

MINING THE BIOMEDICAL LITERATURE FOR DETAILED INFORMATION ABOUT BIOSPECIMENS

Paul Fearn, MBA

PhD Candidate

University of Washington

Biomedical and Health Informatics

BRN Symposium, February 22-23, 2012

Support from NLM training (NIH NLM #T15 LM07442) and ITHS (NIH NCRR RR 025014) grants

PROBLEMS AND MOTIVATION

- Problems
 - Biospecimen information may be “lost in translation” from patient care to results of experiments in biomedical literature
 - Post-hoc statistical correction for batch effects?
 - Loss of statistical power, as well as accuracy and reproducibility of biospecimen experiments
 - High costs of data management ~ Herbek G, Grizzle W @ BRN 2012
 - Motivations
 - Improve accuracy and reproducibility of research using biospecimens
 - Accelerate turnaround time and quality of biomedical research
 - Reduce costs (and/or increase value) of information abstraction
-

INFORMATION FRAMEWORK FOR BIOSPECIMEN RESEARCH



Biomedical Literature (i.e. PubMed)

Public Databases (i.e. GenBank, PDB)

Biological and Biomedical Investigations

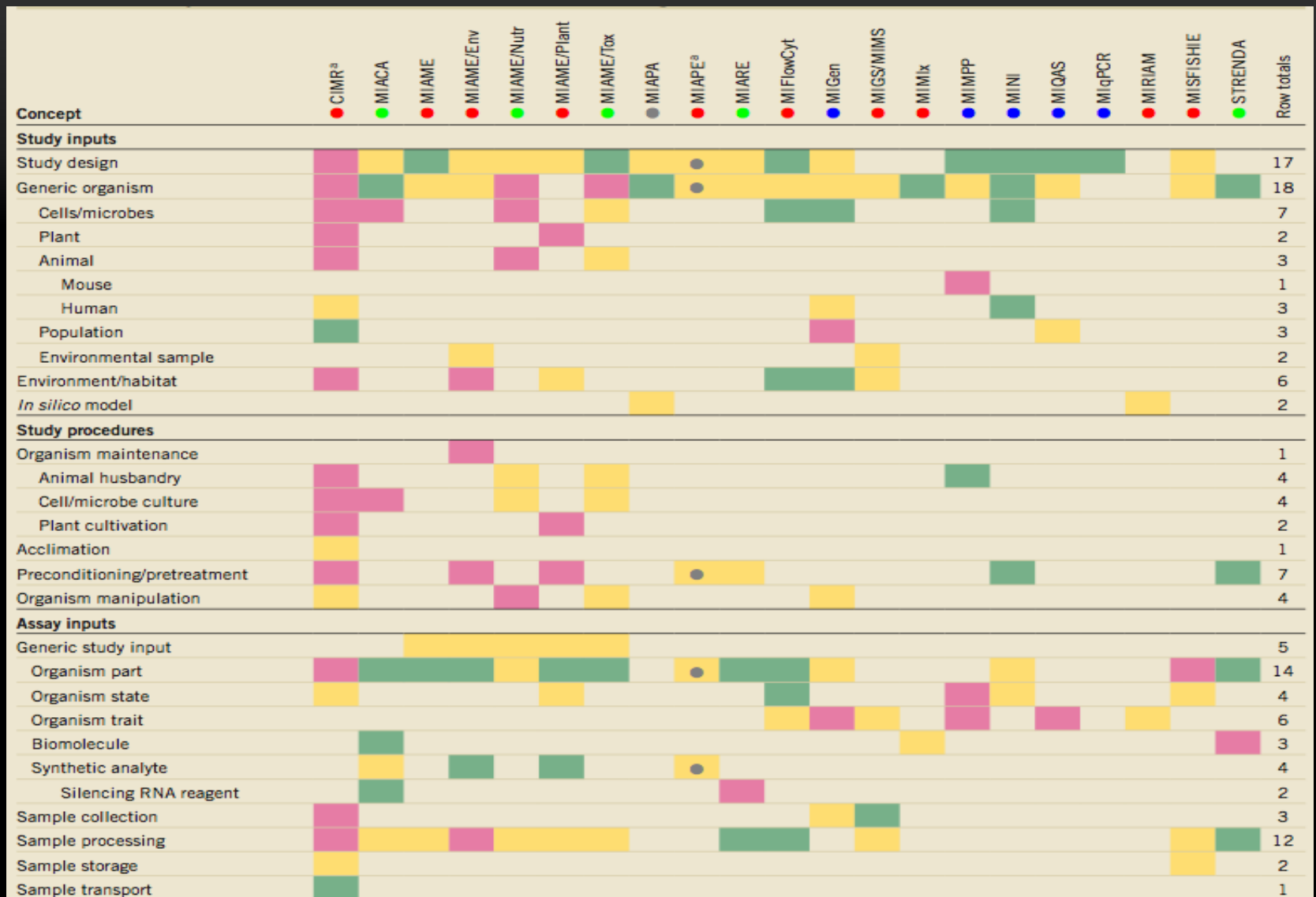
Biorepositories / Biospecimens

Patients / Sources (Donor, Family, Environment)

RELATED WORK IN GUIDELINES AND STANDARDS

- Reporting checklists for biomedical and biological experiments
 - MIAME, MISFISHIE,... MIBBI
 - Biospecimen Reporting for Improved Study Quality (BRISQ)
- Controlled vocabularies, Terminologies, Ontologies
 - OBI, UMLS,...BioPortal
- Database or object models (e.g. FuGE, BioSample)
- Data formats (e.g. ISA-TAB, SPREC)
- Have these information guidelines and standards brought about measurable effects on information reported in literature?
- ~~"Have to show the data to change behavior!"~~ ~ Herbek, G

UNEVEN COVERAGE IN MIBBI CHECKLISTS



Taylor et al. (2008). Promoting coherent minimum reporting guidelines for biological and biomedical investigations: the MIBBI project. *Nature Biotechnology*, 26(8), 889–896. doi:10.1038/nbt.1411

