# Data Structures vs. Study Results:

Confessions of a failed epidemiologist
who had an informatics epiphany

JOHNS HOPKINS
M E D I C I N E

CG Chute, MD DrPH, Bloomberg Distinguished Professor of Health Informatics
April 7, 2015

# Chris's story

- Recognized clinical training as apprenticeship
  - Folklore and anecdote
- Sought methodology training for outcomes research
  - Had not yet heard of "evidence-based medicine"
- DrPH in epidemiology and biostatistics—but no data
- Sought informatics; discovered that data was junk
  - No comparability or consistency, no standards
- Established career in clinical data representation

JOHNS HOPKINS
MEDICINE

# Where did my training go wrong?

- Set up "Health Professionals' Follow-up Study" as graduate student
  - Did thesis on Nurses' Health Study

- Became far more interested in process, data collection, methods, meaning, and data quality
  - Latent informaticist, though I did not know the word

- Rather indifferent to "results" as inferences
  - Not a good sign for a junior epidemiologist

JOHNS HOPKINS
M E D I C I N E

# Why did my training go wrong?

- Exposed to 256-byte programmable calculator in HS

- Became an English major in college

- Imbued in computer science
  - All undergrads had computer accounts
  - Daily user of email (campus) since 1973
  - Lots of CS and applied math courses
  - Directed undergraduate computer consulting program

- It was in the water…
  - Musen, Cimino, Lipman, Butte, Kohane, …

JOHNS HOPKINS
M E D I C I N E

# How many boats did I miss?

- Myopic focus on clinical data generated during the process of care
  - Discount survey data
  - Discount reimbursement data
  - Discount vital statistics
  - Discount environment and exposure
  - Discount occupational health

JOHNS HOPKINS
MEDICINE

# What do I think was going on?

- Healthcare benefits from analyses
- Inferencing methodology is not sufficient
- It's all about the [clinical] DATA
    - Assume universal healthcare
    - Assume complete data capture and availability
- Data remained heterogeneous, non-comparable
- *Informatics emerged as the only path to truth!*

# One of my least-cited articles….

## Invited Commentary

## Invited Commentary: Observational Research in the Age of the Electronic Health Record

Christopher G. Chute*

* Correspondence to Dr. Christopher G. Chute, Department of Health Sciences Research, Mayo Clinic College of Medicine, 200 First Street SW, Rochester, MN 55905 (e-mail: chute@mayo.edu).

Historically, clinical epidemiologic research has been constrained by the costs and time associated with manually identifying cases and abstracting clinical data. In this issue, Carrell et al. (*Am J Epidemiol.* 2014;179(6);749–758) report on their impressive success using natural language processing techniques to correctly identify cases of cancer recurrence among women with previous breast cancer. They report a 10-fold decrease in the need for chart abstraction, though with an 8% loss in case detection. This commentary outlines some recent history associated with the development of "high-throughput clinical phenotyping" of electronic health records and speculates on the impact such computational capabilities may have for observational research and patient consent.

