# Measuring Scholarly Impact and **Beyond**

Ying Ding
Indiana University
dingying@indiana.edu
http://info.slis.indiana.edu/~dingying/index.html

# Scholarly Communication: Challenges and Opportunities in Digital Age

Ying Ding
Indiana University
dingying@indiana.edu
http://info.slis.indiana.edu/~dingying/index.html

# Digital Age



THE SECOND MACHINE AGE

WORK, PROGRESS, AND PROSPERITY IN A TIME OF BRILLIANT TECHNOLOGIES

ERIK BRYNJOLFSSON
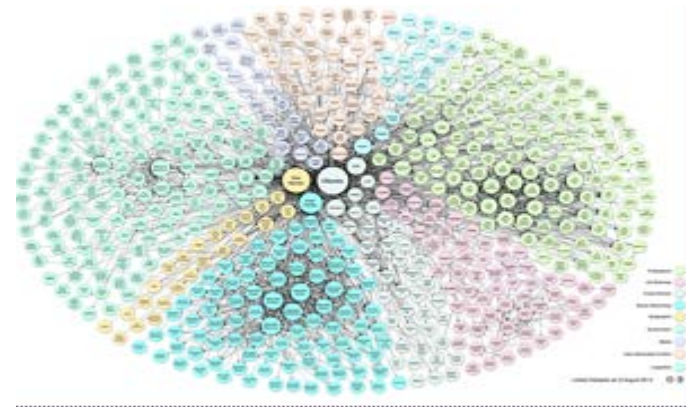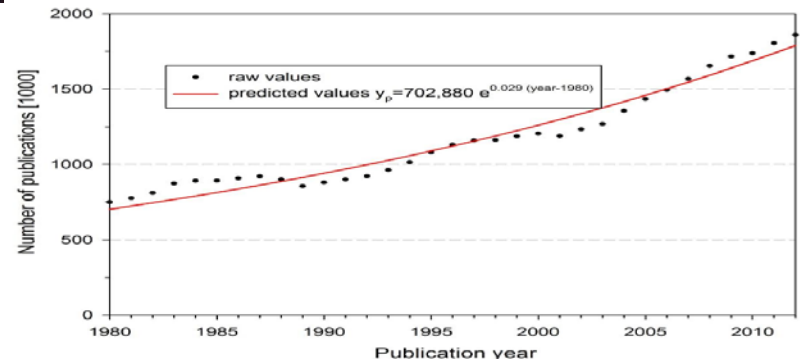ANDREW McAFEE

Google™ self-driving car

**Machine World → Human World**

# Digital Challenges



**Human World → Machine World**

# Confusing: Machine or Human World





**You and Me are currently living in such a confusing time**

# Digital Advantages

- Easy to access data
- Easy to share data (e.g., zero cost to make a copy)
- Powerful computing technologies
- Innovative minds and many eyeballs
- Motivated human capitals

- Digital Transformation:
  - Digital DNA (250M photos uploaded to Facebook daily, >5B have mobile phones)
  - Digital Bonding (we spend more time with our smartphones than with our partner, 119 minute/day)

# Rich Data and Huge Opportunities

- Smart Phone → Smart Home → Smart Robot → Smart City→ Smart … → Smart doctor



- Scholarly communication → Digital Scholarly Communication → Literature-Driven Discovery →Data-Driven Discovery→ Scientific Discovery
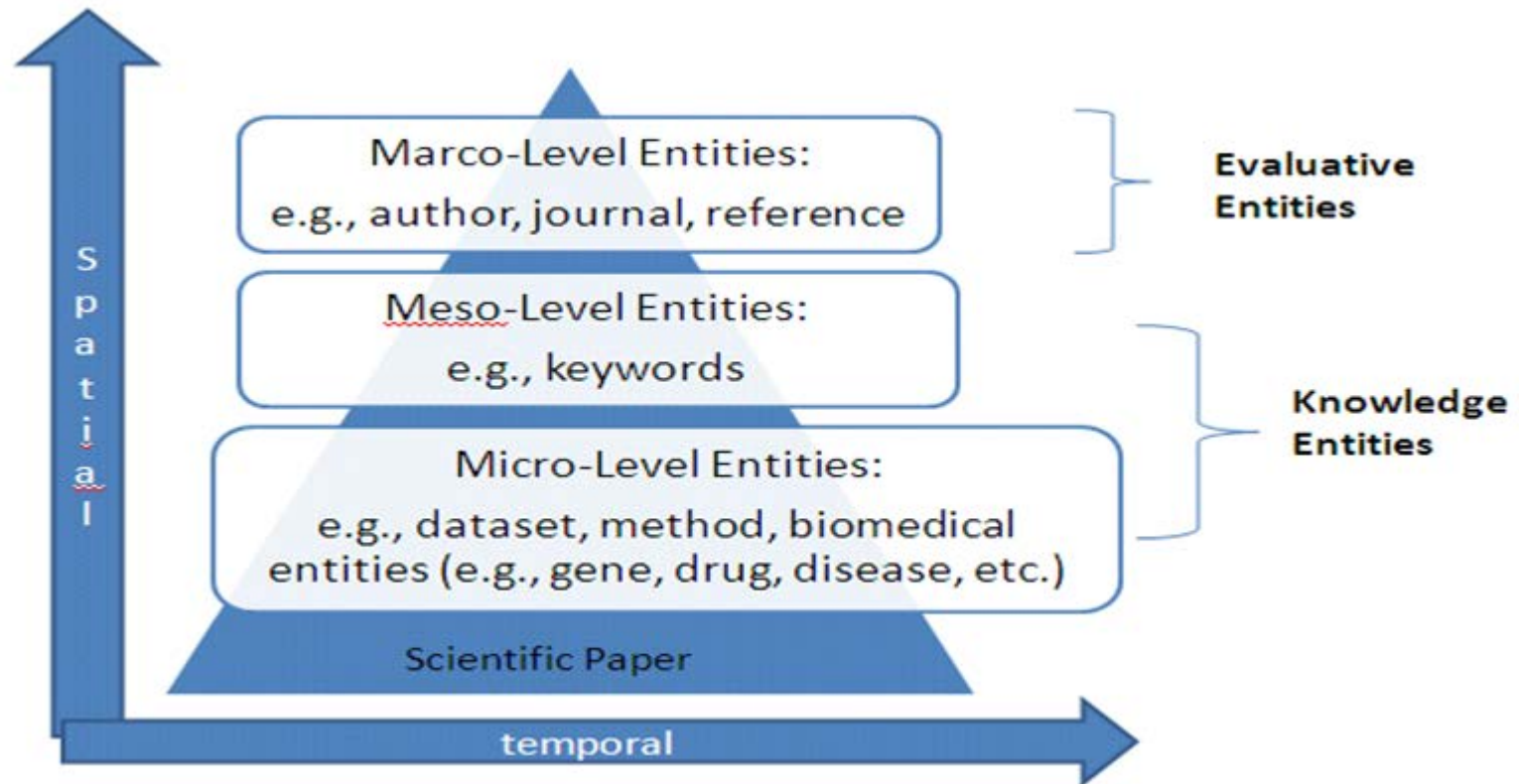
**Data2Knowledge**