ZOOMING IN ON BIOSPECIMEN LIFECYCLE AND PRE-ANALYTIC VARIATION

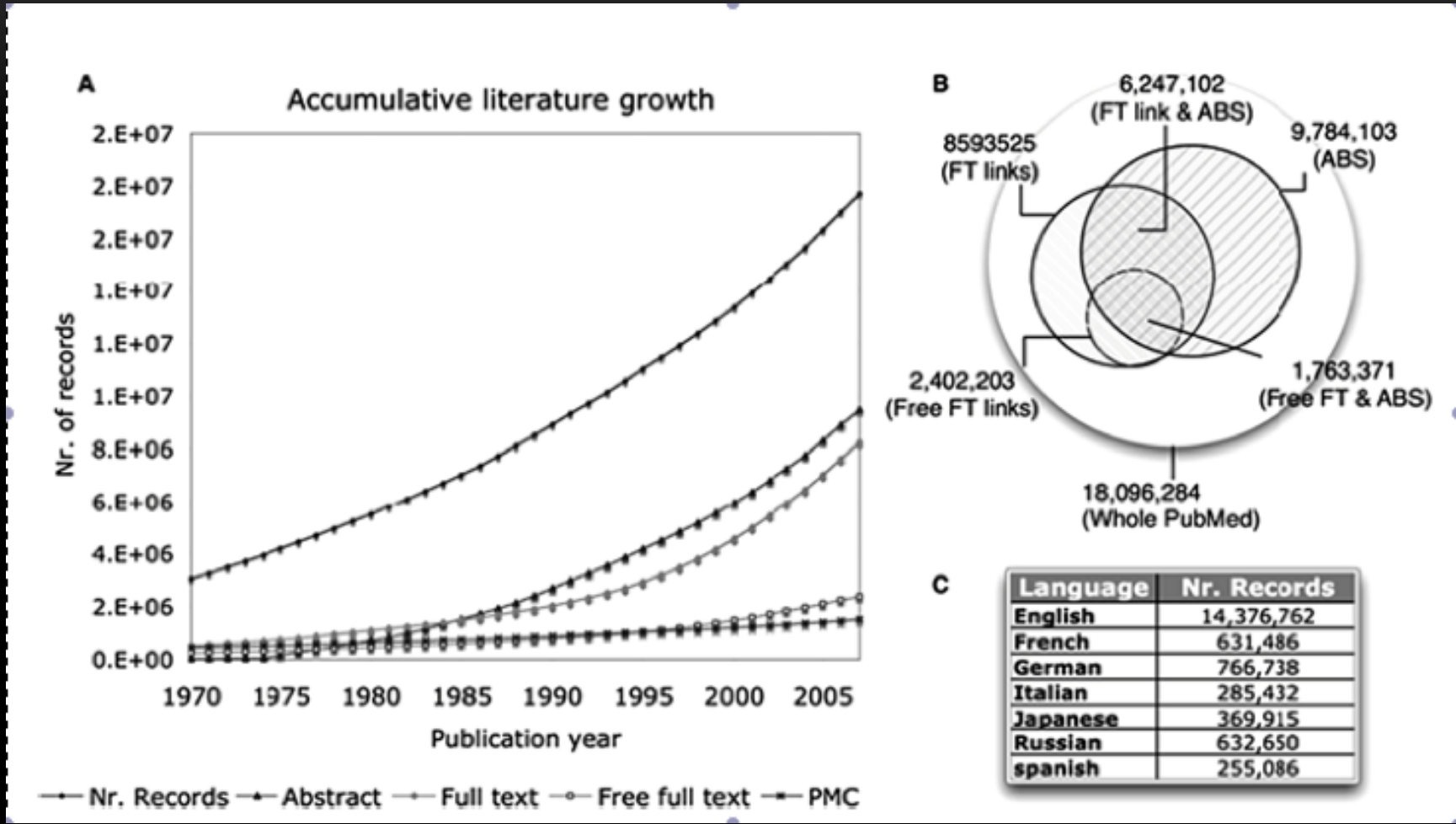


FIGURE 1. The lifecycle of the biospecimen is illustrated. The preanalytical phase of the lifecycle of the biospecimen includes each stage from patient to distribution. Preanalytical variables are addressed in the BRISQ list.

BIOSPECIMEN REPORTING FOR IMPROVED STUDY QUALITY

Data Elements	Examples
<input type="checkbox"/> Biospecimen type <i>Solid tissue, whole blood, or another product derived from a human being</i>	Serum, Urine
<input type="checkbox"/> Anatomical site <i>Organ of origin or site of blood draw</i>	Liver, Antecubital area of the arm
<input type="checkbox"/> Disease status of patients <i>Controls or individuals with the disease of interest</i>	Diabetic, Healthy control
<input type="checkbox"/> Clinical characteristics of patients <i>Available medical information known or believed to be pertinent to the condition of the biospecimens</i>	Pre-menopausal breast cancer patients
<input type="checkbox"/> Vital State of patients <i>Alive or deceased patient when biospecimens were obtained</i>	Postmortem
<input type="checkbox"/> Clinical diagnosis of patients <i>Patient clinical diagnoses (determined by medical history, physical examination, and analyses of the biospecimen) pertinent to the study</i>	Breast cancer
