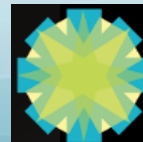


# Big Data and Computational Pathology

The Future of Pathology Informatics

Bruce Levy, MD, CPE  
Associate Chief Health Information Officer  
Associate Professor of Pathology  
University of Illinois at Chicago

THE  
UNIVERSITY OF  
ILLINOIS  
AT  
CHICAGO



**PATHOLOGY  
INFORMATICS  
SUMMIT 2016**

May 23-26, 2016  
Pittsburgh, PA

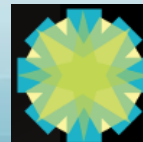
Brought to you by the Association for Pathology Informatics.

# Big Data and Computational ~~Pathology~~ Medicine

The Future of Pathology Informatics

Bruce Levy, MD, CPE  
Associate Chief Health Information Officer  
Associate Professor of Pathology  
University of Illinois at Chicago

THE  
UNIVERSITY OF  
ILLINOIS  
AT  
CHICAGO



**PATHOLOGY  
INFORMATICS  
SUMMIT 2016**

May 23-26, 2016  
Pittsburgh, PA

Brought to you by the Association for Pathology Informatics.

# Notice of Faculty Disclosure

In accordance with ACCME guidelines, any individual in a position to influence and/or control the content of this ASCP CME activity has disclosed all relevant financial relationships within the past 12 months with commercial interests that provide products and/or services related to the content of this CME activity.

The individual below has responded that he/she has no relevant financial relationship(s) with commercial interest(s) to disclose:

**Bruce Levy, MD**

# Objectives

- Explain the Intelligence Cycle
- Define Big Data
- Discuss the emerging field of Computational Pathology and provide examples of how it is already changing pathology and medicine
- Explain why our leadership in Computational Pathology/Medicine is necessary and important

# Where Are We Coming From?

The “great divide” between AP and CP



## Anatomic Pathology

- Descriptive and interpretive
- Morphology-based
- Unstructured

## Clinical Pathology

- Structured and computerized
- Quantitative
- Little interpretation

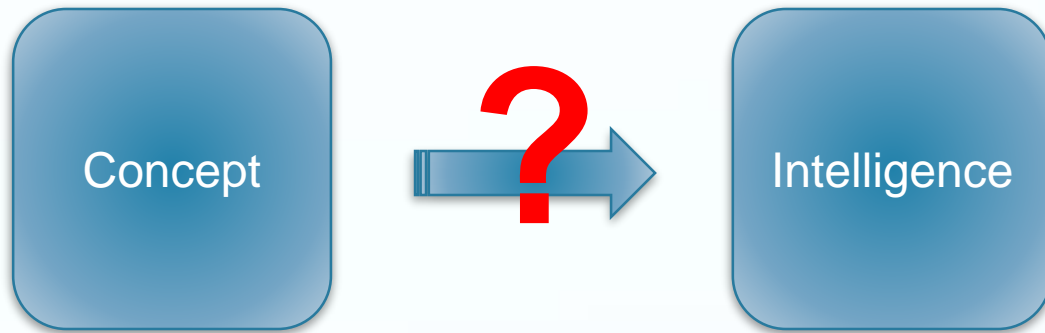


- More structure in AP
- More interpretation in CP
- Integrated reporting



How do we achieve  
intelligence?

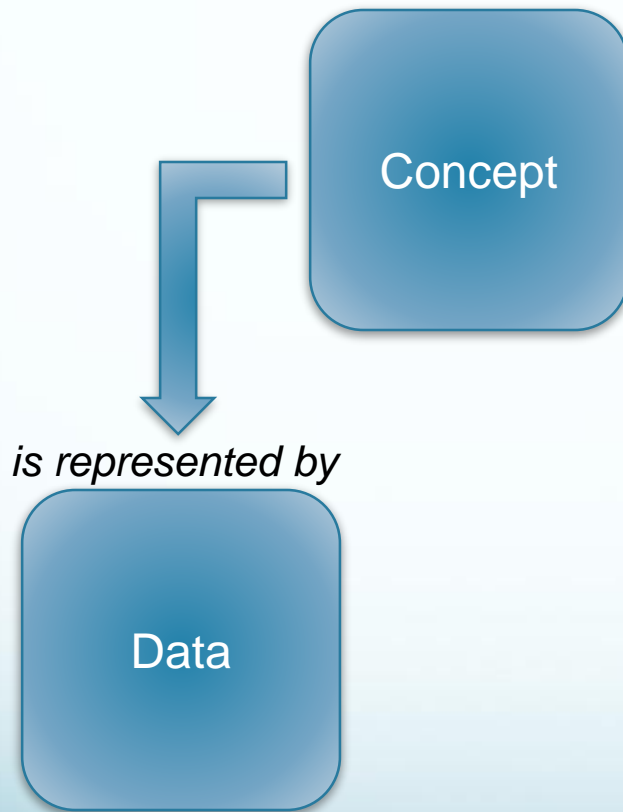
# Got Intelligence?



A concept is the  
representation of an  
idea

Expression of ideas in  
healthcare requires  
detail

# Got Intelligence?



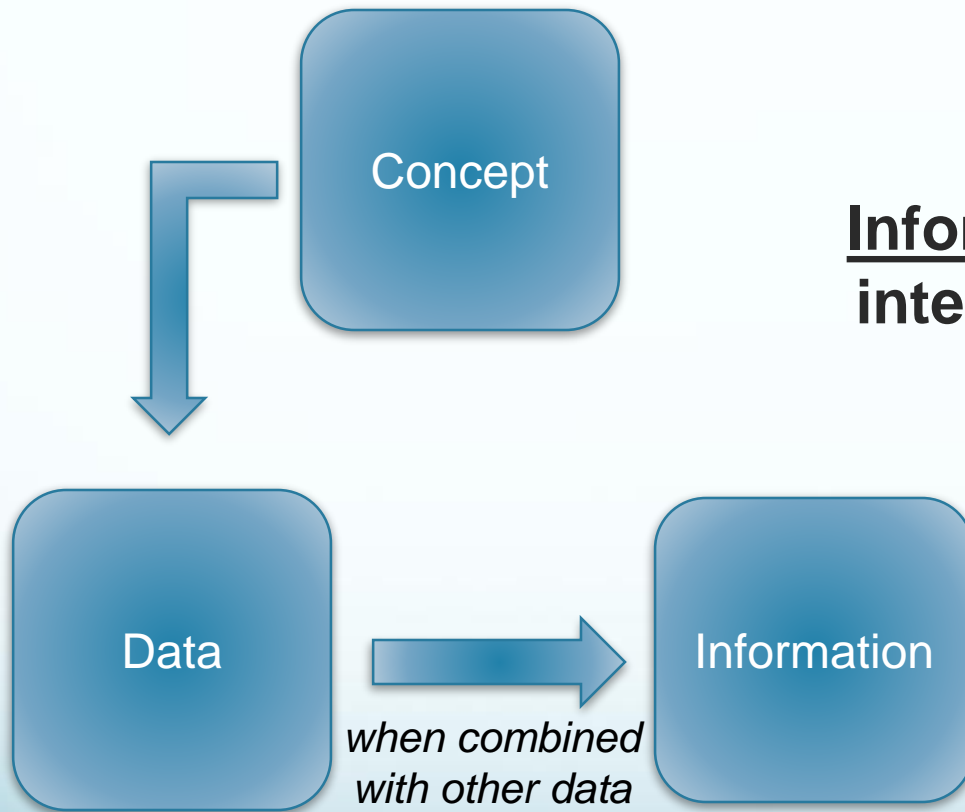
**Data are facts used as a basis for reasoning, discussion or calculation**

