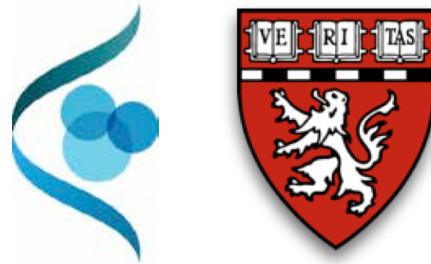


An automatically assembled knowledge graph from
literature-extracted biomedical knowledge
with human-machine dialogue to support discovery



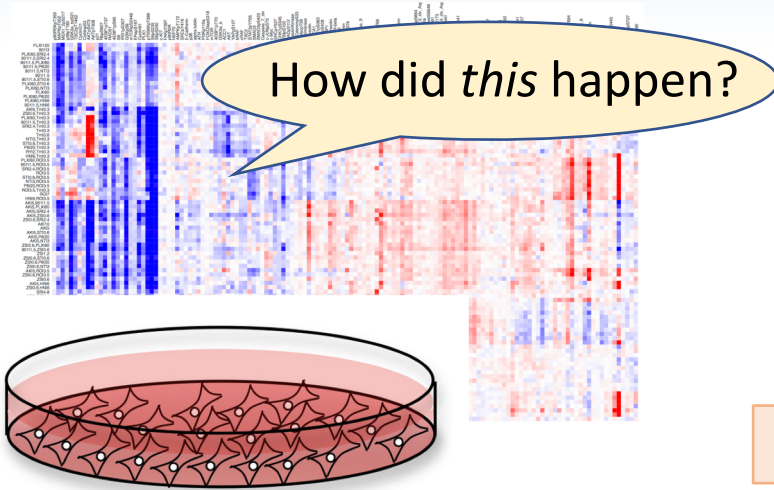
Benjamin M. Gyori
Laboratory of Systems Pharmacology
Harvard Medical School

KGC HCLS

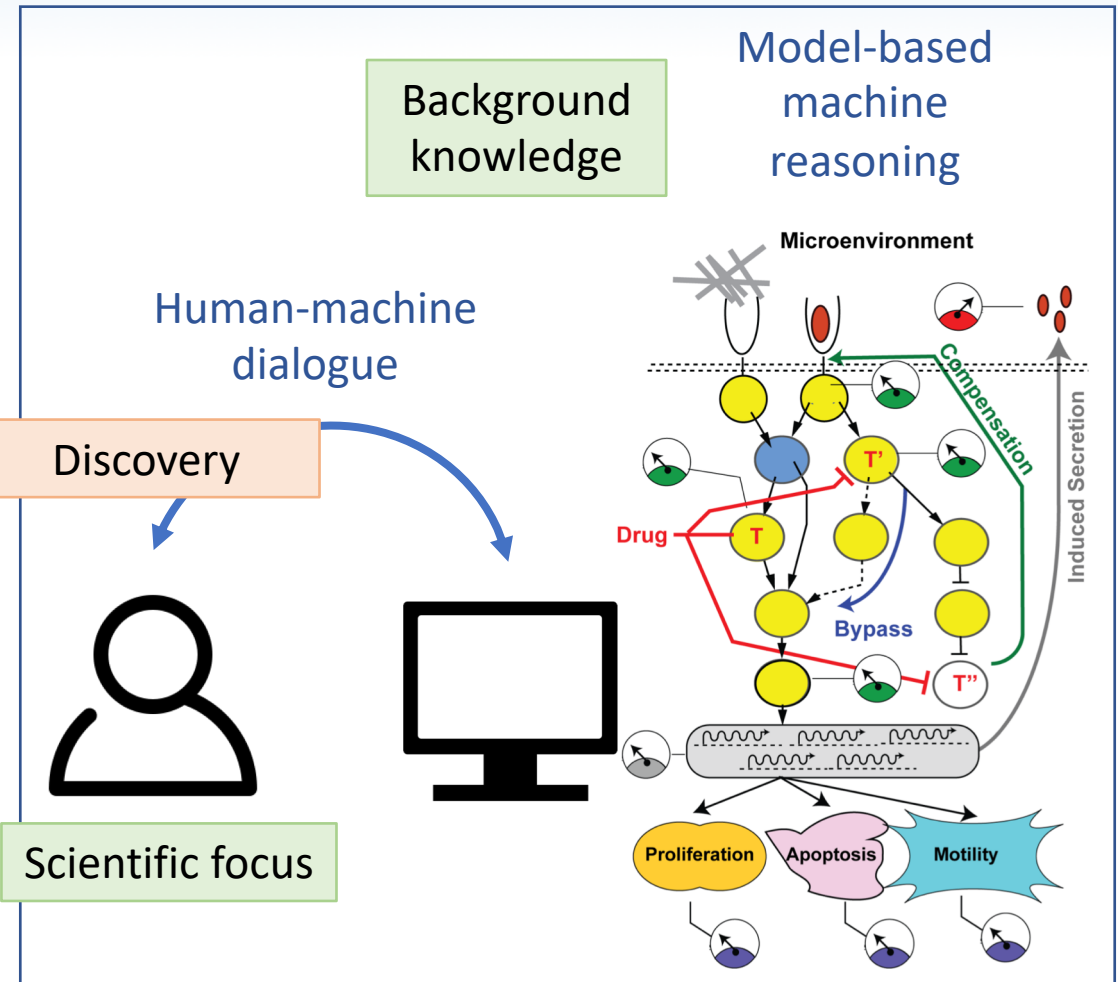
5/2/2022

Human-machine collaboration for discovery and decision making

Scientific discovery

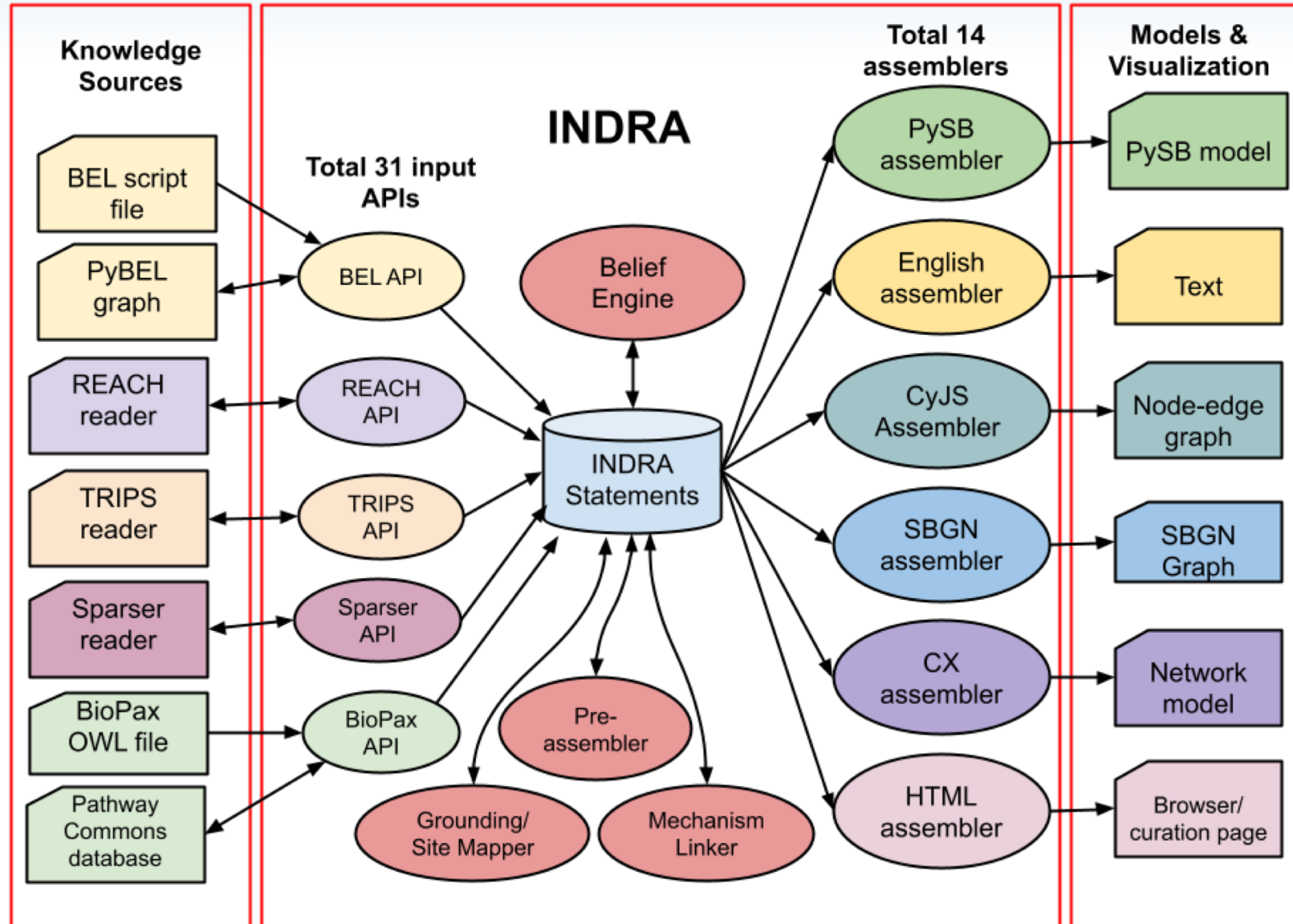
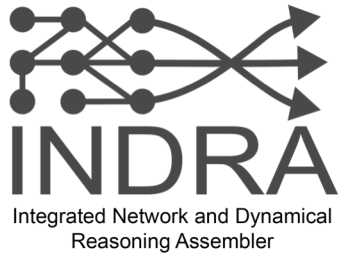


