# Instrumenting the Health Care Enterprise for Discovery Research
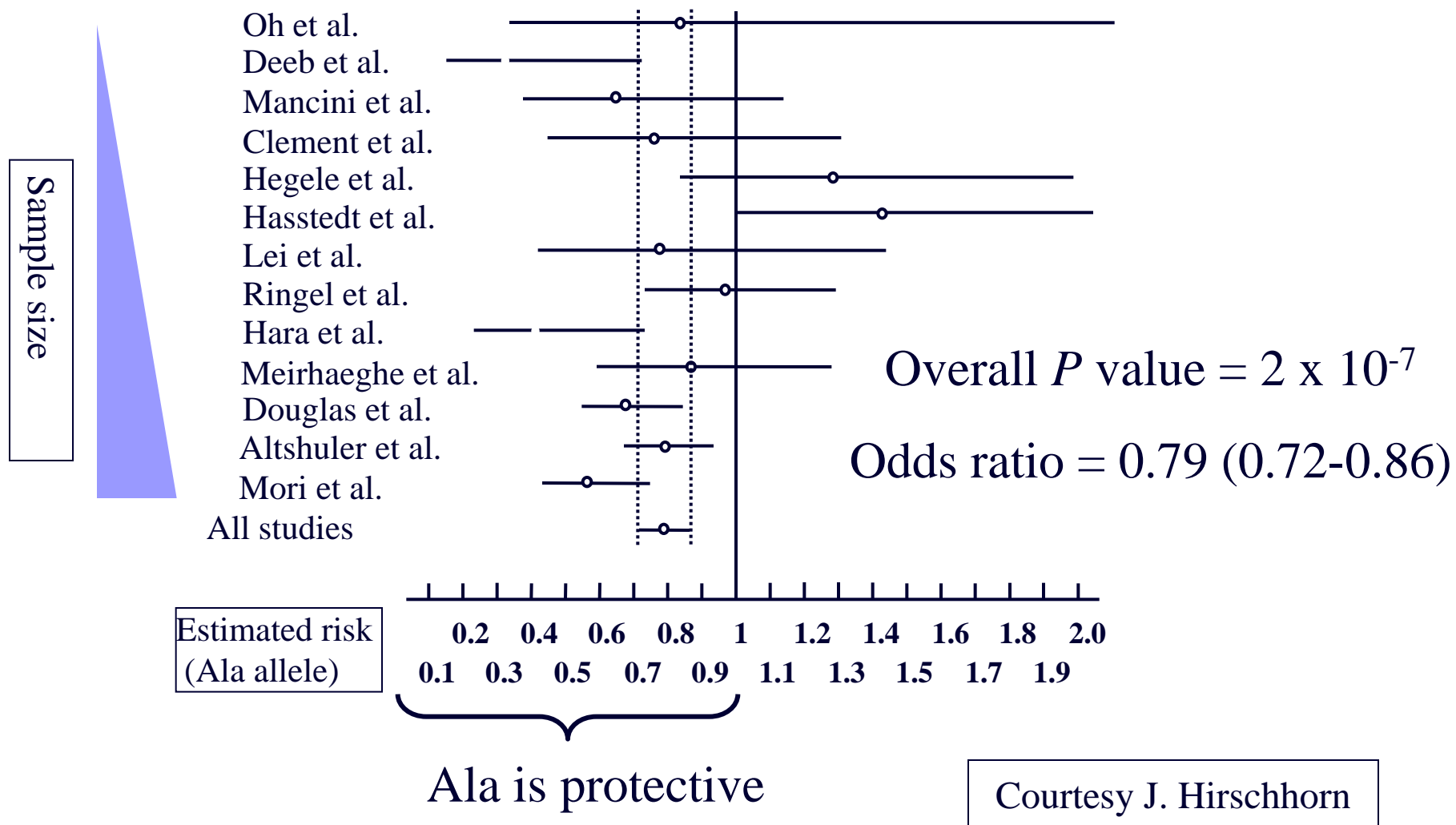
*Shawn Murphy MD, Ph.D.*
*CrEDIBLE  Project Meeting*
*October 15, 2012*

# Conflict of Interest Disclosure (Nothing to Disclose)

Shawn Murphy MD, Ph.D.

Neither I nor members of my immediate family have any financial relationships with commercial entities that may be relevant to this presentation.

# Example: PPARγ Pro12Ala and Diabetes



Sample size

| Study |
|---|
| Oh et al. |
| Deeb et al. |
| Mancini et al. |
| Clement et al. |
| Hegele et al. |
| Hasstedt et al. |
| Lei et al. |
| Ringel et al. |
| Hara et al. |
| Meirhaeghe et al. |
| Douglas et al. |
| Altshuler et al. |
| Mori et al. |
| All studies |

Overall $P$ value $= 2 \times 10^{-7}$

Odds ratio $= 0.79$ (0.72-0.86)

Estimated risk (Ala allele)

0.1  0.2  0.3  0.4  0.5  0.6  0.7  0.8  0.9  1  1.1  1.2  1.3  1.4  1.5  1.6  1.7  1.8  1.9  2.0

Ala is protective

Courtesy J. Hirschhorn

# The Power of Numbers: Efficiently Reaching a Large *N* for clinical studies

- High throughput genotyping
- High throughput phenotyping + sample acquisition

DHHS Secretary's Advisory Committee on Genetics, Health, and Society (SACGHS) argues for the health value of a 500,000 to 1M subject study. Estimated cost: $3,000,000,000

Cost of the pediatric 100,000 study recently launched >> $1B + decades.

# High Throughput Methods for supporting Translational Research
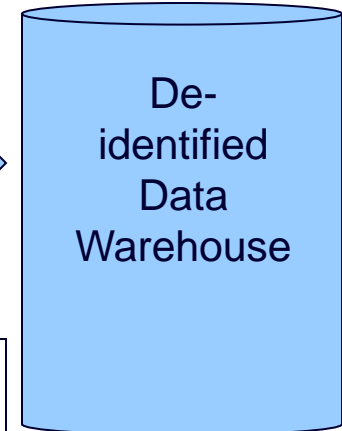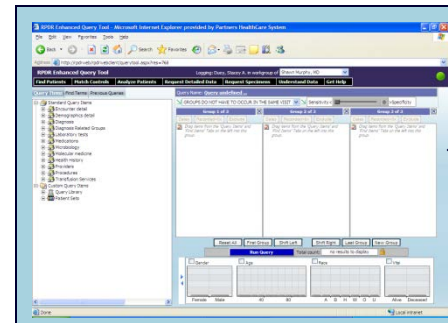
- Set of patients is selected from medical record data in a high throughput fashion

