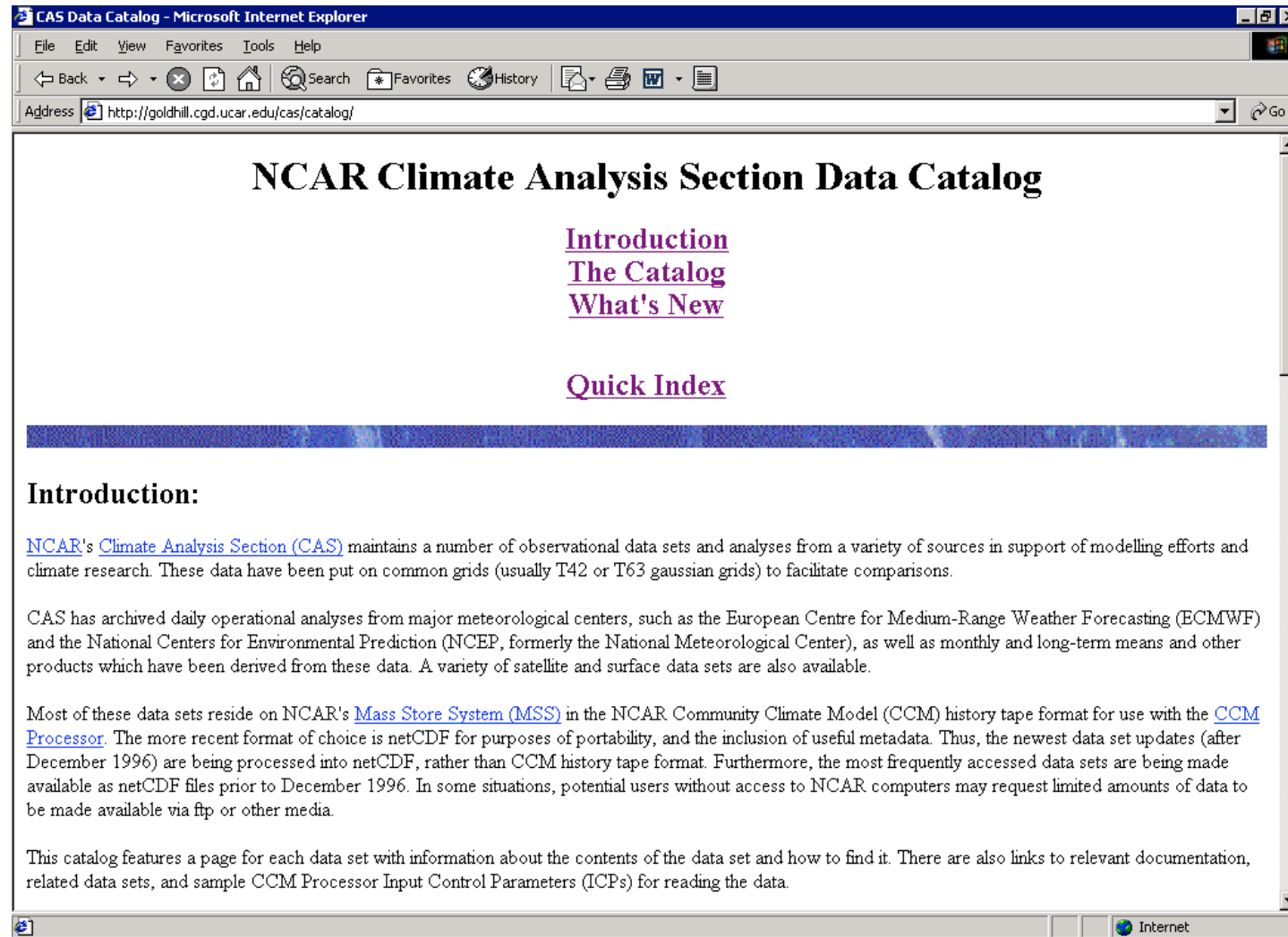# The Analysis & Mining of Globally Distributed Data

## Chapter 1. A Quick Introduction to Data Grids, Data Webs, Semantics Webs, & Distributed Data Mining

Robert Grossman
Laboratory for Advanced Computing
University of Illinois at Chicago

&

Open Data Partners

# 1.1 Background

## Three Fundamental Trends

# Trend 1. Explosion of Data...

# Distributed Exabytes



Source: IDC (1999) "1999 Winchester Disk Drive Market Forecast and Review"

# ... All in the Wrong Format

# The Data Gap



The Data Gap

Total new disk (TB) since 1995

New Ph.D.s

4,000,000
3,500,000
3,000,000
2,500,000
2,000,000
1,500,000
1,000,000
500,000
0

1995    1996    1997    1998    1999

# Trend 2: Most Data is Distributed



Some else's data is more valuable than some else's cycles.

θ  Pearson's Law: The usefulness of a column of data varies as the square of the number of columns it is compared to.

