# GazOntology

## A Tool for Building GATE Gazetteer Lists from Ontologies

Daniel R. Schlegel, Rose Fontana, Adrian Naaktgeboren

Department of Computer Science
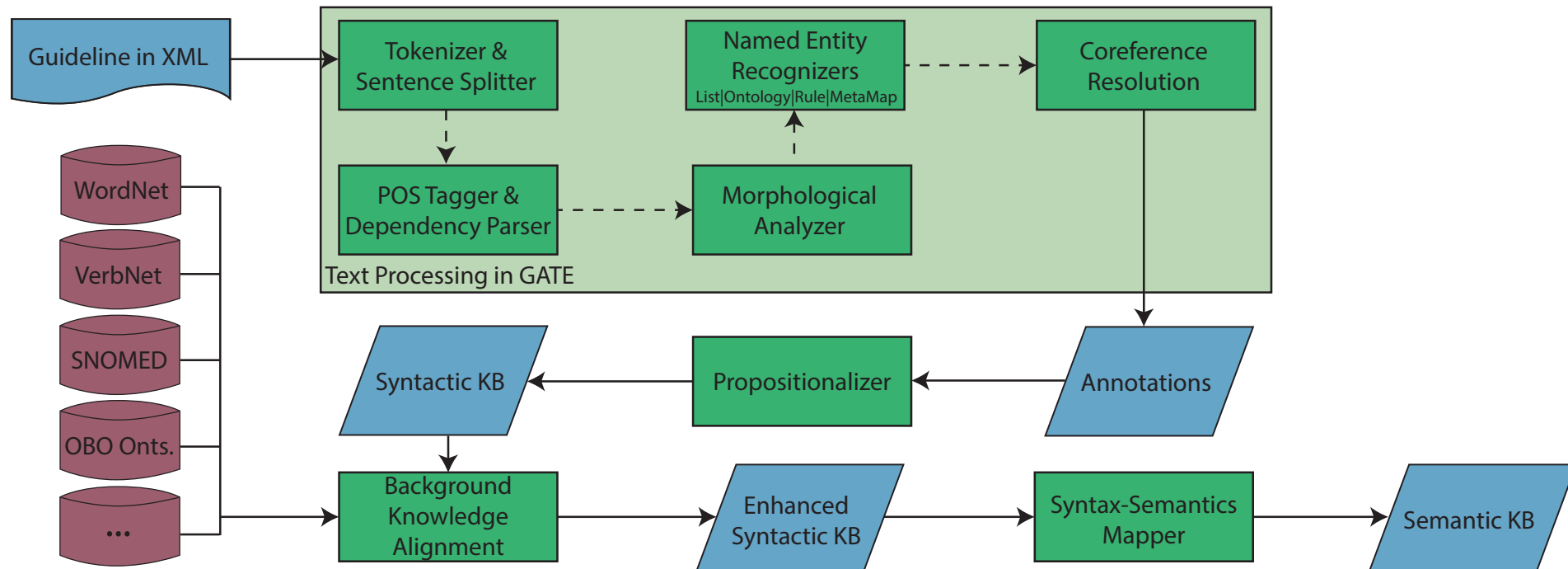
OSWEGO
STATE UNIVERSITY OF NEW YORK

# Outline

- Motivation: Our Application and Philosophy
- GATE and 'Gazetteer' Lists
- Two Examples
- Our Application, Revisited

# Clinical Tractor

- A tool for natural language understanding of clinical practice guidelines
    - Under development
    - Goal: represent the semantics of guidelines for the purpose of automatically generating computer interpretable guidelines.

