# COVID-19 Literature Knowledge Graph Construction and Drug Repurposing Report Generation
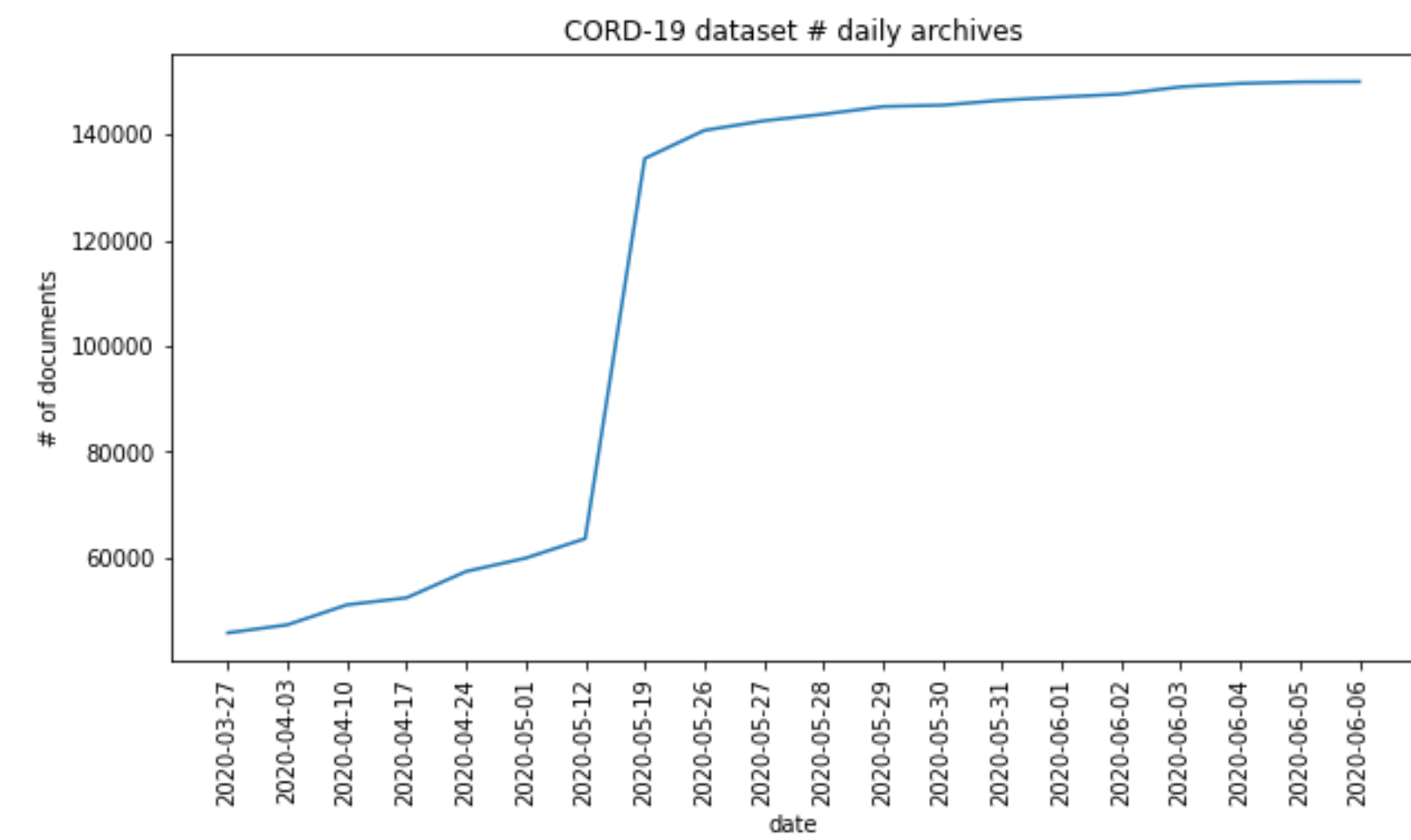
Qingyun Wang[1], Manling Li[1], Xuan Wang[1], Nikolaus Parulian[1], Guangxing Han[2], Jiawei Ma[2], Jingxuan Tu[3], Ying Lin[1], Haoran Zhang[1], Weili Liu[1], Aabhas Chauhan[1], Yingjun Guan[1], Bangzheng Li[1], Ruisong Li[1], Xiangchen Song[1], Yi R. Fung[1], Heng Ji[1], Jiawei Han[1], Shih-Fu Chang[2], James Pustejovsky[3], Jasmine Rah[4], David Liem[5], Ahmed Elsayed[6], Martha Palmer[6], Clare Voss[7], Cynthia Schneider[8], Boyan Onyshkevych[9]