clinical case retrieval; electronic medical records; high-throughput clinical phenotyping; natural language processing
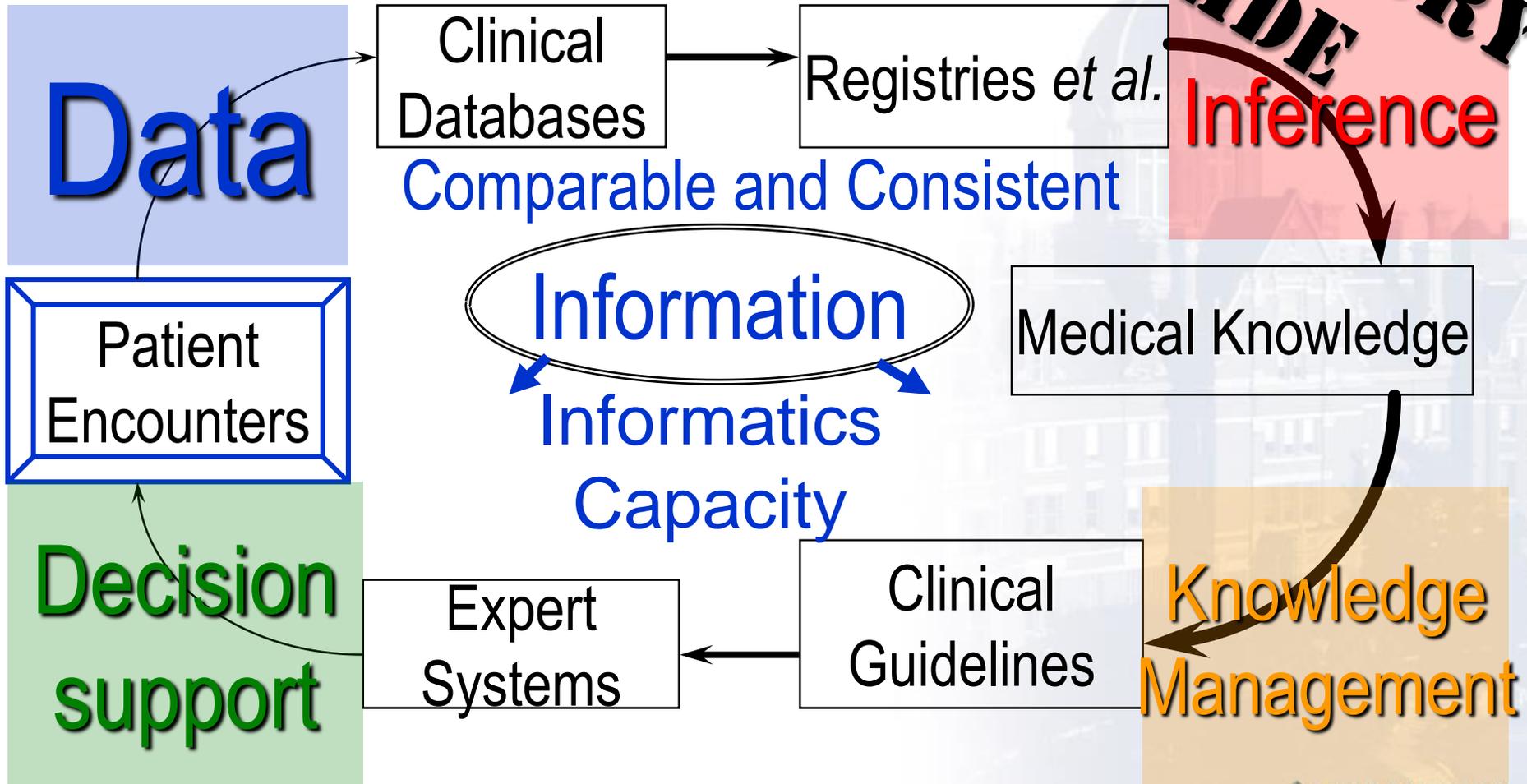
# So, what was the epiphany?
# What are epiphanati?

Within the biomedical data world:

- Comparable and consistent data is prerequisite
- That rests on semantic coherence
  - Classification, Ontology, Terminology, Value Sets
- Semantics must be bound to context
  - Information models, EHR
- Practice late-binding to application schema

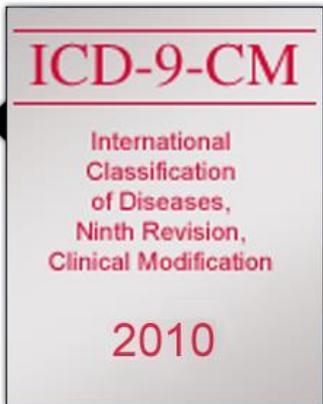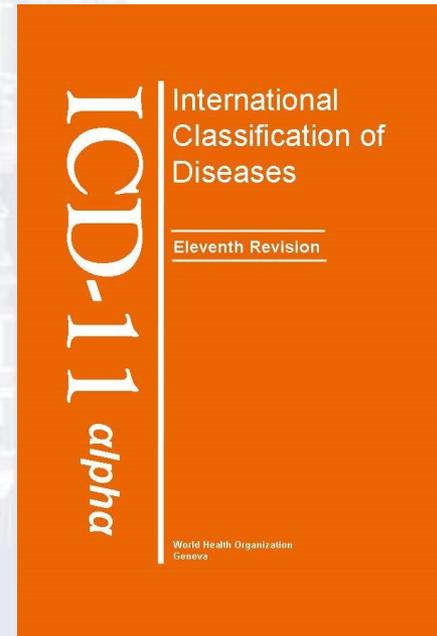# From Practice-based Evidence to Evidence-based Practice