Clinical decision making



Collaborative problem solving

INDRA: Integrated Network and Dynamical Reasoning Assembler



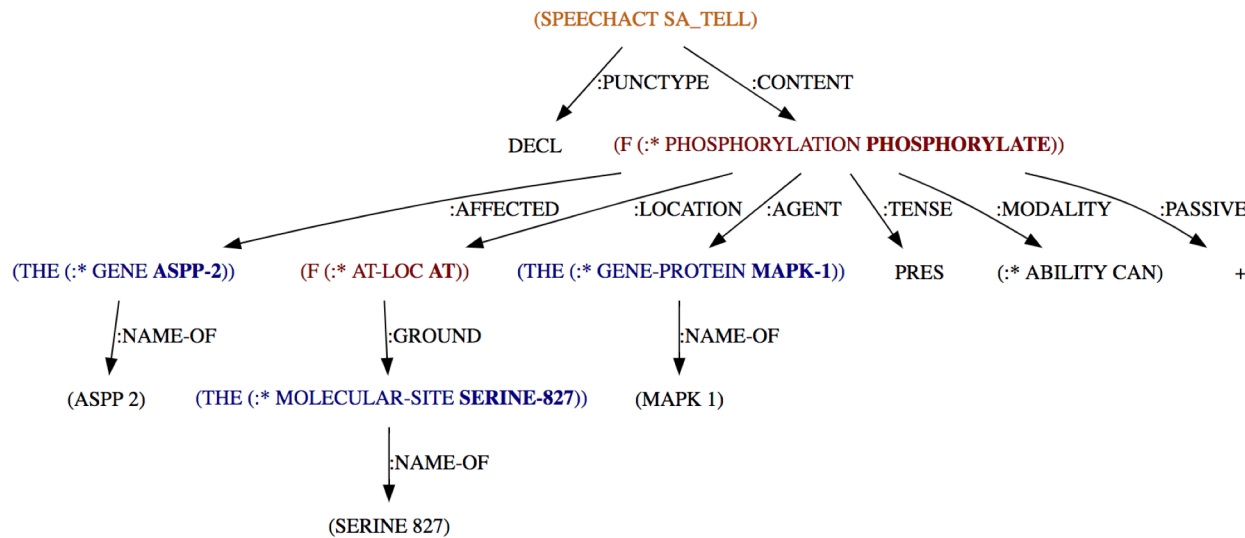
INDRA integrates pathway databases and other structured sources

Database / Exchange format	Module	Reference
PathwayCommons / BioPax	<code>indra.sources.biopax</code>	http://pathwaycommons.org/ http://www.biopax.org/
Large Corpus / BEL	<code>indra.sources.bel</code>	https://github.com/pybel/pybel https://github.com/OpenBEL
Signor	<code>indra.sources.signor</code>	https://signor.uniroma2.it/
BioGRID	<code>indra.sources.biogrid</code>	https://thebiogrid.org/
Target Affinity Spectrum	<code>indra.sources.tas</code>	https://doi.org/10.1101/358978
HPRD	<code>indra.sources.hprd</code>	http://www.hprd.org
TRRUST	<code>indra.sources.trrust</code>	https://www.grnpedia.org/trrust/
Phospho.ELM	<code>indra.sources.phosphoelm</code>	http://phospho.elm.eu.org/
VirHostNet	<code>indra.sources.virhostnet</code>	http://virhostnet.prabi.fr/
CTD	<code>indra.sources.ctd</code>	http://ctdbase.org
DrugBank	<code>indra.sources.drugbank</code>	https://www.drugbank.ca/
OmniPath	<code>indra.sources.omnipath</code>	https://omnipathdb.org/
DGI	<code>indra.sources.dgi</code>	https://www.dgidb.org/
CRoG	<code>indra.sources.crog</code>	https://github.com/chemical-roles/chemical-roles
CREEDS	<code>indra.sources.creeds</code>	https://maayanlab.cloud/CREEDS/
UbiBrowser	<code>indra.sources.ubibrowser</code>	http://ubibrowser.ncpsb.org.cn/



INDRA is integrated with multiple machine reading systems

ASPP2 can be phosphorylated at serine 827 by MAPK1.

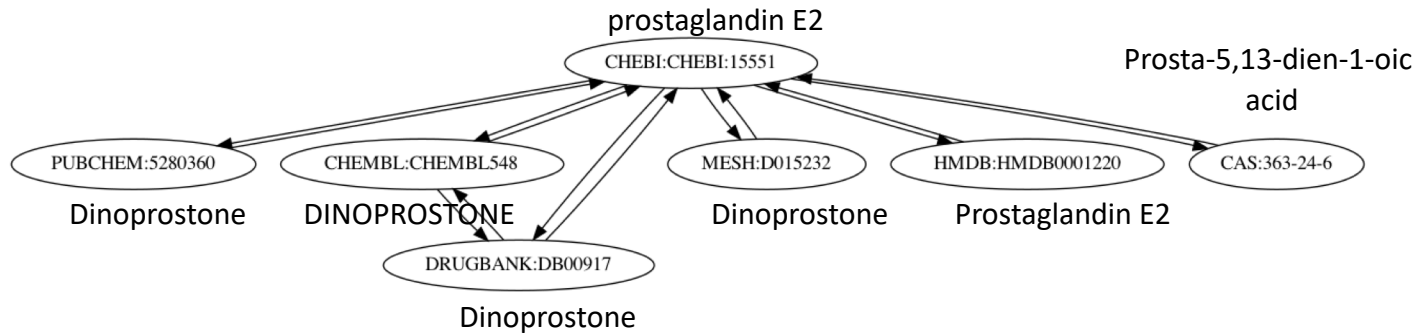


(Allen et al., 2015)

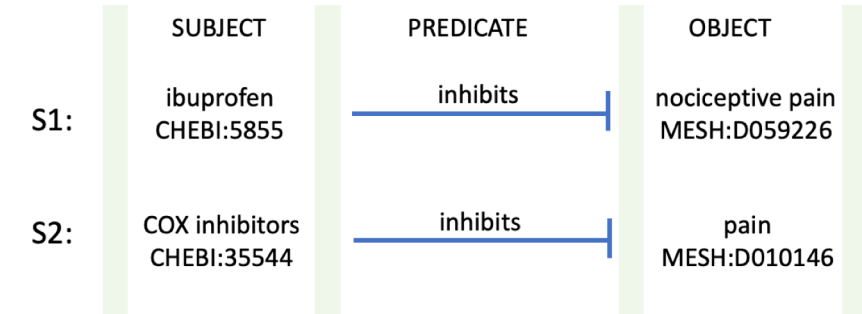
Reader	Module	Reference
TRIPS/DRUM	indra.sources.trips	http://trips.ihmc.us/parser/cgi/drum
REACH	indra.sources.reach	https://github.com/clulab/reach
Sparser	indra.sources.sparser	https://github.com/ddmcdonald/sparser
Eidos	indra.sources.eidos	https://github.com/clulab/eidos
TEES	indra.sources.tees	https://github.com/jbjorne/TEES
MedScan	indra.sources.medscan	https://doi.org/10.1093/bioinformatics/btg207
RLIMS-P	indra.sources.rlimsp	https://research.bioinformatics.udel.edu/rlimsp
ISI/AMR	indra.sources.isi	https://github.com/sgarg87/big_mech_isi_gg
Geneways	indra.sources.geneways	https://www.ncbi.nlm.nih.gov/pubmed/15016385
GNBR	indra.sources.gnbr	https://zenodo.org/record/3459420

INDRA assembly standardizes concepts and finds relationships between causal statements

