











NATIONAL INSTITUTE OF NEUROLOGICAL DISORDERS AND STROKE



Addressing Gaps, Challenges, and Opportunities Related to Metadata Standards

Moving the Field Forward with Technology: A Data Scientist's Perspective

Maria Bagonis, PhD NIDDK Data and Metadata Standards Workshop June 4th, 2025



Defining the Gaps: A Personal Research Story

HUGE disconnect between:

A Predictive AI Model, good enough for a Nature paper vs

A Predictive AI Model that can be productionized and released in the clinic.

Clinical
Predictive Al Model



Lots of High-Quality DATA!

Good Meta Data Standards; Community Cooperation

> No Meta Data Standards; Isolated Labs

Academic Predictive Al Model

Reproducing/Extending a Model Requires Information Regarding:

- Sampling Distributions of the Training Set vs Test Set
- Exact Feature Extraction Processing Pipeline
- The Final Model (Architecture and any Parameter Optimization)









Large, Well-documented, Multi-site Datasets are needed for Biomarker Development and Clinical Translation of Predictive Models To:

- Overcome Small Sample Size Sampling Bias
- Embrace Patient Heterogeneity

Given the magnitude of these tasks, this work CANNOT be realistically achieved by a single lab or even one consortium. It requires a committed community working together effectively as a team.

Led to the Establishment of NIH-funded centralized Repositories:

BRICS









Includes Featured Consortium Datasets

FEATURED DATASE

Adolescent Brain Cognitive Development Study (ABCD)

Use Featured Datasets to quickly access commonly requested datasets

But how do we efficiently query, aggregate data, and synthesize information from these repositories to optimize return on investment?









Challenges: What Makes a Common Data Element 'Common'?

A set of quality CDEs enables a researcher to **quickly and systematically** identify datasets that can be combined and/or compared.

The Logistically Challenging Common Data Element (CDE) Paradox:

CDEs need to be:

- Structured enough to ensure that the necessary information is stored, standardized whenever possible, and communicated effectively, to facilitate dataset aggregation and comparison.
 - i.e. Eliminate Unintentional, Unnecessary Divergences
- Flexible enough to evolve with changing technology and goals:
 - i.e., Keep Intentional, Scientifically Necessary Divergences;

How do we deal with the nuances of this paradox on a large scale?









Opportunities: The Al Large Language Model Revolution

Text Embedding Models and Generative Large Language Models Make it Easier for Computers to Help us Synthesize Detailed, Nuanced Text Quickly, Interactively, and Automatically

Current Al Projects At BRICS to:

- Cross-map and harmonize new measurement variables of interest with existing common data element dictionaries, thereby reducing the introduction of unnecessary redundancy, facilitating automated data aggregation, and ensuring robust clinical data standards that meet interoperability requirements.
- Automatically identify and aggregate relevant cohort-specific datasets across studies and/or repositories that can be combined or compared to answer a specific research question, facilitating machine learning and AI applications, as well as efficient biomarker development.







Avoiding the Data Element (DE) Dictionary

Explosion: Reducing Unintentional Redundancy

BRICs Tool Development: Approaches to Facilitate Data Dictionary Searches



UMLS Code Assignments to Group Common Data Elements

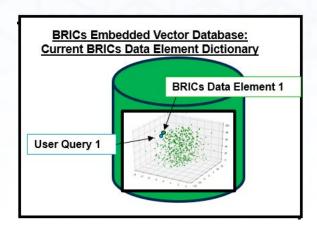
(Standardized Concept Assignments to DEs Via Manual Curation)

UMLS Concept1:

