

Using Literature-Based Knowledge Extraction to Develop a Disease-Specific Ontology for Skin Dysbiosis



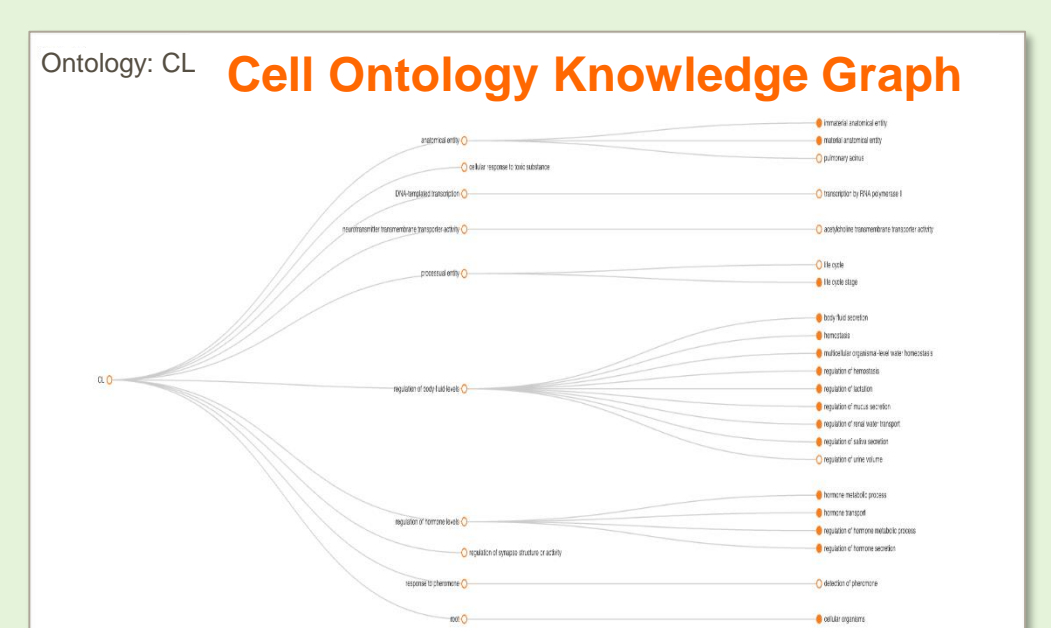
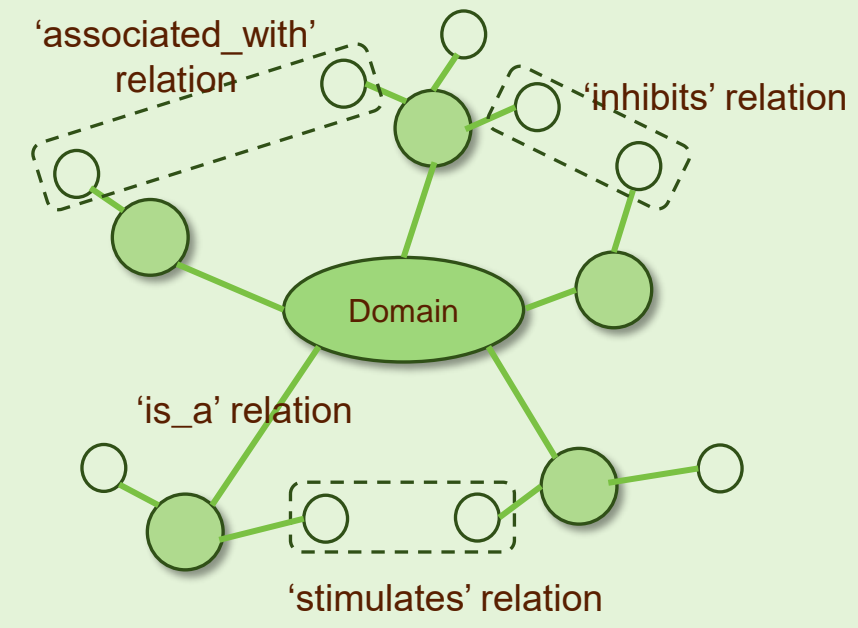
Cynthia J Grondin, Oleg Stroganov, Daniel Cooper, Yaw Nti-Addae
Rancho Biosciences, LLC, Rancho Santa Fe, CA USA

