

# HEALTHCARE DOES HADOOP: AN ACADEMIC MEDICAL CENTER'S FIVE-YEAR JOURNEY

---

Charles Boicey, MS, RN-BC, CPHIMS  
Chief Innovation Officer  
Clearsense



**UC Irvine Health**

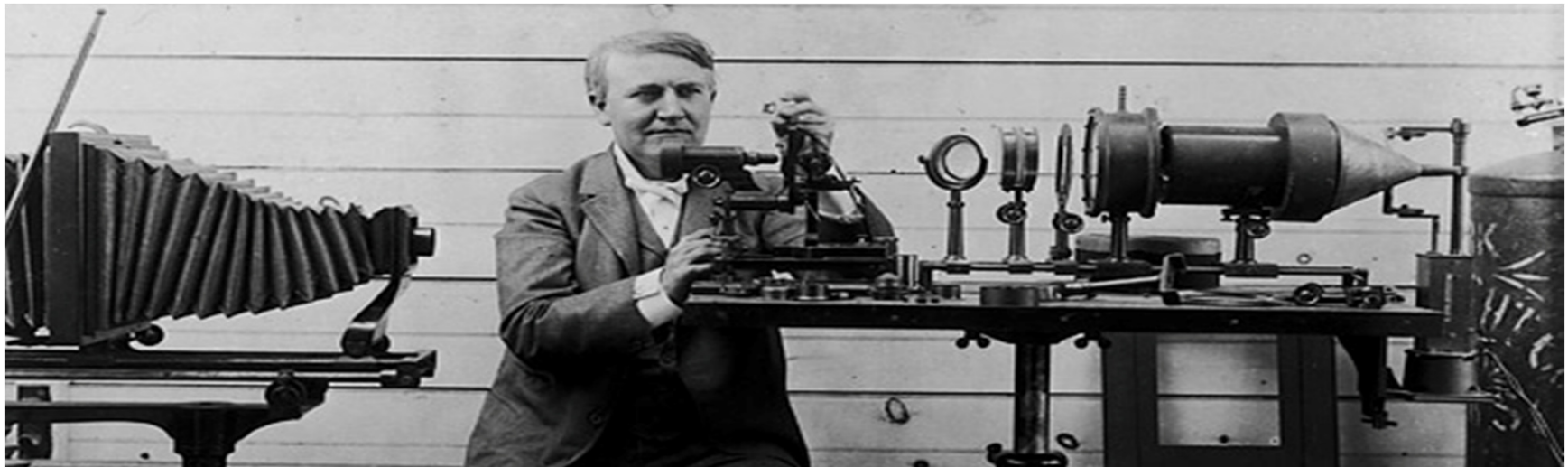


**Stony Brook  
Medicine**

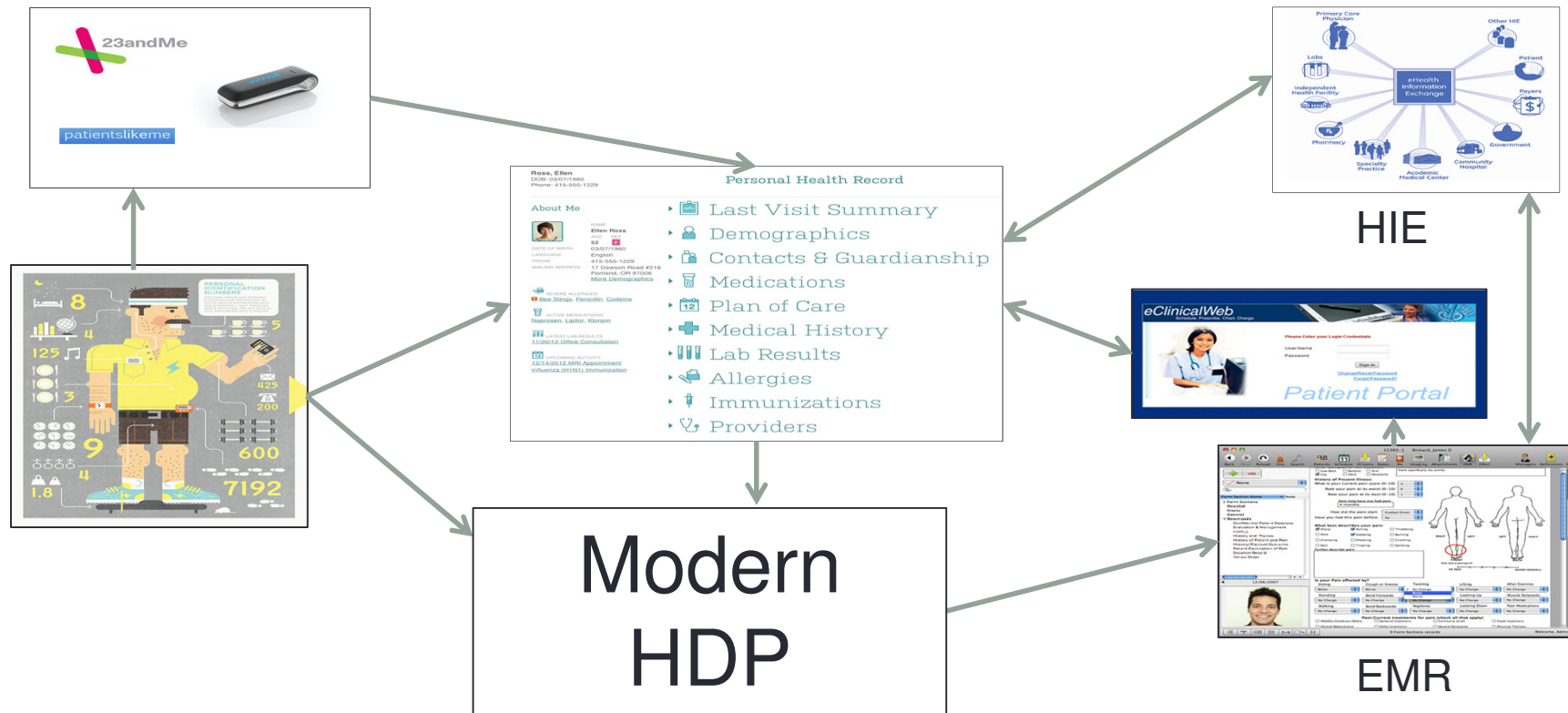


The doctor of the future will give no medicine, but instead will interest his patients in the care of human frame, in diet, and in the cause and prevention of disease.

Thomas Edison (1847 – 1931)



# PHR Centric Health



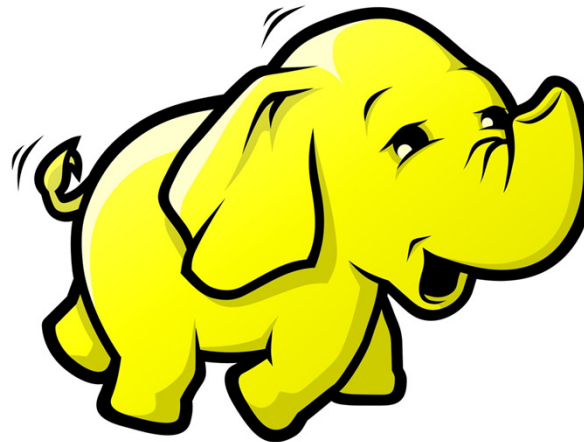
# Early Days - 2010

twitter

LinkedIn

facebook

YAHOO!



**Continuity of Care Document**  
Created On: September 16, 2010

**Patient:** Jeffrey Swann MRN: 0004201  
147 Grove Street  
Wilmington, PA, 17701  
tel: (610) 836-0000 ext 2000