Examples of health data:

- Patient demographics
- Clinical signs and symptoms
- Medical, family and social histories
- Tests ordered and their results
- Problems and diagnoses
- Treatments provided
- Results of treatments
- Provider data

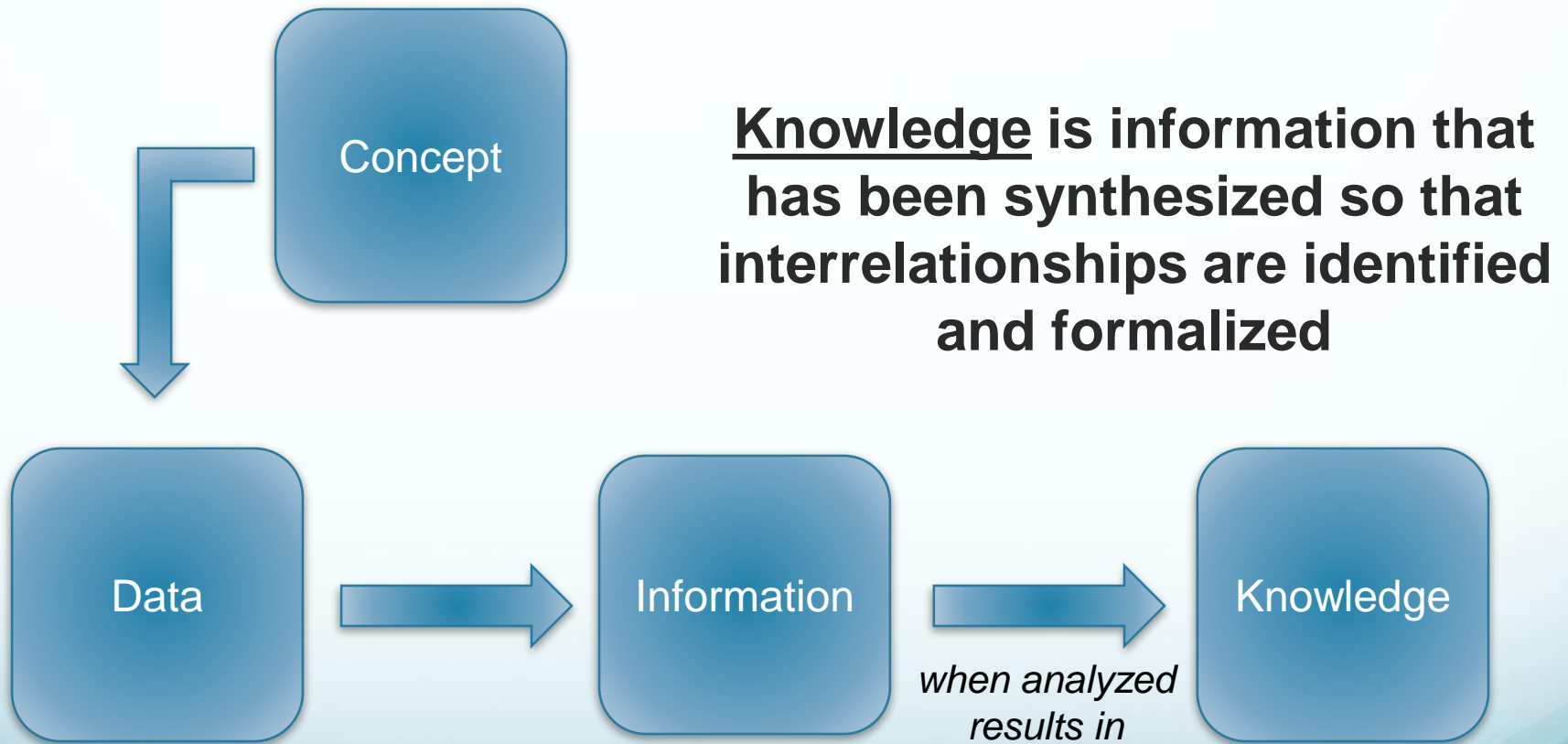


# Got Intelligence?

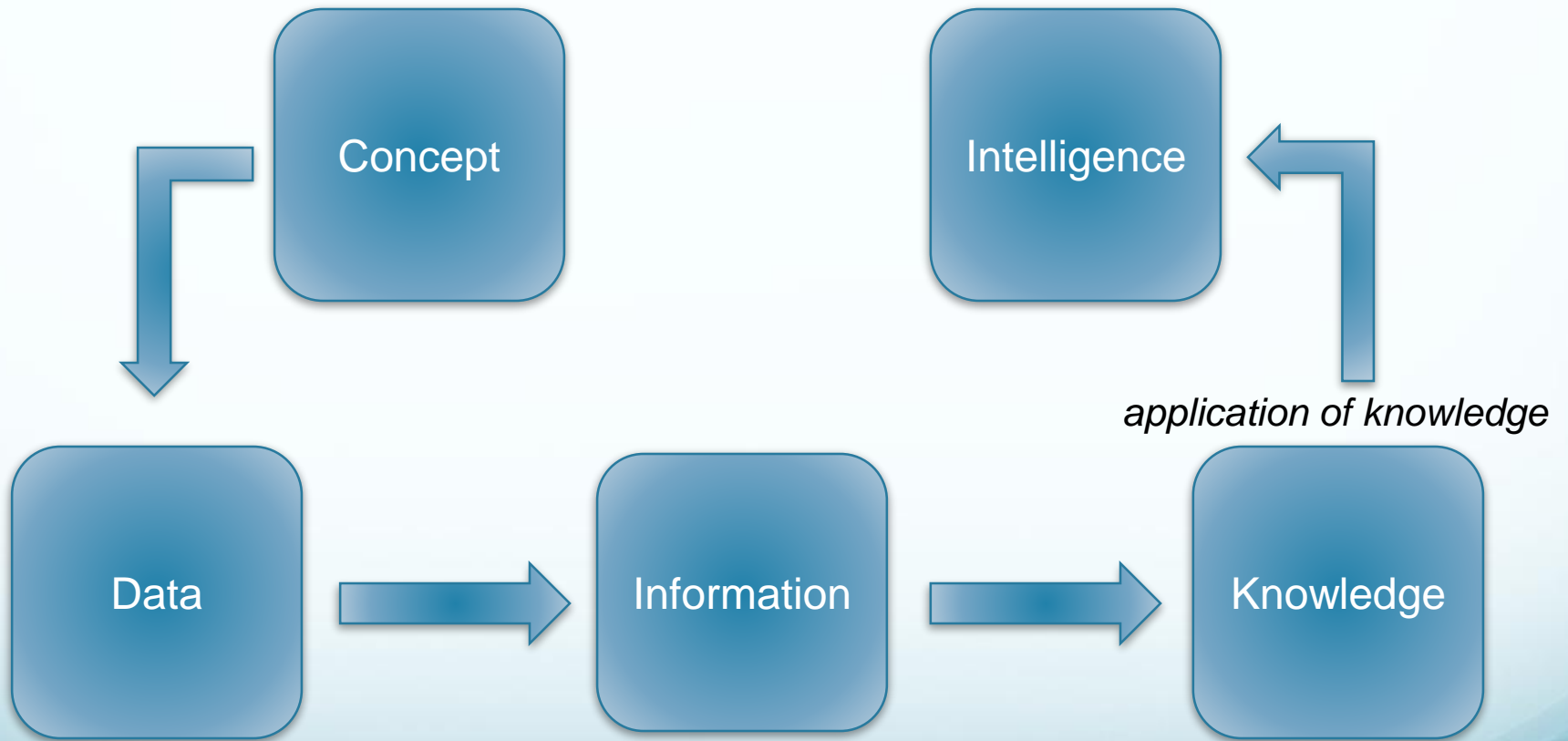


**Information is data that are interpreted, organized and structured**

# Got Intelligence?



# Got Intelligence?



# BIG DATA

A photograph of two men against a black background. The man on the left has dark, curly hair and a mustache, wearing a black shirt. The man on the right has dark hair, a beard, and is wearing a black shirt. Both men have black face paint applied to their faces, resembling stylized data points or binary code. The paint is applied in a way that suggests a digital or data-driven theme, with lines and dots forming abstract shapes around their eyes and noses.

