# Massive Cloud Auditing using Data Mining on Hadoop

## Prof. Sachin Shetty

### CyberBAT Team, AFRL/RIGD

### AFRL VFRP
### Tennessee State University

# Outline

- ➢ Massive Cloud Auditing
- ➢ Traffic Characterization
- ➢ Distributed Data Storage and Online Querying
- ➢ Online Data Mining
- ➢ System Architecture
- ➢ Prototype
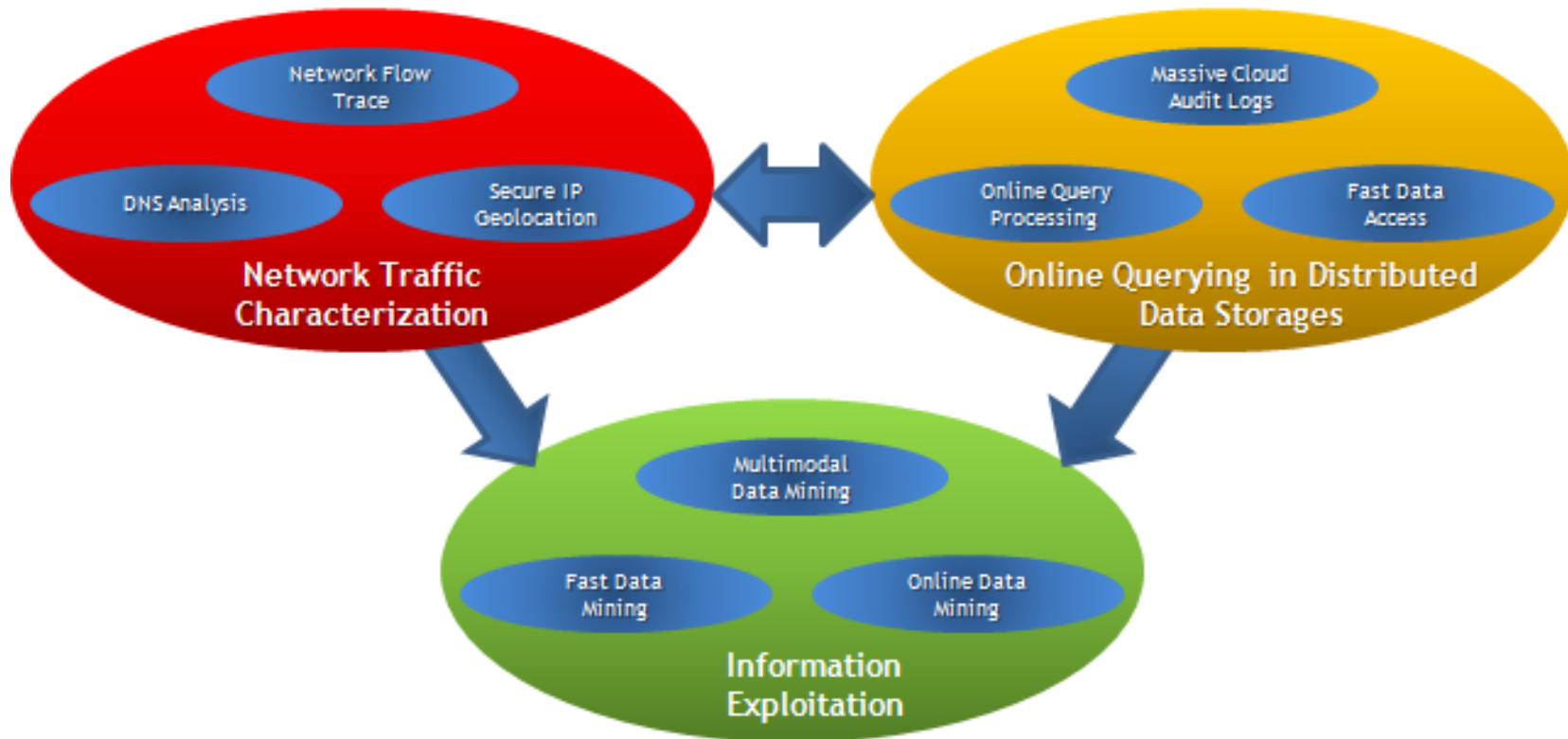- ➢ Work Accomplished
- ➢ Future Work

# Massive Cloud Auditing

- Cloud Auditing of massive logs requires analyzing data volumes which routinely cross the peta-scale threshold.

