# BiG Data Security

## CISC 6640 PRIVACY AND SECURITY IN BIG DATA

Instructor:
**Md Zakirul Alam Bhuiyan**
**Assistant Professor**
Department of Computer and Information Sciences
Fordham University

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# We Have Learned ...

- **Database Security**
- **Relational Databases**
  - Database security models
- **No SQL Databases**
- **Object Based vs. Object Oriented**
- **Overview of Database Vulnerabilities**
  - Common DBMS vulnerabilities
- **Overview of Database topics/issues (indexing, inference, aggregation, polyinstantiation)**
  - Security issues of inference and aggregation
- **Hashing and Encryption**
- **Database access controls (DAC, MAC, RBAC, Clark-Wilson)**
- **Information flow between databases/servers & applications**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# What We Are Going to Learn…

o **Data Administration**

- Big Data Security
- Data Policies, Data Quality
- Data Ownership
- Data Warehousing
- Long Term Archival
- Data Validation
- Data Security (access control, encryption)

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Implications of Big Data Key Aspects

- Security impact of big data dephasing 5Vs
  - Volume, Velocity, Variety, Veracity, Value

- The key aspects increase challenge for security
  - Science and foundations
  - Infrastructure
  - Management
  - Searching and mining
  - Applications.

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Big Data Security

○ **Data security through Hadoop**

- **Challenges centers around the fragmented data security issues, though new tools and technologies are surfacing.**
  - The Kerberos (a network authentication protocol) is a great step toward making Hadoop environments secure.
    - It is still immature

○ **Hadoop File Permissions**

- **Added in HADOOP-1298**
  - Hadoop 0.16, Early 2008
- **Authorization without authentication**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Big Data Security

o **MapReduce ACLs**
  - **Added in HADOOP-3698**
    - Hadoop 0.19, Late 2008
  - **ACLs per job queue**
  - **Set a list of allowed users or groups per operation**
    - Job submission
    - Job administration
  - **No authentication**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Big Data Security

o **Securing a Cluster Through a Gateway**

- **Hadoop cluster runs on a private network**

- **Gateway server dual-homed (Hadoop network and public network)**

- **Users SSH onto gateway**
  - Optionally can create an SSH proxy for jobs to be submitted from the client machine

- **Provides minimum level of protection**

SSH: Secure Shell, a cryptographic network protocol, used to log into a remote machine and execute commands

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Big Data Security

- Why security matters
  - **Prevent Accidental Access**
    - Doesn't require strong authentication
  - **Stop Malicious Users**
    - Security has to get real
    - Hadoop runs arbitrary code
    - Implicit trust doesn't prevent the insider threat
  - **Big data means getting rid of stovepipes**
    - Scalability and flexibility are only 50% of the problem
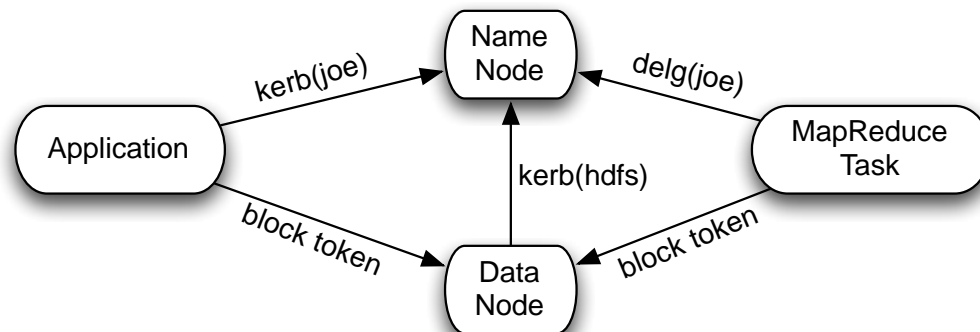    - Trust your data in a multi-tenant environment

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Big Data Security

o **An evolving story**

- **Files**
- **Service-level authorization**
  - Whitelists and blacklists of hosts and users

o **Authentication**

- HADOOP-4487, Hadoop 0.22 and 0.20.205, Late 2010
- Based on Kerberos and internal delegation tokens
  - Provides strong user authentication
  - Also used for service-to-service authentication

