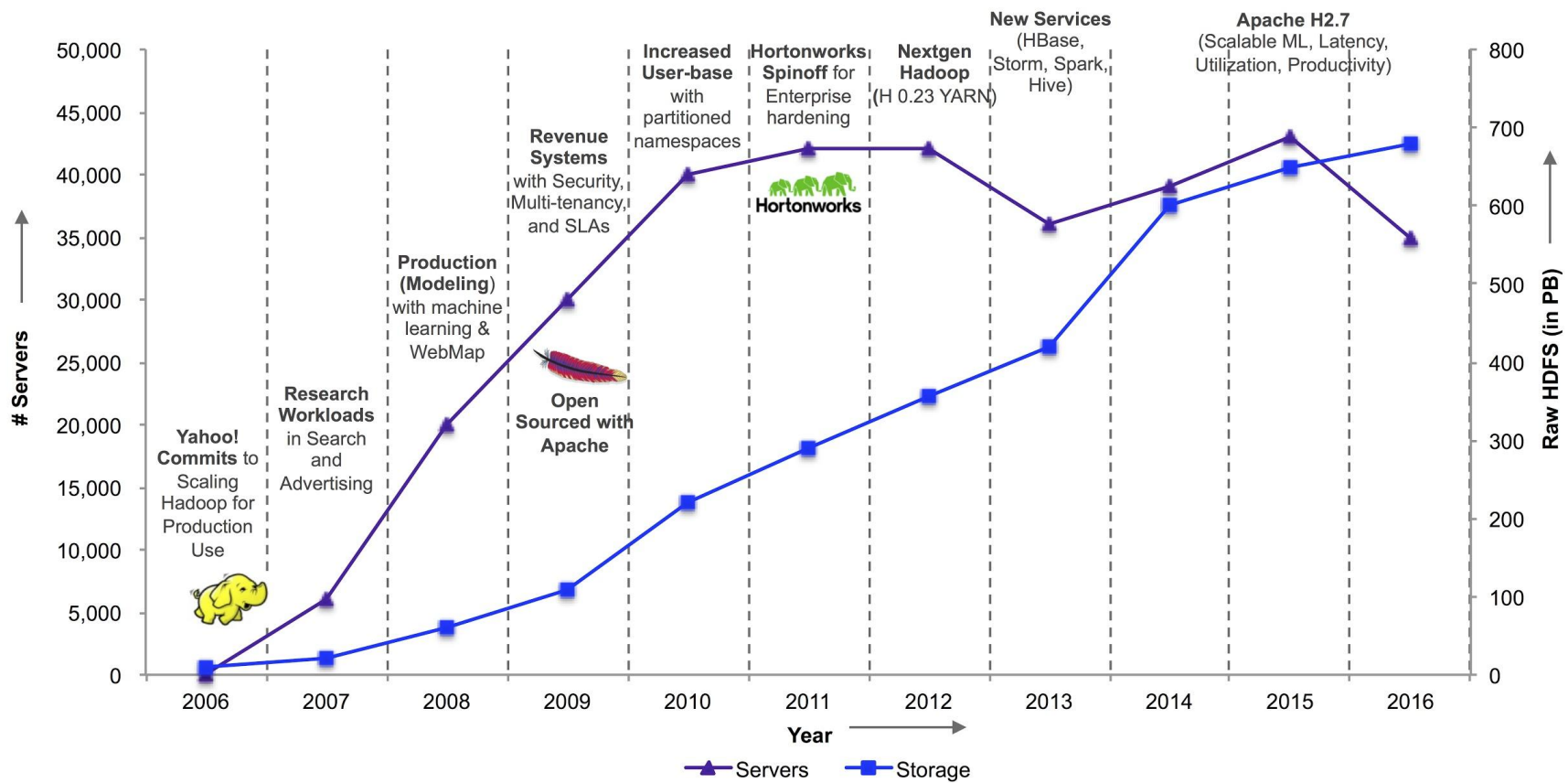# Big-Data Innovation:
## Hadoop, Real-time, and Machine Learning

Andy Feng

Yahoo

# Hadoop at Yahoo

# Real-time Processing at Yahoo

# Big-Data Applications

# Weather App

- ## Beauty
  - › Computational assessed

- ## Relevant
  - › Location
  - › Time
  - › Cloudy
  - › Shower
  - › …

Weather App

Yahoo Weather App

YAHOO!

