# Clinical Data Mining

"If we knew what it was we were doing, it wouldn't be called research."

- Albert Einstein

Dr. Keith J. Dreyer

Partners HealthCare System  ~  Massachusetts General Hospital
Harvard Medical School

# Clinical Data Mining
## Goals

- **Query for Clinical Conditions**
- **Extract Related Information**
- **Knowledge Discovery**

# Clinical Data Mining Challenges

- **Various data sources**

- **Limited Structure**

- **Imbedded Ontology**

# Clinical Data Mining Processes

- **Digital format for all pertinent data**
- **Create structure**
    - **Obtain coded information**
    - **Natural Language Understanding**
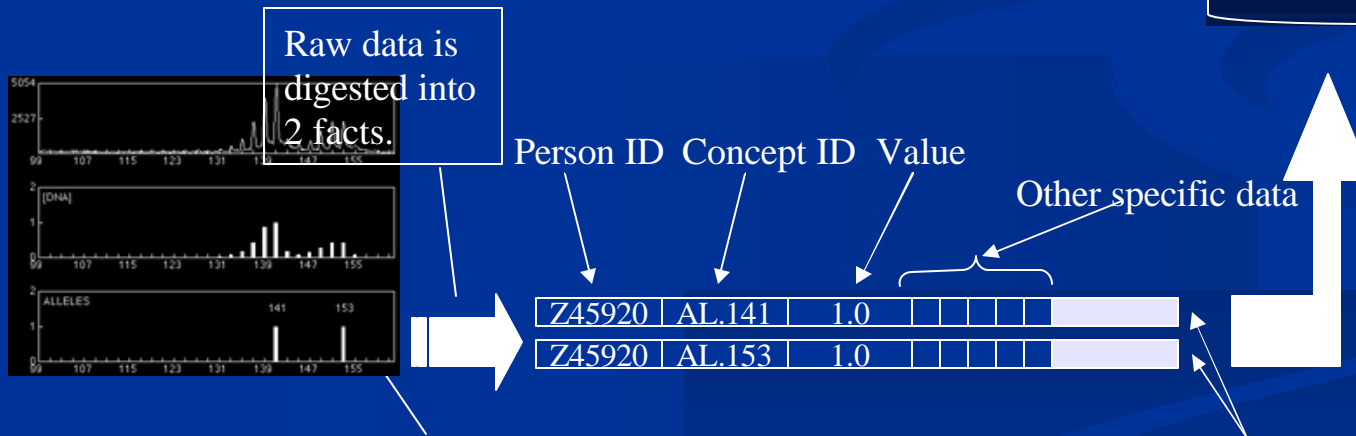- **Create a widely accessible repository**

# Query for Clinical Conditions

# Research Patient Data Registry

**Query construction in web tool**



De-identified Data Warehouse

Research data is digested into facts that can be queried with patient data

Raw data is digested into 2 facts.

Person ID  Concept ID  Value

Other specific data

| Person ID | Concept ID | Value | | | | | | | | |
|-----------|-----------|-------|---|---|---|---|---|---|---|---|
| Z45920 | AL.141 | 1.0 | | | | | | | | |
| Z45920 | AL.153 | 1.0 | | | | | | | | |

raw data goes into the RPDR as an encrypted object

# Research Patient Data Repository

- 1.8 million MGH & BWH patients

- 480 million diagnoses, medications, procedures, laboratories, with patient demographic & visit info

- Authorized use by 750 faculty status users.

- Researcher can construct complex independent queries

- Pilot of radiology image display

- Currently no radiology query capability

# Integrating Radiology Queries

- **Radiology Reports**
- **Radiology Images**

- **Challenge ?   No Structure!**

# Creating Structure

- **2 million patient demographics**
- **6 million radiology interpretations**
- **150 million radiology images**

Current Data Load: 100 TB

# LEXIMER
## Lexicon Mediated Entropy Reduction

- Radiology Language Understanding
- Lexicon Based Hierarchical Decision Trees
  - Trained to remove report noise, thus retain signal
    - Where Signal = $F_T$ and $R_T$ (defined through trained lexicon)
  - Lexicon trained by heuristics for CT and MRI
  - Expanded to all of radiology by modality training sets
    - Training stop point - 95% accuracy
  - Currently the HDTs contain 4,132 nodes

# LEXIMER

## Create an automated method to detect:

$F_T$ - The observation of a positive finding

$R_T$ - A recommendation that further imaging be performed

**Yield** = $F_T$ Exams / Total Exams

**Recommendation Rate** = $R_T$ Exams / Total Exams

# Radiology Report Understanding
## ~ LEXIMER ~

**Phrase Isolation**

**Noise Reduction** ↔ $L_N$

**Signal Extraction** ↔ $L_S$

**Findings**

Bilateral subdural hemorrhages with subarachnoid hemorrhage.

