# Big Data Platform

Lessons Learned in Growing a Big Data Capability for Network Defense

# Who am I?

- Technical Director, Enlighten IT Consulting, a MacAulay-Brown company
- Software Engineering Consultant
- Helped found Apache Rya
- Chief Architect of DoD's Big Data Platform
- Currently working for:
    - Defense Information Systems Agency (DISA)
    - Army Cyber Command
    - US Cyber Command
    - Center for Army Analysis
    - Air Force

# Talk Overview

- DCO Big Data Problem Space

- DoD's Big Data Platform

- Scaling for Big Data

- Multi-Tenancy

- Lessons Learned

# Problem Space

- Huge variety of DCO sensors

- Heterogeneous data formats

- No enterprise standardization on infrastructure

- Petabyte scale storage/retention/analysis requirements

- No single "out of the box" COTS, GOTS, or OSS solution by itself meets the unique DoD cyber security challenges

- Enabling collaborative investigation while eliminating redundant efforts

# Problem Space



CYBERscape · 3Q16 — Momentum

# What is the BDP?

- A cloud-based distributed architecture for ingesting and storing large datasets, building analytics, and visualizing the results.

- Allows critical decisions to be made based on rich and broad data.

- Developed around open source and unclassified  components while leveraging community tech transfer from other DoD entities.

- DISA-controlled software baseline

- RMF accredited with current Authority To Operate in multiple organizations

- 99% open source, specifically integrated to meet DoD's needs

# Data Sources

**1**

Network Security Devices

JRSS  ArcSight
HBSS  Firewalls
Netflow  IDS
Proxies  IPS

Network Management Devices

Hosts Logs
Router Logs
Switch Logs

Application Services

SQL
File Server
Active Directory
ACAS
DNS Logs

Intel/ Operational

NTOC
Incident Handling
OpenPhish

# Connect

**2**

Messaging

Processing

Streaming

# Transform (Enrich)

**3**

Data Fusion

Geospatial

Trends

Persona

# Store

**4**

Structured
Unstructured
Semi-structured
Data

Data Caching

HDFS and Accumulo

# Analyze

**5**

When will it happen?

What happened?

Why did it happen?

How can we prevent it from happening?

Identify Potential Exploit

Develop Hypothesis to Test

Select Relevant Data Sources

Determine Analysis to Perform

Interpret Results

# Visualize/Act

**6**

Data Exploration

Unity

Kibana

20,551

# Big Data Platform Technology Stack

## INGEST AND MESSAGING

- Ingest Pipeline
- Parse
- Canonincalize
- Enrich

- NiFi
- Flume
- Kafka
- Storm

## BUSINESS INTELLIGENCE AND ANALYTICS

- OWF
- R Shiny Server
- Navigator
- Kibana
- Watchtower
- AFD

## DATA PROCESSING

- R
- YARN
- MapReduce2
- Spark

## DATA STORES

- PostgreSQL
- ElasticSearch
- HDFS
- Accumulo
- Kronos
- Content
- AFD
- GEM (RYA)
- Metrics

## BDP CORE

- OS — RHEL
- Core Services — SIMP — Puppet — Java
- Citadel — OpenLDAP — Active Directory — Akamai

## INFRASTRUCTURE

- Bare Metal Commodity Servers
- AWS (GovCloud)
- Azure/HyperV
- Vsphere/ESXi
- OpenStack

## OPERATIONS

### MANAGEMENT

- Puppet
- RDA Deployer
- Slider

### COORDINATION

- Zookeeper
- Consul

### MONITORING

- Grafana
- Overseer

### ALERTING

- Nagios

### WEB SERVICES

- nGINX

# BIG DATA PLATFORM 3.0 ARCHITECTURE

Analytics   Capability Deployments

## DATA PROCESSING
Streaming Analytics

Sensor Events, Reports, Data Feeds

**KAFKA**

**STORM**

| Parse | AsyncServices |
| Canonicalize | |
| Enrich | |

## MANAGEMENT

RDAs

| AFD Service | Metric Service |
| Alerting | Webapps |
| Analytics | |

Consul

Airflow

Overseer

## SCALABLE WEB TIER
User Driven Analytics

**JETTY**

| Services | Widgets |
| Common API | |
| Resolvers | Navigator |
| GEM | |

NGINX

**BUSINESS INTELLIGENCE**

| Cache | Widgets |
| Elastic Search | |
| R/Shiny | Node.js |
| | Kibana |
| | Shiny Apps |