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Big Data Security

o **Encryption with Hadoop**

- **Over the wire encryption for some socket connections**
- **RPC (Remote Proc. Call) encryption added after Kerberos**
- **Shuffle encryption (HTTPS) added in Hadoop 2.0.2-alpha**
- **HDFS block streamer encryption added in Hadoop 2.0.2-alpha**
- **Volume-level encryption for data at rest**

o **Pluggable Authentication**

- **Currently supports username/password authentication backed by ZooKeeper**
- **Authentication info replaced with generic tokens**
- **Supports multiple security level implementations (e.g. Kerberos)**

**FORDHAM UNIVERSITY**
THE JESUIT UNIVERSITY OF NEW YORK

# Big Data Security

o **Application Level**

- **Apache Accumulo is a computer software project that developed a sorted, distributed key/value store based on the BigTable technology from Google**
  - Accumulo often paired with application level authentication/authorization
  - Accumulo users created secure access level per application
  - Each application granted access level of most permitted user
  - Application authenticates users, grabs user authorizations, passes user labels with requests

# Big Data Security

- **Apache Hbase**
  - **HBase is the Hadoop database, a distributed, scalable, big data store**
    - Also based on Google's BigTable
    - Started as a Hadoop contrib project
    - Supports column-level ACLs
    - Kerberos for authentication
    - Discussion and early prototypes of cell-level security ongoing
- **Future: Encryption for Data at Rest**
  - **Need multiple levels of granularity**
  - **Encryption keys tied to authorization labels (like Accumulo labels or HBase ACLs)**
  - **APIs for file-level, block-level, or record-level encryption**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# HBase

○ **HBase is an open-source, distributed, column-oriented database built on top of HDFS based on BigTable!**

- **Designed to operate on top of the Hadoop**
- **Distributed file system (HDFS) or Kosmos File System (KFS, aka Cloudstore) for scalability, fault tolerance, and high availability.**
  - No real indexes
  - Automatic partitioning
  - Scale linearly and automatically with new nodes
  - Commodity hardware
  - Fault tolerance
  - Batch processing

# HBase Security

o **Cell Tags**
  - **All values in HBase are now written in cells, can also carry arbitrary no. of tags such as metadata**

o **Cell ACLs**
  - **Enables the checking of (R)ead, (W)rite, E(X)excute, (A)dmin & (C)reate**

o **Cell Labels**
  - **Visibility expression support via new security coprocessor**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# HBase Security

o **Transparent Encryption**

- **Data is encrypted on disk – HFiles are encrypted when written and decrypted when read**

o **RBAC**

- **Uses Hadoop Group Mapping Service and ACL's to implement**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Encryption

o **Over the wire encryption for some socket connections**

- **RPC encryption added soon after Kerberos**
- **Shuffle encryption (HTTPS) added in Hadoop 2.0.2-alpha, back ported to CDH4 MR1**
- **HDFS block streamer encryption added in Hadoop 2.0.2-alpha**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Encryption

- **Over the wire encryption for some socket connections**
  - **Need multiple levels of granularity**
  - **Encryption keys tied to authorization labels (like HBase ACLs)**
  - **APIs for file-level, block-level, or record-level encryption**
  - **Volume-level encryption for data at rest**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Cell Level Security

o **Cell level ACL means explicit 'RW' access can be set on individual cells when the cell data is put into HBase.**

o **<span style="color:red">Visibility labels</span> allow administrators to associate secure access to cells with visibility labels**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Cell Level Security

o **Labels stored per cell**

o **Labels consist of Boolean expressions (AND, OR, nesting)**

o **Labels associated with each user**

o **Cell labels checked against user's labels with a built-in iterator**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Cell Level Security



Table-Level Security

Column-Level Security

Row-Level Security

Cell-Level Security

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Cell Level Security

o **Data Model**
- **Multi-dimensional, persistent, sorted map**
- **Key/Value store with a twist**
- **A single primary key (Row ID)**
- **Secondary key (Column) internal to a row**
  - Family
  - Qualifier
- **Per-cell timestamp**

| Key | | | | | Value |
|-----|---|---|---|---|-------|
| Row ID | Column | | | Timestamp | |
| | Family | Qualifier | Visibility | | |

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# ○ **Data Administration**

- Big Data: Hadoop / Mongo DB / HBASE
- Data Policies, Data Quality
- Data Ownership
- Data Warehousing
- Long Term Archival
- Data Validation
- Data Security (access control, encryption)

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

- **Data** is any type of stored digital information that is asset.