- Computational and storage requirements of any data analysis methodology will be significantly increased.

- Distributed data mining algorithms and implementation techniques needed to meet scalability and performance requirements entailed in such massive data analyses.

- Current distributed data mining approaches pose serious issues in performance and effectiveness in information extraction of cloud auditing logs.

- Reasons include  scalability, dynamic and hybrid workload, high sensitivity,  and stringent time constraints.

DISTRIBUTION C: Distribution authorized to U.S. Government Agencies and their Contractors, 07/11/11. Other requests for this document must be referred to Air Force Research Laboratory/RIGD, 525 Brooks Rd., Rome NY 13441-4505.

Keesook.Han@rl.af.mil

INFORMATION INSTITUTE WORKSHOP on ASSURING THE CLOUD, 11 July 2011, GRIFISS Institute, Rome, NY    POC: AFRL/RIGD Dr. Keesook J. Han

# Cloud Auditing Challenges



**Network Traffic Characterization**
- Network Flow Trace
- DNS Analysis
- Secure IP Geolocation

**Online Querying in Distributed Data Storages**
- Massive Cloud Audit Logs
- Online Query Processing
- Fast Data Access

**Information Exploitation**
- Multimodal Data Mining
- Fast Data Mining
- Online Data Mining

## *Data Mining Approaches*

# Traffic Characterization

- Cloud Traffic logs accumulated from diverse and geographically disparate sources.

- Sources include stored and live traffic from popular web applications: Web and Email

- Live Packet Capture from packet sniffing tools (Wireshark)

- Honeypot traffic from UTSA comprising of malicious traffic.

- Augment traffic information with IP,DNS, geolocation analysis procured from publicly available datasets and high level network and flow statistics.

DISTRIBUTION C: Distribution authorized to U.S. Government Agencies and their Contractors, 07/11/11. Other requests for this document must be referred to Air Force Research Laboratory/RIGD, 525 Brooks Rd., Rome NY 13441-4505.

Keesook.Han@rl.af.mil

INFORMATION INSTITUTE WORKSHOP on ASSURING THE CLOUD, 11 July 2011, GRIFISS Institute, Rome, NY    POC: AFRL/RIGD Dr. Keesook J. Han

# Traffic Characterization

- IP geolocation from public databases retrieve the name and street address of the organization which registered the address block. For large ISPs the registered street address usually differs from the real location of its hosts.

- Measurement based IP geolocation utilize active packet delay measurements to approximate the geographical location of network hosts.

- Secure IP geolocation to defend against adversaries manipulating packet delay measurements to forge locations

- Traffic characterization will generate massive amount of data. Need for distributed data storage.

# Distributed Data Storage and Online Querying

- Hadoop based distributed data storage and online querying solution.

- **Hadoop**: Software library which supports distributed processing of large data sets across clusters of computers using a simple programming model. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage

- Hbase and Hive Distributed Data Storage

- Chukwa: Distriubted data collection system based on Hadoop system.

Keesook.Han@rl.af.mil
POC: AFRL/RIGD Dr. Keesook J. Han

# Distributed Data Storage and Online Querying

- Relatively high latency in Hadoop system which is not desirable for online query processing.

- Design middleware to support structured and unstructured data and multiple data storage(HIVE and HBASE)

- Design appropriate data structure for the auditing data to achieve quick retrieve and access of large data.

- Develop efficient data collection system to reduce the amount of network disk access.

- Develop algorithms and implementation techniques for both batch processing and online query processing of cloud audit data based on Hadoop distributed system.

DISTRIBUTION C: Distribution authorized to U.S. Government Agencies and their Contractors, 07/11/11. Other requests for this document must be referred to Air Force Research Laboratory/RIGD, 525 Brooks Rd., Rome NY 13441-4505.

Keesook.Han@rl.af.mil
INFORMATION INSTITUTE WORKSHOP on ASSURING THE CLOUD, 11 July 2011, GRIFISS Institute, Rome, NY    POC: AFRL/RIGD Dr. Keesook J. Han

# Online Data Mining

- Develop data mining algorithms that work in a massively parallel and yet online fashion for mining of large data streams

- Reducing time between query submission and obtaining results.

- Overall speed of query processing depends critically on the query response time

- Map-Reduce programming model used for fault-tolerant and massively parallel data crunching .