**Birthdate:** September 24, 1960 **Sex:** Male

**Allergies and Adverse Reactions**

Substance	Adverse Event Type	Reaction	Status	Note
KODINE PHOSPHATE POWDER	Drug allergy	allergic drug reaction	Active	Diarrhea, nausea, vomiting
AMPICILLIN 75 250 MG CAPSULE		prurigoide to adverse reactions	Active	

**Medications**

Medication	Instructions	Start Date	Status
Atorvastatin (LIPITOR 10 MG TABLET)	1 tablet(s), oral, QD	2002-05-05	Active
Protonix (PANTOPRAZOLE 30 MG TABLET)	1 tablet(s), oral, BID	2002-05-05	Active
Formamide (LASIX 20 MG TABLET)	1 tablet(s), oral, BID	2002-05-05	Active
Glibenclamide (GLIBENCLAMIDE 2.5 MG TABLET)	1 tablet(s), oral, QD, AM	2009-09-16	Active

**Problems**

Problem Name	Type	ICD-9-CM	Status
DIABETES UNCOMPLTYPE II UNOSPEC	Disease	250.02	Active
401.9 - HYPERTENSION ESSENTIAL	Symptom	401.9	Active
CAD	Finding	414.01	Chronic
272.4 - HYPERLIPIDEMIA OTHUNSPEC	Condition	272.4	Active

