

1 Knowledge graph analytics platform with 2 LINCS and IDG for Parkinson's disease 3 target illumination

4 **Jeremy J Yang^{1,2,4}, Christopher R Gessner^{1,3}, Joel L Duerksen², Daniel Biber², Jessica L**
5 **Binder⁴, Murat Ozturk^{1,2}, Brian Foote², Robin McEntire², Kyle Stirling^{1,2}, Ying Ding^{1,2,5} and**
6 **David J Wild^{1,2}**

7 ¹*School of Informatics, Computing and Engineering, Indiana University, Bloomington, IN, USA.*

8 ²*Data2Discovery, Inc., Bloomington, IN, USA.*

9 ³*Chemical Abstracts Service, Columbus, OH, USA.*

10 ⁴*Dept. of Internal Medicine Translational Informatics Division, University of New Mexico,*
11 *Albuquerque, NM, USA.*

12 ⁵*School of Information, Dell Medical School, University of Texas, Austin, TX, USA.*

13 Corresponding author: **David Wild**

14 Email address: djwild@indiana.edu

15

16 Abstract

17 Background

18 LINCS, "Library of Integrated Network-based Cellular Signatures", and IDG, "Illuminating the
19 Druggable Genome", are both NIH projects and consortia that have generated rich datasets for
20 the study of the molecular basis of human health and disease. LINCS L1000 expression
21 signatures provide unbiased systems/omics experimental evidence. IDG provides compiled and
22 curated knowledge for illumination and prioritization of novel drug target hypotheses. Together,
23 these resources can support a powerful new approach to identifying novel drug targets for
24 complex diseases, such as Parkinson's disease (PD), which continues to inflict severe harm on
25 human health, and resist traditional research approaches.

26 Results

27 Integrating LINCS and IDG, we built the Knowledge Graph Analytics Platform (KGAP) to
28 support an important use case: identification and prioritization of drug target hypotheses for
29 associated diseases. The KGAP approach includes strong semantics interpretable by domain
30 scientists and a robust, high performance implementation of a graph database and related
31 analytical methods, using Neo4j. Illustrating the value of our approach, we investigated results
32 from queries relevant to PD. Approved PD drug indications from IDG's resource DrugCentral
33 were used as starting points for evidence paths exploring chemogenomic space via LINCS
34 expression signatures for associated genes, evaluated as target hypotheses by integration with
35 IDG. The KG-analytic scoring function was validated against a gold standard dataset of genes
36 associated with PD as elucidated, published mechanism-of-action drug targets, also from
37 DrugCentral. IDG's resource TIN-X was used to rank and filter KGAP results for novel PD
38 targets, and one, SYNGR3 (Synaptogyrin-3), was manually investigated further as a case study
39 and plausible new drug target for PD.

40 Conclusions

41 The synergy of LINCS and IDG, via KG methods, empowers graph analytic methods for the
42 investigation of the molecular basis of complex diseases, and specifically for identification and
43 prioritization of novel drug targets. The KGAP approach enables downstream applications via
44 integration with resources similarly aligned with modern KG methodology. The generality of the
45 approach indicates that KGAP is applicable to many disease areas, in addition to PD, the focus
46 of this paper.

47 Keywords

48 Knowledge graph, graph analytics, systems biology, drug discovery, drug target, druggable
49 genome, Parkinson's disease

50 Background

51 Integration of heterogeneous datasets is often essential for biomedical knowledge discovery,
52 where relevant evidence may derive from diverse subdomains that include but aren't limited to:
53 foundational biology, chemistry, clinical data, and social sciences of epidemiology and health
54 economics. Furthermore, there is a need for judicious selection of data types and datasets to
55 integrate, driven by scientific use cases, guided by applicability, accessibility, and veracity of
56 datasets. Accordingly, we have integrated LINCS and IDG to identify and prioritize novel drug
57 target hypotheses, via KG methods and tools. LINCS content includes assay results from
58 cultured and primary human cells treated with bioactive compounds (small molecules or
59 biologics) or genetic perturbations. IDG data utilized in this study include drug-disease
60 associations, gene-disease associations and bibliometric scores from literature text mining.

61 **Synergy of LINCS and IDG.** LINCS and IDG are NIH Common Fund(1) projects, chosen for
62 integration for specific applicability to drug target discovery. LINCS, Library of Networked Cell-
63 based Signatures, is described as a "System-Level Cataloging of Human Cells Response to
64 Perturbations"(2) and features experimental and computational methodology designed to
65 generate useful biomedical knowledge from a systems/omics framework. As a key example,
66 perturbagens (e.g., small molecules, ligands such as growth factors and cytokines, micro-
67 environments, or CRISPR gene over-expression and knockdowns representing disease
68 phenotypes) are characterized by proteomics, transcriptomics (RNA-seq), or biochemical and
69 imaging readouts. Therefore LINCS provides useful mapping from genome to phenome with
70 direct relevance to biomolecular mechanisms and therapeutic hypotheses. IDG(3), with its
71 Target Central Resource Database (TCRD) and data portal Pharos(4), integrates
72 heterogeneous datasets from IDG experimental centers and diverse external sources with a
73 clear purpose, to illuminate understudied ("dark") protein-coding genes as potential drug targets.
74 IDG is particularly strong in text mining of current biomedical literature and bibliometrics suited
75 for knowledge discovery and evidence evaluation, focused on the sciences of drug discovery
76 research. In addition, another featured resource of IDG is DrugCentral(5). DrugCentral provides
77 information on active ingredients, chemical entities, pharmaceutical products, and their
78 biological targets.