<input type="checkbox"/> Pathology diagnosis <i>Patient pathology diagnoses (determined by macro and/or microscopic evaluation of the biospecimen at the time of diagnosis and/or prior to research use) pertinent to the study</i>	Her2-negative intraductal carcinoma
<input type="checkbox"/> Collection mechanism <i>How the biospecimens were obtained</i>	Fine needle aspiration, Pre-operative blood draw
<input type="checkbox"/> Type of stabilization <i>The initial process by which biospecimens were stabilized during collection</i>	Heparin, On ice
<input type="checkbox"/> Type of long-term preservation <i>The process by which the biospecimens were sustained after collection</i>	Formalin fixation, freezing
<input type="checkbox"/> Constitution of preservative <i>The make-up of any formulation used to maintain the biospecimens in a non-reactive state</i>	10% neutral-buffered formalin, 10 USP Heparin Units/mL
<input type="checkbox"/> Storage temperature <i>The temperature or range thereof at which the biospecimens were kept until distribution/analysis.</i>	-80 °C, 20 to 25 °C
<input type="checkbox"/> Storage duration <i>The time or range thereof between biospecimen acquisition and distribution or analysis.</i>	8 days, 5 to 7 years
<input type="checkbox"/> Shipping temperature <i>The temperature or range thereof at which biospecimens were kept during shipment or relocation.</i>	-170 °C to -190 °C
<input type="checkbox"/> Composition assessment & selection <i>Parameters used to choose biospecimens for the study</i>	Minimum 80% tumor nuclei & maximum 50% necrosis

