

# UNLEASHING DATA INTELLIGENCE WITH INTEL AND APACHE SPARK\*

AI: UNLEASHING THE NEXT WAVE

**MICHAEL GREENE**

VICE PRESIDENT INTEL  
SOFTWARE & SERVICES GROUP  
@greene1of5

# OUR JOURNEY WITH SPARK COMMUNITY SINCE 2015

**TODAY**

Stay Tuned\*\*\*



**INTEL+DATABRICKS+AMPLABS  
COLLABORATION ANNOUNCED**

Streaming SQL Open Sourced

**SCALABLE, HIGH PERFORMANCE APACHE  
SPARK\* ON IA**

WebScale ML Open Sourced

**BRINGING DEEP LEARNING TO APACHE SPARK\***

BigDL Open Sourced

# CONTRIBUTION TO APACHE SPARK

## EXAMPLES

### PERFORMANCE & SECURITY

**4.3X**

MLlib\* with Intel®  
Math Kernel  
Library

**1.28X**

Spark Shuffle File  
Encryption

**1.35X**

Spark\* Shuffle RPC  
Encryption

### SCALABILITY

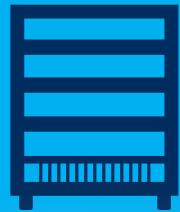
**>10X**

Scalability Improvement  
For Customer Analysis  
Using Word2vec

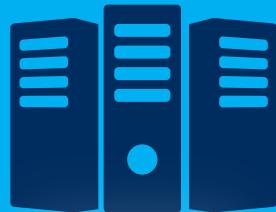
**>70X**

Scalability Improvement For  
Topic Modeling Using Latent  
DIRICHLET ALLOCATION

# THE NEXT BIG WAVE



MAINFRAMES



STANDARDS-  
BASED SERVERS



CLOUD  
COMPUTING



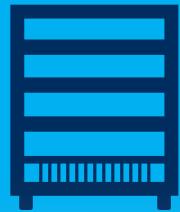
DATA DELUGE  
COMPUTE BREAKTHROUGH  
INNOVATION SURGE

ARTIFICIAL  
INTELLIGENCE

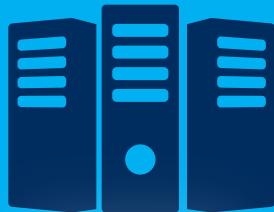
---

AI COMPUTE CYCLES WILL GROW **12X** BY 2020

# THE NEXT BIG WAVE



MAINFRAMES



STANDARDS-  
BASED SERVERS



CLOUD  
COMPUTING



DATA DELUGE  
COMPUTE BREAKTHROUGH  
INNOVATION SURGE

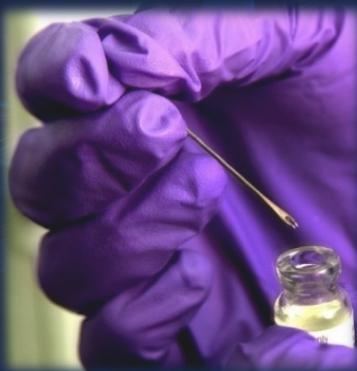
ARTIFICIAL  
INTELLIGENCE

---

AI COMPUTE CYCLES WILL GROW **12X** BY 2020

# AI WILL USHER IN A BETTER WORLD

ON THE SCALE OF THE AGRICULTURAL, INDUSTRIAL AND DIGITAL REVOLUTIONS



## ACCELERATE

Large scale solutions

- Cure Diseases
- Prevent Crime
- Unlock Dark Data

## UNLEASH

Scientific Discovery

- Explore New Worlds
- Decode the Brain
- Uncover New Theories

## EXTEND

Human Capabilities

- Personalize Learning
- Enhance Decisions
- Optimize Time

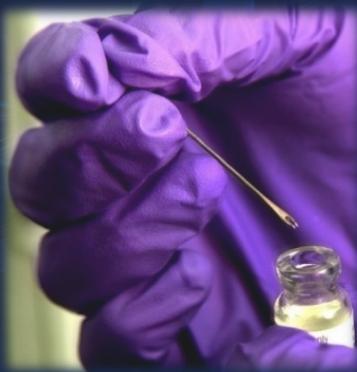
## AUTOMATE

Undesirable Tasks

- Automate Driving
- Save Lives in Danger
- Perform Chores

# AI WILL USHER IN A BETTER WORLD

ON THE SCALE OF THE AGRICULTURAL, INDUSTRIAL AND DIGITAL REVOLUTIONS



## ACCELERATE

Large scale solutions

- Cure Diseases
- Prevent Crime
- Unlock Dark Data

## UNLEASH

Scientific Discovery

- Explore New Worlds
- Decode the Brain
- Uncover New Theories

## EXTEND

Human Capabilities

- Personalize Learning