# Flickr [flickr.com/cameraroll](flickr.com/cameraroll)

# Yahoo Vision Kit



**Image Enrichment Results**

**Dominant Colors**

**Autotags**

outdoor: 0.94 | rock: 0.90 | nature: 0.90 | tower: 0.86 | architecture: 0.86 | lighthouse: 0.86

sea: 0.77 | sunset: 0.74 | serene: 0.72 | hdr: 0.71 | water: 0.71 | landscape: 0.61 | lake: 0.61

waterfront: 0.61 | bay: 0.61 | shore: 0.56 | beach: 0.56 | dusk: 0.55 | sunlight: 0.54

rock formation: 0.53 | lego: 0.51

**Aesthetics**

Meh — Hot damn

**NSFW**

Boring — xxx<sup>xxx</sup>

**Faces**

**Celebrities**

**Similarity**

YAHOO!

# Personalized Content



Content augmentation

User profiling

Recommendation

# Mail Smart Views

# Search



Query intention

Page ranking

Ads matching

# Yahoo Search2Vec: Big Data & Machine Learning
http://bit.ly/28SLdjU

| Bucket Tests | Query Coverage | Auction Depth | Revenue per Search |
|---|---|---|---|
| Simple model vs. Baseline | +1.14% | +2.13% | +7.07% |
| Large model vs. Baseline | +3.61% | +4.57% | +17.12% |

YAHOO!

# Big-Data Innovations

# Technology Stack

| Scalable Database | Scalable Analytics | Scalable Machine Learning | Scalable Deep Learning |
|---|---|---|---|
| Batch Compute: Hadoop | | Iterative Compute: Spark | Real-Time Compute: Storm |
| **CPUs** | | YARN | **GPUs** |
| | | HDFS | |

# Flickr Pipeline: 10B photos, 8M per day

| (1) Feature Engineering @ Scale | (2) Deep Learning @ Scale | (3) Non-deep Learning @ Scale | (4) Apply Model @ Scale | (5) Product Integration @ Scale |
|---|---|---|---|---|

* http://bit.ly/1KIDfof  by Pierre Garrigues, Deep Learning Summit 2015

# System Configuration: Homogenous Clusters

*Network Backplane*

*CPU Servers with JBODs & 10GbE*

*Rack 1*

*Rack 2*

*Rack N*

YAHOO!

# System Configuration: Heterogeneous Clusters
## http://bit.ly/295RNbU

Network Backplane

CPU Servers with JBODs & 10GbE

GPU Servers

100Gbps InfiniBand

Hi-Mem Servers

Rack 1

Rack 2

Rack N

YAHOO!

# Resource Allocation: Early Days



Legend: —— Total   —— Used

Compute Total and Used (TB)

One Month Sample (2015)

**Cluster 1 (2,000 servers)**
*HDFS     12 PB*
*Memory    23 TB*
*Avg. Util:   26%*

**Cluster 2 (3,100 servers)**
*HDFS     21 PB*
*Memory    52 TB*
*Avg. Util:   40%*

**Cluster 3 (5,400 servers)**
*HDFS     36 PB*
*Memory    70 TB*
*Avg. Util:   59%*

YAHOO!

# Resource Allocation: 40% decrease in TCO
## http://bit.ly/295RNbU



Compute Total and Used (TB)

**Consolidated Cluster**
*HDFS     65 PB*
*Memory    240 TB*
*Avg. Util:    70%*

One Month Sample (2016)

Total     Used



**Before**          **After**

**65% increase** in compute capacity

**50% increase** in avg. utilization

18

YAHOO!

# Apache Storm: Real-time Processing
## https://storm.apache.org

# Resource Aware Scheduling: 100% increase utilization

# Apache HBase: NoSQL Database for Hadoop
https://hbase.apache.org



- Several million regions (partitions of tables)
- 800 region servers

YAHOO!

# Support Millions of Regions
http://bit.ly/29hlLFW

**3 Million Regions**

| | Scan Meta | Assignment |
|---|---|---|
| 1 Meta / 1 RS | 56min | 19.79min |
| 32 Meta / 3 RS | 2.91min | 12.56min |

YAHOO!