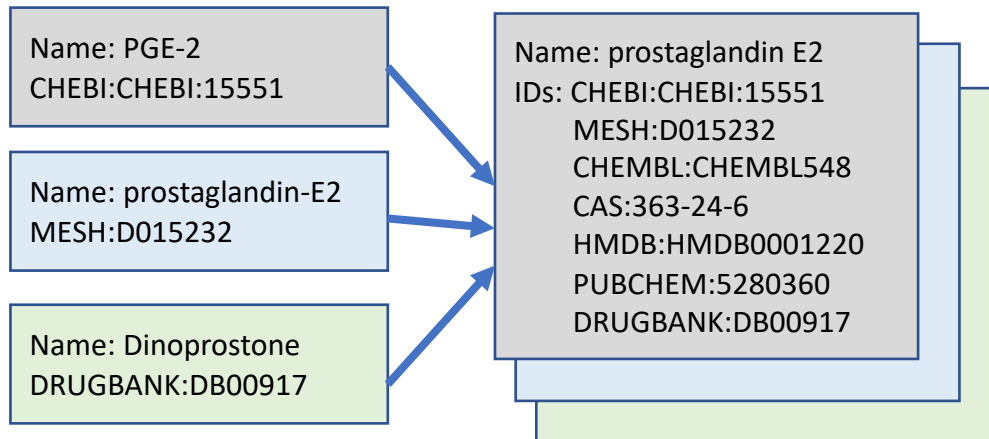
INDRA ontology graph cross-references



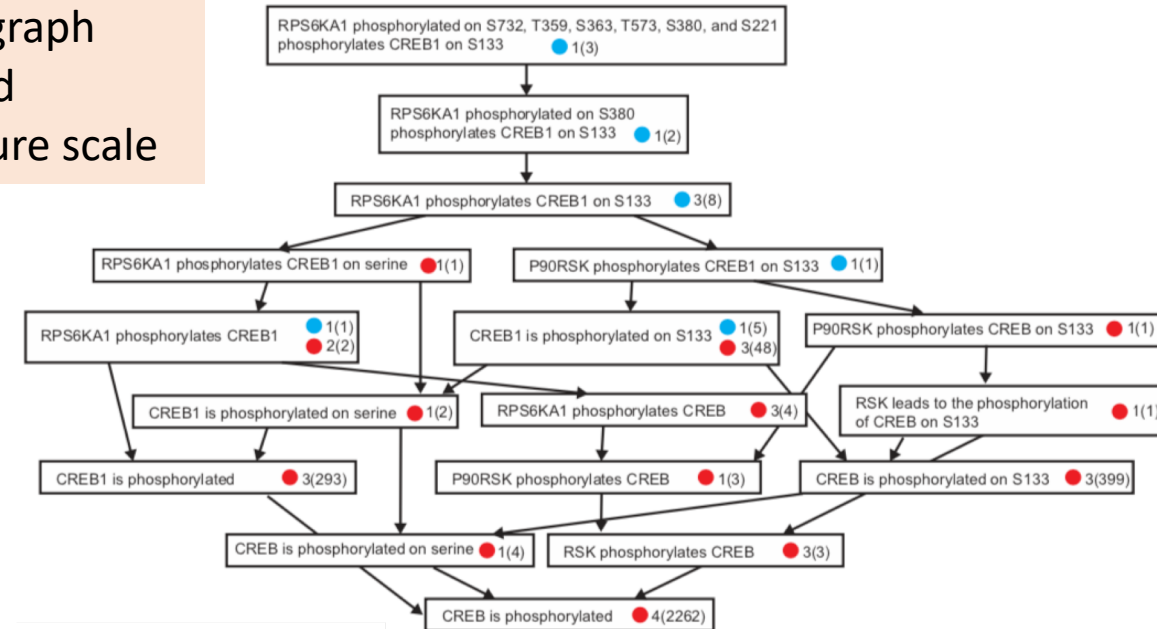
INDRA ontology family relationships induce relationships between statements



Concept normalization from diverse sources based on cross-references



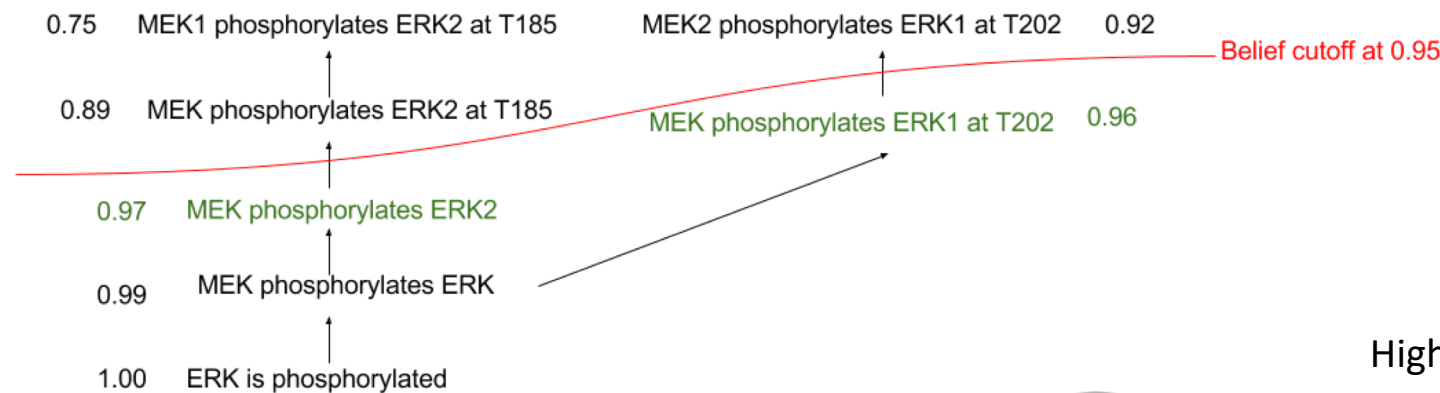
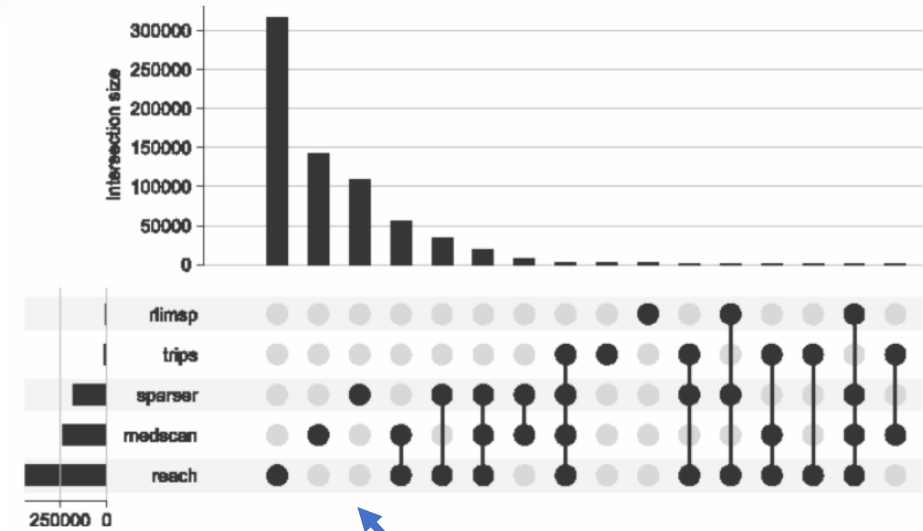
Statement support graph recovered at literature scale



INDRA uses a probability model to determine a belief score

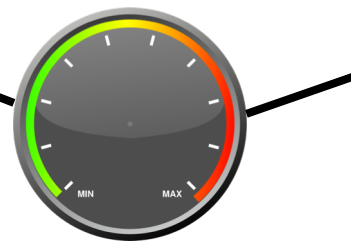
Estimate reliability of Statements probabilistically by:

- Calculating joint probability of an incorrect statement given repeated extractions from different sentences
- Combining results from different readers
- Propagating error estimates through the network of related statements



Machine-learning approaches based on curation data

Low cutoff: allows uncorroborated results into model



High cutoff: the mechanism needs evidence from multiple sentences / systems / with database backup to reach high belief level

Large-scale continuous assembly of causal mechanisms with INDRA

Biomed
literat
text m

Exper
curate
databa

Results

I found statements that are from papers with MeSH ID C000657245, are not only from medscan, and have an agent where HGNC=11876 with role=OBJECT

- ERG affects TMPRSS2 | 295 8
- ACE2 affects TMPRSS2 | 1 72 17
- camostat affects TMPRSS2 | 6 41
- Camostat inhibits TMPRSS2. | 6 37
- Camostat inhibits TMPRSS2. | 6 37

reach	eid	text	id
		However, the mechanism of inhibition of TMPRSS2 by Camostat mesylate, per se, has not been clearly elucidated.	32469279
	eid	Next , using the established assay in 384-well format , we tested the inhibition of TMPRSS2 by camostat , nafamostat and gabexate (Figure 3) .	32596694
reach		Camostat mesylate (a serine protease inhibitor), well documented for the treatment of cancer, pancreatitis and liver fibrosis, has recently been shown to inhibit TMPRSS2 and entry of SARS-CoV-2 virus into the cells (Hoffmann et al., 2020).	32469279
	eid	It was reported the ten-fold reduction of SARS-CoV titers in the Calu-3 cells with the inhibition of TMPRSS2 by camostat (Shulla et al. 2011) .	32546018

Select error type... Optional description (240 chars) Submit

Prior Curations

7/4/2020, 5:18:20 PM	ben.gyori@gmail.com	correct	No text given.	EMMAA
----------------------	---------------------	---------	----------------	-------

Executable models

Our Sources:

Databases

- Causal Bionet
- PathwayCommons
- Signor
- CTD
- BioGRID
- DrugBank
- VirHostNet
- TRRUST
- Phosphosite Plus
- BEL
- HPRD
- Target Affinity Spectrum
- Phospho.ELM

