# Introduction to querying electronic health records for patient cohorts using the i2b2 data platform

## I2b2 and the Query Tool
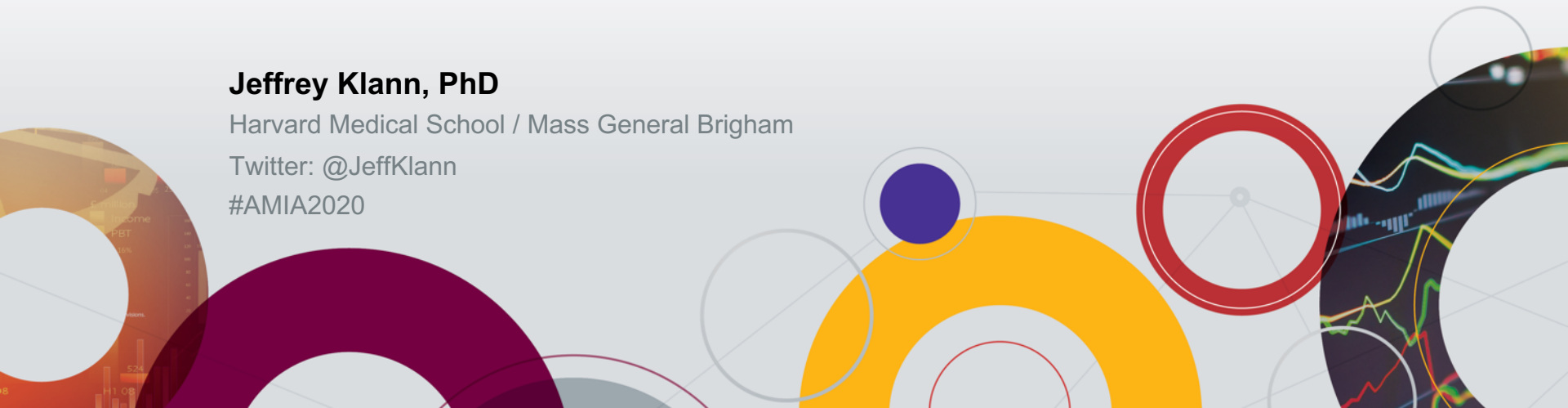
## W16, Session 1: Introduction to i2b2

**Jeffrey Klann, PhD**

Harvard Medical School / Mass General Brigham
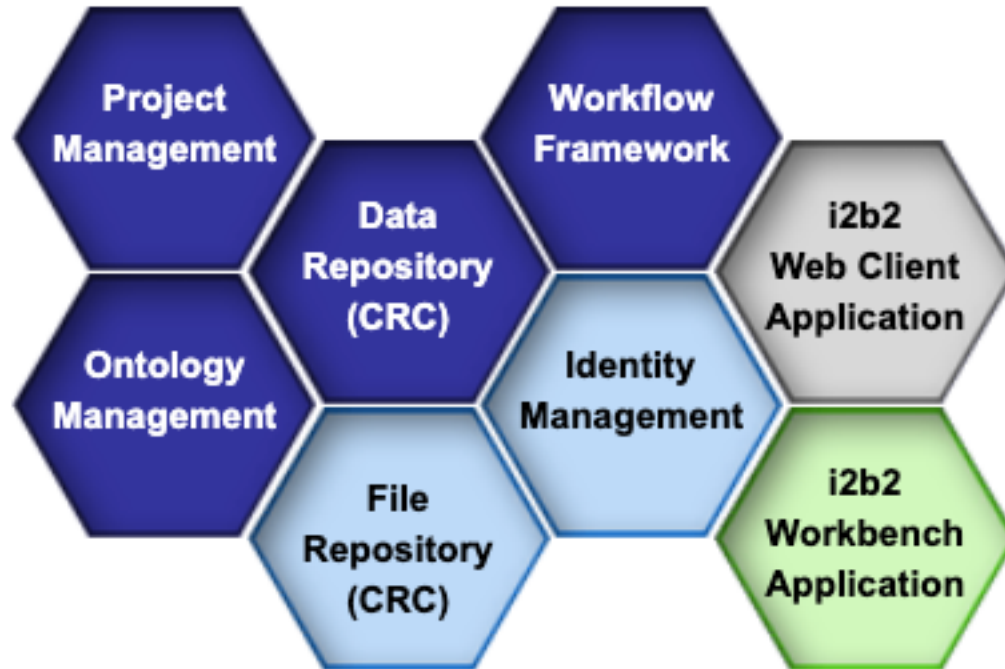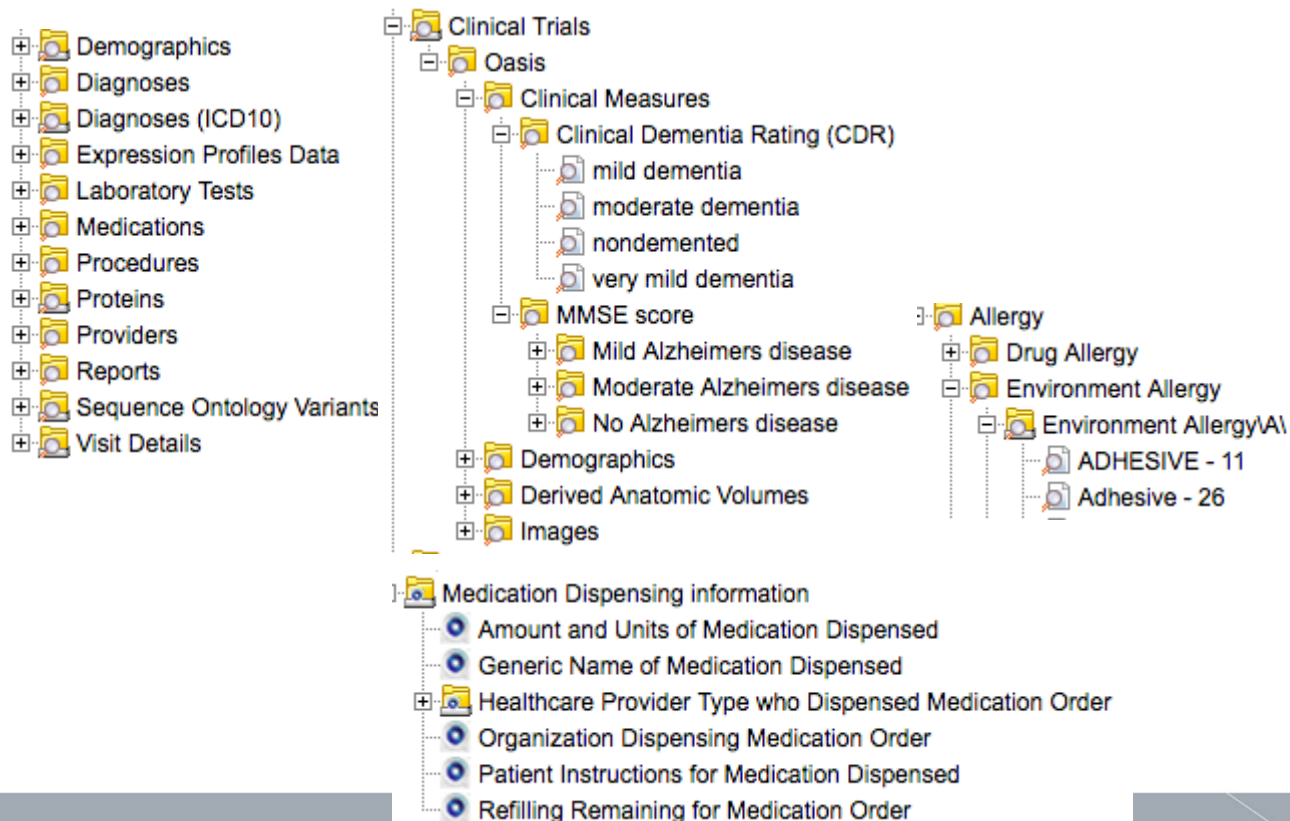
Twitter: @JeffKlann

#AMIA2020

# What is i2b2?

- A Clinical Data Warehousing and Analytics Platform
    - >200 installations worldwide
    - Active user community with many plugins and tools
    - Open-source (Mozilla license)

- Componentized, extensible architecture

- Curated biomedical ontologies

- Flexible star-schema data model

- Graphical query tool

- https://community.i2b2.org/

# The i2b2 "hive"

# Many Ontologies

# The i2b2 Query Tool

# The i2b2 Data Model

# Common Data Model for i2b2 and tranSMART



- **Unified Core Data Model for i2b2 and tranSMART**

- Core tables, fields, keys, indexes, constraints
- Detailed descriptions of each field
- Best practices and examples for different use cases