BRICs Data Element 1 BRICs Data Element 2 BRICs Data Element 3... Etc

Semantic Similarity Searches to Group Common Data Elements

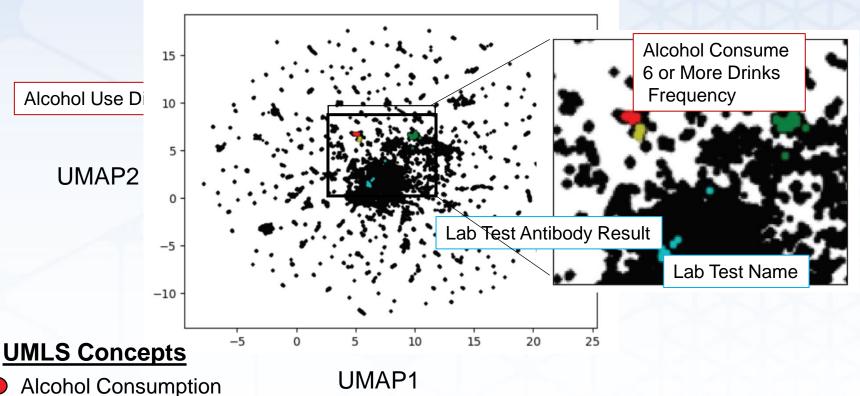
Pre-Trained Large Language Models transform DE text such that DEs with similar meaning cluster in embedding space





Text Embedding Models Automatically **Cluster Common Data Elements by** Concept

LLM Embedding Space: FITBIR Data Dictionary (Vector = 768)



- Illicit Drugs
- **Laboratory Procedures**
- Headache



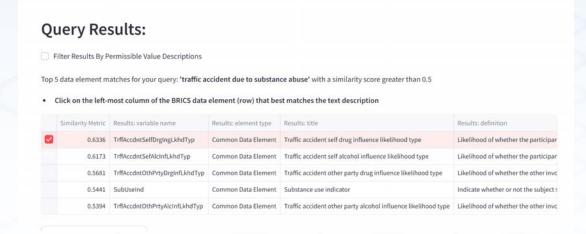




Specific Al Projects in Development at BRICs: Data Element Dictionary Cross-Mapping Tool

AIMS4DE: Artificial Intelligence Meaningful Search for Data Elements Public Beta-Version: https://brics.cit.nih.gov/aims4de/

Currently, searches FITBIR data element dictionary only, but an update allowing searches over all BRICs Instances and additional features will be pushed to the public soon.



Performs well where traditional keyword searches fail

Incorporating the same semantic search method into different UI/Tools:

- <u>Data Mapping and Transformation (DMT) Tool (IP)</u>
- <u>Data Discovery and Dataset Visualization</u>: (IP)



