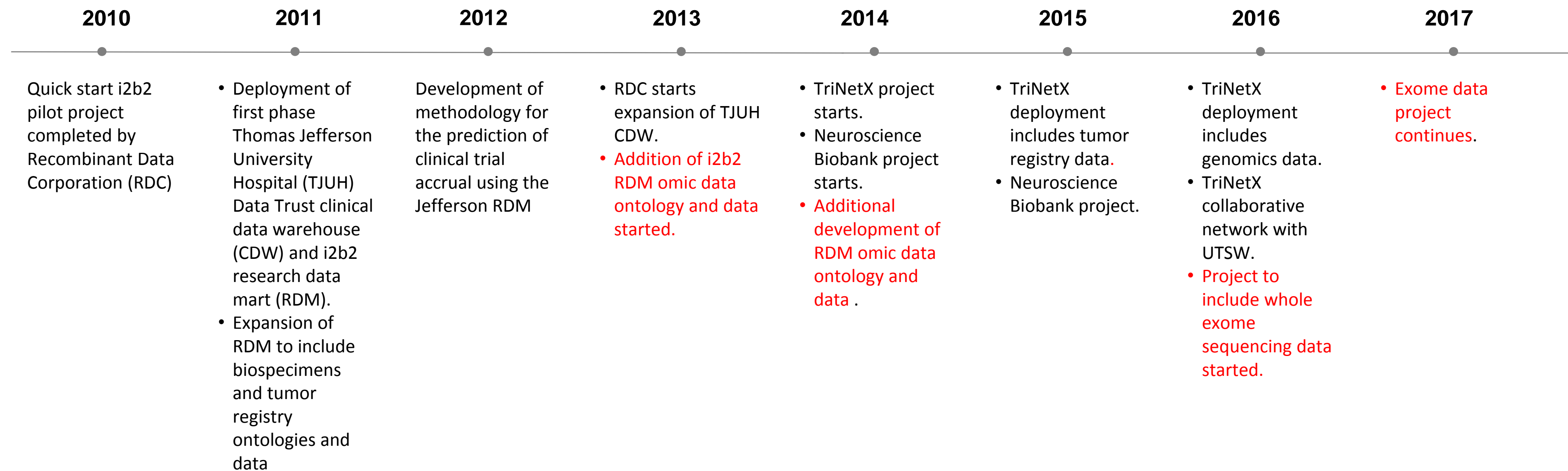# i2b2 Integration and Use of Molecular Diagnostic Genomic Assay Results for Clinical Research

Jack London PhD

Research Professor of Cancer Biology
Thomas Jefferson University
and Informatics Director
Sidney Kimmel Cancer Center
Philadelphia Pennsylvania USA

October 5, 2017
i2b2 2017 European meeting

**Jefferson**
HEALTH IS ALL WE DO

# Jefferson University's Evolution with Research Data Analytics

## Jefferson research informatics timeline

| 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 |
|------|------|------|------|------|------|------|------|
| Quick start i2b2 pilot project completed by Recombinant Data Corporation (RDC) | • Deployment of first phase Thomas Jefferson University Hospital (TJUH) Data Trust clinical data warehouse (CDW) and i2b2 research data mart (RDM).<br>• Expansion of RDM to include biospecimens and tumor registry ontologies and data | Development of methodology for the prediction of clinical trial accrual using the Jefferson RDM | • RDC starts expansion of TJUH CDW.<br>• Addition of i2b2 RDM omic data ontology and data started. | • TriNetX project starts.<br>• Neuroscience Biobank project starts.<br>• Additional development of RDM omic data ontology and data . | • TriNetX deployment includes tumor registry data.<br>• Neuroscience Biobank project. | • TriNetX deployment includes genomics data.<br>• TriNetX collaborative network with UTSW.<br>• Project to include whole exome sequencing data started. | • Exome data project continues. |

Jefferson
HEALTH IS ALL WE DO

# Jefferson's i2b2 research data mart deployment

In addition to **EMR patient data** (demographics, diagnoses, medications, labs, and procedures):

o Comprehensive data set for cancer patients from **Cancer Registry**

- Tumor histology, stage, recurrence, treatment, disease-specific factors

o **"Omic" molecular diagnostic patient data** from Pathology A/P system and outsource vendor system