**Recommendations**

A follow up MRI of the brain is recommended within 7 days to assess progression of hemorrhage.

This study is reviewed with Dr Smith. Standard protocol was used to obtain an MRI of the brain with MRA of the circle of Willis and DWI imaging.

Dizziness and recurrent syncope. Please evaluate the posterior circulation. Comparison is to a CT of the head performed 3 September 99. Comparison is also to a CT performed the day after the MRI on 5 September 1909. Bilateral subdural hemorrhages are present. The right sided subdural hemorrhage appears improved when compared to the prior CT. It has a component extending further posteriorly than appreciated on the CT, appearing to involve the occipital lobe on the right side. The left subdural hemorrhage is worse than it appeared on the initial CT. There is extensive subarachnoid hemorrhage better appreciated on MRI than on CT.

There is no evidence of tentorial subdural hematoma. The subsequent CT did show such a bleed, this must have occurred in the interval between studies. DWI imaging of the brain parenchyma is normal in appearance. There is no evidence of acute infarction. The circle of Willis was imaged with particular attention to the posterior circulation. The right vertebral artery appears prominent. The procedure circulation appears entirely normal. Because imaging was centered on the procedure circulation, the MCA's are not completely evaluated. The ventricular system and CSF spaces do not show evidence of abnormal dilation. The visualized extracranial structures are normal in appearance.

Impression. No evidence of acute infarction on diffusion weighted imaging. Bilateral subdural hemorrhages with subarachnoid hemorrhage. The posterior circulation appears entirely normal. A follow up MRI of the brain is recommended within 7 days to assess progression of hemorrhage.

# Motivation



**MGH Annual Radiology Exam Volume**
1995 to 2003*

550,000
500,000
450,000
400,000
350,000
300,000

199
1999   2000   2001   2002   2003*

# Applications
## Annual Trends



**1995 - 2003 Exam Analysis**
3,605,223 reports, 12.4 hours

# Modality Analysis



| Composite Modality | Subtotal | Pos FX | Pos REC |
|---|---|---|---|
| GRAND TOTAL | 4072636 | 67.68 | 7.28 |
| SUBTOTAL | 4072636 | 67.68 | 7.28 |
| CHEST-X RAY | 1208059 | 72.72 | 5.60 |
| PERIPHERAL-X RAY | 662145 | 61.75 | 4.50 |
| BREAST-MAMMOGRAPHY | 279850 | 24.85 | 6.74 |
| SPINE-X RAY | 188699 | 62.54 | 6.81 |
| ABDOMEN-CT | 186699 | 76.12 | 11.86 |
| PELVIS-ULTRASOUND | 158680 | 85.99 | 13.38 |
| ABDOMEN-X RAY | 156662 | 66.79 | 7.59 |
| HEAD-MR | 134683 | 69.39 | 9.14 |
| HEAD-CT | 134308 | 66.67 | 12.26 |
| CHEST-CT | 127136 | 81.65 | 25.29 |
| CHEST-NUCLEAR MEDICINE | 81814 | 98.18 | 0.69 |
| ABDOMEN-ULTRASOUND | 63186 | 62.82 | 6.92 |
| SPINE-MR | 58097 | 82.18 | 9.53 |
| ABDOMEN-X RAY-BARIUM | 53846 | 64.58 | 3.15 |
| PERIPHERAL-ULTRASOUND | 41088 | 47.85 | 2.04 |
| BODY-NUCLEAR MEDICINE | 40746 | 65.43 | 12.10 |
| PERIPHERAL-MR | 37820 | 87.54 | 6.40 |
| SPINE-CT | 28268 | 61.38 | 8.88 |
| BREAST-ULTRASOUND | 26089 | 62.22 | 17.89 |
| | 23933 | 67.22 | 5.11 |
| HEAD-ANGIOGRAM | 22573 | 79.47 | 5.54 |
| HEAD-ULTRASOUND | 21322 | 62.00 | 7.70 |
| CHEST-ANGIOGRAM | 20680 | 84.59 | 1.07 |
| PERIPHERAL-ANGIOGRAM | 18218 | 83.77 | 1.90 |

