

# The BigData Top100 List Initiative

Chaitan Baru  
San Diego Supercomputer Center

# Background

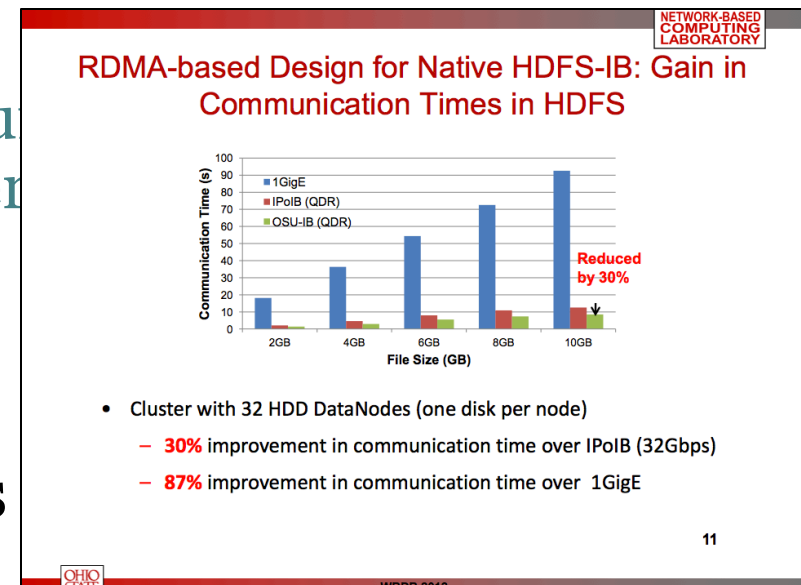
- Workshop series on Big Data Benchmarking (WBDB)
  - First workshop, May 2012, San Jose. Hosted by Brocade.
  - Second workshop, December 2012, Pune, India. Hosted by Persistent Systems / Infosys.
  - Third workshop, July 2013, Xi'an, China. Hosted by Xi'an University.
  - Fourth workshop, October 9-10, 2013, San Jose. Hosted by Brocade. (collocated with IEEE Big Data Conference)
- Selected WBDB papers to be published in two volumes of Springer Verlag LNCS, for 2012 and 2013.
- Paper from First Workshop
  - Setting the Direction for Big Data Benchmark Standards by C. Baru, M. Bhandarkar, R. Nambiar, M. Poess, and T. Rabl, published in *Selected Topics in Performance Evaluation and Benchmarking*, Springer-Verlag
- Article in inaugural issue of Big Data Journal
  - *Big Data Benchmarking and the Big Data Top100 List* by Baru, Bhandarkar, Nambiar, Poess, Rabl, Big Data Journal, Vol.1, No.1, 60-64, Anne Liebert Publications.

# Some principles for good benchmark design

- Self-scaling, e.g. TPC-C
- Comparability between scale factors
  - Results should be comparable at different scales
- Extrapolating Results
  - To larger configurations
- Elasticity and durability
  - Big data systems should be intrinsically elastic and durable
  - TPC runs ACID outside the performance window
- Performance + Price/performance
  - Try to capture price in a simple, intuitive, meaningful way...
  - For **price/performance**: what is the most useful quantity for price?
- Simple to run

# Benchmarks

- Several types of data benchmarks
- Micro-benchmarks
  - E.g. A Micro-benchmark Suite for HDFS Operations on Modern, OSU
- Functional benchmarks
  - Terasort
- Genre-specific benchmarks
  - Graph500
- Application-level benchmarks
  - TPC-C



# How to think about application-level benchmarks?

- Which application domain / application pattern to cover?
- What should the data look like?
- What are the operations on that data?
- E.g. TPC-C (i.e., Order-Entry)
  - A midweight read-write transaction (i.e., New-Order),
  - A lightweight read-write transaction (i.e., Payment),
  - A midweight read-only transaction (i.e., Order-Status),
  - A batch of midweight read-write transactions (i.e., Delivery)
  - A heavyweight read-only transaction (i.e., Stock-Level)
  - Specified in the semantic context, or story-line, of an order processing environment.
- Pushed technologies for transaction processing

# TPC-D

- Decision support use case
- Sales applications
  - Lineitem table + dimension tables
  - Lots of aggregation
  - Views
- Pushed technology on query optimization
- Concurrent query streams
  - Query order different for each stream (specified)
- Side effects
  - Introduced “trickle” update to exclude read-only solutions
  - Included reporting of data load times
- Became “outdated” due to materialized views

# Big Data

- Data at the “edges” of the enterprise
  - Outside of enterprise transaction systems
- “Interactions” data
  - E.g. log data, social network data, etc.
- Used for “event detection”
  - User clicks
  - Device failures
  - Hospital re-admissions
  - ...