# Next Generation of Scholarly Communication

- Newly developed methods allow in-depth analysis of scholarly communication
  - Topic modeling (e.g., Latent Dirichlet Allocation)
  - Information Extraction (e.g., OpenIE)
  - Social Network Analysis (e.g., Community Detection)
- Big data demonstrates the power of connected data to enable knowledge discovery
  - Structured data
  - Unstructured data
  - Social media data
- Digital age incubates transformative innovations
  - Working with domain experts (e.g., biologist, sociologist, historian)
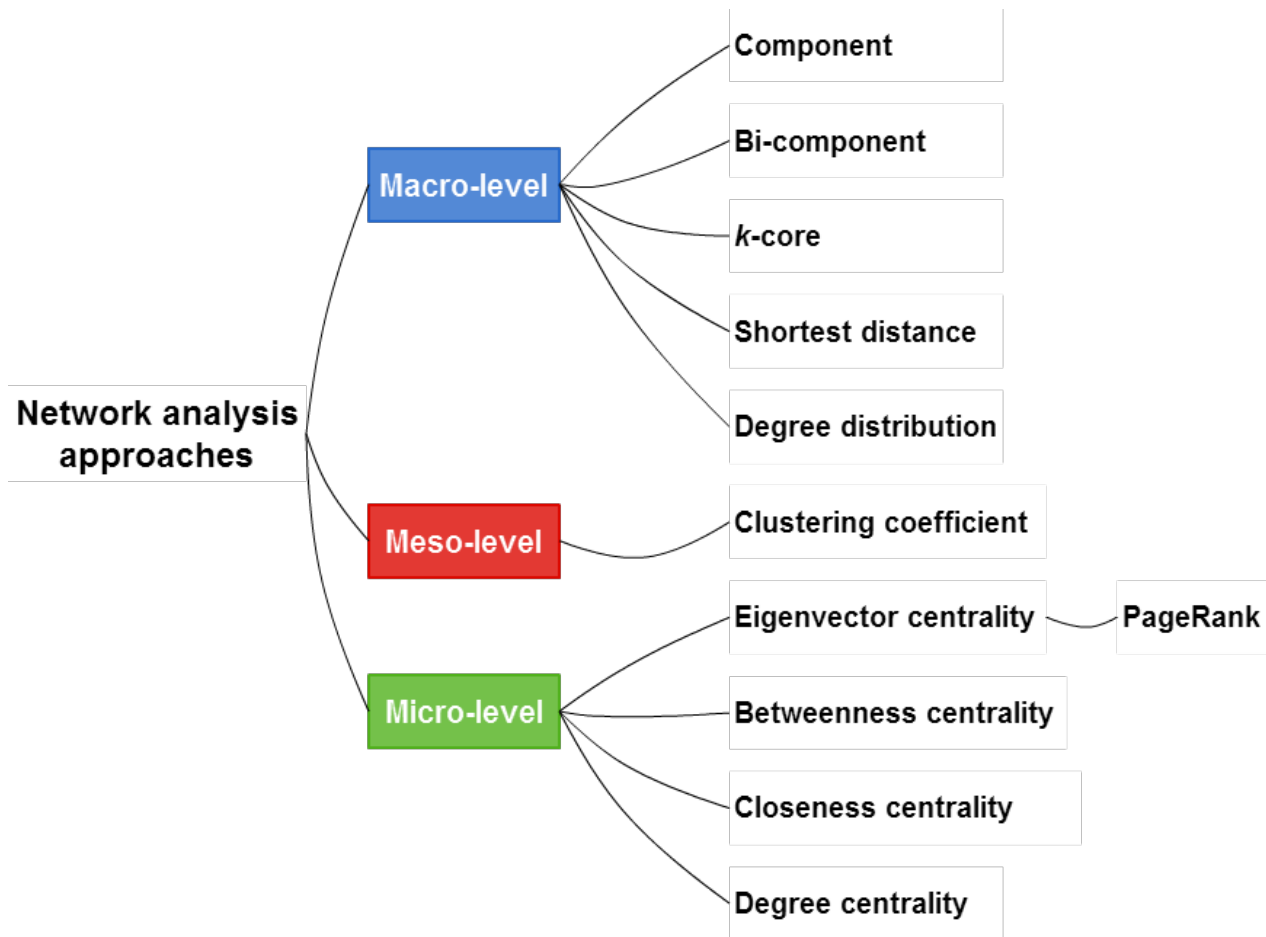  - Computational discovery in science, social science, and humanities

Ding, Y., Rousseau, R., & Wolfram, D. (Eds.) (2014). *Measuring scholarly impact: Methods and practice*. Springer.

# EntityMetrics

**Entitymetrics is defined as using entities (i.e., evaluative entities or knowledge entities) in the measurement of impact, knowledge usage, and knowledge transfer, to facilitate knowledge discovery.**



Ding, Y., Song, M., Han, J., Yu, Q., Yan, E., Lin, L., & Chambers, T. (2013). Entitymetrics: Measuring the impact of entities. PLoS One, 8(8): 1-14.

# EntityMetrics

# PubMed Entities

Drug
Disease
Protein
Pathway
Gene

## Antidiabetic drug metformin induces apoptosis in human MCF breast cancer via targeting ERK signaling.

Malki A, Youssef A.

Biochemistry Department, Faculty of Science, Alexandria University, Alexandria, Egypt. amalky@yahoo.com

### Abstract

Metformin is the most widely used antidiabetic drug for type II diabetes in the world. Recent studies provide clues that the use of metformin may be associated with reduced incidence and improved prognosis of certain cancers, and there is increasing evidence of a potential efficacy of this agent as an anticancer drug. This observation led us to hypothesize that metformin might inhibit human breast cancer cells (MCF-7) growth. Here, we report that metformin induced apoptosis in human breast carcinoma cell lines MCF-7 cells via novel signaling pathway. Metformin induced apoptosis by arresting cells in G1 phase and reducing cyclin D level and increasing the expression of p21 and cyclin E. Molecular and cellular studies indicated that metformin significantly elevated p53 and Bax levels and reduced STAT3 and Bcl-2. Inhibitors of signaling proteins were used to study the mechanism(s) of metformin function. Receptor inhibitor studies indicated that p53 activation was mediated through insulin receptor (IR), not insulin... inhibit... SAPK... All the... has no... strateg

Cite

## Obesity and insulin resistance in breast cancer--chemoprevention strategies with a focus on metformin.

Goodwin PJ, Stambolic V.