- Enhance Decisions
- Optimize Time

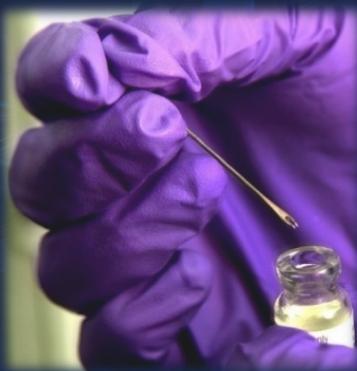
## AUTOMATE

Undesirable Tasks

- Automate Driving
- Save Lives in Danger
- Perform Chores

# AI WILL USHER IN A BETTER WORLD

ON THE SCALE OF THE AGRICULTURAL, INDUSTRIAL AND DIGITAL REVOLUTIONS



## ACCELERATE

Large scale solutions

- Cure Diseases
- Prevent Crime
- Unlock Dark Data

## UNLEASH

Scientific Discovery

- Explore New Worlds
- Decode the Brain
- Uncover New Theories

## EXTEND

Human Capabilities

- Personalize Learning
- Enhance Decisions
- Optimize Time

## AUTOMATE

Undesirable Tasks

- Automate Driving
- Save Lives in Danger
- Perform Chores

# INTEL® AI PORTFOLIO

## EXPERIENCES



## TOOLKITS

Intel® DL  
Training &  
Deployment

Intel® Nervana™  
DL Software &  
Cloud

Intel®  
Computer  
Vision SDK

Intel® GO™  
Automotive  
SDK

Movidius  
Fathom

## FRAMEWORKS



## LIBRARIES



Intel Distribution

Intel® DAAL

Intel® Nervana™ Graph\*  
Intel® MKL MKL-DNN Intel® MLSL

## HARDWARE



Compute



\*



Memory/Storage



Networking



Computer Vision

\*Future

END  
TO  
END  
AI

## AI ON INTEL: UNLEASHING THE NEXT WAVE

# BIG DATA BROUGHT AI TO ENTERPRISE

**Increasing trend of AI workloads  
(ML/DL) in data center**

**Spark has emerged as The Big Data Analytics OS in data center and cloud**

**Spark is evolving to meet the Enterprise Needs - Intel continues to be part of that Journey ...**



# BIGDL

BRINGING DEEP LEARNING TO BIG DATA

## DATA FRAME



## SPARK CORE



Released to open source Dec 30<sup>th</sup> 2016



**BIGDL VIDEO**

# BIGDL: WHAT'S NEW?

## Features Released EOQ1:

- ✓ Python Support
- ✓ Notebook Integration
- ✓ Tensorboard Support
- ✓ Better RNN Support
- ✓ Improved Robustness

## Coming Out EOQ2:

- ✓ Functional API support
- ✓ Tensorflow Model Read/Write
- ✓ Recursive Net Support
- ✓ 3D Convolutions
- ✓ Python 3.5 Support
- ✓ Spark 2.1 Support

Developer Centric Zone: [software.intel.com/bigdl](http://software.intel.com/bigdl)

# BIGDL: WHAT'S NEW?

## Features Released EOQ1:

- ✓ Python Support
- ✓ Notebook Integration
- ✓ Tensorboard Support
- ✓ Better RNN Support
- ✓ Improved Robustness

## Coming Out EOQ2:

- ✓ Functional API support
- ✓ Tensorflow Model Read/Write
- ✓ Recursive Net Support
- ✓ 3D Convolutions
- ✓ Python 3.5 Support
- ✓ Spark 2.1 Support

Developer Centric Zone: [software.intel.com/bigdl](http://software.intel.com/bigdl)

