# Distributed Data Parallel Computing: The Sector Perspective on Big Data

# Robert Grossman

July 25, 2010

Laboratory for Advanced Computing
University of Illinois at Chicago

Open Data Group

Institute for Genomics & Systems Biology
University of Chicago 1

# Part 1.

# Open Cloud Testbed

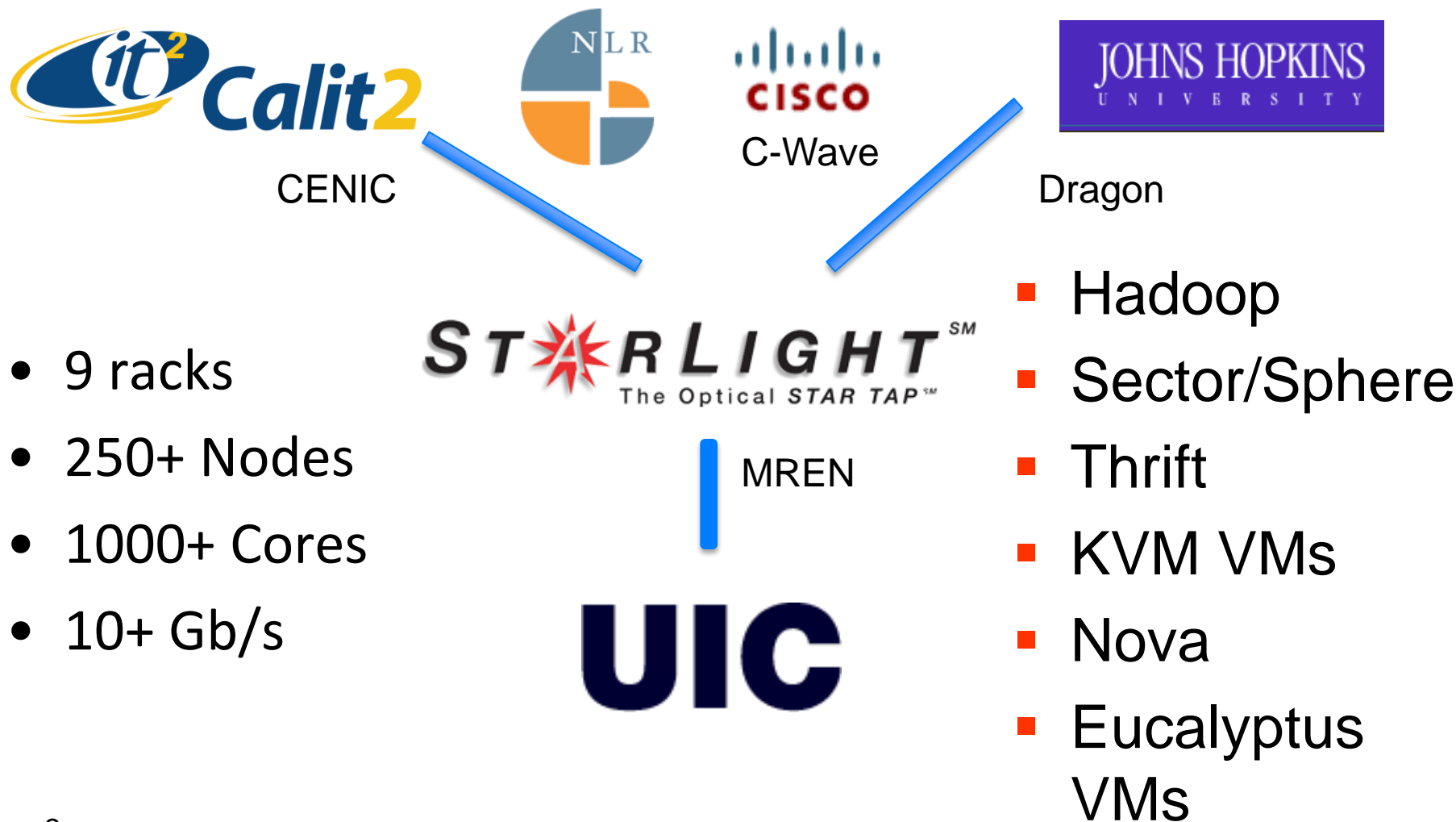**Calit2** — CENIC

**NLR**

**CISCO** C-Wave

**JOHNS HOPKINS UNIVERSITY** — Dragon

**STARLIGHT** ᔆᴹ
The Optical STAR TAP ᔆᴹ

MREN