Department of Medicine, Division of Clinical Epidemiology at the Samuel Lunenfeld Research Institute, Mount Sinai Hospital, Princess Margaret Hospital, University of Toronto; Mount Sinai Hospital, 1284-600 University Avenue, Toronto, Ontario M5G 1X5, Canada. pgoodwin@mtsinai.on.ca
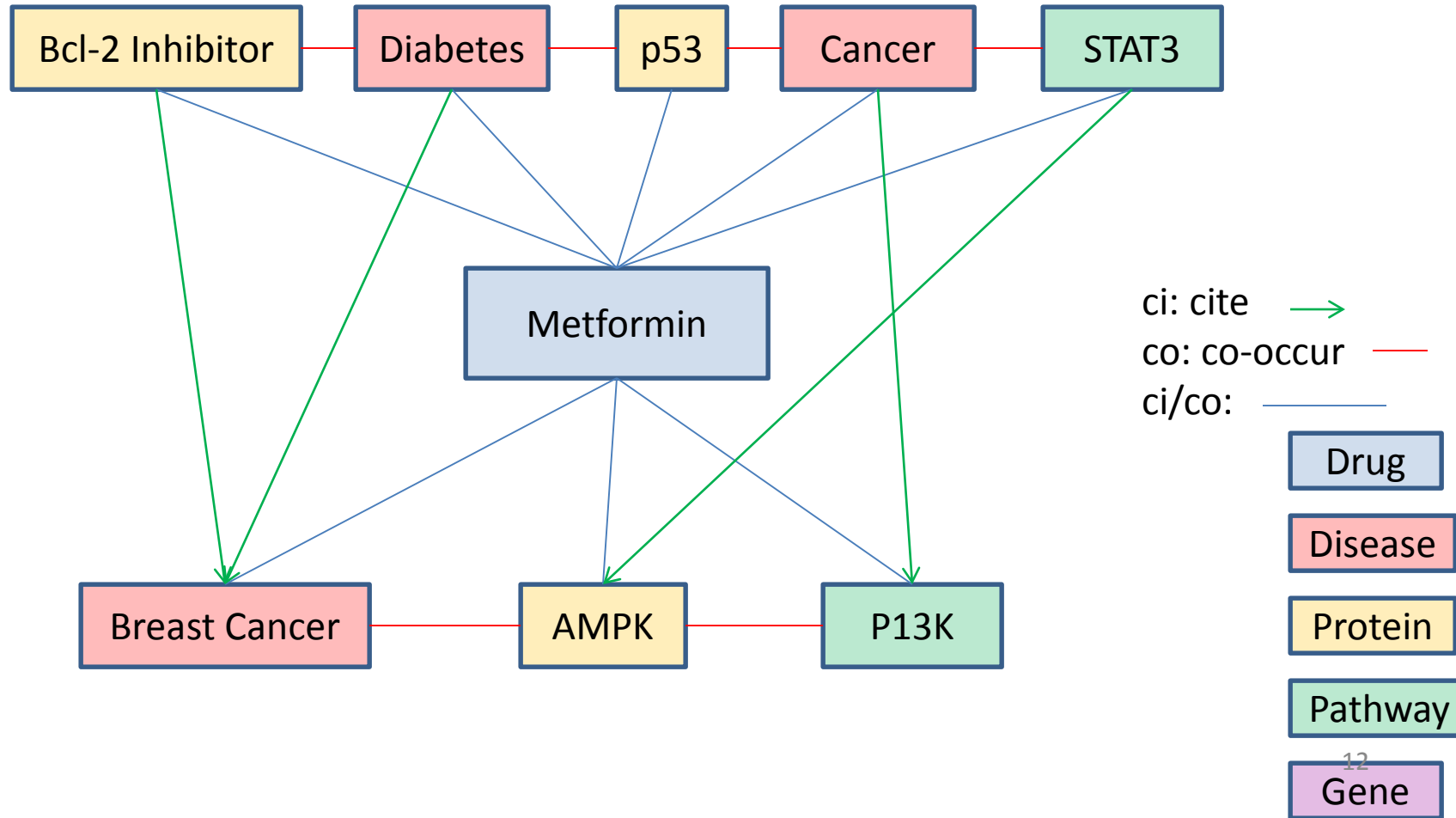
### Erratum in

### Abstract

Obesity and insulin resistance have been associated with breast cancer risk, and breast cancer outcomes. Recent research has focused on insulin as a potential biologic mediator of these effects given frequent expression of insulin/IGF-1 receptors on breast cancer cells which, when activated, can stimulate signaling through PI3K and Ras-Raf signaling pathways to enhance proliferation. Metformin, a commonly used diabetes drug, lowers insulin in non-breast diabetic cancer patients, likely by reducing hepatic gluconeogenesis; it also appears to have potential insulin independent direct effects on tumor cells which are mediated by activation of AMPK with downstream inhibition of mTOR. There is growing epidemiologic, clinical and preclinical (in vitro and in vivo) evidence in keeping with anticancer effects of metformin in breast and other cancers. This has led to the hypothesis that metformin may be effective in breast cancer prevention and treatment. Clinical studies in the neoadjuvant and adjuvant settings are ongoing; additional Phase 2 trials in the metastatic setting and proof of principle studies in the prevention setting are planned.

