# Scientific Text Mining and Knowledge Graphs

**Meng Jiang** and **Jingbo Shang**

University of Notre Dame    University of California, San Diego

mjiang2@nd.edu    jshang@ucsd.edu

# Mining Knowledge from **Big Data**



Big Data

Structured Knowledge & Insights

# Massive Unstructured Text Data



News

Social Media

Business & Finance

Scientific Papers

Medical Records

......

# Massive Unstructured Text Data



News

Social Media
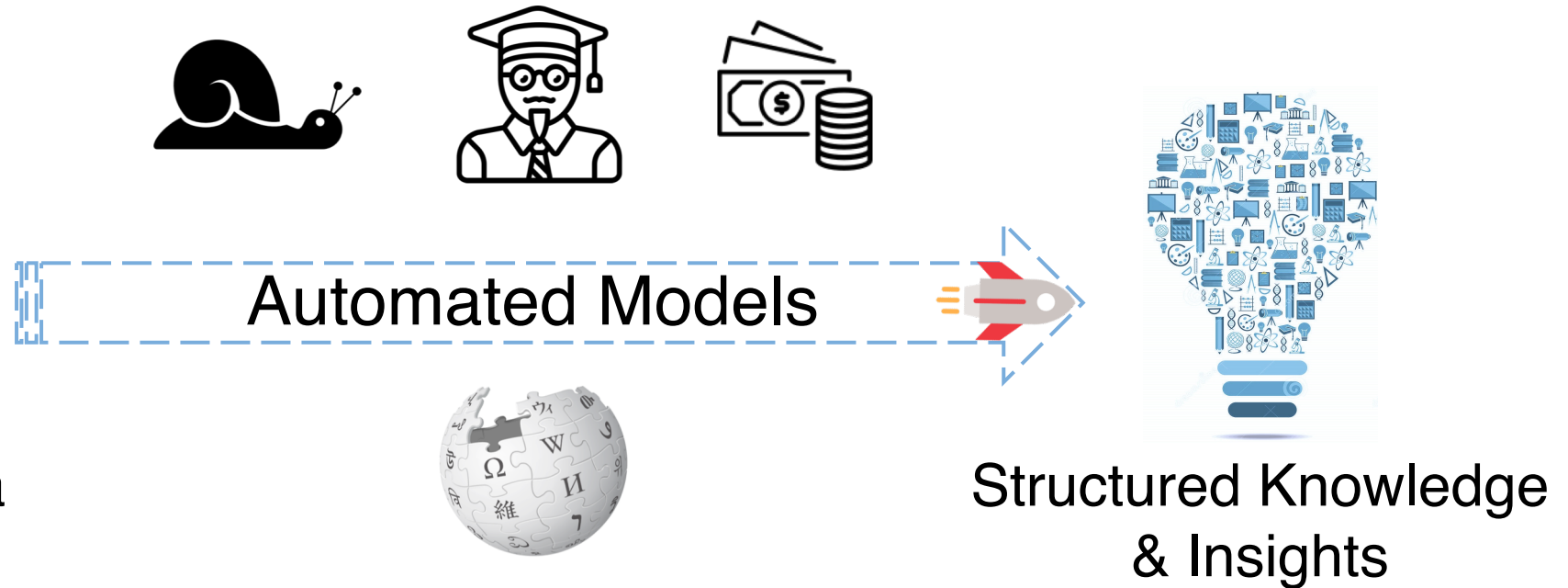
Business & Finance

Scientific Papers

Medical Records

4

# **Goal**: Texts → Knowledge & Insights



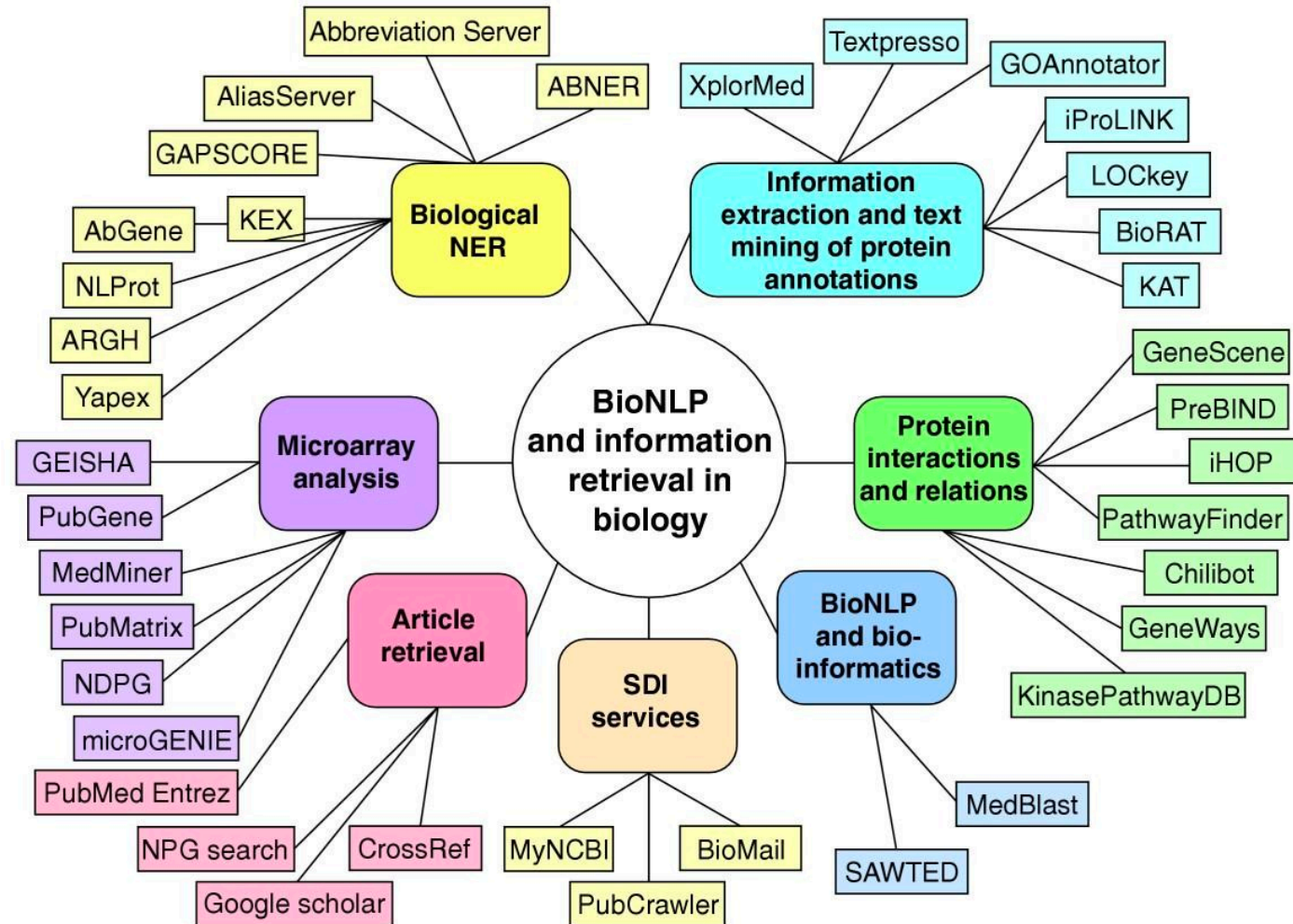❑ **Traditional methods** rely on ***extensive annotations from domain experts***

