# INTEL SOFTWARE

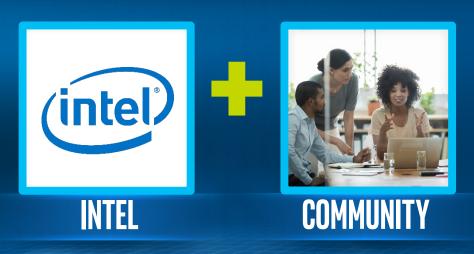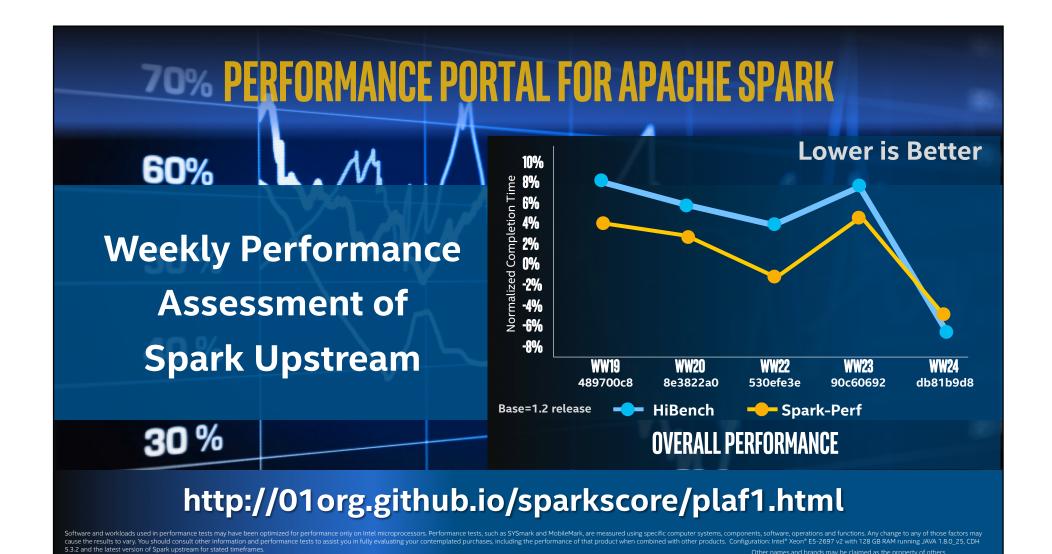When everything computes and connects, Intel software relates.

| Connected Devices | Connective Fabric | Cloud and Data Center |
|---|---|---|



**Ecosystem Enabling**

OPEN INTERCONNECT CONSORTIUM · intel Intel Developer Zone · intel Services · python · Java · OpenJDK · openstack CLOUD SOFTWARE · hadoop · Spark

**OS Enabling**

Android · Chrome · Windows · Apple · Linux

**YOUKU 优酷**

**Data Processing Time Reduced**

**94%**

**Computational Time Reduced**

**92%**

Spark

intel inside™
Xeon®

intel
Math Kernel
Library

**Spark** · **intel® inside™ Xeon®**

**TENCENT 腾讯**

**10X**
MODEL SIZE

**4X**
SPEEDUP

*"With Intel optimized, distributed machine learning on Spark, we were able to scale our model size by >10X, while reduce the training time by ~4x; it helps us to run the distributed machine learning on Xeon clusters and provide better service experience"*

**Hyton Deng**
*Director of Data Infrastructure*

"It sometimes took *weeks to write* new business applications. It was a cumbersome process *involving two different teams* to analyze real time streaming data."

**Xiaohui, Liao**
Senior Development Manager

JD.COM 京东.COM

STRUCTURED QUERIES

ANNOUNCING:

*STREAMING SQL*

FOR APACHE SPARK

STREAM PROCESSING

## Spark Streaming

```
val weblog= KafkaUtils.createStream(…)
 .map(_._2)
 .flapMap(_.split("|")).map(_._2, _)

val category=
  sc.textFile(…).flatMap(_.split("|"))

val result
  = weblog.transform(r =>
r.join(category))
    .map(_._2._2, 1L))
    .reduceByKey(_ + _)

result.print()
```

## Streaming SQL

```
CREATE TABLE weblog ( sourcep_ip STRING,
cookie_id INT, item_id INT …) USING
stream.source.KafkaSource OPTIONS( zkQuorum
"localhost:2181", topic "weblog:1"…)

CREATE TABLE category (item_id INT,
category STRING);

SELECT category, COUNT(*) FROM weblog JOIN
category ON category.item_id =
weblog.item_id GROUP BY category.category;
```

## 48% of developers use SQL
*Source: StackOverflow.com 2015 survey

JOIN US AND CONTRIBUTE!

https://github.com/Intel-bigdata/spark-streamingsql