Reading

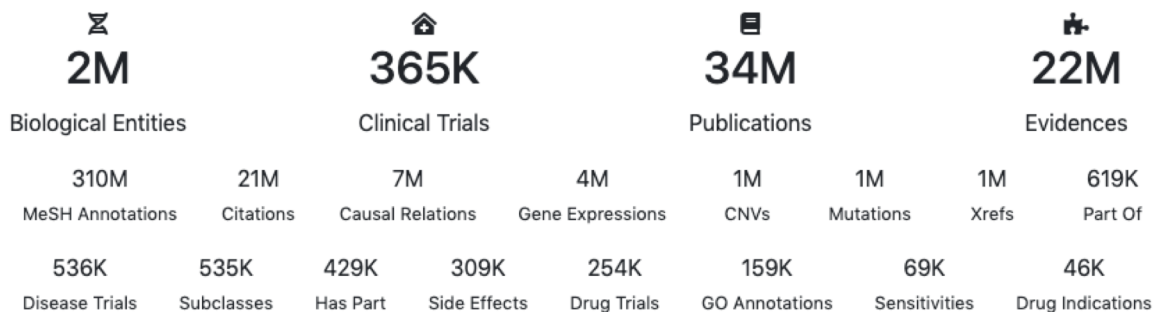
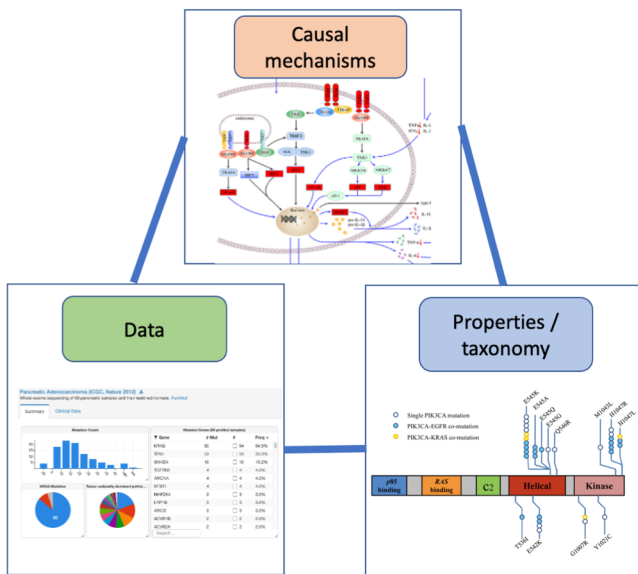
- ISI/AMR
- RLIMS-P
- Eidos
- TRIPS/DRUM
- MedScan
- Sparsen
- REACH

- psp cbn pc bel_lc signor biogrid tas hprd trrust
- ctd vhn pe drugbank omnipath conib crog dgi minerva creeds

knowledge assembly

INDRA Context Graph Extension (CoGEx) – knowledge graph

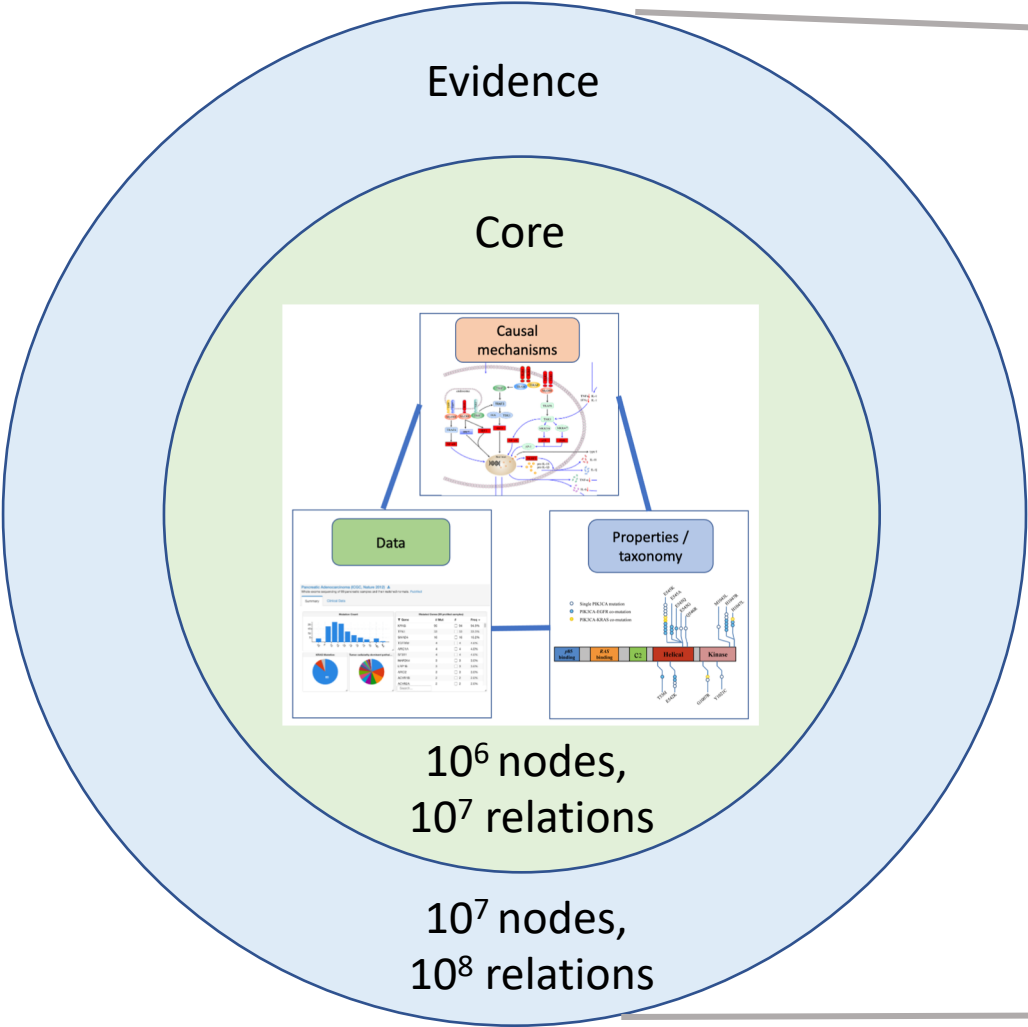
Idea: combine causal mechanisms assembled by INDRA with relations representing data and properties/taxonomy in a knowledge graph



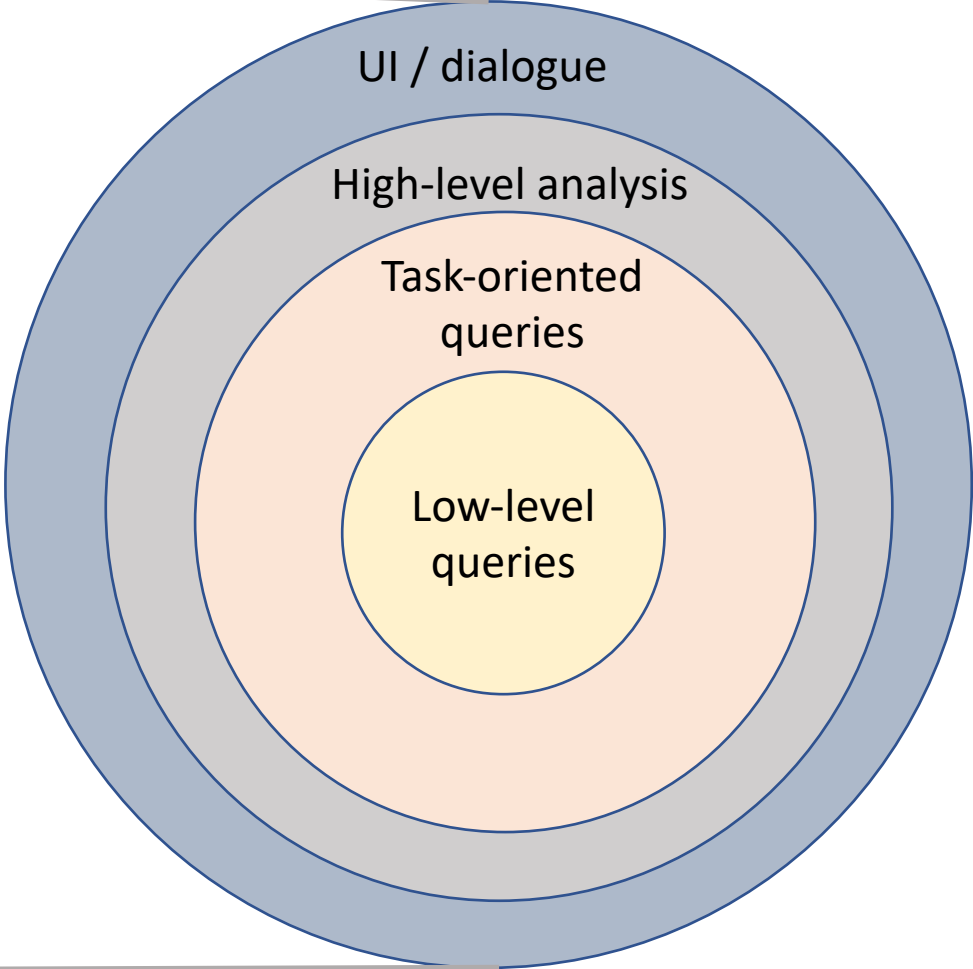
Knowledge source	Relation	Description
INDRA	indra_rel	The source regulates or interacts with the target according to an INDRA Statement.
INDRA Ontology	isa / partof	The source node is a subclass or part of the target node.
Gene Ontology	associated_with	The gene represented by the source is associated with the GO term represented by the target.
Bgee	expressed_in	The gene represented by the source is expressed in the tissue/cell type represented by the target.
CCLE	mutated_in / copy_number_altered_in	The gene represented by the source is mutated or its copy number is altered in the cancer cell line represented by the target.
CCLE	sensitive_to	The cancer cell line represented by the source is sensitive to the drug represented by the target.
ClinicalTrials.gov	tested_in	The drug represented by the source is tested in the clinical trial represented by the target.
ClinicalTrials.gov	has_trial	The disease/condition represented by the source has a clinical trial represented by the target.
ChEMBL	has_indication	The chemical represented by the source has been studied for use against the indication represented by the target.
SIDER	has_side_effect	The chemical represented by the source has a side effect represented by the target.
Reactome	haspart	The pathway represented by the source node contains the gene represented by the target node.
WikiPathways	haspart	The pathway represented by the source node contains the gene represented by the target node.