**Results**

Test	LOINC	Ref	Result	Unit	Scale	Abn	Notes
HDL Cholesterol (d) - 99mg/dl	14645-4	164	146	mg/dl	164		
Total Cholesterol (d) - 206mg/dl	14647-2	164	206	mg/dl	164		
Triglyceride (d) - 1.4mg/dl	14642-9	164	1.4	mg/dl	164		
Fasting Blood Glucose (7) - 100mg/dl	14711-0	177	100	mg/dl	177		
Hemoglobin A1c (d) - 5.7%	14927-8	177	5.7	%	177		
BUN (t) - 10mg/dl	14927-7	164	10	mg/dl	164		
LDL cholesterol (d) - 150mg/dl	2089-1	164	150	mg/dl	164		
Chest X-ray, PA	24644-8						No disease is seen in the lung fields or pleura

```

<div class="view" data-bbox="98 561 326 778">
<pre>
<code>
</code>
</pre>
</div>

```

Naveen Ashish, PhD

# NowTrending 2012

The screenshot shows the NowTrending.HHS.gov website. At the top, there is a navigation bar with the text "Public Health Emergency" and "U.S. Department of Health & Human Services". Below this, the site name "NowTrending.HHS.gov" is displayed along with "beta" and navigation options for "Trends", "by condition", and "by location".

A "Welcome!" message states: "We are tracking disease trends, 140 characters at a time".

A text box explains the challenge: "In March 2012, the Assistant Secretary for Preparedness and Response at the Department of Health and Human Services launched a challenge competition titled Now Trending: #Health in My Community. This contest challenged entrants to create a web-based application that searched open source Twitter data for health topics and delivered analyses of that data for both a specified geographic area and the national level. This website is a result of that contest. The information available below and throughout the website is a tool intended for health departments and other health entities to use in multiple ways such as serving as an indicator of potential health issues emerging in the population, building a baseline of trend data, engaging the public on trending health topics, or cross-referencing other data sources."

A note states: "The data and metrics on this site represent data for up to the last two weeks. Full historical data is being maintained but is not publicly available at this time."

Three key metrics are displayed:

- 4,133,872** tweets gathered from Twitter's Streaming API. All of them match at least one of the 234 condition terms currently tracked across 27 conditions.
- 65,609 (1%)** tweets with a sensor-based location (read more about how we calculate this)
- 2,184,485 (52%)** tweets with a popular user profile location (read more about how we calculate this)

Three data sections are shown below:

- Conditions by Tweet Count:** ebola, acute respiratory illness, std, natural disaster, influenza, common cold, meningitis, gastroenteritis, tuberculosis, polio, pertussis, malaria, dengue, pneumonia, rabies, varicella, tick borne disease, legionnaires disease, tetanus, enterovirus, cholera, anthrax, measles, chagas, mosquito borne disease, mumps, typhoid, yellow fever, smallpox, diphtheria.
- Top 20 Tweet Locations:** california, us, nigeria, nevada, us, los angeles, ca, manhattan, ny, punjab, pakistan, national capital region, republic of the philippines, georgia, us, chicago, il, alaska, us, houston, tx, florida, us, kenya, texas, us, pennsylvania, us, orlando, fl, philadelphia, pa, washington, dc, oklahoma, us, san francisco, ca.
- Top 20 User Locations:** usa, london, #215love, #lcm, philly, atl, to, philly, dmV, to, philly, chi, town, to, philly, atlanta, georgia, new york, philippines, worldwide, uk, canada, california, united states, indonesia, global, nigeria, texas, india.

At the bottom left, a PDF file named "dsrip\_timeline.pdf" is visible. At the bottom right, there is a "Show All" button.

# Current Environment

## Electronic Medical Record

- Not designed to process high volume/velocity data
- Not intended to handle complex operations
  - Such as:
    - Anomaly detection
    - Machine learning
    - Building complex algorithms
    - Pattern set recognition

## Enterprise Data Warehouse

- Suffer from a latency factor of up to 24 hours
- The EDW serves all of the following retrospectively as opposed to in real time
  - Clinicians
  - Operations
  - Quality and research

# Big Data = Interoperability

- Big Data Ecosystem that Supports:
  - Hadoop (HDFS)
  - Hbase
  - Hive
  - Pig
  - MapReduce
  - Mahout
  - MongoDB (NoSQL)
- Neo 4j (Graph Database)
- Relational Data Base
- R
- Spark
- Storm
- Weka

# Big Data = Complete Data

- The Electronic Medical Record is primarily transactional taking feeds from source systems via an interface engine
- The Enterprise Data Warehouse is a collection of data from the EMR and various source systems in the enterprise
- In both cases decisions are made concerning data acquisition
- A Big Data system is capable of ingesting and storing healthcare data in total and in real time



# Modern Healthcare Data Platform

**A healthcare information ecosystem built on “Big Data” technologies should:**

Be capable of serving the needs of clinicians, operations, quality and research  
And should do so in real time and in one environment

**Should be:**

Able to ingest all healthcare generated data both internal and external in native format