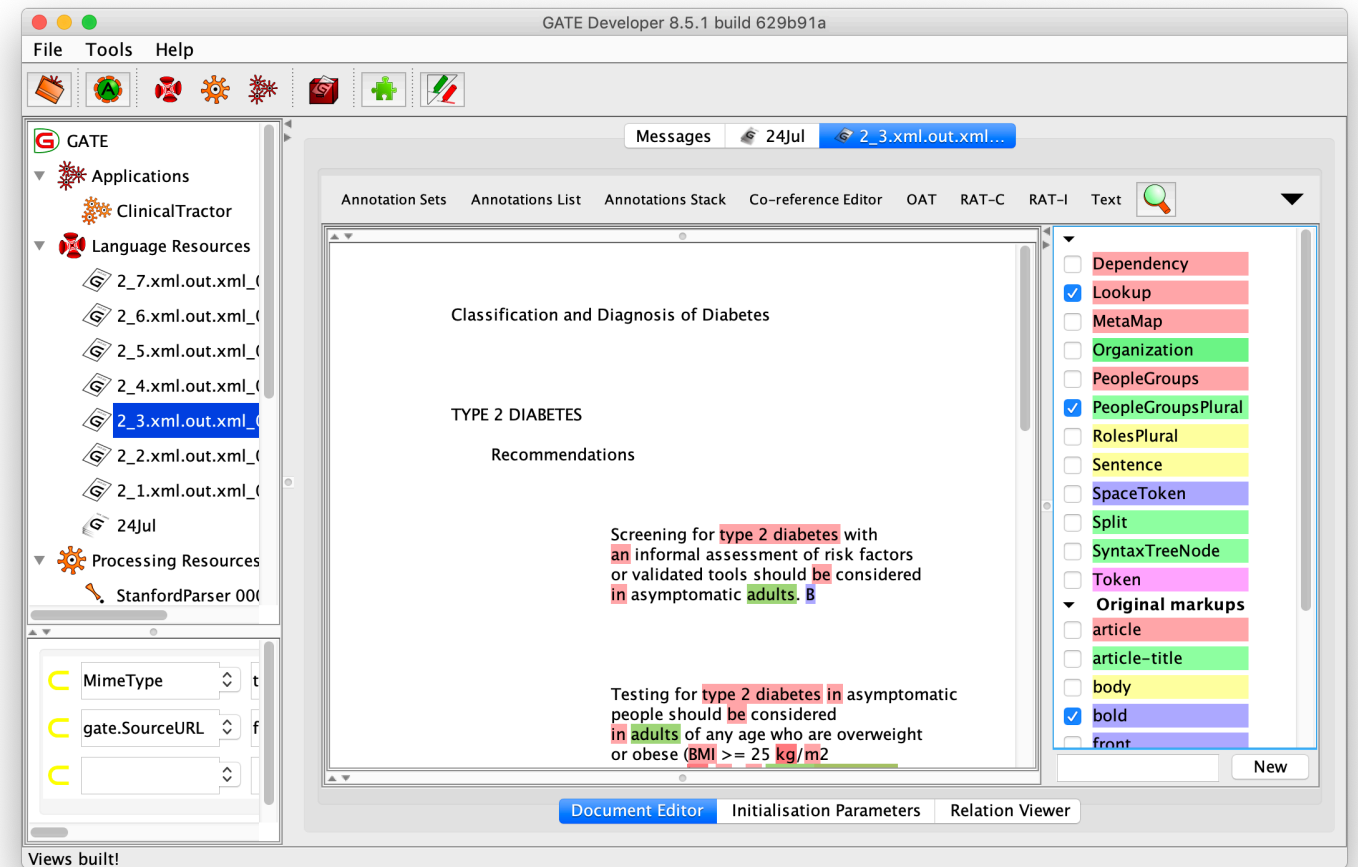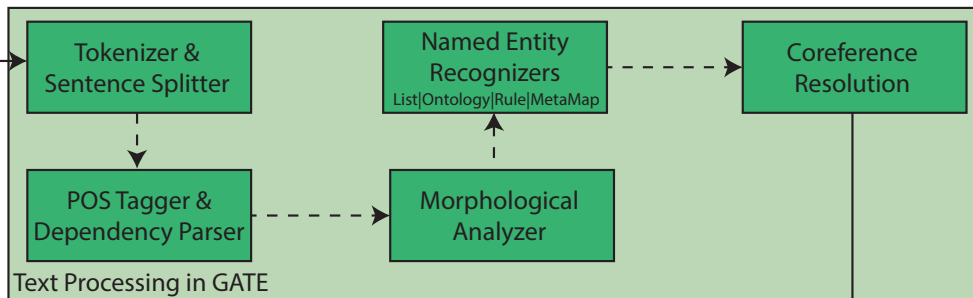# Our Natural Language Understanding Strategy

