ABD310

# AWS re:INVENT

## Big Data, AWS, and Security
## How FINRA Secures Its Big Data and Data Science Platform on AWS

Vincent Saulys
David Yacono

November 29, 2017

re:Invent

aws

---

# Who is **FINRA**?

- **F**inancial **I**ndustry **R**egulatory **A**uthority.
- Our Mission: "Investor Protection—Market Integrity."
- We are a private sector not-for-profit organization authorized by Congress to protect America's investors.
- We do this by:
  - Writing and enforcing rules that govern the activities of 3,800 broker-dealers with 634,000 brokers.
  - Examining firms for compliance with those rules.
  - Fostering market transparency.
  - Educating investors.

And most significant to this discussion:
- FINRA uses big data and data science technologies to detect and analyze fraud, market manipulation, and insider trading across US capital markets.

re:Invent

aws

## FINRA Technology

# Innovating
**to protect investors and ensure market integrity**

FINra

Over
**25 PETABYTES**
of storage

Market
reconstruction
containing
**TRILLIONS**
of nodes and edges

Up to
**75 BILLION**
events
per day

BUY

SELL

3 - 4 years
**QUERYABLE**
data online

4 years
**ARCHIVAL**
storage

---



# Session takeaways

FINra

- Learn to be **realistic** in your risk assessment of using the cloud.

- Learn about Amazon Web Services and its foundational security controls and practices relevant to **safeguarding** your big data workloads.

- See how FINRA **securely enables** our data scientists, and other big data projects, in AWS by achieving a balance between productivity and security.
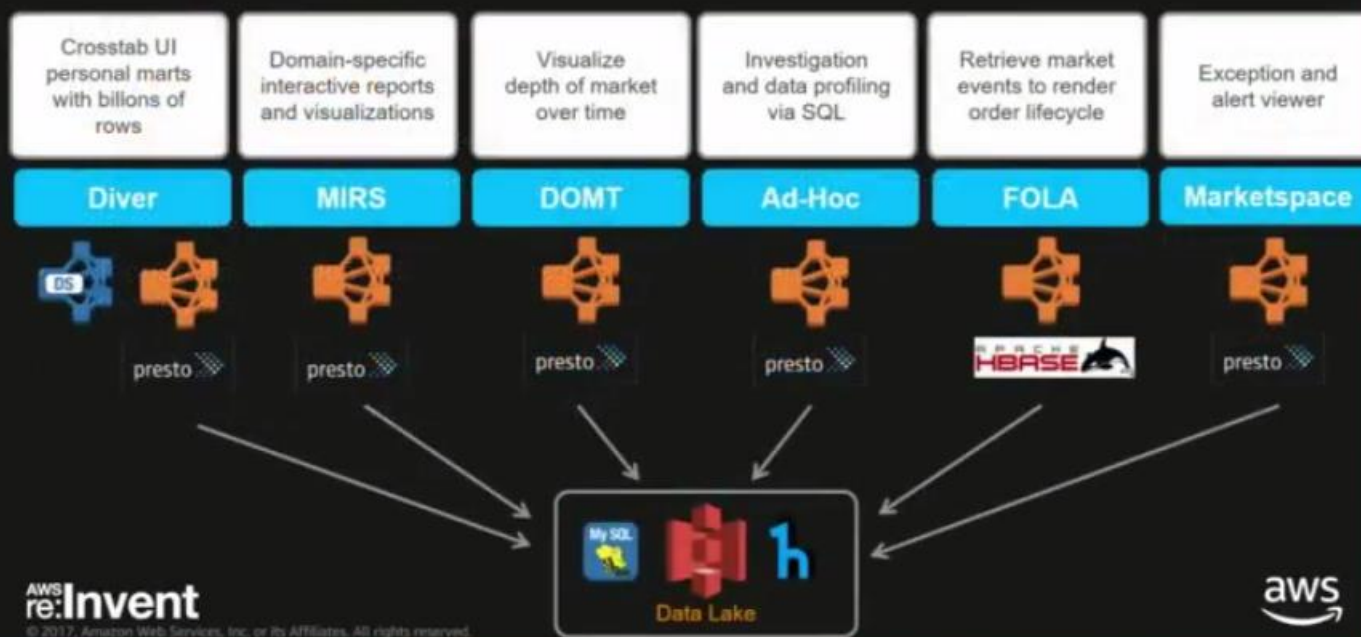
# Data science needs