- **Security** is about the protection of assets.

- Prevention
  - **Measures taken to protect your assets from being damaged.**

- Detection
  - **Measures taken to allow you to detect when an asset has been damaged, how it was damaged and who damaged it.**

- Reaction
  - **Measures that allow you to recover your assets.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

o **Data confidentiality, integrity, and availability.**

o **Accountability is audit information that is kept and protected so that security actions can be traced to the responsible party.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

o **Audit Standards**

- **Data Security is subject to several types of audit standards and verification.**
  - The most common are ISO 17799, ISO 27001-02, PCI, ITIL, SAS-70, HIPPA, SOX

- **Security Administrators are responsible for creating and enforcing a policy that forms to the standards that apply to their organizations business.**

- **IT certification audits are generally carried out by 3rd party accounting firms.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

○ **Security Policies**

- **A security policy is a comprehensive document**
  - That defines a companies' methods for prevention, detection, reaction, classification, accountability of data security practices and enforcement methods.
    - It generally follows industry best practices as defined by ISO 17799,27001-02, PCI, ITIL, SAS-70, HIPPA , SOX or a mix of them.

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

o **Security Policies**

- **A security policy is a comprehensive document**
  - The security policy is the key document in effective security practices.
    - Once it has been defined it must be implemented and modified and include any exceptions that may need to be in place for business continuity.
  - All users need to be trained on these best practices with continuing education at regular intervals.

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

o **Tools to Secure Data**

- **Data needs to be classified in the security policy according to its sensitivity.**

- **Once this has taken place, the most sensitive data has extra measures in place to safeguard and ensure its integrity and availability.**

- **All access to this sensitive data must be logged.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

o **Tools to Secure Data**

- **Secure data is usually isolated from other stored data.**
  - Controlling physical access to the data center or area where the data is stored.
  - Active or Open Directory is a centralized authentication management system that is available to companies to control and log access to any data on the system.
- **Encryption of the sensitive data is critical before transmission across public networks.**

# Data Policies

o **Tools to Monitor Secure Data**

- **Walk around and look for passwords in the open.**
- **Event Viewer / Log Files**
- **Intrusion Detection/ Protection systems (IDS/IPS) such as SNORT.**
- **These will alert Administrators of suspicious data flows.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

o **Tools to Document Data Security**

- **Microsoft Visio is the standard for drawing network maps.**
  - These maps allow a detailed overview of the system and how it is functions.
  - They also allow the spotting of weak points of security and flaws in design that can impact reliability or continuity of the data to the end user.

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

o **End User Education**

- **All relevant security polices must be clearly explained to the end users.**

- **A clear explanation of the consequences for violating these polices must also be explained.**

- **The end user needs to sign a document acknowledging that they understand the policies and consequences for violating these policies.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Policies

o **Enforcement**

- **Must obtain executive authority to enforce security policy.**
- **Systematic approach of warnings and punishments.**
- **Coordinate with HR to document continued issues with staff.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

o **Data Administration**

- Big Data: Hadoop / Mongo DB / HBASE
- Data Policies, Data Quality
- Data Ownership
- Data Warehousing
- Long Term Archival
- Data Validation
- Data Security (access control, encryption)

# Data Quality

o **Data in the real world is of low quality (dirty!)**

- Poor data across businesses and the government costs the U.S. economy $3.1 trillion dollars a year.

- Poor data can cost businesses 20%–35% of their operating revenue.

- Bad data or poor data quality costs US businesses $600 billion annually.