- But Map-Reduce implementations work only in batch mode and do not allow stream processing or exploiting of preliminary results.
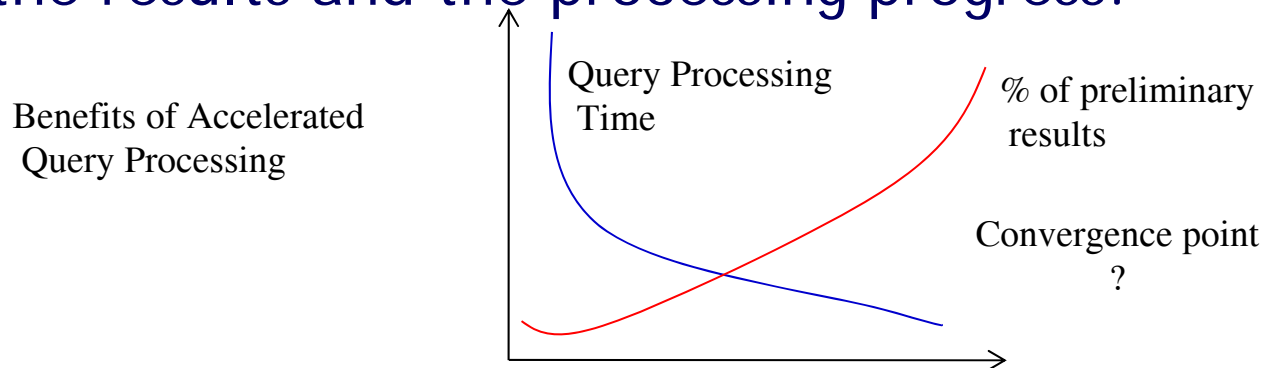
# Online Data Mining

- Online Aggregation: Reduce query processing time by monitoring of preliminary results of an aggregation query along with estimates on the result's error intervals.

- Proposed Solution: Develop game theoretic model to combine Map-Reduce parallelization and online aggregation to reduce query processing time.

- Lack of memory to store preliminary data mining results when processing streams, and expensive computation for growing size of preliminary result grows and frequent input updates

- Both problems can be solved by limiting the number of "historical" results to be stored.

# Online Data Mining

- Balance the risks of using a preliminary result against the benefits of an accelerated computation, two questions need to be answered:

- "What is the likelihood that the query result is likely to change if all data is processed?" and "How much time can be saved if the current result is used?".

- Game theoretic approach to answer these questions by providing mechanisms for monitoring the convergence of the results and the processing progress.

Benefits of Accelerated
Query Processing

Query Processing
Time

% of preliminary
results

Convergence point
?

Risk of Using Preliminary Results

DISTRIBUTION C: Distribution authorized to U.S. Government Agencies and their Contractors, 07/11/11. Other requests for this document must be referred to Air Force Research Laboratory/RIGD, 525 Brooks Rd., Rome NY 13441-4505.