# BIGDL: WHAT'S NEW?

## Features Released EOQ1:

- ✓ Python Support
- ✓ Notebook Integration
- ✓ Tensorboard Support
- ✓ Better RNN Support
- ✓ Improved Robustness

## Coming Out EOQ2:

- ✓ Functional API support
- ✓ Tensorflow Model Read/Write
- ✓ Recursive Net Support
- ✓ 3D Convolutions
- ✓ Python 3.5 Support
- ✓ Spark 2.1 Support

Developer Centric Zone: [software.intel.com/bigdl](http://software.intel.com/bigdl)

# BIGDL: WHAT'S NEW?

## Features Released EOQ1:

- ✓ Python Support
- ✓ Notebook Integration
- ✓ Tensorboard Support
- ✓ Better RNN Support
- ✓ Improved Robustness

## Coming Out EOQ2:

- ✓ Functional API support
- ✓ Tensorflow Model Read/Write
- ✓ Recursive Net Support
- ✓ 3D Convolutions
- ✓ Python 3.5 Support
- ✓ Spark 2.1 Support

Developer Centric Zone: [software.intel.com/bigdl](http://software.intel.com/bigdl)

# BIGDL READY FOR WIDE ADOPTION

cloudera®

databricks



Alibaba Cloud  
aliyun.com

amazon  
web services™

**CRAY**  
THE SUPERCOMPUTER COMPANY



Microsoft  
Azure



Lightbend

UnionPay  
银联



# BIGDL READY FOR WIDE ADOPTION

cloudera®

**“The integration of BigDL with Cloudera Data Science Workbench allows organizations to leverage deep learning libraries on CPU architecture and an easy way to create native Spark data science pipelines and integrate them with BigDL and other Spark/Hadoop components.”**