**MANDATORY SLIDE**

**Data**

Clinical Databases → Registries *et al.*

**Inference**

Comparable and Consistent

Patient Encounters

*Information*

Medical Knowledge

Informatics Capacity

**Decision support**

Expert Systems ← Clinical Guidelines

**Knowledge Management**

JOHNS HOPKINS MEDICINE

# Coherent semantics

# Content vs. Structure
# Semantics is intertwined with structure

Heart Disease

An Information Model

Family History

Isomorphic

A Terminology Model

Family history of heart disease

JOHNS HOPKINS
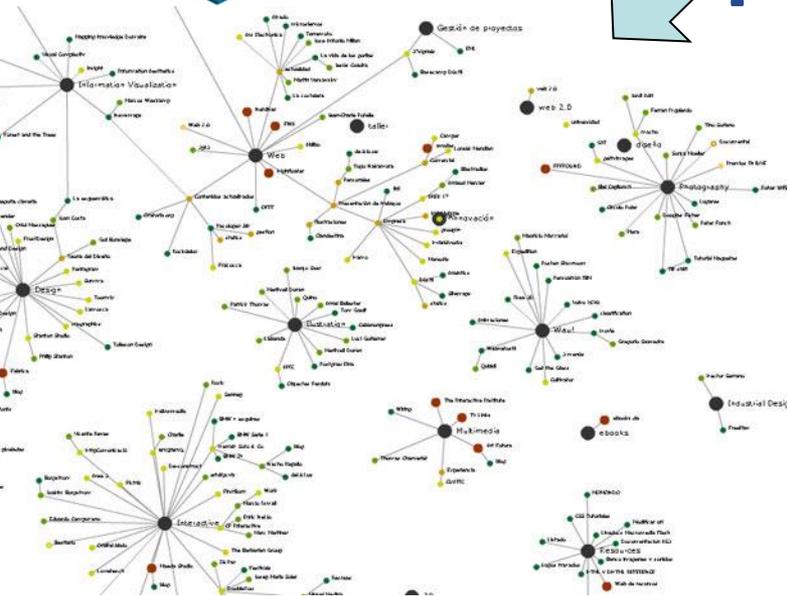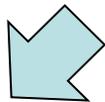MEDICINE

# Discrete data elements
# Just-in-time model binding
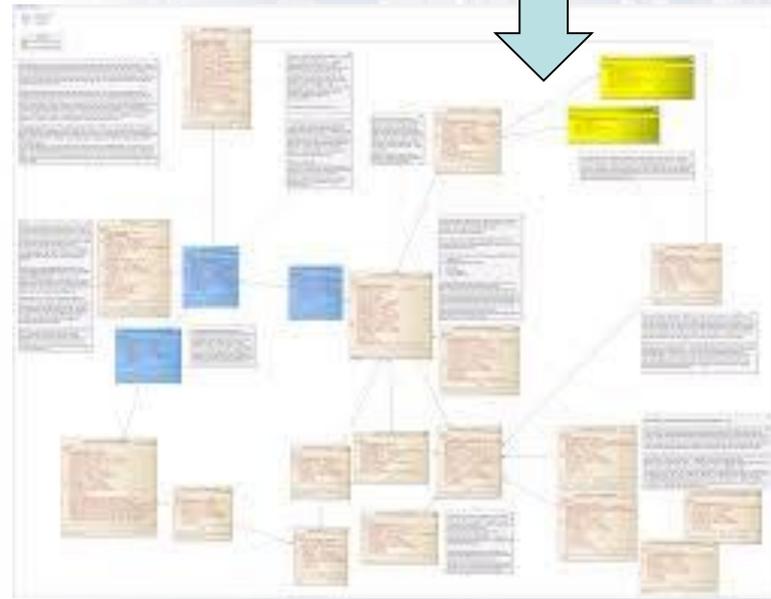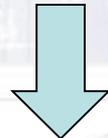
**LEGO PIECES**

## CIMI Archetypes

- Demographics
- Observations
- Medications
- Procedures
- ...

## Data Marts

- Registries
- Protocols
- Studies
- Cohorts
- ...

**VS.**

# What does any of this mean for Hopkins?

- Promote principle of "clinical data as a first-rank resource"
- Pursue the implications
  - Data governance, security, curation
  - Informatics critical mass, development, application
- Propose extension beyond Hopkins to community
  - Population health

# The "data lake"

Establish repository of clinical data

- Invoke NOSQL accumulation of data elements
- Leverage Accumulo/Topaz (Armstrong Institute)
- Leverage EPIC data warehouse
- Incorporate departmental data sources
  - Include original content and metadata
  - Capture waveforms and raw signals
  - Integrate claims data
- Incrementally normalize to canonical form

JOHNS HOPKINS
M E D I C I N E

# Maryland as a Population Laboratory

Many unique features

- CMS waiver among hospitals
- Successful emergence of CRISP
  - Framework for collaboration
- Goal of federated data repositories
  - Build on "data lake" technologies
  - Participants have secure silos

JOHNS HOPKINS
MEDICINE

# Where is this going?

- Outstanding opportunity, talent, material
- Hopkins must embrace clinical data
- Collaborate with University resources
- Collaborate with community partners
- Enable unprecedented discovery
- Rewind Chris's story

Normalized data➜Analyses➜Evidence➜Practice