memegenerator.net

What is 'Big Data' from the  
point of view of Pathology and  
Laboratory Medicine?



## "Big Data" in Laboratory Medicine

Moderators: Nicole V. Tolan<sup>1\*</sup> and M. Laura Parnas<sup>2</sup>

Experts: Linnea M. Baudhuin,<sup>3</sup> Mark A. Cervinski,<sup>4</sup> Albert S. Chan,<sup>5</sup> Daniel T. Holmes,<sup>6</sup> Gary Horowitz,<sup>7</sup>  
Eric W. Klee,<sup>8</sup> Rajiv B. Kumar,<sup>9</sup> and Stephen R. Master<sup>10</sup>



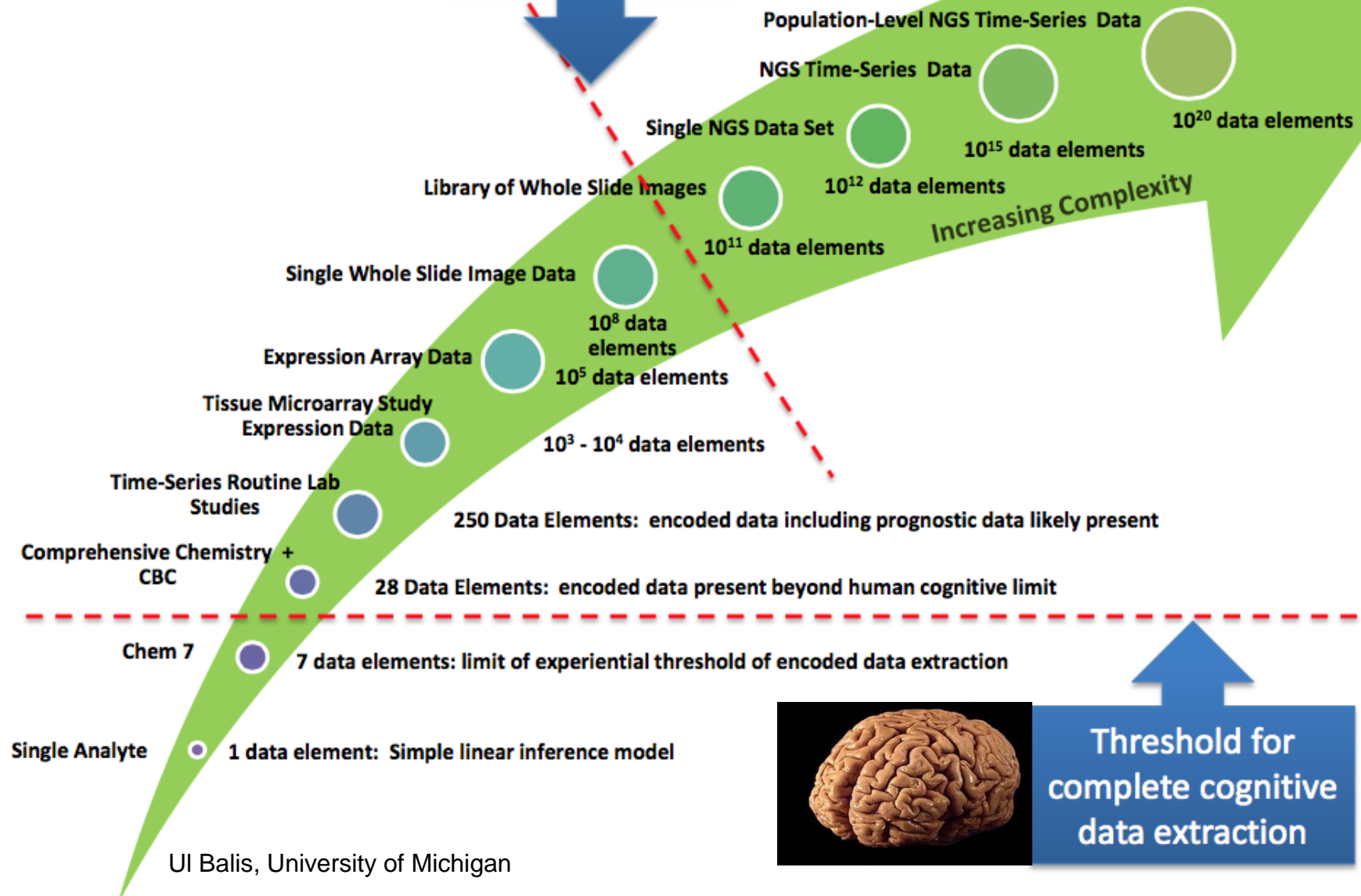
Beth Israel Deaconess Medical Center  
Sutter Health Shared Laboratory  
Mayo Clinic  
Dartmouth-Hitchcock Medical Center  
Sutter Health  
St. Paul's Hospital, Vancouver  
Stanford  
New York Presbyterian

