# Knowledge Graph-driven Tabular Data Discovery from Scientific Documents

Vijay S. Kumar[1], Varish Mulwad[2], Jenny Weisenberg Williams[1], Tim Finin[3], Sharad Dixit[1], Anupam Joshi[3]

1: GE Research, Niskayuna, NY, USA, 2: GE Research, Bengaluru, KA, India, 3: University of Maryland, Baltimore County, Baltimore, MD, USA
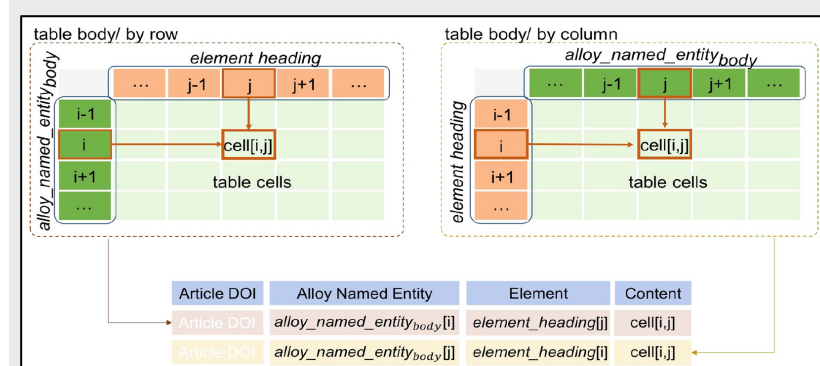
## Augmenting generative AI-driven search with valuable nuggets from tabular data sources is an emerging need

### Intelligence Report Generation & Enhancement
✓ Gather intelligence from more information sources
✓ Strengthen analysis reports with tabular data

**Incorporate Effective Visual Presentations When Feasible**

C-14. Analysts should present intelligence in a visual format to clarify an analytical conclusion and to complement or enhance the presentation of intelligence and analysis. In particular, visual presentations should be used when information or concepts, such as spatial or temporal relationships, can be conveyed better in graphic form, such as tables, flow charts, and images coupled with written text. Visual presentations

### AI-assisted Scientific Research
✓ Augment understanding & discovery from literature
  ✓ Assemble datasets in low-data domains
  (alloy materials discovery, technology forecasting)

**Data Sets and Associated Data Creation/Preparation Tools** (*NSF APTO*)

Data: e.g., aggregate historical data from lab notebooks and academic journals from 1730 to 2010 on telecommunication technologies' bandwidth, latency, and power requirements.

## Despite prevalence of tables in technical documents, limited focus on tabular data discovery in scientific domains



| Dataset | Document Type / Source | Domain | Corpus size | # tables |
|---|---|---|---|---|
| ChemTables | Patents / USPTO | Chemical | 1,000 | 788 |
| ArxivPapers | Preprints / arXiv | ML | 104,723 | 277,996 |
| ProCure (this work) | Papers & preprints / PubMed Central OA | Biomedical, clinical | 62,777 | 120,417 |

| Dataset | Downstream Task |
|---|---|
| PubTables-1M | Table detection, Table structure recognition |
| ChemTables | Table classification |
| ArxivPapers | Table extraction and segmentation |
| SciGen | Reasoning-aware table-to-text generation |
| TAT-QA | Question-answering over tables and text |
| S2abEL | Entity Linking for scientific tables |

## Scientific Tables are Hard!

### Domain-specific Entities

Similar to open data … less text, more numbers … with ranges, multi-value cells

### High Structural Heterogeneity

Row and column headers … sub-columns … abridged header cells

| Characterization | System Count | Precision | Recall |
|---|---|---|---|
| Tables with Header Rows | 113,582 | 1.00 | 0.94 |
| Tables with Header Columns | 48,733 | 1.00 | 0.55 |
| Tables with Concise Header Rows | 36,182 | 0.84 | 0.94 |
| Tables with Multi-level Header Rows | 32,169 | 1.00 | 0.97 |
| Tables with ONLY Numeric Data Cells | 12,969 | 1.00 | 0.83 |
| Tables with Concise Body | 40,158 | 0.97 | 0.67 |
| Horizontal Tables | 21,863 | 0.95 | 0.50 |
| Vertical Tables | 7205 | 0.91 | 0.62 |

### Lack of Information Reliability

**PMC7979515**: *SARS-CoV-2 Infection is Effectively Treated and Prevented by EIDD-2801*

**PPR230896**: *Efficacy and Safety of Ivermectin for Treatment and Prophylaxis of COVID-19 Pandemic*

Our automated rule-based structural characterization of 120,000+ tables shows high variability amongst scientific tables