Abstract Highlights

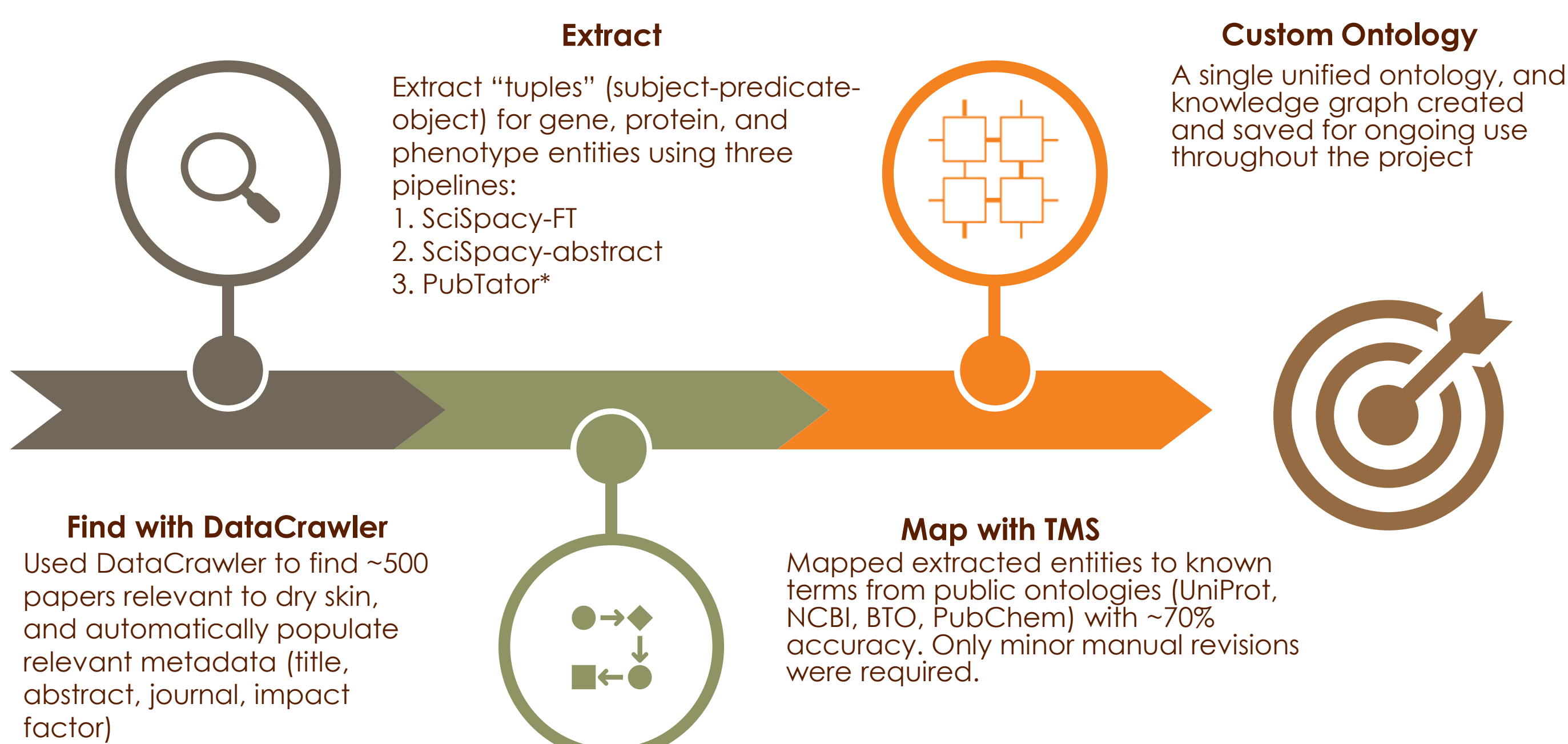
- Custom ontology needed for Skin Dysbiosis**
 - Identify relevant sources from peer-reviewed literature using Rancho's DataCrawler
 - Extract tuples from the scientific literature with 'subject-predicate-object' triplets
 - Map, clean and harmonize terms to existing ontologies using Rancho's TMS tool
 - Incorporate MESH, NCIT, GO, CL and BTO, definitions, synonyms and identifiers from PubChem, ChEBI, UniProt, and NCBI Taxonomy
 - Add new terms with uniquely generated identifiers
 - Generate Biolink predicates using OpenAI's GPT-4 model to assist relationship classification
 - Build SD-Ontology using R from cleaned terms using linked ontology hierarchies and Biolink predicates
 - n-triples (*.nt), turtle (*.ttl) and OWL (*.owl) formats using rdf conversion tool rapper and Protégé
- Custom SD-Ontology**
 - Provides insights on relationships among terms
 - Increases understanding of the pathophysiology of this disease
 - Advances strategies for improved treatment
 - Provides proof-of-concept of successful utilization of Rancho's DataCrawler and TMS tools as innovative and effective methods to identify and prioritize peer-reviewed literature from public databases and clean disparate terms to facilitate knowledge extraction and organization to advance understanding of relationships across various disciplines