INDRA CoGEx Integrated analysis approach

Content



Queries / analysis

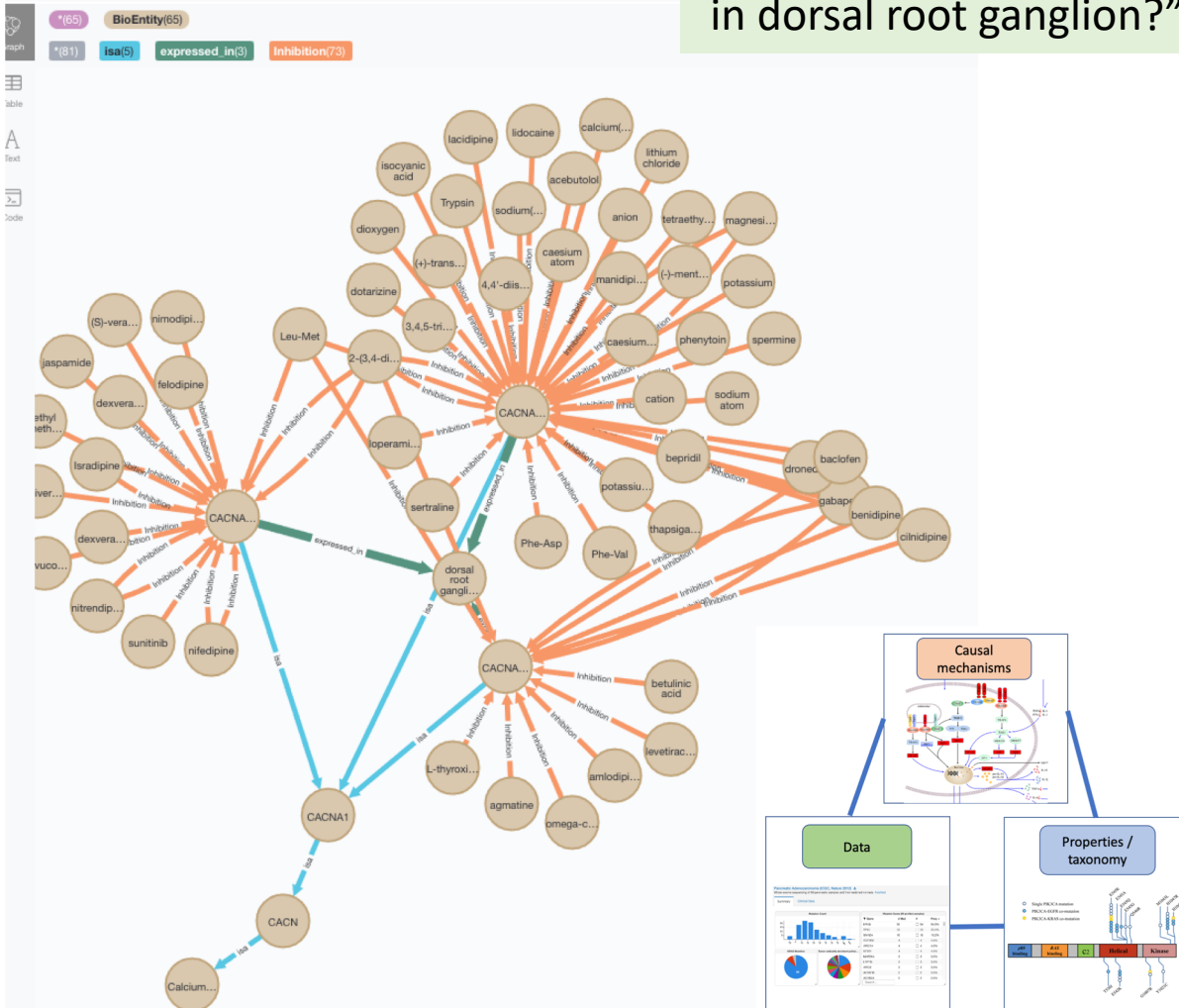


Querying the INDRA CoGEx KG programmatically

```

1 MATCH p1=(n)-[:isa*1..]->({id: 'fplx:Calcium_channels'})
2 MATCH p2=(n)-[:expressed_in]->({id: 'uberon:0000044'})
3 MATCH p3=(d)-[stmt:Inhibition]->(n)
4 MATCH p4=(d)-[:isa*1..]->({id: 'chebi:24431'})
5 RETURN p1, p2, p3
6 LIMIT 5000
    
```

“What drugs target calcium channels that are expressed in dorsal root ganglion?”



discovery.indra.bio/apidocs

INDRA CoGEx Query API ^{1.0}

[Base URL: /]
/swagger.json

REST API for INDRA CoGEx queries

CoGEx Queries Queries for INDRA CoGEx

POST /api/get_diseases_for_trial Return the diseases for the given trial

POST /api/get_drugs_for_side_effect Return the drugs for the given side effect

POST /api/get_drugs_for_target Return the drugs targeting the given protein

POST /api/get_drugs_for_trial Return the drugs for the given trial

POST /api/get_evidences_for_mesh Return the evidence objects for the given MESH term

POST /api/get_evidences_for_stmt_hash Return the matching evidence objects for the given statement

POST /api/get_evidences_for_stmt_hashes Return the matching evidence objects for the given statements

POST /api/get_genes_for_go_term Return the genes associated with the given GO term

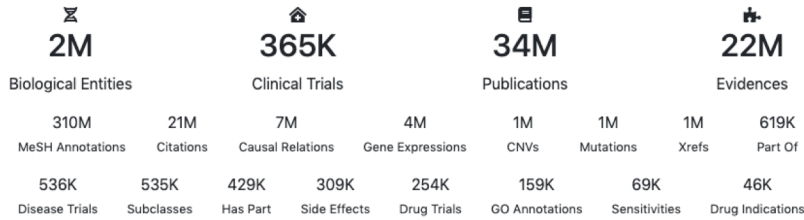
POST /api/get_genes_for_pathway Return the genes for the given pathway

POST /api/get_genes_in_tissue Return the genes in the given tissue

INDRA Biomedical Discovery Engine using the CoGEx KG

INDRA Biomedical Discovery Engine

The INDRA Biomedical Discovery Engine is built on INDRA CoGEx, a graph database integrating causal relations, ontological relations, properties, and data, assembled at scale automatically from the scientific literature and structured sources.



Apps and Services using INDRA CoGEx

Discrete Gene List Analysis

Performs gene set enrichment analysis using INDRA CoGEx.

[Click to Start](#)

Signed Gene List Analysis

Perform signed gene set enrichment analysis using INDRA CoGEx and the Reverse Causal Reasoning algorithm.

[Click to Start](#)

Continuous Gene Set Enrichment Analysis

Perform gene set enrichment analysis on continuous data using INDRA CoGEx.

[Click to Start](#)

Discrete Metabolite List Analysis

Perform metabolite set enrichment analysis using INDRA CoGEx.

[Click to Start](#)

Directed Curation

GO Pathway Curator

Curate networks induced by genes associated with a given

Topic Curator

Curate statements from publications annotated with the given MeSH term

Publication Curator

Curate statements from a given publication.

Subnetwork Curator

Curate statements for

[Chat with INDRA CoGEx](#)

Gene List Analysis

Query

[HGNC:11133 \(SNAP29\)](#) [HGNC:4910 \(HIFA1\)](#) [HGNC:2734 \(DDX1\)](#) [HGNC:1116 \(BSG\)](#) [HGNC:11432 \(STX17\)](#) [HGNC:18348 \(TICAM1\)](#)
[HGNC:23785 \(PIKFYVE\)](#) [HGNC:4931 \(HLA-A\)](#) [HGNC:1697 \(CD74\)](#) [HGNC:16404 \(IL17F\)](#) [HGNC:5432 \(IFNAR1\)](#) [HGNC:29517 \(DHX58\)](#)
[HGNC:11876 \(TMPPRS2\)](#) [HGNC:9788 \(RAB7A\)](#) [HGNC:28948 \(TMEM41B\)](#) [HGNC:4983 \(HMOB1\)](#) [HGNC:6727 \(LYGE\)](#)
[HGNC:11985 \(TOMM70\)](#) [HGNC:6395 \(KPNA2\)](#) [HGNC:14374 \(NLRP1\)](#) [HGNC:18669 \(PAL51\)](#) [HGNC:12647 \(VAMP8\)](#)
[HGNC:20593 \(VPS39\)](#) [HGNC:10819 \(SGTA\)](#) [HGNC:8912 \(PHB\)](#) [HGNC:6122 \(IRF7\)](#) [HGNC:613 \(APOE\)](#) [HGNC:9253 \(PPIA\)](#)
[HGNC:5985 \(IL17RA\)](#) [HGNC:6118 \(IRF3\)](#) [HGNC:6378 \(KLRD1\)](#) [HGNC:11120 \(SMPD1\)](#) [HGNC:6374 \(KLRG1\)](#) [HGNC:10498 \(S100A8\)](#)
[HGNC:11878 \(TMPPRS4\)](#) [HGNC:4962 \(HLA-E\)](#) [HGNC:8668 \(FURIN\)](#) [HGNC:6120 \(IRF5\)](#) [HGNC:6021 \(IL6ST\)](#) [HGNC:11584 \(TBK1\)](#)
[HGNC:12713 \(VPS41\)](#) [HGNC:30308 \(PHB2\)](#) [HGNC:6769 \(SMAD3\)](#) [HGNC:20820 \(TPCN2\)](#) [HGNC:6019 \(IL6R\)](#) [HGNC:18873 \(IFIH1\)](#)
[HGNC:4932 \(HLA-B\)](#) [HGNC:8004 \(NRP1\)](#) [HGNC:2637 \(CTSL\)](#) [HGNC:1119 \(BST2\)](#) [HGNC:18358 \(IL17RC\)](#) [HGNC:6148 \(ITGAL\)](#)
[HGNC:7087 \(CIITA\)](#) [HGNC:5433 \(IFNAR2\)](#) [HGNC:11849 \(TLR3\)](#) [HGNC:5981 \(IL17A\)](#) [HGNC:6018 \(IL6\)](#)

