# Daniel Himmelstein, Leo Brueggeman & Sergio Baranzini present *Repurposing drugs on a heterogeneous network*
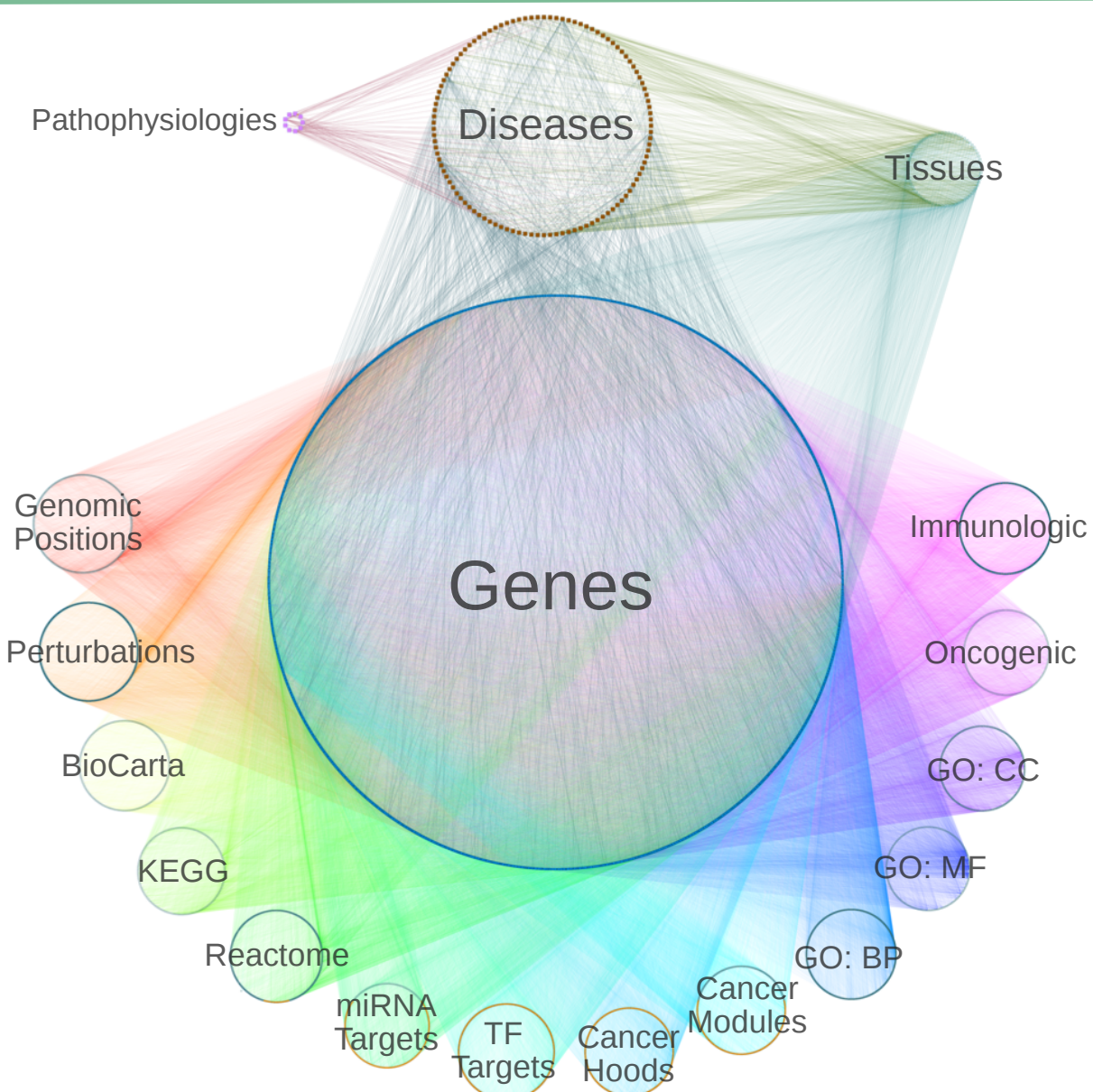
Last year, we introduced *heterogeneous