Importance of Ontologies and Knowledge Graphs

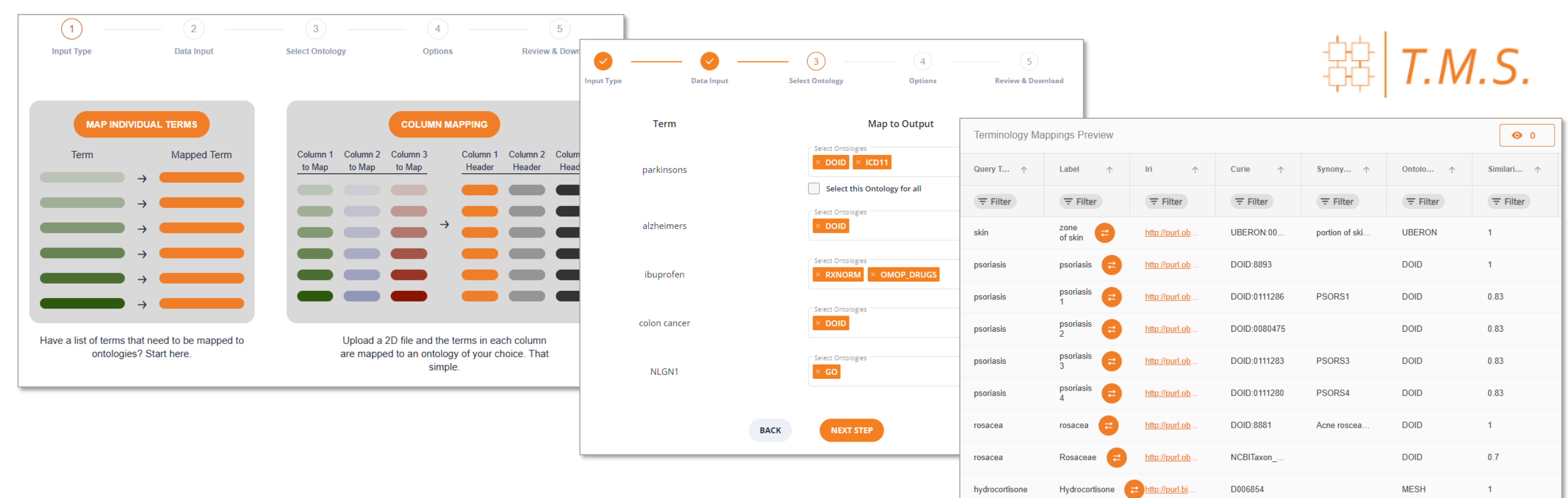
- Ontologies**
 - Knowledge classification of a domain, where the relationships between concepts are formally defined and logically related
 - Text (human readable) and Logical (machine readable) definitions
 - Hierarchical arrangement of defined terms and relationships
 - Centralizes and harmonizes data
 - Facilitates computational reasoning
 - Promotes logical inferences and sophisticated data queries
 - Facilitates knowledge extraction for therapeutics
 - Over 700 biomedical ontologies in BioPortal
 - 262 ontologies in EMBL-EBI Ontology Lookup Service (8,584,670 classes, 44,633 properties, 687,082 individuals (updated Sep 2024))
 - Language formats in RDFS, OBO, or OWL
 - Organize, Filter and Connect Data to Suggest New Relationships
- Knowledge Graphs (KGs)**
 - Coupling of hierarchical knowledge with direct relationships
 - Results in entity "Nodes" connected to relationships "Edges"
 - Stored and presented as triplet subject-->predicate-->object
 - While entities are primary information for other modalities, edges are key for KGs
 - Enables complex data and logical flow queries
 - Enables data-centric hypothesis testing