BUT, RAPID GROWTH IN LITERATURE...



LOTS OF PLACES WHERE SLIPS OR ERRORS CAN OCCUR...



Biomedical Literature (i.e. PubMed)

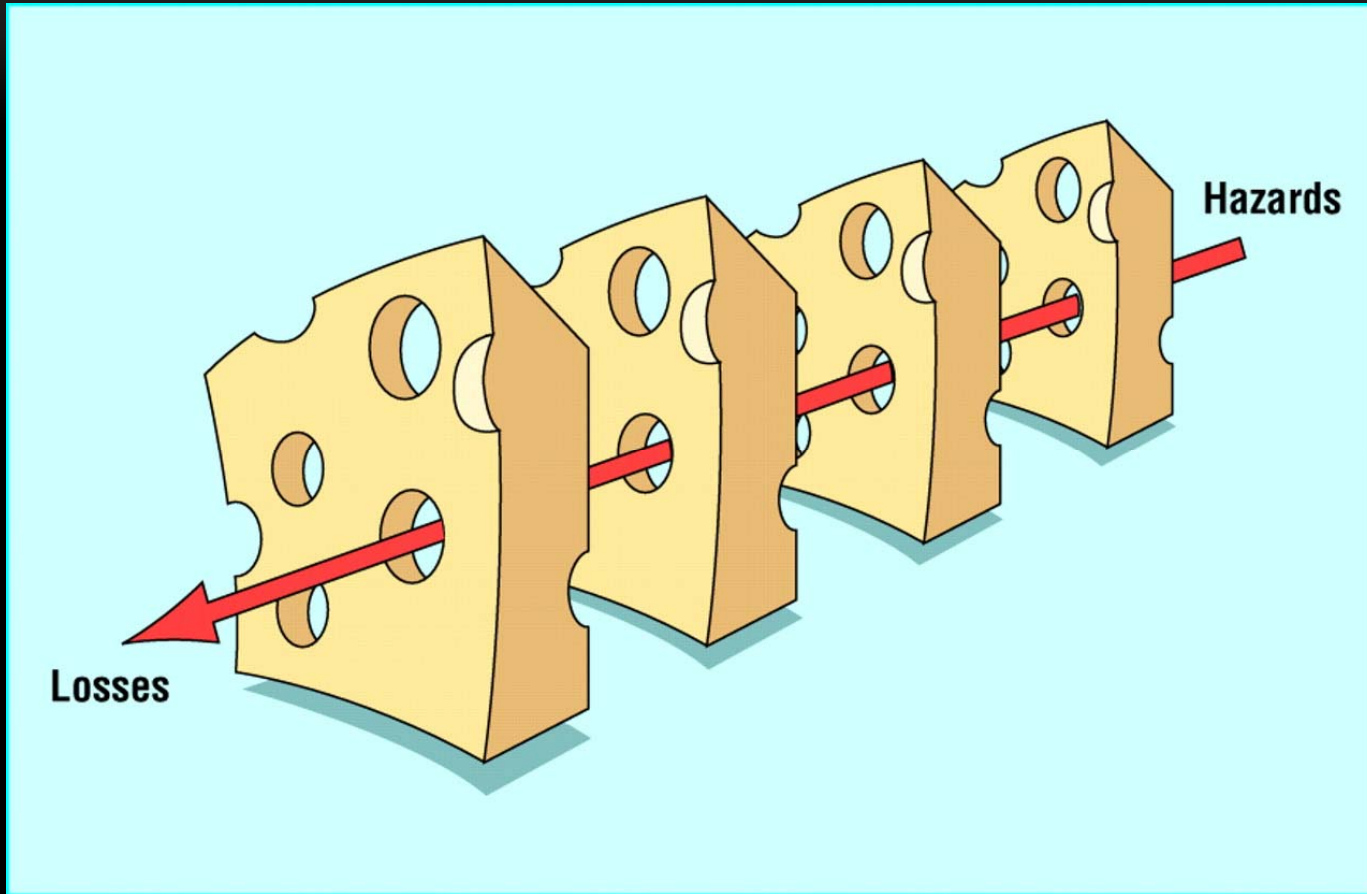
Public Databases (i.e. GenBank, PDB)