- Currently > 350 genes with > 4,100 mutations (both in-house and outsourced Foundation Medicine results)

o **Biospecimen annotation** from biobanking system

- Specimen anatomic origin, class, type, pathology, slide images

## More than 125 million observations on 2.9 million patients, refreshed weekly

Jefferson
HEALTH IS ALL WE DO

**Research and applications**

## Design-phase prediction of potential cancer clinical trial accrual success using a research data mart

Jack W London,[1,2] Luanne Balestrucci,[3] Devjani Chatterjee,[1] Tingting Zhan[4]

[1]Kimmel Cancer Center, Thomas Jefferson University, Philadelphia, Pennsylvania, USA
[2]Department of Cancer Biology, Thomas Jefferson University, Philadelphia, Pennsylvania, USA
[3]Jefferson Graduate School of Biomedical Sciences, Thomas Jefferson University, Philadelphia, Pennsylvania, USA
[4]Department of Pharmacology & Experimental Therapeutics, Thomas Jefferson University, Philadelphia, Pennsylvania, USA

**Correspondence to**
Dr Jack London, Kimmel Cancer Center, Thomas Jefferson University, 233 S. 10th Street, Room 808 BLSB, Philadelphia, PA 19107, USA; Jack.london@jefferson.edu

**ABSTRACT**
**Background** Many cancer interventional clinical trials are not completed because the required number of eligible patients are not enrolled.
**Objective** To assess the value of using a research data mart (RDM) during the design of cancer clinical trials as a predictor of potential patient accrual, so that less trials fail to meet enrollment requirements.
**Materials and methods** The eligibility criteria for 90 interventional cancer trials were translated into i2b2 RDM queries and cohort sizes obtained for the 2 years prior to the trial initiation. These RDM cohort numbers were compared to the trial accrual requirements, generating predictions of accrual success. These predictions were then compared to the actual accrual performance to evaluate the ability of this methodology to predict the trials' likelihood of enrolling sufficient patients.
**Results** Our methodology predicted successful accrual (specificity) with 0.969 (=31/32 trials) accuracy (95% CI 0.908 to 1) and predicted failed accrual (sensitivity) with 0.397 (=23/58 trials) accuracy (95% CI 0.271 to 0.522). The positive predictive value, or precision rate, is 0.958 (=23/24) (95% CI 0.878 to 1).
**Discussion** A prediction of 'failed accrual' by this methodology is very reliable, whereas a prediction of accrual success is less so, as causes of accrual failure other than an insufficient eligible patient pool are not considered.
**Conclusions** The application of this methodology to cancer clinical design would significantly improve cancer clinical research by reducing the costly efforts expended initiating trials that predictably will fail to meet accrual

As important as interventional clinical trials are in translational research, these studies may never accrue the statistically required number of participants to complete the study's research plan. An Institute of Medicine (IOM) report on cancer cooperative group trials found that 40% were never completed because of failure to achieve minimum accrual goals.[1] The IOM report states, 'The ultimate inefficiency is a clinical trial that is never completed because of insufficient patient accrual, and this happens far too often.' These non-accruing trials are often kept open for many months before closure, consuming personnel resources in their setup and operation at a significant cost to institutions, without providing any return in definitive research findings. Furthermore, while many of these trials register zero patients, others accrue some patients, resulting in thousands of patients nationwide who are recruited to unproductive research studies.[2] A number of studies have investigated barriers to clinical trial accrual, and reported various physician-related and patient-related obstacles.[3–9] Physician barriers cited include inadequate reimbursement, lack of support resources, the irrelevance of available studies to the practice population, and treatment preferences. Patient barriers cited include concerns and uncertainty about treatments, treatment preferences, unavailability of an appropriate trial, lack of awareness of trials, and transportation and other logistical constraints. These cited studies all have focused on accrual issues occurring *after* trial activation. Recently, however, Schroen *et al*[10] have

---

In 2012 we **used our research data mart** to determine recent cohort sizes of patients satisfying the eligibility criteria for a number of open trials at SKCC.

The study's hypothesis was that the RDM could have predicted whether a sufficient number of patients could be recruited for the trial during the trial's design phase.