11

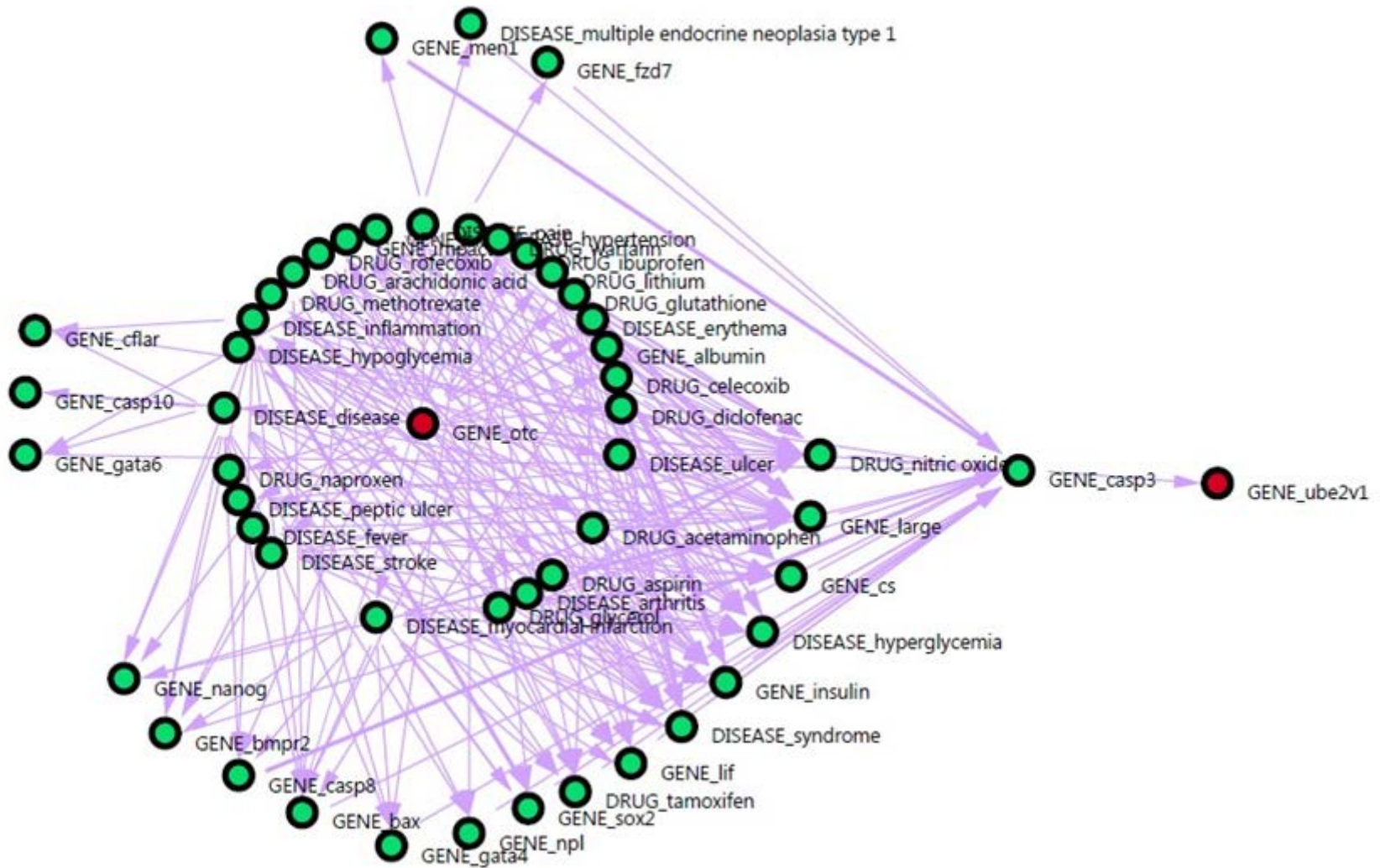# Entity Graph

- Heterogeneous Entity Graph

# Metformin related entity-entity citation network
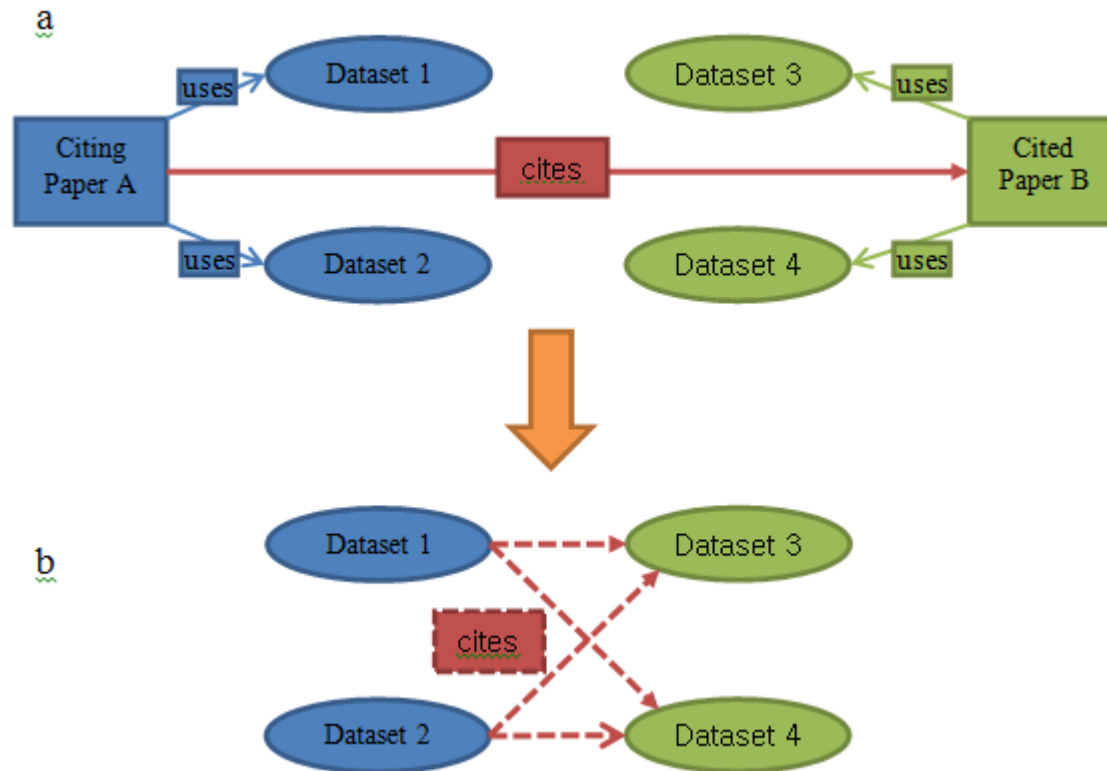
Table 4 In-degree centrality (top 20)

| Rank | Disease | Drug | Gene | All Entities |
|---|---|---|---|---|
| 1 | DISEASE_disease | DRUG_glycerol | GENE_large | DISEASE_disease |
| 2 | DISEASE_erythema | DRUG_arachidonic acid | GENE_insulin | DRUG_glycerol |
| 3 | DISEASE_syndrome | DRUG_calcium | GENE_impact | DISEASE_erythema |
| 4 | DISEASE_death | DRUG_cholesterol | GENE_set | DRUG_arachidonic acid |
| 5 | DISEASE_hypertension | DRUG_nitric oxide | GENE_tnf | GENE_large |
| 6 | DISEASE_obesity | DRUG_potassium | GENE_lep | DISEASE_syndrome |
| 7 | DISEASE_inflammation | DRUG_glutathione | GENE_hr | DISEASE_death |
| 8 | DISEASE_diabetes mellitus | DRUG_ester | GENE_ca2 | GENE_insulin |
| 9 | DISEASE_necrosis | DRUG_dexamethasone | GENE_camp | GENE_impact |
| 10 | DISEASE_insulin resistance | DRUG_norepinephrine | GENE_met | DISEASE_hypertension |

Data: 4,770 articles retrieved from PubMed Central with 134,844 references, and 1,969 bio-entities (i.e., 880 genes, 376 drugs, and 713 diseases)

# Metformin related entity-entity citation network

# Dataset



Yu, Q., **Ding, Y.**, Song, M., Song, S., Liu, J., & Zhang, B. (forthcoming). Tracing database usage: Detecting main paths in database link networks. *Journal of Informetrics*.

