



Massive Database Analysis on the Cloud with D4M

Jeremy Kepner

William Arcand, William Bergeron, Chansup Byun, Matthew Hubbell, Ben Landon, Andrew McCabe, Craig McNally, Peter Michaleas, Andrew Prout, Tony Rosa, Delsey Sherrill, Albert Reuther, and Chuck Yee

This work is sponsored by the Department of the Air Force under Air Force contract FA8721-05-C-0002. Opinions, interpretations, conclusions and recommendations are those of the author and are not necessarily endorsed by the United States Government.



Outline

A large blue arrow pointing to the right, with a gradient from light blue on the left to dark blue on the right. It contains the text 'Introduction' in white.

- **Introduction**

- Technologies

- Results

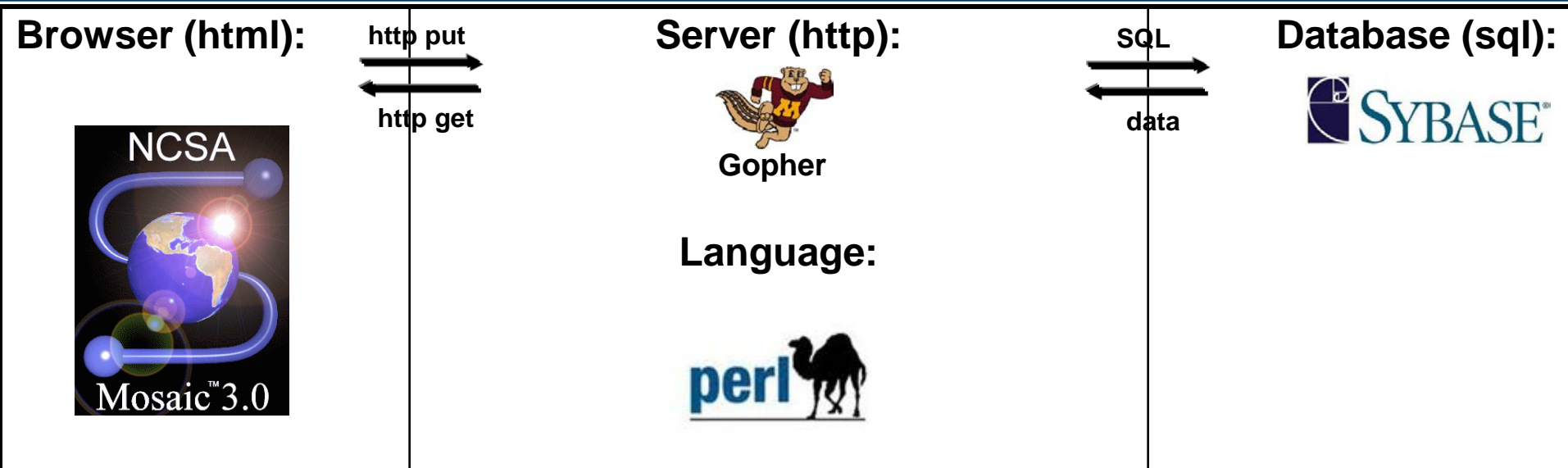
- Demo

- Summary

- *Webolution*
- *As is, is OK*
- *D4M*

Primordial Web

Kepner & Beaudry 1992, Visual Intelligence Corp (now GE Intelligent Platforms)

