# Hadoop Patterns & Practices

**11.30.2012**

**George Chu** (朱金生)
Sr. Director of Engineering
Hadoop, Cloud Services & Mobile

# Hadoop Today

# About Yahoo!

- Aspire to make the world's daily habits more inspiring & entertaining

- Focus on building highly personalized experiences that connect people to what matters most to them

- Connect advertisers and partners with the audiences who build their business

**Email** – 35M hours per day, 190M user engagements per day

**News** & **Information** – Personalized news, sports, finance, weather, etc.

**Photos** – Reimaging Flickr to make it faster, more beautiful & social

**Search** – Innovating search across platforms (Direct Display, Axis)

**Personalization** – 13M different versions of the Homepage tailored for users' distinct interests

**Across platforms** – Online, mobile, TV, second screen
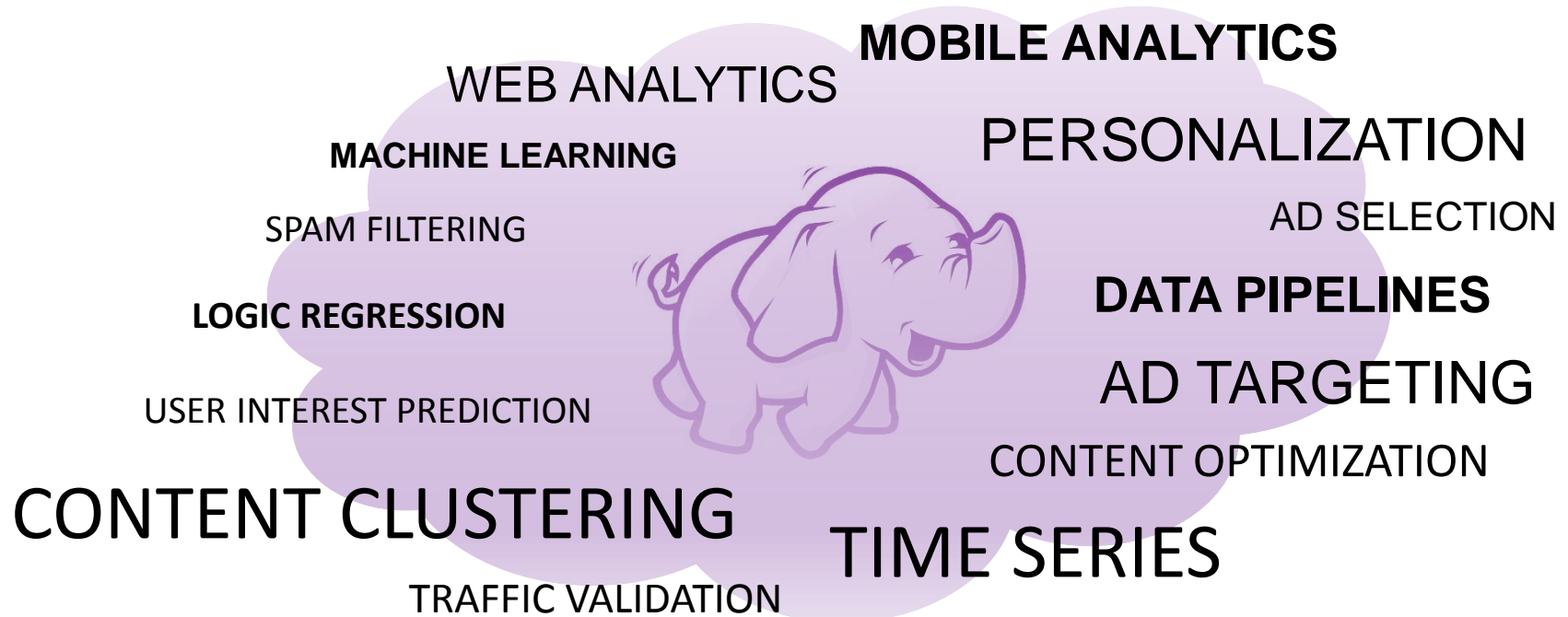
**Trusted** & **transparent**

**Personalized user experiences** help advertisers target the people that matter

**Global scale**

**Innovative platforms** (mobile, ad networks, video, media, partnerships, etc.)

**YAHOO!**

# Hadoop at Yahoo!

**Hadoop** is behind every click at **Yahoo!**, turning data into insights and making content and ads relevant for our consumers. You will quickly learn from this presentation about how **Yahoo!** is leveraging the cloud, scaling the core, and expanding the ecosystem.

MOBILE ANALYTICS

WEB ANALYTICS

MACHINE LEARNING

PERSONALIZATION

AD SELECTION

SPAM FILTERING

DATA PIPELINES

LOGIC REGRESSION

AD TARGETING

USER INTEREST PREDICTION

CONTENT OPTIMIZATION

CONTENT CLUSTERING

TIME SERIES

TRAFFIC VALIDATION

YAHOO!