[1]University of Illinois at Urbana-Champaign [2]Columbia University [3]Brandeis University [4]University of Washington [5]University of California, Los Angeles [6]Colorado University [7]Army Research Lab [8]QS2 [9]Defense Advanced Research Projects Agency

hengji@illinois.edu, hanj@illinois.edu, sc250@columbia.edu
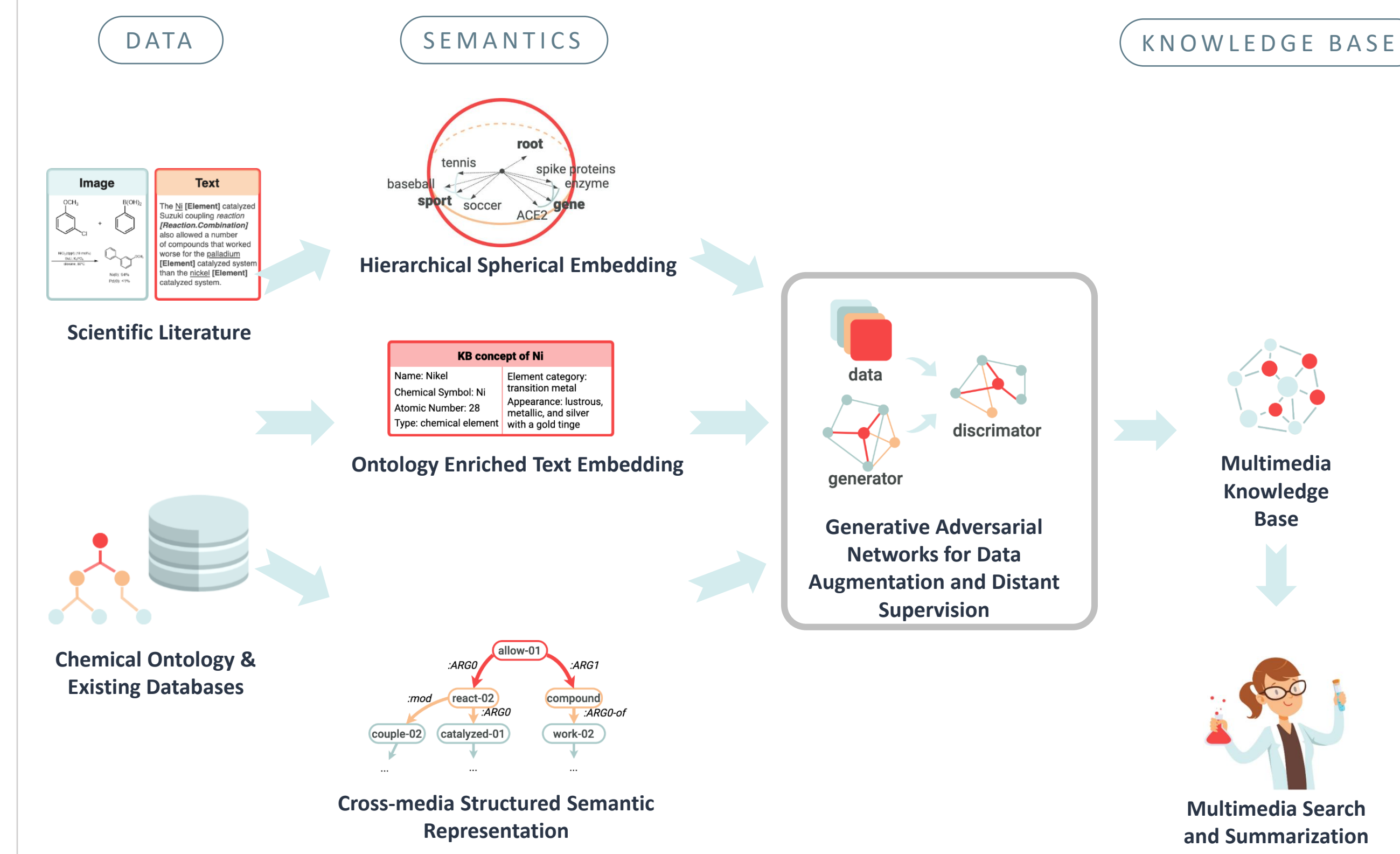
## Motivation