Keesook.Han@rl.af.mil

INFORMATION INSTITUTE WORKSHOP on ASSURING THE CLOUD, 11 July 2011, GRIFISS Institute, Rome, NY          POC: AFRL/RIGD Dr. Keesook J. Han
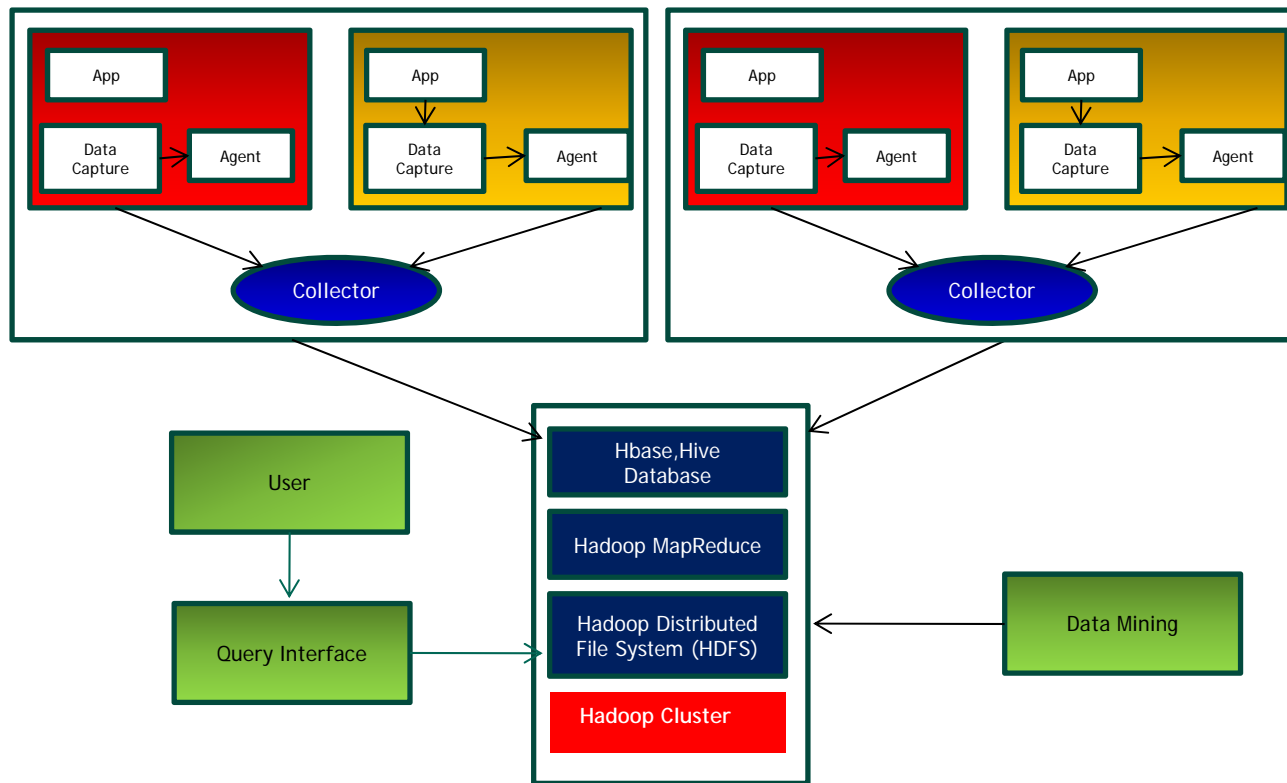
# System Architecture

# Prototype

DISTRIBUTION C: Distribution authorized to U.S. Government Agencies and their Contractors, 07/11/11. Other requests for this document must be referred to Air Force Research Laboratory/RIGD, 525 Brooks Rd., Rome NY 13441-4505.

Keesook.Han@rl.af.mil

INFORMATION INSTITUTE WORKSHOP on ASSURING THE CLOUD, 11 July 2011, GRIFISS Institute, Rome, NY    POC: AFRL/RIGD Dr. Keesook J. Han

# Work Accomplished

- Malicious traffic available from UTSA's cyber security monitoring system.

- Developed prototype of IP and DNS analysis of malicious and normal traffic

- Design of secure IP geolocation algorithm to defend against attacks on delay measurements.

- Implemented Hadoop based prototype for collecting and storage of network flow trace.

- Developed prototype of web user interface to retrieve data stored in Hadoop based system.

- Enhanced the portability of the Chukwa's data collection process by supporting Windows based cloud users.

DISTRIBUTION C: Distribution authorized to U.S. Government Agencies and their Contractors, 07/11/11. Other requests for this document must be referred to Air Force Research Laboratory/RIGD, 525 Brooks Rd., Rome NY 13441-4505.

Keesook.Han@rl.af.mil

INFORMATION INSTITUTE WORKSHOP on ASSURING THE CLOUD, 11 July 2011, GRIFISS Institute, Rome, NY    POC: AFRL/RIGD Dr. Keesook J. Han

# Future Works

- Develop secure IP geolocation algorithm to defend against attacks on measurement based IP geolocation

- Develop Hadoop based algorithms and implementation techniques for distributed data storage and online query processing

- Design and evaluate middleware to support structured and unstructured data and multiple databases

- Design efficient data collection system to reduce the amount of network disk access.

- Develop game theoretic model for refining convergence monitoring, and develop tools to handle concept drift in data.

DISTRIBUTION C: Distribution authorized to U.S. Government Agencies and their Contractors, 07/11/11. Other requests for this document must be referred to Air Force Research Laboratory/RIGD, 525 Brooks Rd., Rome NY 13441-4505.

Keesook.Han@rl.af.mil

INFORMATION INSTITUTE WORKSHOP on ASSURING THE CLOUD, 11 July 2011, GRIFISS Institute, Rome, NY    POC: AFRL/RIGD Dr. Keesook J. Han