# Using p2p systems to translate routine tissue collection and clinical encounters into robust research discovery



PSL
Pathology Specimen Locator

Andrew_McMurry@hms.harvard.edu

# Translational Research

**Routine care delivery → Robust research discovery**

<u>use cases</u>

1. Share routinely collected <span style="color:red">human tissues</span> for biomarker discovery and high-throughput validation

2. Share <span style="color:red">experimental outcomes</span> derived from tissue processing, with an emphasis on genomic measurements

## NCI Vision 2001

➤ **Millions of Paraffin Embedded Tissues**
  ◆ Biomarker Discovery / Validation:

  DNA fragments of up to 400 bp and RNA fragments of up to 150 nucleotides can be routinely isolated for mutation detection, SNP analysis, detection of translocation, and microRNA quantification. Pathology services and screens, TMA construction, …

➤ **Smaller Collections of Fresh / Frozen Tissues**
  ◆ DNA/RNA Microarrays, chip-chip, chip-seq, etc.

for Translational Research Requiring Human Specimens

# Gene Expression in Fixed Tissues and Outcome in Hepatocellular Carcinoma

*Results* The expression-profiling method for formalin-fixed, paraffin-embedded tissue was highly effective: samples from 90% of the patients yielded data of high quality, including samples that had been archived for more than 24 years. Gene-expression profiles of tumor ti

# Sharing Human Tissues for Discovery and Validation

- **Challenges**

  - How to link routine pathology databases for research?
    - *Local Control* → each hospital is a "peer" on the network

  - How to ensure patient privacy in accordance with HIPAA?
    - *Local Control* → anonymization and statistical aggregates

  - How to engender hospital participation?
    - *Local Control* → hospitals remain owners of specimens and stewards of patient data

# How It Works

1. **Link existing databases**
   - Extract from existing hospital systems
   - Transform the data into common HIPAA-safe vocabulary
   - Load into locally controlled "SPIN peer" with deidentified ID

2. **Protect Patient Privacy per HIPAA**
   - <u>Anonymized</u>:        Case Counts / Aggregates

   - <u>Limited</u>:        When authorized for individual cases

   - *PHI is rarely used, and only with IRB from each hospital.*

3. **Hospital Control**
   - No central governing body or server
   - Peers (hospital) remains in control over disclosures at all times

# (1) Linking routine care systems



- <u>Extract</u> from routine care delivery systems
  - ➤ Databases *or* XML

- <u>Transform</u> free text reports
  - ➤ "Scrub" patient identifiers (per HIPAA)
  - ➤ NLP (autocode) into controlled vocabularies such as UMLS

- <u>Load</u> into the hospital controlled PEER database
  - ➤ Assign a randomly generated ID to each case

# Transforming Free Text: "Scrubber"

# BMC Medical Informatics and Decision Making

BioMed Central

Software

**Open Access**

# Development and evaluation of an open source software tool for deidentification of pathology reports

Bruce A Beckwith[*1,2], Rajeshwarri Mahaadevan[2], Ulysses J Balis[2,3] and Frank Kuo[2,4]

Address: [1]Department of Pathology, Beth Israel Deaconess Medical Center, 330 Brookline Ave., Boston, MA, USA, [2]Department of Pathology, Harvard Medical School, 25 Shattuck Street, Boston, MA, USA, [3]Department of Pathology, Massachusetts General Hospital, 55 Fruit Street, Boston, MA, USA and [4]Department of Pathology, Brigham & Women's Hospital, 75 Francis Street Boston, MA, USA

Email: Bruce A Beckwith[*] - bruce_beckwith@bidmc.harvard.edu; Rajeshwarri Mahaadevan - rajeshwarri@yahoo.com; Ulysses J Balis - balis@helix.mgh.harvard.edu; Frank Kuo - fkuo@partners.org

[*] Corresponding author

**http://spin.chip.org/software.html**

# (2) Protecting Patient Privacy

Increasing levels of investigator access commensurate with authorization by the hospital & investigator demonstrated need.

| Use Case | *Statistical Queries > 90%* | *Non Identifying < 10%* | *PHI < 1%* |
|---|---|---|---|
| Tissue Sharing | Feasibility Studies | Case Selection | Clinical Data |
| Public Health | Automated Analysis | Investigation | Emergencies |
| Genomic Studies | Significant Markers | Case Selection | Genotypes |