# Main Path

PLACE is a database of motifs found in plant cis-acting regulatory DNA elements
PlantCARE is a database of plant cis- acting regulatory elements, enhancers and
repressors. The motifs are collected in this database as well,

# Entity Citation Network vs. Entity Co-Occurrence Network

- Gene Gene Co-Occurrence Network (GG) vs. Gene Cite Gene Network (GCG)
  - The GCG network shares many genes with the GG network and as a result is a competitive complement to the GG network
  - Using gene relationships based on citation relation extends the assumption of gene interaction being limited to the same article and opens up a new opportunity to analyze gene interaction from a wider spectrum of datasets.
  - 1,149 gene pairs from GCG were found in GG. A total of 164 pairs out of 1,149 were not found in GG before 2005, but were found in GCG before 2005. In particular, the PARK2 and PINK1 gene pair ranks fifth by co-occurrence frequency in the GG network, implying the gene pair has highly been studied since 2005

Song, M., Han, N., Kim, Y., Ding, Y., & Chambers, T. (2013). Discovering implicit entity relation with the gene-citation-gene network. *PLoS One*, 8(12), e84639

# Big Data in Life Sciences

- There is now an incredibly **rich resource of public information** relating compounds, targets, genes, pathways, and diseases. Just for starters there is in the public domain information on:
    - **69 million compounds** and **449,392 bioassays** (PubChem)
    - **59 million compound bioactivities** (PubChem Bioassay)
    - **4,763 drugs** (DrugBank)
    - **9 million protein sequences** (SwissProt) and 58,000 3D structures (PDB)
    - **14 million human nucleotide sequences** (EMBL)
    - **22 million life sciences publications** - 800,000 new each year (PubMed)
    - Multitude of other sets (drugs, toxicogenomics, chemogenomics, metagenomics …)

- Even more important are the **relationships between these entities.** For example a chemical compound can be linked to a gene or a protein target in a multitude of ways:
    - Biological assay with percent inhibition, IC50, etc
    - Crystal structure of ligand/protein complex
    - Co-occurrence in a paper abstract
    - Computational experiment (docking, predictive model)
    - Statistical relationship
    - System association (e.g. involved in same pathways cellular processes)

Wild, D. J., Ding, Y., Sheth, A. P., Harland, L., Gifford, E. M., & Lajiness, M. S. (2012). System chemical biology and the Semantic Web: What they mean for the future of drug discovery research. *Drug Discovery Today* (impact factor=6.422), 17(9-10), 469-474.

# Chem2Bio2RDF

- NCI Human Tumor Cell Lines Data
- PubChem Compound Database
- PubChem Bioassay Database
- PubChem Descriptions of all PubChem bioassays
- Pub3D: A similarity-searchable database of minimized 3D structures for PubChem compounds
- Drugbank
- MRTD: An implementation of the Maximum Recommended Therapeutic Dose set
- Medline: IDs of papers indexed in Medline, with SMILES of chemical structures
- ChEMBL chemogenomics database
- KEGG Ligand pathway database
- Comparative Toxicogenomics Database
- PhenoPred Data
- HuGEpedia: an encyclopedia of human genetic variation in health and disease.



**31m chemical structures**
**59m bioactivity data points**
**3m/19m publications**
**~5,000 drugs**



Chen, B., Dong, X., Jiao, Dazhi, Wang, H., Zhu, Q., Ding, Y. and Wild, D. (2010). Chem2Bio2RDF: A semantic framework for linking and mining chemogenomic and systems chemical biology data. *BMC Bioinformatics*, 2010, 11, 255.

| primary classes | description | sample instance data sources | # of sample instances |
|---|---|---|---|
| SmallMolecule | a small bioactive molecule | PubChem, ChEBI | 15509 |
| Drug | a chemical used in the treatment, cure, prevention, or diagnosis of disease | DrugBank, PharmGKB, TTD | 6544 |
| Protein | a physical entity consisting of a sequence of amino acids | Uniprot, HGNC, GOA | 12242 |
| BioAssay | an experiment to measure the effects of some substance on target, cell, or a living organism | PubChem BioAssay, ChEMBL, BindingDB, DPSP | 26861 |
| Disease | any condition that causes pain, dysfunction, distress or social problems | OMIM, DO | 8724 |
| SideEffect | undesired effect from a medicine | SIDER | 1385 |
| Literature | a scientific article | Medline | 28392 |
| Pathway | a set or series of biological interactions | KEGG, Reactome | 347 |
| Interaction — DrugDrugInteraction | a drug affects the activity of another drug | DrugBank, DCDB | 9690 |
| Interaction — ProteinProteinInteraction | two or more proteins bind together | HPRD, DIP, BioGrid | 54345 |
| Interaction — DrugInducedSideEffect | a drug interaction that results in side effect | SIDER | 61102 |
| Interaction — DrugTreatment | the use of drug to treat disease | Diseasome | 812 |
| Interaction — ChemicalProteinInteraction | genomic response to chemical compounds | ChEMBL, BindingDB, DPSP Ki, TTD, BindingMOAD, DrugBank, CTD, MATADOR, ArrayExpress, KEGG | 47282 |

Chen, B., Ding, Y., & Wild, D. J. (2012). Improving integrative searching of systems chemical biology data using semantic annotation. *Journal of Cheminformatics*, 4:6 (doi:10.1186/1758-2946-4-6).