# Applications
## Modality Analysis



Modality Fᴛ & Rᴛ Analysis
1,070 Reports, 24 Seconds

# Extraction of Related Information

# Extract Practice Pattern Inconsistencies

- **Ordering practices and yield of all referring physicians**

- **Recommendation practices of all interpreting radiologists**

- **Provide online access to individual and group level statistics**

# Practice Pattern Analysis

# Diagnostic Decisions Evaluation



Clinician Ordering Analysis
MRI Knee - 5 Clinicians

Positive Findings

100%

0%

Recommenation Rate

100%

# Diagnostic Decisions Evaluation



Radiologist Interpretation Report Analysis
CT Chest – 8 Radiologists

# Diagnostic Decisions Evaluation



Clinician Ordering Analysis
MRI Brain – 10 Clinicians

# Indication ? Examination Evaluation

Thunderclap Headaches ? Head Imaging Analysis
males <u>under</u> 40 years old

# Signal Classification

Phrase Isolation

Noise Reduction $\leftrightarrow$ **$L_N$**

Signal Extraction $\leftrightarrow$ **$L_S$**

Findings

> Bilateral subdural hemorrhages with subarachnoid hemorrhage.

Recommendations

> A follow up MRI of the brain is recommended within 7 days to assess progression of hemorrhage.

Classification $\leftrightarrow$ **$L_C$**

Concept Schema

**Post Coordinated Expressions**

Semantic Type: IS_FINDING
Etiology:  Unspecified
Finding: Hemorrhage
  Location: Subdural, Subarachnoid
  Side:  Bilateral

Semantic Type: IS_RECOMMEND
 Action:  Brain MRI
  Time:  7 Days

# Extraction of Related Information

# Extraction of Related Information

# Knowledge Discovery

# Knowledge Management

- **No standards for existing knowledge**
- **Difficulties in representing knowledge**
- **Difficulties obtaining new knowledge**

# Ontology

- The hierarchical structuring of knowledge about things through sub-categorization according to their essential qualities.

# Create Structure



- **<u>Classifying Interpreted Findings</u>**
  - **LEXIMER** (Lexicon Medicated Entropy Reduction)
  - **Codification - CPT, ICD-9, RADLEX, SNOMED**

- **<u>Image Feature Extraction</u>**
  - **Modern CAD Applications**
    - **Digital Mammography**
      - **Micro-Calcification Detection**
    - **CT Chest**
      - **Lung Nodule Detection**
  - **Future Application Examples**
    - **Brain CT, MRI**
      - **Sagittal asymmetry index, Tumor detection**
    - **Cardiac MRI, US, NM**
      - **Myocardial wall thickness quantification, HCM**
    - **Future Imaging Modalities**
      - **Molecular, Fusion, 3D Imaging**

- **<u>Essential for knowledge discovery</u>**

# Clinical Data System Integration

Research Patient Data Repository

Clinical Data Repositories

Structured Image Repository

Genetic Data Repositories

Structured Feature Sets

Research Image Archive

PACS

PACS

PACS

IMAGING MODALITIES

# Imaging and Genetics

- **Genotype**
  - An individual's genetic composition.
  - Identified by genetic analysis.

  **+**

- **Penetrance, Developmental and Environmental**

  **=**

- **Phenotype**
  - The features of health and disease expressed throughout life.
  - Can be identified, in large part, by medical imaging.

# Genetics 101

- **Allele**
  - **One of at least two forms of any individual gene.**

- **SNP** (**Single Nucleotide Polymorphism**)
  - **Small DNA regions that vary between individuals.**

- **It is these differences that represent underlying disease susceptibility and drug responsiveness.**

# Disease Knowledge Discovery

- **Understand diseases in individuals via medical imaging**
  - **Identify phenotypic expression suggestive of disease**
  - **Explore genetic analysis alternatives**
    - **Identify a potential genetic penetrance from allele or SNP**
  - **Evaluate therapy effectiveness**