**Should be:**

A platform for advanced analytics such as early detection of sepsis & hospital acquired conditions

Be enabled to predict potential readmissions

Leverage complex algorithms and be a machine learning platform



# Architecture Guiding Principles

- Architecture to minimize encumbrance on IT staff
- Ability to store all healthcare data in native form and complete
- Use of supported open source code
- Ensure architectural compatibility with commercial applications

# Infrastructure

- Low Cost of Entry & Scalable
  - Open Source
  - Commodity Hardware
    - UCI Hadoop Ecosystem
      - 10 nodes
      - 5 terabytes
    - Yahoo Hadoop Ecosystem
      - 60K nodes
      - 160 petabytes
  - Cloud Ready

# Data Sources

- Legacy Systems
  - Print to Text or Delimited String
- All HL7 Feeds (EMR source systems)
- All EMR Initiated Data (Stored Procedures)
- Device Data (in one minute intervals)
- Physiological Monitors (HL7)
- Ventilators (HL7)
- Smart Pumps
- Social Media (POC)
- Healthcare Organization Sentiment Analysis
- Patient Engagement
- Home Monitoring (POC)
- Real Time Location System (RFID)
- Hospital Sensors

# Newer Data Sources

- External Streaming Device Data
- Wearables
- Home Devices
- Social Media
- Geographic Information System (GIS) Data
- Omic Data
- Open Data
  - [www.data.gov](http://www.data.gov)
- Adverse Drug Event
  - [www.researchae.com](http://www.researchae.com)
- ***Internet of Things (IoT)***
  - Telematics
  - 5G

# Use Cases

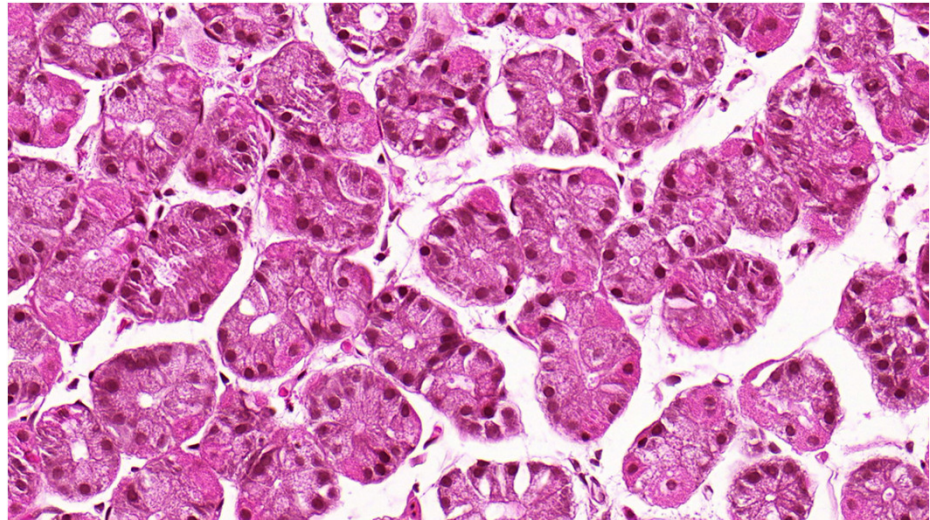
- Legacy System Retirement
- ***Patient Condition Changes***
  - ***RRT***
- ***Early Sepsis Detection***
- ***Environmental Response***
- Real Time Nursing Unit Utilization
  - Staffing and Resource Allocation
- Social Media Sentiment Analysis
- Research
- Cohort Discovery
- Data Science
- Clinician Aware Applications
- Patient Monitoring External to Traditional Healthcare Setting
- ***Event Driven Care & Real Time Quality Monitoring***
- ***Personal Health Record***

# Future Use Cases

- Ventilator Management
  - Vent dashboard in EMR
- Hospital Acquired Infections (HAI)
- VTE Surveillance
- Sensium Vitals Digital Patch
- Patient-Generated Data
  - Home Devices (Scale, Vital Signs, Glucose)
  - Exercise & Diet (Fit Bit, Jawbone, Nike)
- Combining Phenotype Data with Genotype Data
- Patient Threat Analysis
- Edge and Vertices Analysis
  - Patient caregivers and outcomes

# Imaging Analytics

- NIH Funded U24 Grant
- Joel Saltz, PhD



- This project is to develop, deploy, and disseminate a suite of open source tools and integrated informatics platform that will facilitate multi-scale, correlative analyses of high resolution whole slide tissue image data, spatially mapped genetics and molecular data for cancer research.