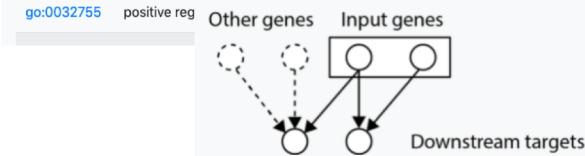
Over-Representation Analysis

[GO](#) [Reactome](#) [WikiPathways](#) [INDRA Downstream](#) [INDRA Upstream](#)

These results are acquired by running over-representation analysis using Fisher's exact test and correcting using `fdr_bh` and $\alpha=0.05$ on the genes annotated to terms in the [Gene Ontology](#) via the [Gene Ontology Annotations Database](#).

Show entries Search:

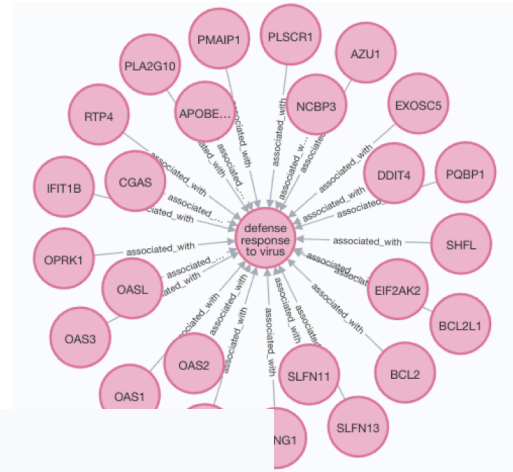
Over-Representation Analysis
[GO](#) [Reactome](#) [WikiPathways](#) [INDRA Downstream](#) [INDRA Upstream](#)



These results are acquired by running over-representation analysis using Fisher's exact test and correcting using `fdr_bh` and $\alpha=0.05$ on genes causally downstream in one step from all entities in the INDRA Database.

Show entries Search:

CURIE	Name	p-value	q-value
go:0006955	immune response	4.90e-41	1.47e-36
go:0001816	cytokine production	1.44e-39	2.16e-35
mesh:D007239	Infections	1.80e-37	1.80e-33
go:0006914	autophagy	1.25e-36	9.41e-33
hgnc:5438	IFNG	2.53e-36	1.52e-32



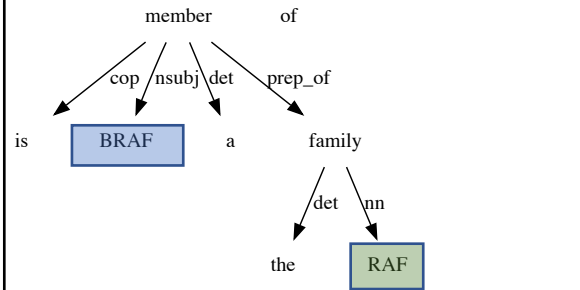
Human-machine dialogue

Example: Is BRAF a member of the RAF family?

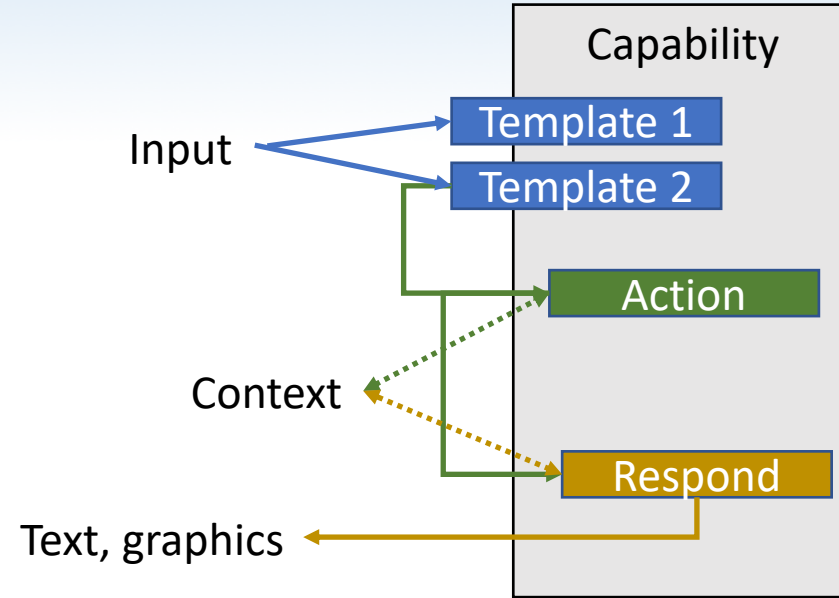
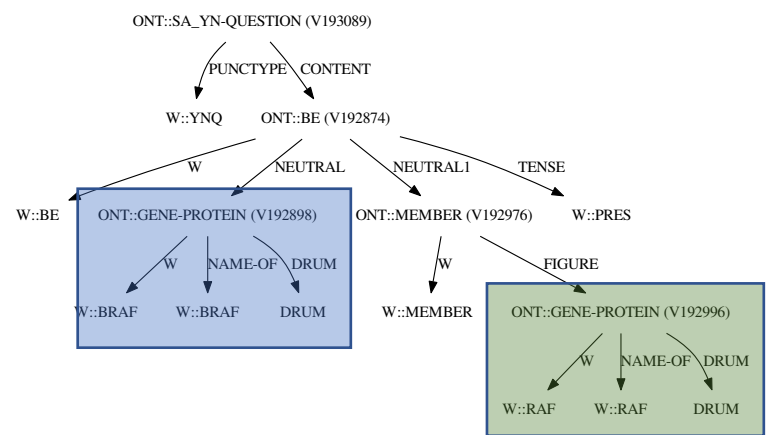
Regular expression templates

is **BRAF** a member of the **RAF** family

Syntactic dependency graph



Semantic interpretation graph

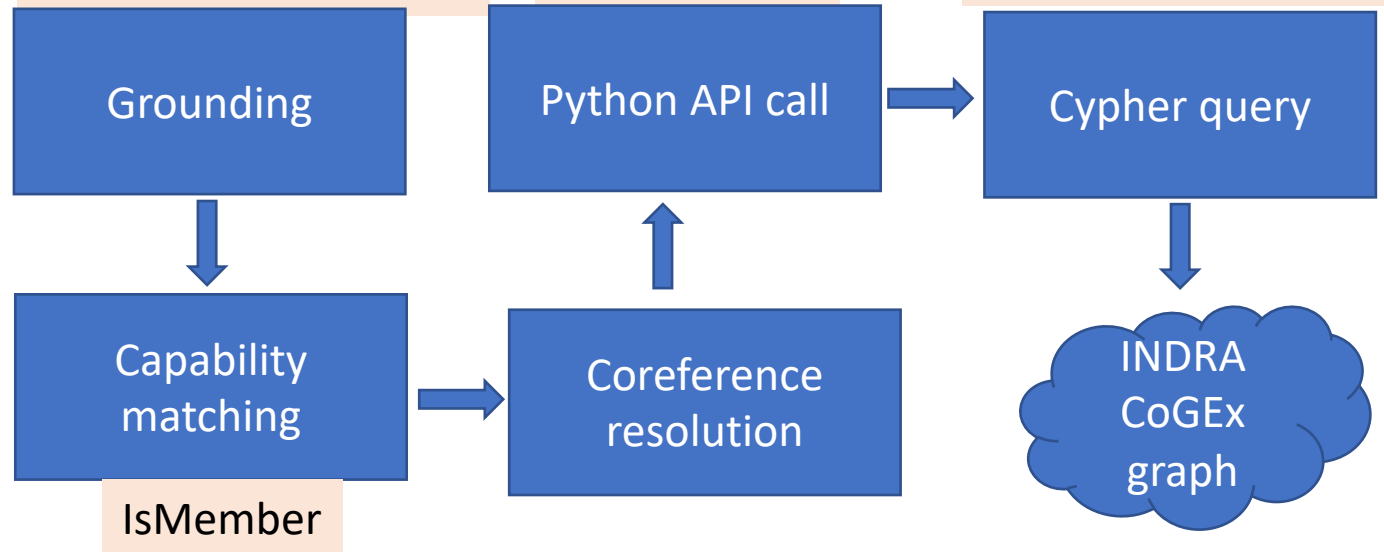


“BRAF” -> HGNC:1097

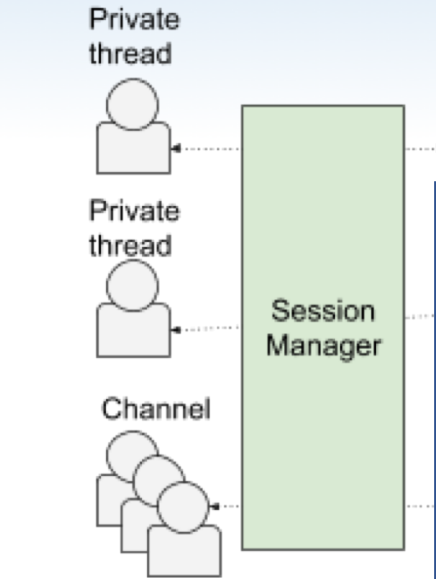
“RAF” -> FPLX:RAF

isa_or_partof

MATCH gene:BioEntity ...