# Big Data - Application characteristics

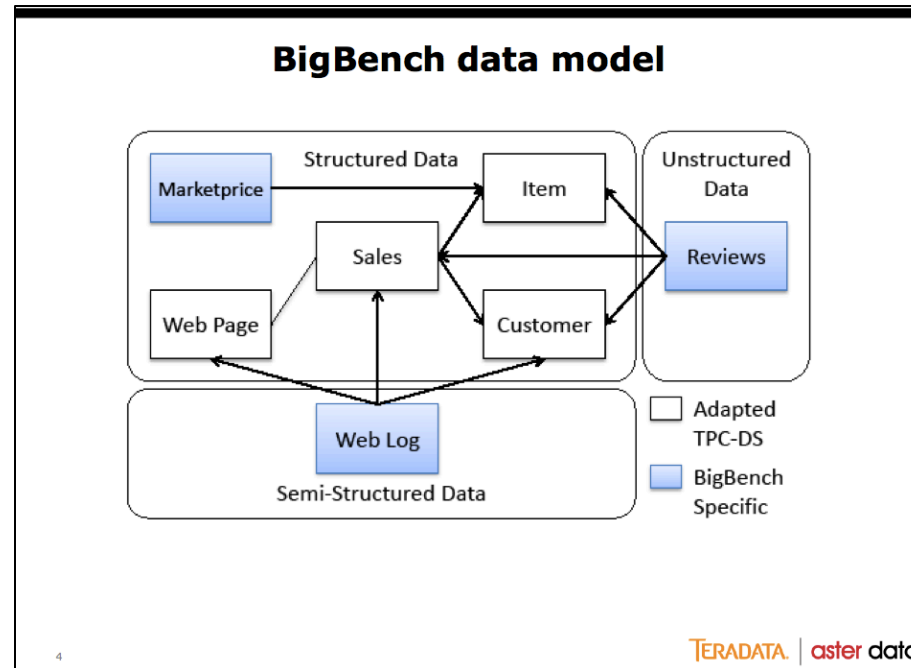
- Data characteristics
  - Data from multiple sources
  - Loose, flexible schema
  - Data requires structuring
  - ETL / ELT
- Application characteristics
  - Pipeline of processing
  - Running models with data



# BigBench: A proposal

- By Ghazal et al: Teradata, Oracle, U.of Toronto, InfoSizing
- Derives from TPC-DS
- TPC-DS:
  - Multiple snowflake schemas with shared dimensions
  - 24 tables with an average of 18 columns
  - 99 distinct SQL 99 queries with random substitutions
  - More representative skewed database content
  - Sub-linear scaling of non-fact tables
  - Ad-hoc, reporting, iterative and extraction queries
  - ETL-like data maintenance

# BigBench Data Model



- Workload = Set of queries
  - On structured, semistructured, unstructured data
  - Data mining, ML

# Alternative Proposal: Data Analytics Pipeline

- Sequence of processing steps
- Feed data from one to the next
- Different operations at each step
- From ETL/ELT to SQL to ML and PA
- Data generation of event-based data

# Data Analytics Pipeline: Quest for Typical Workload

- Tune systems for broadly applicable workloads
- Benchmarks most relevant if representative
- Designing optimized systems: Make common tasks fast, other tasks possible

# Is There a Typical Big Data Pipeline Workload?

- Big Data systems are characterized by flexibility
- Multiple Interfaces: SQL, MapReduce, Streaming, GraphLab,...
- Workloads evolving rapidly

# Encouraging Early Results

- Analysis of characteristics of 1M+ real Hadoop jobs on production clusters at Yahoo, 100+ features
- Identified 8 Job types
- Verified with GridMix 3
- Characterization of Hadoop Jobs Using Unsupervised Learning, Sonali Aggarwal, Shashank Phadke & **Milind Bhandarkar**, in 2010 IEEE Second International Conference on Cloud Computing Technology and Science, Indianapolis, Indiana, December 2010, <http://doi.ieeecomputersociety.org/10.1109/CloudCom.2010.20>