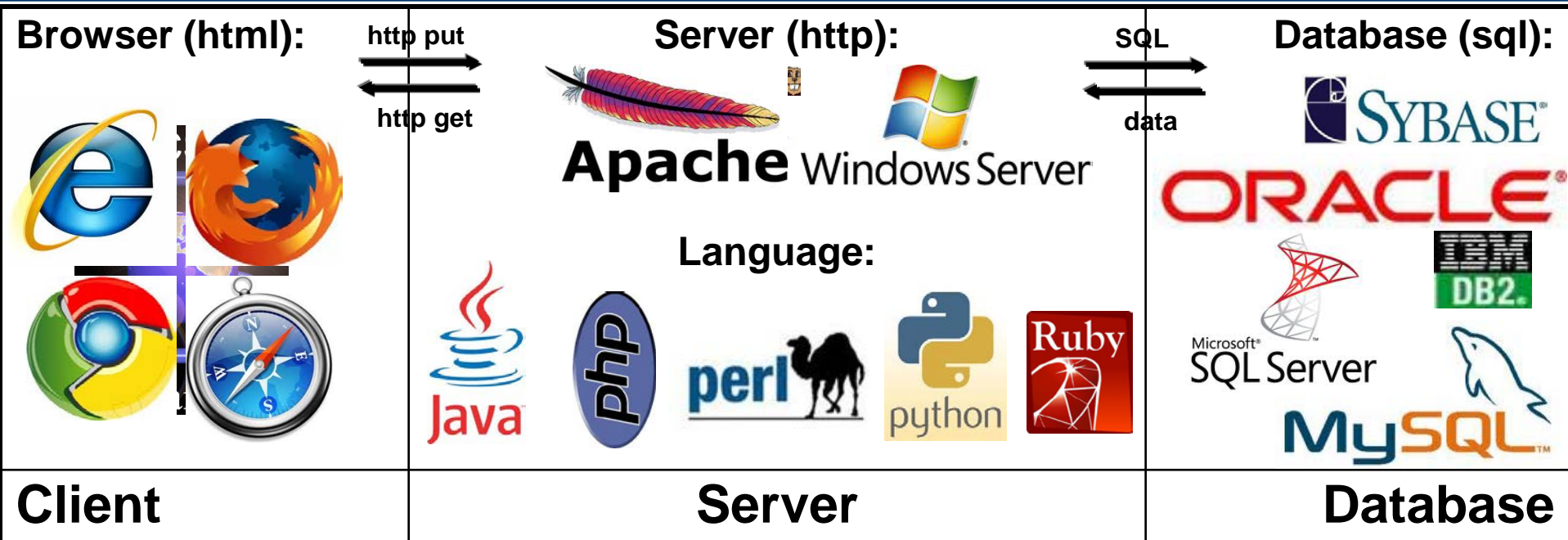


Client	Server	Database
---------------	---------------	-----------------



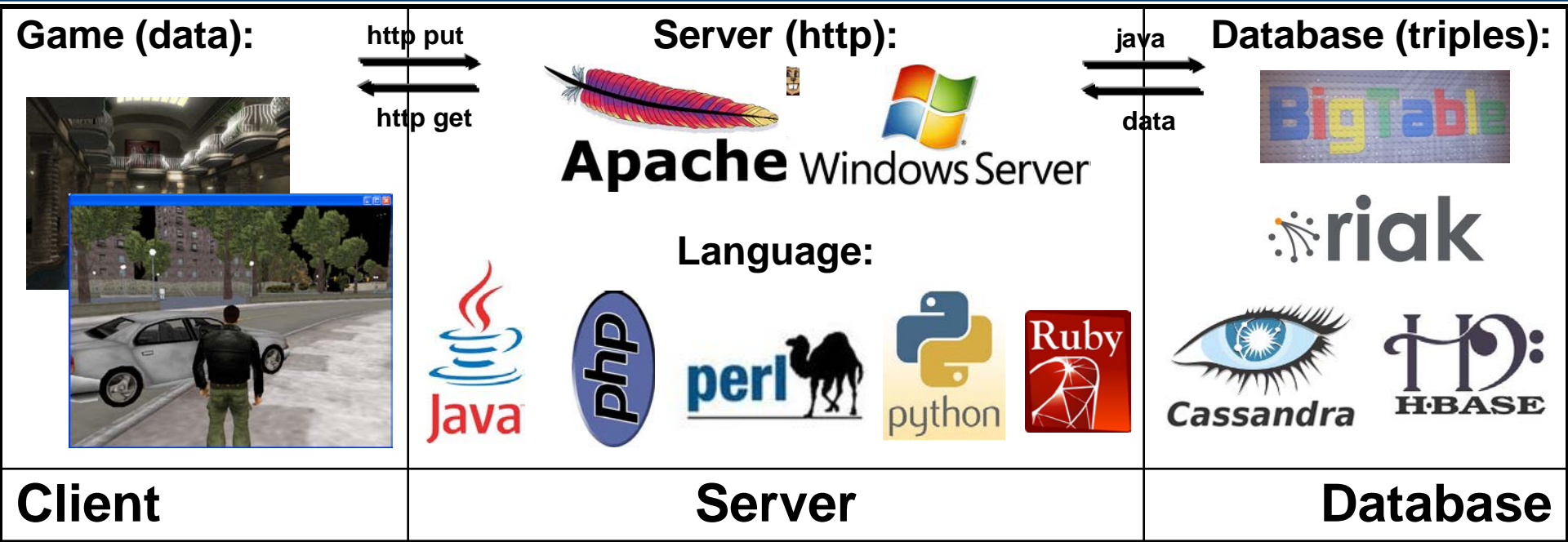
- Browser GUI? HTTP for files? Perl for analysis? SQL for data?
- A lot of work just to view data.
- Won't catch on.

Cambrian Web



- Browser GUI? HTTP for files? Perl for analysis? SQL for data?
- A lot of work to view a little data.
- ~~Won't catch on.~~

Modern Web



http put
←
http get

java →
←
data



- **Game GUI! HTTP for files? Perl for analysis? Triples for data!**
- **A lot of work to view a lot of data.**
- **Great view. Massive data.**

Modern Web



Client

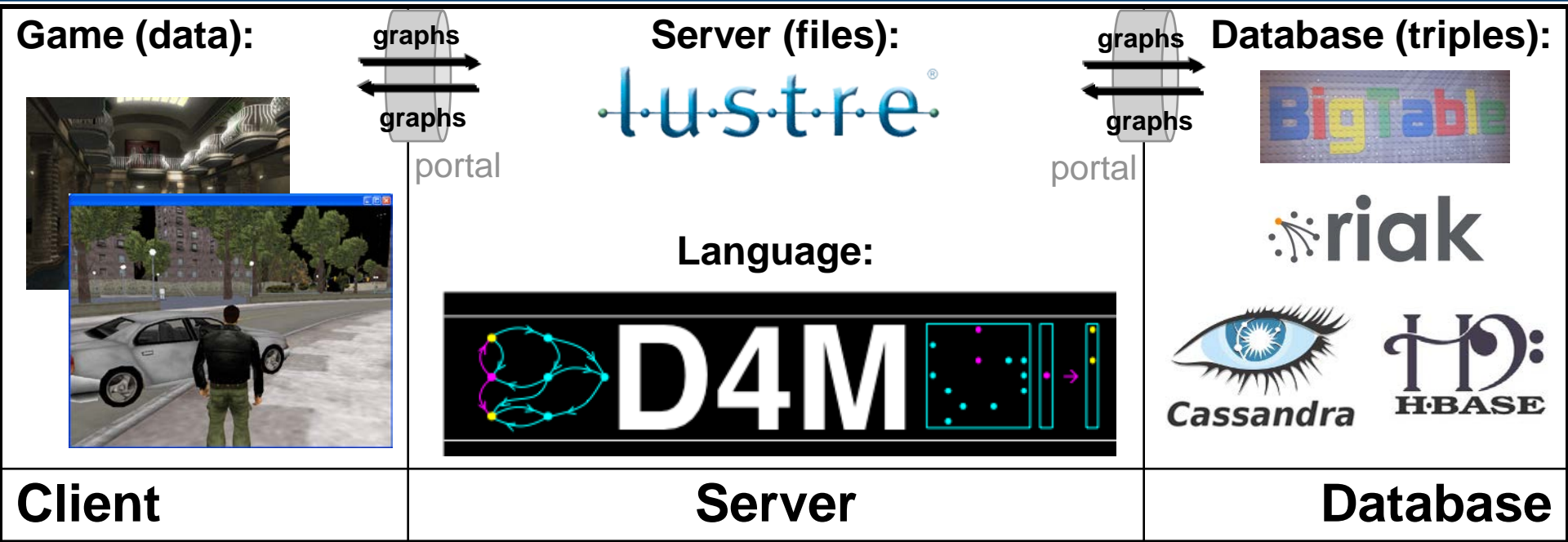
Server

Database



- **Game GUI! HTTP for files? Perl for analysis? Triples for data!**
- **A lot of work to view a lot of data. Missing middle.**
- **Great view. Massive data.**

Future Web?

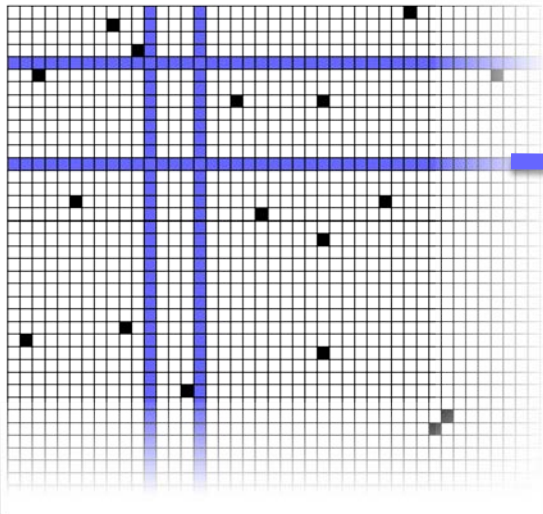


- **Game GUI! Fileserver for files! D4M for analysis! Triples for data!**
- **A little work to view a lot of data. Securely.**
- **Great view. Massive data.**

D4M: “Databases For Matlab”

Triple Store

Distributed Database



Triple store are high performance distributed databases for heterogeneous data

D4M

Dynamic
Distributed
Dimensional
Data
Model

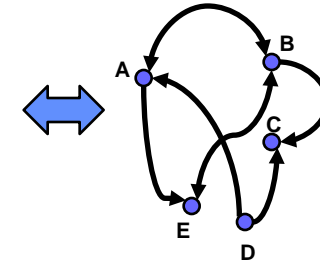
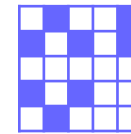
Query:

Alice
Bob
Cathy
David
Earl

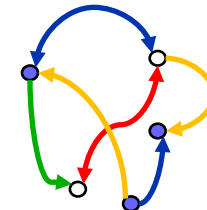
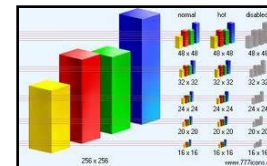


Associative Arrays

Numerical Computing Environment



A D4M query returns a sparse matrix or graph from Cloudbase...