**UIC**

- 9 racks
- 250+ Nodes
- 1000+ Cores
- 10+ Gb/s

- Hadoop
- Sector/Sphere
- Thrift
- KVM VMs
- Nova
- Eucalyptus VMs

# Open Science Data Cloud

Open Cloud Consortium

YAHOO!

NLR

CISCO

JOHNS HOPKINS UNIVERSITY

iT² Calit2

STARLIGHT℠
The Optical STAR TAP℠
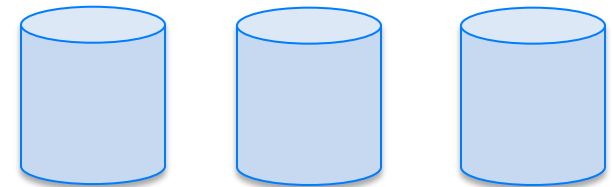
sky cloud

NSF OSDC PIRE Project – Working with 5 international partners (all connected with 10 Gbps networks).
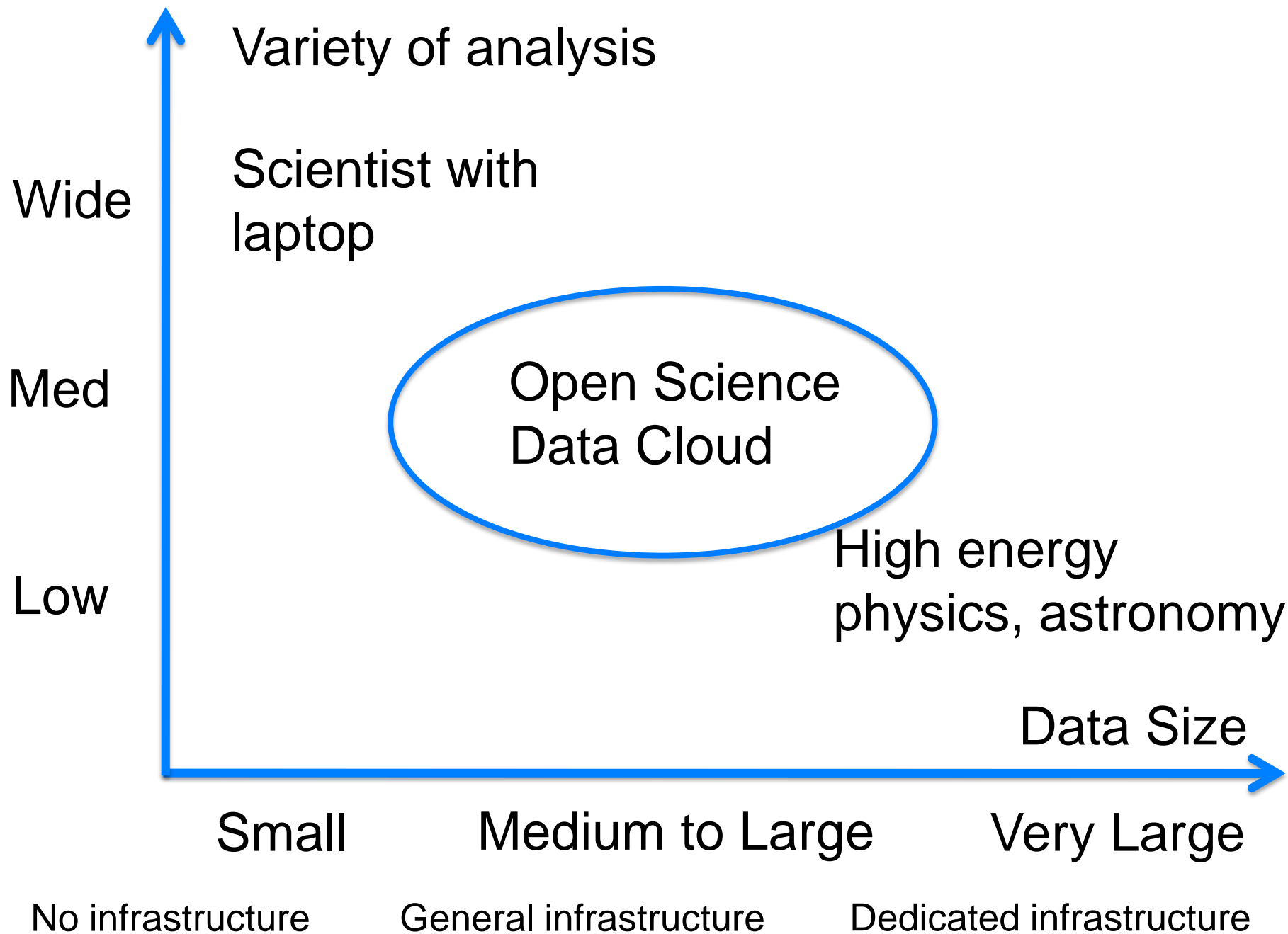
UIC