- **Understand diseases in individuals via genetic analysis**
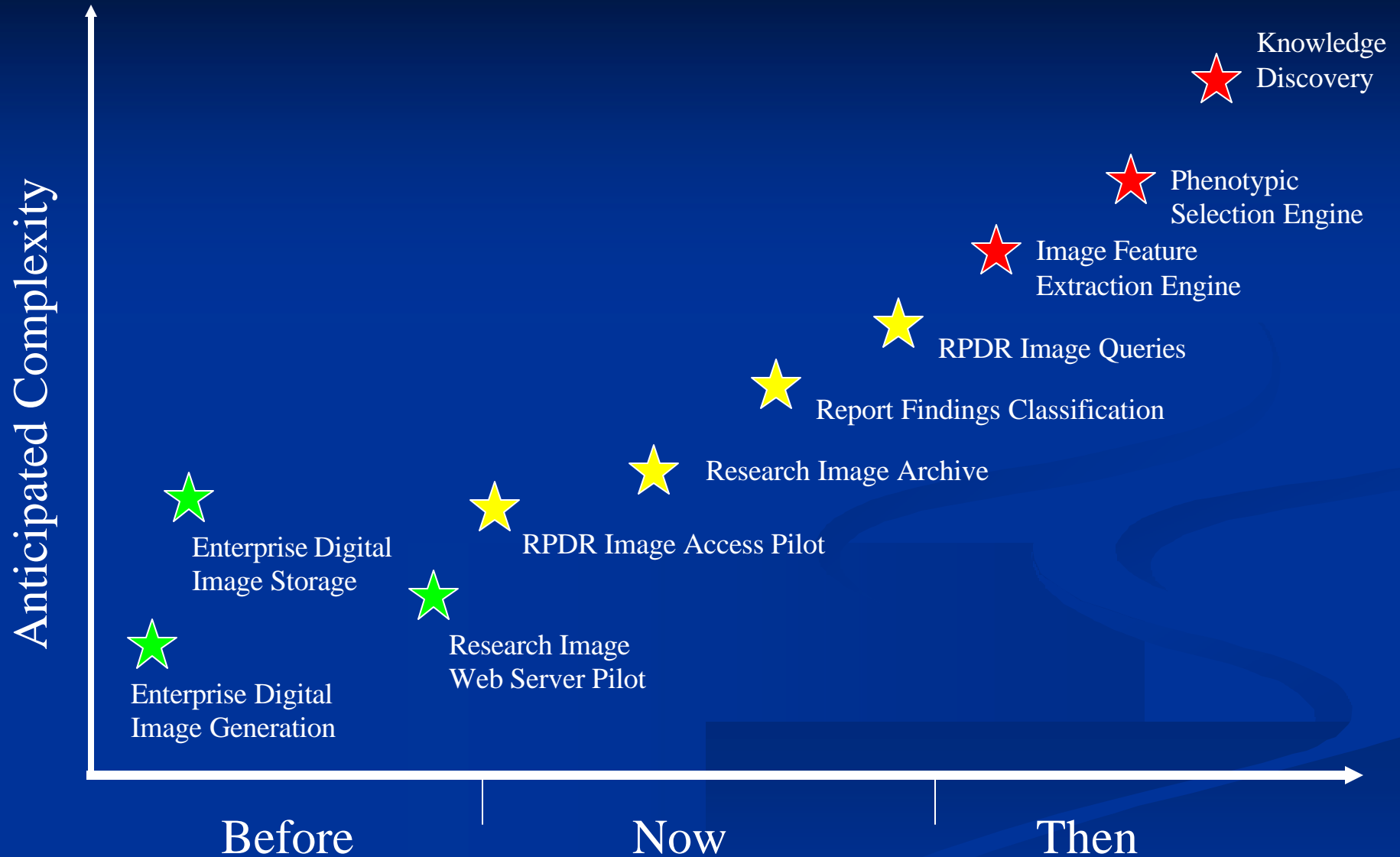  - **Identify genotypic expression suggestive of disease**
  - **Classify the pathologic spectrum of expression states**
    - **Standard Imaging - Pathologic Anatomic and Physiologic States**
    - **Molecular Imaging – Specifically Targeted Expression States**
  - **Create specific imaging work-ups for genotypes**
  - **Evaluate therapy effectiveness**

- **Discover new diseases in populations using genetic and image analysis**
  - **Utilize extensive combinations of genotypic and phenotypic data**
    - **Perform clustering algorithms to identify classification possibilities**
  - **Map genetic states to possible phenotypic outcomes**
  - **Map phenotypic presentations to possible genetic states**
  - **Monitor therapy alternatives**

# Summary Timeline



Anticipated Complexity (vertical axis)

Knowledge Discovery

Phenotypic Selection Engine

Image Feature Extraction Engine

RPDR Image Queries

Report Findings Classification

Research Image Archive

RPDR Image Access Pilot

Enterprise Digital Image Storage

Research Image Web Server Pilot

Enterprise Digital Image Generation

Before        Now        Then

# Clinical Data Mining

"If we knew what it was we were doing, it wouldn't be called research."

- Albert Einstein

## Dr. Keith J. Dreyer
Partners HealthCare System  ~  Massachusetts General Hospital
Harvard Medical School