- Data discovery and exploration.
- Bring disparate sources of data together.
- Semantic understanding of the data sets.
- Ease of use: Enable users without having to understand underlying data infrastructure.
- Safeguard information with a high degree of security and least privileges access.
- Model migration from research to prototype to production.
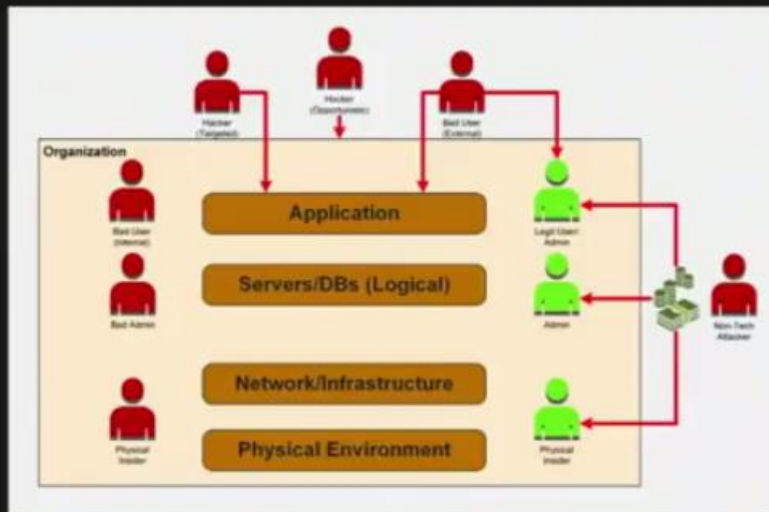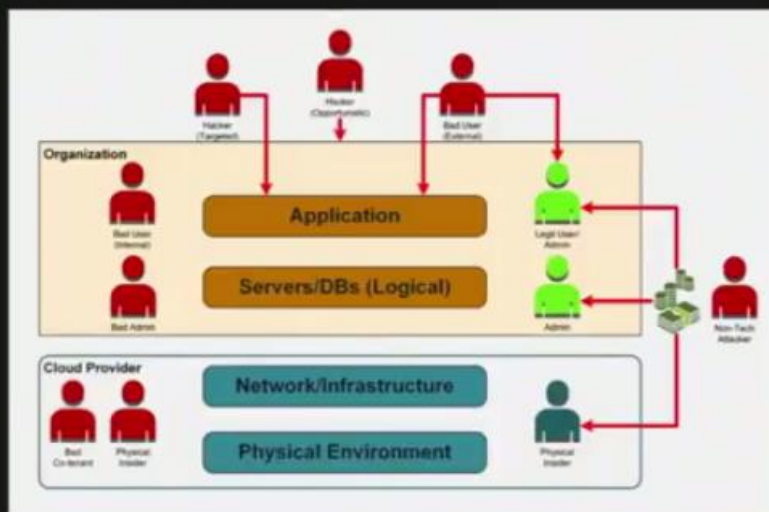- Avoid time spent on environment administration.

# Interactive big data portfolio

| Crosstab UI personal marts with billons of rows | Domain-specific interactive reports and visualizations | Visualize depth of market over time | Investigation and data profiling via SQL | Retrieve market events to render order lifecycle | Exception and alert viewer |
|---|---|---|---|---|---|
| Diver | MIRS | DOMT | Ad-Hoc | FOLA | Marketspace |
| presto | presto | presto | presto | APACHE HBASE | presto |

Data Lake

These are all the query engines that we built and the applications that we have built on top the engines. We built our own data lake (~25 Petabytes) using S3, we also have our OSS data catalog called Herd that helps store the metadata about our data including all previous data history. We generally use Presto (and HBase) to query the data lake from the many apps we build.

There are security needs for big data.

# Evaluating the risk

## Key factors

**Compare risk of alternatives**
- Cloud vs. legacy data center
  - The "idealized zero-risk scenario" is unrealistic.
  - Legacy data centers have risk!

**Evaluate risk in context**
- "Cloud" risks get the most press.
- Many existing risks are unchanged by adoption of cloud.
- Many existing risks are significantly "riskier" than new "cloud" risks.
- The "unknown" creates a powerful perception of risk.

**Cybersecurity is only one dimension of risk**
- Operational, legal, and opportunity risks as well

aws

---

# Shared responsibility model plus

Security "ABOVE" the cloud
- All the security controls you're already using.

Security "IN" the cloud
- Controls to supplement Cloud Service Provider (CSP) controls.

