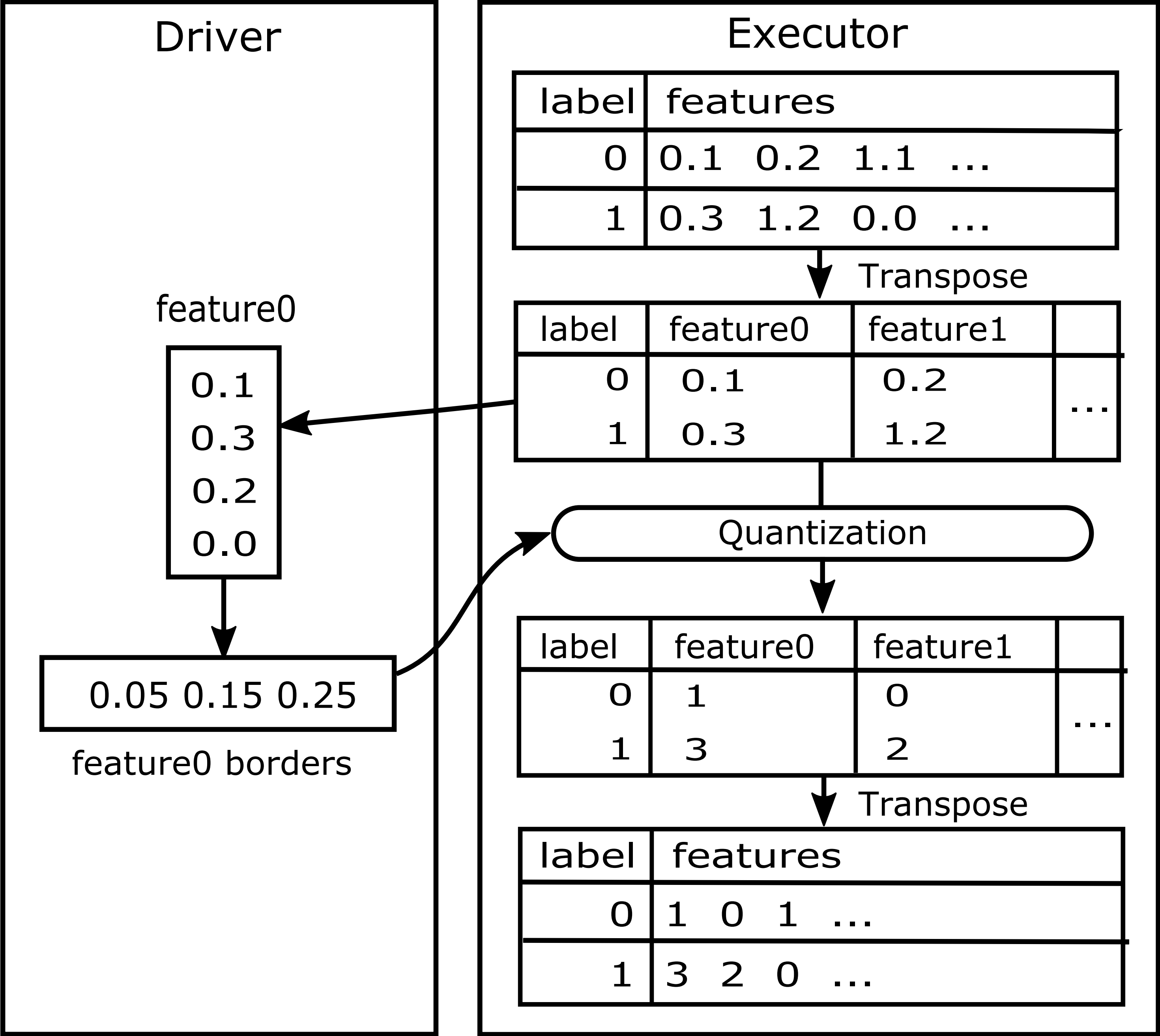
# Training Workflow



# Application Workflow

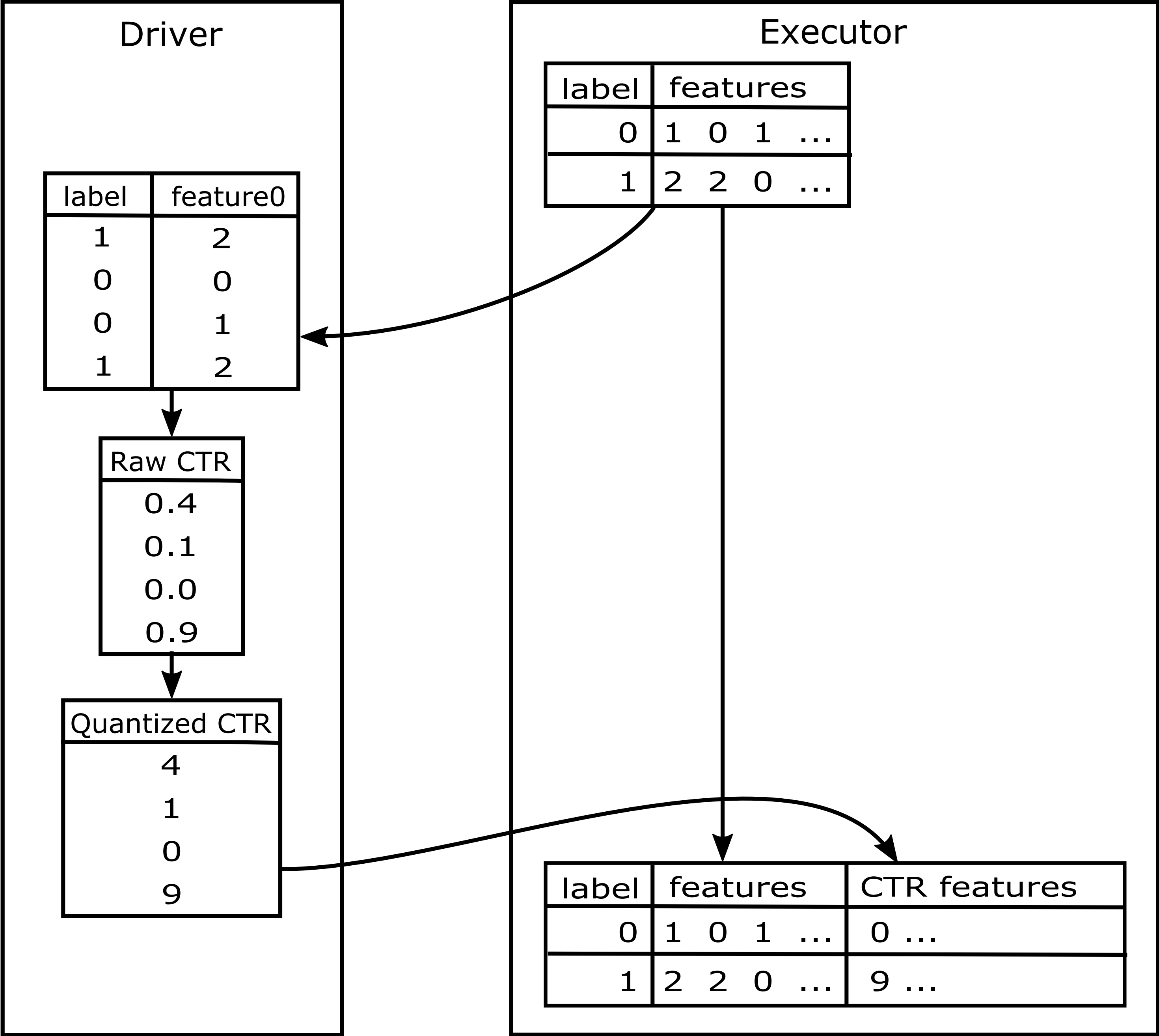


# Quantization



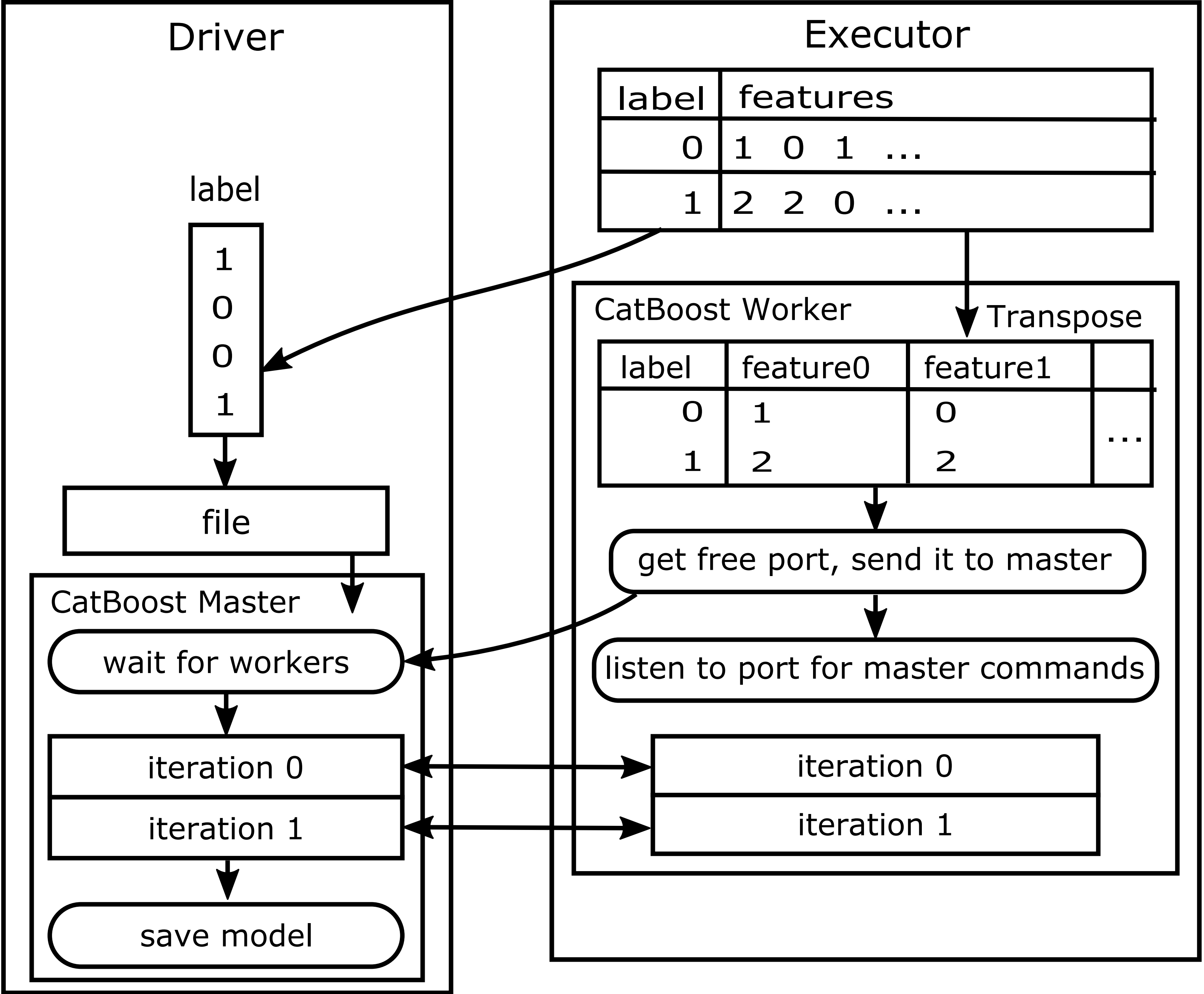
# CTRs

## Precalculation





# Training



# The End



- › CatBoost website: <https://catboost.ai/>
- › CatBoost documentation: <https://catboost.ai/docs>
- › CatBoost on GitHub: <https://github.com/catboost>
- › CatBoost for Apache Spark home:  
<https://github.com/catboost/catboost/tree/master/catboost/spark/catboost4j-spark>