

Name:	Chi Yuan
Affiliation:	Columbia University
Email:	cy2465@cumc.columbia.edu
Presentation type (s):	Software Demonstration

## Criteria2Query: Automatically Transforming Clinical Research Eligibility Criteria Text to OMOP Common Data Model (CDM)-based Cohort Queries

Chi Yuan, MS<sup>1,4</sup>, Patrick B. Ryan, PhD<sup>1,2,3</sup>, Yixuan Guo, MA<sup>1</sup>,  
Peng Jin, MS<sup>1</sup>, Kang Tian, MA<sup>1</sup>, Chunhua Weng, PhD<sup>1</sup>

<sup>1</sup>Department of Biomedical Informatics, Columbia University, New York, NY, USA;

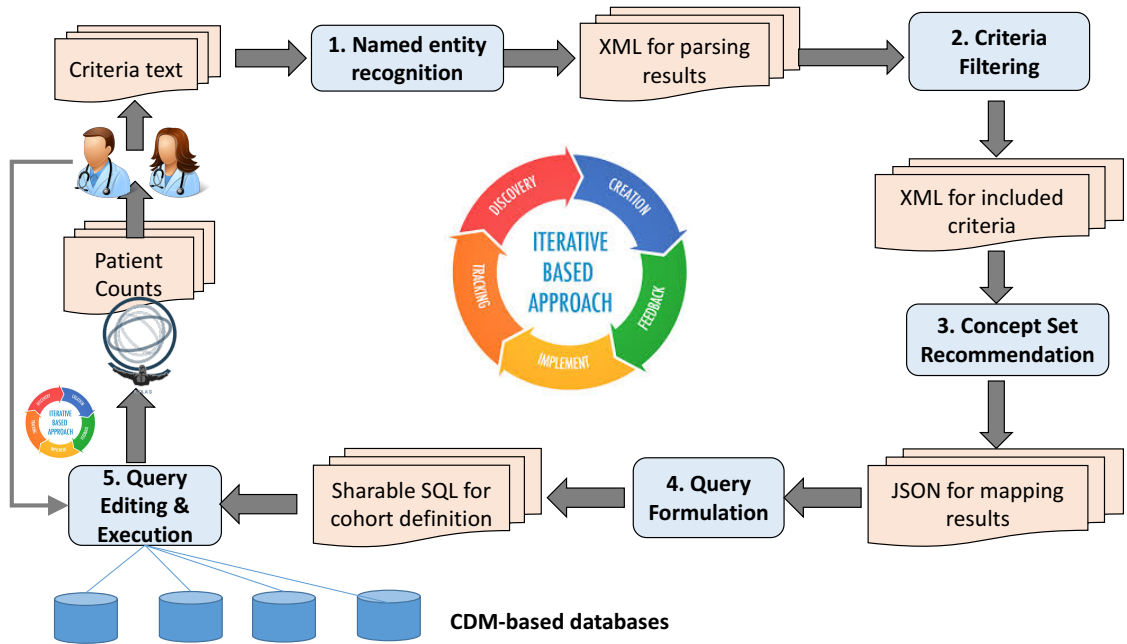
<sup>2</sup>Observational Health Data Sciences and Informatics, New York, NY, USA;

<sup>3</sup>Janssen Research & Development, LLC, Titusville, NJ, USA;

<sup>4</sup>School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing, Jiangsu, P.R. China

Patient recruitment has been a persistent barrier to clinical and translational research<sup>1</sup>. Over 50% of studies fail due to difficulties in recruitment. The universal adoption of electronic health records makes it possible to query distributed, large clinical databases for prescreening potentially eligible patients for performing feasibility assessment or cohort identification for clinical studies. However, interpretations of clinical research eligibility criteria may vary from site to site or from person to person, which can lead to incompatible or mismatched patient queries and compromise the integrity of multi-site studies. In addition, encoding of eligibility criteria into database queries often involves laborious human effort, which is costly, not scalable, and often prohibitive for clinicians and researchers, who do not necessarily understand informatics and know how to map criteria concepts to their data representations. This study aims to build a natural language processing system to transform eligibility criteria text into standards-based cohort identification queries that are sharable and configurable by end users, who can then focus their energy on feasibility assessment and iterative refinement of the criteria based on dynamic feedback with patient counts during eligibility criteria design.