# “Big Data” from our perspective

- “Anything that challenges an institution’s computational infrastructure”
  - “Too much information to process or store on a single computer”
- “I can’t do the analysis in Excel”
- “Any data set that challenges or exceeds an individual’s ability to manually evaluate all data points for clinical relevance”



## Supercomputing Threshold





# “Big Data” from our perspective

- “The kind of analytics that Google does”

# Detecting influenza epidemics using search engine query data

Jeremy Ginsberg<sup>1</sup>, Matthew H. Mohebbi<sup>1</sup>, Rajan S. Patel<sup>1</sup>, Lynnette Brammer<sup>2</sup>, Mark S. Smolinski<sup>1</sup> & Larry Brilliant<sup>1</sup>

<sup>1</sup>Google Inc. <sup>2</sup>Centers for Disease Control and Prevention

“One way to improve early detection is to monitor health-seeking behavior in the form of online web search queries, which are submitted by millions of users around the world each day. Here we present a method of analyzing large numbers of Google search queries to track influenza-like illness in a population. Because the relative frequency of certain queries is highly correlated with the percentage of physician visits in which a patient presents with influenza-like symptoms, we can accurately estimate the current level of weekly influenza activity in each region of the United States, with a reporting lag of about one day.”



# “Big Data” from our perspective

- “The kind of analytics that Google does”
- “Information we get from our large numbers of patients, samples and analytes in the lab”
  - “Includes all of the non-result data”
- “Assessment of massive amounts of information from multiple electronic sources in unison...to reveal otherwise unrecognized patterns”

# Big Data Definition – 3 V's

- “Big data” is high-**V**olume, -**V**elocity and –**V**ariety information assets that demand cost-effective, innovative forms of information processing for enhanced insight and decision making  
Doug Laney (Gartner) 2001
- The fourth ‘v’ is **V**eracity

## 40 ZETTABYTES

[ 43 TRILLION GIGABYTES ]

of data will be created by 2020, an increase of 300 times from 2005



## Volume SCALE OF DATA

It's estimated that **2.5 QUINTILLION BYTES**

[ 2.3 TRILLION GIGABYTES ]

of data are created each day

Most companies in the U.S. have at least **100 TERABYTES**

[ 100,000 GIGABYTES ]

of data stored

# The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015  
**4.4 MILLION IT JOBS**  
will be created globally to support big data,  
with 1.9 million in the United States



As of 2011, the global size of data in healthcare was estimated to be

**150 EXABYTES**

[ 161 BILLION GIGABYTES ]



**30 BILLION  
PIECES OF CONTENT**

are shared on Facebook every month



## Variety DIFFERENT FORMS OF DATA

By 2014, it's anticipated there will be

**420 MILLION  
WEARABLE, WIRELESS  
HEALTH MONITORS**

**4 BILLION+  
HOURS OF VIDEO**  
are watched on  
YouTube each month



**400 MILLION TWEETS**  
are sent per day by about 200  
million monthly active users



The New York Stock Exchange captures

**1 TB OF TRADE  
INFORMATION**

during each trading session



## Velocity ANALYSIS OF STREAMING DATA

Modern cars have close to **100 SENSORS**  
that monitor items such as  
fuel level and tire pressure



By 2016, it is projected there will be

**18.9 BILLION  
NETWORK  
CONNECTIONS**

— almost 2.5 connections  
per person on earth



**1 IN 3 BUSINESS  
LEADERS**

don't trust the information  
they use to make decisions



**27% OF  
RESPONDENTS**

in one survey were unsure of  
how much of their data was  
inaccurate

## Veracity UNCERTAINTY OF DATA

Poor data quality costs the US  
economy around

**\$3.1 TRILLION A YEAR**



# Other V's

- Variability – Not variety. Variability is when the meaning is changing (rapidly)
- Visualization – Challenging. How is data organized and analyzed in a manner that is easy to understand and read?
- Value – The data has to generate value. After all, the collection, storage and analysis of data costs a large amount of money. Data needs to be turned into intelligence

# Big Data Definition

- “Big data” is high-Volume, -Velocity and –Variety information assets that demand **cost-effective, innovative forms of information processing** for enhanced insight and decision making

Doug Laney (Gartner) 2001

# Big Data Definition

- “Big data” is high-Volume, -Velocity and –Variety information assets that demand cost-effective, innovative forms of information processing for **enhanced insight and decision making**

Doug Laney (Gartner) 2001



# Computational Pathology



# Computational Pathology

- Terms ‘computational’ and ‘pathology’ have been used together in literature as far back as 1980
- First defined by Fuchs and Buhmann in 2011
  - “Computational Pathology investigates a complete probabilistic treatment of scientific and clinical workflows in general pathology.”
  - “It combines experimental design, statistical pattern recognition and survival analysis within a unified framework to answer scientific and clinical questions in pathology.”

# 2014 Definition

“An approach to diagnosis that incorporates multiple sources of raw data; extracts biologically and clinically relevant information from these data; uses mathematic models at the molecular, individual and population levels to generate diagnostic inferences and predictions; and presents this clinically actionable knowledge to customers through dynamic and integrated reports and interfaces, enabling physicians, patients, laboratory personnel and other health care system stakeholders to make the best possible medical decisions.”

# 2014 Definition

“An approach to diagnosis that incorporates **multiple sources of raw data**; extracts biologically and clinically relevant information from these data; uses mathematic models at the molecular, individual and population levels to generate diagnostic inferences and predictions; and presents this clinically actionable knowledge to customers through dynamic and integrated reports and interfaces, enabling physicians, patients, laboratory personnel and other health care system stakeholders to make the best possible medical decisions.”

# 2014 Definition

“An approach to diagnosis that incorporates multiple sources of raw data; **extracts biologically and clinically relevant information** from these data; uses mathematic models at the molecular, individual and population levels to generate diagnostic inferences and predictions; and presents this clinically actionable knowledge to customers through dynamic and integrated reports and interfaces, enabling physicians, patients, laboratory personnel and other health care system stakeholders to make the best possible medical decisions.”

# 2014 Definition

“An approach to diagnosis that incorporates multiple sources of raw data; extracts biologically and clinically relevant information from these data; **uses mathematic models at the molecular, individual and population levels to generate diagnostic inferences and predictions**; and presents this clinically actionable knowledge to customers through dynamic and integrated reports and interfaces, enabling physicians, patients, laboratory personnel and other health care system stakeholders to make the best possible medical decisions.”