# HBase Multi-tenancy & Isolation
http://bit.ly/28YK4a7



Cluster Workload

Group A  RS  RS  ···  RS

Group B  RS  RS  ···  RS

Group C  RS  RS  ···  RS

Group D  RS  RS  ···  RS

Group E  RS  RS  ···  RS

Region Server Groups

Isolated Group Workloads

# Apache Omid: Transactions for NoSQL DB
*http://omid.incubator.apache.org*

ACID
Transacti
ons

- *Multi-row/multi-table transactions*

- *Snapshot isolation*

- *Lock-free*

YAHOO!

# Omid: Transactional HBase App

*TransactionManager tm = HBaseTransactionManager.newBuilder().build();*

**Transaction tx = tm.begin();**

*Put row1 = new Put(Bytes.toBytes("EXAMPLE_ROW1"));*
*row1.add(family, qualifier, Bytes.toBytes("VALUE_1"));*
*tt.put(**tx,** row1);*

*Put row2 = new Put(Bytes.toBytes("EXAMPLE_ROW2"));*
*row2.add(family, qualifier, Bytes.toBytes("VALUE_2"));*
*tt.put(**tx,** row2);*

**tm.commit(tx);**

YAHOO!

# Omid: Snapshot Isolation
http://bit.ly/28YZvQF

| TxId | Time Overlap | Spatial Overlap (WriteSet) | | | |
|------|--------------|------|------|------|------|
| T1 | | R1 | R3 | | |
| T2 | | | R2 | | R4 |
| T3 | | | | R3 | R4 |
| T4 | | R1 | R2 | R3 | R4 |

- **T2 overlaps in** time **with** T1 & T3, **but**
  - Spatially T1 ∩ T2 = ∅      ➔ T1 and T2 can both commit.
  - Spatially T2 ∩ T3 = { R4 }      ➔ T2 or T3 will abort.

# Machine Learning: Search2Vec
## http://bit.ly/28SLdjU

**Input: Training Data**

S1: gas_caps gas_cap_replacement_for_car slc_679f037d
gas_door_replacement_for_car slc_466145af1 fuel_door_covers
adid_28540536 slc_348709d7 autozone_auto_parts adid_33183157
auoto_zone slc_8dcdab5d slc_58f979b6

S2: hoka_running shoe_reviews adid_22830711 hoka_shoes_for_bad_feet
hoka_shoes_amazon slc_231gfaw zappos_shoes slc_7c126f71
hoka_walking_shoes

S3: king_tut king_tut_exhibit king_tut_exhibit_seattle_2015 slc_726y6j51
charlies_seattle adid_55774014

**Output:**
**Numeric Vectors**
**Of queries & ads**

Vector("san jose weather") ≈
Vector("weather 95113") ≈
Vector(ad123)

YAHOO!

# Search2Vec: Cosine similarity b/w ads and queries

**adid_243609**

ad vector:

| 0.2 | 1.1 | 7.2 | 0.8 | 3.1 |
|-----|-----|-----|-----|-----|

similarity=0.931

query vector:

| 0.2 | 1.2 | 6.8 | 0.7 | 3.2 |
|-----|-----|-----|-----|-----|

**mystery_games**

| ad metadata | **ad id**: 243609 |
|---|---|
| | **bidterm id**: 341454 |
| | **ad title**: Host a Fun Murder Mystery Party |
| | **ad description**: Huge selection of fun murder mystery games for all ages, groups. #1 site for instant downloads and boxed sets of exciting murder mystery party games |
| generated keywords | murder mystery |
| | murder mystery games |
| | mystery games |
| | murder mystery party |
| | mystery party |
| | free murder mystery games for parties |
| | detective games |
| | murder mystery game |
| | how to host a mystery party for kids |
| | murder mystery dinner |
| | friends game night |
| | murder mystery parties at home |
| | murder mystery dinner party |
| | … |

YAHOO!

