

# Systematic explanation of drug effects and gene dependency correlations using mechanistic networks assembled from literature mining and databases

John A. Bachman<sup>1</sup>, Benjamin M. Gyori<sup>1</sup>, Klas Karis<sup>1</sup>, and Peter K. Sorger<sup>1</sup>

<sup>1</sup>Laboratory of Systems Pharmacology, Harvard Medical School, Boston, MA, USA

Availability: <https://github.com/sorgerlab/indra> Funding: DARPA ASDF program, ARO contract W911NF018-1-0124



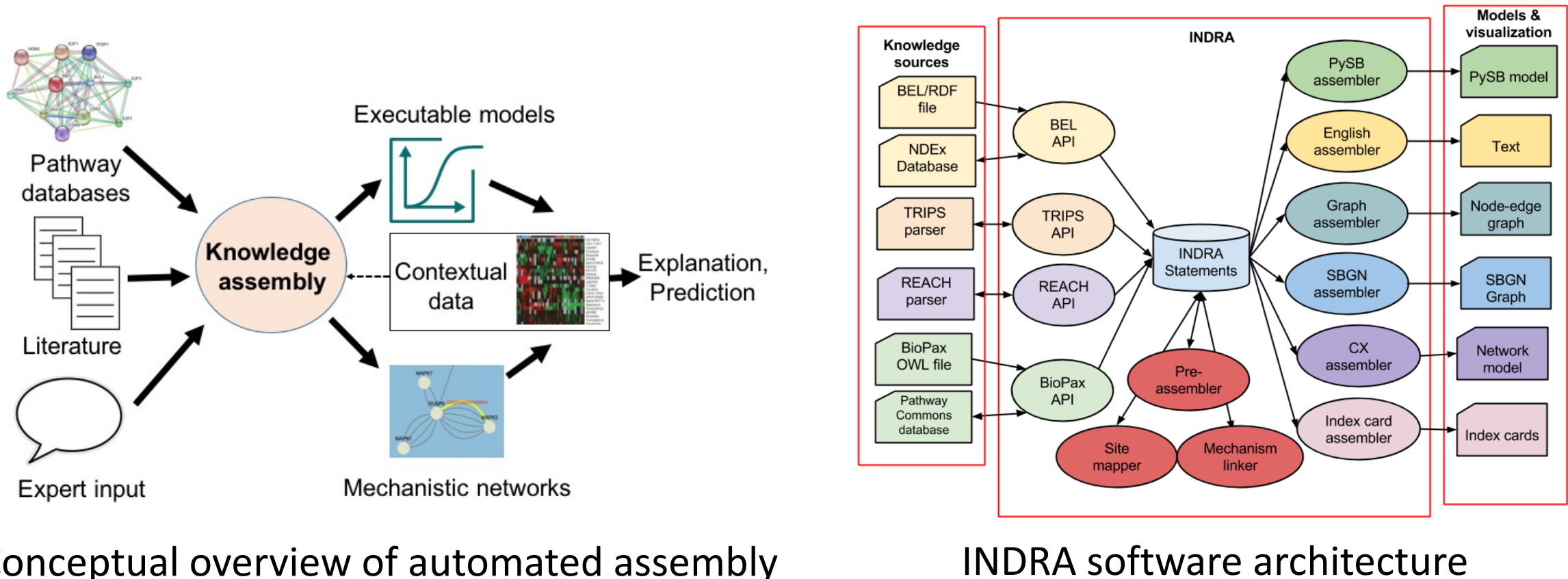
## INTRODUCTION

Biology is currently grappling with the challenge of integrating large datasets with a body of scientific knowledge that has grown too large and complex for any single scientist to read or understand. The size and scope of large-scale data compendia (“big data”) has created analytical challenges that require new solutions, ideally ones that make use of aggregated scientific knowledge. The capacity of modern experimental methods to generate data about biological processes has surpassed the ability of existing informatics approaches to generate meaningful mechanistic explanations. Mechanistic systems biology models could potentially address this gap, but model construction remains a labor-intensive process requiring both biological knowledge and modeling expertise. As a result, modeling studies remain fairly small in scope and are disconnected from genome-scale research. For mechanistic models to attain the necessary scope, methods for the automated assembly and analysis of large models from available knowledge sources will be required. Here we describe the use of the Integrated Network and Dynamical Reasoning Assembler (INDRA)<sup>1</sup> to assemble mechanistic facts from databases and literature into different types of models for explanation of large datasets.

## RESULTS

### System architecture and approach

The Integrated Network and Dynamical Reasoning Assembler (INDRA)<sup>1</sup> automatically assembles mechanistic models from pathway databases, literature, and expert knowledge expressed in natural language. INDRA draws on three existing natural language processing systems<sup>4,5,6</sup> and uses a modular architecture to build different types of models from a variety of sources.