...for statistical signal processing or graph analysis in MATLAB

D4M binds Associative Arrays to Triple Store, enabling rapid prototyping of data-intensive cloud analytics and visualization



Outline

- Introduction

- **Technologies**

- *Cloud Software*
- *D4M*

- Results

- Demo

- Summary

Recap: Cloud Computing Concepts

Data Intensive Computing

- **Compute architecture for large scale data analysis**
 - Billions of records/day, trillions of stored records, petabytes of storage
 - Google File System 2003
 - Google MapReduce 2004
 - Google BigTable 2006
- **Design Parameters**
 - Performance and scale
 - Optimized for ingest, query and analysis
 - Co-mingled data
 - Relaxed data model
 - Simplified programming
- **Community:**



Utility Computing

- **Compute services for outsourcing IT**
 - Concurrent, independent users operating across millions of records and terabytes of data
 - IT as a Service
 - Infrastructure as a Service (IaaS)
 - Platform as a Service (PaaS)
 - Software as a Service (SaaS)
- **Design Parameters**
 - Isolation of user data and computation
 - Portability of data with applications
 - Hosting traditional applications
 - Lower cost of ownership
 - Capacity on demand
- **Community:**



Recap: Cloud Computing Concepts

Data Intensive Computing

- Compute architecture for large scale data analysis
 - Billions of records/day, trillions of stored records, petabytes of storage
 - Google File System 2003
 - Google MapReduce 2004
 - Google BigTable 2006
- Design Parameters
 - Performance and scale
 - Optimized for ingest, query and analysis
 - Co-mingled data
 - Relaxed data model
 - Simplified programming
- Community:



Utility Computing

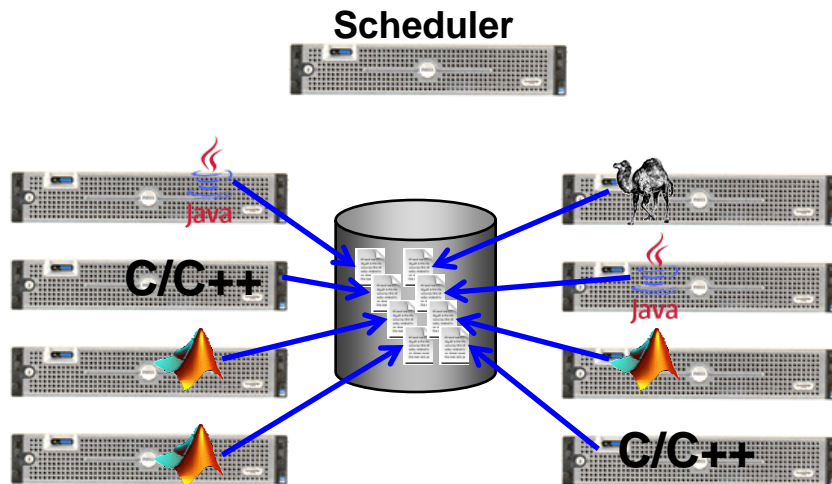
- Compute services for outsourcing IT
 - Concurrent, independent users operating across millions of records and terabytes of data
 - IT as a Service
 - Infrastructure as a Service (IaaS)
 - Platform as a Service (PaaS)
 - Software as a Service (SaaS)
- Design Parameters
 - Isolation of user data and computation
 - Portability of data with applications
 - Hosting traditional applications
 - Lower cost of ownership
 - Capacity on demand
- Community:



Advantages of Data Intensive Cloud

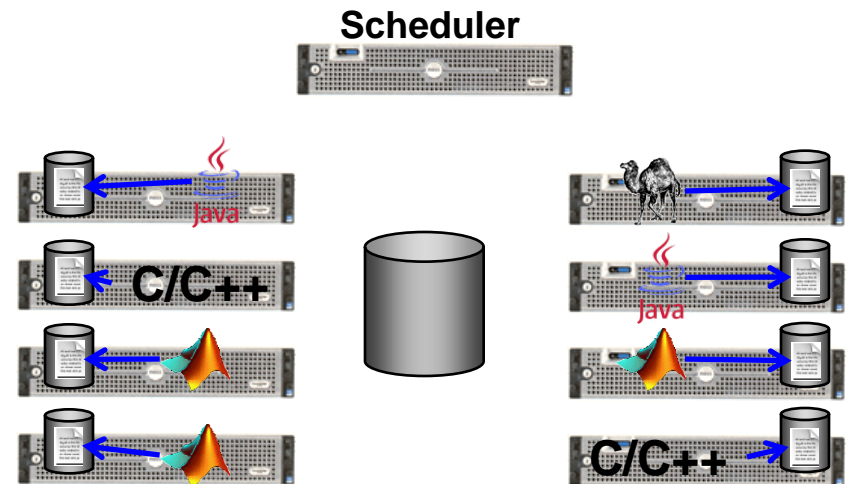
Traditional:

Data from central store to compute nodes



Cloud:

Data replicated on nodes, computation sent to nodes



- **Cloud computing moves computation to data**
 - **Good for applications where time is dominated by reading from disk**
- **Replaces expensive shared memory hardware and proprietary database software with cheap clusters and open source**
 - **Scalable to hundreds of nodes**