- Use off-the-shelf tools and data for the standard NLP pipeline
  - Make customizations as needed for the domain

- When mapping syntax to semantics:
  - Be as general as possible
  - Avoid specific cases wherever possible
  - Make judicious use of background knowledge to maintain generality

We use GATE to coordinate the standard NLP pipeline. We want it to provide us "hooks" into background knowledge we can use later.

# GATE: General Architecture for Text Engineering

- Packages several tools for NLP
  - Notably for us today, named entity recognition tools
- Easy to use interface
  - But can run files in batch
- Extendable to add additional *processing resources*



In Clinical Tractor…

Text Processing in GATE

# Gazetteer Lists

# GazOntology

- Benefit of biomedical domain: lots of ontologies!
  - How can we leverage them in GATE?

- Build Gazetteer lists from ontology terms
  - Use labels and synonyms
  - Retain mapping to IRI

# A First Example…

"Blood glucose rather than A1C should be used to diagnose the acute onset of type 1 diabetes in individuals with symptoms of hyperglycemia." –Standards of Care in Diabetes 2017

**Class: hyperglycemia**

   **Term IRI:** http://purl.obolibrary.org/obo/DOID_4195

**Annotations**

- **database_cross_reference:** UMLS_CUI:C0020456; NCI:C26797; ICD10CM:R73.9; MESH:D006943; SNOMEDCT_US_2018_03_01:80394007
- **has_obo_namespace:** disease_ontology
- **id:** DOID:4195
- **in_subset:** http://purl.obolibrary.org/obo/doid#NCIthesaurus

**Class Hierarchy**

Thing
   + disease
      + disease of metabolism
         + acquired metabolic disease
            + carbohydrate metabolism disease
               + glucose metabolism disease
                  - triosephosphate isomerase deficiency
                  - hyperinsulinism
                  + diabetes mellitus
                  - hypoglycemia
                  - hyperglycemia
                     - glucose intolerance