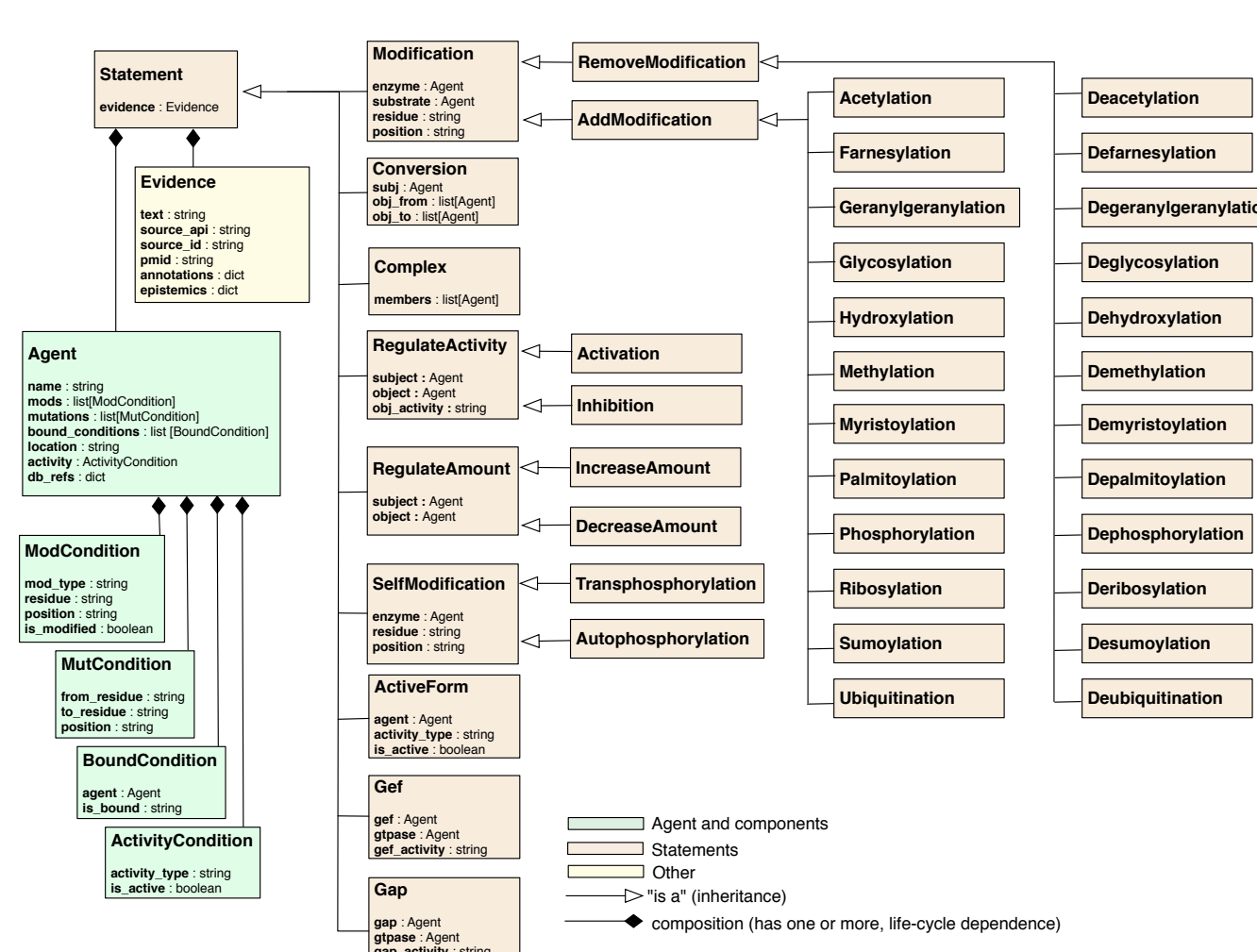


Conceptual overview of automated assembly

INDRA software architecture

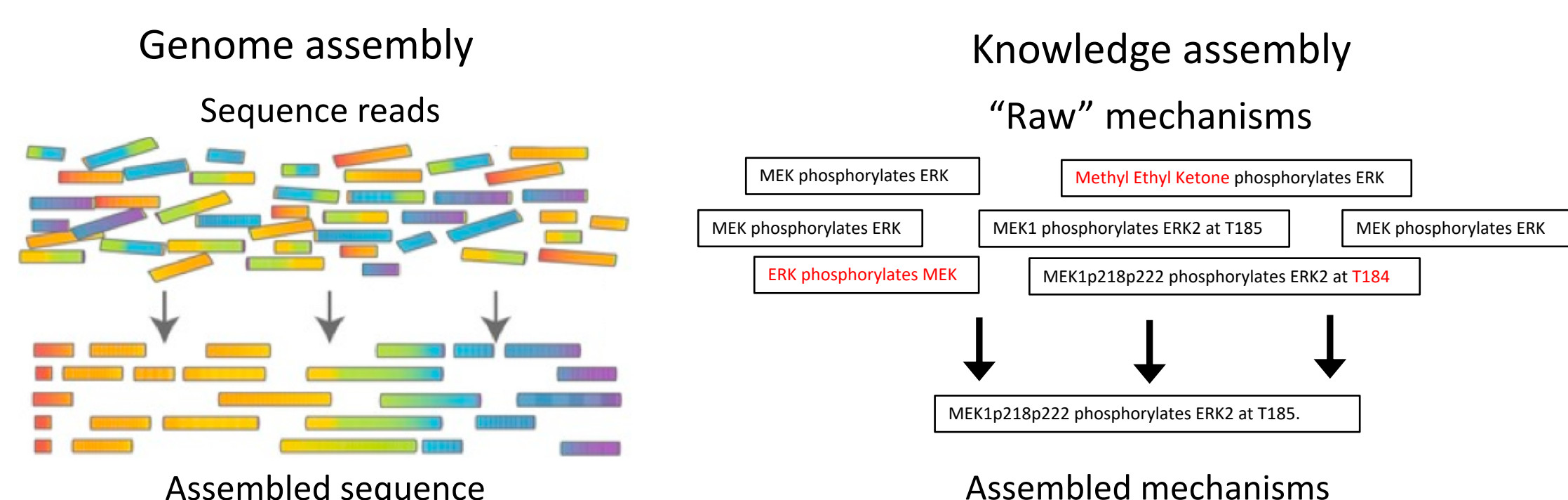
### Mechanisms are normalized into Statements

Mechanisms extracted from each source format are normalized into Statements, an SBO-compatible internal representation, where they are processed to remove errors, identify overlaps, and estimate reliability. Statements are designed to correspond in both specificity and ambiguity to descriptions of biochemistry as found in text (e.g., “MEK1 phosphorylates ERK2”, rather than a detailed reaction mechanism). The representation currently encompasses post-translational modifications, chemical conversions, protein expression and degradation, and generic activation/inhibition relationships.