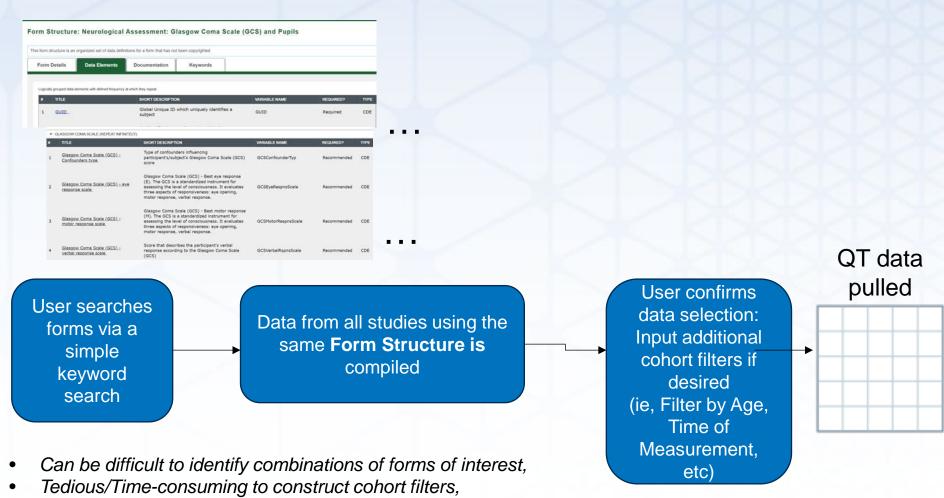






Standardized Form Structures Facilitate Dataset Aggregation

Form Structure: Collection of Common Data Elements



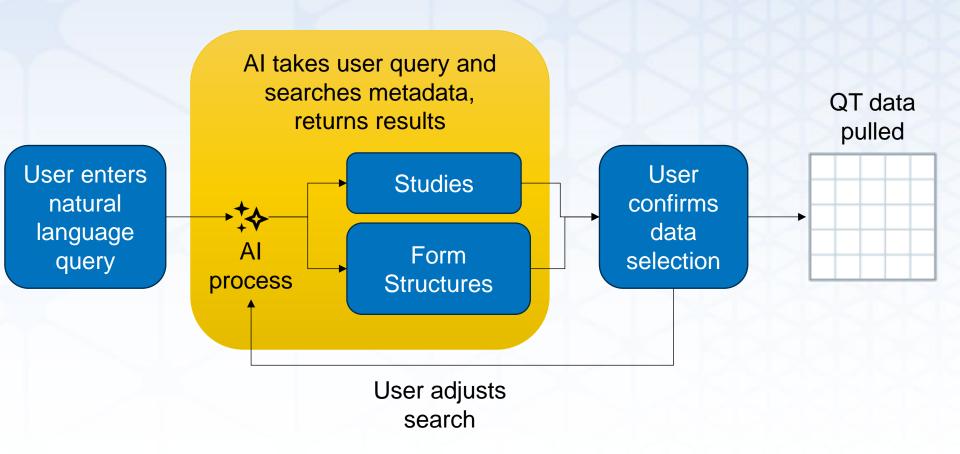




No way to quickly check for sampling distributions in the aggregated dataset, without downloading



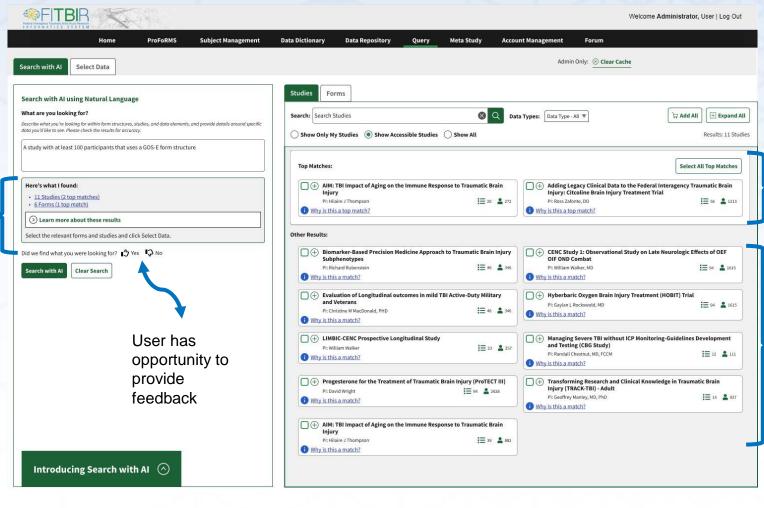
Specific Al Projects in Development at BRICs: Al for Dataset Queries







Specific Al Projects in Development at BRICs: Al for Dataset Queries



Al shows results, sorted into Studies & Forms, and provides summary of results

Top
matches
highlighted,
with one
click ability
to select

Other results listed below



Expandable accordion to view full summary

Search with AI using Natural Language What are you looking for? Describe what you're looking for within form structures, studies, and data elements, and provide details around specific data you'd like to see. Please check the results for accuracy: A study with at least 100 participants that uses a GOS-E form structure Here's what I found: 11 Studies (2 top matches) 6 Forms (1 top match) Learn more about these results Select the relevant forms and studies and click Select Data. Did we find what you were looking for? 13 Yes 13 No Search with AI Clear Search

Specific Al Projects in Development at BRICs: Al for Dataset Queries

Search with Al using Natural Language

What are you looking for?

Describe what you're looking for within form structures, studies, and data elements, and provide details around specific data you'd like to see. Please check the results for accuracy.

A study with at least 100 participants that uses a GOS-E form structure

Here's what I found:

- · 11 Studies (2 top matches)
- · 6 Forms (1 top match)

(^) Learn more about these results

Based on your query for data on anxiety in elderly subjects, here is a comprehensive response:

Primary Form Structure:

- · Form: Generalized Anxiety Disorder (GAD-7)
- Form ID: 296
- Justification: This form is specifically designed to screen and measure anxiety symptoms, making it ideal for analyzing anxiety in elderly populations.

Primary Data Element:

- · Data Element: 'GAD7WryTooMchScale'
- Reason: Directly measures anxiety symptoms using the GAD-7 scale, with a specific focus on worrying - a key indicator of anxiety particularly relevant to elderly subjects.