Distributed Cloud File Systems on TX-2500 Cluster

Service Nodes

Shared network storage



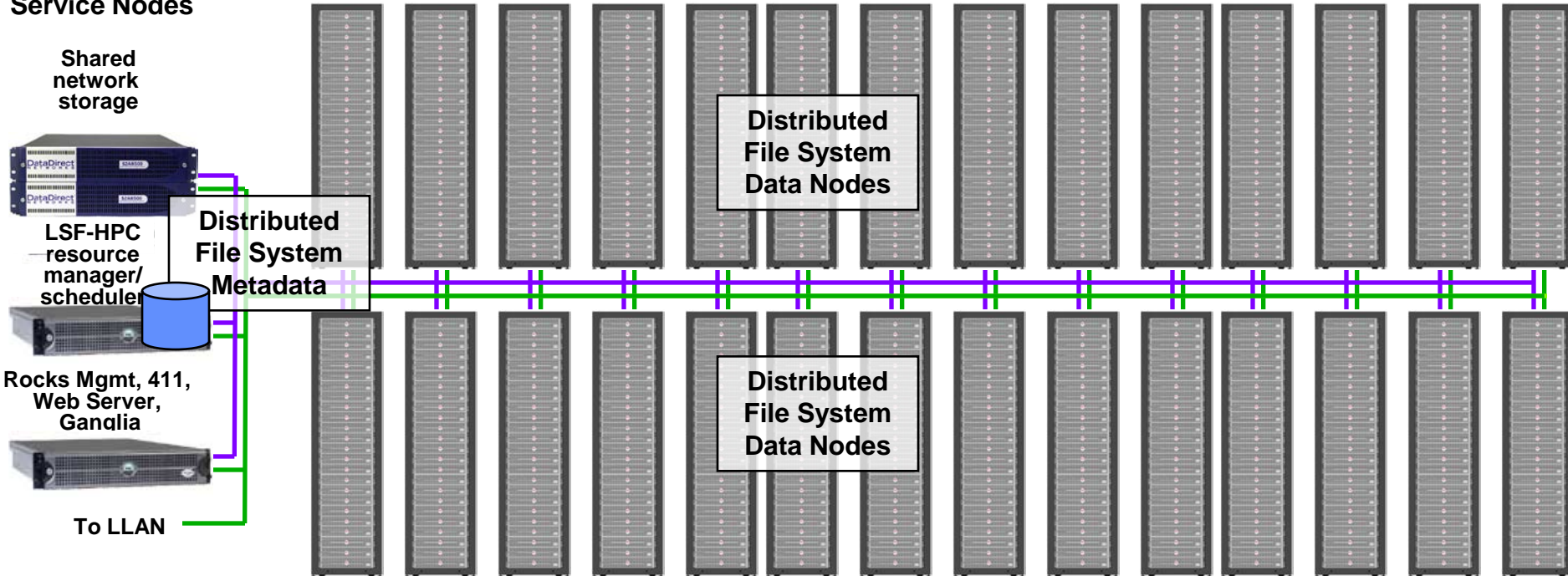
LSF-HPC resource manager/scheduler



Rocks Mgmt, 411, Web Server, Ganglia



To LLAN



432 **DELL** PowerEdge 2850



Dual 3.2 GHz EM64-T Xeon (P4)
 8 GB RAM memory
 Two Gig-E Intel interfaces
 Infiniband interface
 Six 300-GB disk drives

- 432+5 Nodes
- 864+10 CPUs
- 3.4 TB RAM
- 0.78 PB of Disk
- 28 Racks

MIT-LL Cloud	Hadoop DFS	Sector
Number of nodes used	350	350
File system size	298.9 TB	452.7 TB
Replication factor	3	2



Outline

- Introduction

- **Technologies**

- *Cloud Software*
- *D4M*

- Results

- Demo

- Summary

Multi-Dimensional Associative Arrays

- Extends associative arrays to 2D and mixed data types

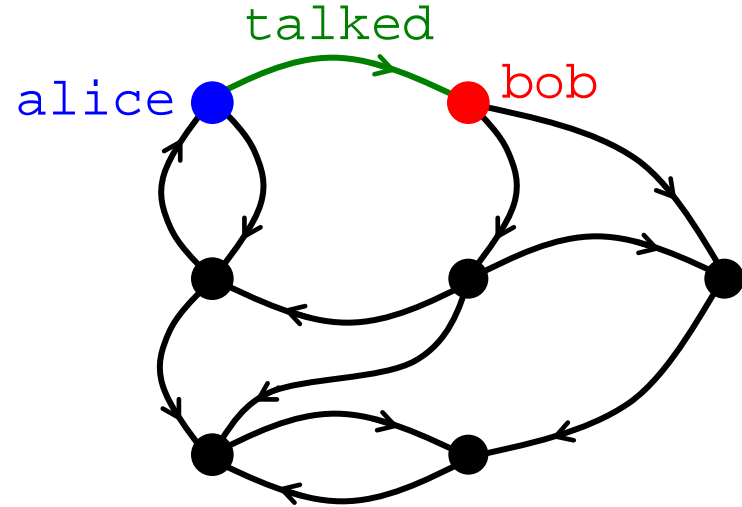
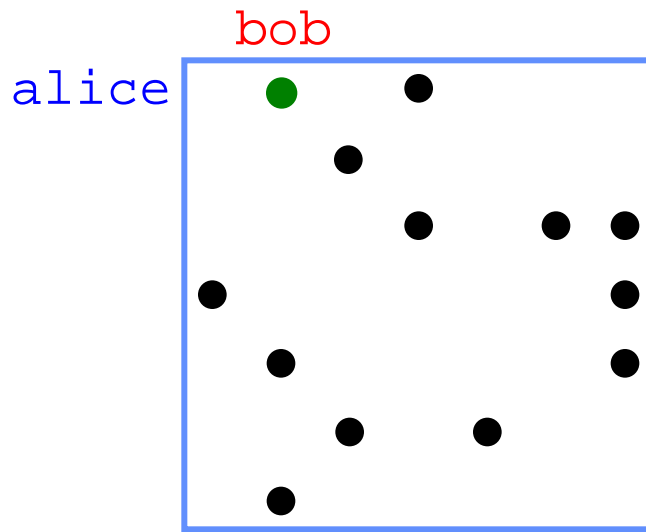
`A('alice ', 'bob ') = 'talked '`

or `A('alice ', 'bob ') = 47.0`

- Key innovation: 2D is 1-to-1 with triple store

`('alice ', 'bob ', 'talked ')`

or `('alice ', 'bob ', 47.0)`



- Associative arrays unify four viewpoints into one concept



Composable Associative Arrays

- **Key innovation: mathematical closure**
 - all associative array operations return associative arrays

- **Enables composable mathematical operations**

$A + B$ $A - B$ $A \& B$ $A | B$ $A * B$

- **Enables composable query operations via array indexing**

`A('alice bob ', :)` `A('alice ', :)` `A('al* ', :)`
`A('alice : bob ', :)` `A(1:2, :)` `A == 47.0`

- **Simple to implement in a library (~2000 lines) in programming environments with: 1st class support of 2D arrays, operator overloading, sparse linear algebra**

- **Complex queries with ~50x less effort than Java/SQL**
- **Naturally leads to high performance parallel implementation**

Associative Array Algebra

- Keys and values are from the infinite strict totally ordered set \mathbb{S}
- Associative array $A(\mathbf{k}) : \mathbb{S}^d \rightarrow \mathbb{S}$, $\mathbf{k}=(k^1, \dots, k^d)$, is a partial function from d keys (typically 2) to 1 value, where

$$A(\mathbf{k}_i) = v_i \quad \text{and} \quad \emptyset \text{ otherwise}$$
- Binary operations on associative arrays $A_3 = A_1 \oplus A_2$, where $\oplus = \cup_{f()}$ or $\cap_{f()}$, have the properties
 - If $A_1(\mathbf{k}_i) = v_1$ and $A_2(\mathbf{k}_i) = v_2$, then $A_3(\mathbf{k}_i)$ is

$$v_1 \cup_{f()} v_2 = f(v_1, v_2) \quad \text{or} \quad v_1 \cap_{f()} v_2 = f(v_1, v_2)$$
 - If $A_1(\mathbf{k}_i) = v$ or \emptyset and $A_2(\mathbf{k}_i) = \emptyset$ or v , then $A_3(\mathbf{k}_i)$ is

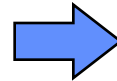
$$v \cup_{f()} \emptyset = v \quad \text{or} \quad v \cap_{f()} \emptyset = \emptyset$$

- High level usage dictated by these definitions
- Deeper algebraic properties set by the collision function $f()$
- Frequent switching between “algebras”

Universal “Exploded” Schema

Input Data

Time	src_ip	domain	dest_ip
2001-01-01	a		a
2001-01-02	b	b	
2001-01-03		c	c



Triple Store Table: Ttranspose

	2001-01-01	2001-01-02	2001-01-03
src_ip/a	1		
src_ip/b		1	
domain/b		1	
domain/c			1
dest_ip/a	1		
dest_ip/c			1



	src_ip/a	src_ip/b	domain/b	domain/c	dest_ip/a	dest_ip/c
2001-01-01	1				1	
2001-01-02		1	1			
2001-01-03				1		1