--Charles Zedlewski, Sr. Vice President, Products, Cloudera

# BIGDL READY FOR WIDE ADOPTION



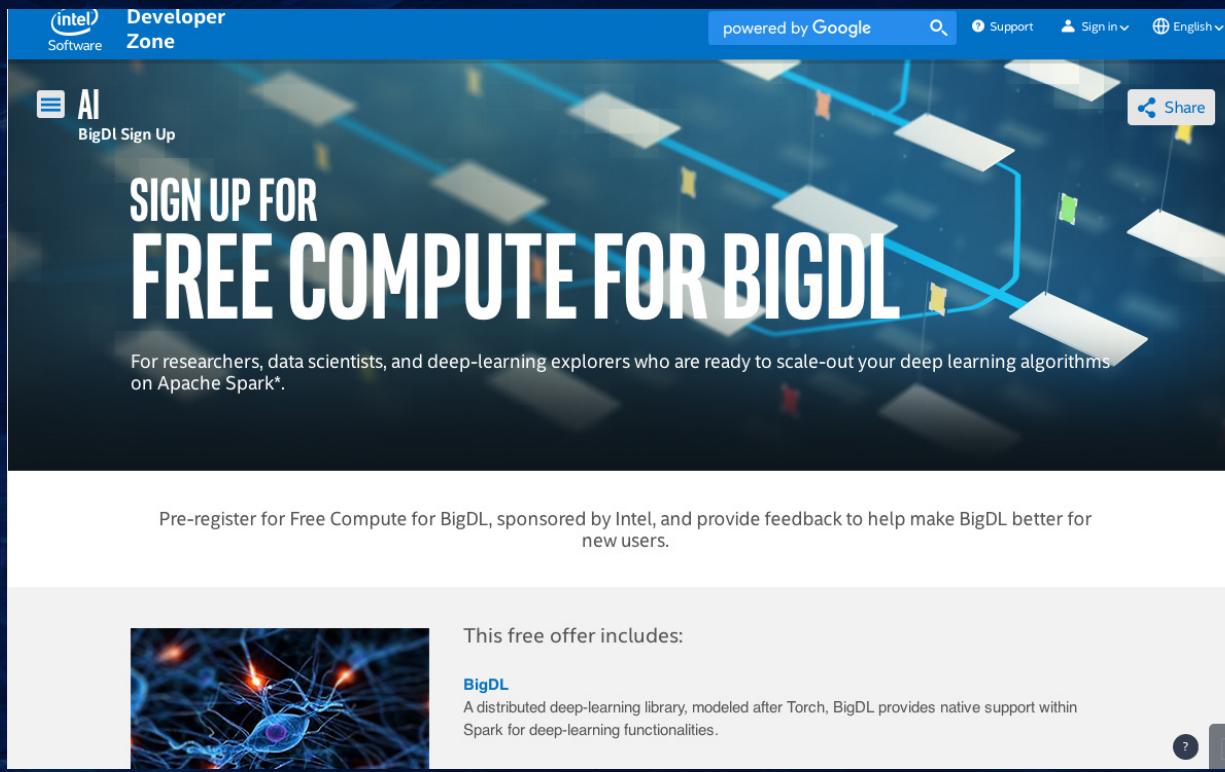
Lightbend

**“BigDL is featured in our sample app, and it's our preferred library to date for DL.”**

--Mark Brewer, President & CEO, Lightbend, Inc

# FREE COMPUTE FOR BIGDL

Visit <https://software.intel.com/bigdlcompute>

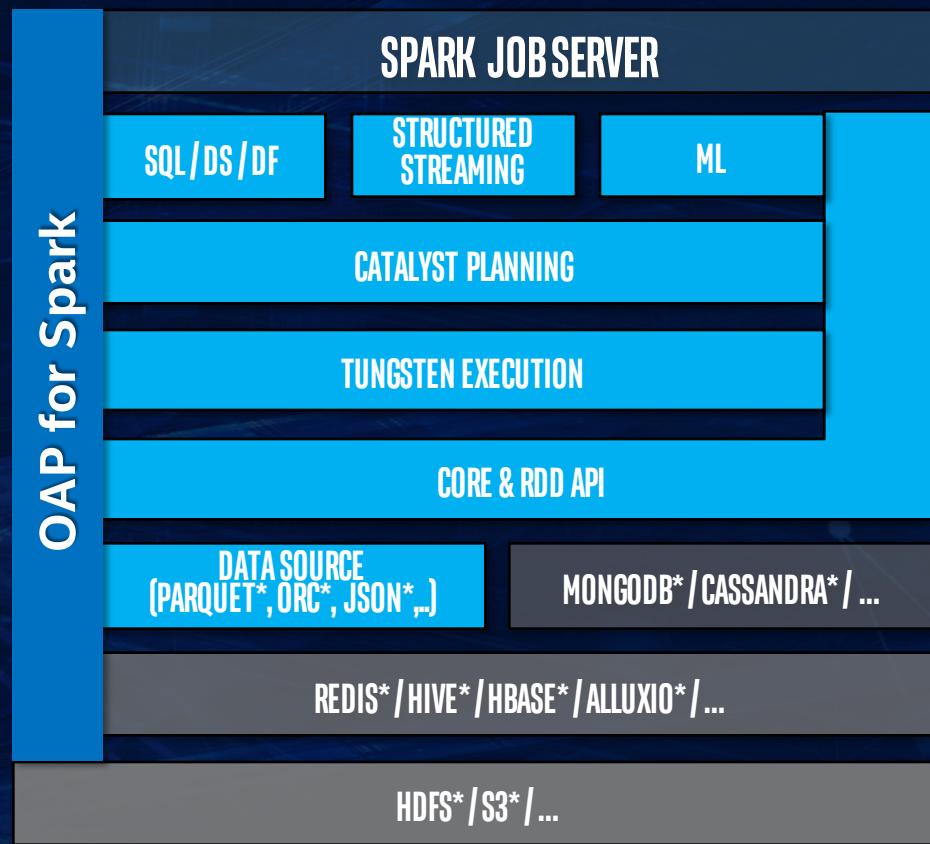


The screenshot shows the Intel Software Developer Zone homepage. At the top, there's a navigation bar with the Intel logo, "Developer Zone", "powered by Google", a search icon, "Support", "Sign in", and "English". Below the header, a large banner features a blue-toned abstract background with floating white rectangular shapes and a network of blue lines. The text "SIGN UP FOR FREE COMPUTE FOR BIGDL" is prominently displayed in white. Below the banner, a sub-headline reads: "For researchers, data scientists, and deep-learning explorers who are ready to scale-out your deep learning algorithms on Apache Spark\*." A "BigDL Sign Up" button is visible. The main content area below the banner contains a paragraph about pre-registering for free compute, followed by a section titled "This free offer includes:" which lists "BigDL" with a brief description. There are also small icons for a question mark and a zero count.