THE UNIVERSITY OF CHICAGO
BIOLOGICAL SCIENCES

Bionimbus (biology & health care)

# Part 2
# What's Different About
# Data Center Computing?

Data center scale computing provides storage and computational resources at the scale and with the reliability of a data center.

# The Datacenter as a Computer

*An Introduction to the Design of Warehouse-Scale Machines*

**Luiz André Barroso and Urs Hölzle**
Google Inc.

A very nice recent book by
Barroso and Holzle

# Scale is new

# Elastic, Usage Based Pricing Is New

costs the same as

1 computer in a rack for 120 hours

120 computers in three racks for 1 hour

# Simplicity of the Parallel Programming Framework is New

A new programmer can develop a program to process a container full of data with less than day of training using MapReduce.

# Elastic Clouds



# Large Data Clouds



Goal: Minimize cost of virtualized machines & provide on-demand.

# HPC



Goal: Maximize data (with matching compute) and control cost.

Goal: Minimize latency and control heat.

2003
10x-100x

data
science

1976
10x-100x

1670
250x

simulation
science

1609
30x

experimental
science

|  | **Databases** | **Data Clouds** |
|---|---|---|
| Scalability | 100's TB | 100's PB |
| Functionality | Full SQL-based queries, including joins | Single keys |
| Optimized | Databases optimized for safe writes | Data clouds optimized for efficient reads |
| Consistency model | ACID (Atomicity, Consistency, Isolation & Durability) | Eventual consistency |
| Parallelism | Difficult because of ACID model; shared nothing is possible | Parallelism over commodity components |
| Scale | Racks | Data center |

|  | **Grids** | **Clouds** |
| --- | --- | --- |
| Problem | Too few cycles | Too many users & too much data |
| Infrastructure | Clusters and supercomputers | Data centers |
| Architecture | Federated Virtual Organization | Hosted Organization |
| Programming Model | Powerful, but difficult to use | Not as powerful, but easy to use |

# Part 3
# How Do You Program A Data Center?

# How Do You Build A Data Center?



- Containers used by Google, Microsoft & others
- Data center consists of 10-60+ containers.



Microsoft Data Center, Northlake, Illinois

# What is the Operating System?

VM 1 ... VM 5

VM 1 ... VM 50,000

Data Center Operating System

workstation

- Data center services include: VM management services, VM fail over and restart, security services, power management services, etc.

# Architectural Models:
# How Do You Fill a Data Center?

on-demand computing instances

| App | App | ... | App |
|-----|-----|-----|-----|

large data cloud services

| App | App | App | App | App |
|-----|-----|-----|-----|-----|

| Cloud Data Services (BigTable, etc.) | Quasi-relational Data Services | App | App |
|---|---|---|---|

| Cloud Compute Services (MapReduce & Generalizations) | App | App |
|---|---|---|