Ontology Development Workflow



3. Map to Existing Ontologies - Clean- Harmonize



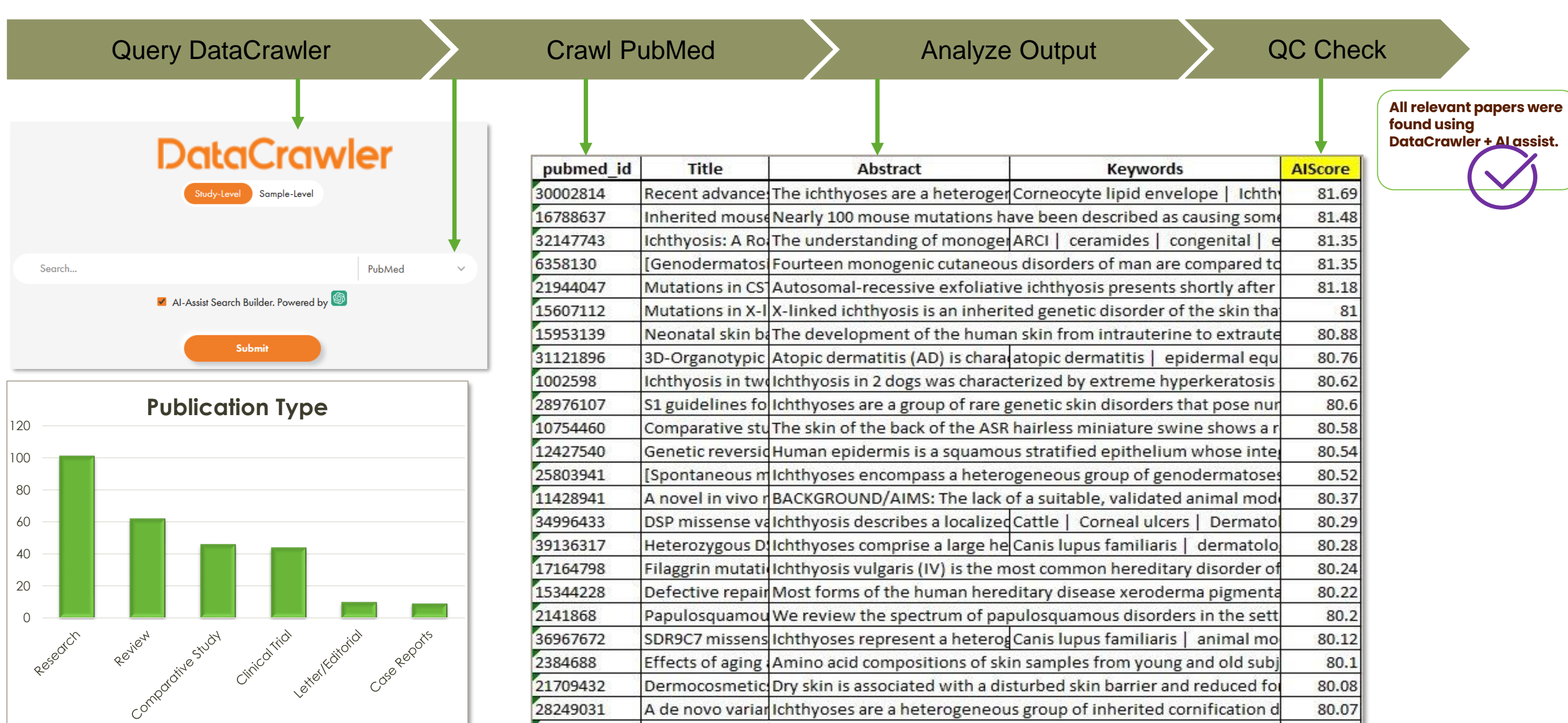
TMS Supported Ontologies and Dictionaries