# Big Data Sources

- Events
  - Direct - Human Initiated
  - Indirect - Machine Initiated
- Software Sensors (Clickstreams, Locations)
- Public Content (blogs, tweets, Status updates, images, videos)

# User Modeling

- Objective: Determine user interests by mining user activities
- Large dimensionality of possible user activities
- Typical user has sparse activity vector
- Event attributes change over time



# User Modeling Pipeline

- Data Acquisition
- Sessionization
- Feature and Target Generation
- Model Training
- Offline Scoring & Evaluation
- Batch Scoring & Upload to serving

# Data Acquisition

- Batched and collected at the edge
- Loaded incrementally
- Simple ETL/ELT
- Append / Replace

# Denormalization / Cleansing

- Augment raw events data with attributes
- Look up dictionaries
- Multi-way Joins with dimension tables

# Find Sessions with Target Events

- User-Actions of Interest
  - Clicks on Ads & Content
  - Site & Page visits
  - Conversion Events
    - Purchases, Quote requests
    - Sign-Up for membership etc

# Feature Selection

- Summary of user activities over a time-window
- Aggregates, moving averages, rates over various time-windows
- Incrementally updated

# Join Targets and Features

- Target rates very low: 0.01% ~ 1%
- First, construct targets
- Filter user activity without targets
- Join feature vector with targets

# Model Training

- Regressions
- Boosted Decision Trees
- Naive Bayes
- Support Vector Machines
- Maximum Entropy modeling

# Offline Scoring & Evaluations

- Apply model weights to features
- Pleasantly parallel
- Sort by scores and compute metrics
- Evaluate metrics



# Batch Scoring

- Apply models to features from all user activity
- Upload scores to serving systems

## 5 Different Classes

- Tiny (10K Entities, 100GB)
- Small (100K Entities, 1TB)
- Medium (1M Entities, 10 TB)
- Large (10M Entities, 100 TB)
- Huge (1B Entities, 1PB)

# Technical Issues - 1

- Audience: Who is the audience for such a benchmark?
  - Marketing, Internal Use, Academic Use
- Application: What is the application that should be modeled?
  - Abstractions of a data pipeline, e.g. Internet-scale business
- Single benchmark spec: Is it possible to develop a single benchmark to capture characteristics of multiple applications?
  - Single, multi-step benchmark, with plausible end-to-end scenario

## Technical Issues - 2

- Component vs. end-to-end benchmark. Is it possible to factor out a set of benchmark “components”, which can be isolated and plugged into an end-to-end benchmark?
  - The benchmark should consist of individual components that ultimately make up an end-to-end benchmark
- Paper and Pencil vs Implementation-based. Should the implementation be specification-driven or implementation-driven?
  - Start with an implementation and develop specification at the same time

# Technical Issues - 3

- Reuse. Can we reuse existing benchmarks?
  - Leverage existing work and built-up knowledgebase
- Benchmark Data. Where do we get the data from?
  - Synthetic data generation: structured, semistructured, unstructured data
- Innovation or competition? Should the benchmark be for innovation or competition?
  - Successful competitive benchmarks will be used for innovation

# Moving Forward...

- BigData Top100 will select a benchmark specification for ranking systems
- Two alternatives for workload specification
  - **BigBench**: based on TPC-DS
    - Extended with semistructured and unstructured data and operations on those data. See 1<sup>st</sup> WBDB and 2<sup>nd</sup> WBDB program online for slides.
  - **Data Analytics Pipeline (DAP)**
    - Proposed end-to-end pipeline, from data ingestion to predictive modeling
  - **Track [bigdatatop100.org/benchmarks](http://bigdatatop100.org/benchmarks) for info**
- Use Kaggle for competitions
  - Propose data generation programs
  - Propose / evaluate operations for steps in the Data Analytics Pipeline

# Deadlines...

- March 15, 2013
  - BigBench turned into first draft of a benchmark spec
  - First draft of data and operations for the Data Analytics Pipeline
- April 15, 2013
  - Receive comments from the community
- May 1, 2013:
  - Short paper submission deadline for 3<sup>rd</sup> WBDB, July 16-17, Xi'an, China.
  - By 3<sup>rd</sup> WBDB, July 2013: Proposals for metrics, execution rules, audit rules and reporting rules.
- July 16, 2013
  - Reference implementations of BigBench and Data Analytics Pipeline
- August 31, 2013: Release of a benchmark specification