Cloud Storage Services

# Instances, Services & Frameworks

| | | | |
|---|---|---|---|
| | | | |
| **many instances** | | Hadoop DFS & MapReduce Amazon's SQS Azure Services | Microsoft Azure Google AppEngine | VMWare Vmotion… |
| **single instance** | Amazon's EC2 | S3 | | |
| | instance (IaaS) | service | framework (PaaS) | operating system |

# Some Programming Models for Data Centers

- Operations over data center of disks
  - MapReduce ("string-based")
  - Iterate MapReduce (Twister)
  - DryadLINQ
  - User-Defined Functions (UDFs) over data center
  - SQL and Quasi-SQL over data center
  - Data analysis / statistics functions over data center

# More Programming Models

- Operations over data center of memory
  - Memcached (distributed in-memory key-value store)
  - Grep over distributed memory
  - UDFs over distributed memory
  - SQL and Quasi-SQL over distributed memory
  - Data analysis / statistics over distributed memory

# Part 4.  Stacks for Big Data

# The Google Data Stack



- The Google File System (2003)
- MapReduce: Simplified Data Processing… (2004)
- BigTable: A Distributed Storage System… (2006)

# Map-Reduce Example

- Input is file with one document per record

- User specifies map function
  - key = document URL
  - Value = terms that document contains

("doc cdickens",
"it was the best of times")

→ map

"it", 1
"was", 1
"the", 1
"best", 1

# Example (cont'd)

- MapReduce library gathers together all pairs with the same key value (shuffle/sort phase)
- The user-defined reduce function combines all the values associated with the same key

key = "it"
values = 1, 1

key = "was"
values = 1, 1

key = "best"
values = 1

key = "worst"
values = 1

reduce

"it", 2
"was", 2
"best", 1
"worst", 1

# Applying MapReduce to the Data in Storage Cloud

map/shuffle

reduce

# Google's Large Data Cloud

Applications

Compute Services          Google's MapReduce

Data Services             Google's BigTable

Storage Services          Google File System (GFS)

Google's Stack

# Hadoop's Large Data Cloud

| | |
|---|---|
| **Applications** | |
| **Compute Services** | Hadoop's MapReduce |
| **Data Services** | NoSQL Databases |
| **Storage Services** | Hadoop Distributed File System (HDFS) |

Hadoop's Stack

# Amazon Style Data Cloud



Load Balancer

Simple Queue Service

SDB

EC2 Instances

EC2 Instances

S3 Storage Services

# Evolution of NoSQL Databases

- Standard architecture for simple web apps:
  - Front end load balanced web servers
  - Business logic layer in the middle
  - Backend database
- Databases do not scale well with very large numbers of users or very large amounts of data
- Alternatives include
  - Sharded (partitioned) databases
  - master-slave databases
  - memcached

# NoSQL Systems

- Suggests No SQL support, also Not Only SQL

- One or more of the ACID properties not supported

- Joins generally not supported

- Usually flexible schemas

- Some well known examples: Google's BigTable, Amazon's S3 & Facebook's Cassandra

- Several recent open source systems

# Different Types of NoSQL Systems

- Distributed Key-Value Systems
  - Amazon's S3 Key-Value Store (Dynamo)
  - Voldemort

- Column-based Systems
  - BigTable
  - HBase
  - Cassandra

- Document-based systems
  - CouchDB

# Cassandra vs MySQL Comparison

- MySQL > 50 GB Data
  Writes Average : ~300 ms
  Reads Average : ~350 ms

- Cassandra > 50 GB Data
  Writes Average : 0.12 ms
  Reads Average : 15 ms

Source: Avinash Lakshman, Prashant Malik, Cassandra Structured Storage System over a P2P Network, static.last.fm/johan/nosql-20090611/cassandra_nosql.pdf

# CAP Theorem

- Proposed by Eric Brewer, 2000
- Three properties of a system: consistency, availability and partitions
- You can have at most two of these three properties for any shared-data system
- Scale out requires partitions
- Most large web-based systems choose availability over consistency

Reference: Brewer, PODC 2000; Gilbert/Lynch, SIGACT News 2002

# Eventual Consistency

- All updates eventually propagate through the system and all nodes will eventually be consistent (assuming no more updates)

- Eventually, a node is either updated or removed from service.