# Yahoo!'s commitments to Hadoop & community

## Leveraging the Cloud

Yahoo! operates one of the world's largest private cloud infrastructure, handling **15B page visits per month** & **100B events per day** from more than **700M unique monthly users**. Content handled by the Yahoo! Cloud has grown to **more than 200PBs**, with **50TBs of additional data collected daily**.

The Hadoop project is an integral part of Yahoo!'s cloud infrastructure, and is at the heart of many of Yahoo!'s important business processes like Yahoo!'s next generation display advertising system, which **optimizes forecasting and pricing for over 24B ads served daily.**

Hadoop at Yahoo! works in concert with Yahoo!'s other cloud services such as **Edge Services** (built with standard Yahoo! Technologies like **Yahoo! Traffic Server), Data Servicing Containers, Distributed Structured** & **Unstructured Storage Services**, and **Data Highway** (Yahoo!'s event collection & delivery platform) for collecting, storing, processing, managing and analyzing data.

## Scaling the Core

Yahoo!'s technical leadership has taken Hadoop from a science project to a mainstream big data technology serving thousands of companies around the world. Yahoo! continues to be a key contributor across all areas of Hadoop, including **Hadoop 0.23** with next generation **MapReduce** and **HDFS federation**.

Yahoo! currently operates the largest production deployment of Hadoop clusters in the world made up of **over 42,000 servers** which process **10.8M jobs per month** and **140PBs of data**, stretching the scale and stability limits of core Hadoop MapReduce and HDFS.

Like **GridMix** (benchmark for Hadoop clusters) and **Vaidya** (rule-based performance diagnostics), Yahoo! continues to develop key tools that improve the overall stability and usability of Hadoop. Tools in the pipeline include **Groundhog** (Pig record and playback for regression testing), **QuAREH** (replay workloads and model utilization), **C3** (compute capacity calculator), and **Anarchy Ape** (fault injection into Hadoop clusters).

## Expanding the Ecosystem

Yahoo! continues to make the Hadoop ecosystem stronger, working closely with key collaborators in the Hadoop community and helping to drive more users and contributors to Hadoop. Yahoo! remains significantly invested in code, resources, and adoption of technology to further ensure a strong and vibrant Hadoop community.

Yahoo! has contributed a majority of the current Hadoop code and other related projects such as Pig (language to express data transformation), **Oozie** (workflow scheduling & coordination system), **Zookeeper** (centralized service for highly reliable distributed coordination), and **HCatalog** (unified table & schema management).

Yahoo! has also adopted other Hadoop stack components such as **Hive** and **HBase** from the open source community to solve additional use-cases, and looks forward to extending the capabilities of these components and contributing the development back to the community.

# Hadoop Operational Statistics | May, 2012

| | |
|---|---:|
| Number of (centrally-managed) nodes | **42,000** |
| | **Production** |
| | **Research** (Ad-hoc usage) |
| Cluster types | **Sandbox** (Release validation) |
| | **Innovation** (Dev, QE, Benchmarking) |
| | **Data Loading** |
| Maximum nodes per cluster | **4,000** |
| HDFS | **>350** petabytes |
| Compute Slots (Map and Reduce Slots) | **>9M** slot hours available / day |
| Jobs submitted | **>10M** / month |
| Number of monthly unique users (submitting jobs in a month) | **>1000** / month |
| Average daily unique users (submitting jobs each day) | **>300** / day |

# Hadoop usage at Yahoo!

| **Search Assist** | **Behavioral Targeting** |
|---|---|

**Problem**
Related concepts appear close together in text corpus to assist users with search term

**Problem**
Quickly make complex decisions to serve the right ads to the right customer by targeting billions of impressions per day across one of the largest ad networks in the world



**Solution**
Hadoop helps Yahoo! process 1B web pages of about 10K bytes each (10 TB of input data) to create the output list of related words
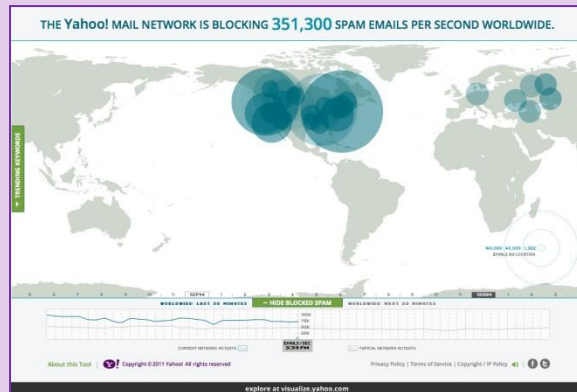
**Solution**
Hadoop helps Yahoo! process declared data and recent activity to segment users and determine the right ads to serve in milliseconds

YAHOO!