Security "OF" the cloud
- CSP provides these controls. Customer due diligence through third-party risk management.

| Standard Cyber Security Controls (IDS, Secure SDLC, IAM) | | | |
|---|---|---|---|
| Customer Data | | | |
| Platform, Applications, Identity & Access Management | | | |
| Operating System, Network & Firewall Configuration | | | |
| Client-side Data Encryption & Data Integrity Authentication | Server-side Encryption (File System and/or Data) | Network Traffic Protection (Encryption/Integrity/Identity) | |
| Compute | Storage | Database | Networking |
| AWS Global Infrastructure | Regions | Availability Zones | Edge Locations |

aws

# Foundational security controls

- AWS security best practices
- Access management
  - Authentication
  - Authorization/Separation of Duties (SoD)
  - Policy enforcement and oversight
- Logging/monitoring/alerting/UEBA
  - Controls for Economic DoS
- Network architecture
- Encryption and Key Management (KMS)
- Governance

---

# Access mgmt: Human

We have access granted by roles per environment. No IAM role has Delete role on Production assets like S3 buckets or EC2 instances. We follow the principle of least entitlement where an individual is given the minimum set of role entitlements that they need to get the job done.

We use an OSS application credential vault called Credstash that implements KMS encryption of the secrets and stores them in DynamoDB. The resource (DB admin) creates the secrets and passwords, they store it in Credstash, automation on the user's system (Python or Java client) pulls those secrets at deploy time without human intervention, the target system is IAM role restricted for seeing the credentials.

This is a sample deployment architecture we use at FINRA, we use up to 4 AZs.

We use SSL in-transit encryption with all the AWS service endpoints we have, we also use encryption at rest and S3 and KMS encryption.

# Securing the services—Amazon EC2

### AMI updates

- Created at least monthly
  - Plus out-of-cycle for critical security patches
- Start: Latest Amazon AMI
- Harden: Remove unneeded packages, update remaining packages (security patches) [Yum], apply compliance modules [Puppet]
- Extend: Install common tools (AWS CLI, Puppet agent, Splunk agent, Trend Micro agent, etc.);
- Snapshot new AMI.

### Security Groups

- Goals:
  - Narrowly crafted (microsegmentation),
  - Policy of least privilege,
  - Separation of Duties
- Challenges: Many groups to manage!
- Solution: FINRA Portus

### Strictly Limited Access

- Goal: No access to production.
- Reality: Occasional prod access may be needed
- Modified Goals: Temporary, just-in-time access
  - Restricted by IAM Role, AWS Tag
  - Approved and Logged
- Solution: FINRA Gatekeeper

aws

---

# Security group mgmt: Portus

**FINRA-developed Centralized Security Groups Management tool for Developers and the FINRA Cyber & Information Security Department**

- Cyber & Information Security DEFINE security policies
- Developers SELF MANAGE AGS security groups

- Maximizes FLEXIBILITY for developers
- SIMPLIFIES administration for InfoSec

Portus

aws

Portus dashboard



Portus dashboard

Security Policy for each type of logical system

**Portus** dashboard

Policy is a set of Whitelisted Rules (inbound and outbound)



**Portus** dashboard

Only these Whitelisted Rules are allowed in Security Groups

# Portus **developer** dashboard



Access only to the AGS owned by developer (priv_aws_<AGS>_dev_d AD group)

# Portus **developer** dashboard



Select AGS from drop down list

Gatekeeper is another system that we built that allows that Just in time, limited production environment access to an EC2 instance.

This further minimizes the opportunity for credentials to be exposed.



The Herd system helps us to manage access to the information to the data lake in S3, it uses metadata about the data like namespace, object definition, format, partition, encryption type, etc.

# Securing big data services—Amazon EMR

- AWS::EMR::SecurityConfiguration

- At-rest encryption
  - Local volumes (LUKS), HDFS encryption

- In-transit encryption
  - Inter-node: Spark/Presto/Hive (TLS). HBase in 2018?
  - EMR-S3: EMRFS/TLS

- Logging—Splunk agent in bootstrap:
  - JobCode ID, project code logged

- Access Controls
  - No access to underlying cluster. (App layer AuthN/AuthZ)
  - Gatekeeper for admin access (rare)