### The assembly challenge

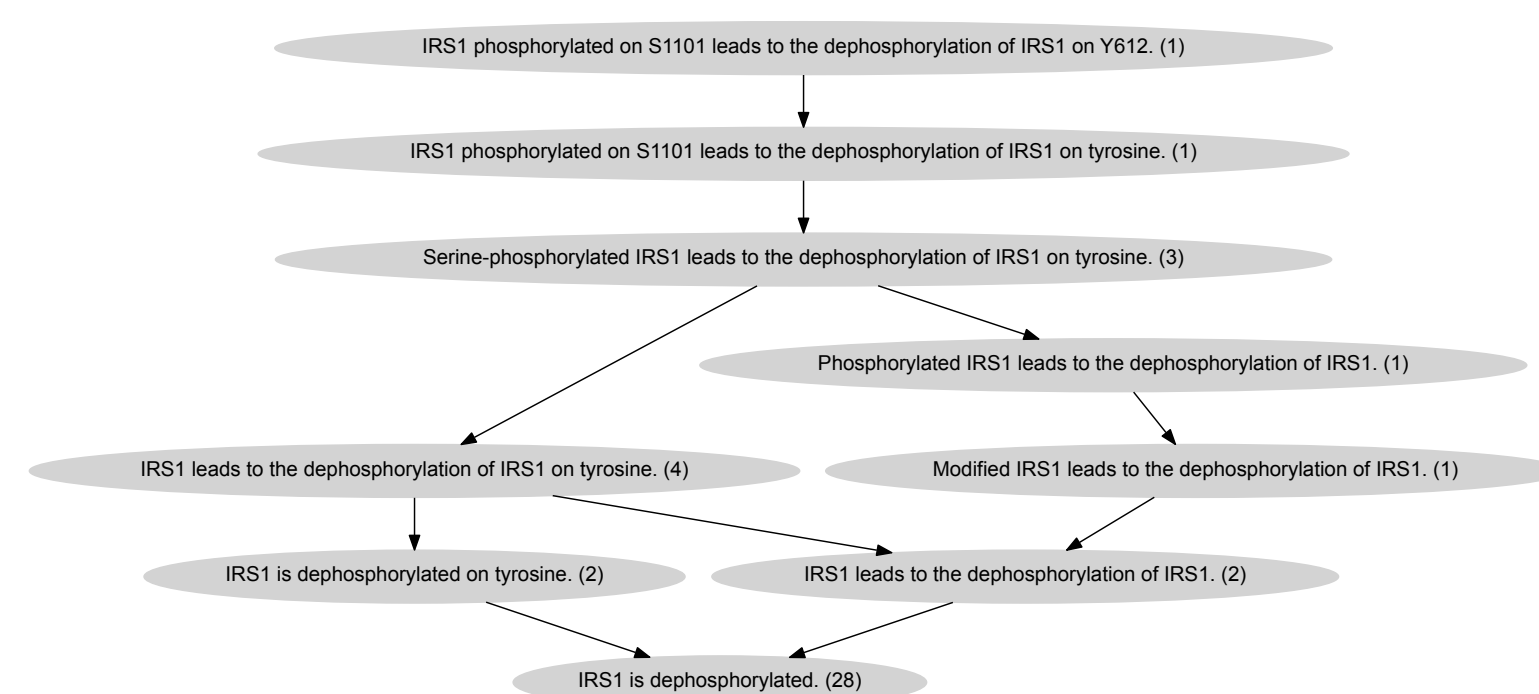
Assembly of a large number of mechanistic facts is analogous to genome assembly: databases and literature yield a large number of redundant, partially overlapping facts that may contain errors. Mechanisms must be corrected and “aligned” in order to produce a set of facts suitable for generating a non-redundant, non-degenerate model.



### Identifying relationships between mechanisms

A key challenge in assembling detailed mechanistic networks is that a single mechanism may be described at different levels of specificity among the literature and various databases. Reconciling these overlapping mechanisms is essential to eliminate spuriously distinct edges in the assembled model. Using hierarchical ontologies of protein modification types, activity types, and the protein family information provided in Bioentities, INDRA implements duplicate removal, hierarchy-based redundancy resolution, and other forms of error correction and mechanism linking.

Relations can be organized into a hierarchy based on their specificity



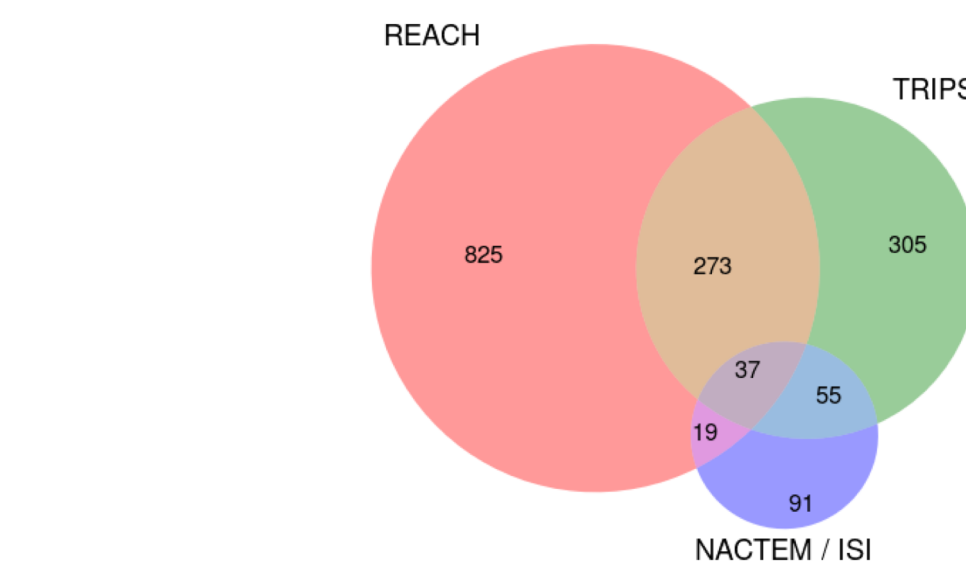
### Estimating the reliability of extracted mechanisms