# Example: ENSO & Cholera



El Nino Data at NCAR

Cholera Data at WHO

# Trend 3. Bandwidth is a Commodity.



For the first time, Gigabytes can be moved in minutes.

# Trend 3. Bandwidth is a Commodity.

For the first time, Gigabytes can be moved in minutes.

# Gigabytes can be Moved in Minutes

# Example 1: Data Grids – National Virtual Observatory

# NASA SkyView

# Example 2: Data Webs – Molecular Data Space

# Example 3:
# Semantic Web – DAML

# Example 4: Data Mining – Sky Survey Catalog

θ Goal: To predict class (star or galaxy) of sky objects, based on survey images (from Palomar Observatory)

– 3000 images with 23K x 23K pixels/image.

θ Approach:

– Partition the image & create 40 features

– Build a classification model

– Success Story: Could find 16 new high red-shift quasars, some of the farthest objects that are difficult to find.

# Classifying Galaxies

**Early**



**Class:**
- **Stages of Formation**

**Intermediate**



**Attributes:**
- **Image features,**
- **Characteristics of light waves received, etc.**

**Late**



**Data Size:**
- **72 million stars, 20 million galaxies**
- **Object Catalog: 9 GB**
- **Image Database: 150 GB**

From Fayyad, et.al., Advances in Knowledge Discovery and Data Mining, 1996

# 1.2 Four Different Philosophies for Working with Distributed Data

## Data Grids, Data Webs, Semantic Webs, & Distributed Data Mining

# Data Grids –
# Categorical Imperative



grid

θ How can we *interoperate* distributed supercomputers?

# Data Webs –
# Categorical Imperative

θ How can we *explore* other people's data?

# Semantic Web – Categorical Imperative

```
<City rdf:ID="AABS">
   <geolocationCode>AABS</geolocationCode>
   <name>AABENRAA</name>
   <installationTypeCode>CTY</installationTypeCod
   <primeGeoloc rdf:resource="#AABQ" />
   <longitude>0092600E</longitude>
   <latitude>550300N</latitude>
</City>
```

DAML Geofile

What is the distance between Chicago and Baltmore?

θ   How can we extend the web to support *knowledge?*

# Data Mining – Categorical Imperative

| Longitude | Latitude | Time | Cloud Cover |
|---|---|---|---|
| 19.6875 | -12.557755 | 120.0 | 0.43481043 |
| 75.9375 | -12.557755 | 120.0 | 0.9641479 |
| 132.1875 | -12.557755 | 120.0 | 0.82385314 |
| 188.4375 | -12.557755 | 120.0 | 0.91212153 |
| 244.6875 | -12.557755 | 120.0 | 0.8691945 |
| 300.9375 | -12.557755 | | 0.36265105 |
| 357.1 | | | 0.52140427 |
| 5 | | | 0.9674745 |
| 1 | | | 0.9179723 |

Are there any patterns relating cloud cover & biodiversity?

θ How can we find *patterns* in distributed data?

# Technologies for Global Data



**Object**

| | View | Mine/Discover | Compute |
|---|---|---|---|
| Knowledge | Digital Libraries | Knowledge Mining | Semantic Web |
| Data | Web-based databases | Data Webs | Data Grids |
| Files | Persistent Archives | Distributed Data Mining | Grids |

**Action**

# Data Grids vs. Data Webs

Browsing &
Casual
Exploration

Data Webs
- Searching
- Exploration
- Casual correlation

Collaborations

Data Grids
- Security
- Authorization
- Scheduling

Distributed
Computer

Web Based
Computing

# Data Grids, Data Mining & Data Webs

|  | Data Grid | Distributed Data Mining | Data Web |
|---|---|---|---|
| Goal | distributed computation | distributed data mining | data explor. & mining |
| Services | authorization, security, resources | building models, transforming data, etc. | publishing, merging, & correlating columns |
| Protocol | TCP, GridFTP | TCP | DWTP, … |
| Platform | dist. clusters | server | dist. cluster |

# Semantic Web vs. Data Web

|  | Document Web | Semantic Web | Data Web |
|---|---|---|---|
| Protocol | HTTP | HTTP, SOAP | DWTP, SOAP |
| Languages | HTML, XML | XML, RDF | XML, PMML … |
| Action | keyword search | RDF inferences | correlate and mine |
| Platform | server | server | server, cluster |

# What is a Petabyte?

θ HSS Camp:
- $10^{15}$ bytes
- Tertiary storage, data migration, data staging, …

θ Data Grid Camp:
- Thousand TB data sets with AAA
- Security, authorization, task scheduling, replication management, …

θ Data Web Camp:
- Hundred Million 10 MB / Million 1 GB open data sets
- Discovery, correlation, normalization, transform, …

# Overview of Course

1. Introduction
2. Introduction to Data Mining
3. Protocols and Stacks
4. Web Services
5. Data Grids
6. Data Webs

# Sources for Some of the Slides

θ Semantic web and web services: Isabel Cruz,
  SC 02 Tutorial Notes

θ Data grids: Introduction to Grid Computing
  and Globus Toolkit, www.globus.org