Human-machine dialogue



ben.gyori 1:51 PM
what are the members of calcium channels?

clare-test APP 1:52 PM
I know about 41 members of Calcium_channels: CACNA1H, CACNA1I, CACNA1S, CACNA2D1, CACNA2D2, CACNA2D3, CACNA2D4, CACNA2D5, CACNA2D6, CACNA2D7, CACNA2D8, CACNA2D9, CACNA2D10, CACNA2D11, CACNA2D12, CACNA2D13, CACNA2D14, CACNA2D15, CACNA2D16, CACNA2D17, CACNA2D18, CACNA2D19, CACNA2D20, CACNA2D21, CACNA2D22, CACNA2D23, CACNA2D24, CACNA2D25, CACNA2D26, CACNA2D27, CACNA2D28, CACNA2D29, CACNA2D30, CACNA2D31, CACNA2D32, CACNA2D33, CACNA2D34, CACNA2D35, CACNA2D36, CACNA2D37, CACNA2D38, CACNA2D39, CACNA2D40, CACNA2D41, CACNA2D42, CACNA2D43, CACNA2D44, CACNA2D45, CACNA2D46, CACNA2D47, CACNA2D48, CACNA2D49, CACNA2D50, CACNA2D51, CACNA2D52, CACNA2D53, CACNA2D54, CACNA2D55, CACNA2D56, CACNA2D57, CACNA2D58, CACNA2D59, CACNA2D60, CACNA2D61, CACNA2D62, CACNA2D63, CACNA2D64, CACNA2D65, CACNA2D66, CACNA2D67, CACNA2D68, CACNA2D69, CACNA2D70, CACNA2D71, CACNA2D72, CACNA2D73, CACNA2D74, CACNA2D75, CACNA2D76, CACNA2D77, CACNA2D78, CACNA2D79, CACNA2D80, CACNA2D81, CACNA2D82, CACNA2D83, CACNA2D84, CACNA2D85, CACNA2D86, CACNA2D87, CACNA2D88, CACNA2D89, CACNA2D90, CACNA2D91, CACNA2D92, CACNA2D93, CACNA2D94, CACNA2D95, CACNA2D96, CACNA2D97, CACNA2D98, CACNA2D99, CACNA2D100.

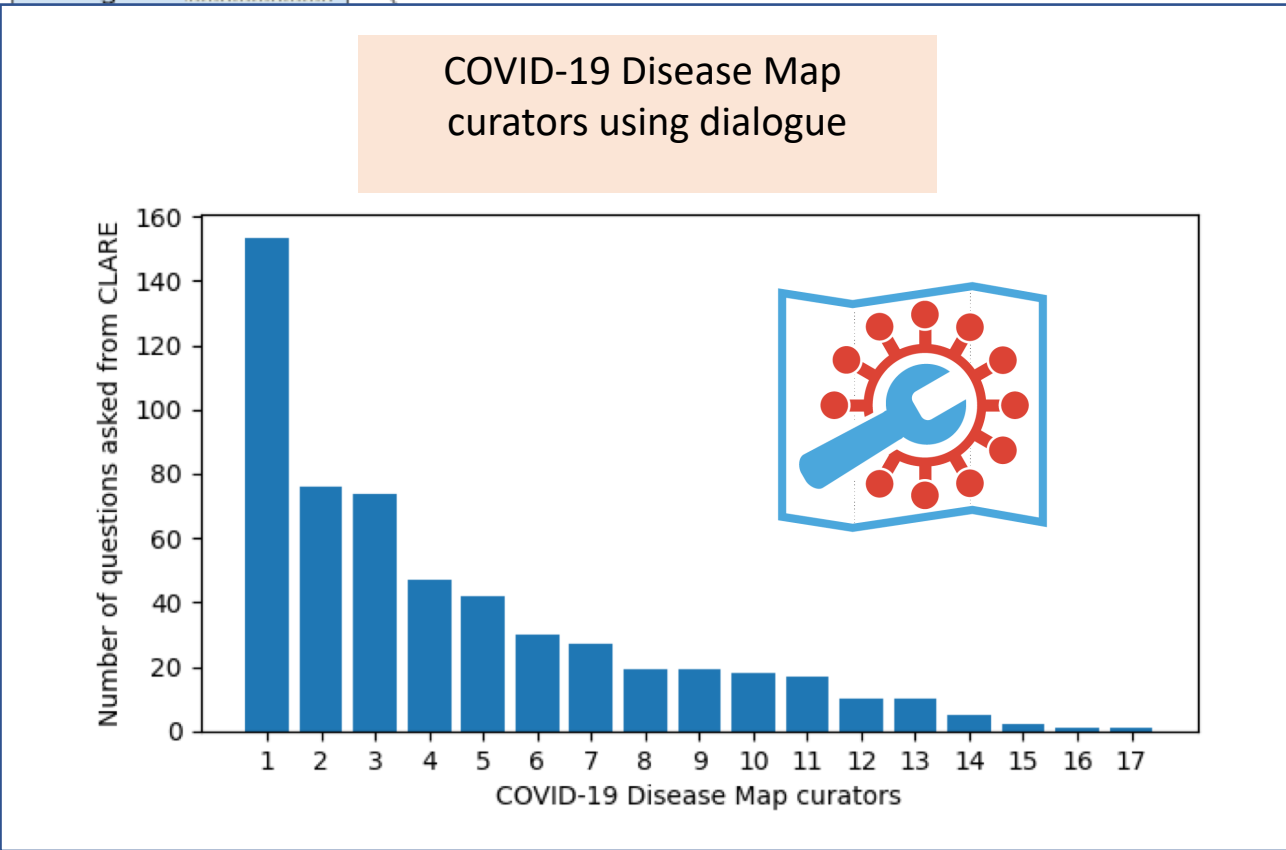
ben.gyori 1:52 PM
are any of those expressed in dorsal root ganglion?