Disease/Phenotype	Multi-Modal	Assay	Strain	Other
DOID HPO ICDO3 ICD10CM ICD11 MONDO ORDO	BTO EFO FMA GARD MeSH NCIT OBA OMIM OMOP SNOMEDCT UMLS MedDRA	BAO CHMO OBI	NCBITAXON RS	AFO CDISC GO HGNC LOINC PATO PR UO
		Drug CHEBI RXNORM VO	Tissue/Cell BTO CL CLO UBERON	

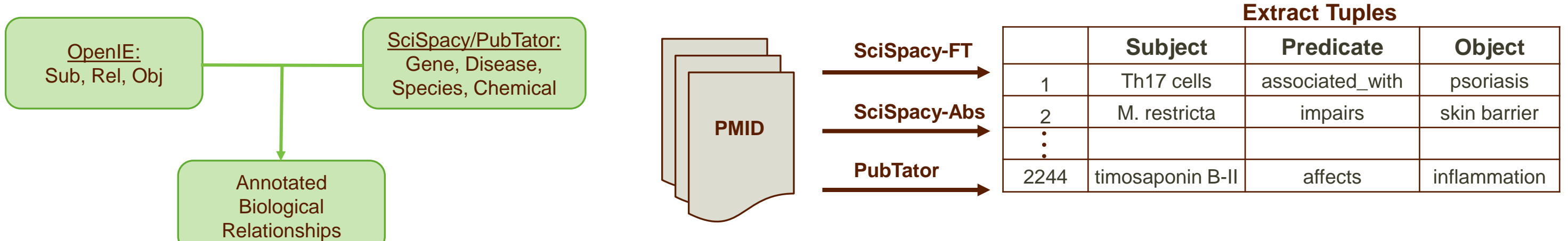
Bold = Ontologies used in Skin Dysbiosis Ontology

- Use functions and macros to clean/remove non-ascii and non-printable characters extracted from literature