**Dereferenable URI**

**Browsing**

**PlotViz: Visualization**

**Bio2RDF**

**Chem2Bio2RDF**

**LODD**

**Others**

RDF
Triple store

**Cytoscape Plugin**

**Linked Path Generation and Ranking**
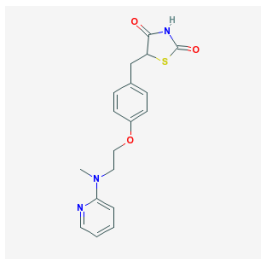
**SPARQL ENDPOINTS**

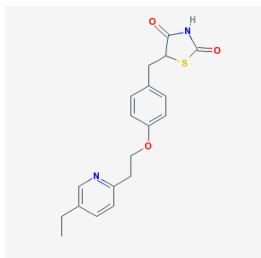**Third party tools**

Dijkstra's algorithm

# Thiazolinediones (TZDs) – revolutionary treatment for type II Diabetes



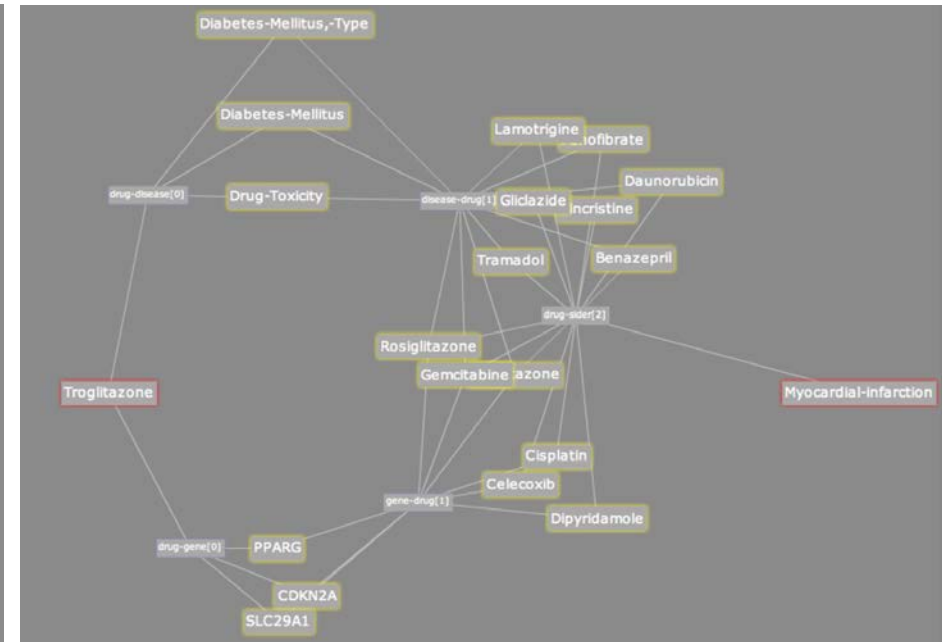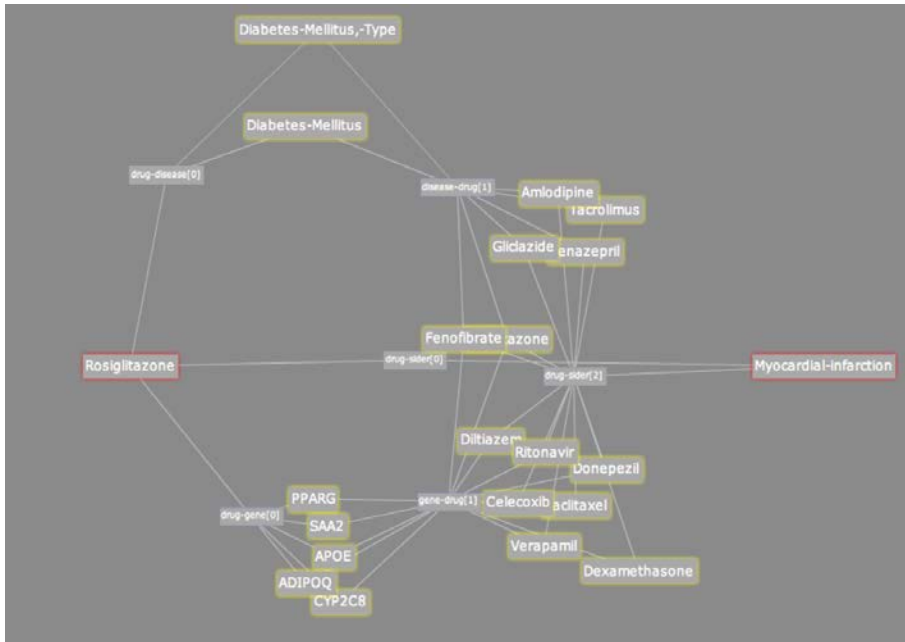Troglitazone (Rezulin): withdrawn in 2000 (liver disease)



Rosiglitazone (Avandia): restricted in 2010 (cardiac disease)



Rosiglitazone bound into PPAR-γ

Pioglitazone: ???? (does decrease blood sugar levels,  was associated with bladder tumors and has been withdrawn in some countries.)

**PPARG**: TZD target

**SAA2:** Involved in inflammatory response implicated in cardiovascular disease (Current Opinion in Lipidology 15,3,,269-278 2004)

**APOE**: Apolipoprotein E3 essential for lipoprotein catabolism. Implicated in cardiovascular disease.

**ADIPOQ:** Adiponectin involved in fatty acid metabolism. Implicated in metabolic syndrome, diabetes and cardiovascular disease

**CYP2C8:** Cytochrome P450 present in cardiovascular tissue and involved in metabolism of xenobiotics

**CDKN2A:** Tumor suppression gene

**SLC29A1:** Membrane transporter

# Semantic Prediction
# http://chem2bio2rdf.org/slap



Chen, B., Ding, Y., & Wild, D. (2012). Assessing Drug Target Association using Semantic Linked Data. *PLoS Computational Biology*, 8(7): e1002574. doi:10.1371/journal.pcbi.1002574,

# Example: Troglitazone and PPARG

# Topology is important for association
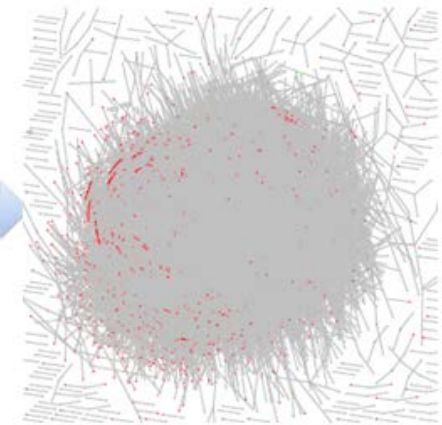
# Semantics is important for association

# SLAP Pipeline



| PubChem | ChEBI | DrugBank |
| --- | --- | --- |
| UniProt | UniProtKB-GOA | HGNC |
| SIDER | OMIM | KEGG |
| HPRD | ChEMBL | TTD |
| BindingDB | CTD | PDSP |

(a) Raw Data Sets

(b) Ontological level schema

(c) Semantic Linked Data

Path finding

(f) Statistical Models

1. Edge weight:

$$p(e(i \rightarrow j)) = \frac{1}{\sum_k \sum_{n=1}^{n=1} R_{i,n} == R_{i,j}}$$

2. Path score:

$$p(P_l(t \rightarrow s)) = p(P_l(e_{m \rightarrow m-1}, ..., e_{3 \rightarrow 2}, e_{2 \rightarrow 1})) = \prod_{i=1}^{m-1} e_{i+1 \rightarrow i}$$

$$log(p(P_l(t \rightarrow s))) = \sum_{i=1}^{m-1} log(e_{i+1 \rightarrow i})$$

3. Association score

$$raw\ score(s,t) = \sum_{l}^{n} \frac{log(p(P_l)) - \theta(log(P_l))}{\sigma(log(P_l))}$$

Association Assessing

Drug: Troglitazone

Target :PPARG

(e) Significant Paths between Two Nodes

Path filtering

Drug: Troglitazone

Target :PPARG

(d) Paths (length <4) between Two Nodes

Integration

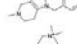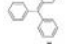Semantic annotation
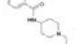
# Cross-check with SEA

- SEA analysis (Nature 462, 175-181, 2009) predicts 184 new compound-target pairs, 30 of which were experimentally tested
- 23 of these pairs were experimentally validated (<15uM) including 15 aminergic GPCR targets and 8 which crossed major receptor classification boundaries
- 9 of the aminergic GPCR target pairings were correctly predicted by SLAP (p<0.05) – for the other 6 compounds were not present in our set
- 1 of the 8 cross-boundary pairs was predicted

# Assessing drug similarity from biological function

- Took 157 drugs with 10 known therapeutic indications, and created SLAP profiles against 1,683 human targets
- Pearson correlation between profiles > 0.9 from SLAP was used to create associations between drugs
- Drugs with the same therapeutic indication unsurprisingly cluster together
- Some drugs with similar profile have different indications – potential for use in drug repurposing?

# Data2Knowledge platform...

**Data2Knowledge**

| AMiner | PMiner | SLAP | ... |
|--------|--------|------|-----|
| Mining knowledge from articles:<br>• Researcher profiling<br>• Expert search<br>• Topic analysis<br>• Reviewer suggestion | Mining knowledge from patents:<br>• Competitor analysis<br>• Company search<br>• Patent summarization | Mining drug discovery data<br>• Predicting targets<br>• Repurposing drugs<br>• Heterogeneous graph search | Mining more data... |

# AMiner

- Research profiling
- Integration
- Interest analysis

- Topic analysis
- Course search
- Expert search

- Association
- Disambiguation
- Suggestion

- Geo search
- Collaboration recommendation

# What is PMiner?

- Current patent analysis systems focus on search
  - Google Patent, WikiPatent, FreePatentsOnline

- PMiner is designed for an ***in-depth*** analysis of patent activity at the topic-level
  - Topic-driven modeling of patents
  - Heterogeneous network co-ranking
  - Intelligent competitive analysis
  - Patent summarization

* Patent data:
> 3.8M patents
> 2.4M inventors
> 400K companies
> 10M citation relationships

* Journal data:
> 2k journal papers
> 3.7k authors

The crawled data is increasing to >300 Gigabytes.