# 2014 Definition

“An approach to diagnosis that incorporates multiple sources of raw data; extracts biologically and clinically relevant information from these data; uses mathematic models at the molecular, individual and population levels to generate diagnostic inferences and predictions; and **presents this clinically actionable knowledge to customers through dynamic and integrated reports and interfaces**, enabling physicians, patients, laboratory personnel and other health care system stakeholders to make the best possible medical decisions.”

# 2014 Definition

“An approach to diagnosis that incorporates multiple sources of raw data; extracts biologically and clinically relevant information from these data; uses mathematic models at the molecular, individual and population levels to generate diagnostic inferences and predictions; and presents this clinically actionable knowledge to customers through dynamic and integrated reports and interfaces, enabling physicians, patients, laboratory personnel and other health care system stakeholders to **make the best possible medical decisions.**”



# Potential Benefits

1. Minimize surprise of diseases
  - a. Likelihood of disease before onset
  - b. Disease trend before complications
2. Improve quality and efficiency of care
3. Hub for data-related research
4. Better selection of patients for clinical trials

# Comp Path Use-Cases

- Molecular Pathology
- Test utilization
- Hematology
- Surgical Pathology
- Microbiology/Infectious disease
- Visualization of data/SAGE

# Test Utilization

- Study of patients in Calgary over 1 year
- Six common lab tests (cholesterol, Hb A1C, TSH, vitamin B12, vitamin D and ferritin)
- Looked at repeat rates at 3, 6 and 12 months
- Found 16% inappropriately repeated tests
- Excess costs of \$2.2 million (Canadian)
- This could be built as rules in CPOE systems

# Hematology

- Analyzing populations of CBC data:
  1. Predict the onset of anemia before patients have clinical disease

Higgins and Mahadevan. *Proc Natl Acad Sci* 2010 ;23;107(47):20587-92

1. Identify myelodysplastic syndrome in unselected outpatient populations

Raess et al. *Am J Hematol.* 2014;89(4):369–374

# Surgical Pathology

- Computational image analysis to differentiate breast intraductal proliferative lesions
  - Differentiate UDH from DCIS
  - Stratify nuclear grade in DCIS
- 116 breast bxs to build classification models
- Applied to 51 breast biopsies
  - AUC of 0.86 for differentiation
  - AUC of 0.98 for low grade v. high grade


# Antibiogram

- Antibiograms
  - Microorganism specific
  - Single antibiotic
  - Not patient specific
- Use computational pathology and informatics to change frame of reference
  - from will this drug work for this bug
  - to will this regimen work for this patient

# Weighted Incidence Syndromic Combination Antibigram

## WISCA

1. Reorganize from organism-antibiotic table to patient-infection-regimen

Organism	Antibiotic 1		Patient	Syndrome	Regimen 1
<i>E. coli</i>	R		1	ABI	1
<i>E. coli</i>	S		2	ABI	0
<i>E. coli</i>	S		3	ABI	0
<i>E. coli</i>	S		4	ABI	0
<i>E. coli</i>	S		5	ABI	0
<i>E. coli</i>	S		6	ABI	1
<i>E. coli</i>	S		7	ABI	0
<i>E. coli</i>	S		8	ABI	1
<i>E. coli</i>	S		9	ABI	0

# WISCA

1. Reorganize from organism-antibiotic table to patient-infection-regimen
2. Create a model that uses clinical and previous test data to predict the probability of coverage of each regimen
3. Input new patient's characteristics and WISCA engine outputs probabilities for regimens
4. Models automatically recalculated based on continuously collected data



# WISCA for UTI

Characteristics	Base case	patient 1	patient 2	patient 3	patient 4	patient 5	patient 6	patient 7	patient 8
Age	45	45	75	75	75	75	75	75	75
Gender	F	M	F	F	F	F	F	F	F
Hospital	1	1	1	4	1	1	1	1	1
≥4 positive UCx in prior yr	no	no	no	no	yes	no	no	no	no
Albumin <2.5 g/dL	no	no	no	no	no	yes	no	no	no
≥2 recent hospitalizations <sup>a</sup>	no	no	no	no	no	no	yes	yes	yes
Cephalosporins in last 30 d	no	no	no	no	no	no	no	yes	no
FQ in last 30 d	no	no	no	no	no	no	no	no	yes
Antibiotic Regimens	Calculated probability of Coverage								
TMP-SMX	71%	61%	73%	63%	60%	71%	75%	66%	56%
Cefazolin	77%	60%	81%	78%	78%	75%	79%	66%	61%
Ertapenem	84%	66%	85%	90%	81%	80%	79%	69%	60%
Ciprofloxacin	85%	80%	85%	72%	77%	82%	88%	90%	56%
Ceftriaxone	85%	67%	86%	85%	84%	82%	83%	74%	65%
Ceftazidime	87%	75%	89%	87%	88%	85%	84%	80%	66%
Ampicillin + Gentamicin	93%	91%	93%	90%	89%	91%	94%	93%	87%
Ampicillin-Sulbactam	94%	85%	95%	95%	92%	91%	95%	90%	88%
Ceftazidime+ vancomycin	96%	94%	96%	95%	95%	94%	96%	94%	89%
Pip-Tazo	97%	95%	96%	97%	95%	94%	95%	93%	88%
Meropenem	98%	95%	97%	98%	96%	95%	96%	95%	90%

# Visualizing BIG Data

Including Whole-Slide Images

# Scalable Adaptive Graphics Environment

- Access, display, share and collaborate in real time
  - Data-intensive information
  - Variety of resolutions and formats
  - Multiple simultaneous user input
  - Multiple drivers, no passengers
- Cloud based and web browser enabled
- Open source