## Hadoop for Healthcare

### ANALYSIS

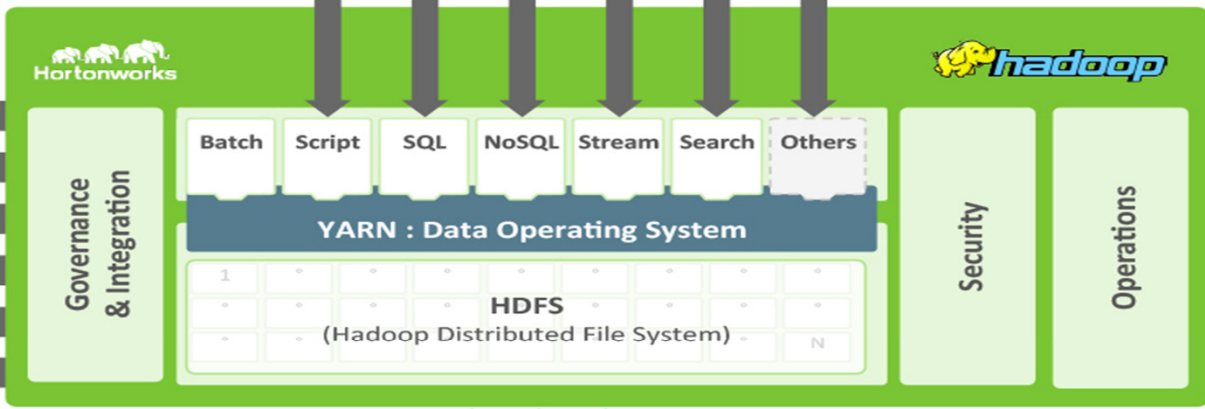
Cohort discovery  
 Predicting readmission  
 Detection of sepsis pathways  
 Analyzing test variances  
 Rapid bedside response

Tracking patient wait times  
 Home health monitoring  
 Chronic disease management  
 Patient scorecards



### DATA REPOSITORIES

- EDW
- Surgical Data Mart
- Diagnosis Data Mart
- Quality Data Mart
- Clinical Info Data Mart
- Neo4j