Automated Models

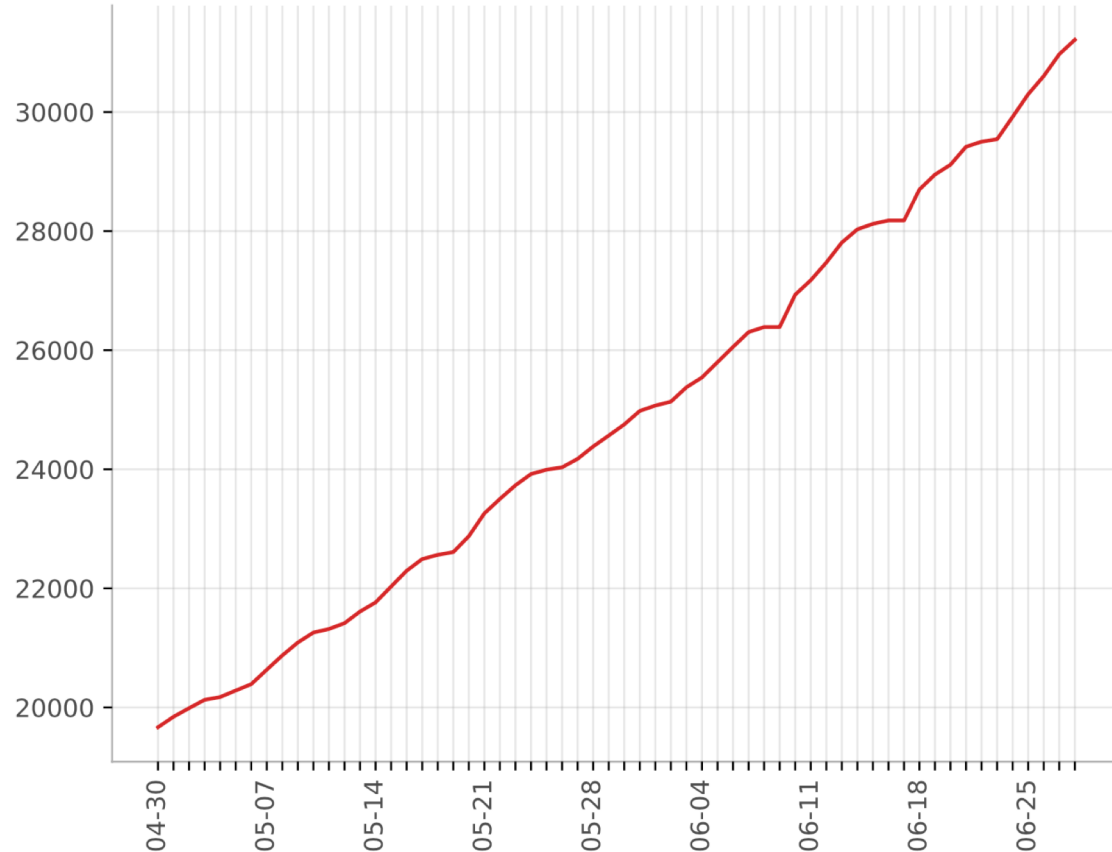Unstructured Text Data

Structured Knowledge & Insights

❑ **This tutorial** focuses on an ***automatic*** way

❑ "Automatic" – using public knowledge bases only

# An Example: BioNLP

# COVID-19 Research



**The Growing Number of COVID-19 Papers at PubMed**

Angiotensin-converting enzyme 2 **GENE_OR_GENOME** ( ACE2 **GENE_OR_GENOME** ) as a SARS-CoV-2 **CORONAVIRUS** receptor: molecular mechanisms and potential therapeutic target. SARS-CoV-2 **CORONAVIRUS** has been sequenced [3]. A phylogenetic **EVOLUTION** analysis [3, 4] found a bat **WILDLIFE** origin for the SARS-CoV-2 **CORONAVIRUS**. There is a diversity of possible intermediate hosts for SARS-CoV-2 **CORONAVIRUS**, including pangolins **WILDLIFE**, but not mice **EUKARYOTE** and rats **EUKARYOTE** [5]. There are many similarities of SARS-CoV-2 **CORONAVIRUS** with the original SARS-CoV **CORONAVIRUS**. Using computer modeling, Xu *et al*. [6] found that the spike proteins **GENE_OR_GENOME** of SARS-CoV-2 **CORONAVIRUS** and SARS-CoV **CORONAVIRUS** have almost identical 3-D structures in the receptor binding domain that maintains Van der Waals forces **PHYSICAL_SCIENCE**. SARS-CoV spike proteins **GENE_OR_GENOME** has a strong binding affinity to human ACE2 **GENE_OR_GENOME**, based on biochemical interaction studies and crystal structure analysis [7]. SARS-CoV-2 **CORONAVIRUS** and SARS-CoV spike proteins **GENE_OR_GENOME** share identity in amino acid sequences and ……

**Scientific Named Entity Recognition and Typing**

Wang et al. "COVID-19 Literature Knowledge Graph Construction and Drug Repurposing Report Generation"

# Knowledge Graphs in COVID-19 Research

# Drug Repurposing

# Two Chapters

- Chapter 1: Mining Structures from Scientific Text
  - Phrase mining
  - Concept recognition (Named entity recognition) — Shang
  - Language models
  - Relation and attribute extraction
  - Conditional statement extraction — Jiang
  - Experimental evidence extraction
- Chapter 2: Constructing and Learning Scientific Knowledge Graphs
  - Taxonomy construction — Shang
  - Knowledge graph construction
  - Learning KG for literature search — Jiang
  - Learning KG for scientific text generation