**Superclasses & Asserted Axioms**
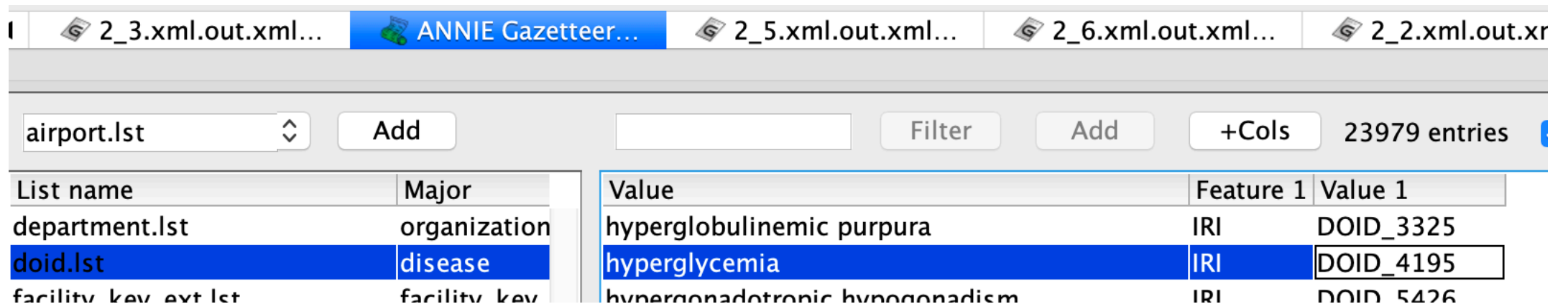
- glucose metabolism disease

# Using GazOntology

- Download the tool and the ontology of interest

- Run the tool

  ```
  java -jar GazOntology-1.0.1.jar -f doid.owl -o doid.lst
  ```
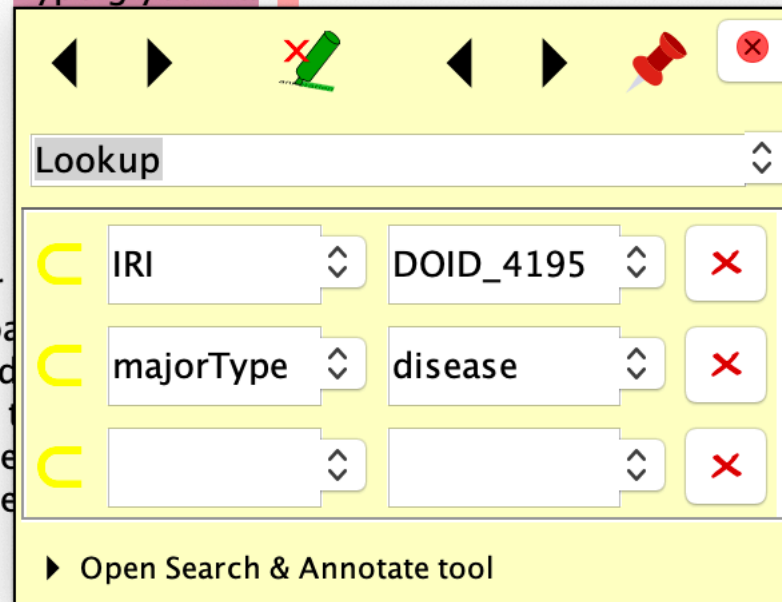
- Add the Gazetteer list to GATE

# After Running GATE…

# One Advantage: Add Synonyms Easily

- Adding missing synonyms is simple in the GATE interface
  - Add them in the ontology Gazetter list directly, or make a new list for additions
  - Likely easier to update ontology versions this way than editing the ontology directly

"Screening for type 2 diabetes with an informal assessment of risk factors or validated tools should be considered in asymptomatic adults."