aws

---

# Securing the architecture

**Data sanitization**
- One-way hash/tokenization
  Preserves ability to associate related records by the sensitive data element, search on tokenized values
- Format-preserving encryption
  Preserves ability to associate related records, some limited ability to operate on data (search, sort, categorize)
- Generalization, subsetting
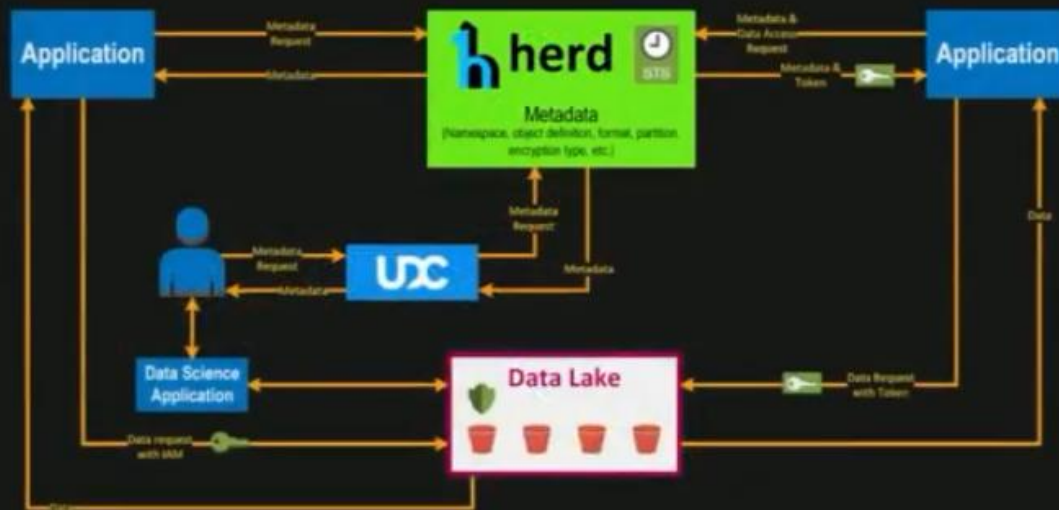- De-identification
  - Be wary of re-identification strategies

**Limit Credential Exposure**
- IAM Role-based access is ideal
- Secrets Management (Credstash)

**Make Security Easy**
- Internal mirrors of external resources preserves isolation
- Empower users, managers with utilization/cost information, necessary entitlements to provide oversight.

aws

We moved all our data into S3 and provide access to it through the various SQL engines like Redshift, Presto, Hive, and added access to on-premise data in databases and files via the VPC.

We now move the compute into the cloud by providing EC2 instances with the SQL engines for querying and working with the S3 data using some Python libraries.

# What went **wrong**?

**Needs driven by technology**
- IT: Reduce costs
- Users: Need more compute

**Secure but inflexible**
- Local machines were more flexible
- Install any package and experiment

**Data availability**
- On-premises databases not reachable

**Setup still required**
- Driver configuration to connect to databases

**Technology in the way**
- Technology required to install any new package

hello?

---

# **U**niversal **D**ata **S**cience **P**latform

FINRA

Amazon Redshift
Presto
Hive
Spark
Data Catalog

S3

CRAN
PyPI

Internet

VPC

Files

FINRA Databases

# What went **right**?

**Completely self service, no Technology administration**
- Users select UDSP version and machine capacity

**Users associated to groups**
- Users manage their instances
- AWS billing tags and machine selection choices to group

**Create, stop, terminate (delete)**
- Managers can administer their teams' instances

**Dashboard to monitor resource usage**
- Stop instances from the dashboard

**Reports for historical usage**

**USAGE INCREASED**

**20 FOLD!!!!**

aws

---

# Recap

- Be **realistic** in your risk assessment. The security risks in using your own data center are equal to or more than going to the cloud.

- Use of strong **foundational cloud security controls** is paramount.

- AWS provides controls which, when properly applied, **balance** productivity and security.

aws

# Related FINRA presentations

## 2017 re:INVENT

- **SID326 - AWS Security State of the Union**
  Steve Schmidt, chief information security officer of AWS, addresses the current state of security in the cloud. As part of this presentation, John Brady (CISO of FINRA) shares the FINRA journey to the cloud. Wednesday, Nov 29, 12:15 p.m. – 1:15 p.m. MGM, Level 3, Premier Ballroom 316

- **FSV307 - Capital Markets Discovery: How FINRA Runs Trade Analytics and Surveillance on AWS**
  The FINRA analytics platform unlocks the value in capital markets data by accelerating trade analytics and providing a foundation for machine learning at scale. Monday, Nov 27, 10:45 a.m. – 11:45 a.m. Venetian, Level 5, Palazzo P

- **ENT328 - FINRA's Managed Data Lake: Next-Gen Analytics in the Cloud**
  The Financial Impact Regulatory Authority (FINRA) Technology Group has changed its customers' relationships with data by creating a managed data lake Thursday, Nov 30, 1 p.m. – 2 p.m. MGM, Level 3, Premier Ballroom 319

- **DEV335 - Manage Infrastructure Securely at Scale and Eliminate Operational Risks**
  Managing AWS and hybrid environments securely and safely while having actionable insights is an operational priority and business driver for all customers. Thursday, Nov 30, 4 p.m. – 4 p.m. Venetian, Level 2, Venetian E

## 2016 re:INVENT

- **BDM203: Building a Secure Data Science Platform on AWS**