Biological and Biomedical Investigations

Biorepositories / Biospecimens

Patients / Sources (Donor, Family, Environment)

JAMES REASON'S SWISS CHEESE MODEL



Clinic/OR
Pathology
Repository
Laboratory

HIERARCHY OF INTERVENTIONS TO CHANGE PRACTICE, FROM STRONGEST (1) TO WEAKEST (4)

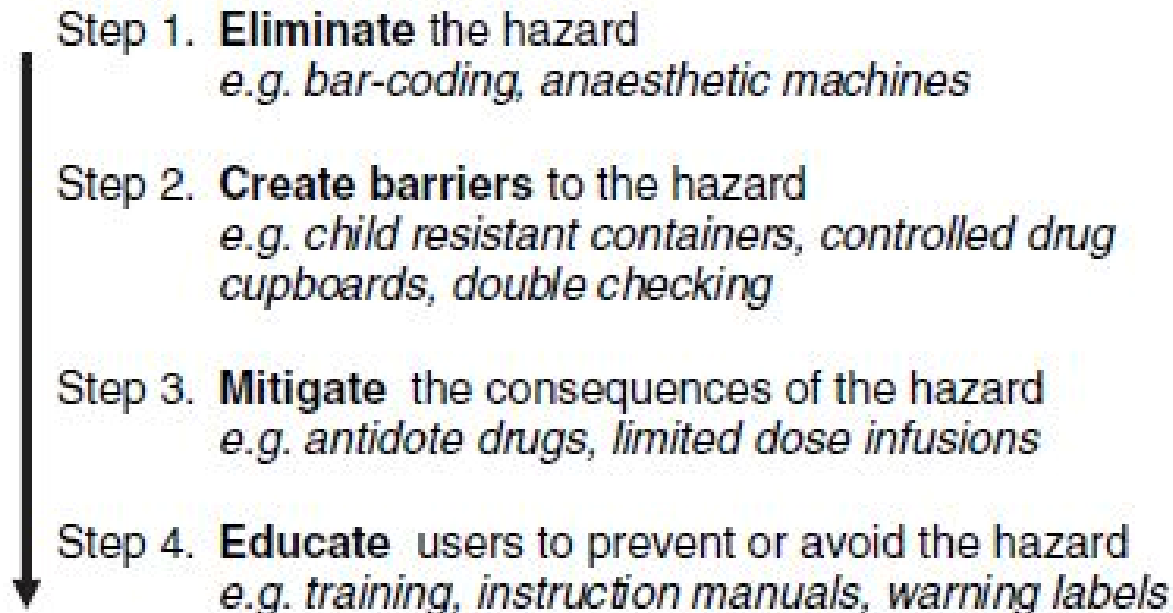


Figure 2

Hierarchy of interventions to improve safety (adapted from Hale & Glendon 1987).