- Investigators explore phenotypes of these patients using i2b2 tools and a translational team developed to work specifically with medical record data

- Distributed networks cross institutional boundaries for phenotype selection, public health, and hypothesis testing

- Tissues of these patients can be made available for genomic and biochemical analysis

# High Throughput Methods for supporting Translational Research

- <span style="color:red">Set of patients is selected from medical record data in a high throughput fashion</span>

- Investigators explore phenotypes of these patients using i2b2 tools and a translational team developed to work specifically with medical record data

- Distributed networks cross institutional boundaries for phenotype selection, public health, and hypothesis testing

- Tissues of these patients can be made available for genomic and biochemical analysis

# Research Patient Data Registry exists at Partners Healthcare to find patient cohorts for clinical research

**Query construction in web tool**

## 1) Queries for aggregate patient numbers

- Warehouse of in & outpatient clinical data
- 6.0 million Partners Healthcare patients
- 1.5 billion diagnoses, medications, procedures, laboratories, & physical findings coupled to demographic & visit data
- Authorized use by faculty status
- Clinicians can construct complex queries
- Queries cannot identify individuals, internally can produce identifiers for (2)



De-identified Data Warehouse

Z731984X
Z74902XX
...
...

Encrypted identifiers

## 2) Returns identified patient data

- Start with list of specific patients, usually from (1)
- Authorized use by separate IRB Protocols
- Returns contact and PCP information, demographics, providers, visits, diagnoses, medications, procedures, laboratories, microbiology, reports (discharge, LMR, operative, radiology, pathology, cardiology, pulmonary, endoscopy), and images into a Microsoft Access database and text files.

0000004
2185793
...
...

OR

0000004
2185793
...
...

Real identifiers

# Security and Patient Confidentiality of Step 1

- All patients at Partners are added
  - HIPAA notification that their data may be used for research upon registration.

- RPDR data is anonymized at the Query Tool.
  - Aggregated numbers are obfuscated to prevent identification of individuals; automatic lock out occurs if pattern suggests identification of an individual is being attempted.

**A Security Architecture for Query Tools used to Access Large Biomedical Databases**

Shawn N. Murphy, MD, Ph.D. and Henry C. Chueh, MD, M.S.

Laboratory of Computer Science, Massachusetts General Hospital, Boston, MA.

- Queries done in Query Tool available for review by RPDR team, a user lock out will specifically direct a review.

- De-identified data warehouse is a "Limited Data Set" by HIPAA
  - Medical record numbers are encrypted and obvious identifiers are removed from data.

- Concept of "established medical investigator" is promoted by classification as a faculty sponsor.

# Security and Patient Confidentiality of Step 2

- Only studies approved by the Institutional Review Board (IRB) are allowed to receive identified data.

- Queries may be set up by workgroup member, but faculty sponsor on IRB protocol must directly approve all queries that return identified data.

- Special controls exist when distributing data regarding HIV antibody and antigen test results, substance abuse rehab programs, and genetic data, due to specific state and federal laws.

- Queries that return identified data are reviewed (retrospectively) by the IRB.

# 2011's usage of RPDR

- 2,733 registered users, 457 new in 2011

- 462 teams gathering data for research studies

- 1852 detailed patient data sets returned to these teams, containing data of 7.8 million patient records.

- From a survey of 153 teams
    - Importance of the data received from the RPDR was evaluated in relation to the study it was supporting.
    - The adequacy of the match of a patient profile that could be obtained through the RPDR query tool was estimated.

- $94-136 million total research support critically dependent on RPDR from patient data received throughout life of funding.

- ~300 data marts were created to support hospital operations, representing about 80 million patient records

**Usefulness of Detailed Data**
*106 Total Responses*

Not Useful 15%
Critical 43%
Useful 42%

**% of Patients Who Fit Required Profile**
*105 Total Responses*

< 10% 19%
> 75% 33%
25% - 50% 26%
50% - 75% 22%

# Organizing data in the Clinical Data Warehouse

## Star schema

**Concept DIMENSION**
concept_key
concept_text
search_hierarchy

**Patient-Concept FACTS**
patient_key
concept_key
start_date
end_date
practitioner_key
encounter_key
value_type
numeric_value
textual_value
abnormal_flag

**Encounter DIMENSION**
encounter_key
encounter_date
hospital_of_service

**Patient DIMENSION**
patient_key
patient_id (encrypted)
sex
age
birth_date
race
deceased
ZIP

**Pract . DIMENSION**
practitioner_key
name
service

## Binary Tree

start
search

.16
150
6.0
.06
1500 million

# FINDING PATIENTS



Query items

Person who is using tool

Query construction

Results - broken down by number distinct of patients

# MATCHING PATIENTS



Previous query items

Case set construction

Control set construction

Estimate set size and run program

# Obtaining Data Extracts

**RPDR Detailed Data Request Wizard -- Web Page Dialog**

Using Partners IRB#2002P000381 (Research Patient Data Registry (RPDR)) to obtain data from the RPDR

You are logged in as Duey, Stacey A. in workgroup Shawn Murphy, MD

## Please enter your IRB protocol.

Partners IRB (required): 2002P000381

Title: Research Patient Data Registry (RPDR)

Status: Active - Ongoing

Newton Wellesley Hospital IRB:

Spaulding Rehabilitation Hospital IRB:

North Shore Medical Center IRB: NSM 2008-786 demo

Title:

Status:

Options for returned set of patients:

☑ Exclude Partners Healthcare employees

☐ Create a static set of patients from this query that can be used in other RPDR queries

☑ Rerun the base query shown above to obtain a fresh set of patients

| Help | < Back | **Step 3** | Next > | Cancel |

# RPDR Detailed Data Request Wizard -- Web Page Dialog

Using Partners IRB#2002P000381 (Research Patient Data Registry (RPDR)) to obtain data from the RPDR

You are logged in as Duey, Stacey A. in workgroup Shawn Murphy, MD

## Select the types of data that should be returned from the RPDR
## Only data allowed by your protocol should be chosen

(Identified data sets will always return a set of identified patient medical numbers)