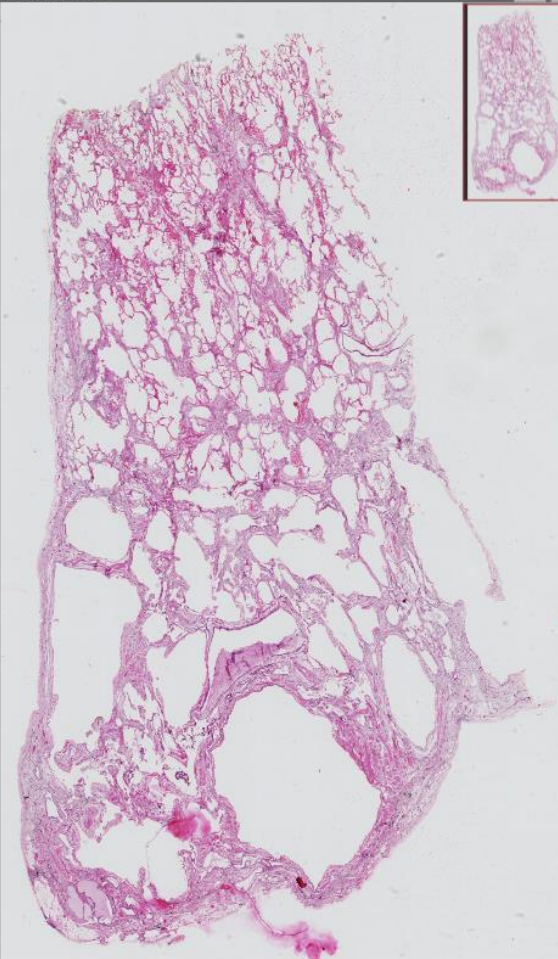
# Medical Education

4:13pm bplevy.path.uic.edu

Pathology Wall

Sevi Egan

Histology Viewer



## Asbestosis Case History.pdf

### Case 1 – PULMONARY ASBESTOSIS

#### Clinical Information:

This 68-year-old man was employed as an installer in the insulation industry in the city of Chicago for 25 years. He was forced to retire at the age of 63 because of increasingly severe dyspnea and chronic recurrent pulmonary infections. Chest X-ray studies revealed a diffuse reticular pattern consistent with interstitial pulmonary fibrosis, most severe in the lower lobes, as well as evidence of pleural thickening and calcified pleural plaques. Marked finger clubbing was noted on physical exam. Pulmonary function studies performed several months prior to death demonstrated decreased total lung capacity, forced vital capacity, and FEV1; normal FEV1 / FVC ratio; decreased forced mid-expiratory flow rate; and greatly diminished CO diffusing capacity. The patient died in severe respiratory and right cardiac failure five years after retirement.

#### Gross Pathology:

At autopsy, the lungs were small, firm, and contracted. The visceral and parietal pleurae were diffusely thickened and showed pearly white patches measuring up to 0.6 cm in thickness. Pleural adhesions were present. Sections of the lungs revealed multiple cystic spaces, especially at the bases. Interstitial and subpleural scarring was also evident.

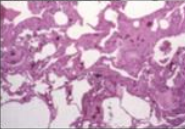
#### Microscopic Pathology:

- Thickening and fibrosis of visceral pleura.
- Thickening and fibrosis of alveolar septa (interstitial fibrosis) with loss of normal alveolar architecture.
- Dilatation and coalescence of alveolar spaces (honeycombing).
- Reactive hyperplasia of alveolar pneumocytes.
- Yellow-brown ferruginous bodies in alveolar spaces and interstitial tissues.
- Medial hypertrophy and sclerosis of pulmonary arterioles and arteries.