**Jefferson**
HEALTH IS ALL WE DO

# Jefferson University's Evolution with Genomics Data Analytics

<u>Precision Medicine's targeted therapies were the driver for integrating genomic data with other patient data.</u>

o Clinical (trial) research is focused on determining genomic markers indicating diagnostic specificity, and predicative of treatment response and outcomes for personalized patient care.

o Inclusion of genomics data necessary for trial design and feasibility analyses.

o Integrated genomics/clinical data also a requirement for hypothesis generation (i.e., biomarker discovery).

Jefferson.
HEALTH IS ALL WE DO

## Consider a clinical trial for locally advanced colon patients "Neoadjuvant Treatment of Colon Cancer (NCT01108107)"

This study will investigate the effect of preoperative combination chemotherapy in patients with locally advanced colon cancer with mutation in the KRAS, BRAF or PIK3CA gene

inclusion criteria
• stage 2 or 3 colon cancer
• KRAS, BRAF, and/or PIK3CA mutation testing determined in a CLIA-certified lab

exclusion criteria
• Clinically significant cardiovascular disease (including myocardial infarction, unstable angina, symptomatic congestive heart failure).

Jefferson
HEALTH IS ALL WE DO

# Genomic Clinical Data Sources

"in-house" Jefferson Pathology Department
clinical molecular diagnostic NGS panel assays

"outsourced" Foundation Medicine, Inc.
clinical molecular diagnostic NGS panel assays

read-only access to patient results database

XML result files electronically sent to Jefferson

E-T-L to Jefferson i2b2 RDM

# Jefferson University's Evolution with Genomics Data Analytics

Initial effort to integrate genomic data with other clinical data started with developing an i2b2 ontology for patients' molecular diagnostic gene panel results.

All available sequencing data elements were included

- gene
- chromosome
- Start/end position
- reference/alternate allele
- DNA/amino acid change
- mutation type
- COSMIC id
- read depth
- aletrnate allele frequency
- ...
- ...



EVERYTHING BUT...... The Kitchen Sink

Jefferson
HEALTH IS ALL WE DO

# Molecular diagnostic assays – gene panels

Discussions with pathologists and oncologists brought the realization that – as in many other situations – the needed data classification is a function of the specific use case being addressed.

For **molecular diagnostic targeted gene panel data**, **used for clinical research cohort identification**, the necessary data elements are:

- tissue sampled
- gene tested
- mutations expressed  (using standard mutation nomenclature based on coding DNA reference sequences and protein-level amino acid sequences – "c." and "p."), or wild type).

Such as, **colon tumor, BRAF, p.V600E, c.1799T>A**

**These data elements correspond to clinical trial eligibility criteria**, and are usually sufficient for mutation identification in databases/literature so that further information can be obtained.

Jefferson.
HEALTH IS ALL WE DO

# SKCC's i2b2 "omic" ontology

# Other genomic data considerations

o When extracting data from some molecular diagnostic sources, thresholds need to be set for the alternate allele frequency and percent tumor of valid variant results.

o The genes tested for an assay need to be known, since all genes that are wild type are not necessarily reported. (Pertinent negatives are sometimes reported; also results for ordered genes may only be reported.)

o The molecular diagnostic ontology (i.e., taxonomy) is not static. Variants not previously reported may occur at any time; genes may be added to test panels at any time (ideally you would be alerted in advance to test panel changes).

Jefferson
HEALTH IS ALL WE DO

# Integration of genomic data with other patient data

For **hypothesis generation** – perhaps for **biomarker discovery** – further data capabilities may also be needed.

- inclusion of high-dimensional data

- statistical pipeline ("R") and other analysis tools (genomic data viewer)

- additional of public data sets (TCGA)

This type of environment found on tranSMART.

# Targeted therapies also a driver for data sharing

o About half of therapeutic oncology trials investigate targeted approaches.

o The small patient sub-populations needed for these studies increase the likelihood that multi-institutional trials are required (and that networked, inter-operative data analytics platforms are available).

## Work now "in progess"

Since last summer, we are working on how to process and integrate whole exome sequencing data, particularly from a biomarker discovery perspective.

Jefferson.
HEALTH IS ALL WE DO

# The Jefferson i2b2 RDM team

Jack London, PhD
SKCC Informatics Director

Stephen Peiper, MD
Chair, Department of Pathology

Robert Stapp, MD
Jefferson Department of Pathology

Chirayu Goswami, MS
Bioinformatician

John Reber
Systems Administrator/Programmer

Funding:  NCI/NIH grant # P30CA056036

Jefferson.
HEALTH IS ALL WE DO