Triple Store Table: T

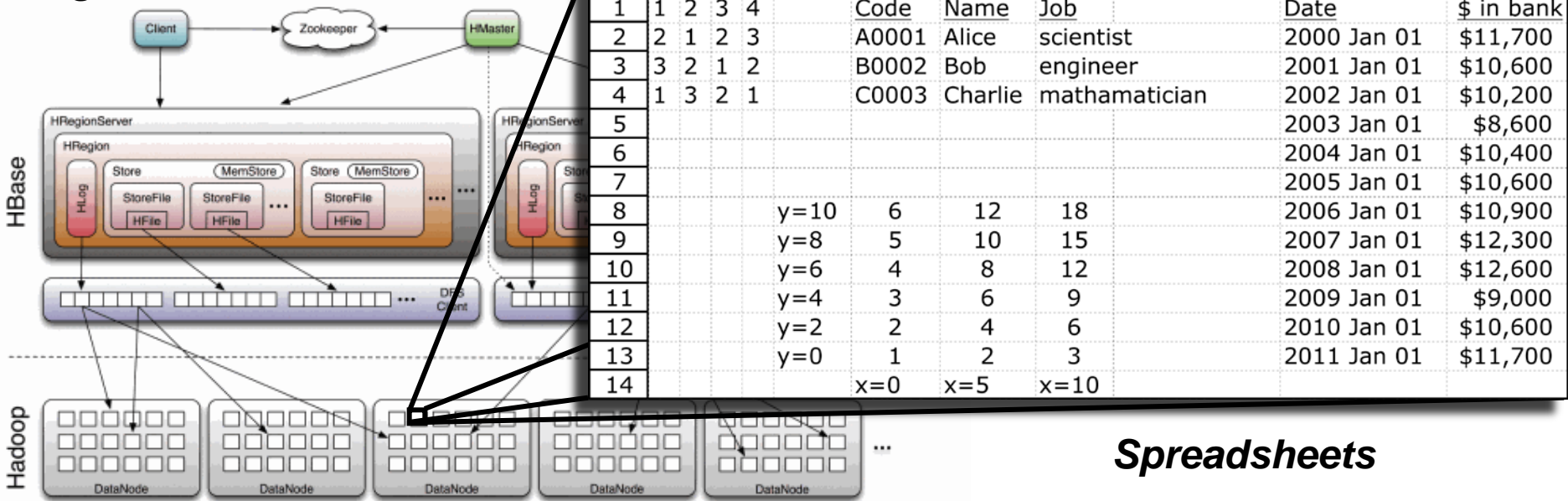
Key Innovations

- Handles all data into a *single* table representation
- Transpose pairs allows quick look up of *either* row or column



Unifies Spreadsheets and Big Tables

Big Tables



Spreadsheets

- **Spreadsheets** are the most commonly used analytical structure on Earth (100M users/day?)
- **Big Tables** (Google, Amazon, Facebook, ...) store most of the analyzed data in the world (Exabytes?)
- **Simultaneous** diverse data: strings, dates, integers, reals, ...
- **Simultaneous** diverse uses: matrices, functions, hash tables, databases, ...



Outline

- Introduction
- Technologies



- **Results**

- *Network Monitoring*
- *Text Query*
- *Insert Performance*

- Demo
- Summary

Stats Diagram

Triple Store Table: T

Row	Key (time)	src_ip/a	src_ip/b	src_ip/c	src_ip/d	domain/a	domain/b	domain/c	domain/d	dest_ip/a	dest_ip/b	dest_ip/c	dest_ip/d	Recv/a	Recv/b	Recv/c	Recv/d	Recv/e	
1	2001-10-01 01 01 00																		
2	2001-10-01 01 02 00																		
3	2001-10-01 01 03 00																		
4	2001-10-01 01 04 00																		
5	2001-10-01 01 05 00																		
6	2001-10-01 01 06 00																		

Associative Array: A

- Copy a set of rows from T into associative array A
- Perform the following statistical calculations on A
 - Column count: how many times each column appears in A
 - Column type count: how many times each column type appears in A
 - Column covariance: how many times a each pair of columns in A appear in the same row together
 - Column covariance: how many times a each pair of column types in A appear in the same row together

• Good for identifying column types, gaps, clutter, and correlations



Stats Implementation

- **Define a set of rows**

```
r = '2001-01-01 01 02 00,2001-01-01 01 03 00, 2001-01-01 01 04 00,'
```

- **Copy rows from table to associative array and convert '1' to 1**

```
A = double(logical(T(r,:)))  
A = A(:, 'src_ip/ * ,domain/ * ,dest_ip/ * ,')
```

- **Find popular columns counts**

```
sum(A,1) > 200
```

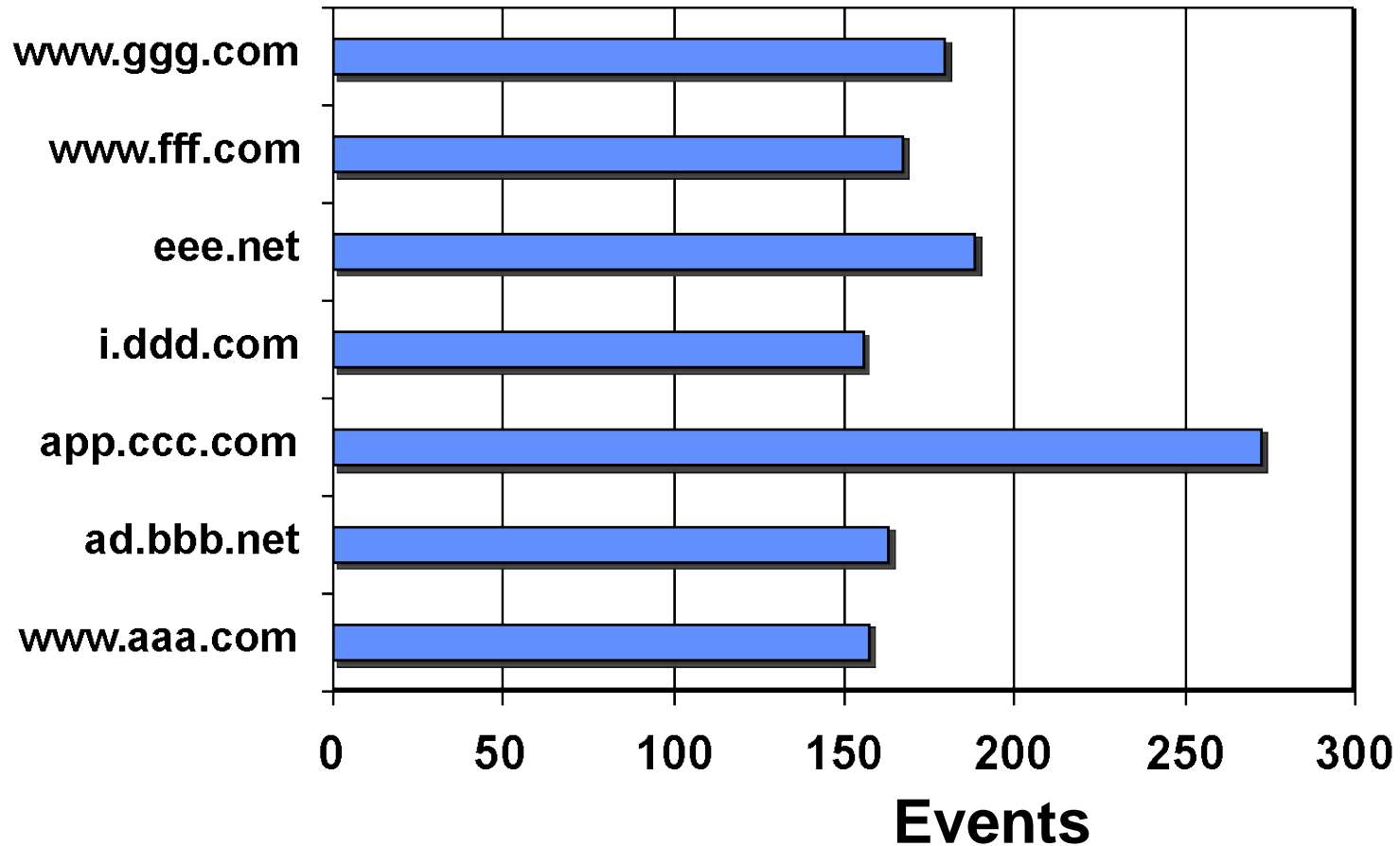
- **Find popular pairs**

```
A' * A > 200      or      sqIn(A) > 200
```