# Common Data Model for i2b2 and tranSMART

Single schema, documented, MSSQL/Oracle/Postgres

| | Table Name | CONCEPT_PATH | Data Type: MSSQL | Data Type: Oracle | Data Type: Postgres | Primary Key | Foreign Key | Core | Admin | Future | Values | Description | History |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Table Name | CONCEPT_PATH | Data Type: MSSQL | Data Type: Oracle | Data Type: Postgres | Primary Key | Foreign Key | Core | Admin | Future | Values | Description | History |
| 2 | CONCEPT_DIMENSION | CONCEPT_PATH | varchar(700) | varchar2(700) | varchar(700) | Y | | Y | | | | A path that delineates the concept's hierarchy | v1.0 |
| 3 | CONCEPT_DIMENSION | CONCEPT_CD | varchar(50) | varchar2(50) | varchar(50) | | | Y | | | | A code that represents the coded value (e.g. diagnosis, procedure, etc.) | v1.0 |
| 4 | CONCEPT_DIMENSION | NAME_CHAR | varchar(2000) | varchar2(2000) | varchar(2000) | | | | | | | The name of the concept | v1.0 |
| 5 | CONCEPT_DIMENSION | CONCEPT_BLOB | varchar(max) | clob | text | | | | | | | Holds any raw or miscellaneous data that exists, often encrypted PHI or additional information in a parseable format like XML | v1.0 |
| 6 | CONCEPT_DIMENSION | UPDATE_DATE | datetime | date | timestamp | | | | Y | | | Date the row was updated by the source system Obtained from the source system | v1.0 |
| 7 | CONCEPT_DIMENSION | DOWNLOAD_DATE | datetime | date | timestamp | | | | Y | | | Date the data was downloaded from the source system | v1.0 |
| 8 | CONCEPT_DIMENSION | IMPORT_DATE | datetime | date | timestamp | | | | Y | | | Date the data was imported into the CRC | v1.0 |
| 9 | CONCEPT_DIMENSION | SOURCESYSTEM_CD | varchar(50) | varchar2(50) | varchar(50) | | | | Y | | | A coded value for the data source system | v1.0 |
| 10 | CONCEPT_DIMENSION | UPLOAD_ID | int | number(38) | integer | | | | Y | | | A numeric id given to the upload | **v1.3** |
| 11 | OBSERVATION_FACT | ENCOUNTER_NUM | int | number(38) | integer | Y | | Y | | | | i2b2 patient visit number | v1.0 |
| 12 | OBSERVATION_FACT | PATIENT_NUM | int | number(38) | integer | Y | | Y | | | | i2b2 patient number | v1.0 |
| 13 | OBSERVATION_FACT | CONCEPT_CD | varchar(50) | varchar2(50) | varchar(50) | Y | | Y | | | | Code for the observation of interest (i.e. diagnoses, procedures, medications, lab tests) | v1.0 |
| 14 | OBSERVATION_FACT | PROVIDER_ID | varchar(50) | varchar2(50) | varchar(50) | Y | | Y | | | | Practitioner or provider id | v1.0 |
| 15 | OBSERVATION_FACT | START_DATE | datetime | date | timestamp | Y | | Y | | | | Starting date-time of the observation | v1.0 |
| 16 | OBSERVATION_FACT | MODIFIER_CD | varchar(100) | varchar2(100) | varchar(100) | Y | | Y | | | | Code for modifier of interest (i.e. "ROUTE", "DOSE"). Note that the value columns are often used to hold the amounts such as "100" (mg) for the modifier of DOSE or "PO" for the modifier of ROUTE. | v1.0 |
| 17 | OBSERVATION_FACT | INSTANCE_NUM | int | number(18) | integer | Y | | Y | | | | Encoded instance number that allows more than one modifier to be provided for each CONCEPT_CD. Each row will have a different MODIFIER_CD but a similar INSTANCE_NUM. | **v1.4** |
| 18 | OBSERVATION_FACT | VALTYPE_CD | varchar(50) | varchar2(50) | varchar(50) | | | Y | | | N = Numeric T = Text (enums / short messages) B = Raw Text (notes / reports) | Format of the concept | v1.0 |
| | | | | | | | | | | | Used in conjunction with VALTYPE_CD = "T" or "N" | | |

# Common Data Model for i2b2 and tranSMART

- Draft version available for download!

- On https://community.i2b2.org/wiki/

# I2b2/tranSMART Bundle

- Run I2b2 and tranSMART on the same database!

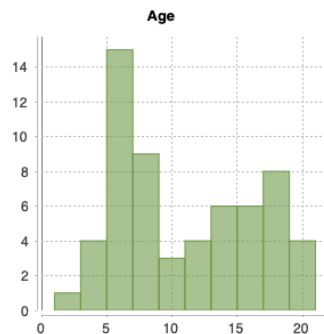- Detailed documentation with step-by-step instructions.