*Read more about this work here*
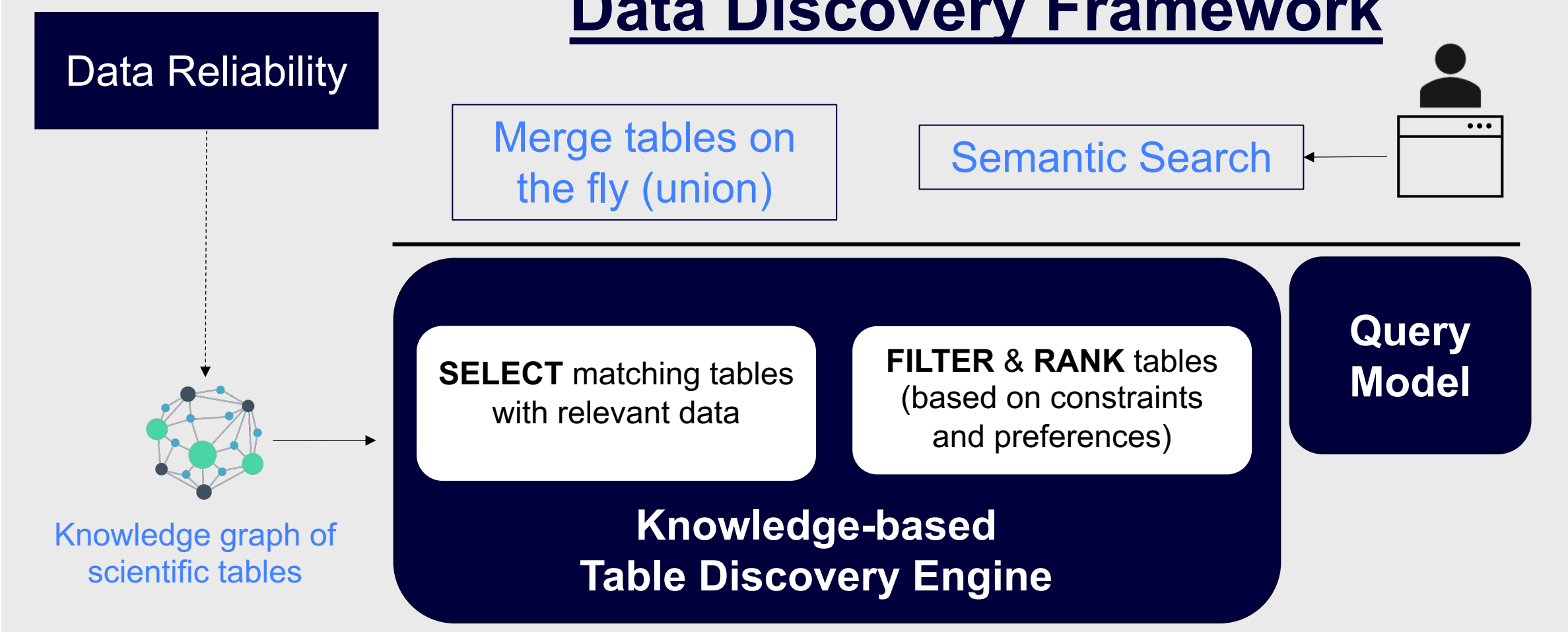
## A semantics-driven approach to tabular data discovery

Extract tables from scientific documents → Represent scientific tables as semantically annotated linked data in knowledge graph (KG) → Discover tabular data from KG under contextual constraints

Estimate information reliability and augment KG

### Data Discovery Framework

Data Reliability

Merge tables on the fly (union)

Semantic Search

**SELECT** matching tables with relevant data

**FILTER** & **RANK** tables (based on constraints and preferences)

**Query Model**

**Knowledge-based Table Discovery Engine**

Knowledge graph of scientific tables

- Synthesize (relational) tabular response to semantic search requests
- Specify diverse set of contextual constraints
- Auto-generation of SPARQL query/ies corresponding to search request
- Extend with preliminary on-the-fly table generation capability

## Table Discovery Prototype System



**ProCure Data Discovery**

Enter list of search terms / Upload file

country — Mapped to Q6256: country
vaccine — Mapped to Q134808: vaccine
trial — Mapped to Q30612: clinical trial

ProCure Search | Advanced Search | I'm Feeling Lucky | Reset

Searching for tabular objects of the form:

| Q6256 | Q134808 | Q30612 |
|---|---|---|
| ... | ... | ... |

**Reliability Metrics for Table: PMC7350246_Table_5**

| PMCID | PMC7350246 |
|---|---|
| PROVENANCE_RELIABILITY_METRIC | 0.523411 |
| PLACE_OF_ORIGIN | 1.0 |
| PUBLICATION_AVENUE | 0.046821 |

Retrieved 2 original results (0.3 seconds)
Retrieved 1 fused results (1.7 seconds)

▼ Result Constraints:
1. ☑ Table must have caption? 2. Return All types of tables 3. Time range: mm/dd/yyyy
4. Coverage constraints: 2 / 0 (Min. # of matching header cells / Min # rows in matching table)
5. Reliability constraints: 0.2 <= Rel_PROV <= 1

▼ Result Ranking Preferences:
# of matching header cells in table — Highest-first — second (Sort by / Preference order)