- **Find domains with many dest IPs**

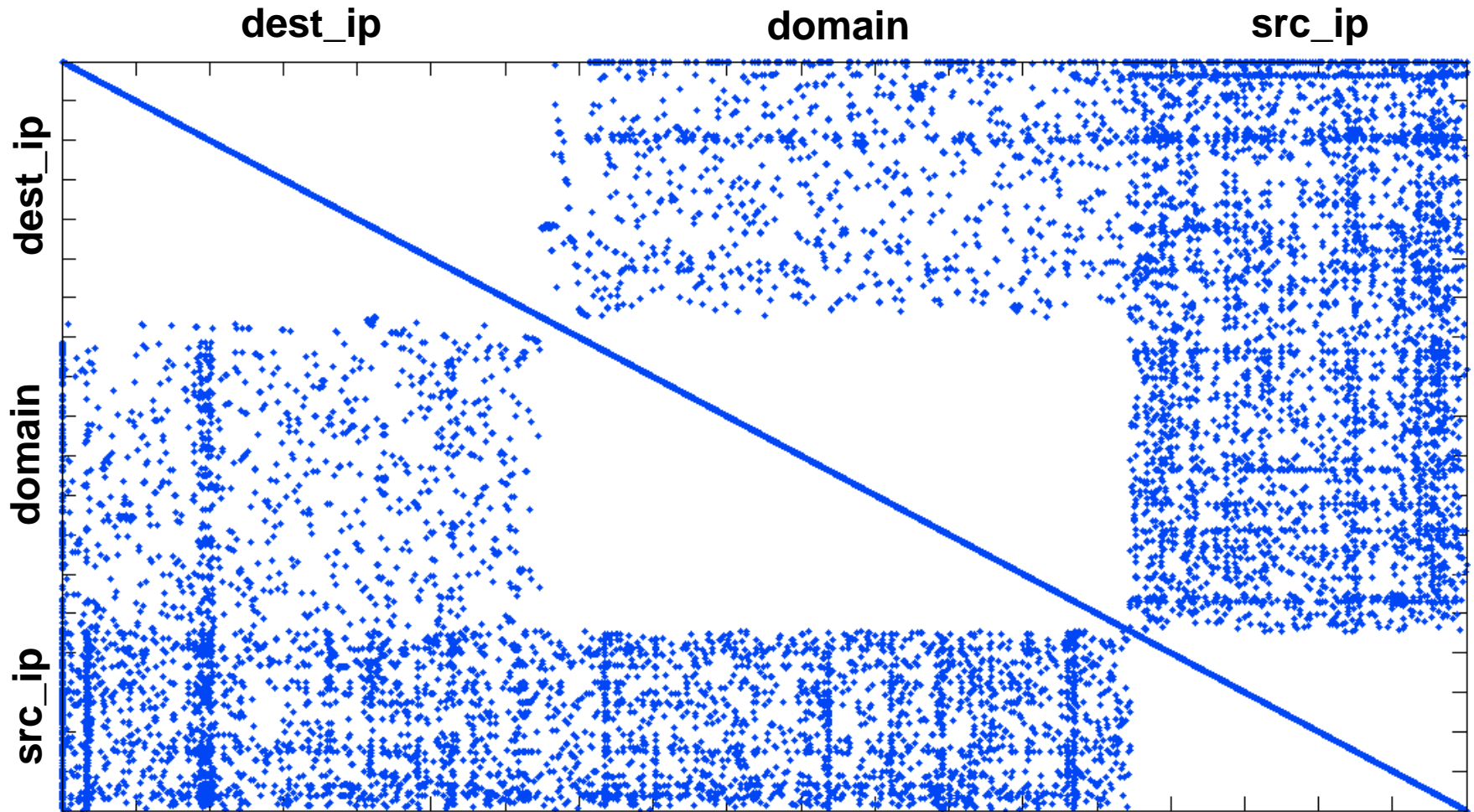
```
sum(double(logical(sqIn(A))),2) > 3
```

Count



- Very easy to get elementary count info necessary for finding clutter and anomalies

Covariance



- Adjacency matrix a natural result of covariance calculation



Facet Search

STRUCTURED KNOWLEDGE SPACE

Search Upload File Status Dictionary Update

DOCUMENT SEARCH

afghanistan Search

PEOPLE

- AHMAD SHAH MASOOD (60)
- ABDUL RASHID (55)
- OSAMA BIN LADEN (14)
- CHRIS BIRD (12)
- ALEXANDER LE (1)
- 15 more...

LOCATIONS

- AFGHANISTAN (297)
- KABUL (147)
- PAKISTAN (134)
- TAJIKISTAN (66)
- MOSCOW (64)
- 15 more...

ORGANIZATIONS

- UNITED NATIONS (105)
- AFGHAN ISLAMIC PRESS (31)
- NORTH ATLANTIC TREATY ORGANI... (25)
- THE TALIBAN (22)
- UNITED NATIONS HIGH COMMISSI... (17)
- UNITED NATIONS SECURITY COUN... (17)
- AL QAEDA (16)
- INTERNATIONAL RED CROSS (16)
- CENTRAL INTELLIGENCE AGENCY (15)
- UNITED STATES ARMY (14)
- 10 more...

SELECTOR B

SELECTOR C

SELECTOR D

SELECTOR E

SELECTOR F

SELECTOR G

TEXT

- BECKY.SIU@MND-B.ARMY.MIL (1)
- BRAD.ANDERSON@IRAQ.CENTCOM.M...
- BRIAN.SMITH@UNDP.ORG (1)
- BUSCHM@SAF-HQ.NATO.INT (1)
- COSTAP@CFC.AFGN.ARMY.MIL (1)
- 15 more...