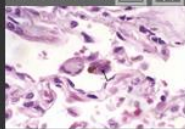
Asbestosis-Gross1.jpg



Asbestosis-Mic1.jpg



Asbestosis-Mic2.jpg



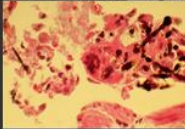
Asbestosis-Mic3.jpg



Asbestosis-Asb1.jpg



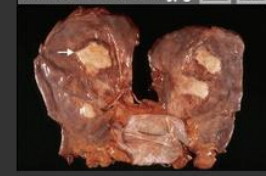
Asbestosis-Mic4.jpg



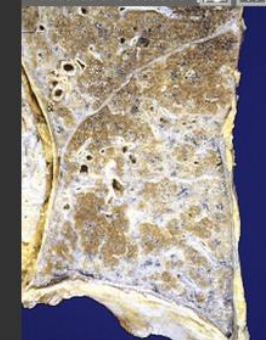
Asbestosis-GrossCut1.jpg



Asbestosis-Gross2.jpg

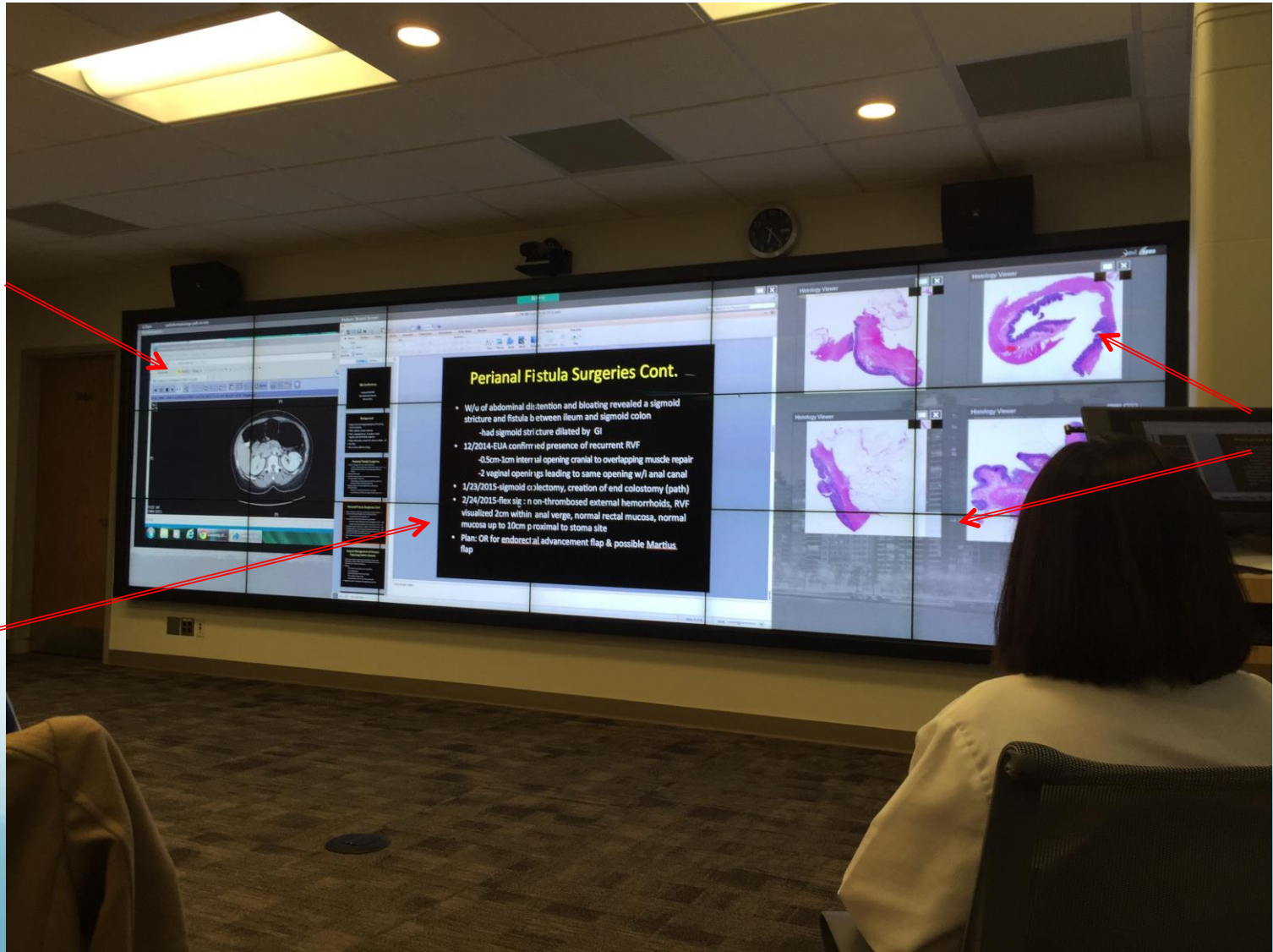


Asbestosis-GrossCut2.jpg





# Multidisciplinary Case Conference

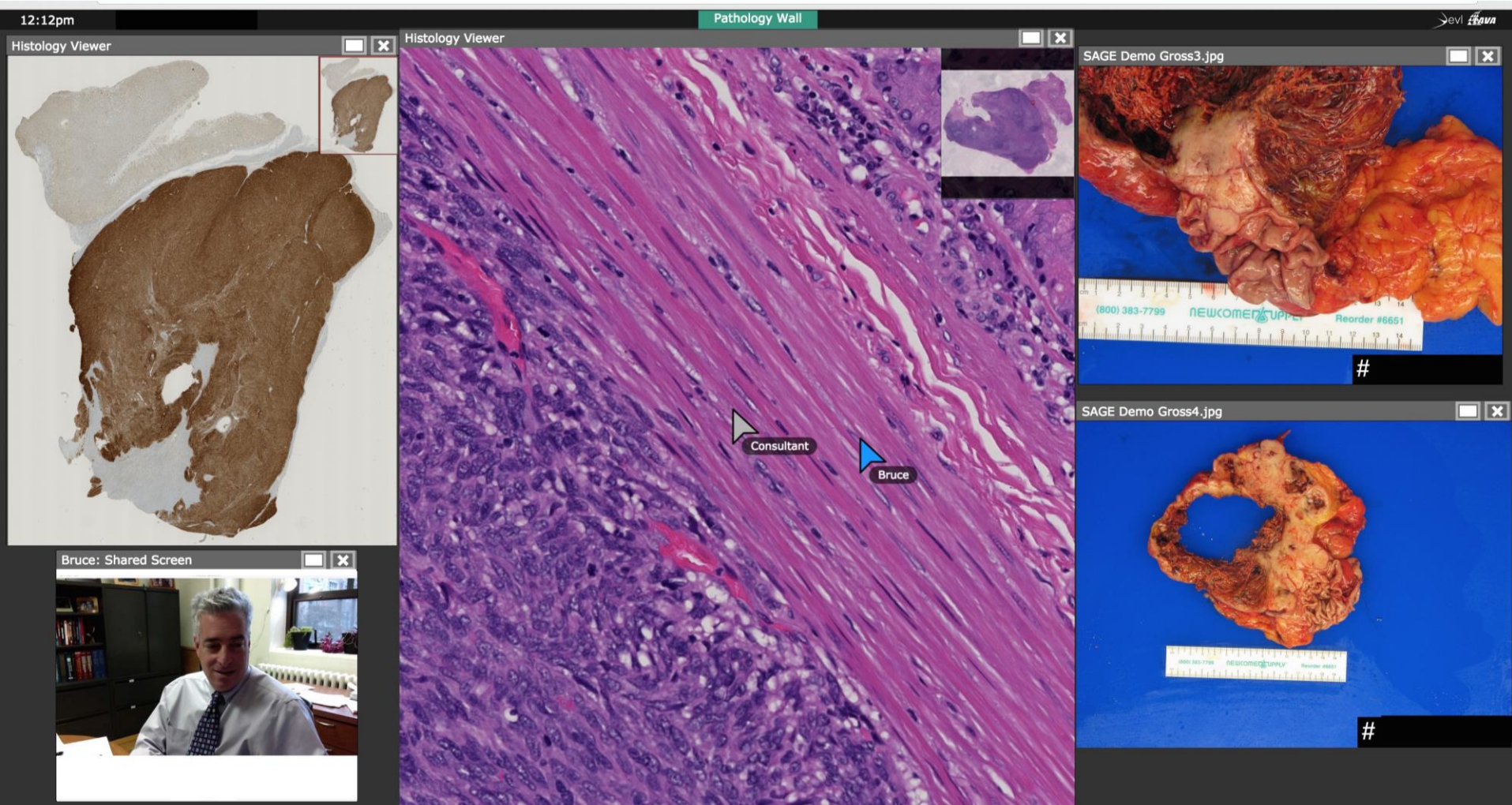


Radiology  
PACS

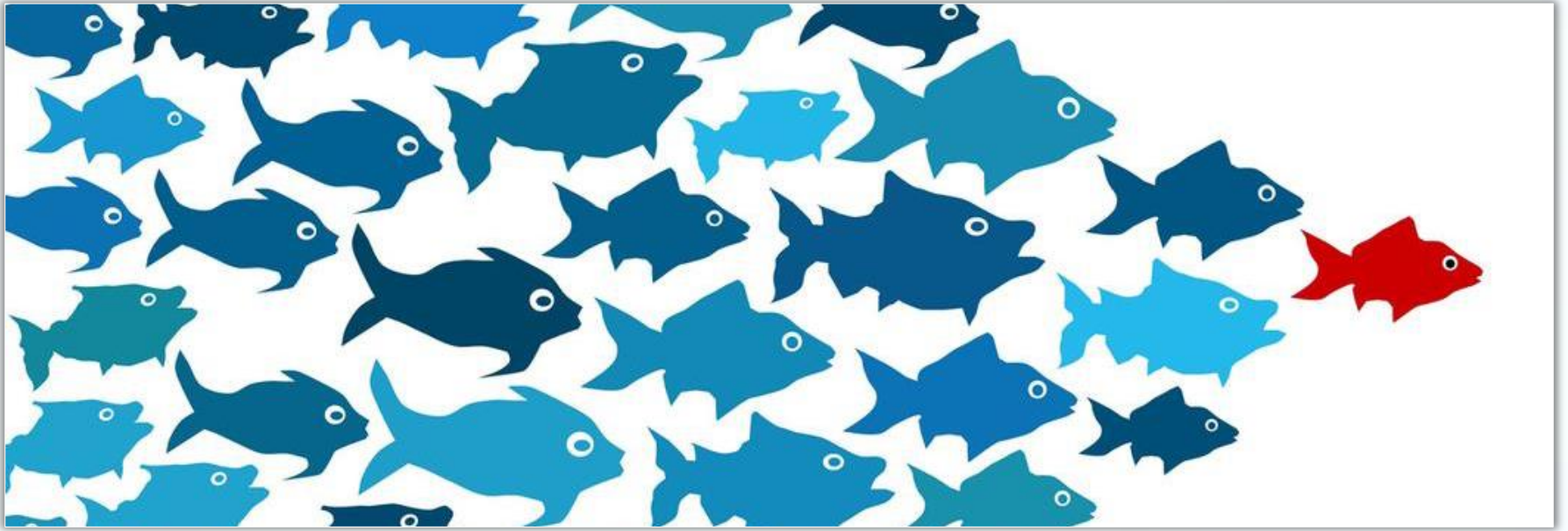
Clinical  
Info

WSIs

# Medical Consultation







# Leadership of Computational Medicine

We need to recognize and  
assume our proper place in the  
leadership of health care