**Class: type 2 diabetes mellitus**

**Term IRI:** http://purl.obolibrary.org/obo/DOID_9352

- **has_exact_synonym:** type II diabetes mellitus; non-insulin-dependent diabetes mellitus; NIDDM
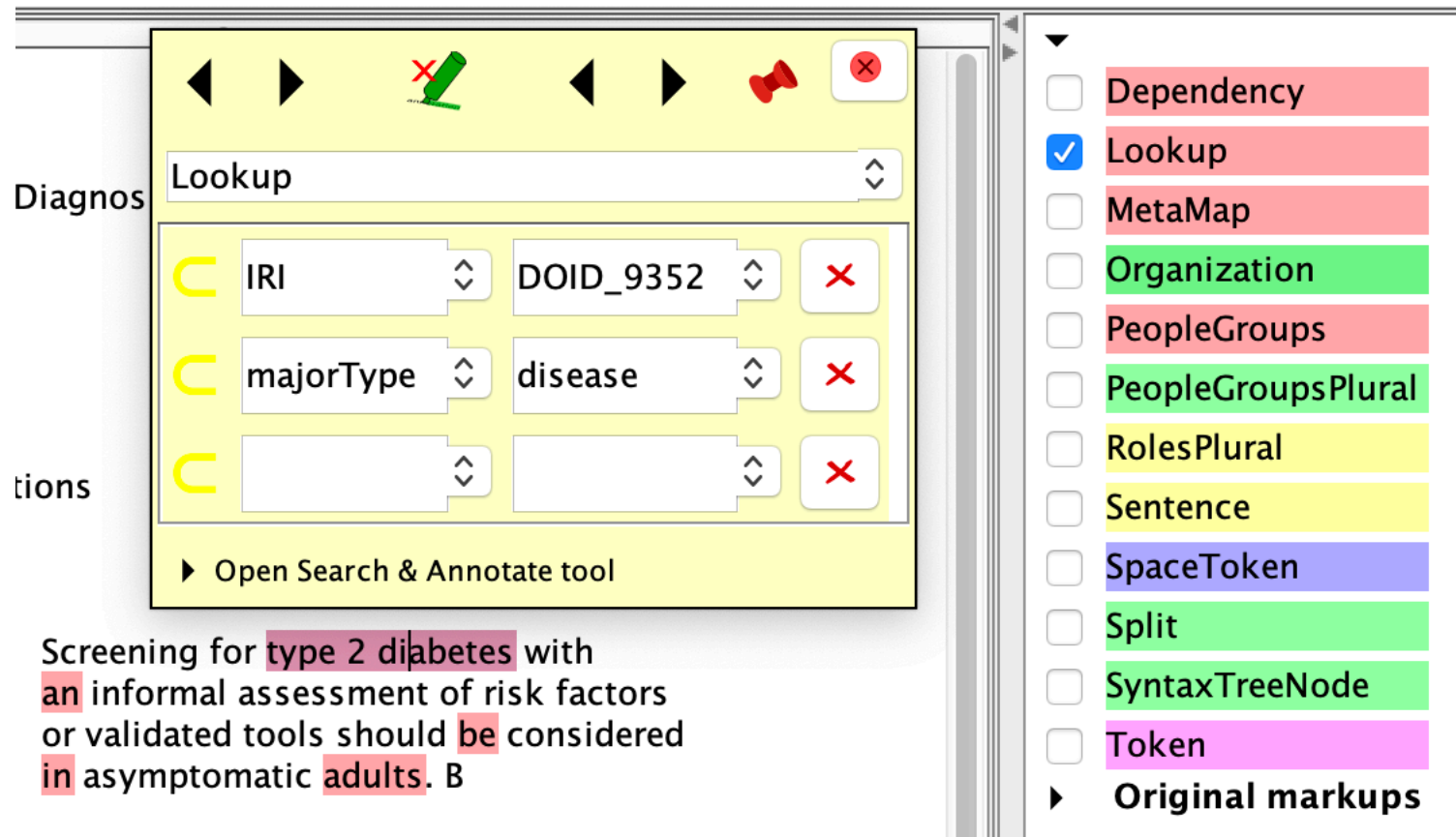
| Messages | ANNIE Gazetteer... |

| airport.lst | Add | | Filter | Add | +Cols | 23979 entrie |

| List name | Major | Minor | Language | Annotation type |
| --- | --- | --- | --- | --- |
| department.lst | organization | government | | Lookup |
| doid.lst | disease | | | Lookup |
| facility.lst | facility | building | | Lookup |