The system Criteria2Query consists of five modules (**Figure 1**): i.e., text parsing, criteria filtering, terminology standards-based concept mapping, automatic query formulation, and query execution with dynamic feedback generation for users. Criteria parsing is supported by an open-source machine learning parser, EliIE<sup>2</sup>, which outputs XML representations for recognized entities and their temporal relationships using the Observational Medical Outcomes Partnership (OMOP) common data model (CDM) V5<sup>3</sup>. Criteria (e.g., “willing to sign the consent” and “able to walk 5 miles on treadmill”) that cannot be queried within EHR for prescreening purposes are filtered out based on heuristics or empirical knowledge. Concepts and their relationships in retained criteria are then mapped to the Observational Health Data Sciences and Informatics (OHDSI) controlled clinical vocabularies to obtain concept IDs for each concept and to create structured cohort definitions, which are further translated into cohort queries in JSON or SQL formats using a public OHDSI Web API (<https://github.com/OHDSI/WebAPI>). These queries can be executed against any OMOP CDM-based patient database to prescreen potentially eligible patients.



**Figure 1.** The behind-the-scenes parsing pipeline that enables iterative eligibility criteria design.

Criteria2Query Support

Transforming Eligibility Criteria Text to Cohort Queries

---

ClinicalTrials.gov Identifier:  [Get From ClinicalTrials.gov](#)

**Inclusion Criteria:(one criterion per line)** **Exclusion Criteria:(one criterion per line)**

Patients have type 2 diabetes and Alzheimer's Disease

[↩ Parsing](#)   [Format ↘](#)

#	Inclusion Criteria:	EHR Status
0	Patients have <span style="border: 1px solid red; padding: 2px;">type 2 diabetes</span> <span style="border: 1px solid red; padding: 2px;">AND</span> <span style="border: 1px solid red; padding: 2px;">Alzheimer's Disease</span>	YES
#	Exclusion Criteria:	EHR Status
No matching records found		

[Searching](#)

**Figure 2.** Snapshot of free-text parsing page.

Criteria2Query

ConceptSet Recommendation

**type 2 diabetes** [Sync with ATLAS](#) [Create New ConceptSet](#)

<input type="checkbox"/>	id	title
<input type="checkbox"/>	499	Type 2 Diabetes NEw
<input type="checkbox"/>	2737	test Type 2 diabetes
<input checked="" type="checkbox"/>	3218	type 2 diabetes mellitus
<input type="checkbox"/>	24616	Type 2 Diabetes Mellitus
<input type="checkbox"/>	103068	type 2 diabetes mellitus
<input type="checkbox"/>	101231	SK Type 2 Diabetes Mellitus
<input type="checkbox"/>	917966	Exclusion concepts for Type 2 diabetes
<input type="checkbox"/>	917882	Type 2 diabetes mellitus - inclusion codes from EMIF

**Alzheimer's Disease** [Sync with ATLAS](#) [Create New ConceptSet](#)

<input type="checkbox"/>	id	title
<input checked="" type="checkbox"/>	105252	test-Alzheimer's disease
<input type="checkbox"/>	105254	yctest0406Alzheimer's disease

[GenerateJSON](#)

Figure 3. Snapshot of concept set mapping page.