network edge prediction* (HNEP) to underline predict disease-associated genes.
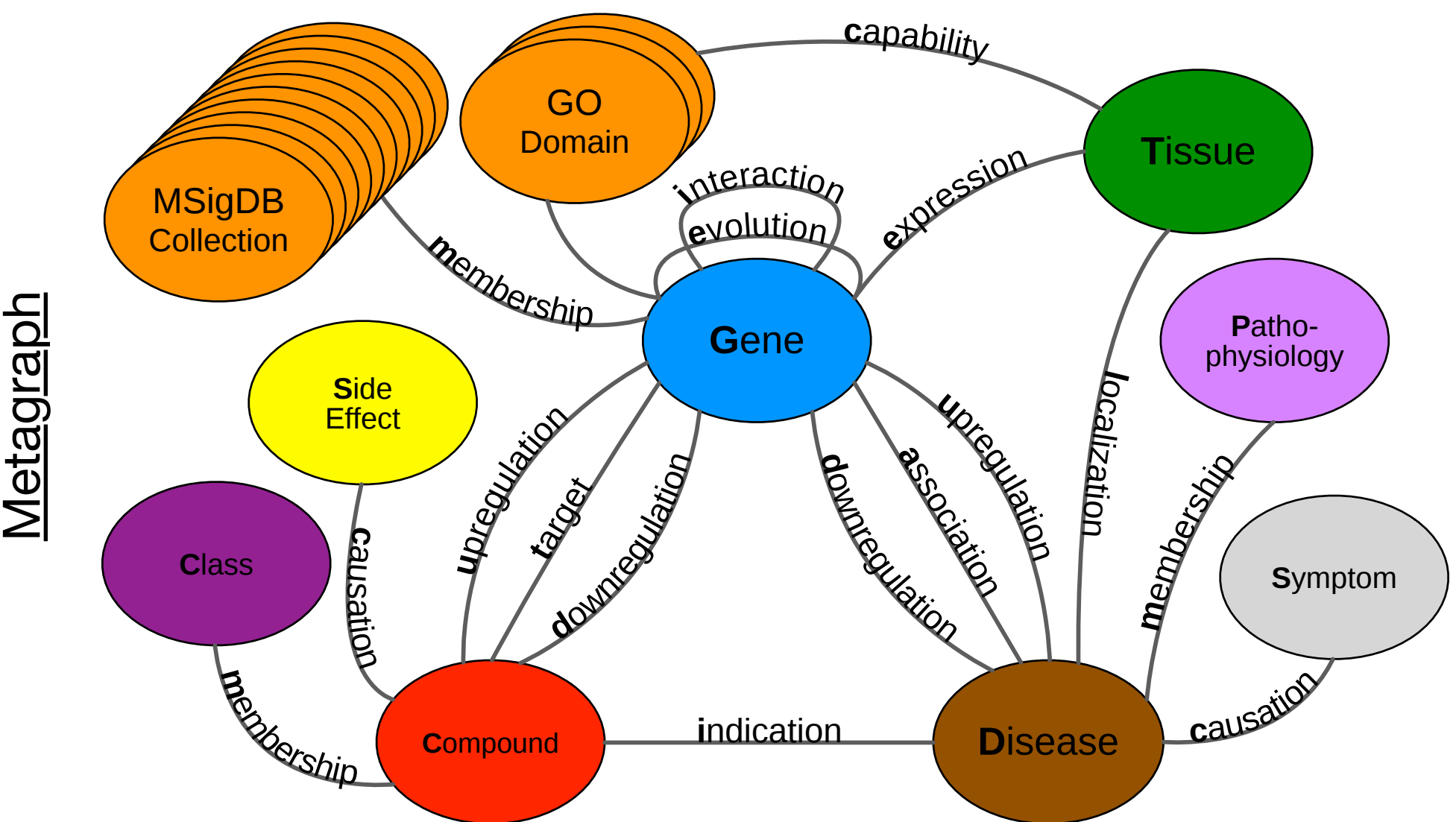
Heterogeneous networks contain multiple node and edge types.