- **Quantity**
  - More than **140K** paper are published about coronavirus by June 13, 2020.
- **Quality**
  - Many research results are **redundant**, **complementary** or even **conflicting** with each other
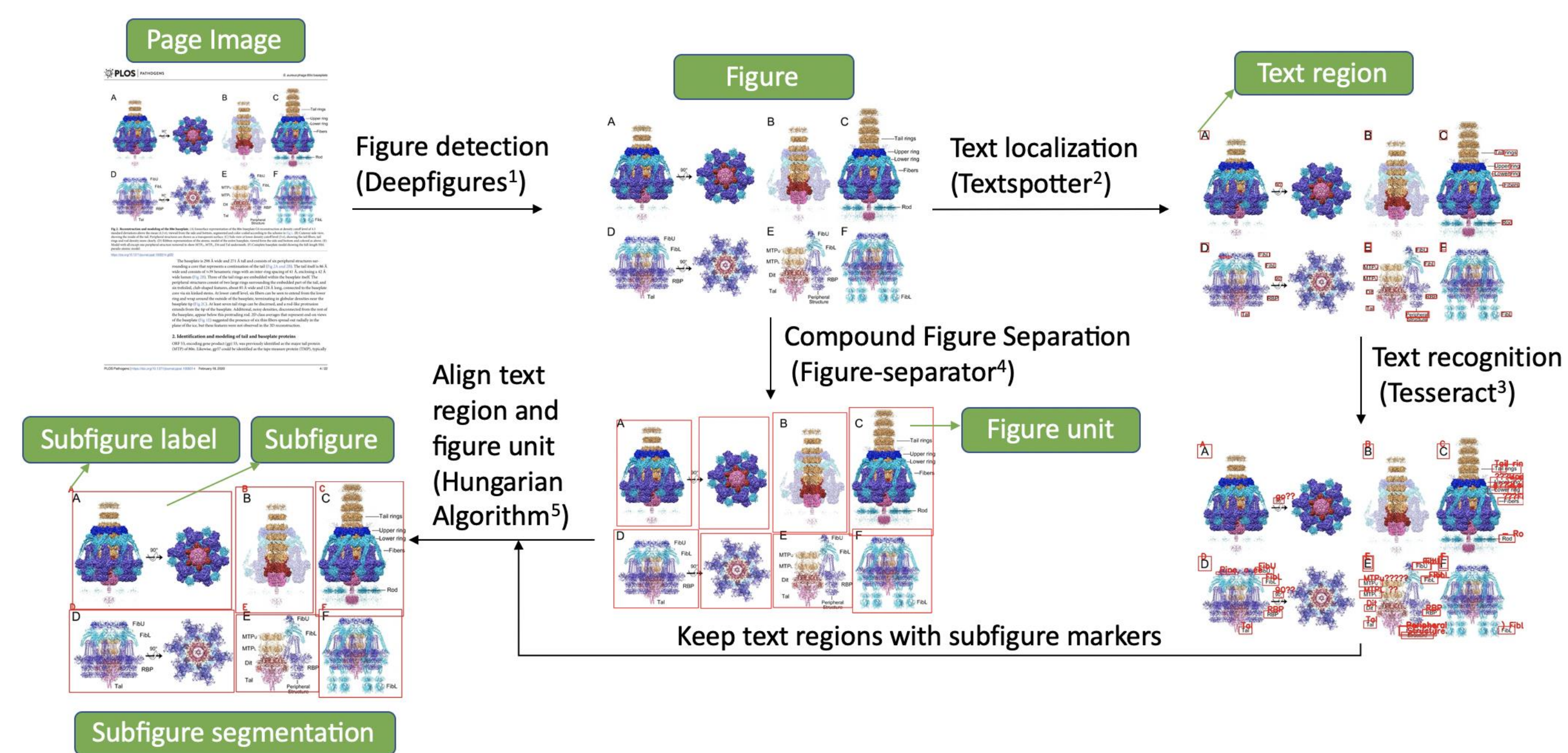