clare-test APP 1:53 PM
Out of those, CACNA1A, CACNA1B, CACNA1C, CACNA1D, CACNA1E, CACNA1F, CACNA1G, CACNA1H, CACNA1I, CACNA1J, CACNA1K, CACNA1L, CACNA1M, CACNA1N, CACNA1O, CACNA1P, CACNA1Q, CACNA1R, CACNA1S, CACNA1T, CACNA1U, CACNA1V, CACNA1W, CACNA1X, CACNA1Y, CACNA1Z, CACNA2A, CACNA2B, CACNA2C, CACNA2D, CACNA2E, CACNA2F, CACNA2G, CACNA2H, CACNA2I, CACNA2J, CACNA2K, CACNA2L, CACNA2M, CACNA2N, CACNA2O, CACNA2P, CACNA2Q, CACNA2R, CACNA2S, CACNA2T, CACNA2U, CACNA2V, CACNA2W, CACNA2X, CACNA2Y, CACNA2Z, CACNA3A, CACNA3B, CACNA3C, CACNA3D, CACNA3E, CACNA3F, CACNA3G, CACNA3H, CACNA3I, CACNA3J, CACNA3K, CACNA3L, CACNA3M, CACNA3N, CACNA3O, CACNA3P, CACNA3Q, CACNA3R, CACNA3S, CACNA3T, CACNA3U, CACNA3V, CACNA3W, CACNA3X, CACNA3Y, CACNA3Z, CACNA4A, CACNA4B, CACNA4C, CACNA4D, CACNA4E, CACNA4F, CACNA4G, CACNA4H, CACNA4I, CACNA4J, CACNA4K, CACNA4L, CACNA4M, CACNA4N, CACNA4O, CACNA4P, CACNA4Q, CACNA4R, CACNA4S, CACNA4T, CACNA4U, CACNA4V, CACNA4W, CACNA4X, CACNA4Y, CACNA4Z, CACNA5A, CACNA5B, CACNA5C, CACNA5D, CACNA5E, CACNA5F, CACNA5G, CACNA5H, CACNA5I, CACNA5J, CACNA5K, CACNA5L, CACNA5M, CACNA5N, CACNA5O, CACNA5P, CACNA5Q, CACNA5R, CACNA5S, CACNA5T, CACNA5U, CACNA5V, CACNA5W, CACNA5X, CACNA5Y, CACNA5Z, CACNA6A, CACNA6B, CACNA6C, CACNA6D, CACNA6E, CACNA6F, CACNA6G, CACNA6H, CACNA6I, CACNA6J, CACNA6K, CACNA6L, CACNA6M, CACNA6N, CACNA6O, CACNA6P, CACNA6Q, CACNA6R, CACNA6S, CACNA6T, CACNA6U, CACNA6V, CACNA6W, CACNA6X, CACNA6Y, CACNA6Z, CACNA7A, CACNA7B, CACNA7C, CACNA7D, CACNA7E, CACNA7F, CACNA7G, CACNA7H, CACNA7I, CACNA7J, CACNA7K, CACNA7L, CACNA7M, CACNA7N, CACNA7O, CACNA7P, CACNA7Q, CACNA7R, CACNA7S, CACNA7T, CACNA7U, CACNA7V, CACNA7W, CACNA7X, CACNA7Y, CACNA7Z, CACNA8A, CACNA8B, CACNA8C, CACNA8D, CACNA8E, CACNA8F, CACNA8G, CACNA8H, CACNA8I, CACNA8J, CACNA8K, CACNA8L, CACNA8M, CACNA8N, CACNA8O, CACNA8P, CACNA8Q, CACNA8R, CACNA8S, CACNA8T, CACNA8U, CACNA8V, CACNA8W, CACNA8X, CACNA8Y, CACNA8Z, CACNA9A, CACNA9B, CACNA9C, CACNA9D, CACNA9E, CACNA9F, CACNA9G, CACNA9H, CACNA9I, CACNA9J, CACNA9K, CACNA9L, CACNA9M, CACNA9N, CACNA9O, CACNA9P, CACNA9Q, CACNA9R, CACNA9S, CACNA9T, CACNA9U, CACNA9V, CACNA9W, CACNA9X, CACNA9Y, CACNA9Z, CACNA10A, CACNA10B, CACNA10C, CACNA10D, CACNA10E, CACNA10F, CACNA10G, CACNA10H, CACNA10I, CACNA10J, CACNA10K, CACNA10L, CACNA10M, CACNA10N, CACNA10O, CACNA10P, CACNA10Q, CACNA10R, CACNA10S, CACNA10T, CACNA10U, CACNA10V, CACNA10W, CACNA10X, CACNA10Y, CACNA10Z, CACNA11A, CACNA11B, CACNA11C, CACNA11D, CACNA11E, CACNA11F, CACNA11G, CACNA11H, CACNA11I, CACNA11J, CACNA11K, CACNA11L, CACNA11M, CACNA11N, CACNA11O, CACNA11P, CACNA11Q, CACNA11R, CACNA11S, CACNA11T, CACNA11U, CACNA11V, CACNA11W, CACNA11X, CACNA11Y, CACNA11Z, CACNA12A, CACNA12B, CACNA12C, CACNA12D, CACNA12E, CACNA12F, CACNA12G, CACNA12H, CACNA12I, CACNA12J, CACNA12K, CACNA12L, CACNA12M, CACNA12N, CACNA12O, CACNA12P, CACNA12Q, CACNA12R, CACNA12S, CACNA12T, CACNA12U, CACNA12V, CACNA12W, CACNA12X, CACNA12Y, CACNA12Z, CACNA13A, CACNA13B, CACNA13C, CACNA13D, CACNA13E, CACNA13F, CACNA13G, CACNA13H, CACNA13I, CACNA13J, CACNA13K, CACNA13L, CACNA13M, CACNA13N, CACNA13O, CACNA13P, CACNA13Q, CACNA13R, CACNA13S, CACNA13T, CACNA13U, CACNA13V, CACNA13W, CACNA13X, CACNA13Y, CACNA13Z, CACNA14A, CACNA14B, CACNA14C, CACNA14D, CACNA14E, CACNA14F, CACNA14G, CACNA14H, CACNA14I, CACNA14J, CACNA14K, CACNA14L, CACNA14M, CACNA14N, CACNA14O, CACNA14P, CACNA14Q, CACNA14R, CACNA14S, CACNA14T, CACNA14U, CACNA14V, CACNA14W, CACNA14X, CACNA14Y, CACNA14Z, CACNA15A, CACNA15B, CACNA15C, CACNA15D, CACNA15E, CACNA15F, CACNA15G, CACNA15H, CACNA15I, CACNA15J, CACNA15K, CACNA15L, CACNA15M, CACNA15N, CACNA15O, CACNA15P, CACNA15Q, CACNA15R, CACNA15S, CACNA15T, CACNA15U, CACNA15V, CACNA15W, CACNA15X, CACNA15Y, CACNA15Z, CACNA16A, CACNA16B, CACNA16C, CACNA16D, CACNA16E, CACNA16F, CACNA16G, CACNA16H, CACNA16I, CACNA16J, CACNA16K, CACNA16L, CACNA16M, CACNA16N, CACNA16O, CACNA16P, CACNA16Q, CACNA16R, CACNA16S, CACNA16T, CACNA16U, CACNA16V, CACNA16W, CACNA16X, CACNA16Y, CACNA16Z, CACNA17A, CACNA17B, CACNA17C, CACNA17D, CACNA17E, CACNA17F, CACNA17G, CACNA17H, CACNA17I, CACNA17J, CACNA17K, CACNA17L, CACNA17M, CACNA17N, CACNA17O, CACNA17P, CACNA17Q, CACNA17R, CACNA17S, CACNA17T, CACNA17U, CACNA17V, CACNA17W, CACNA17X, CACNA17Y, CACNA17Z.

ben.gyori 1:53 PM
are there any drugs targeting any of those?

clare-test APP 1:53 PM
I found the following drugs for these targets:

- CACNA1C: nifedipine, nitrendipine, Verapamil
- CACNA2D1: gabapentin
- RYR1: EPIRYANODINE
- RYR2: ALADORIAN



using INDRA CoGEX. using INDRA CoGEX and the Reverse Causal enrichment analysis continuous data usir

Discovery Engine

a graph database integrating causal scale automatically from the scientific

4M	22M
Publications	Evidences
1M	1M
1M	619K
159K	69K
159K	46K
159K	46K

Annotations

Minimize Chat Window

Ben: what does sorafenib inhibit?

CLARE: Overall, I know that sorafenib can inhibit, for instance, cell population proliferation, angiogenesis, RAF, ERK, apoptotic process, Carcinoma, Hepatocellular, VEGFR, BRAF, RAF1, and VEGF. ([View statements](#))

Ben: which of those are kinases?

CLARE: Out of those, BRAF, RAF1, KDR, FLT3, RET, PDGFRB, KIT, FLT4, MTOR, FLT1, EGFR, MAPK14, PDGFRA, ABL1, MAP3K20, CSF1R, DDR2, FGFR1,

Enter Message Send

Sequential exploration of causal mechanisms, properties and data through dialogue

bgyori 6:06 PM
what proteins might lead to the development of melanoma?

clare APP 6:06 PM
It looks like BRAF has a mutation frequency of 52% in melanoma so it could be a useful target.

bgyori 6:06 PM
is it a kinase?

clare APP 6:06 PM
Yes, BRAF is a kinase.

bgyori 6:06 PM
what does it phosphorylate?

clare APP 6:06 PM
Overall, I know that BRAF can phosphorylate, for instance, MEK, MAP2K1, ERK, MAP2K2, BRAF, MAPK7, PAX3, Pax3, BAD, and RAF1.

You can browse the full set of 177 statements here:
<https://s3.amazonaws.com/indrabot-results/49c73701-b9a8-4c12-97ad-9718f70df344.html>

Data

Properties /
taxonomy

Causal
mechanisms

bgyori 8:40 PM
do any of those activate EMT?

clare APP 8:41 PM
Out of those, I know that for instance, AKT, NFKappaB, ERK, STAT3, EGFR, SRC, MAPK, NANOG, MET, and TERT can activate epithelial to mesenchymal transition.

You can browse the full set of 72 statements here:
<https://s3.amazonaws.com/indrabot-results/2d48594c-ac79-4067-93bc-64e471962989.html>

bgyori 8:44 PM
are there any drugs that target BRAF?

clare APP 8:44 PM
I know that ALW-II-38-3, AZ-628, dabrafenib, Encorafenib, GDC-0879, GSK-718429A, GW-5074, L-779450, LY3009120, PD0325901, PD173955, PLX-4720, regorafenib, SB590885, sorafenib, staurosporine, TAK-632, and vemurafenib target BRAF.

bgyori 8:45 PM
are there any clinical trials attempting sorafenib?

clare APP 8:45 PM
I found the following clinical trials for sorafenib:

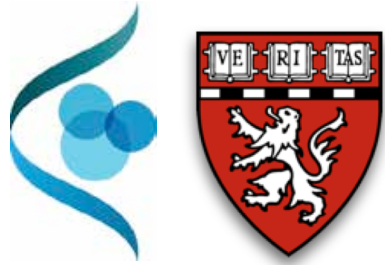
- preoperative trial of sorafenib in combination with cisplatin followed by paclitaxel for early stage breast cancer - terminated
- sorafenib and vinorelbine in treating women with stage iv breast cancer - completed
- phase i/ii trial of letrozole and sorafenib in post-menopausal hormone-receptor + breast cancer - active, not recruiting
- sorafenib and paclitaxel in treating patients with metastatic breast cancer - terminated

Causal
mechanisms

Causal
mechanisms

Data

Acknowledgments



Peter Sorger
Laura Maliszewski
Catherine Luria

Charlie Hoyt
Klas Karis
Diana Kolusheva
Samuel Bunga



John Bachman
Albert Steppi
Patrick Greene



Website: indralab.github.io
Twitter: @IndraSysBio

Funding

DARPA Young Faculty Award
(W911NF2010255)