# Feasibility Study:
## ascertain if there are enough samples available

# Case Selection and Retrieval

## My Specimen Requests

| Status | Peer(s) | Project | Tracking # |
|---|---|---|---|
| OPEN | MGH, BIDMC | Multi tumor TMA | 7c7130c0-e7ed-4b39-bde3-01593b7e5287 |
| NEEDS IRB | MGH, BIDMC | Microarray, gene expression liver | d2481460-2021-4ce2-b27a-9a3481fde11a |
| RETREIVING BLOCKS | MGH, BIDMC | Microarray, gene expression lung | 2e051b9b-e65c-414e-a54f-27f83ef9a8dd |
| COMPLETED | MGH, BIDMC | Microarray, gene expression white blood cells | c96a220c-91f7-4a6d-a5fe-89cce0f7781f |

## Participating Peer Institutions

[ BIDMC Rules and Pricing ] xxxxx Paraffin Specimens Available

[ CHB Rules and Pricing ] xxxxxx Paraffin Specimens Available

# (3) Hospitals remain in control

- Each hospital (Peer) chooses **who** to share with



- And **what** to share (Path Reports, ED Visits, ... )

# Sites Participating in National Demonstration

1. **Brigham & Women's Hospital***
2. **Beth Israel Deaconess Medical Center***
3. Cedars-Sinai Medical Center
4. **Dana-Farber Cancer Institute***
5. **Children's Hospital Boston***
6. **Harvard Medical School***
7. **Massachusetts General Hospital***
8. National Institutes of Health
9. National Cancer Institute
10. Olive View Medical Center
11. Regenstrief Institute
12. University of California at Los Angeles Medical Center
13. University of Pittsburgh Medical Center
14. VA Greater LA Healthcare System

**\* Participate in ongoing "Virtual Specimen Locator" collaboration**

# Sites Participating in National Demonstration

# Overview



1. Extract pathology data from existing systems into a locally controlled *peer node database* and identifiers codebook

Identifiers Codebook

Existing Systems

Peer

supernode

Peer

Peer

Increasing Levels of Investigator Access

2. Feasibility studies — Statistical/Public Level Detail

Query Interface

3. Select cases — DeIdentified diagnosis

4. Request Specimens — Authorized Re-Link

Membership verification and Role Based Access

# Sharing Human Tissues for Discovery and Validation

- ## Results
  - ➤ National prototype including HMS, UCLA, Indiana, UPMC, …
  - ➤ Live Production instance at HMS including 4 hospitals
  - ➤ Developed Open Source Tools
  - ➤ caBIG adopted caTIES from SPIN
  - ➤ Influenced Markle's Common Framework federated query
  - ➤ TMA construction using specimens from four sites

## SPIN: Sharing Human Tissues for Discovery and Validation

# Harvard hopes database will speed cancer cures

The Boston Globe

By Liz Kowalczyk, Globe Staff | November 21, 2005

Since World War II, many cancer patients who have had surgery at a Harvard-affiliated teaching hospitals have left a small piece of their tumor to science.

These clumps of human cells have been frozen in liquid nitrogen or preserved in paraffin blocks the size of small Post-it notes -- and they now fill giant freezers and floor-to-ceiling shelves in hospital basements and off-site warehouses.

The value of this tissue trove has soared in recent years with the successful cataloging of humans genes. Researchers need to study hundreds of specimens to find genetic mutations, proteins, and other molecules linked to cancer, in hopes of developing new medicines and tests to diagnose cancer early and help customize treatment for individual patients.

# Sharing Human Tissues for Discovery and Validation

**Editorial** Comments

JAMIA

*Editorial* ■

# Lessons Learned from the Shared Pathology Informatics Network (SPIN): A Scalable Network for Translational Research and Public Health

Michael J. Becich, MD, PhD

# Sharing Genomic Results for Association Studies?

- ## Motivation:
  - ➤ Enable Phenotype – Genotype association studies for Autism Spectrum Disorders
  - ➤ Integrative genomics across multiple measurement modalities such as DNA->RNA (EQTL)

- ## New Challenges:
  - ➤ Privacy Policy: genotypes are clearly identifiable
  - ➤ Resources: storage, processing, network load for SNP data
  - ➤ Multiple Testing and False Discovery

# Sharing Genomic Results for Association Studies?

- **Policy Challenges**

## GENETICS

# No Longer De-Identified

Amy L. McGuire[1]* and Richard A. Gibbs[2]

As DNA sequencing becomes more afford able and less time-consuming, scientists are adding DNA banking and analysis to research protocols, resulting in new disease-specific DNA databases. A major ethical and policy question will be whether and how much information about a particular individual's DNA sequence ought to be publicly accessible.

Without privacy protection, public trust will be compromised, and the scientific and medical potential of the technology will not be realized.

## Sharing Genomic Results for Association Studies?

- **Technical Challenges**



**http://en.wikipedia.org/wiki/Hard_drive**

## Sharing Genomic Results for Association Studies?

- ## Multiple Testing Challenges

CLINICIAN'S CORNER
## The Incidentalome

### A Threat to Genomic Medicine

Isaac S. Kohane, MD, PhD; Daniel R. Masys, MD; Russ B. Altman, MD, PhD

*JAMA.* 2006;296:212-215.