Even state-of-the-art NLP and text mining algorithms have limited accuracy, with roughly 20-30% of extracted relations representing a misinterpretation of the corresponding sentence (“reader error”). Given empirical estimates of the per-sentence error rate for different readers, INDRA’s BeliefEngine component aggregates results to estimate the overall probability that a relation is the result of reader error. It accomplishes this by:

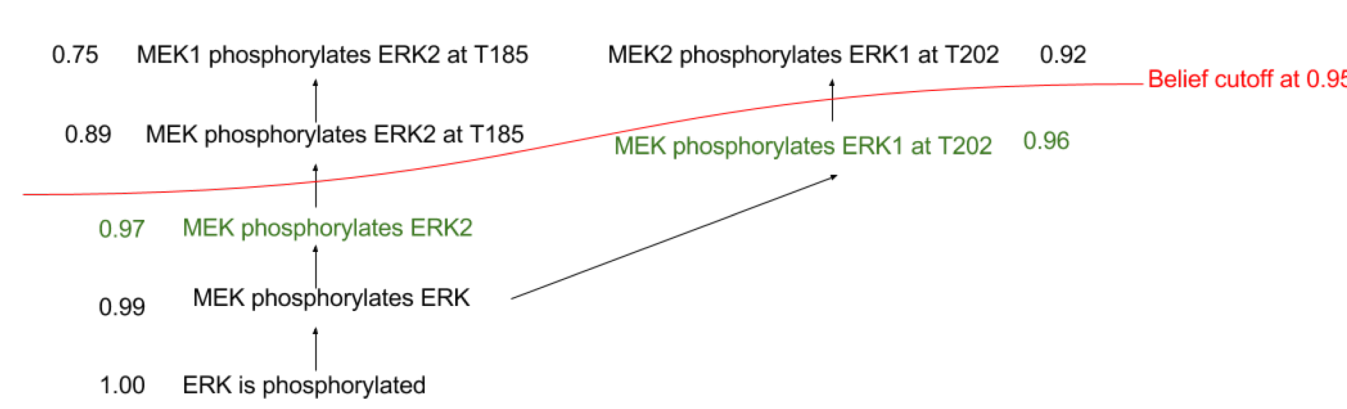
- 1) aggregating evidence from multiple sentences read by the *same reader*
- 2) aggregating results from different reading algorithms on the *same sentence*
- 3) propagating error estimates through the network of related statements

Mechanisms can then be filtered with a precision threshold (e.g., 95% confidence).

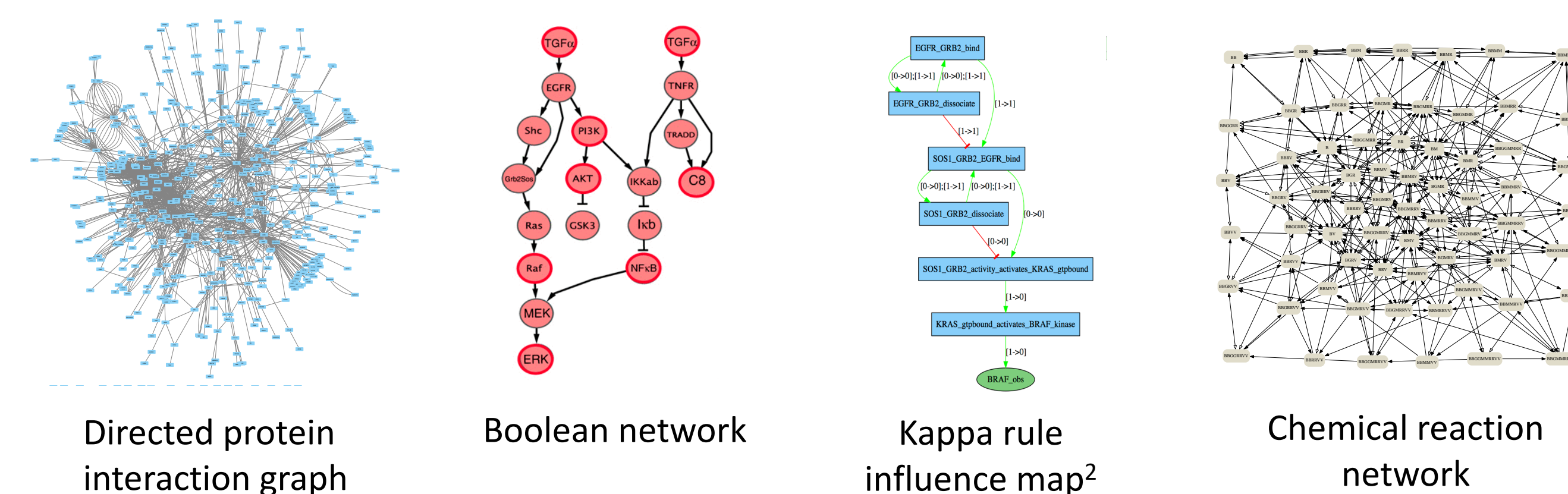
Reading systems produce partially overlapping extractions