J. Tang, B. Wang, Y. Yang, P. Hu, Y. Zhao, X. Yan, B. Gao, M. Huang, P. Xu, W. Li, and A. K. Usadi. PatentMiner: Topic-driven Patent Analysis and Mining. KDD'12, pp. 1366-1374.

# Semantic Publishing

- Turn literature knowledge into actionable data to generate more powerful knowledge



**http://nanopub.org/**

**Paper 1, 2, …**

AAAAAAAAAAAAAAAAAAAAAAAAA **Gene1**

AAAAAAAAAAAAAAAAAAAAAAAAA **Gene1**
AAAAAAA(XYZ, 2002), AAAAAAA
**Disease2**AAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAA

Paper n: Gene 1, Disease 2,  (XYZ, 2002)

……

Paper 2: Gene 1, Disease 2,  (XYZ, 2002)
Paper 1: Gene 1, Disease 2,  (XYZ, 2002)

Demonstrating strong evidence of the connection/relationship between concept Gene 1 and Disease 2, or concept GENE and DISEASE

38

# Knowledge Graph

# Semantic Search



**The 5 Steps of Google's Semantic Search**

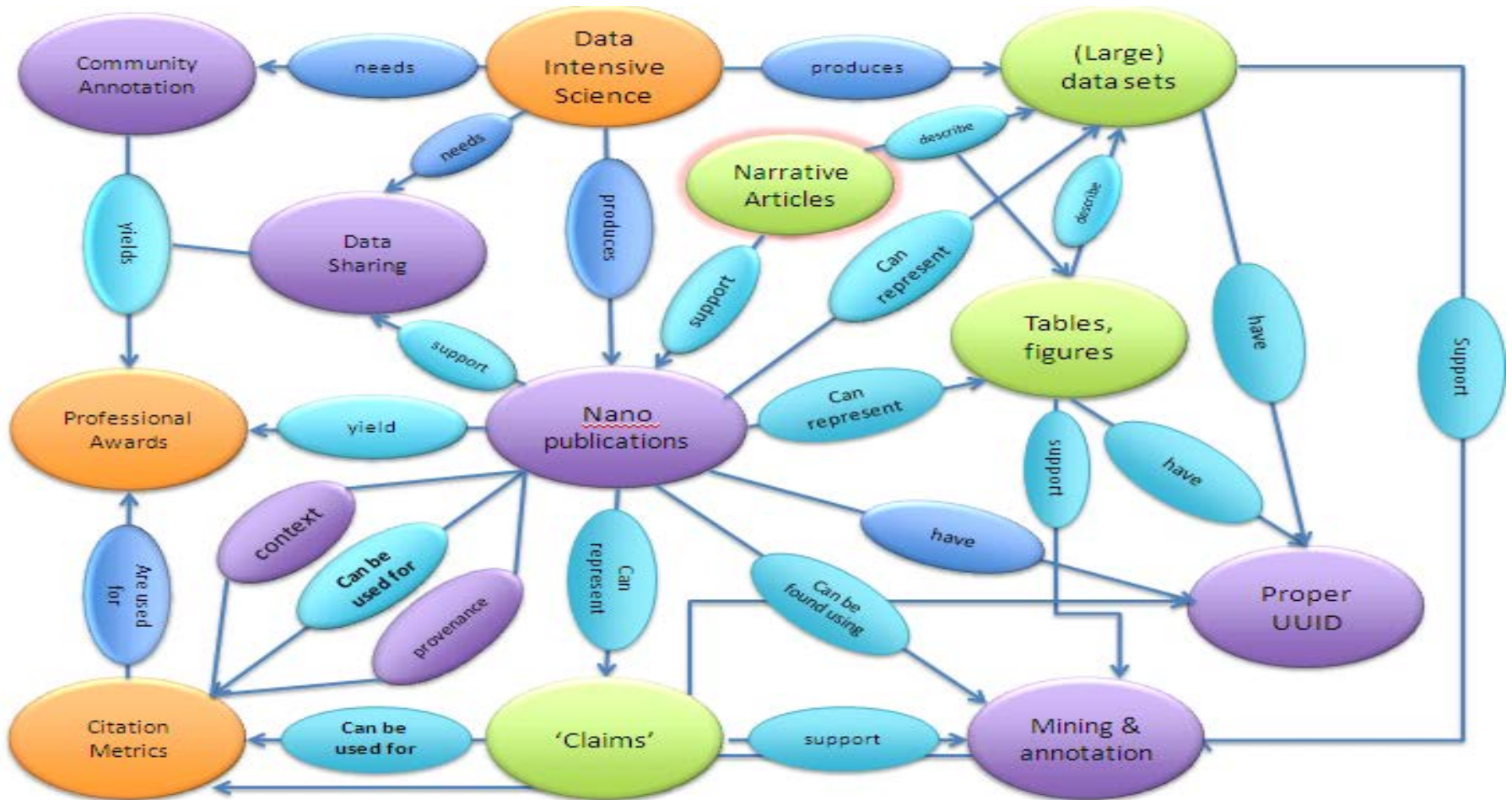Index the entire web. Discover new relationship data

Create comprehensive list of trusted sources. Use information to verify new data.

Knowledge Graph

Fact Extraction

Web Questions

Use social signals and user-generated data to mine facts and create entities.

Ask trusted Humans

Solve Issues

Ask questions to discover facts you don't yet have.