### TRADITIONAL SOURCES

- LEGACY EMR
- FINANCIAL
- RADIOLOGY
- PHARMACY POS
- PACS
- RTLS
- CLINICAL TRIALS
- TRANSCRIPTIONS

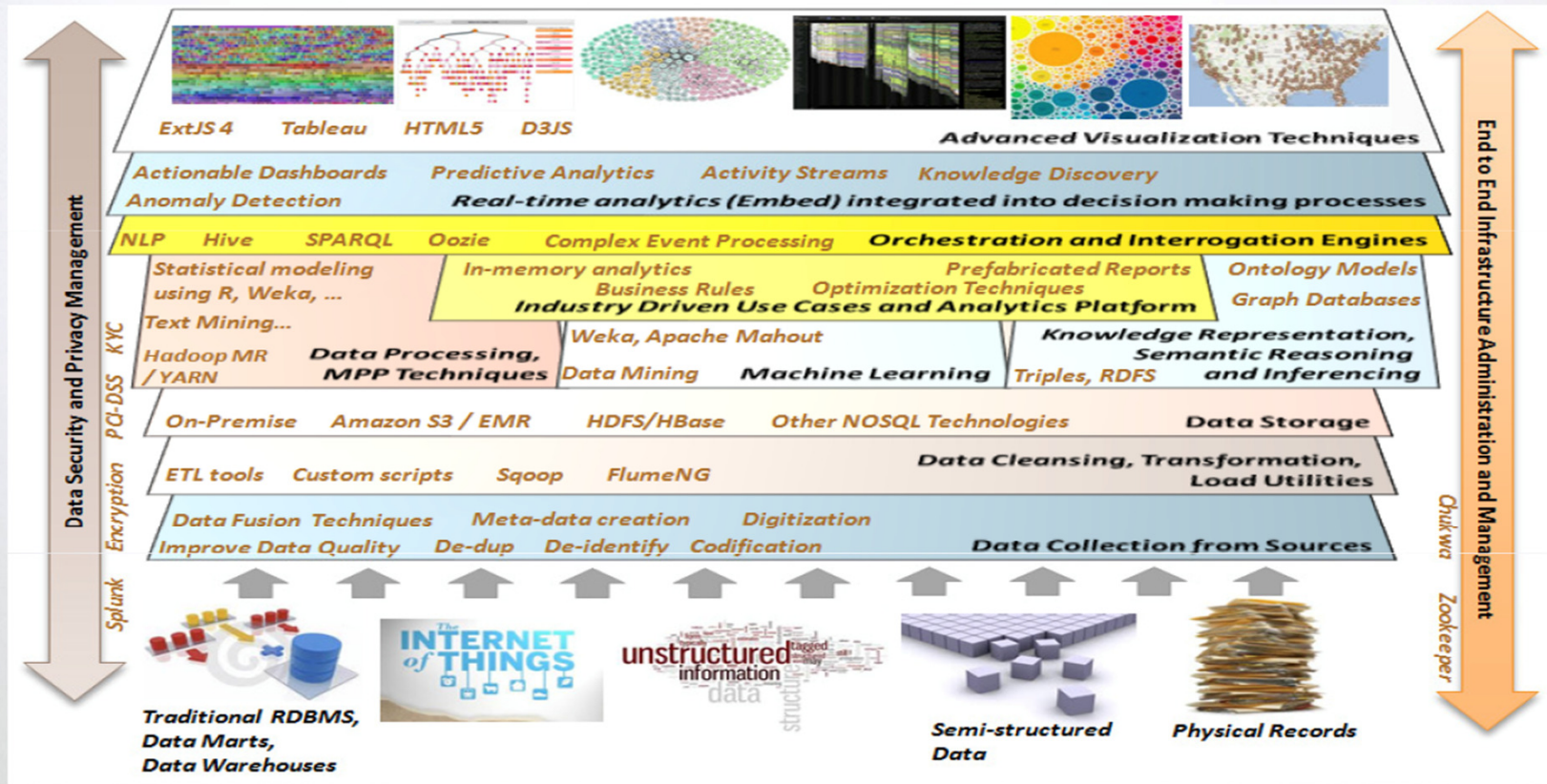
### EMERGING & NON-TRADITIONAL SOURCES

- SOCIAL MEDIA
- MEDICATION
- HOME DEVICES
- DEVICE INTEGRATION
- LABORATORY
- BIO REPOSITORY
- GENOMICS
- QUANTIFIED SELF

# FOSS Driven Protean

- Is a centrally-hosted, instrumented “Smart and Connected” platform servicing real time business event streams using high-speed MPP Compute and Storage Grids
- Primarily based on the concepts and principles of Event Driven Architecture (EDA), Complex Event Processing (CEP) and Multi-Agent-Systems (MAS)
- Support for high speed data ingestion - Structured and Unstructured (Textual)
- Core Advanced Analytics enabled through Model Building, Data Mining and Machine Learning techniques (Supervised and Unsupervised)
- Context modeling creation across Time-Space-Value dimensions
- Enables creation of a Central Enterprise Data Refinery to enable “Source of Truth” for transactional information within the Healthcare Enterprise

A reference architecture blueprint for realizing the Big Data platform leveraging Free and Open Source Software (FOSS)....The platform has been deployed successfully across 4 large client implementations across various business domains....



# FHIR – The “Public API” for Healthcare?

## FHIR = Fast Health Interoperability Resource

- Emerging HL7 Standard (DSTU 2 soon)
- More powerful & less complex than HL7 V3



## ReSTful API

- ReST = Representational State Transfer – basis for Internet Scale
- Resource-oriented rather than Remote Procedure Call (nouns > verbs)
- Easy for developers to understand and use

## FHIR Resources

- Well-defined, simple snippets of data that capture core clinical entities
- Build on top of existing HL7 data types
- Resources are the “objects” in a network of URI reference links

# SMART Platform – Open Specification for Apps

- “Substitutable Medical Apps”
  - Kohane/Mandl – NEJM (2009)
- A SMART App is a Web App
  - HTML5 + JavaScript
  - Remote or embedded in EHR
  - URL passes context & FHIR link
- EHR Data Access via FHIR
- OAuth2 / OIDC for security