- Can be implemented with Gossip protocol

- Amazon's Dynamo popularized this approach

- Sometimes this is called BASE (**B**asically **A**vailable, **S**oft state, **E**ventual consistency), as opposed to ACID

# Part 5.  Sector Architecture

# Design Objectives

1. Provide Internet scale data storage for large data

    – Support multiple data centers connected by high speed wide networks

2. Simplify data intensive computing for a larger class of problems than covered by MapReduce

    – Support applying User Defined Functions to the data managed by a storage cloud, with transparent load balancing and fault tolerance

# Sector's Large Data Cloud

Applications

Compute Services — Sphere's UDFs

Data Services

Storage Services — Sector's Distributed File System (SDFS)

Routing & Transport Services — UDP-based Data Transport Protocol (UDT)

Sector's Stack

39

# Apply User Defined Functions (UDF) to Files in Storage Cloud

map/shuffle reduce

UDF

# UDT



udt.sourceforge.net

Sterling Commerce

Movie2Me

Globus

Nifty TV

Power Folder

UDT has been downloaded 25,000+ times

# Alternatives to TCP – Decreasing Increases AIMD Protocols



$$x \leftarrow x + \alpha(x) \qquad \text{increase of packet sending rate x}$$

$$x \leftarrow (1 - \beta)\, x \qquad \text{decrease factor}$$

# System Architecture

User account
Data protection
System Security

Metadata
Scheduling
Service provider

System access tools
App. Programming
Interfaces

Security Server

Masters

Clients

SSL

SSL

Data

UDT
Encryption optional

slaves

slaves

Storage and
Processing

|                      | **Hadoop DFS**             | **Sector DFS**              |
| -------------------- | -------------------------- | --------------------------- |
| Storage Cloud        | Block-based file system    | File-based                  |
| Programming Model     | MapReduce                  | UDF & MapReduce             |
| Protocol             | TCP                        | UDP-based protocol (UDT)    |
| Replication          | At write                   | At write or period.         |
| Security             | Not yet                    | HIPAA capable               |
| Language             | Java                       | C++                         |

|  | MapReduce | Sphere |
|---|---|---|
| Storage | Disk data | Disk & in-memory |
| Processing | Map followed by Reduce | Arbitrary user defined functions |
| Data exchanging | Reducers pull results from mappers | UDF's push results to bucket files |
| Input data locality | Input data is assigned to nearest mapper | Input data is assigned to nearest UDF |
| Output data locality | NA | Can be specified |

# Terasort Benchmark

|  | 1 Rack | 2 Racks | 3 Racks | 4 Racks |
|---|---|---|---|---|
| Nodes | 32 | 64 | 96 | 128 |
| Cores | 128 | 256 | 384 | 512 |
| Hadoop | 85m 49s | 37m 0s | 25m 14s | 17m 45s |
| Sector | 28m 25s | 15m 20s | 10m 19s | 7m 56s |
| Speed up | 3.0 | 2.4 | 2.4 | 2.2 |

Sector/Sphere 1.24a, Hadoop 0.20.1 with no replication on Phase 2 of Open Cloud Testbed with co-located racks.

# MalStone



sites □  entities →

$d_{k-2}$  $d_{k-1}$  $d_k$

time →

# MalStone Benchmark

| | MalStone A | MalStone B |
|---|---|---|
| Hadoop | 455m 13s | 840m 50s |
| Hadoop streaming with Python | 87m 29s | 142m 32s |
| Sector/Sphere | 33m 40s | 43m 44s |
| Speed up (Sector v Hadoop) | 13.5x | 19.2x |