1. Identify Relevant Literature



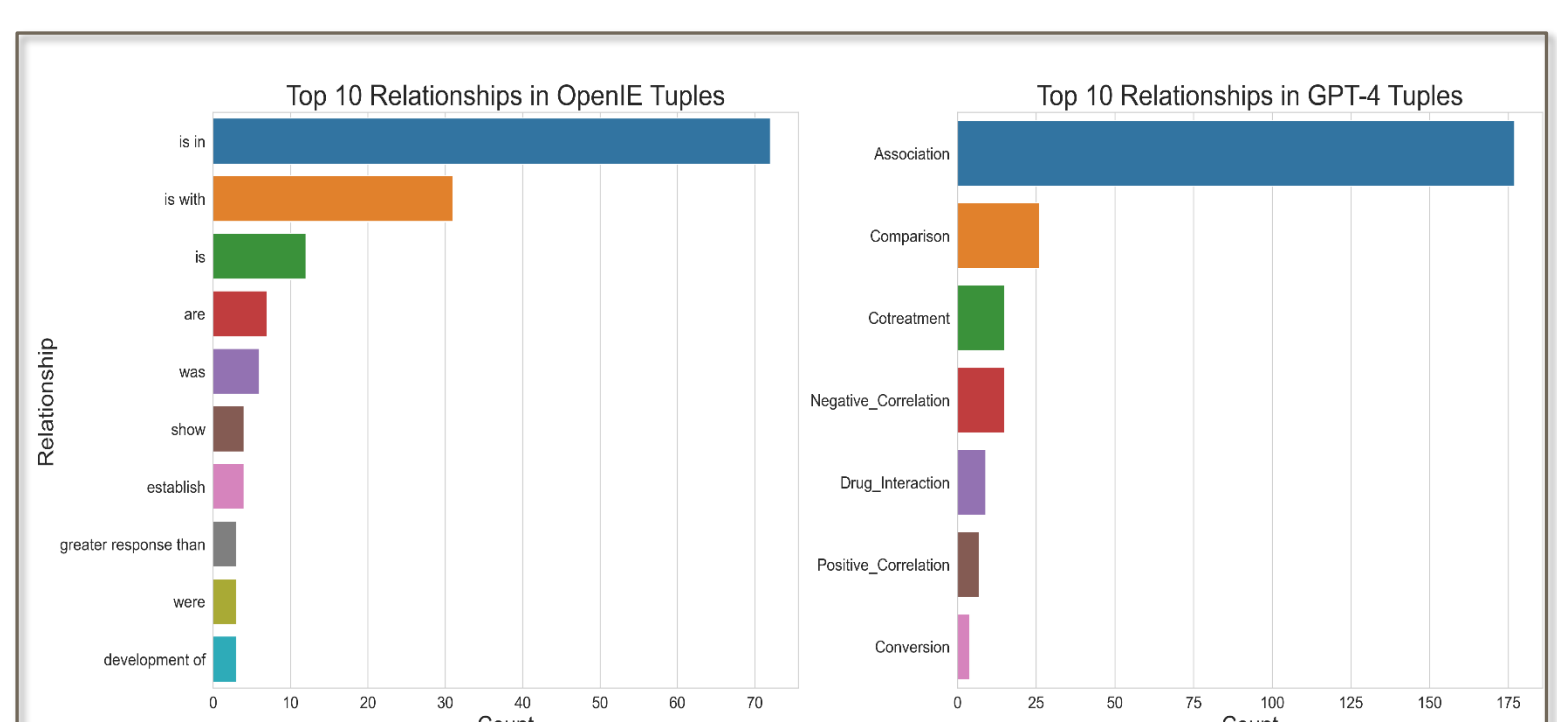
2. Extract Tuples



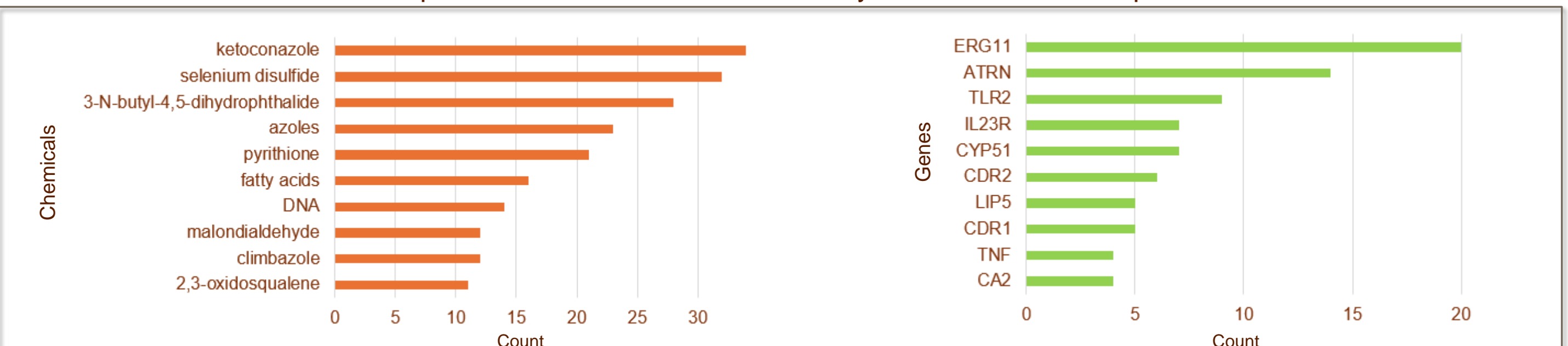
Compare OpenIE and other non-GPT based NLP tools to GPT-4 for biological/biomedical relationship extraction

- GPT-4 pipeline with OpenAI function calling and Few-Shot prompting is cautious in labeling relationships and entities as seen in the high precision
- In this task, biggest error is in recall; improve by allowing the models to extract multiple relationships per sentence, fine-tuning, prompt engineering

Metric	RE	Sub_NER	Obj_NER	NER
Precision	0.86	0.94	0.97	0.96
Recall	0.20	0.46	0.44	0.45
F1	0.32	0.62	0.61	0.62



Top Chemicals and Genes in Skin Dysbiosis Extracted Tuples



4. Build Ontology and Knowledge Graphs

- Skin Dysbiosis Ontology was built in R, with ontology represented in n-triples format (*.nt) and (*.nt) format was converted to turtle (*.ttl) and owl (*.owl) formats using rdf conversion tool rapper and Protégé UI
- Skin Dysbiosis Ontology combined standard ontologies and an additional ontology with terms which were missing in standard ontology