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Quality

o **Data in the real world is of low quality (dirty!)**

o **When data can be of high quality?**

- **The data is accurate**
- **The data is stored according to data type**
- **The data has integrity**
- **The data is consistent**
- **The databases are well designed**
- **The data is not redundant**
- **The data follow business rules**
- **The data corresponds to established domains**
- **The data is timely**
- **The data is well understood**
- **etc.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Quality

o **Data in the real world is of low quality (dirty!)**

- **Inconsistent**
  - Containing discrepancies in codes or names
  - e.g., Age="42" Birthday="03/07/1997"
  - e.g., Was rating "1,2,3", now rating "A, B, C"
  - e.g., discrepancy between duplicate records
- **Incomplete**
  - Lacking attribute values, lacking certain attributes of interest, or containing only aggregate data
  - e.g., occupation=""
- **Noisy**
  - Containing errors or outliers
  - e.g., Salary="-10"

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Quality

o **Data in the real world is dirty**

- **Illegal values**
- **Violated attribute dependencies**
- **Uniqueness violation**
- **Referential integrity violation**
- **Missing values**
- **Misspellings**
- **Cryptic values**
- **Embedded values**
- **Misfielded values**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Quality

o **Data in the real world is dirty**

- **Word transpositions**
- **Duplicate records**
- **Contradicting records**
- **Wrong references**
- **Overlapping data/matching records**
- **Name conflicts**
- **Structural conflicts**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Quality

o **Assessment of existing data quality**

- **Programs that abnormally terminate with data exceptions**
- **Clients who experience errors/anomalies**
- **Clients who do not know or are confused about what the data actually means**
- **Data that cannot be shared due to lack of integration**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# ○ **Data Administration**

- Big Data: Hadoop / Mongo DB / HBASE
- Data Policies, Data Quality
- Data Ownership
- Data Warehousing
- Long Term Archival
- Data Validation
- Data Security (access control, encryption)

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Ownership

o **Data Owner Definition:**

- **A person with statutory or operational authority for specified information (e.g., supporting a specific business function) and responsibility for establishing the controls for its generation, collection, processing, access, dissemination, and disposal**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Ownership

- **Grants** access to the Information System under his/her responsibility.

- **Classifies** Digital Data based on Data sensitivity and risk.

- **Backs up** Data under his/her responsibility in accordance with risk management decisions and secures back up media.

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

o **Data Administration**

- Big Data: Hadoop / Mongo DB / HBASE
- Data Policies, Data Quality
- Data Ownership
- **Data Warehousing**
- Long Term Archival
- Data Validation
- Data Security (access control, encryption)

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Warehouse

o **Definition**

- An enterprise structured repository of **subject-oriented**, **time-variant**, **historical data** used for information retrieval and decision support. The data warehouse stores atomic and summary data. [Oracle Data Warehouse Method]

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Why Do We Need Data Warehouses?

o **Consolidation of information resources**

- **Improved query performance**

- **Separate research and decision support functions from the operational systems**

- **Foundation for data mining, data visualization, advanced reporting and OLAP tools**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Why Do We Need Data Warehouses?

# What Is a Data Warehouse Used for?

o **Knowledge discovery**

- **Making consolidated reports**
- **Finding relationships and correlations**
- **Examples**
  - Banks identifying credit risks
  - Insurance companies searching for fraud
  - Medical research

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Warehouse Properties

○



Subject Oriented

Equity Plans, Shares, Insurance, Savings

Integrated

Relational databases, flat
on-line transaction record
E.g., Hotel price: current
tax, breakfast covered, et

Data Warehouse

Load, Read
No: Modification, update,

Non Volatile

Time Variant

| Time | Data |
|------|------|
| Jan-97 | January |
| Feb-97 | February |
| Mar-97 | March |

# How Do Data Warehouses Differ From Operational Systems?

- Operational
  - **Online transaction and query processing**
  - **Day-to-day operations**
  - **Users and system orientation: customer oriented**
    - Online transaction processing (OLTP)
- Data warehouse
  - **Serve users or knowledge workers in the role of data analysis and decision making**
  - **Decision support: operations as defined or needed**
  - **Market-oriented**
    - Online analytical processing (OLAP)

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# How Do Data Warehouses Differ From Operational Systems?

| Property | Operational DB | Data Warehouse |
|---|---|---|
| Response time | Sub seconds to seconds | Seconds to hours |
| Operations | Historic data | Primarily read only |
| Nature of Data | 30-60 days | Snapshots over time |
| Data Organization | Applications | Subject, time |
| Size | Small to large | Large to very large |
| Data Source | Operational, Internal | Operational, Internal, External |
| Activities | Processes | Analysis |

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# How Do Data Warehouses Differ From Operational Systems?