Sector/Sphere 1.20, Hadoop 0.18.3 with no replication on Phase 1 of Open Cloud Testbed in a single rack.  Data consisted of 20 nodes with 500 million 100-byte records / node.
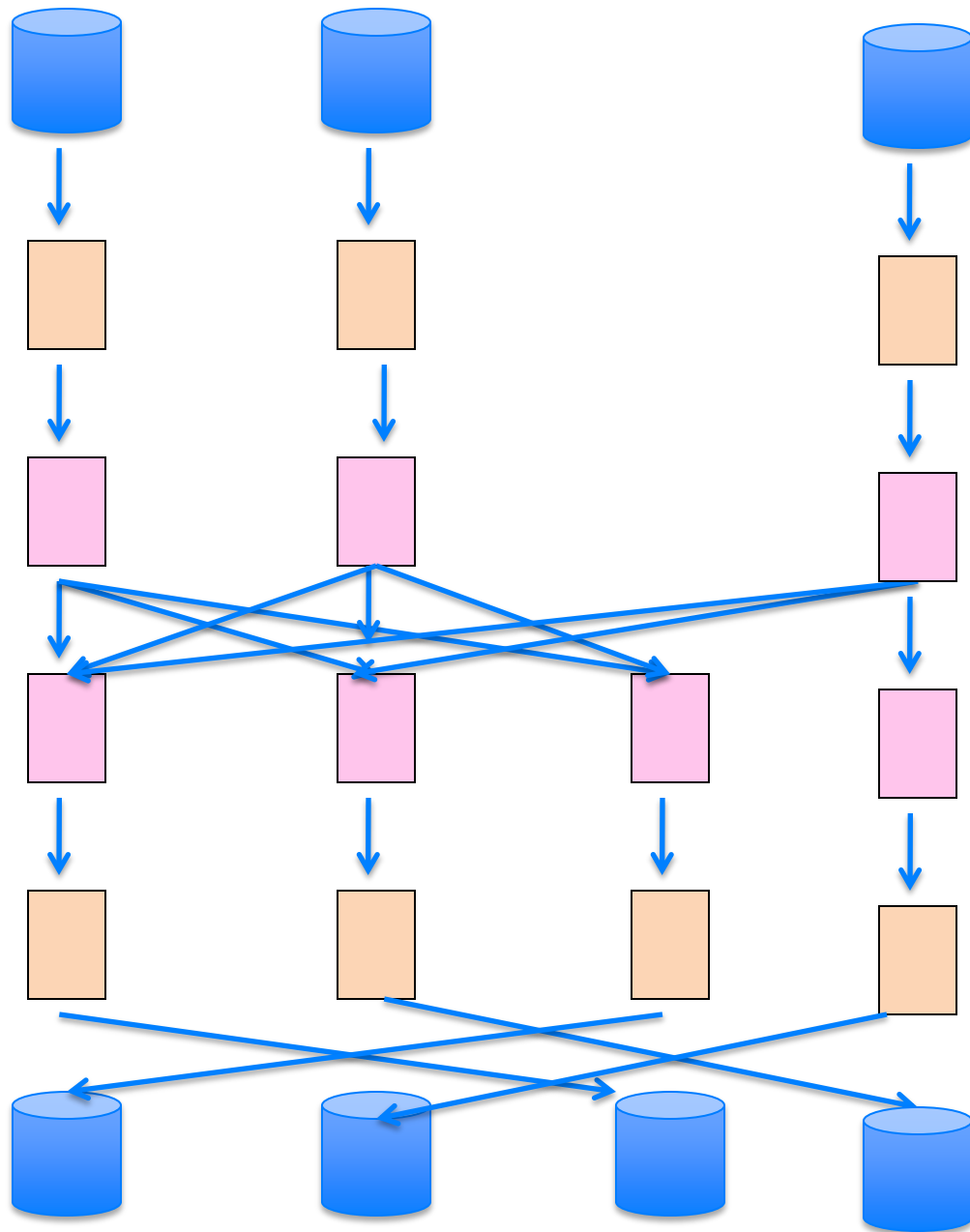
Disks

Input Segments

UDF

Bucket Writers

Output Segments

Disks

- Files not split into blocks
- Directory directives
- In-memory objects

# Sector Summary

- Sector is fastest open source large data cloud
  - As measured by MalStone & Terasort
- Sector is easy to program
  - UDFs, MapReduce & Python over streams
- Sector does not require extensive tuning
- Sector is secure
  - A HIPAA compliant Sector cloud is being launched
- Sector is reliable
  - Sector supports multiple active master node servers

# Part 6.
# Sector Applications

# App 1: Bionimbus



## Public Cistrack Data

[Drosophila Chromatin Time Course] - modENCODE

[Drosophila Insulator] - modENCODE

[All Data]- modENCODE

Browse Flynet

Browse and download public Cistrack data

-[By Experiment]

-[By File]

-[By Experimental Unit]

## Lookup Cistrack accession number

File ▾ [                    ] search

## Welcome to Cistrack

Cistrack supports data distribution for the Drosophila modENCODE cis-regulatory project (NHGRI contract U01HG004264) and the Chicago Center for Systems Biology (NIGMS grant P50GN081892).