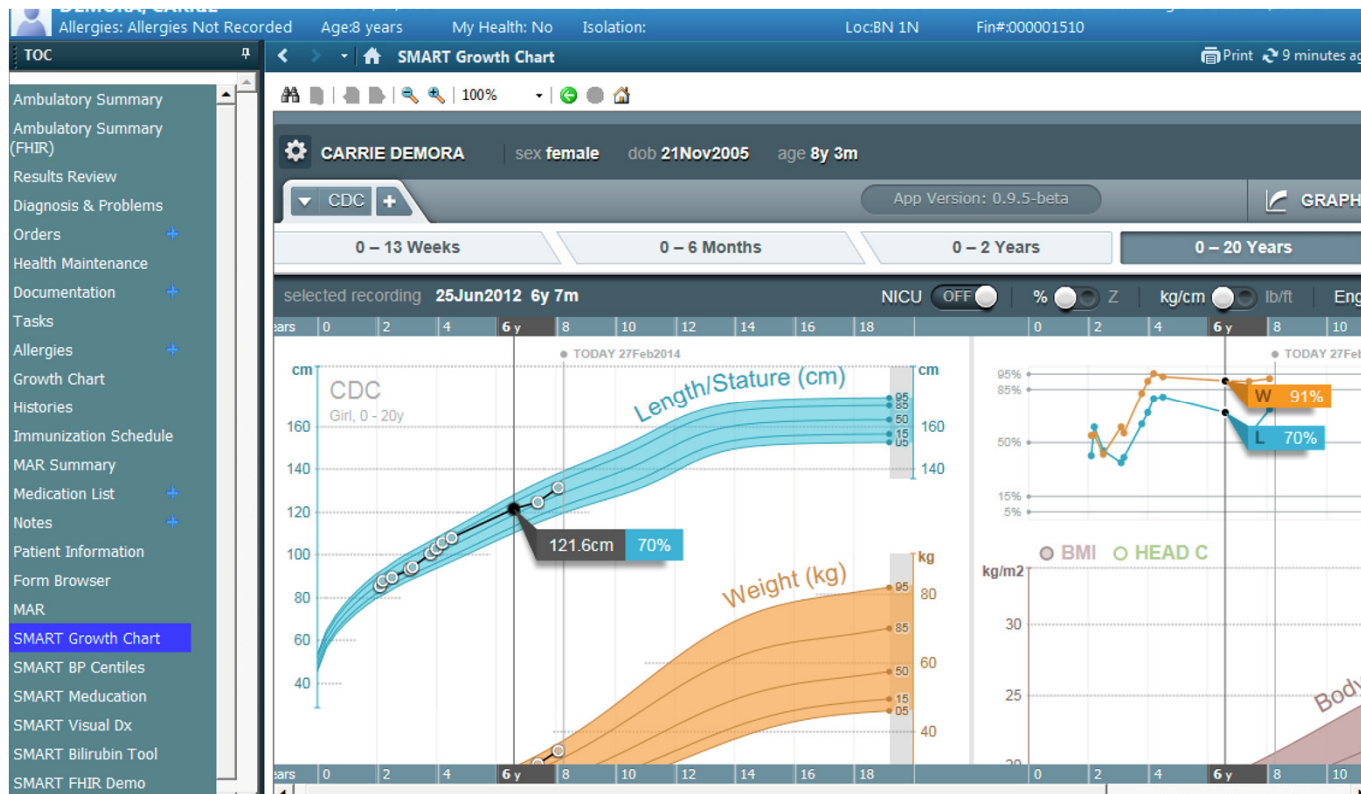
## Some SMART Hotbeds



SMART



# Boston Childrens: SMART Growth Chart



Huff, S., McCallie, D HIMSS 2015





# DSRIP

- 8 billion dollar grant (Medicaid waiver) from CMS to NY State
  - 25% reduction over five years in avoidable hospitalizations and ER visits in the Medicaid and uninsured population
  - Collaborative effort to implement innovative projects focused on
    - System transformation
    - Clinical improvement
    - Population health improvement



## 5 Year Goals

- Create integrated Suffolk County care delivery system for 387K lives anchored by safety net providers
- Engage partners across the care delivery spectrum to create a countywide network of care
- After five years, transition this network to an ACO which will contract with insurance providers on an at risk basis

# Suffolk Care Collaborative IT Architecture

Suffolk County Providers

EMRs or clinical Information System

Stony Brook Medicine

EMRs or clinical Information System

Suffolk County PPS Population Management Tools

Registries	Care Plans	Workflow	Med Adherence	Mobility
------------	------------	----------	---------------	----------

Suffolk County PPS Patient Portal

eForms	Patient Wellness	Alerts	Mobile Monitoring	Patient Education	Clinical Records	Collaboration
--------	------------------	--------	-------------------	-------------------	------------------	---------------

Suffolk County Big Data Platform

Predictive Analytics	Event Engine	Structured Data	Financial Data	Legacy Data
Machine Learning	NLP	Unstructured Data	Wearables Data	Social Data
Anomaly Detection	Rules	Device Data	HL7/CCD	Open Data

Suffolk county PPS Master Patient Index (MPI)