## Our Goals



## Coarse-grained Text Extraction

- **Entity Extraction + Entity Linking**
  - Extract entities from unstructured texts, link entity mentions to external biomedical ontologies including **Comparative Toxicogenomics Database (CTD)** and obtain **Medical Subject Headings (MeSH) IDs**
- **Relation Extraction**
  - Extract 133 relation types including *Gene–Chemical–Interaction Relationships, Chemical–Disease Associations, Gene–Disease Associations, Chemical–GO Enrichment Associations* and *Chemical–Pathway Enrichment Associations*
- **Event Extraction**
  - Extract 13 Event types and the roles of entities involved in these events, including *Gene expression, Transcription, Localization, Protein catabolism, Binding, Protein modification, Phosphorylation, Ubiquitination, Acetylation, Deacetylation, Regulation, Positive regulation,* and *Negative regulation*

## Fine-grained Text Extraction

- **Fine-grained Knowledge Element**
  - Fine-grained entity extraction for **75** entity types (Xuan Wang and Jiawei Han, 2020), including many COVID-19 specific new entity types (e.g., *coronaviruses, viral proteins, evolution, materials, substrates,* and *immune responses*)
  - So we will be able to answer questions that include fine-grained knowledge elements such as "Which *amino acids* in glycoprotein (a spike protein of COVID-19) are most related to Glycan (CHEMICAL)?"
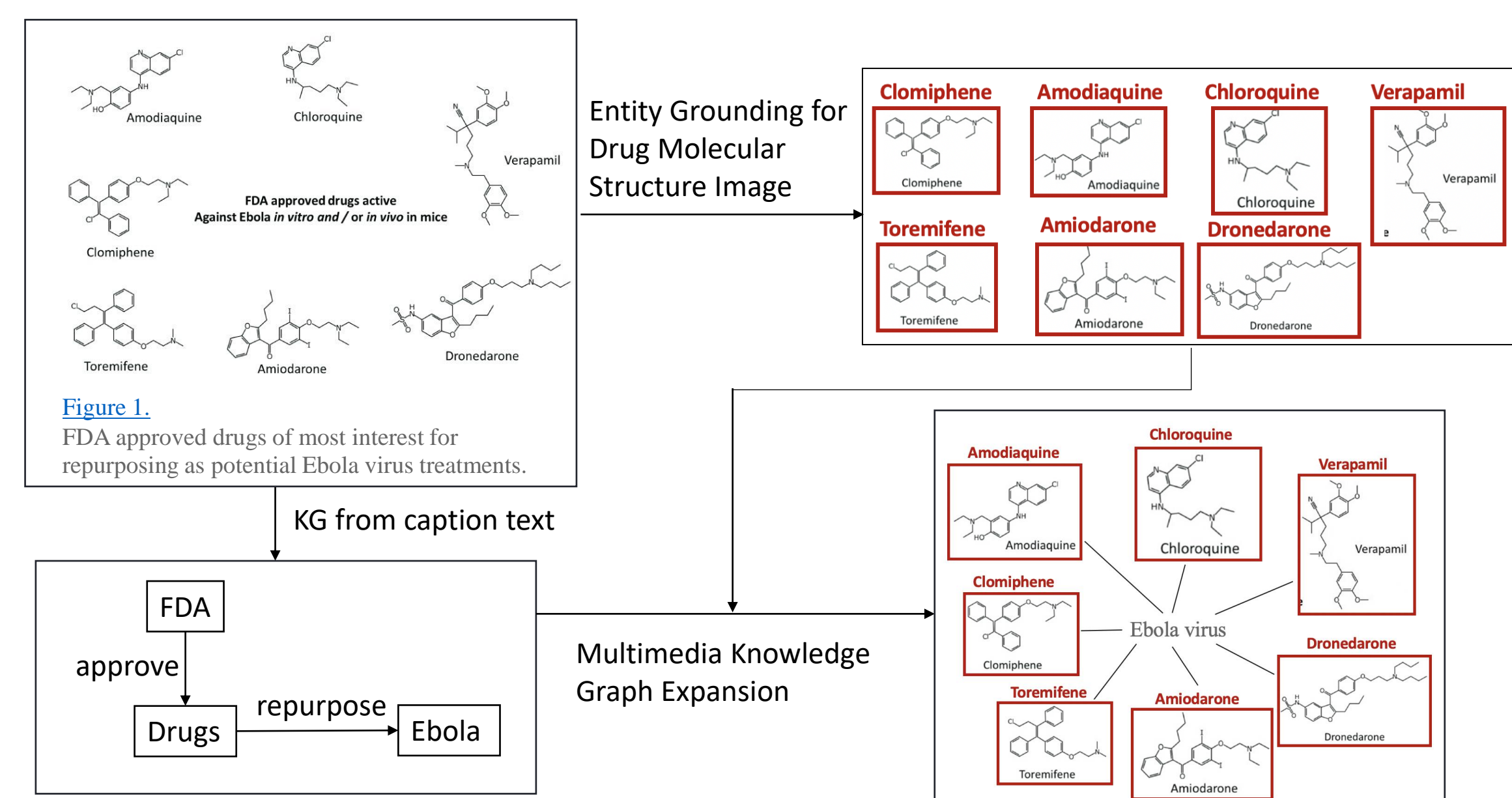