Reliability estimates are propagated through the specificity hierarchy



### Model representations for statically identifying causal paths



More false positive paths (less stringent context) vs. More false negative paths (more stringent context)

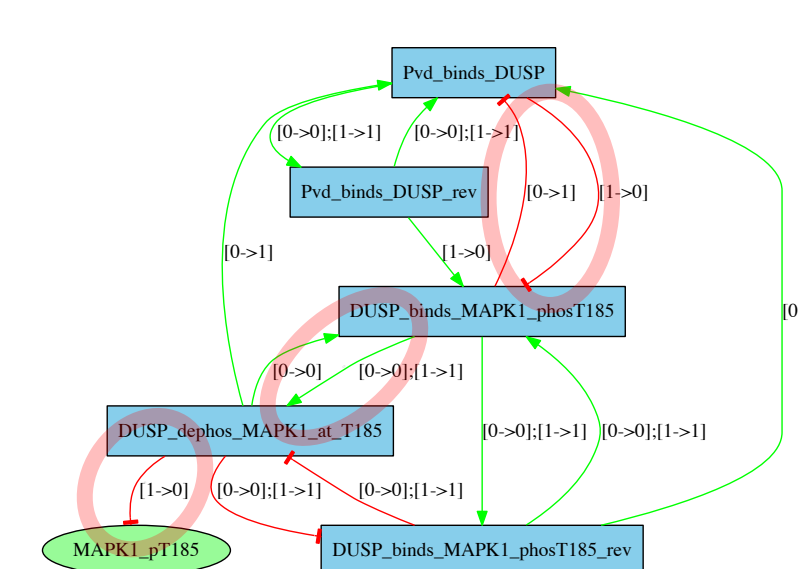
In directed interaction graphs, the relatively limited causal context leads to an explosion of paths between any two proteins. This leads to many false positive paths and makes identification of long causal chains difficult (or even intractable) in large networks.

### Generating explanations from the Kappa<sup>2</sup> rule influence map

The Kappa influence map captures detailed context while avoiding the combinatorial explosion of chemical species. Paths are obtained by:

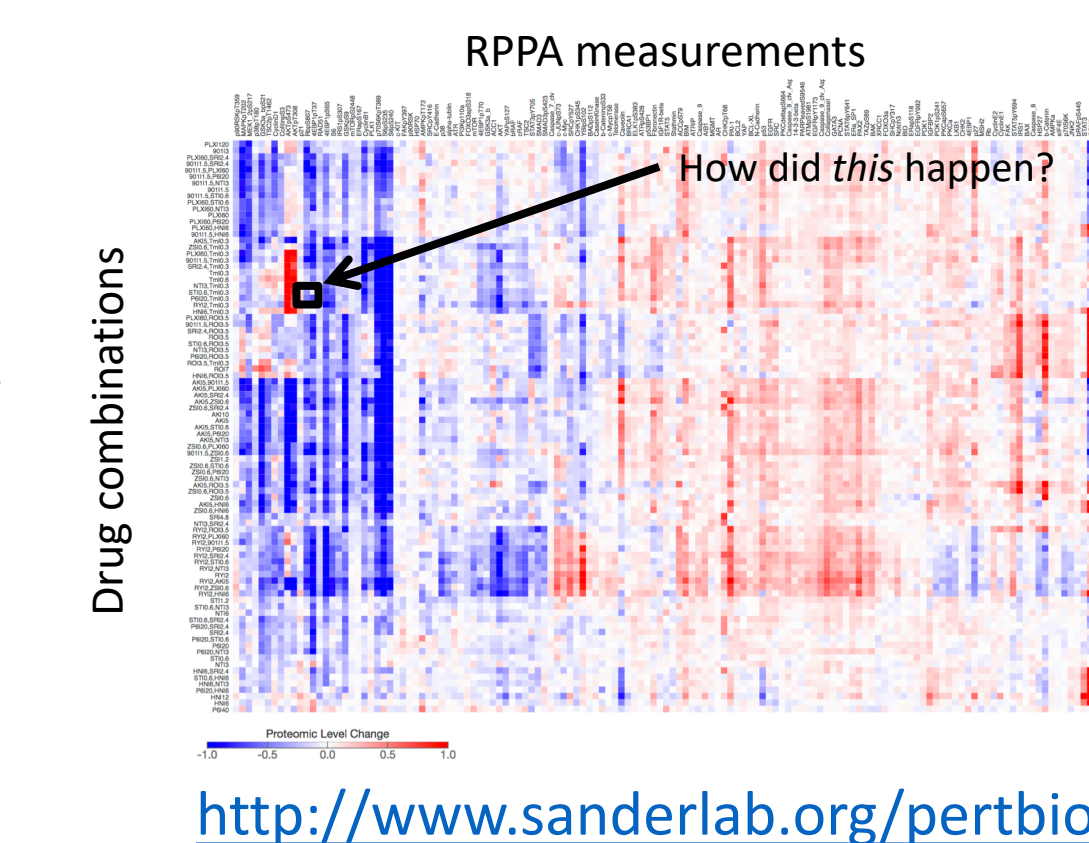
- identifying rules whose activity is increased by the abundance of the subject (e.g., drug)
- searching for a path to an observable representing the object (e.g., a measured protein) with the appropriate overall polarity
- scoring paths by whether the signs of measured intermediate nodes are correctly predicted