# Search2vec: 300M Unique Queries/Ads/Links



- Computation on servers
  - › (1) Negative sampling
  - › (1) Compute gradient: X*Y
  - › (3) Adjust vectors: Y=aX+Y

- Daily training enabled
  - › weeks ➔ hours

**YAHOO!**

# Deep Learning: Previous Practice

# Scalable Deep Learning for Big Data

**Hadoop/Spark Cluster**

Spark application program:
(1) Prepare datasets
(2) DL training & Test
(3) Apply DL model

Hadoop Datasets

YAHOO!

# CaffeOnSpark Open Sourced



github.com/yahoo/CaffeOnSpark

- Released in Feb. 2016
  - Apache 2.0 license
- Distributed deep learning
  - GPU or CPU
  - Ethernet or InfiniBand
- Easily deployed on public cloud or private cloud

YAHOO!

# CaffeOnSpark: One Notebook http://bit.ly/1REZ0cN

## Logistic Regression using MLlib

```
In [57]: from pyspark.mllib.linalg import Vectors
         from pyspark.mllib.regression import LabeledPoint
         from pyspark.mllib.classification import LogisticRegressionWithLBFGS
```

```
In [58]: data = f.map(lambda row: LabeledPoint(row.label[0], Vectors.dense(row.ip1)))
```

```
In [59]: lr = LogisticRegressionWithLBFGS.train(data, numClasses=10, iterations=10)
```

```
In [60]: predictions = lr.predict(data.map(lambda pt : pt.features))
```

```
In [61]: predictions.take(5)
```

Out[61]: [7, 2, 1, 0, 4]

```
         |00000001|    [1.0]|[1.3683326, -0.0,...|[2.0906663, 1.048...|[2.0]|
         |00000002|    [1.0]|[1.5641443, -0.0,...|[-0.773368, 10.61...|[1.0]|
         |00000003|    [1.0]|[-0.0, 1.9505613,...|[16.46351, -6.917...|[0.0]|
         |00000004|    [1.0]|[0.5979191, 0.075...|[-0.48371825, -2....|[4.0]|
         +--------+--------+--------------------+--------------------+-----+
```

```
In [45]: dl_train_source = DataSource(sc).getSource(cfg,True)
```

```
In [46]: cos.train(dl_train_source)
```

# Datasets for YOU

# Data, Data, Data.

- ## Machine learning is all about data
  - More data ➔ Better model
- ## Industrial researchers work on large-scale data
  - Academic researchers need data, too ☹

# Yahoo Webscope Datasets

| Categories | datatype | # of datasets |
|---|---|---|
| L: Language Data | l | 25 |
| G: Graph and Social Data | g | 8 |
| R: Ratings and Classification Data | r | 10 |
| A: Advertising and Market Data | a | 4 |
| C: Competition Data | c | 3 |
| S: Computing Systems Data | s | 4 |
| I: Image Data | i | 5 |

- *Ex. https://webscope.sandbox.yahoo.com/catalog.php?datatype=r*

- Non-commercial use by academics and scientists: 8,000+

YAHOO!

# L5 - Query Vectors

- **8M** query vectors trained using our search2vec system
- May serve as a testbed for query rewriting task
  - › IR research
  - › Word and sentence similarity task in NLP research

| Method | oAUC | Macro NDCG@5 |
|---|---|---|
| word2vec | 0.817 | 0.929 |
| search2vec | **0.880** | **0.959** |

YAHOO!

# I3 - Flickr Creative Commons

- 100M Flickr photos & videos

  › Largest public multimedia dataset ever released

  - ImageNet … 14M images

- Metadata

  › Tags, title, des & geo

  - 68M photos have at least one tag

  › Comments, favorites, social network data queried via Flickr API



**longing** by Robert-Crouse Baker

**Tags:** Ocean Beach, San Francisco, Pacific, Surf, surfing, waves, green, water, man, surfer, January, océan, le surf, vagues, mer, sea
**Description:** "My soul is full of longing for the secret of the sea,
and the heart of the great ocean sends a thrilling pulse through me." — Henry Wadsworth

# Summary

- Big-data technologies and continued innovation are critical for Yahoo business.

- Yahoo continues open source contribution ever since our donation of Hadoop to Apache
  › Apache Spark/Storm/Hbase/Omid, CaffeOnSpark
  › Webscope datasets

YAHOO!

# Thanks!
yahoohadoop.tumblr.com
bigdata@yahoo-inc.com

YAHOO!