**Detail Data Items**

- ☐ Allergy Data from PEAR (Partners Enterprise Allergy Repository)
- ☐ Demographic Data
- ☐ Identifying Patient Information - not available for Limited Data Sets
- ⊞ 📁 LMR (Longitudinal Medical Record)
- ⊞ 📁 Medications, Diagnoses and Procedures
- ⊟ 📁 Patient Clinical Reports- not available for Limited Data Sets
  - ☐ Cardiology Reports
  - ☐ Discharge Summaries
  - ☐ Endoscopy Reports
  - ☐ Microbiology Data
  - ☐ Operative Notes
  - ☐ Pathology Reports
  - ☐ Pulmonary Reports
  - ☐ Radiology Reports
  - ☐ Transfusion Data, Blood Bank Data
- ☐ Top three providers for each patient

| Help | < Back | **Step 9** | Next > | Cancel |

# High Throughput Methods for supporting Translational Research

- Set of patients is selected from medical record data in a high throughput fashion

- Investigators explore phenotypes of these patients using i2b2 tools and a translational team developed to work specifically with medical record data

- Distributed networks cross institutional boundaries for phenotype selection, public health, and hypothesis testing

- Tissues of these patients can be made available for genomic and biochemical analysis

**The National Center for Biomedical Computing entitled Informatics for Integrating Biology and the Bedside (i2b2), what is it?**

- Software for explicitly organizing and transforming person-oriented clinical data to a way that is optimized for clinical genomics research
    - Allows integration of clinical data, trials data, and genotypic data
- A portable and extensible application framework
    - Software is built in a modular pattern that allows additions without disturbing core parts
    - Available as open source at https://www.i2b2.org

# i2b2 Cell: The Canonical Software Module

Business Logic

i2b2

Data Access

Data Objects

HTTP XML

(minimum: RESTful)

# An i2b2 Environment (the Hive) is built from i2b2 Cells



"Hive" of software services provided by i2b2 cells

A

B

C

Data Model

Data Repository Cell

local

remote

# I2b2 Software components are distributed as open source

# Set of patients is selected through Enterprise Repository and data is gathered into a data mart



EDR

Selected patients

Data directly from EDR

Data from other sources

Data imported specifically for project

Project Specific Phenotypic Data

Automated Queries search for Patients and add Data

# Data is available through the i2b2 Workbench

# Team support for Projects

# NLP (and comedy) is not pretty



SOCIAL HISTORY: The patient is married with four grown daughters, **uses tobacco**, has wine with dinner.

**Smoker**

PRINCIPAL DIAGNOSIS:    LEFT LOWER LOBE PNEUMONIA

SOCIAL HISTORY:  The patient is a **nonsmoker**.  No alcohol.

**Non-Smoker**

SOCIAL HISTORY:  **Negative for tobacco**,  alcohol, and IV drug abuse.

PAST MEDICAL HISTORY:  (1) Hip fracture.  (2) Bronchiectasis.

BRIEF RESUME OF HOSPITAL COURSE:
63 yo woman with COPD, **50 pack-yr tobacco (quit 3 wks ago),** spirometer...

**Past Smoker**

ALLERGIES:  (1) Aspirin.  (2) Ciprofloxacin.  (3) Penicillin.

SOCIAL HISTORY: The patient lives in rehab, married. **Unclear smoking** history from the admission note…

**???**

PHYSICAL EXAMINATION:  Temperature 97.2, pulse 80, respirations 20,
blood pressure 160/63, oxygen saturation 95% on room air.  HEENT:  Normocephalic and atraumatic.  Pupi

HOSPITAL COURSE:  ... It was recommended that she receive …We also added Lactinax, oral form of Lac**tobac**illus  acidophilus to attempt a repopulation of her gut.

**Hard to pick**

Chest x-ray revealed hyperinflated lungs and...

HOSPITAL COURSE:  The patient was seen and evaluated by the

SH: widow,lives alone,2 children,no **tob/**alcohol.

**Hard to pick**

The patient was discharged home on 8/18/99 to finish a five day course of azithromycin

Ln 44 Col 1      274     WR      Rec Off No Wrap DOS INS NUM

# Investigator Review

# *Can We Trust the Phenotypes?*

## Validation Study (N = 185)

- Evaluate case and control algorithms compared to gold standard of diagnostic interview by expert clinician
- Recruit cases and controls as defined by informatics algorithm
- Interview by clinicians blinded to ascertainment group
- Recruited patients with depression or schizophrenia to enhance blinding

Jordan Smoller MD, ScD and team

# Train classification algorithms

1. Over 300 words/phrases (features) were identified using chart review

2. Important features were selected for model using adaptive LASSO shrinkage

Tianxi Cai PhD and team



# of selected features = 29

ORIGINAL ARTICLE

# Using electronic medical records to enable large-scale studies in psychiatry: treatment resistant depression as a model

R. H. Perlis[1,2]*, D. V. Iosifescu[1,3], V. M. Castro[4], S. N. Murphy[5], V. S. Gainer[4], J. Minnier[6], T. Cai[6], S. Goryachev[4], Q. Zeng[7], P. J. Gallagher[2], M. Fava[1], J. B. Weilburg[1], S. E. Churchill[8], I. S. Kohane[9] and J. W. Smoller[2]

Use NLP to define cohorts of treatment-resistant and treatment-responsive depression

Specificity: 95%
AUC > 85%



| Clinical Status | Model | Specificity | Sensitivity | Precision | AUC |
|---|---|---|---|---|---|
| Depressed | Billing Codes | 0.95 | 0.09 (0.03) | 0.57 (0.14) | 0.54 (0.02) |
| Depressed | NLP | 0.95 | 0.42 (0.05) | 0.78 (0.02) | 0.88 (0.02) |
| Depressed | NLP + Billing Codes | 0.95 | 0.39 (0.06) | 0.78 (0.02) | 0.87 (0.02) |
| | | | | | |
| Well | Billing Codes | 0.95 | 0.06 (0.02) | 0.26 (0.27) | 0.55 (0.03) |
| Well | NLP | 0.95 | 0.37 (0.06) | 0.86 (0.02) | 0.85 (0.02) |
| Well | NLP + Billing Codes | 0.95 | 0.39 (0.07) | 0.85 (0.02) | 0.86 (0.02) |

# Research Investigator Workflow enabled by mi2b2



Query is done
To find patients

Derive new
data from images

**Use i2b2**

Request
Images with
Accession #'s

Study
Images

BIRN/XNAT

mi2b2

Images
Retrieved
from Clinical
PACS

# *White matter abnormalities associated with treatment-resistant depression*

- Scans collected as part of routine clinical care
- NLP identified cohort with treatment outcomes and lack of diagnosed brain pathology on MRI
- Diffusion tensor imaging in 150 pts with best data



*Age-related decline in white matter integrity increases with treatment resistant depression*