“All Physicians are  
Leaders”



We need to recognize and  
assume our proper place in the  
leadership of health care

## 70% Rule

“All Physicians are  
Leaders”

We need to recognize and  
assume our proper place in the  
leadership of health care

## 70% Rule

“All Pathologists are  
Leaders”

# Operational Leadership

- Electronic Health Record
  - Participate in EHR selection/implementation
    - Leverage experiences with LIS vendors
  - Be involved with EHR committees and groups
  - Develop EHR expertise among lab professionals
- Clinical Decision Support and test utilization
- Document and communicate lab's contributions and successes
- Institutional Leadership



# Educational Leadership

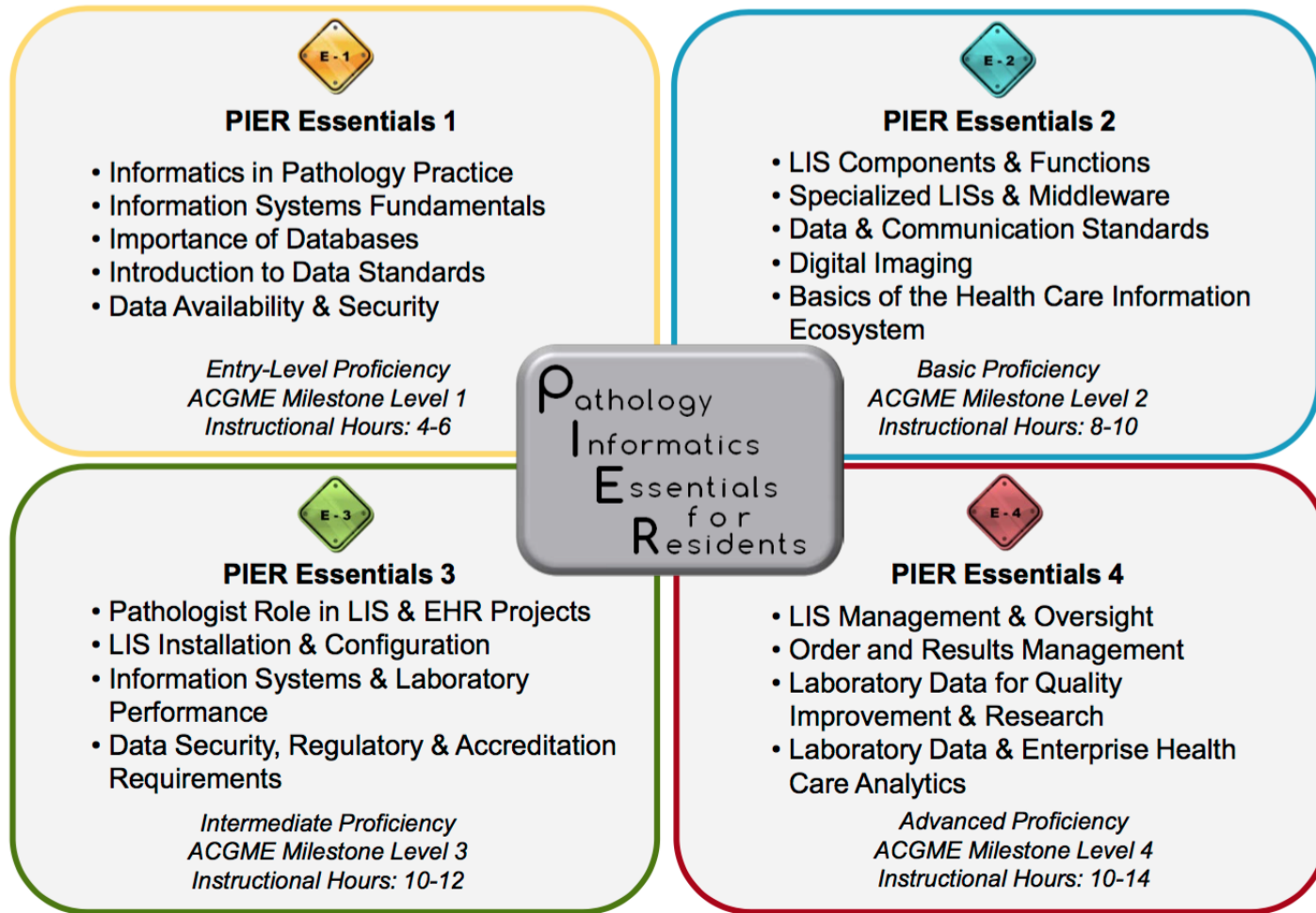
Undergraduate  
and Graduate  
Medical Education

# Undergraduate Med Ed

- Evolution of undergraduate medical education from the Flexner model is an opportunity for pathology, not a challenge
- Pathologists should lead sections of the systems-based model of medical school education
- Required lab medicine rotation
- Informatics education for medical students

# PIER: Model for all Residents

## PIER Scope and Sequence



# Clinical Informatics Fellowships



# Summary

- Intelligence in health care depends on big data
- Computational Medicine is heavily reliant on Pathology and Informatics
- Leadership of Computational Medicine is ours if we are willing to work for it
- Leadership requires us to be willing to look beyond the walls of our labs and understand our critical role in the care of patients and populations



A word cloud featuring the word "THANK YOU" in large, bold, black capital letters. Surrounding it are numerous other words in various sizes and orientations, representing different languages and scripts for "Thank You".

Words visible in the word cloud include:

- THANK YOU
- GRACIAS
- ARIGATO
- SHUKURIA
- JUSPAXAR
- DANKSCHEEN
- TASHAKKUR ATU
- YAQHANYELAY
- BIYAN
- SHUKRIA
- TINGKI
- SUKSAMA
- EKHMET
- MAAKE
- GRAZIE
- MEHRBANI
- PALDIES
- BOLZİN
- MERCI
- GOZAIMASHITA
- EFCHARISTO
- KOMAPSUMNIDA
- MAKETAI
- MINMONCHAR
- SPASSIBO
- SNACHALHUYA
- NUHUN
- CHALTU
- WABEEJA
- MAITEKA
- HUI
- YUSPAGARATAM
- DHANYABAAD
- ANIMA
- ATTO
- MERSI
- SPASIBO
- DENKAUJA
- NENACHALHYA
- UNALCHEESH
- HATUR
- GUL
- EKOJU
- SIKOMO
- BAIKYA
- TAVTAPUCH
- MEDAWAGSE
- FAKAAUE
- AGUYJE
- SAICO
- MERASTAWHY
- GAEJTHO
- LAH
- UNALCHEESH

Questions?