# Hands On : PySparkling Water
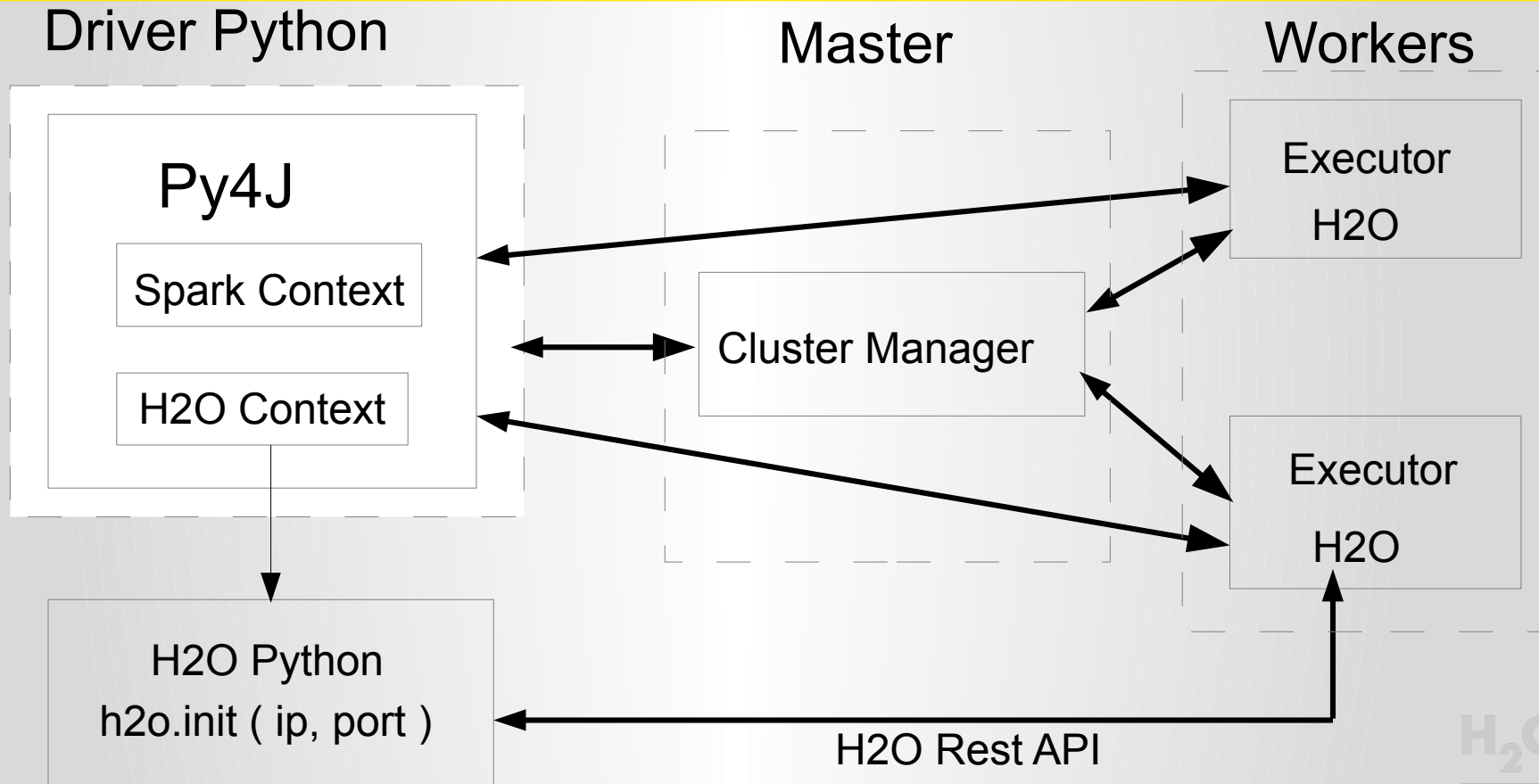
- By Nidhi Mehta

# What is PySparkling Water

PySparkling Water = Python + Spark + $H_2O$

Python + Sparkling Water

H$_2$O WORLD

# PySparkling Architecture

# Demo Workflow

Aim: Build a model to predict Arrest for Chicago crime dataset

- Import Chicago Crime Dataset

- Combine Crime data with Census and Weather data.

- Build a model to predict whether an arrest was made

- Predict on a test dataset

H2O WORLD

# Pre Requisites to run the demo

- Install Spark-1.5.1

- Install and Build Sparkling Water-1.5.6

( ./gradlew build -x check )

- Install H2O-3.6.0.3

- Install H2O-python

( sudo pip install h2o-3.6.0.3-py2.py3-none-any.whl )

# Command to Start/Access PySparking Water Cluster

1)

Set spark environment by specifying SPARK_HOME

export SPARK_HOME = Path_to_Spark_dir

2)

- To run from Python notebook-

IPYTHON_OPTS="notebook" Path_to_Sparkling_dir/bin/pysparkling

- To run from regular Python shell

Path_to_Sparkling_dir/bin/pysparkling

H<sub>2</sub>O WORLD

# Let's Run the Demo!

# Why use PySparkling

- Automatic Parallelization and less lines of code

- Much Faster on big data - uses H2O's rest API calls to connect to H2O Cluster
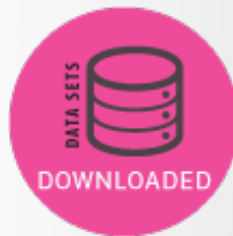
H₂O

WORLD

# Thank You

# What do these stickers mean?

I have Sparkling Water Installed

**H2O** INSTALLED

I have H2O installed

I have Python installed

INSTALLED

DATA SETS DOWNLOADED

I have the H2O World data sets

**Pick up stickers or get install help at the information booth**

H2O WORLD