- Core analytic of SKS
- Gives keyword distribution of a set of documents that share a common keyword(s)
 - Provides useful guide to what keyword to select next
- Currently implemented with several hundreds of lines of Java/SQL
- Associative array implementation has 1 line

Facet Search Algorithm

	NY	DC	IMF	UN	Alice	Bob	Carl
a.txt		●		●			
b.doc			●				
c.pdf				●		●	●
d.htm	●	↓		↓		↓	↓
e.ppt		●		●			●
f.txt			●		●		
g.doc		●					
		1		2		1	2

- Associative array relates documents to place, org and person entities

$$A(x,y) : \mathbf{S}^{N \times M} \rightarrow \mathbf{R}$$

- Facets $y_1=UN$, $y_2=Carl$

- Documents that contain both

$$\underline{A}(:,y_1) \ \& \ \underline{A}(:,y_2)$$

- Entity counts in the above set of documents obtained via matrix multiply

$$(\underline{A}(:,y_1) \ \& \ \underline{A}(:,y_2))^t \ A$$



Graph Analysis/Graph500 Benchmark

- Scalable benchmark specified by graph community
- Goal
 - Stress parallel computer architecture
- Key data
 - Very large Power law graph
- Kernels
 - Data Generator
 - Ingest

GraphAnalysis.org: High Performance Computing for solving large-scale graph problems

HPC Graph Analysis

- [Home](#)
- [News](#)
- [Benchmark](#)
- [Results](#)
- [Publications](#)
- [People](#)
- [Links](#)
- [Contact](#)

Benchmark

Overview

We present a graph theory benchmark representative of computational kernels in computational biology, complex network analysis, and national security. This benchmark is based on the HPCS Scalable Synthetic Compact Applications graph analysis (SSCA#2) benchmark. SSCA#2 is characterized by integer operations, a large memory footprint, and irregular memory access patterns. It has multiple kernels accessing a single data structure representing a weighted, directed multigraph. In addition to a kernel to construct the graph from the input tuple list, there are three additional computational kernels to operate on the graph. Each of the kernels requires irregular access to the graph's data structure, and it is possible that no single data layout will be optimal for all four computational kernels.

References

- SSCA#2 v2.0 Specification: [pdf doc](#)
- Sequential code: [C](#)
- Parallel code: [C/OpenMP](#)
- Matlab version: [ssca2v2-061107.tar.gz](#)

Home Complete Results

GRAPH 500

The Graph 500 List

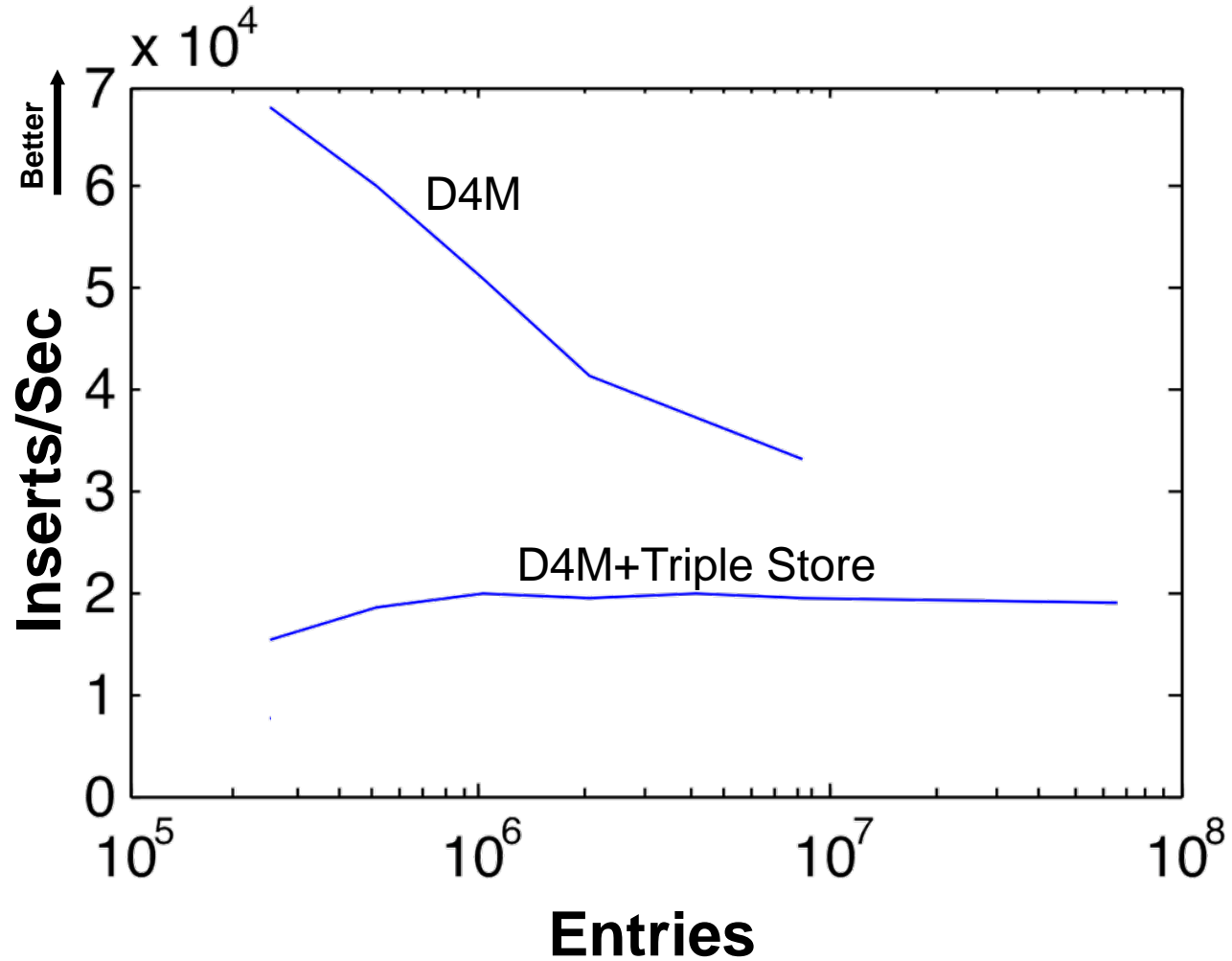
Top 10 (June 2011)

Rank	Machine
1	Intrepid (BG/P, 32768 nodes/ 131072 cores)
2	Jugene (IBM, 32k nodes)
3	Lomonosov (MPP, 4096 nodes/ 8192 cores)