## Image Processing and Cross-media Entity Grounding

- **Automatic Figure Extraction and Subfigure Segmentation**
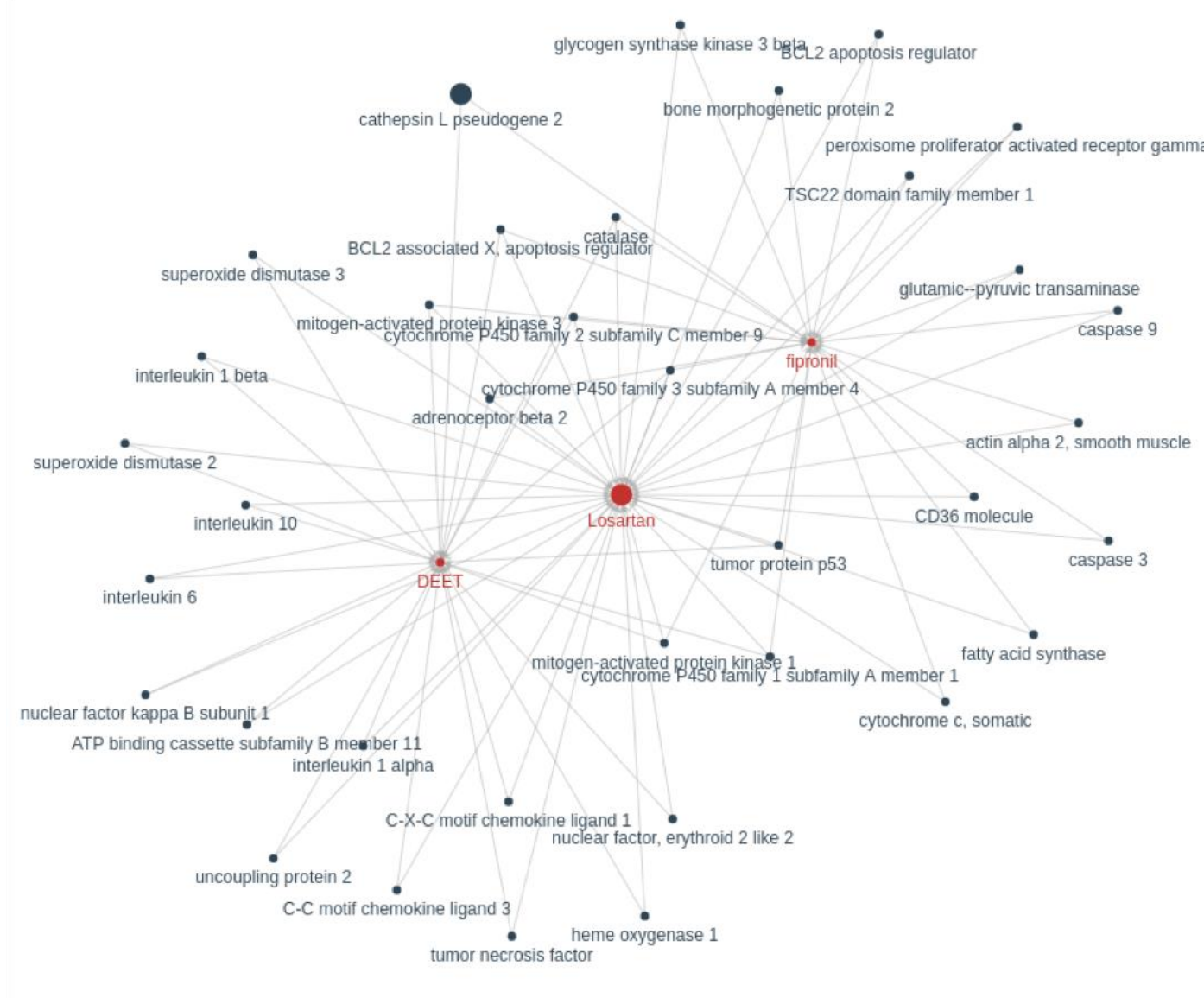  - The figure shown here is from (Kizziah et al., 2020)