Causal path for “Pervanadate increases MAPK1 phosphorylation”

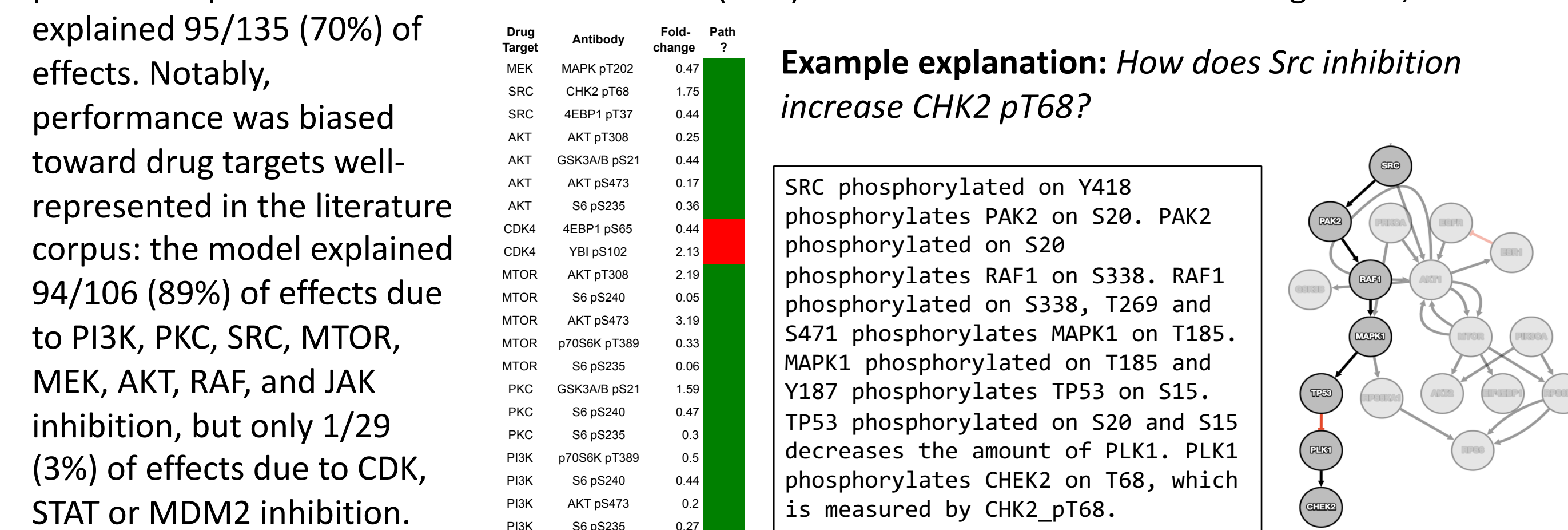


### Use case 1: interpreting phosphoproteomic data

To evaluate the ability of INDRA to systematically generate explanations of high-throughput data, we assembled a rule-based executable model to explain a previously published dataset of the phospho-proteomic response of a melanoma cell line to 12 different drugs.<sup>3</sup> A rule-based model containing 221 proteins and 1451 rules was assembled from mechanisms extracted from databases and ~95,000 publications (abstracts and full texts). Static analysis of the rule influence map provided by Kappa identified possible mechanistic paths linking drug targets to experimentally observed effects on phosphoprotein abundances.



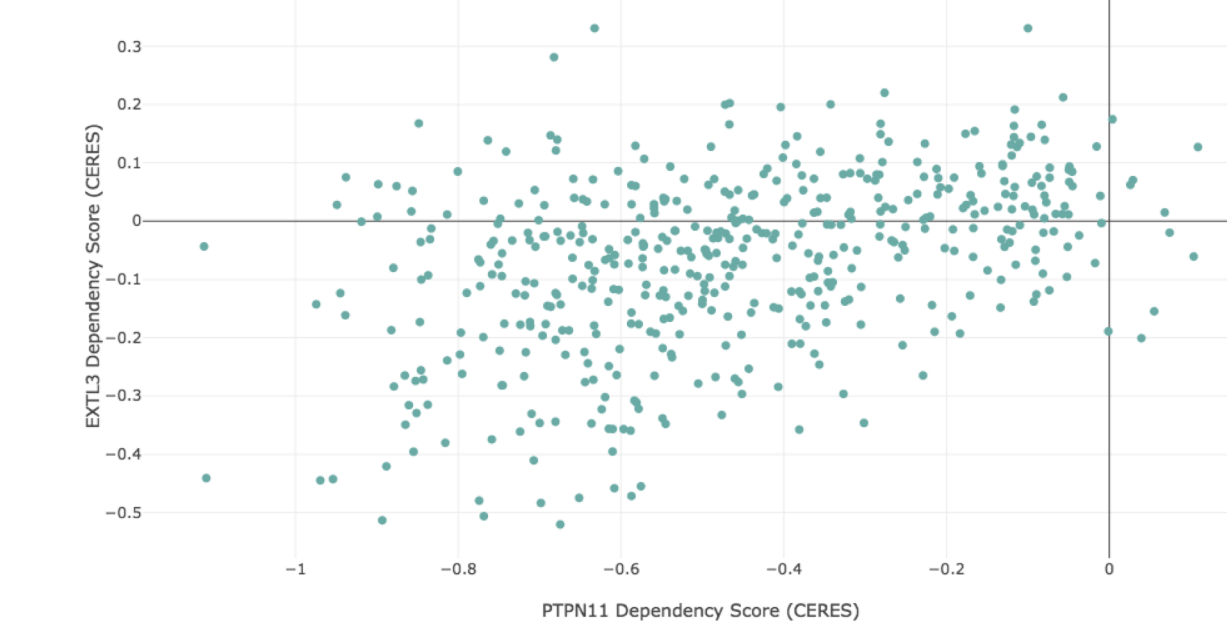
**Performance:** For the largest effects in the data (>50% fold-change) the model generated biochemically plausible explanations for 20 of the 22 effects (91%). For effects at the 20% fold-change level, the model explained 95/135 (70%) of effects. Notably, performance was biased toward drug targets well-represented in the literature corpus: the model explained 94/106 (89%) of effects due to PI3K, PKC, SRC, MTOR, MEK, AKT, RAF, and JAK inhibition, but only 1/29 (3%) of effects due to CDK, STAT or MDM2 inhibition.