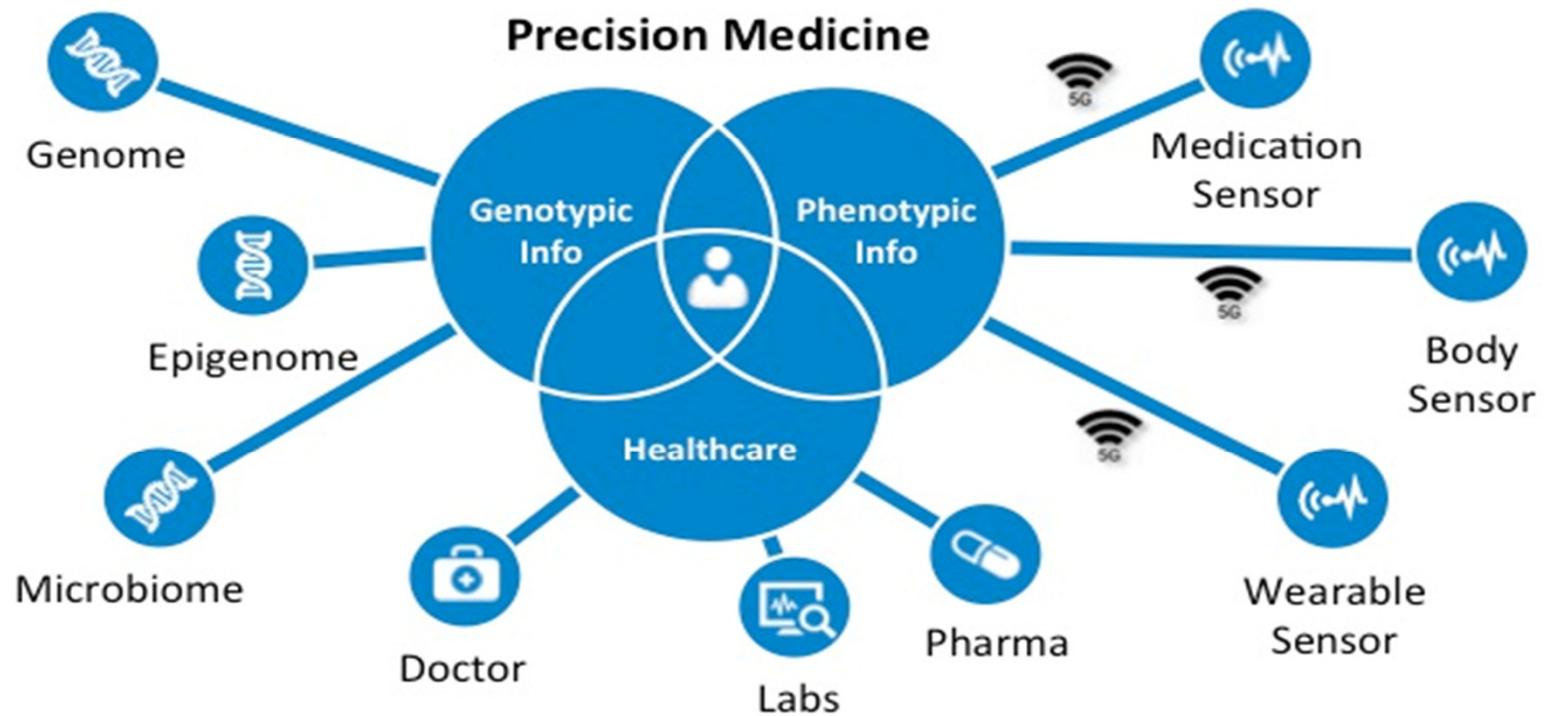
Suffolk county PPS Health Information Exchange (HIE)

E-HNLI RHIO (HIE)

Clinical Data for Patient Care



# The Future of Genomic Medicine



Gavin Stone, edico genome 5G Summit May, 14, 2015

# New Team Members

- Data Scientist
- Developers
- Cognitive and Behavioral Psychology
- User Experience
- Human & Computer Interaction
  - Devices
  - Wearables
- ***Patients & Family***



# Trends: Big Data

- **Definition:** Evolving
- **Creation & Management:** Distributed and augmented
- **Information Governance:** Shared
- **Meaningful Analysis:** Beyond PnL, Reporting, Connections, Correlations, Pattern Recognition, Machine Learning, Natural Language Processing
- **Business Requirements:** Blank Page; We don't know what we want we will figure it out once we look at the data, the data will lead the way, AKA, Data Science

## Trends: Healthcare

- Content Analytics – Suggestive Analytics\* – Prescriptive Analytics
- Imaging Analytics
- Moving Analytics out of the EMR Environment
- Graph Data Mart
  - Edge and Vertices Analysis
- Omic & Phenotype Combines
- Sentiment Analysis



# Takeaways

- Underpinning platforms may change but concept is here to stay, abstract where possible.
- Machine learning will lead to the evolution of Data Science and eventual use of AI in Healthcare.
- Get used to source now, ask questions later: Healthcare evolves with data and it is not a point in time construct any longer.
- Get used to working with constant change, disruptive trends and something new that will make your “frameworks” obsolete.



# Contact Me @

Charles Boicey

[cboicey@uci.edu](mailto:cboicey@uci.edu)

[charles.boicey@stonybrookmedicine.edu](mailto:charles.boicey@stonybrookmedicine.edu)

[cboicey@clearsense.com](mailto:cboicey@clearsense.com)

1+904-373-0831

@N2InformaticsRN



**UC Irvine Health**



**Stony Brook  
Medicine**