MORE "GARBAGE IN / GARBAGE OUT " EVERY DAY...

- "Have to show the data to change behavior!" ~ Herbek, G



BIOSPECIMEN RESEARCH DATABASE

The screenshot shows the website header with the National Cancer Institute logo and name, and the U.S. National Institutes of Health logo and website address. Below this is the OBBR logo and name. There are navigation buttons for 'Launch NCI Best Practices' and 'Launch caHUB'. A search bar is present with a 'Search' button. A navigation menu includes 'About OBBR', 'About NCI Best Practices', 'Biospecimen Research Network', 'caHUB', 'News and Events', and 'Resources'. The main heading is 'Biospecimen Research Database'. On the left, there is a sidebar with links: 'Home', 'Search', 'Quick Search', 'Simple Search', 'Advanced Search', 'Experimental Factor Search', and 'Curator Login'. The main content area shows 'Search Results' with '18 Study(s) Found' and 'Page 1 of 1'. A 'Modify Search' button is visible. The first search result is for a study by Ahrens Kim, Braylan Raul, Almasri Nidal, Foss Robin, and Rimsza Lisa. The abstract text is: 'IgH PCR of zinc formalin-fixed, paraffin-embedded non-lymphomatous gastric samples produces artifactual "clonal" bands not observed in paired tissues unexposed to zinc formalin'. The journal is 'J Mol Diagn', 2002, Vol. 4, Page 159. There is a 'PubMed' icon and '1 Study(s) Found' below it. A dashed box highlights the following text: 'Specimen:Tissue /Stomach /Formalin /Other diagnoses / Platforms: DNA - DNA Sequencing / Sequencing of the nonreproducible bands using primers spanning the IgH joining region confirmed they contained rearranged IgH sequence, indicating the amplicons are not due to nonspecific primer binding.' Below this is another search result by Williams C, Pontén F, Moberg C, Söderkvist P, Uhlén M, Pontén J, Sitbon G, and Lundeberg J. The abstract text is: 'A high frequency of sequence alterations is due to formalin fixation of archival specimens.' The journal is 'Am J Pathol', 1999, Vol. 155, Page 1467.

1150 articles and growing ~ Bass, P

Features / annotations abstracted

Add level of evidence (# of patients, # of samples)?

Find related articles for animal biospecimens?

CAN WE USE TEXT MINING TO ASSESS REPORTING, AND CREATE BARRIERS TO INADEQUATE REPORTING?

- Document classification
 - Classify existing articles according analyte, and technology platform and other BRD labels / categories
 - Classify new manuscripts according to BRD labels
 - Information extraction
 - Named entity recognition of BRISQ variables
 - Describe relations between entities (more complex features)
 - Complementary to synoptic reporting (Herbeck, G) and detailed prospective data capture (e.g. Moore, H with OpenClinica)
-

LABELS / ANNOTATIONS ABSTRACTED FROM ARTICLES AND ENTERED INTO BRD...

Label Field	Label Value
Biospecimen Type	Tissue, Blood... [e.g. dash et al 12414521]
Diagnosis	Neoplastic - Carcinoma
Biospecimen Location	Prostate
Preservation Type	OCT
Analyte	RNA
Technology Platform	DNA Microarray
Experimental Factor	Biospecimen Acquisition, Cold Ischemia Time
Experimental Factor	DNA Microarray Specific, Targeted nucleic acid
Summary of Findings	Using microarray, 0.6% of genes (61 genes, 41 of which were named) were upregulated in prostate tissues stored at room temperature for 1 h or longer compared to 0 h controls. Genes displaying elevated expression included several early response genes (early growth response 1, EGR-1; jun B proto oncogene, junB; jun D proto oncogene, jun D; activating transcription factor 3, ATF3). The degree of upregulation was variable and ranged between minute increases to nearly 2-fold upregulation (EGR-1, junB).