Brief Introduction

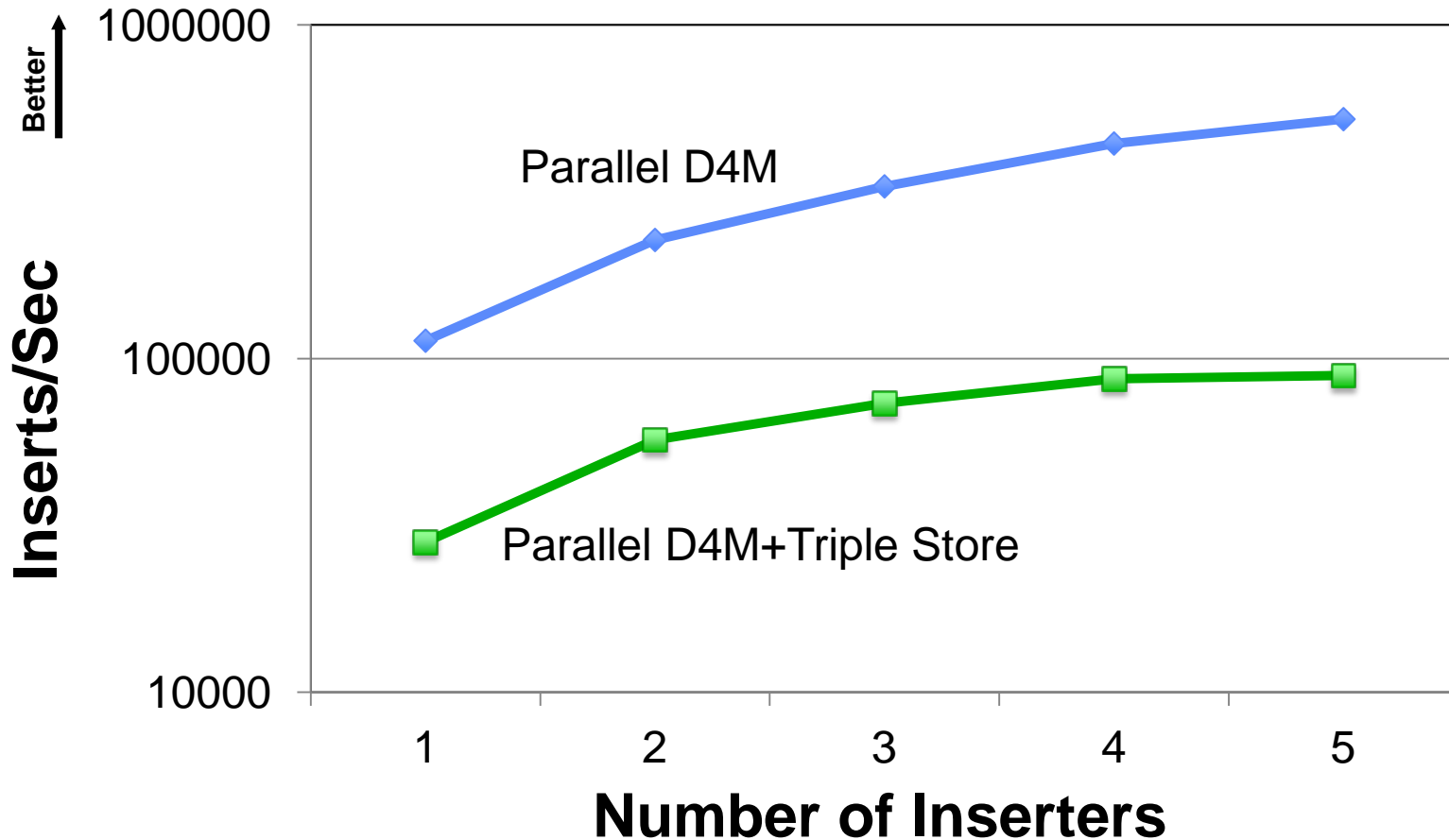
Data intensive supercomputer applications are increasingly important for HPC workloads, but are ill-suited for platforms designed for 3D physics simulations. Current benchmarks and performance metrics do not provide useful information on the suitability of supercomputing systems for data intensive applications. A new set of benchmarks is needed in order to guide the design of hardware architectures and software systems intended to support such applications and to help procurements. Graph algorithms are a core part of many analytics workloads.

Single Node Insert Rate





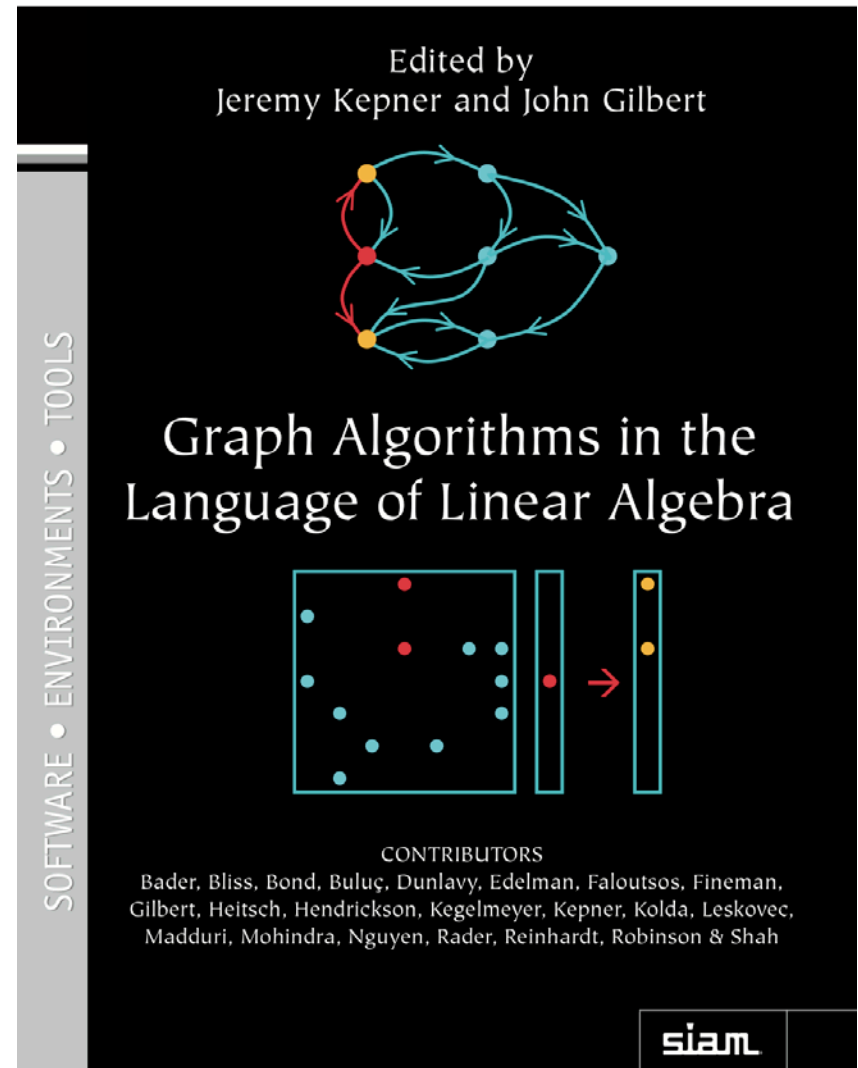
Parallel D4M + Single Node Triple Store Graph500 insert rate



- Parallel D4M provides 500K memory inserts/sec on 1 node
- Parallel D4M + Triple Store provides 100K disk inserts/sec on 1 node
- Parallel D4M enables exploiting full power of Triple Store from D4M

Reference

- **Book: “Graph Algorithms in the Language of Linear Algebra”**
- **Editors: Kepner (MIT-LL) and Gilbert (UCSB)**
- **Contributors**
 - Bader (Ga Tech)
 - Bliss (MIT-LL)
 - Bond (MIT-LL)
 - Dunlavy (Sandia)
 - Faloutsos (CMU)
 - Fineman (CMU)
 - Gilbert (UCSB)
 - Heitsch (Ga Tech)
 - Hendrickson (Sandia)
 - Kegelmeyer (Sandia)
 - Kepner (MIT-LL)
 - Kolda (Sandia)
 - Leskovec (CMU)
 - Madduri (Ga Tech)
 - Mohindra (MIT-LL)
 - Nguyen (MIT)
 - Rader (MIT-LL)
 - Reinhardt (Microsoft)
 - Robinson (MIT-LL)
 - Shah (UCSB)





Summary

- **Web evolution has resulted in a new class of technologies for**
 - **Display (game interfaces)**
 - **Analysis (D4M)**
 - **Storage (triple stores)**

- **D4M is a novel technology that allows complex analytics to be implement with significantly less effort than traditional approaches**

- **D4M is built on composable associative arrays which admit linear algebraic manipulation**