TURN YOUR IDEA  
INTO REALITY

# INTRODUCING OAP FOR SPARK\*

Optimized Analytics Package for Spark\* Platform – Accelerating Spark Queries!



<https://github.com/Intel-bigdata/OAP>

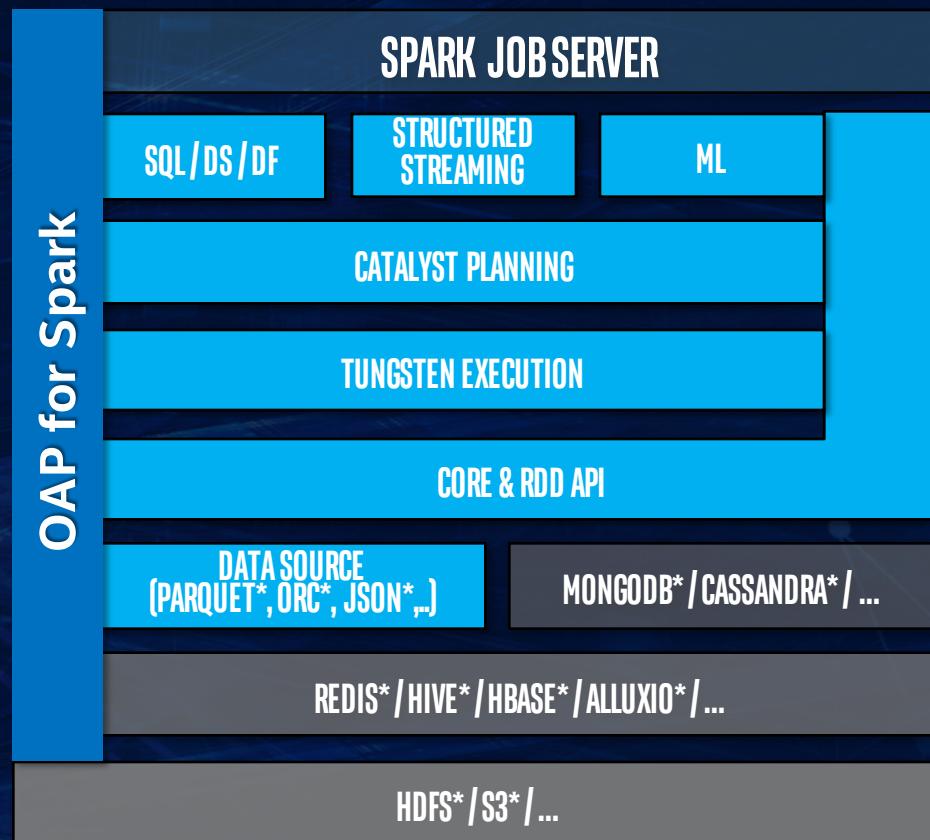
**“OAP for Spark is quite fit for Baidu’s data analytics requirements, and brings 1.5X-5X performance gain for ad-hoc query. We’d like to dive into the OAP open source community with Intel for more significant acceleration in the future releases, to unleash the power of new hardware platforms.”**

--Lin Xiaodong, Director of Baidu Infrastructure Department



# INTRODUCING OAP FOR SPARK\*

Optimized Analytics Package for Spark\* Platform – Accelerating Spark Queries!



<https://github.com/Intel-bigdata/OAP>

**“OAP for Spark is quite fit for Baidu’s data analytics requirements, and brings 1.5X-5X performance gain for ad-hoc query. We’d like to dive into the OAP open source community with Intel for more significant acceleration in the future releases, to unleash the power of new hardware platforms.”**

--Lin Xiaodong, Director of Baidu Infrastructure Department





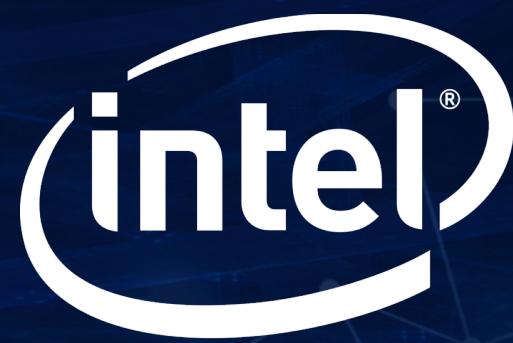
**BETTER TOGETHER  
LET'S COLLABORATE!**

[SOFTWARE.INTEL.COM/BIGDL](http://SOFTWARE.INTEL.COM/BIGDL)