*Medial fornix shows strongest effect*



HC r= -0.420***
MDD F-Rem r= -0.132
MDD P-Rem r= -0.240
MDD N-Rem r= -0.595***

Hoogenboom et al. World J Biol Psychiatry, 2012

# Ontology-driven data organization allows simplistic data models that paste together

# High Throughput Methods for supporting Translational Research

- Set of patients is selected from medical record data in a high throughput fashion

- Investigators explore phenotypes of these patients using i2b2 tools and a translational team developed to work specifically with medical record data

- Distributed networks cross institutional boundaries for phenotype selection, public health, and hypothesis testing

- Tissues of these patients can be made available for genomic and biochemical analysis

# i2b2 Implementations

CTSA's

- Boston University
- Case Western Reserve University (*including Cleveland Clinic*)
- Children's National Medical Center (GWU), Washington D.C.
- Duke University
- Emory University (*including Morehouse School of Medicine and Georgia Tech* )
- Harvard University (*including Beth Israel Deaconness Medical Center, Brigham and Women's Hospital, Children's Hospital Boston, Dana Farber Cancer Center, Joslin Diabetes Center, Massachusetts General Hospital*)
- Medical University of South Carolina
- Medical College of Wisconsin
- Oregon Health & Science University
- Penn State MIlton S. Hershey Medical Center
- Tufts University
- University of Alabama at Birmingham
- University of Arkansas for Medical Sciences
- University of California Davis
- University of California, Irvine
- University of California, Los Angeles*
- University of California, San Diego*
- University of California San Francisco
- University of Chicago
- University of Cincinnati (including *Cinncinati Children's Hospital Medical Center*)
- University of Colorado Denver (including *Children's Hospital Colorado*)
- University of Florida
- University of Kansas Medical Center
- University of Kentucky Research Foundation
- University of Massachusetts Medical School, Worcester
- University of Michigan
- University of Pennsylvania (including *Children's Hospital of Philadelphia*)
- University of Pittsburgh (including their *Cancer Institute*)
- University of Rochester School of Medicine and Dentistry
- University of Texas Health Sciences Center  at Houston
- University of Texas Health Sciences Center at San Antonio
- University of Texas Medical Branch (Galveston)
- University of Texas Southwestern Medical Center at Dallas
- University of Utah
- University of Washington
- University of Wisconsin - Madison (*including Marshfield Clinic*)
- Virginia Commonwealth University
- Weill Cornell Medical College

Academic Health Centers (does not include AHCs that are part of a CTSA):

- Arizona State University
- City of Hope, Los Angeles
- Georgia Health Sciences University, Augusta
- Hartford Hospital, CN
- HealthShare Montana
- Massachusetts Veterans Epidemiology Research and Information Center (MAVERICK), Boston
- Nemours
- Phoenix Children's Hospital
- Regenstrief Institute
- Thomas Jefferson University
- University of Connecticut Health Center
- University of Missouri School of Medicine
- University of Tennessee Health Sciences Center
- Wake Forest University Baptist Medical Center

HMOs:

- Group Health Cooperative
- Kaiser Permanente

*International:*

- Georges Pompidou Hospital, Paris, France
- Hospital of the Free University of Brussels, Belgium
- Inserm U936, Rennes, France
- Institute for Data Technology and Informatics (IDI), NTNU, Norway
- Institute for Molecular Medicine Finland (FIMM)
- Karolinska Institute, Sweden
- Landspitali University Hospital, Reykjavik, Iceland
- Tokyo Medical and Dental University, Japan
- University of Bordeau Segalen, France
- University of Erlangen-Nuremberg, Germany
- University of Goettingen, Goettingen, Germany
- University of Leicester and Hospitals, England (Biomed. Res. Informatics Ctr. for Clin. Sci)
- University of Pavia, Pavia, Italy
- University of Seoul, Seoul, Korea

Companies:

- Johnson and Johnson (TransMART)
- GE Healthcare Clinical Data Services

# Aggregating across 4 hospitals, 3 i2b2 instances
## SHRINE (Shared Research Informatics Network) = Distributed Queries

# Clinical data in SHRINE

- 10 years (2001-2011)
- 4 hospitals
- 6 million total patients
- >1 billion medical observations
  - Demographics
  - Diagnoses                    (ICD9-CM)
  - Medications        (RxNorm)
  - Labs                    (LOINC)

# SHRINE

## Navigate Terms | Find Terms

- SHRINE
  - Demographics
  - Diagnoses
    - Certain conditions originating in the perinatal period
    - Complications of pregnancy, childbirth, and the puerperium
    - Congenital anomalies
    - Diseases of the blood and blood-forming organs
    - Diseases of the circulatory system
    - Diseases of the digestive system
    - Diseases of the genitourinary system
    - Diseases of the musculoskeletal system and connective tissue
    - Diseases of the nervous system and sense organs
    - Diseases of the respiratory system
    - Diseases of the skin and subcutaneous tissue
    - Endocrine, nutritional, and metabolic diseases and immunity disorders
    - Infectious and parasitic diseases
    - Injury and poisoning
    - Mental Illness
      - Adjustment disorders
      - Alcohol-related disorders
      - Anxiety disorders
      - Attention deficit, conduct, and disruptive behavior disorders
      - Delirium, dementia, and amnestic and other cognitive disorders

## Query Tool

Query Name: Pervasi-0-9 yea@00:44:10

| Group 1 | | | Group 2 | | | Group 3 | | |
|---|---|---|---|---|---|---|---|---|
| Dates | Occurs > 0x | Exclude | Dates | Occurs > 0x | Exclude | Dates | Occurs > 0x | Exclude |

**Group 1:** Pervasive developmental d

**Group 2:**
- 0-9 years old
- 10-17 years old
- 18-34 years old

**one or more of these** AND **one or more of these** AND **drag a term to here**

Autism ▼    Info    Request New Topic

Run Query    New Query    Print Query    2 Groups    New Group

## Previous Queries

- Pervasi-0-9 yea@00:44:10 [9-27-2010] [kohane]
- Perva-0-9 y-PHENO@17:57:23 [9-26-2010] [kohane]
- PDD-0-34@17:40:42 [9-26-2010] [kohane]
- Perva-Male-Schiz@16:27:52 [9-26-2010] [kohane]
- AI+PDD-0-34@17:47:17 [9-26-2010] [kohane]
- medicated-ppd-34@16:41:28 [9-26-2010] [kohane]
- =34@16:39:02 [9-26-2010] [kohane]
- Pervasi-0-9 yea@16:37:13 [9-26-2010] [kohane]