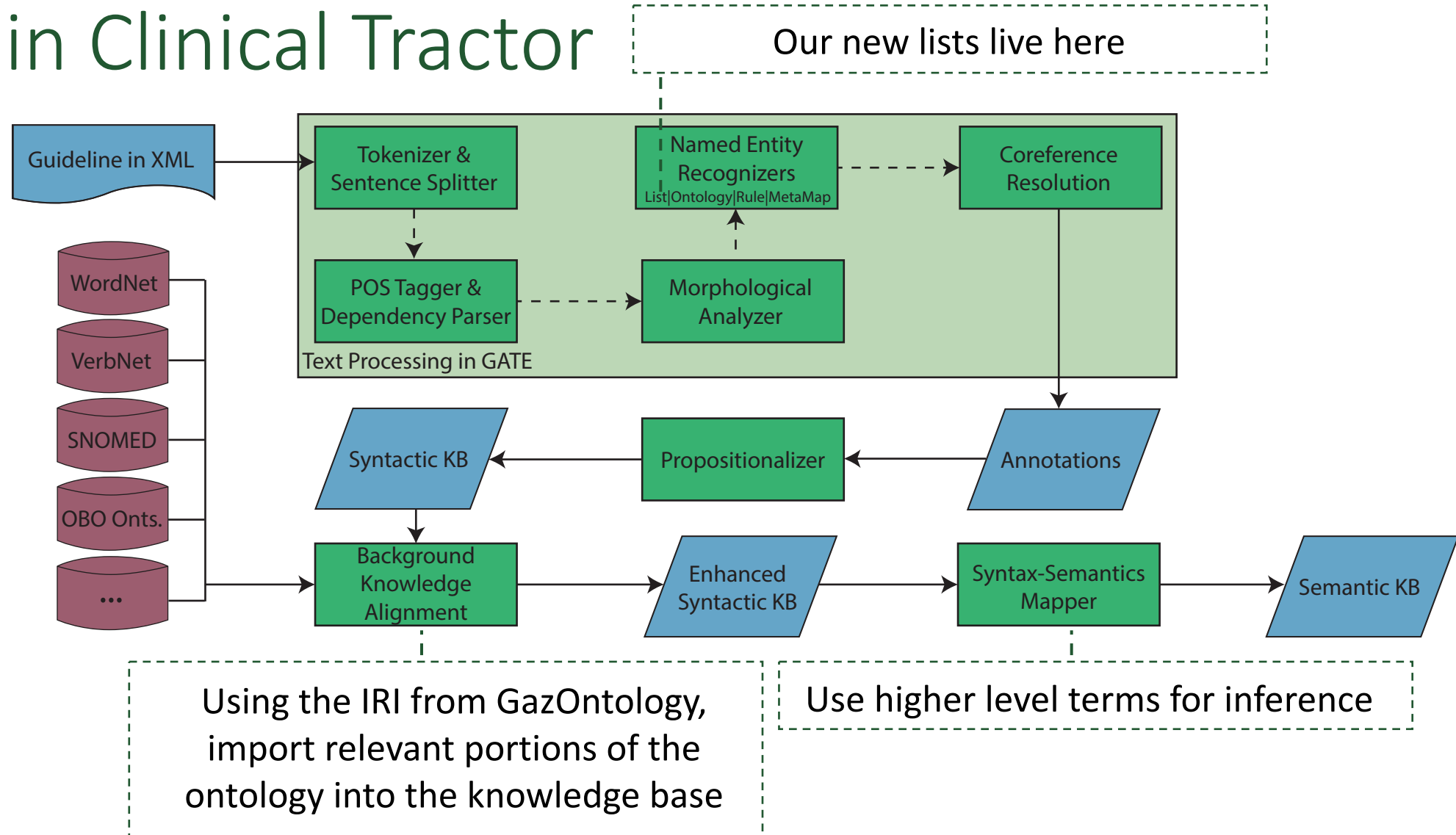| Value | Feature 1 | Value 1 |
| --- | --- | --- |
| type 1 papillary adenoma of the kidney | IRI | DOID_6258 |
| type 2 diabetes | IRI | DOID_9352 |
| type 2 diabetes mellitus | IRI | DOID_9352 |

# Easy to Test: Just Re-run the Gazetteer

# Use in Clinical Tractor



Our new lists live here

Guideline in XML → Text Processing in GATE: Tokenizer & Sentence Splitter → POS Tagger & Dependency Parser → Morphological Analyzer → Named Entity Recognizers (List|Ontology|Rule|MetaMap) → Coreference Resolution

WordNet, VerbNet, SNOMED, OBO Onts., ...

Annotations → Propositionalizer → Syntactic KB → Background Knowledge Alignment → Enhanced Syntactic KB → Syntax-Semantics Mapper → Semantic KB

Using the IRI from GazOntology, import relevant portions of the ontology into the knowledge base

Use higher level terms for inference

# In Clinical Tractor…



S: (v) **screen#1 (screen%2:41:01::)**, test#2 (test%2:41:01::) (test or examine for the presence of disease or infection) *"screen the blood for the HIV virus"*

**Class: disease**

**Term IRI: http://purl.obolibrary.org/obo/DOID_4**

Now we can write rules which rely on the *theme* of a *Testing* action being a *disease*.

"Screening for type 2 diabetes …"

14

# More about GazOntology

- Open source: https://github.com/oswegonlu/GazOntology

- Several options:
  - Include / exclude broad or related synonyms
  - Include only a subhierarchy of the ontology
  - Include only a single level of the ontology

- We're open to adding features! Just submit a bug report!

# Thank You!

GazOntology is available at:

https://github.com/oswegonlu/GazOntology