### Use case 2: explaining gene dependency correlations

We are evaluating the effectiveness of our text-mined network of ~7.5 million unique interactions mined from the literature to identify mechanisms mediating correlations between gene knockouts in the Broad Cancer Dependency Map dataset (<http://depmap.org>)

Example: correlation of PTPN11 (Shp2) with EXTL3: 0.42



Of the 1.7 million correlations (absolute value > 0.3) we found that while only 6,043 (0.3%) could be explained by a known mechanistic link, 22% could be explained with a single intermediate linking gene, either as a common downstream target or part of a pathway.

Type	Count	Percent
(All correlations)	155,417,265	
Correlations > 0.3	1,763,551	100.0
Direct link	6,043	0.3
Pathway or shared target	380,404	21.6
Correlations explained	386,447	21.9

For a subset of 48,000 correlations involving 340 cancer-related genes, we found that a higher proportion (75%) of the correlations can be explained with a single intermediate linking gene.

We are building a browser, the INDRA DepMap Explainer, to explore these explanations.

**Example Explanation: EXTL3 → FGF2 → PTPN11**

Support for EXTL3->FGF2:

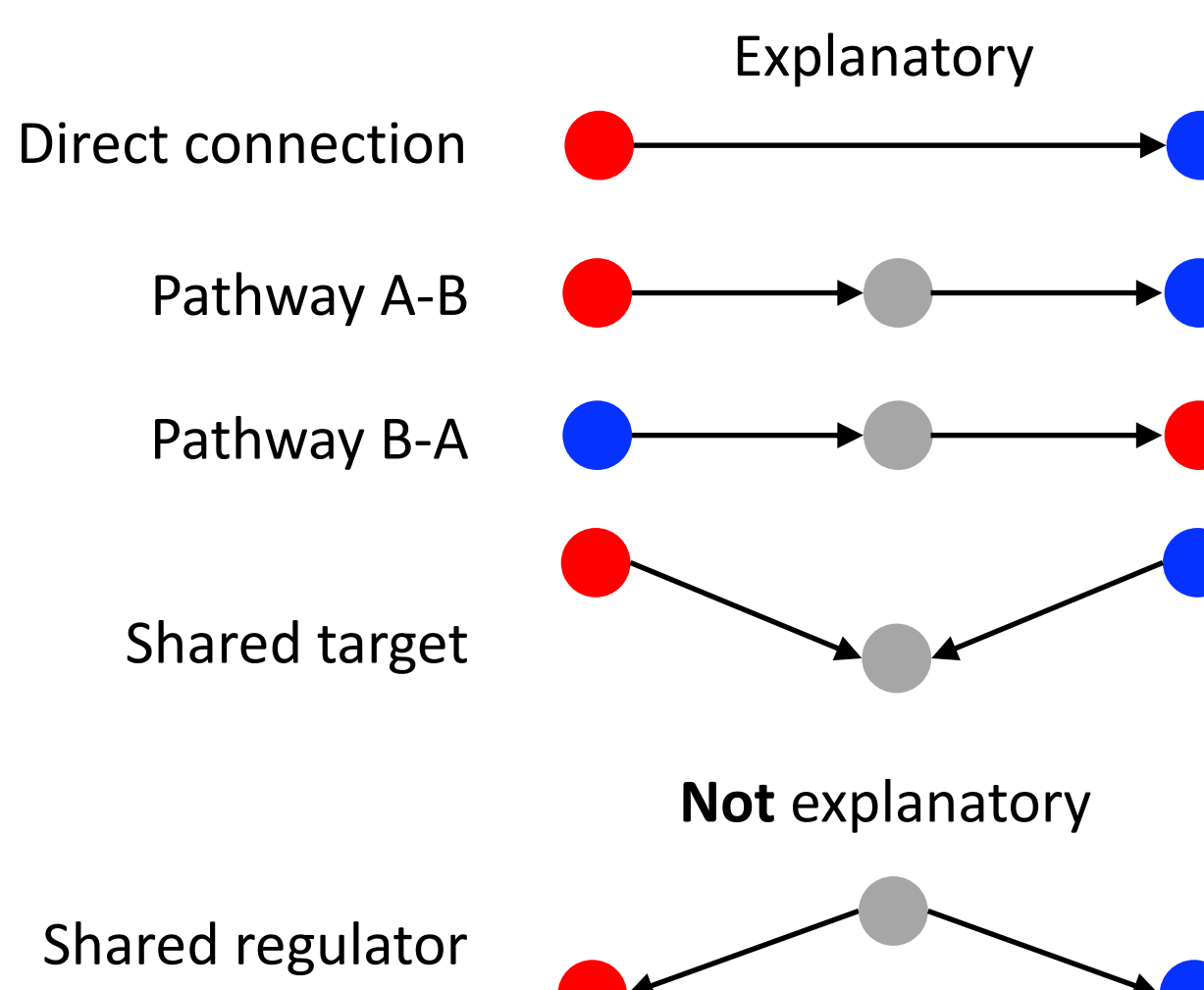
"These results indicate that by altering HS chain length and composition, EXTL3 mutations potentiate FGF2 signaling, thereby contributing to the pathophysiology of the skeletal dysplasia observed in the patients." (PMID 28148688)