- **Expanding KG through Subfigure Segmentation and Cross-modal Entity Grounding**
  - The figure shown here is from (Ekins and Coffee, 2015)



Figure 1. FDA approved drugs of most interest for repurposing as potential Ebola virus treatments.
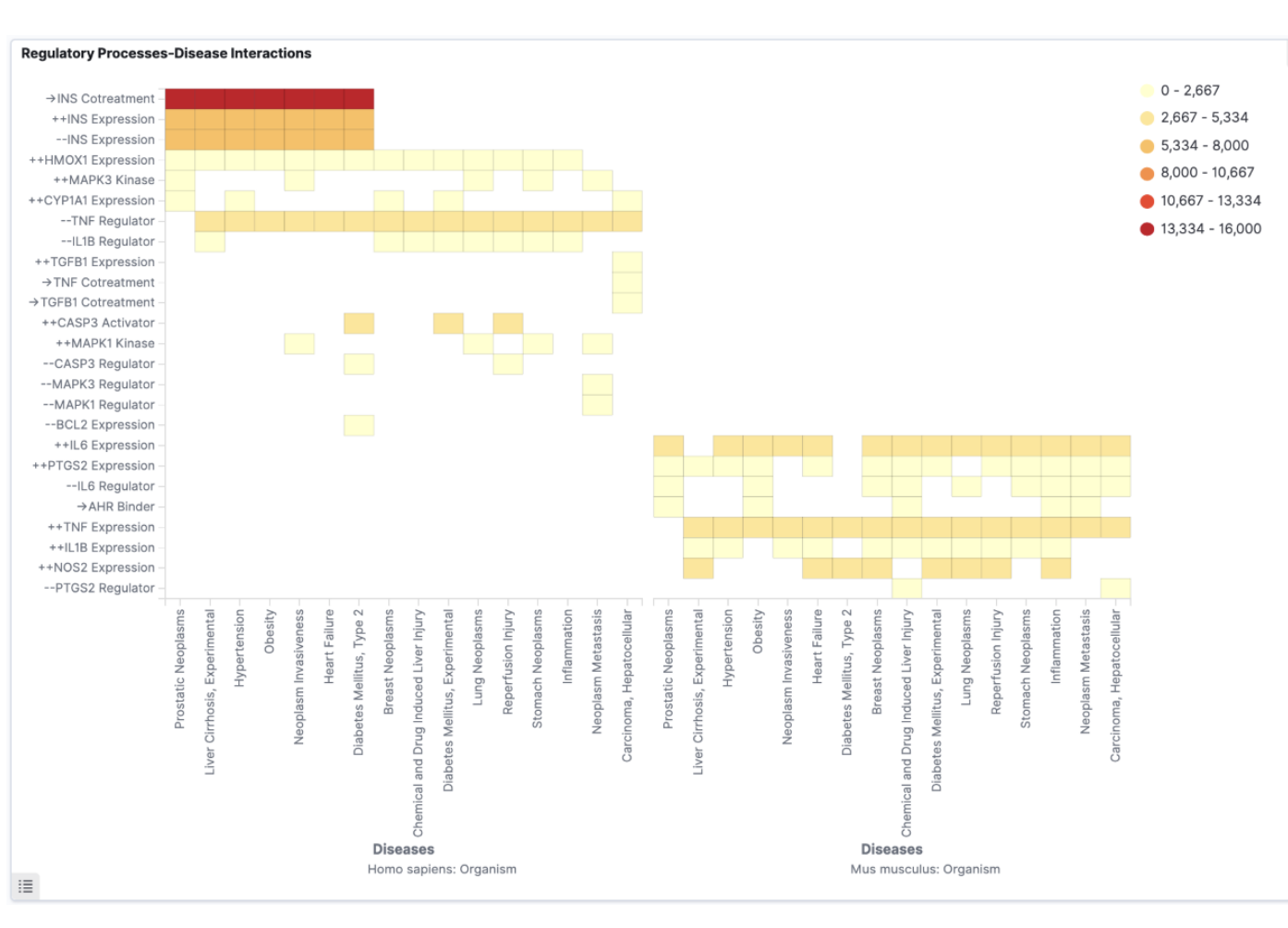
## KG Visualization

- **Constructed KG Connecting Losartan and Cathepsin L pseudogene2**
  - Where **red** nodes represent *chemicals*, **grey** nodes represent *genes*, and edges represent gene-chemical relations



- **Regulatory Processes-Disease Interactions Heatmap**



## Knowledge-driven Question Answering

- **Limitations of State-of-the-art Question Answering**
  - Fully rely on *word-level* or *sentence-level* semantic meaning matching
  - Questions are limited to non-experts (e.g., "Corona Virus Update?") or too high-level (e.g., "What is known about transmission, incubation, and environmental stability?")
- **What We Need**
  - Install a scientific brain (**KG**) for QA
  - Preliminary Results

| Question | # of Answers | Example Answers |
|---|---|---|
| Which genes are related to COVID-19? | 687 | AP2 associated kinase 1, myeloperoxidase, thioredoxin |
| Which chemicals are related to COVID-19? | 3,142 | acetoacetic acid, Chlorine, Zymosan |
| Which diseases are the most similar to COVID-19? | 4 | Enteritis, Transmissible, of Turkeys; Feline Infectious Peritonitis; Gastroenteritis, Transmissible, of Swine; Severe Acute Respiratory Syndrome |
| Which genes are related to COVID-19 that can be transferred from its similar diseases? | 2,168 | DEK proto-oncogene, nuclear receptor corepressor 1 |
| Which chemicals are related to COVID-19 that can be transferred from its similar diseases? | 327 | Ampicillin, Quercetin, Zoledronic Acid |

- **EvidenceMiner** with Query: "CORONAVIRUS cause DISEASEORSYNDROME"



## Case Study

- **KG Statistics**
  - **50,864** Gene nodes, **7,230** Disease nodes, **9,123** Chemical nodes, **1,725,518** chemical-gene links, **5,556,670** chemical-disease links, and **7,7844,574** gene-disease links
- **Sample Questions and Answers**
  - *Current indication: what is the drug class? What is it currently approved to treat?*
    - Results for Benazepril



  - *Was the drug identified by manual or computation screen?*