# Hadoop usage at Yahoo! | Continued

| **Mail Anti-Spam** | **Membership Anti-Abuse** |
|---|---|

**Problem**

Yahoo! Mail delivers 5.6 billion email a day across 300 million mailboxes. Users want to see emails from friends and family in the inboxes, from the people who matter the most... not from spammers and phishers

**Problem**

Membership processes 2.22M new registrations (127 M logins) every day! Abuse taints metrics, and parsing out abusive vs. legitimate user is an ongoing challenge
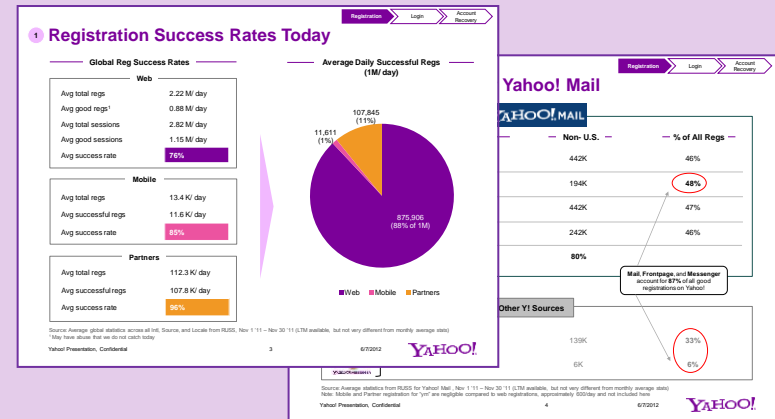


**Solution**

Hadoop helps Yahoo! block 20.5 billion spam emails per day through machine learning on the grid. SpamGuard in conjunction with Hadoop has reduced spam by 60%

**Solution**

Hadoop helps Yahoo! detect abusive registrations through machine learning on the grid

YAHOO!

# Hadoop usage at Yahoo! | Continued

| **Content Agility** | **Personalization** |
|---|---|
| **Problem** | **Problem** |
| Properties had siloed approaches for CMS, front-end development and editorial. A common solution was needed for the entire content network to bring agility to Yahoo! properties | Increase engagement by showing the right content to users with input from science & human editors |



**Solution**
Leverage Content Agility as the single, grid-based, highly scalable CMS. Lego provides reusable UI modules and shared tools

**Solution**
Personalization requires a real-time feedback loop across properties, leveraging user interests, intent, and context to optimize user engagement. Hadoop/HBase is leveraged for modeling (item/ user) and async processing, Hive for analytics

YAHOO!

# Hadoop Tomorrow

# Trends to address

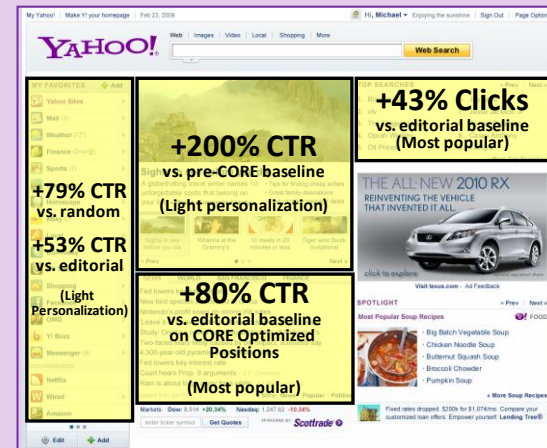| | |
|---|---|
| **T1** | **Data ingestion**<br>▪ Data is doubling every 18 months (50-60% annual growth) with greater access from multiple touch points<br>▪ Analysis has become more important than storage and retrieval with increasing information velocity, volume, and variety |
| **T2** | **Thriving open-source community**<br>▪ Strong community contributing to Hadoop from Yahoo!, eBay, Facebook, LinkedIn, Twitter, Hortonworks, Cloudera, etc.<br>▪ Seen as the preferred solution for big data analytics with several interchangeable components e.g. Cassandra/ HBase |
| **T3** | **New use-cases**<br>▪ From increasing sales and user engagement to managing risk and fraud, detecting stock market patterns, predicting mortgage default rates to civil infrastructure and telco churn management<br>▪ From back-office offline analytics to customer-facing 24x7 production systems |
| **T4** | **Traditional vendor's dilemma – Extend or Fork**<br>▪ Consolidating advanced SQL and NoSQL space (EMC, IBM, HP, and Terradata making their moves with Greenplum, Netezza, Vertica, and Aster Data)<br>▪ Hadoop-based distribution in the hopes of pushing appliance and service sales (Greenplum HD Data Computing Appliance, IBM InfoSphere BigInsights, Informatica 9.1 Big Data Integration platform)<br>▪ No winning solution or approach to processing big data yet (transactional or non-transactional for processing and management) |
| **T5** | **Integrated & real-time processing**<br>▪ Need for an integrated customer solution that can handle transactions, processing, search, and analytics across the entire data set<br>▪ Increased interest in knowing what is happening right now vs. offline processing implies solution support for simultaneous read and write with real-time response |