"Reduced expression of either EXT1, EXT2, or EXTL3 decreased heparan sulfate biosynthesis, and consequently suppressed the FGF2 dependent proliferation of mouse L fibroblasts." (PMID 29305908)

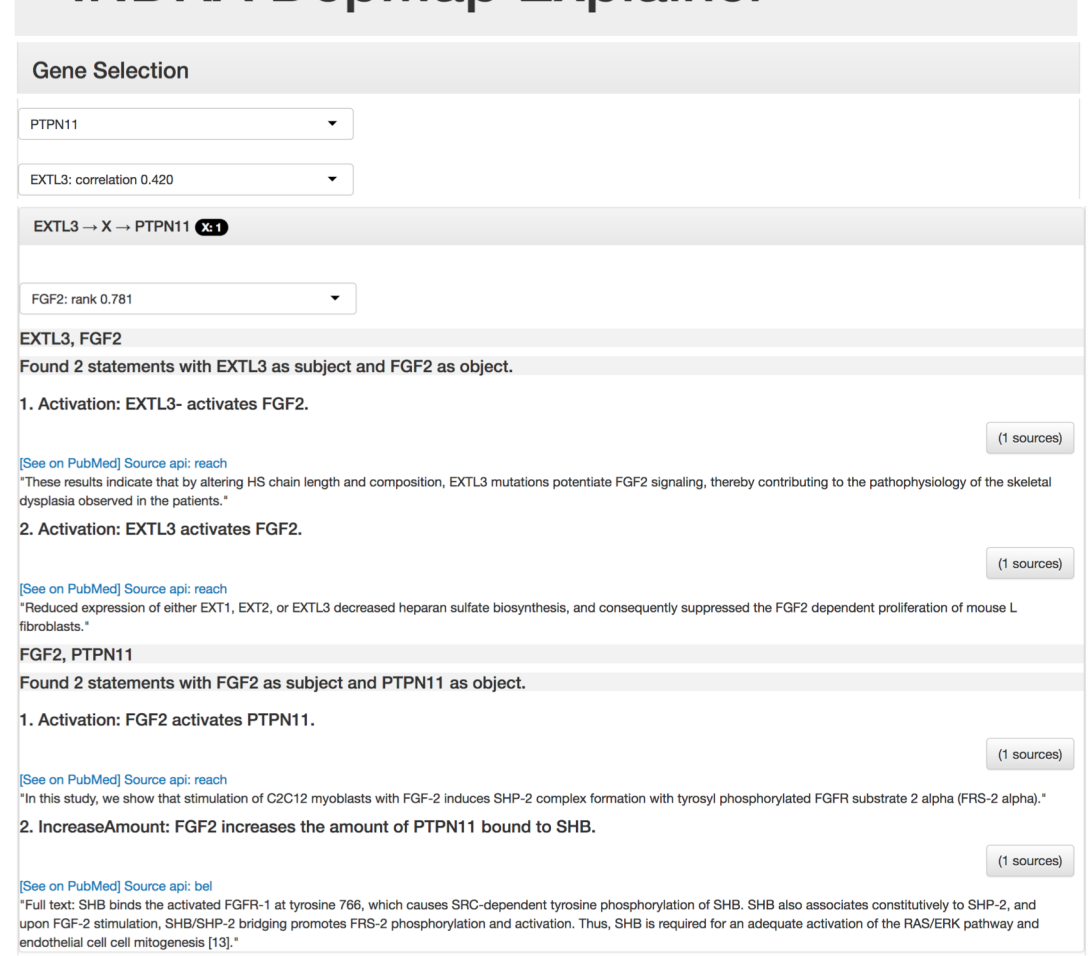
Support for FGF2 -> PTPN11:

"In this study, we show that stimulation of C2C12 myoblasts with FGF-2 induces SHP-2 complex formation with tyrosyl phosphorylated FGFR substrate 2 alpha (FRS-2 alpha)."

Network configurations



INDRA DepMap Explainer



## REFERENCES

1. B. M. Gyori, J. A. Bachman, K. Subramanian, J. L. Muhlich, L. Galescu, and P. K. Sorger. "From word models to executable models of signaling networks using automated assembly." *bioRxiv*, 2017.
2. V. Danos, J. Feret, W. Fontana, R. Harmer, and J. Krivine. "Rule-Based Modeling of Cellular Signaling." *Concurrency Theory (CONCUR)*, Lecture Notes in Computer Science, 4703:17–41, 2007.
3. E. J. Molinelli, A. Korukut, et al. "Perturbation biology: inferring signaling networks in cellular systems." *PLoS Computational Biology*, 9(12):e1003290, Dec 2013.
4. J. Allen, W. de Beaumont, L. Galescu, and C. M. Teng. "Complex event extraction using DRUM." 2015.
5. M. A. Valenzuela-Escarrega, G. Hahn-Powell, T. Hicks, and M. Surdeanu. "A domain-independent rule-based framework for event extraction." In *Proc. 53rd Annual Meeting of the ACL-IJCNLP*, 2015.
6. D. McDonald et al., "Extending Biology Models with Deep NLP over Scientific Articles." *Workshops at the 30th AAAI Conference on Artificial Intelligence*, 2016.
7. C. F. Lopez\*, J. L. Muhlich\*, J. A. Bachman\*, and P. K. Sorger. "Programming biological models in python using PySB." *Molecular Systems Biology*, 9(1):646–646, Apr 2014.