## Query Status

**Finished Query: "Pervasi-0-9 yea@00:44:10"**

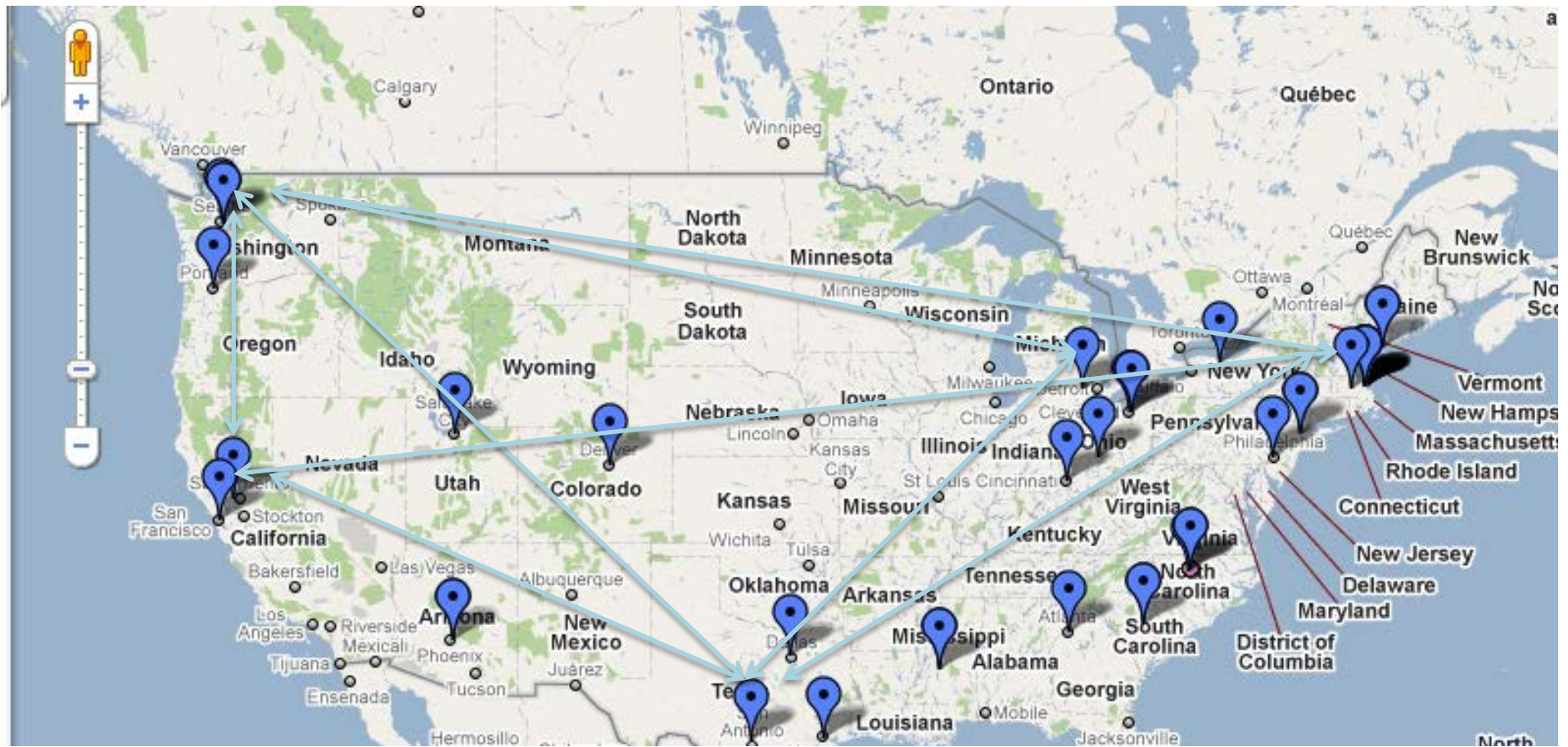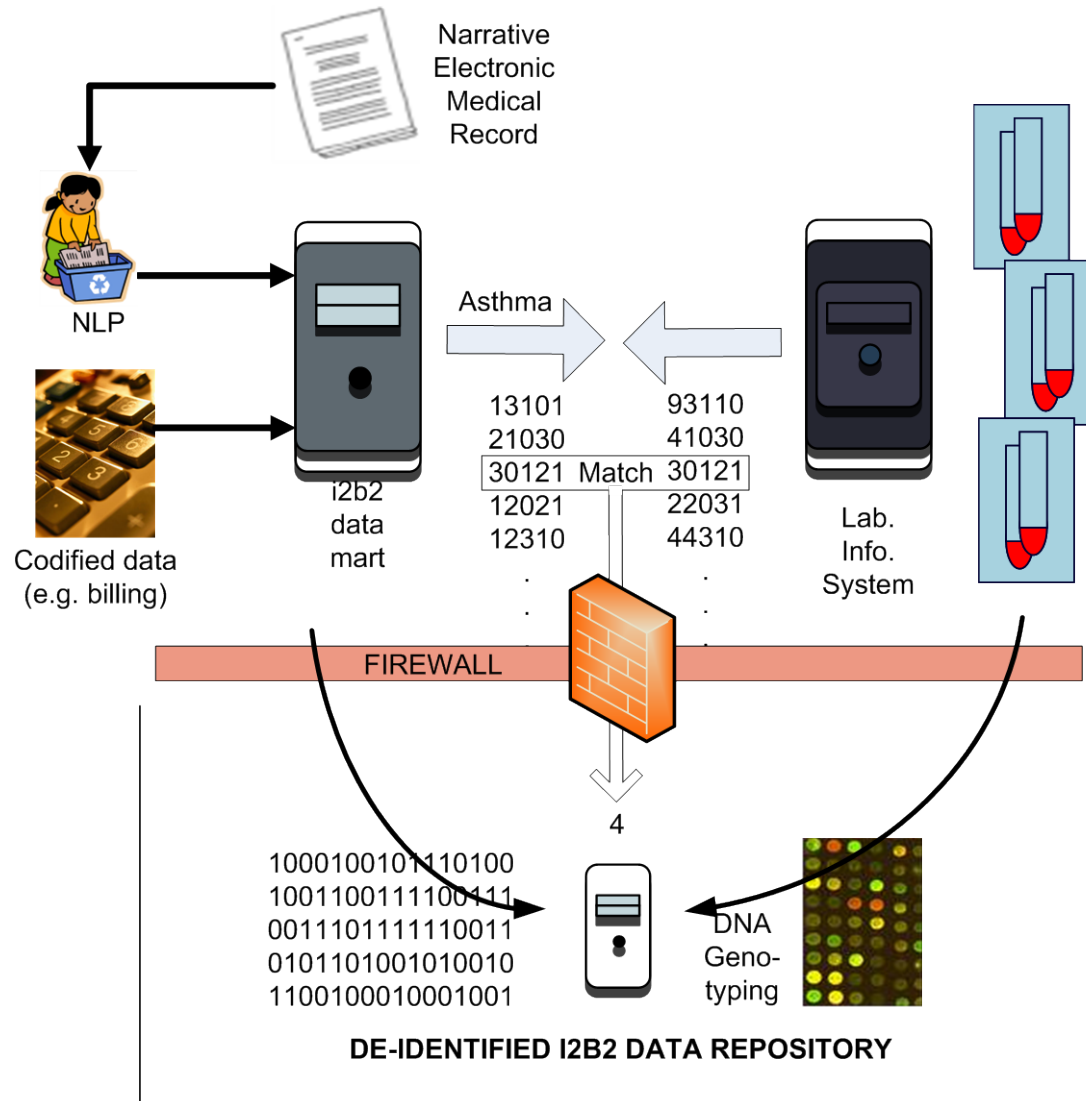| | | |
|---|---|---|
| BIDMC - 141 ±3 patients | FINISHED | [78.7 secs] |
| CHB - 9103 ±3 patients | FINISHED | [78.7 secs] |
| Partners - 5134 ±3 patients | FINISHED | [78.7 secs] |