Source: McKinsey Quarterly, Gartner Research, Yahoo! Research

**YAHOO!**

# Paradigm shifts | Cost/Behavioral angle

| | |
|---|---|
| **Technological**<br>▪ The ability to effectively process multi-petabytes of data across thousands of inexpensive computing resources | Compute |
| **Economic**<br>▪ The cost of data acquisition has practically gone to zero<br>▪ The cost of data storage is approaching zero | Data |
| **Social & Mobile**<br>▪ Consumer's increasing comfort in pushing user-generated content for sharing<br>▪ Consumer's increasing demand for pulling both public and private information, anywhere, anytime, on every form-factor | Latency |

# Paradigm shifts | Technology angle

|         | Scenario | Characteristic |
|---------|----------|----------------|
| **Compute** | **Offload** | Off critical resources |
|         | **Load-balance** | Better resource utilization |
|         | **Batch** | Faster time to completion |
|         | **Speculate** | Real-time, immersive experience |
| **Data** | **Acquire** | Public or proprietary |
|         | **Cure** | Transformative steps prior to compute |
|         | **Aggregate** | Compose structured & unstructured |
|         | **Version** | Storage management |
| **Latency** | **Cache** | Access anywhere, anytime, on any form-factor |
|         | **Stage** | Move large-scale data across storage tiers |
|         | **Distribute** | Locality of reference |
|         | **Transact** | Stateful, resuming where you've left off |

**YAHOO!**®

# New value-props

- Run both compute- (math & statistical) as well as data-intensive workloads

# HPC vs. Hadoop | Architectural comparison

## HPC

**Computational**
(Supercomputing)
**C++/Java Programming**
**MPI**
(Numeric Solutions to Partial Differential Equations)
**OpenMP**
(Image Processing, Rendering, Transcoding)
**SOA**
(Black-Scholes, Monte Carlos, Correlation Matrices)

**On-premise Enterprise Data**
(in Multi-Terabyte range)
**Structured, Floating Points**

**Closely-coupled Architecture**
(Scale to 1000 nodes)
**Fine-grain Scheduling**
**Multiple Network Topologies**
(Compute Node secure isolation)
**Low Latency**
(InfiniBand)
**Mission Critical High Availability**
**Enterprise-grade**
(Security, Versioning, Servicing)

## Hadoop

**Data Intensive**
(NoSQL querying)
**Java/Pig Scripting**
**Map-Reduce**
(Web Analytics, Ads Targeting, Spam Filtering, Web Traffic Validation)

**Web Data in Cloud**
(in Multi-Petabyte range)
**Unstructured, Semi-numerical**

**Loosely-coupled Architecture**
(Scale to 4000 nodes)
**Coarse-grain Scheduling**
**Single Network Topology**
(RPC wire protocol)
**Higher Latency**
**No High Availability**

- *Run MR workloads*
- *Build a pathway to data in cloud*
- *Scale to a higher node count*
- *Add GP-GPU support*

- *Add fine-grain scheduling*
- *Run non-MR workloads*
- *Be enterprise-grade*
- *Add GP-GPU support?*

# HPC vs. Hadoop | HPC workload characterization

**Data Types**

**Data Parallelism**

Traditional, transacted database
Economic simulation
Circuit simulation
Game-tree exploration
Particle movement under fields
Lagrangian fluid dynamics (fluid element decomposition)
Finite-element stress-strain analysis (crash test)
N-body problem
Euleran fluid dynamics (spatial decomposition)
Dense matrix multiplication & factorization
SETI
Integer factorization
Gene matching

**Compute Traits**

Real-time
Physical simulation
Machine learning
Number crunching
Non/Semi-numerical
Boolean
Discrete mathematics
Low/High-precision continuous

**Arrows** indicate a direction toward generally increasing suitability or scalability to cloud migration

# Putting it all together | High-level data workflow

| Data | Compute | Latency |
|---|---|---|
| (Data collection) | (Offline asynchronous processing) | (Online synchronous serving) |

**Acquire Data**

Sensor Data
- PC

**Cure Data** (if raw)

Correctness
Linear algebra
Arithmetical
Statistical

Web crawl

Social graphs

3rd party content

Email

**Run Model on Data**

Offload
Load-balance
Batch
Speculate

**Aggregate & Version Data**

**Stage Data**

Ads Serving

**Distribute & Cache**

Cloud Serving

Edge Serving

**Transact** (if stateful)

YAHOO!