"In the genomic era, **the lack of prior probabilities** regarding the clinical import of each genetic variant **creates the likelihood** of a large proportion of **false positives**, if genetic testing is not placed on a systematic quantitative basis."

## Sharing Genomic Results for Association Studies?

- **Solution: what worked before?**

  ➢ Link genomic test results to the clinical data in a spin peer

  ➢ Protect patient privacy with anonymization and statistical aggregation techniques

  ➢ Engender participation by reasserting local ownership of microarray data and stewardship of patient privacy

# Sharing Genomic Results for Association Studies?

webplink

| Home | My Studies | Create Study | Search Studies | Tutorial | Contact | Create Account | Login | Logout |

**Analyze Data Set**

Analysis to Perform*:  Quantitative  ?

Model*:  Logistic  ?

Phenotype*:  Language Function  ?

Linear/Logistic: The basic
association test for a trait based on
comparing allele frequencies
between phenotypes applying a
linear or logistic regression tests.
Genotypic: Generates two extra tests
per SNP, the dominance deviation
component from the additive model
or a 2 df joint test of both additive
and dominance.

**Analysis Options**

Phenotype Column Number:  1

Missingness Per Marker (less than):  0.1  ?

Missingness Per Individual (less than):  0.1  ?

Minor Allele Frequency (greater than):  0.01  ?

Hardy Weinberg Equalibrium:  ?

Dominant: To specify a model
assuming full dominance for the
minor allele
Recessive: To specify a model
assuming full recessive for the minor
allele

*field is required

Submit

webplink

View Cohort Distribution Graph     View -Log 10 P-Value Graph     View Correlation Call Rate Graph     View Plink Analysis Log File

New P-Value: [　　　　] (Submit)     Format: CSV ▼ (Export)



**Distribution of Phenotype**

**[7]Language Delay : logistic recessive,geno<.1,mind<.1,maf<.01,hwe<.05**

New P-Value: [          ]  (Submit)    Format: [ CSV ⇕ ] (Export)



### Rate of snp-phenotype correlations having pValue < 0.05

chr1  chr2  chr3  chr4  chr5  chr6  chr7  chr8  chr9  chr10  chr11  chr12  chr13  chr14  chr15  chr16
chr17  chr18  chr19  chr20  chr21  chr22  chr23  chr25

# Applying lessons learned:
# Common Architecture

**JAMIA** Editorial Comments

*Editorial* ■

## Lessons Learned from the Shared Pathology Informatics Network (SPIN): A Scalable Network for Translational Research and Public Health

MICHAEL J. BECICH, MD, PHD

# Lessons Learned: IRBs and political will

➢ Statistical level queries easy are OK by IRBs

➢ Difficulty arises going to the next step
  ◆ HIPAA limited data set
  ◆ PHI

➢ ANY use of patient data for research imposes SOME risk

➢ Minimize risk, show that research benefit is overwhelmingly in the best interest of patients

# Lessons Learned: mapping heterogeneous DBs



**VS**

***Start SMALL : Grow the number of common terms!***

# Lessons Learned: <u>mapping heterogeneous DBs</u>

1. Request for Capabilities & Statistics (What is available?)

2. Availability limits scope of the vocabulary

3. Which BIG questions can be asked with only a few identifiers?
   - Pathology:        age, gender, collection,   free text "diagnosis"
   - Public Health:    age, gender, location,    free text "complaint"
   - CTSA:             age, gender, …………,     free text mining

4. Parallel tracks:  autocoding and standard vocabulary approach
   - Different low hanging fruit: diagnosis *vs* MRN

5. Quick End-To-End lifecyles
   - Question, development, research, new question

# Summary

## Addressed 3 pervasive issues:

➤ Linking routine care systems for robust research

➤ Protecting patient privacy

➤ Engendering participation among hospitals

## Use Cases

➤ Routinely collected human tissues for biomarker discovery and high-throughput validation

➤ Genomic measurements derived from tissue sharing

# Collaborators & Acknowledgements

- **Biospecimen Sharing Community**
  - Too many to list!
  - http://spin.chip.org/community.html

- **Public Health Surveillance**
  - http://chip.org/ihl

- **ASD Genotype Phenotype Associations**
  - Developers: Mike Banos , Gregory Polumbo
  - Investigators: Alexa McCray , Dennis Wall, Amanda Sedgewick
  - Collaborator: Shaun Purcell (plink author)

- **Special Thanks**
  - Advisors: Zak Kohane & Ken Mandl
  - Investigators: Kamila Naxerova & Alal Eran