# 2012

# Performing Clinical trials "in-silico"

- Performing an observational, phase IV study is an expensive and complex process that can be potentially modeled in a retrospective database using groups of patients available with large amounts of well organized medical data.

- Fundamental problems complicate this approach:
  - Patients drift in and out of the healthcare system.  Sophisticated statistical models using adequate control populations are necessary to compensate for the drift.
  - Confounding variables may not be found in the database.  Natural language processing may be needed to extract the confounders from textual reports to allow confounders to be exposed.
  - Unknown missing data disrupts typical statistical approaches.
  - Biases in the data can easily mislead the investigator to false conclusions; data exploration and visualization tools are needed to expose these kinds of potential problems.

# High Throughput Methods for supporting Translational Research

- Set of patients is selected from medical record data in a high throughput fashion

- Investigators explore phenotypes of these patients using i2b2 tools and a translational team developed to work specifically with medical record data

- Distributed networks cross institutional boundaries for phenotype selection, public health, and hypothesis testing

- Tissues of these patients can be made available for genomic and biochemical analysis

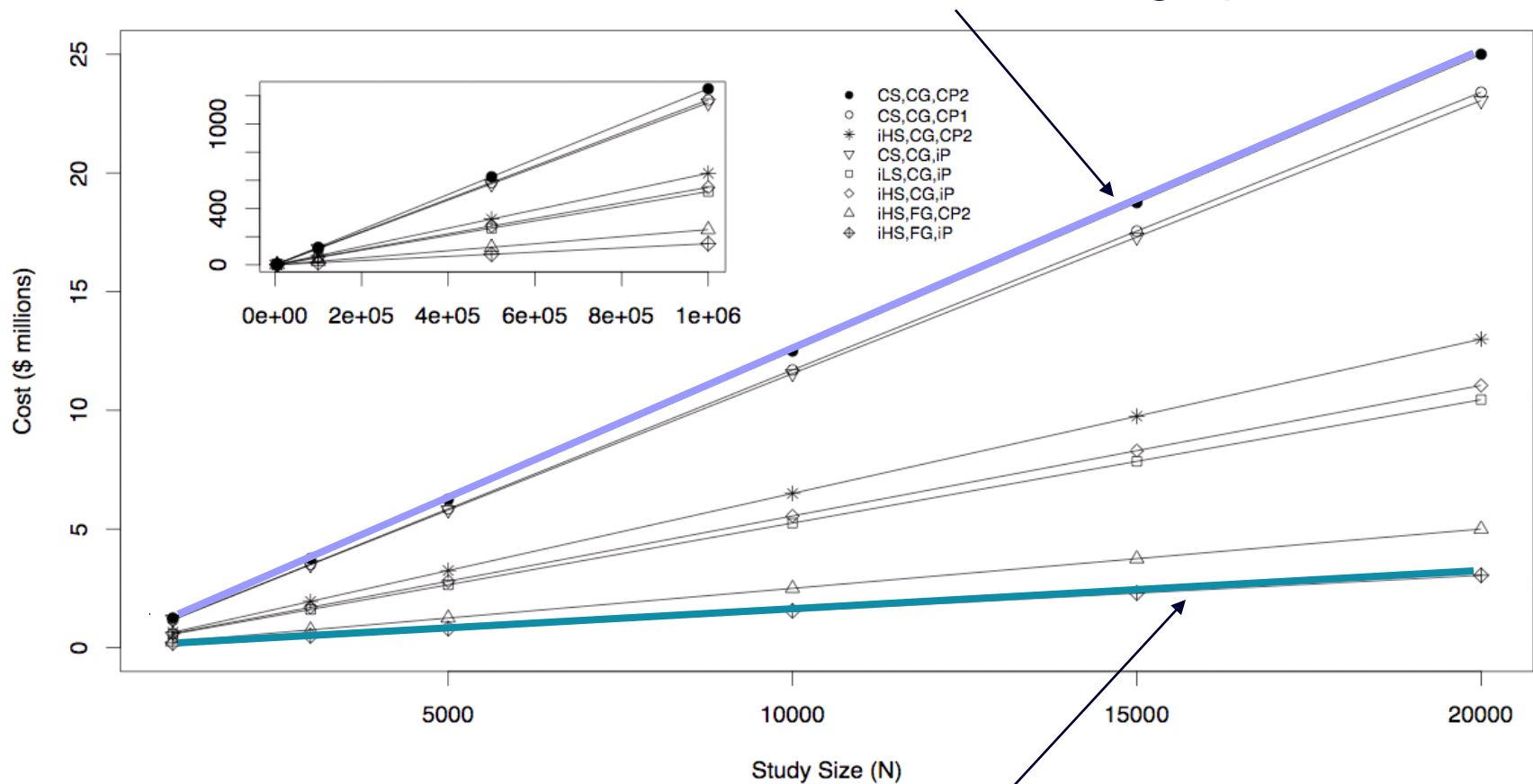# Genotype samples and compare to controls

# Cost and time benefit of Instrumenting with Sample Collection for Modest-size Study with 10,000 subjects (cases + controls)