[SOFTWARE.INTEL.COM/AI](http://SOFTWARE.INTEL.COM/AI)

INTEL BOOTH: 301

WOMEN IN BIG DATA LUNCH/PANEL:  
JUNE 7 IN ROOM 2014



experience  
what's inside™

# LEGAL DISCLAIMERS

- Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Learn more at intel.com, or from the OEM or retailer.
- No computer system can be absolutely secure.
- Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. Consult other sources of information to evaluate performance as you consider your purchase. For more complete information about performance and benchmark results, visit <http://www.intel.com/performance>

## Configurations:

- 4.3X for Spark MLlib thru Intel Math Kernel Library (MKL)
  - Spark-Perf (same for before and after): 9 nodes each with Intel® Xeon® processor E5-2697A v4 @ 2.60GHz \* 2 (16 cores, 32 threads); 256 GB ; 10x SSDs; 10Gbps NIC
- 19x for HDFS Erasure Coding in micro workload (RawErasureCoderBenchmark) and 1.25x in Terasort, plus 50+% storage capacity saving and higher failure tolerance level.
  - RawErasureCoderBenchmark (same for before and after): single node with Intel® Xeon® processor E5-2699 v4 @ 2.20GHz \*2 (22 cores, 44 threads); 256GB; 8x HDDs; 10Gbps NIC
  - Terasort (same for before and after): 10 nodes each with Intel® Xeon® processor E5-2699 v4 @ 2.20GHz \*2 (22 cores, 44 threads); 256GB; 8x HDDs; 10Gbps NIC
- 5.6x for HBase off heap read in micro workload (PE) and 1.3x in real Alibaba production workload
  - PE (same for before and after): Intel® Xeon® Processor X5670 @ 2.93Hz \*2 (6 cores, 12 threads); RAM: 150 GB; 1Gbps NIC
  - Alibaba (same for before and after): 400 nodes cluster with Intel® Xeon® processors
- 1.22x Spark Shuffle File Encryption performance for TeraSort and 1.28x for BigBench
  - Terasort (same for before and after): Single node with Intel® Xeon® Processor E5-2699 v3 @ 2.30GHz \*2 (18 cores, 36 threads); 128GB; 4x SSD; 10Gbps NIC
  - BigBench (same for before and after): 6 nodes each with Intel® Xeon® Processor E5-2699 v3 @ 2.30GHz \*2 (18 cores, 36 threads); 256GB; 1x SSD; 8x SATA HDD 3TB, 10Gbps NIC
- 1.35X Spark Shuffle RPC encryption performance for TeraSort and 1.18x for BigBench
  - Terasort (same for before and after): 3 nodes each with Intel® Xeon® Processor E5-2699 v3 @ 2.30GHz \*2 (18 cores, 36 threads); 128GB; 4x SSD; 10Gbps NIC
  - BigBench (same for before and after): 5 nodes, 1x head node: Intel® Xeon® Processor E5-2699 v3 @ 2.30GHz \*2 (18 cores, 36 threads); 384GB; 1x SSD; 8x SATA HDD 3TB, 10Gbps NIC. 4x worker nodes: each with Intel® Xeon® processor E5-2699 v4 @ 2.20GHz \*2 (22 cores, 44 threads); 384GB; 1x SSD; 8x SATA HDD 3TB, 10Gbps NIC.
- 10X scalability for Word2Vec E5-2630v2 \* 2, 128 GB Memory, 12x HDDs; 1000Mb NIC (14 nodes)
- 70X scalability for LDA (Latent Dirichlet Allocation)
  - Intel Xeon E5-2630v2 \* 2, 288GB Memory, SAS Raid5, 10Gb NIC

Intel, the Intel logo, Xeon, Xeon phi, Lake Crest, etc. are trademarks of Intel Corporation in the U.S. and/or other countries.

\*Other names and brands may be claimed as the property of others.

© 2017 Intel Corporation