Our network contained 40,343 nodes (of 18 types) and 1,608,168 edges (of 19 types).



## HNEP Method



A. Metagraph:

B. Metapaths for G ---a--- D :

C. Hypothetical graph:

D. Calculating and weighting path counts:

$$DWPC(metapath) = \sum_{path \in Paths} PDP(path)$$

$$PDP(path) = \prod_{d \in D_{path}} d^{-w}$$

## Media

Forthcoming in *PLOS Computational Biology*
preprint on *bioRxiv* [doi:10.1101/011569]

Predictions online at het.io



Unlocking THE GENETICS of Complex Diseases: GWAS and Beyond
By Kristin Sainani

Featured in Stanford's *Biomedical Computation Review*

---

Now in 2015, we will use this data integration approach to repurpose drugs on a heterogeneous network.

## Planning the Network Construction



### Metagraph

### Nodes

| Type | Resource |
|------|----------|
| Compound | DrugBank |
| Disease | Disease Ontology |
| Gene | Entrez Gene |
| Tissue | Uberon |
| Gene Set | MSigDB |
| Side Effect | UMLS |
| Pathophysiology | Manual |
| Symptom | MeSH |

Standardized terminologies:
- provide a scalable framework for data integration
- prevent redundancy
- enable semantic data

### Edges

| Source | Target | Type | Resource |
|--------|--------|------|----------|
| Compound | Disease | Indication | MEDI |
| Compound | Disease | Indication | LabeledIn |
| Compound | Gene | Expression | LINCS |
| Compound | Side Effect | Causation | SIDER 2 |
| Compound | Side Effect | Causation | OFFSIDES |
| Disease | Gene | Target | ChEMBL |
| Disease | Gene | Association | GWAS Catalog |
| Disease | Gene | Expression | STAR-GEO |
| Disease | Pathophysiology | Membership | Manual |
| Disease | Symptom | Causation | Human symptoms--disease network |
| Gene | Gene | Interaction | Human Interactome Project |
| Gene | Gene | Interaction | The Incomplete Interactome |
| Gene | Gene | Evolution | Evolutionary Rate Covariation |
| Gene | Gene Set | Membership | MSigDB |
| Gene | Tissue | Expression | GNF Gene Expression Atlas |

Ideal resources are:
- high-throughput
- systematic
- unbiased
- aggregately diverse

---

And you can follow in realtime and get paid to participate.

## ThinkLab

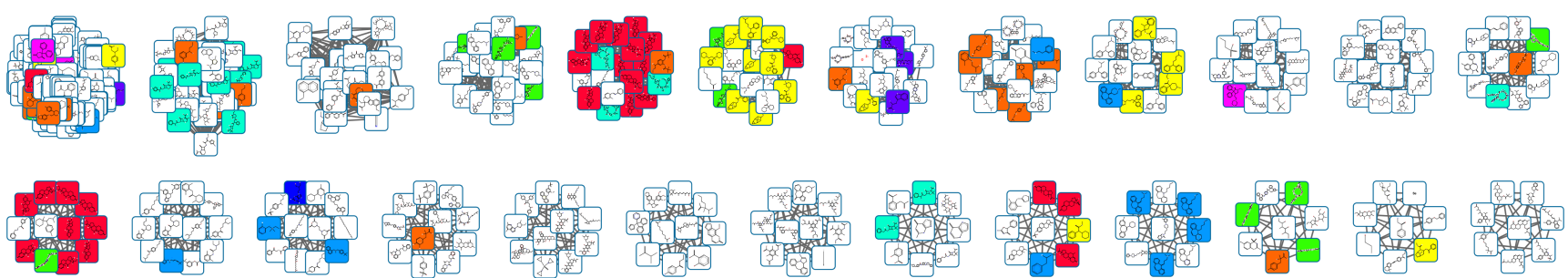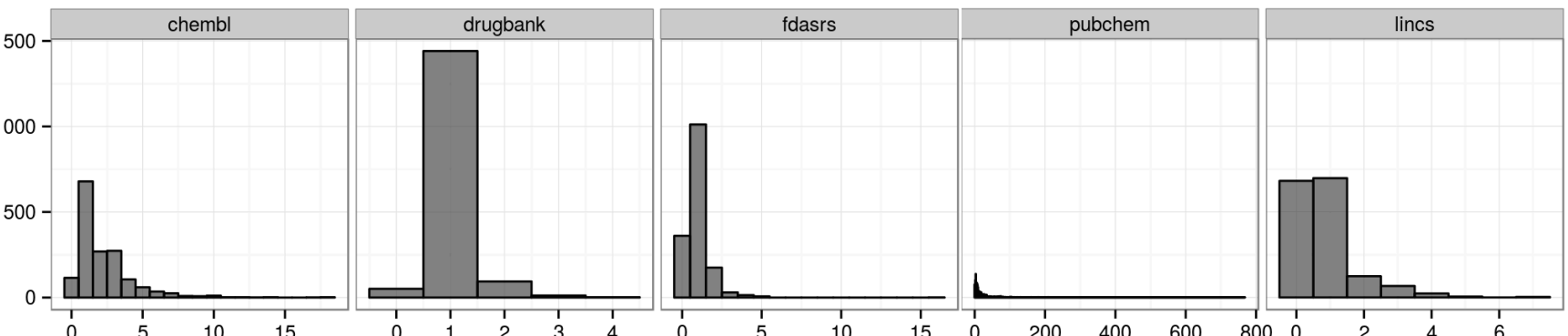thinklab.com/p/rephetio
doi:10.15363/thinklab.4



**ThinkLab** is:
- massively collaborative — all are welcome
- open science — all content is CC-BY
- incentivized — contributions are rewarded
- productive — scientific markdown editor
- efficient — code and results public upon commit

## Results (as of March 2015)

We analyzed **SIDER 2** and investigated its strengths and weaknesses as well as pharmacological utility.



Side-effect similarity modules were concordant with structural similarity modules (colored).

git.dhimmel.com/SIDER2

We created a user-friendly service to retrieve **Gene Ontology annotations** with optional propagation.

| Propagated | | Unpropagated |
|---|---|---|
| Entrez | | Symbol |
| All Genes | | Protein-coding Genes |

git.dhimmel.com/gene-ontology

We mapped compound vocabularies to DrugBank using **UniChem** to enable fuzzy matching.



Number of matches to each approved small molecule in DrugBank

git.dhimmel.com/drugbank/unichem-map.html