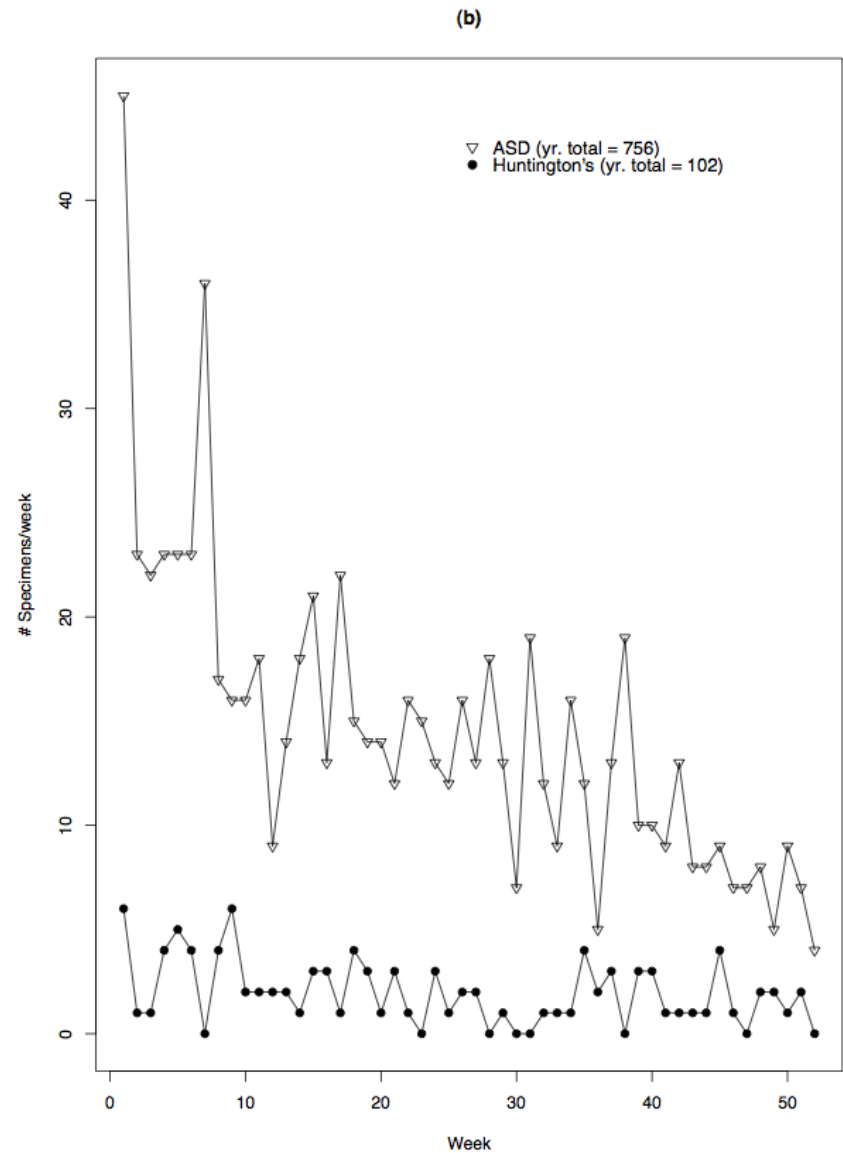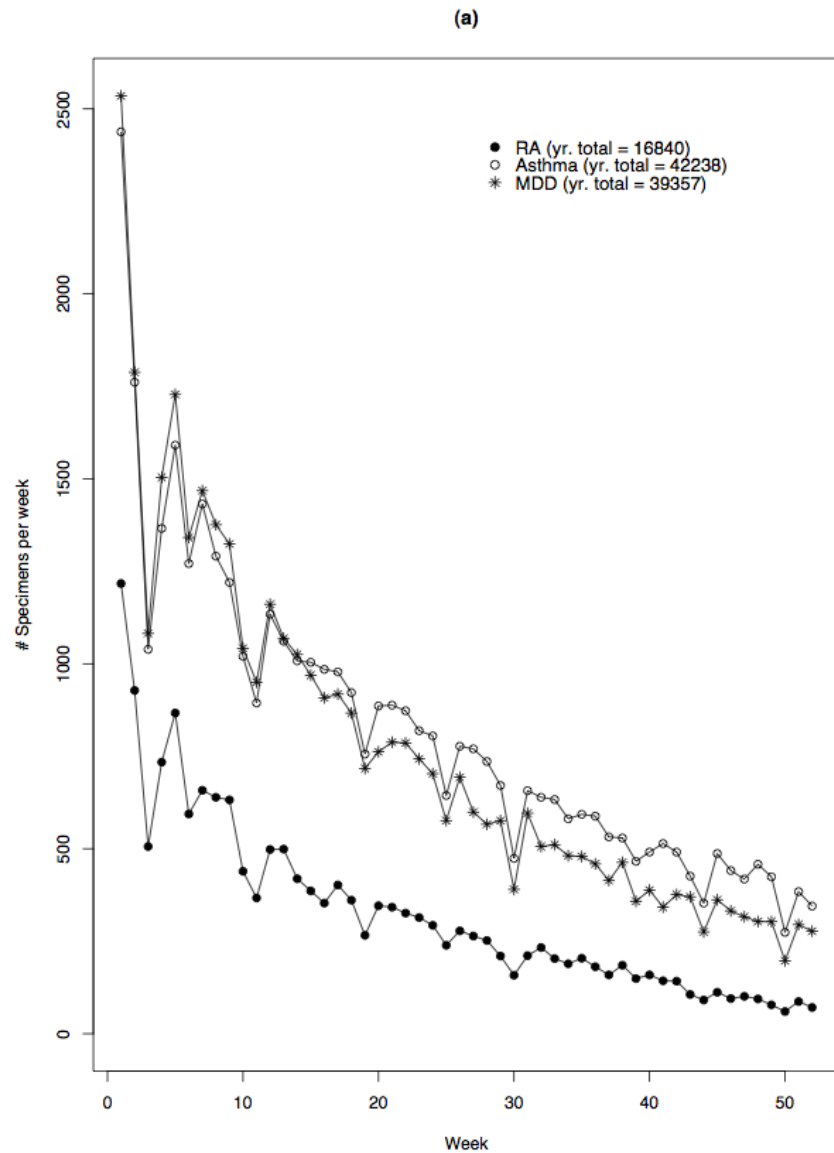| Old vs. **New** | Cost ($) | Time |
|---|---|---|
| 1 chart review per patient (CP1) | $20 | 15 minutes/subject |
| **High-throughput phenotyping (iP) through RPDR and i2b2** | **$50K Total** | **1 month total (conservative high estimate)** |
| Sample acquisition through primary care provider (CP) | $650 | 3-5 subjects/week[1] |
| **High-throughput sample acquisition through RPDR and BETR/Crimson.** | **$20** | **50-200 subjects /week[2]** |