Genomics Analysis Bundle

# tranSMART

I2b2 and tranSMART on the same database

# i2b2 Community Projects

https://community.i2b2.org/

# Accrual to Clinical Trials (ACT) Data Network

# ACT Research Network





Query for sites with adequate data

Curated, mapped ontology

# Genomic Data in i2b2

# Supports Multiple Data Models

# 1.7.13: Community Contributions?

- Develop a new feature!
  - Submit a GitHub pull request
  - Attest to the DCO (Developer Certificate of Origin)
  - **Ideas:** Import/export, client usability, new authentication methods, refactoring, better demo data

- Clean up the Community Projects

- Improve the Documentation

# Synthea/COVID/ACT Demodata

- Synthea synthetic patient data

  – 1.5 million patients

- Mapped to ACT Ontology

  – Demographics, PX, RX, DX, Labs

- Enhanced with COVID Lab Values

  – 20,000 patients



ACT COVID-19
ACT Demographics
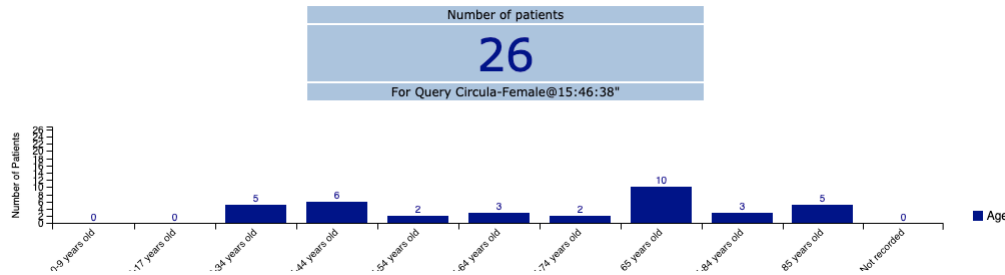ACT Diagnoses ICD-10-CM
  A00-B99 Certain infectious and parasitic diseases (A00-B99) - 46,774
    ACT_2.0.1_UMLS_2018AA
  C00-D49 Neoplasms (C00-D49) - 58,861
  D50-D89 Diseases of the blood and blood-forming organs and certain disorders involving the immune mechanism (D50-D89) - 803
  E00-E89 Endocrine, nutritional and metabolic diseases (E00-E89) - 80,143
  F01-F99 Mental, Behavioral and Neurodevelopmental disorders (F01-F99) - 15,265
  G00-G99 Diseases of the nervous system (G00-G99) - 53,295
  H00-H59 Diseases of the eye and adnexa (H00-H59) - 23
  H60-H95 Diseases of the ear and mastoid process (H60-H95) - 102,542
  I00-I99 Diseases of the circulatory system (I00-I99) - 185,838
  J00-J99 Diseases of the respiratory system (J00-J99) - 1,065,566
  K00-K95 Diseases of the digestive system (K00-K95) - 138,452
  L00-L99 Diseases of the skin and subcutaneous tissue (L00-L99) - 35,124
  M00-M99 Diseases of the musculoskeletal system and connective tissue (M00-M99) - 144,610
  N00-N99 Diseases of the genitourinary system (N00-N99) - 24,494
  O00-O9A Pregnancy, childbirth and the puerperium (O00-O9A) - 89,893
  P00-P96 Certain conditions originating in the perinatal period (P00-P96) - 17,272
  Q00-Q99 Congenital malformations, deformations and chromosomal abnormalities (Q00-Q99) - 967
  R00-R99 Symptoms, signs and abnormal clinical and laboratory findings, not elsewhere classified (R00-R99) - 754,734
  S00-T88 Injury, poisoning and certain other consequences of external causes (S00-T88) - 435,540
  V00-Y99 External causes of morbidity (V00-Y99) - 2,651
  Z00-Z99 Factors influencing health status and contact with health services (Z00-Z99) - 589,783
ACT Diagnoses ICD10-ICD9
ACT Diagnoses ICD-9-CM
ACT Laboratory Tests
ACT Laboratory Tests (Provisional)
ACT Medications Alphabetical
ACT Medications VA Classes
ACT Procedures HCPCS

# Excercises

# Exercise Set 1: Query

1a) Using the ontology browser and find terms, find how many females had circulatory system disease and create a bar chart by age.



1b) How many patients had circulatory system disease and a CPK test >=400 u/l. (Answer: 5 patients.)

1c) Using the query from 1a, find the number of patients with a diagnosis before 11/2005. (Answer: 20).

# Exercise 1 cont'd

1d) Repeat 1c, but now find patients with three or more occurrences of circulatory system disease. (Answer: 13)

1e) How many patients from 1d used a cardiovascular agent medication? (Answer: 6)

# Exercise 1 cont'd

1f) Now view these patients as a timeline.

# Exercise 2: Data Export

Export the previous query using the ExportXLS plugin.

# Exercise 2: Temporal query

Among patients over 65, how many have a diagnosis of asthma prior to their first ever prescription of bronchodilators?

Answer: 6



Now view this as a timeline as well.