CAPTURE ANNOTATION AT THE SOURCE

biobanks defined the **time of excision** as the **time of removal** of the biospecimen from the operating table as documented in the operating room by a nurse. The **time of cryopreservation** was defined as the time that the biospecimen was **stored at -70°C or in a liquid nitrogen vapor freezer**.

Queries of the databases at MBTB and the BCCA-TTR were done by each biobank's informatics staff. These queries returned the excision time and cryopreservation time for each breast tumor biospecimen. Biospecimens for which either of these times were unknown were excluded from these analyses. Data were sorted into four groups: **<30 min, 31 to 60 min, 61 to 120 min, and >120 min between excision and cryopreservation**. These data were graphed using Microsoft Excel (Microsoft Corpo-

PREVIOUS WORK IN TEXT MINING

Pathology reports – caTIES and CAP protocols

Extraction of experiment results

- Protein-protein interactions
- Genome annotation
- Relation extractions (facts or events)
 - Phenotypes
 - Species
 - Pathology

PubMed abstracts, PMC full-text, and supplementary online materials

Krallinger M, Valencia A, Hirschman L. Linking genes to literature: text mining, information extraction, and retrieval applications for biology. *Genome Biol* 2008;9 Suppl 2:S8.

Haeussler M, Gerner M, Bergman CM. Annotating genes and genomes with DNA sequences extracted from biomedical articles. *Bioinformatics* 2011 Apr.;27(7):980–986.

TEXT MINING FIRST STEPS FOR THIS DOMAIN

Word	Base Form	Part-of-Speech	Chunk	Named Entity
HAX-1	HAX-1	NN	B-NP	B-protein
associates	associate	VBZ	B-VP	O
with	with	IN	B-PP	O
cortactin	cortactin	NN	B-NP	B-protein
in	in	IN	B-PP	O
the	the	DT	B-NP	O
apical	apical	JJ	I-NP	O
membrane	membrane	NN	I-NP	O
of	of	IN	B-PP	O
hepatocytes	hepatocyte	NNS	B-NP	B-cell_type
.	.	.	O	O
Word	Morphology	Grammar	Syntax	Semantics

Figure 2

Main natural language processing levels, from word tokenization to semantics. The different processing layers for a given example sentence are shown here. This example is based on the output generated by the GENIA tagger: DT, determiner; IN, preposition or subordinating conjunction; JJ, adjective; NN, Noun (singular or mass); NNS, Noun (plural); VBZ, Verb (third person singular present). The B/I/O terminology refers to begin phrase (B), internal to phrase (I), and outside of phrase (O).

EXPLORING THE VALUE PROPOSITION

- Market based research on text mining for biospecimen information
 - Journal publishers and editors
 - Investigators using human biospecimens
 - Funding of research with biospecimens
- More effective methods to improve practices
 - Education and training is a great first step, but may not sufficient to change individual practices
 - Accreditation (e.g. CAP) based on evaluation is key
 - Structured data capture is expensive and puts burden on pathology, clinical staff, etc; NLP can help

HOW DO WE GET THIS GOING?

- Build annotation of "corpus" and text mining into pipeline for acquiring, curating and evaluating biospecimen information from literature
 - Build text mining support for data abstraction and data entry, to increase performance or reduce costs over time
 - Explore potentially fundable text mining / BioNLP research for this new area
-

CONCLUSIONS

- Text mining can potentially be used to extract and evaluate biospecimen information from literature
 - May be a tool to measure and improve practices (i.e. adherence to SOPs and reporting guidelines)
 - Implementing text mining may require shifts in data abstraction/entry practices
-

ACKNOWLEDGEMENTS

- Peter Tarczy-Hornoch, MD
 - Meliha Yetisgen-Yildiz, PhD
 - John Castle, ScD
 - Helen Moore, PhD
 - PNW Prostate Cancer SPORE
 - Annika Havnaer
 - Kelly Engel, PhD
-