Data Science Analytics

## DATA STORAGE & ANALYTICS

**ACCUMULO**

Content | Metrics | Kronos | Analyics Datasets | GEM (RYA)
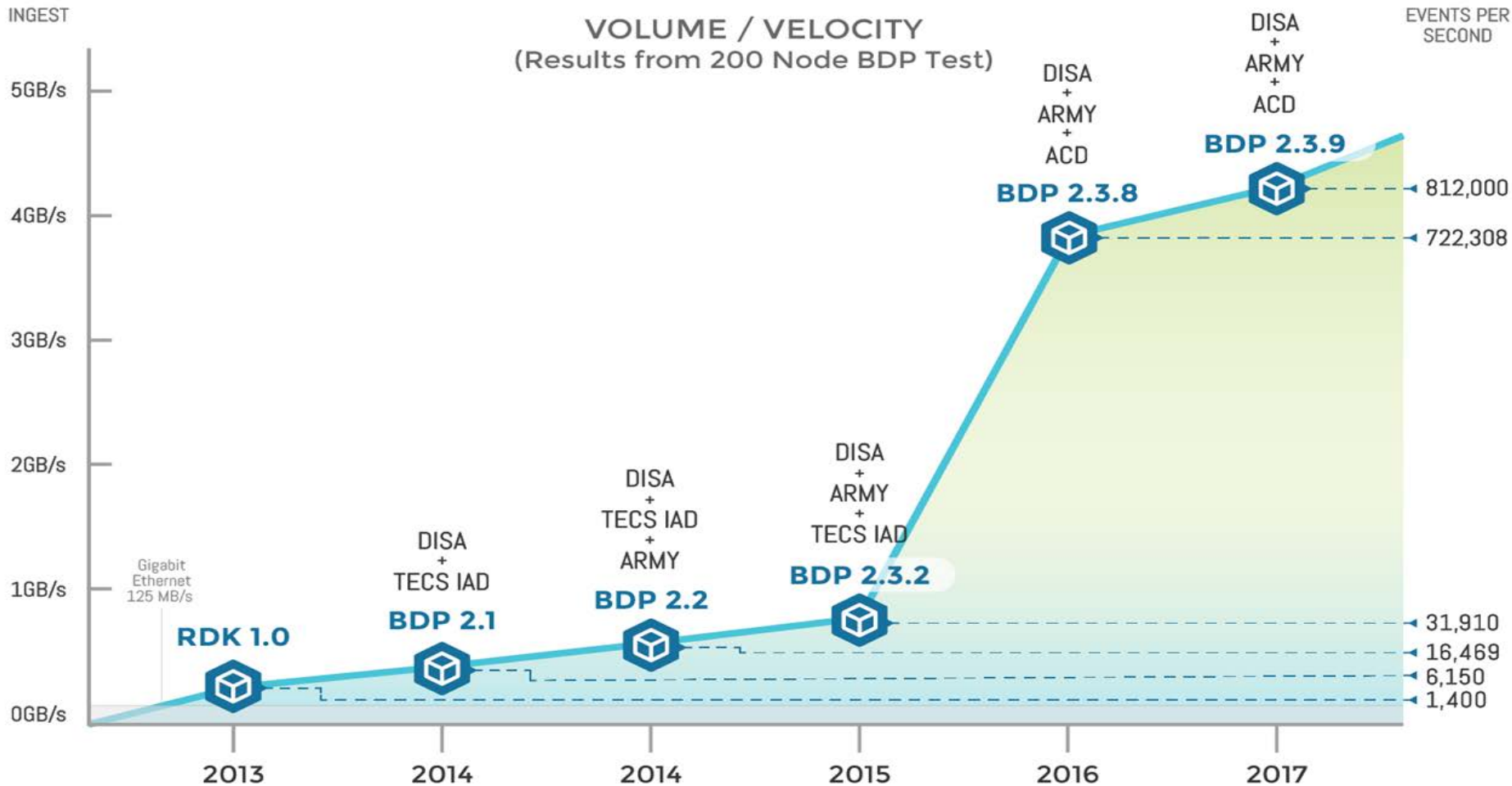
**HADOOP W/ YARN**

MR/2 Analytics | R on Data Nodes | Spark Analytics

Enrichment Stores

**CITADEL**

# Scaling for Volume and Velocity



**VOLUME / VELOCITY**
(Results from 200 Node BDP Test)

INGEST

- 5GB/s
- 4GB/s
- 3GB/s
- 2GB/s
- 1GB/s
- 0GB/s

EVENTS PER SECOND

Gigabit Ethernet 125 MB/s

**RDK 1.0**

**BDP 2.1**
DISA + TECS IAD

**BDP 2.2**
DISA + TECS IAD + ARMY

**BDP 2.3.2**
DISA + ARMY + TECS IAD

**BDP 2.3.8**
DISA + ARMY + ACD — 722,308

**BDP 2.3.9**
DISA + ARMY + ACD — 812,000

- 31,910
- 16,469
- 6,150
- 1,400

2013  2014  2014  2015  2016  2017

# Multi Tenancy (Learning to share)

- HDFS / Accumulo (Storage)
- Analytics
    - Spark
    - Streaming- Kafka/Storm
    - RShiny
- Web Applications
    - Jetty
    - NodeJS
- Microservices
    - Spring/Java/NodeJS
- Ingest

# Lesson Learned:
# It's all about the data

- Don't underestimate the difficulty of collecting and sharing data

- End user analytic questions have to drive data priorities

- You can't wait to start collecting data until you need to use it

- *Just enough* normalization will allow unplanned correlations to emerge

- Data from many vantage points increases the value (but analysts need to understand the vantage point of each)

# Lesson Learned:
# Use commercial cloud infrastructure

- It lets your engineering teams focus on your problems not on infrastructure

- It provides "just in time" capacity that reduces costs in the long run

- It has a refresh rate that is much more frequent than traditional in-house data centers

- It reduces barriers for data transport and acquisition

# Lesson Learned:
# Standardize your platform early, but evolve it

- Organizations can share security accreditation

- Shared data structures will encourage correlations

- Be willing to change and evolve, without reinventing everything every time

- Create and document APIs that encourage reuse

- Leverage a community to share costs

# Lesson Learned:
# Analytics need to scale

- Need to run on commodity hardware (if you can fit all your data into memory, you don't have big data)

- Need to be parallelizable

- Need to handle preemption (half your job may be killed at any moment to make way for higher priority tasks)

- Need to be secure (can't open ports, store passwords; need to handle data security controls)

# Lesson Learned:
# You need to optimize your load

- Use batch ingest

- Cache data near the web tier

- Adjust the allocation of resources to your mission (YARN is great, but it needs to be managed)

- Test with real world datasets (size and variety)

- Understand the computational costs of your analytics before deploying them

# Questions?