Some browsers experience a problem when downloading public data. If you are asked to login when accessing public data, please close the web page and try again. This problem will be fixed with the next release.

## Cistrack Users

Login to Cistrack or [Register]

Upload data files

Annotate uploaded files with metadata

Cistrack Wiki

www.bionimbus.org

# App 2. Sector Application: Cistrack & Flynet

Cistrack Web Portal & Widgets

Cistrack Elastic Cloud Services

Cistrack Database

Analysis Pipelines & Re-analysis Services

Ingestion Services

Cistrack Large Data Cloud Services

# App 3: Bulk Download of the SDSS

| Source | Destin. | LLPR* | Link | Bandwidth |
|--------|---------|-------|------|-----------|
| Chicago | Greenbelt | 0.98 | 1 Gb/s | 615 Mb/s |
| Chicago | Austin | 0.83 | 10 Gb/s | 8000 Mb/s |



**SDSS Data Distribution using Sector and UDT**

search [ ] (go)

Overview
Download Instructions
Software
Nodes Status
Downloading Records
Documentation
Technical Contact

**Related Links**
NCDM
SDSS
UDT
Sector
Teraflow Testbed

**Overview**

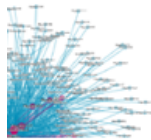Using this web site, you can download the Sloan Digital Sky Survey (SDSS) data if you have access to a high speed wide area network. For example, if your organization is attached to the National Lambda Rail or Internet2's Abilene Network, then you should be able to download the entire SDSS BESTDR5 catalog data set in less than five hours.

In general, it can be quite challenging to use effectively the available bandwidth over a wide area, high performance network. This project uses the UDP-based Data Transfer Protocol or UDT, which has been developed by the National Center for Data Mining (NCDM) at the University of Illinois at Chicago to make effective use of the bandwidth available from high performance wide area networks.

The project is supported by the National Science Foundation through the grant SCI II: The TeraFlow Project: High Performance Flows for Mining Large Distributed Data Archives, Award SCI-0430781.

**Sloan Digital Sky Survey (SDSS)**

The SDSS is systematically mapping a quarter of the entire sky, producing a detailed image of it, and determining the positions and absolute brightness of more than 100 million celestial objects. It is also measuring the distances to a million of the nearest galaxies, giving us a three-dimensional picture of the universe through a volume one hundred times larger than that explored to date. SDSS is also recording the distances to 100,000 quasars — the most distant objects known — giving us unprecedented knowledge of the distribution of matter to the edge of the visible universe.