Criteria2Query

JSON Result

[Check it On ATLAS](#) [Translate to SQL](#)

```
{
  "ConceptSets": [
    {
      "id": "3218",
      "name": "type 2 diabetes mellitus",
      "expression": {
        "items": [
          {
            "concept": {
              "CONCEPT_ID": "201826",
              "CONCEPT_NAME": "Type 2 diabetes mellitus",
              "STANDARD_CONCEPT": "S",
              "INVALID_REASON": "V",
              "CONCEPT_CODE": "44054006",
              "DOMAIN_ID": "Condition",
              "VOCABULARY_ID": "SNOMED",
              "CONCEPT_CLASS_ID": "Clinical"
            }
          }
        ]
      },
      "Finding": {
        "INVALID_REASON_CAPTION": "Valid",
        "STANDARD_CONCEPT_CAPTION": "Standard",
        "isExcluded": false,
        "includeDescendants": false,
        "includeMapped": false
      },
      "id": "105252",
      "name": "test-Alzheimer's disease",
      "expression": {
        "items": [
          {
            "concept": {
              "CONCEPT_ID": "378419",
              "CONCEPT_NAME": "Alzheimer's disease",
              "STANDARD_CONCEPT": "S",
              "INVALID_REASON": "V",
              "CONCEPT_CODE": "26929004",
              "DOMAIN_ID": "Condition",
              "VOCABULARY_ID": "SNOMED",
              "CONCEPT_CLASS_ID": "Clinical"
            }
          }
        ]
      },
      "Finding": {
        "INVALID_REASON_CAPTION": "Valid",
        "STANDARD_CONCEPT_CAPTION": "Standard",
        "isExcluded": false,
        "includeDescendants": true,
        "includeMapped": false
      }
    }
  ],
  "PrimaryCriteria": {
    "CriteriaList": [
      {
        "ConditionOccurrence": {}
      }
    ],
    "ObservationWindow": {
      "PriorDays": 0,
      "PostDays": 0,
      "PrimaryCriteriaLimit": {
        "Type": "First"
      },
      "AdditionalCriteria": {
        "Type": "ALL",
        "CriteriaList": [
          {
            "Criteria": {
              "ConditionOccurrence": {
                "CodesetId": "3218",
                "StartWindow": {
                  "Start": {
                    "Coeff": -1,
                    "End": {
                      "Coeff": "1"
                    },
                    "Occurrence": {
                      "Type": "2",
                      "Count": "1"
                    }
                  }
                },
                "Criteria": {
                  "ConditionOccurrence": {
                    "CodesetId": "105252",
                    "StartWindow": {
                      "Start": {
                        "Coeff": -1,
                        "End": {
                          "Coeff": "1"
                        },
                        "Occurrence": {
                          "Type": "2",
                          "Count": "1"
                        }
                      }
                    },
                    "DemographicCriteriaList": [],
                    "Groups": [],
                    "QualifiedLimit": {
                      "Type": "First",
                      "ExpressionLimit": {
                        "Type": "First",
                        "InclusionRules": [],
                        "CensoringCriteria": []
                      }
                    }
                  }
                }
              }
            }
          }
        ]
      }
    }
  }
}
```

Figure 4. Snapshot of JSON result page.

Patient Count :138

Please select one patient database:

SYNPUF 1K

Get patient count

SQL

```
CREATE TABLE #Codesets (  
  codeset_id int NOT NULL,  
  concept_id bigint NOT NULL  
)  
;  
  
INSERT INTO #Codesets (codeset_id, concept_id)  
SELECT 3218 as codeset_id, c.concept_id FROM (select distinct l.concept_id FROM  
(  
  select concept_id from @cdm_database_schema.CONCEPT where concept_id in (201826)and invalid_reason is null  
) l  
) C;  
INSERT INTO #Codesets (codeset_id, concept_id)  
SELECT 105252 as codeset_id, c.concept_id FROM (select distinct l.concept_id FROM  
(  
  select concept_id from @cdm_database_schema.CONCEPT where concept_id in (378419)and invalid_reason is null  
UNION select c.concept_id  
from @cdm_database_schema.CONCEPT c  
join @cdm_database_schema.CONCEPT_ANCESTOR ca on c.concept_id = ca.descendant_concept_id  
and ca.ancestor_concept_id in (378419)  
and c.invalid_reason is null  
) l  
) C;
```

**Figure 5.** Snapshot of SQL and Patient count page.