= $6.7 million/study   vs.   $250 thousand/study

# Escalating cost and time benefit of Instrumenting with Sample Collection
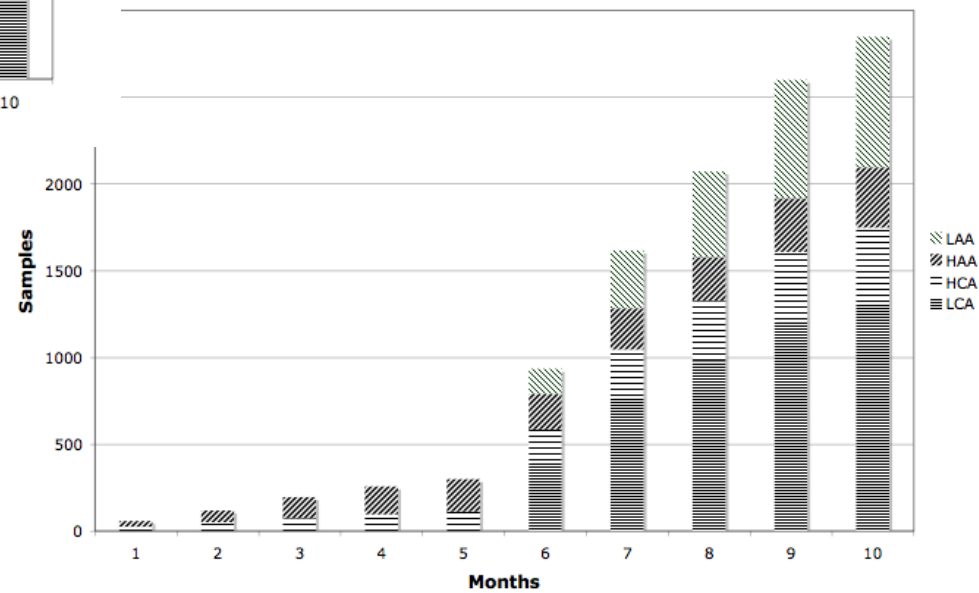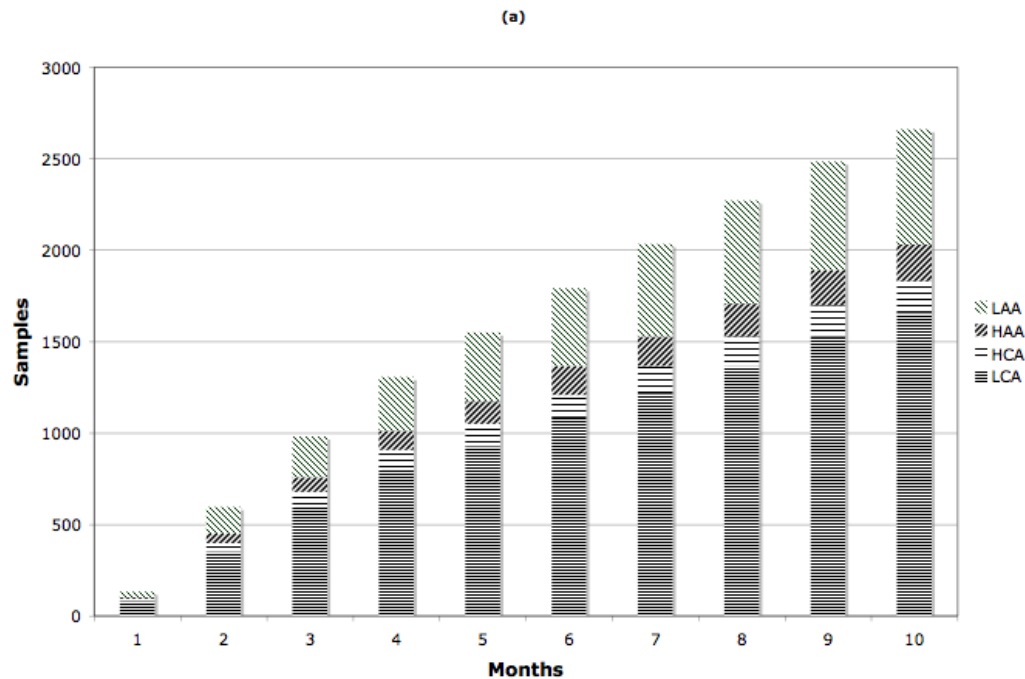
Previous model for collecting specimens



New model for collecting specimens

# Accrual Rates

# Meeting Expectations

# Seven important factors enabled by i2b2 platform

- 1) Enables enterprise-wide repurposing of health care data for research
- 2) Enables extensible software architecture for developers
- 3) Extends EHR research so that data may be shared among sites
- 4) Enables natural language processing
- 5) Provides method for materializing scientific method for EHR-based investigations
- 6) Extends EHR research so that data may be shared among sites and samples may be obtained
- 7) Provides platform for Clinical Trials "in silico"

# Collaborators

- RPDR
  - Eugene Braunwald
  - John Glaser
  - Diane Keogh
  - Henry Chueh

- I2b2
  - Isaac Kohane
  - Susanne Churchill
  - Michael Mendis
  - Nich Wattanasin
  - Vivian Gainer
  - Lori Phillips
  - Wensong Pan
  - Janice Donahue
  - William Simons (SHRINE)
  - Doug McFadden (SHRINE)
  - Christopher Herrick (mi2b2)
  - David Wang (mi2b2)
  - Bill Wang (mi2b2)

- Sample Acquisition
  - Lynn Bry
  - Natalie Boutin

- **Depression Driving Biology and Pharmacovigilance Projects:**
- Roy Perlis/Jordan Smoller/Dan Iosifescu (PIs)
  - Victor Castro
  - Caitlin Clements
  - Wouter Hoogenboom,
  - Martha Shenton
  - Patience Gallagher
  - Stefanie Block
  - Alison Hoffnagle

- **International Cohort Collection for Bipolar Disorder:**
- Jordan Smoller/Pamela Sklar (PIs)
  - Roy Perlis
  - Victor Castro
  - Alison Hoffnagle
  - Sydney Weill
  - Mireya Nadal-Vicens
  - Niels Rosenquist
  - April Hirschberg
  - Alisha Pollastri
  - Jane Erb
  - Shaun Purcell
  - Nadia Solovieff