The SDSS completed its first phase of operations — SDSS-I — in June, 2005. Over the course of five years, SDSS-I imaged more than 8,000 square degrees of the sky in five band passes, detecting nearly 200 million celestial objects, and it measured spectra of more than 675,000 galaxies, 90,000 quasars, and 185,000 stars. These data have supported studies ranging from asteroids and nearby stars to the large scale structure of the Universe.
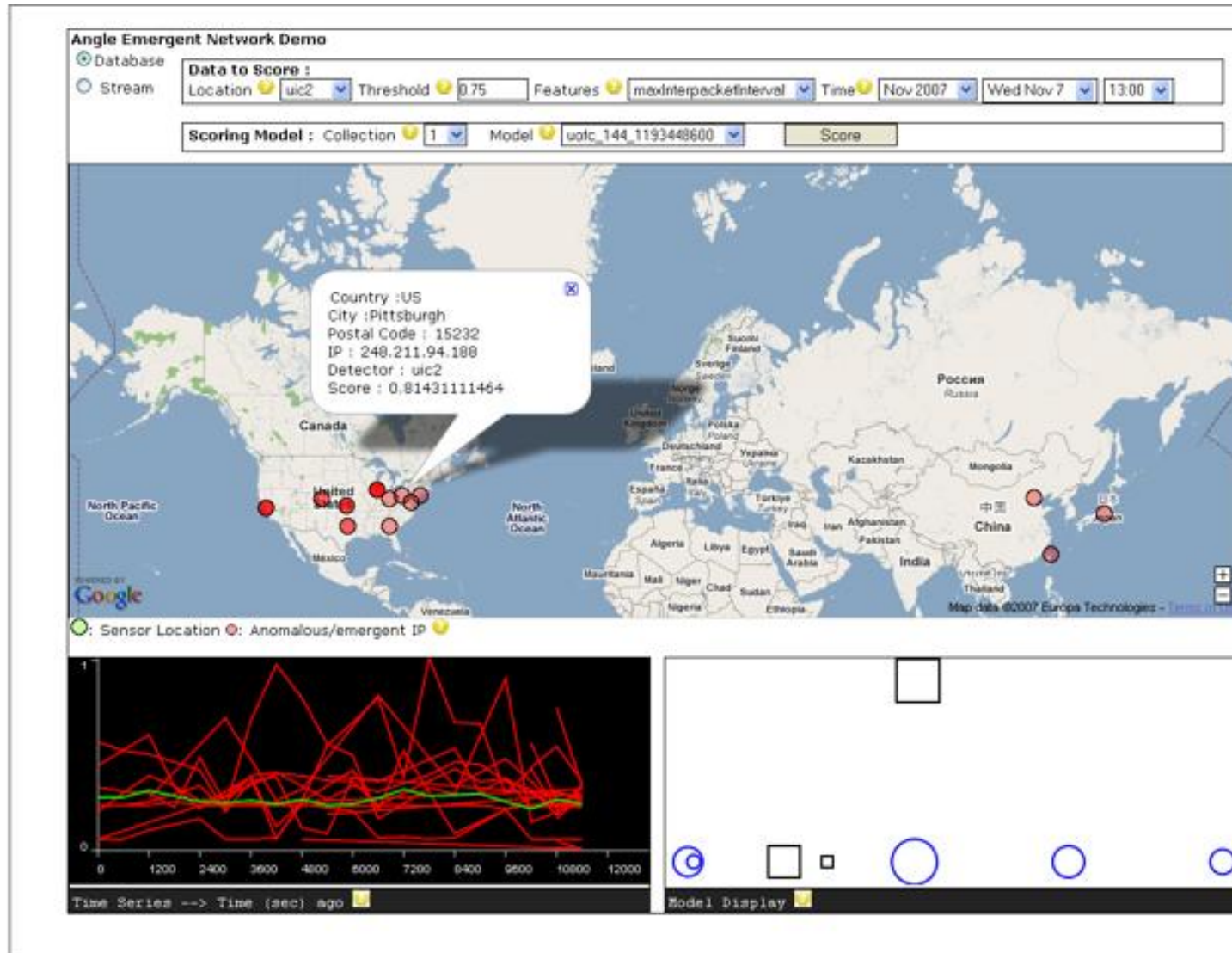
The most recent data product is DR6, which was released on June, 2007.

- LLPR = local / long distance performance
- Sector LLPR varies between 0.61 and 0.98

Recent Sloan Digital Sky Survey (SDSS) data release is 14 TB in size

# App 4: Anomalies in Network Data

# Sector Applications

- Distributing the 15 TB Sloan Digital Sky Survey to astronomers around the world (with JHU, 2005)

- Managing and analyzing high throughput sequence data (Cistrack, University of Chicago, 2007).

- Detecting emergent behavior in distributed network data (Angle, won SC 07 Analytics Challenge)

- Wide area clouds (won SC 09 BWC with 100 Gbps wide area computation)

- New ensemble-based algorithms for trees

- Graph processing

- Image processing (OCC Project Matsu)

# Credits



- Sector was developed by Yunhong Gu from the University of Illinois at Chicago and verycloud.com

# For More Information

For more information, please visit

sector.sourceforge.net

rgrossman.com (Robert Grossman)

users.lac.uic.edu/~yunhong (Yunhong Gu)