Filter Criteria:

- · Data Element: 'AgeYrs'
- · Filter Expression: Subjects aged 65 years and older
- Justification: Allows precise selection of elderly participants by capturing age in full years

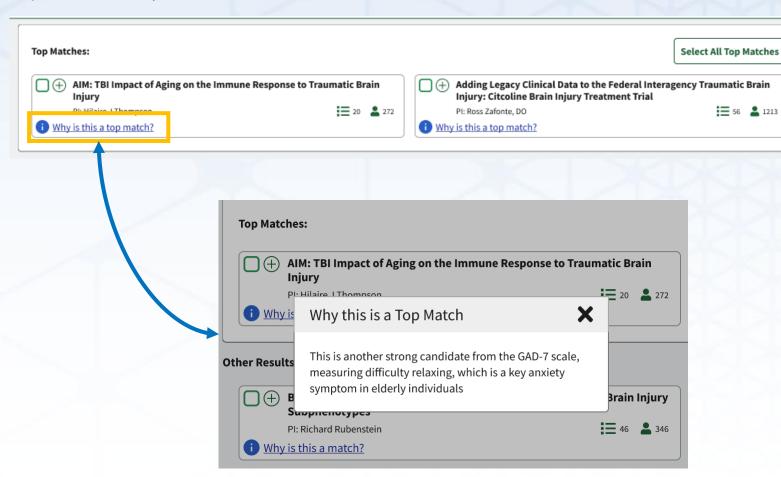
Relevant Studies: The following studies were identified as potentially containing relevant anxiety data:

- 1. FITBIR-STUDY0000259
- 2. FITBIR-STUDY0000242
- 3. FITBIR-STUDY0000243
- 4. FITBIR-STUDY0000244
- FITBIR-STUDY0000245FITBIR-STUDY0000247
- 7. FITBIR-STUDY0000250



Specific Al Projects in Development at BRICs: Al for Dataset Queries

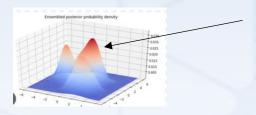
See explanation of why a result is a match





What can be done moving forward: Technical Implementation

 Clustering of unharmonized data element dictionaries from current repositories (Ex: Open Data Commons for Spinal Cord Injury) using text embedding models



High Frequency
Concept Cluster – Which variables can
be collapsed to a single
Common Data Element

- Implementation of Interactive, Collaborative Platforms for CDE curation that utilize the automated organization afforded by Al Models.
- Create an Al interactive personal assistant to help researchers explore the repository and plan for their next study:
 - Suggesting common data elements/forms for their research goals
 - Recommending other studies to consider using those data elements.
 - Automatically aggregating and suggesting potential datasets of use.
 - Providing dataset visualizations and flagging potential sampling bias of a given dataset.
- Consider federated systems for data storage
- Design new tools, workflows, and standards for archiving image feature extraction and predictive machine learning pipelines:
 - ML Experiment Tracking (mlflow, etc)
 - Apps for Easy Archival and Reproduction of Processing Pipelines: Single record containing the entire data workflow (from raw data to published figures) https://brainlife.io/about/ Pestilli Lab









What can be done moving forward: Governance and Culture

- Organize a structured process for additional cross-talk and coordination among the Common Data Element/Repository Leadership Groups across NIH to promote efficiency and eliminate unintentional redundancy in efforts.
- Encourage the practice of good software development skills in the context of a medical regulatory landscape earlier in the research cycle.
- Re-examine traditional funding paradigms when building long-term, collaborative infrastructure.
- Create a system that rewards scientists professionally for collaboration, data sharing, and adherence to high-quality data standards. Encourage respect for team science.
- Implement more formalized, well-documented, iterative 'market' research when developing data repositories and tools.
- **Design and automatically monitor new key performance indicators** (beyond # of papers) to obtain more quantitative statistics for collaboration and technology transfer.





True efficiency arises through collaboration, coordination, recognizing and celebrating human genetic heterogeneity, and collectively building effective infrastructure that provides us all with the bandwidth for enhanced creativity.