**Each step leads to a more semantic web**

(C) davidamerland.com

# Digital Discovery



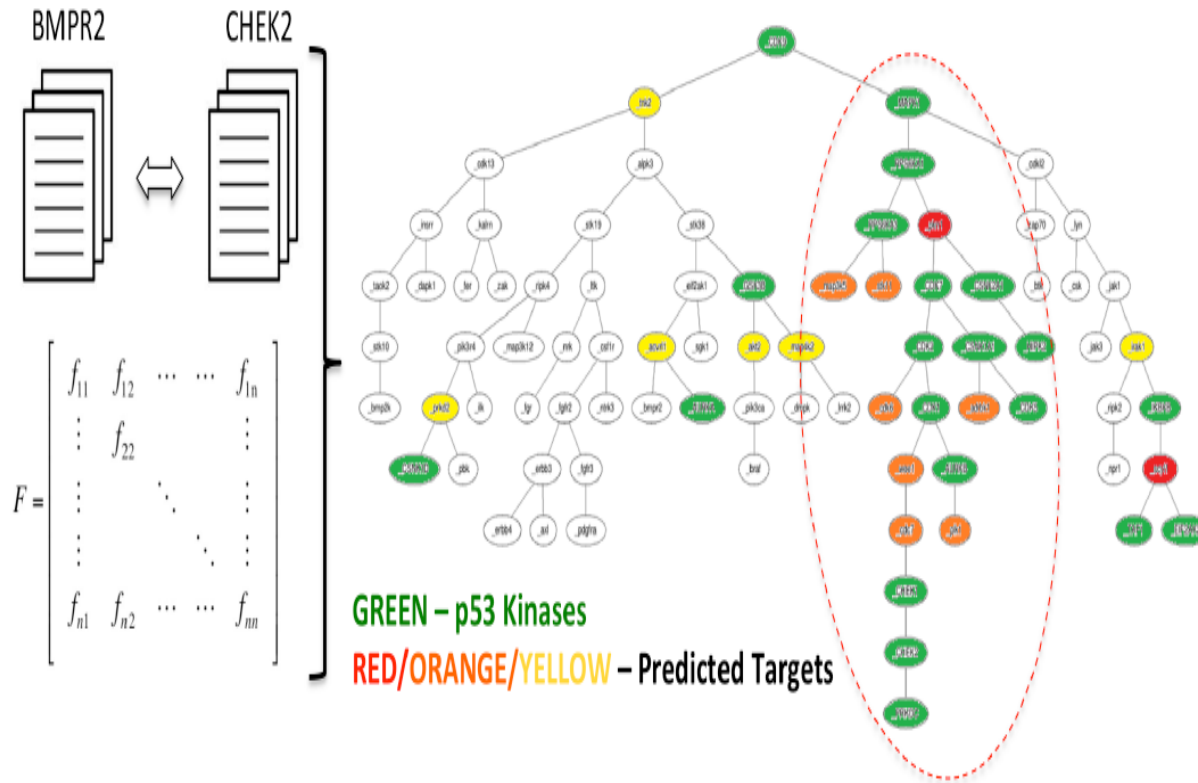Figure 1 Kinases are clustered based on their literature distance. The clustered p53 kinases (green) suggest new kinases that may also phosphorylate p53.
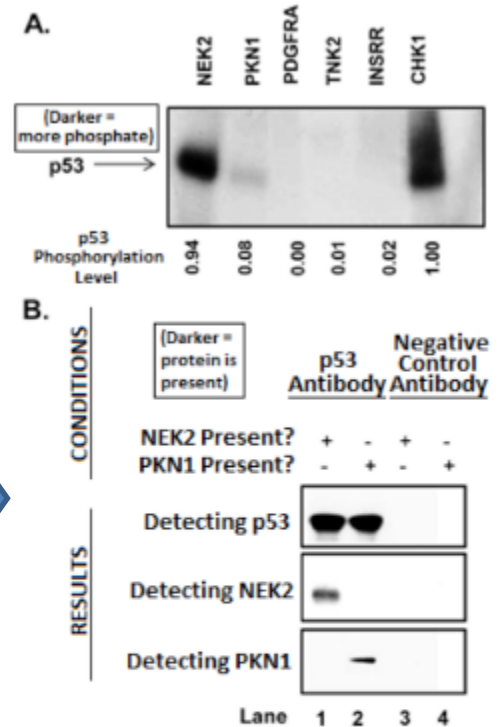
GREEN – p53 Kinases

RED/ORANGE/YELLOW – Predicted Targets

Figure 6 Experimental validation of candidate p53 kinases as *bona fide* p53 kinases. (A) *In vitro* kinase assay demonstrates phosphorylation of p53 by top ranked candidate kinases PKN1 and NEK2. Relative levels of p53 phosphorylation are indicated for each kinase normalized to positive control CHK1. Though the signal is weak for PNK1, subsequent experiments lend further support. (B) PKN1 and NEK2 shown to interact with p53 *in vivo*. A p53 antibody isolates p53 and any proteins bound to it. Antibodies detect the presence of candidate kinases in this isolate.
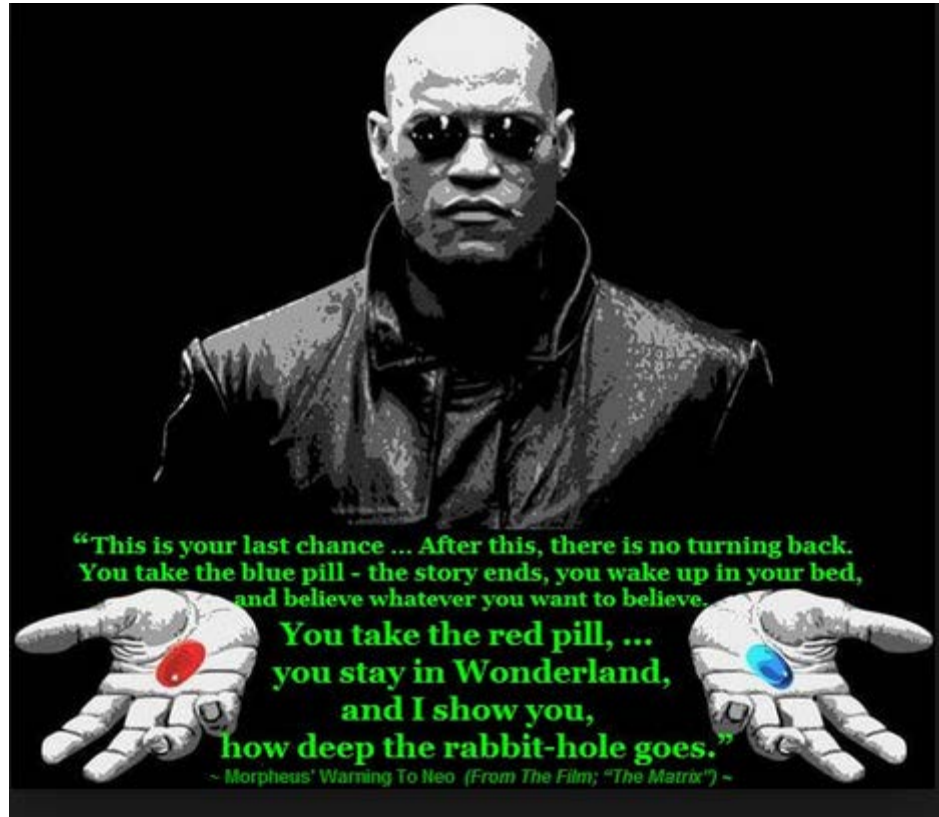
Spangler, S., Wilkins, A.D., Bachman, B. J., et al.(2014). Automated hypothesis generation based on mining scientific literature. Proceedings of the 20th *ACM SIGKDD*, August 24-27, 2014, New York, USA.

# More

- Data-Driven Discovery
  - Medicine
  - Neuroscience
  - Math and Material science
  - Social science (education, business, poverty reduction)
  - Digital humanities and Arts (distant reading, digital painting, digital recipes, computational creativity (story, joke and poetry generation)
- Digital Creativity

# Future of Knowledge

**Questions?**



The Matrix