| OLTP Examples | OLAP Examples |
|---|---|
| Which product is mostly purchased (or sold)? | Amazon analyzes purchases by its customers to come up with an individual screen with products of likely interest to the customer. |
| Answering queries from a Web interface, sales at cash registers, e.g., selling particular medicine. | Analysts at Wal-Mart look for items with increasing sales in some region. |

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Decision Support

o **Used to manage and control business**

o **Data is historical or point-in-time**

o **Optimized for inquiry rather than update**

o **Use of the system is loosely defined and can be ad-hoc**

o **Used by managers and end-users to understand the business and make judgements**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Warehouse Models

o **Enterprise warehouse**

- Collects all of the information about subjects spanning the entire organization

o **Data mart**

- A subset of corporate-wide data that is of value to a specific groups of users.
- Its scope is confined to specific, selected groups, such as marketing data mart

o **Virtual warehouse**

- A set of views over operational databases
- Only some of the possible summary views may be materialized

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Warehouse Products

- Computer Associates -- CA-Ingres
- Hewlett-Packard -- Allbase/SQL
- Informix -- Informix, Informix XPS
- Microsoft -- SQL Server
- Oracle – Oracle
- Red Brick -- Red Brick Warehouse
- SAS Institute -- SAS
- Software AG     -- ADABAS
- Sybase     -- SQL Server, IQ, MPP

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Security in the Data Warehouse

o **Unfortunately, most data warehouses are built with little or no consideration given to security during the development phase.**

- **Basic security concepts**
- **Physical security**
- **Stand-alone or shared security**
- **Remote access**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Security in the Data Warehouse

○ **Achieving proactive security requirements of DW is a seven-phase process:**

1) **identifying data**

2) **classifying data**

3) **quantifying the value of data**

4) **identifying data security vulnerabilities**

5) **identifying data protection measures and their costs**

6) **selecting cost-effective security measures, and**

7) **evaluating the effectiveness of security measures. These phases are part of an enterprise-wide vulnerability assessment and management program.**

# Security in the Data Warehouse

- o **Controlling Access to Warehouse Data**
- o **Role-based Access Control**
- o **Row-Level Security**
- o **Virtual Private Database (VPD)**
  - **Column-Level Virtual Private Database**
  - **Policy Types for Added Performance**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# o **Data Administration**

- Big Data: Hadoop / Mongo DB / HBASE
- Data Policies, Data Quality
- Data Ownership
- Data Warehousing
- Long Term Archival
- Data Validation
- Data Security (access control, encryption)

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Long Term Archival

- **Long Term Data that can be entered and/or checked over a longer time frame, using collaborative arrangements**
- **Archiving is for long-term storage of data that is not required immediately.**
- **More often than not it is never required again but it kept just in case.**
- **Data is often removed from a system and stored separately.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Long Term Archival

o **Consider the following example:**

- **A school records data about pupils' performance every year. If they continued to collect data, even after pupils had left school, the system resources would soon diminish. Instead, records about pupils are removed from the system once they leave.**

- **However, some data may be archived such as average test scores and achievement rates. The data is not needed immediately but may be useful to keep for the future.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# o Data Administration

- Big Data: Hadoop / Mongo DB / HBASE
- Data Policies, Data Quality
- Data Ownership
- Data Warehousing
- Long Term Archival
- Data Validation
- Data Security (access control, encryption)

# Data Validation

o **It allows us to provide data with confidence in its accuracy, and we can consistently provide this data by implementing thorough security.**

o **Data Integrity**

- **Validity, consistency, and accuracy of the data in a database.**
  - Table-level
  - Field-level
  - Relationship-level
  - Business Rules

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Validation

o **The process of determining if an update to a value in a table's data cell is within a preestablished range or is a member of a set of allowable values.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Data Validation

○ **Methods for checking for data validity**

- **Visual/manual**
- **Aggregation**
- **Reviewers guide**
- **Auto data checks**
- **Record counts**
- **Spell checks**
- **Have data provider review**

FORDHAM UNIVERSITY
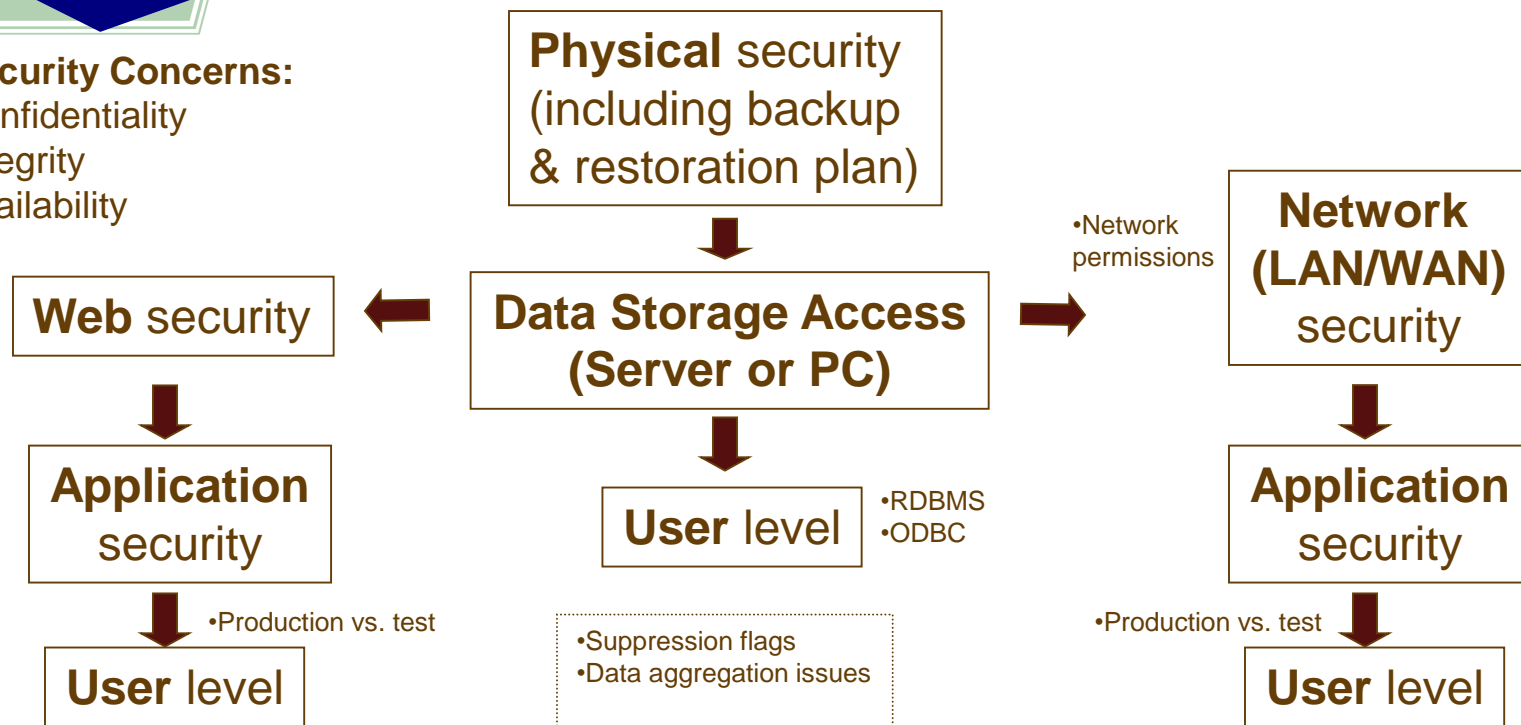THE JESUIT UNIVERSITY OF NEW YORK

# Data Validation

## ○ **Secure Systems for Data Validation**

### Data Security Considerations

The contact in my state is:_____

**Security Concerns:**
Confidentiality
Integrity
Availability

**Physical** security (including backup & restoration plan)

**Web** security

**Data Storage Access (Server or PC)**

**Network (LAN/WAN)** security

•Network permissions

**Application** security

**User** level
•RDBMS
•ODBC

**Application** security

•Production vs. test

•Production vs. test

**User** level

•Suppression flags
•Data aggregation issues

**User** level

**Database security prevents unauthorized person(s) from viewing, destroying or altering data within the database.**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK

# Coming Attraction…

o **Reading: Slides for quiz**

o **Assignment 1**

o **Next class**

- **Proposal Presentaiton**

FORDHAM UNIVERSITY
THE JESUIT UNIVERSITY OF NEW YORK