79 We have combined the complementary strengths of LINCS and IDG to derive new drug
80 discovery insights via these synergies, particularly drug target hypotheses associated with
81 diseases, phenotypes and cell lines. LINCS provides the comprehensive systems/omics view,
82 while IDG focuses on relevant and robust evidence for drug target knowledge and validation.
83 Table 1 summarizes key concepts and the semantic linkage of LINCS and IDG. Previously,
84 LINCS studies have been designed to identify novel drug targets(6), and IDG has included
85 LINCS as a source for selected data(7). However, our approach provides a new KG method and
86 platform to integrate and add value from these sources for an urgent and unmet scientific use
87 case.

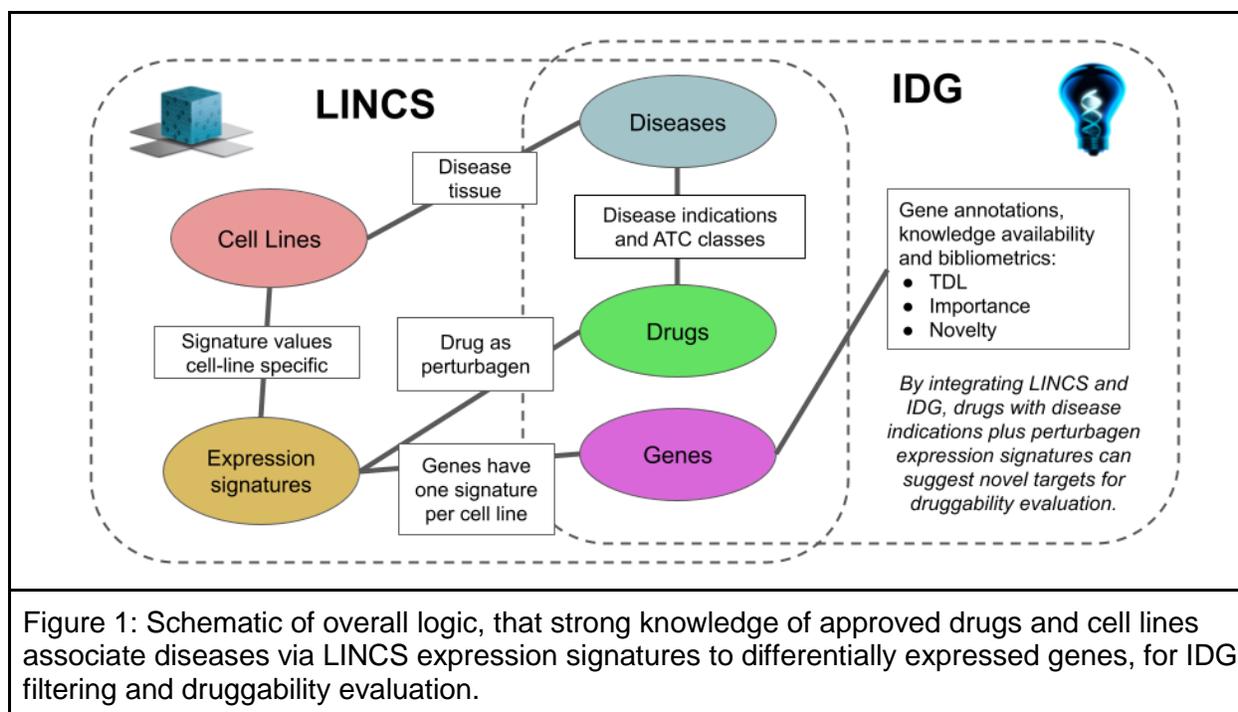
88 **Table 1:** Key concepts from LINCS and IDG

LINCS key concepts	IDG key concepts
--------------------	------------------

<p>L1000: Broad Institute microarray platform and LINCS project involving landmark genes.</p> <p>Landmark Gene: approximately 1000 genes representing the full human genome, from which a maximum set of gene responses can be inferred.</p> <p>Cell Lines: Primary cells and human patient induced pluripotent stem cells (iPSCs) representing diverse tissues and disease states (including 'healthy' controls).</p> <p>Perturbagen: A perturbagen may be chemical (drug-like), biologics, or genetic.</p> <p>Signatures: Various assays (transcriptomics, proteomics, or image analysis) are integrated together to identify patterns of common networks and cellular responses leading to predictive mechanistic interpretations.</p>	<p>Target: Drug targets normally defined as human proteins or equivalently protein coding genes.</p> <p>Mechanism of action: Known, published, biochemical mechanism of pharmacological effect.</p> <p>Disease: Working definition is anything which could be the clinical indication for a drug.</p> <p>Ligand: Small molecule with known or hypothesized relationship with one or more targets.</p> <p>Target Development Level (TDL): Simple measure of illumination: Tdark, Tbio, Tchem and Tclin.</p> <p>Target Importance and Novelty: Bibliometric measures from TIN-X for target prioritization.</p>
LINCS-IDG semantic linkage	
<p>LINCS centers utilize in-depth gene and protein expression assays to generate signatures directly mappable to IDG protein targets. Disease and phenotype ontology mapping is a community challenge with useful, working solutions such as OMOP and UMLS. LINCS perturbagens include rigorously defined chemical entities and small molecule drugs included in IDG resource DrugCentral. Therefore, LINCS's vast number of human cell lines and experimental chemical perturbation datasets, combined with IDG's protein targets (gene and protein IDs) and DrugCentral active pharmaceutical ingredients (drug compounds) databases, provide a tightly integrated combination resource for drug target discovery.</p>	

89

90



91

92 Results

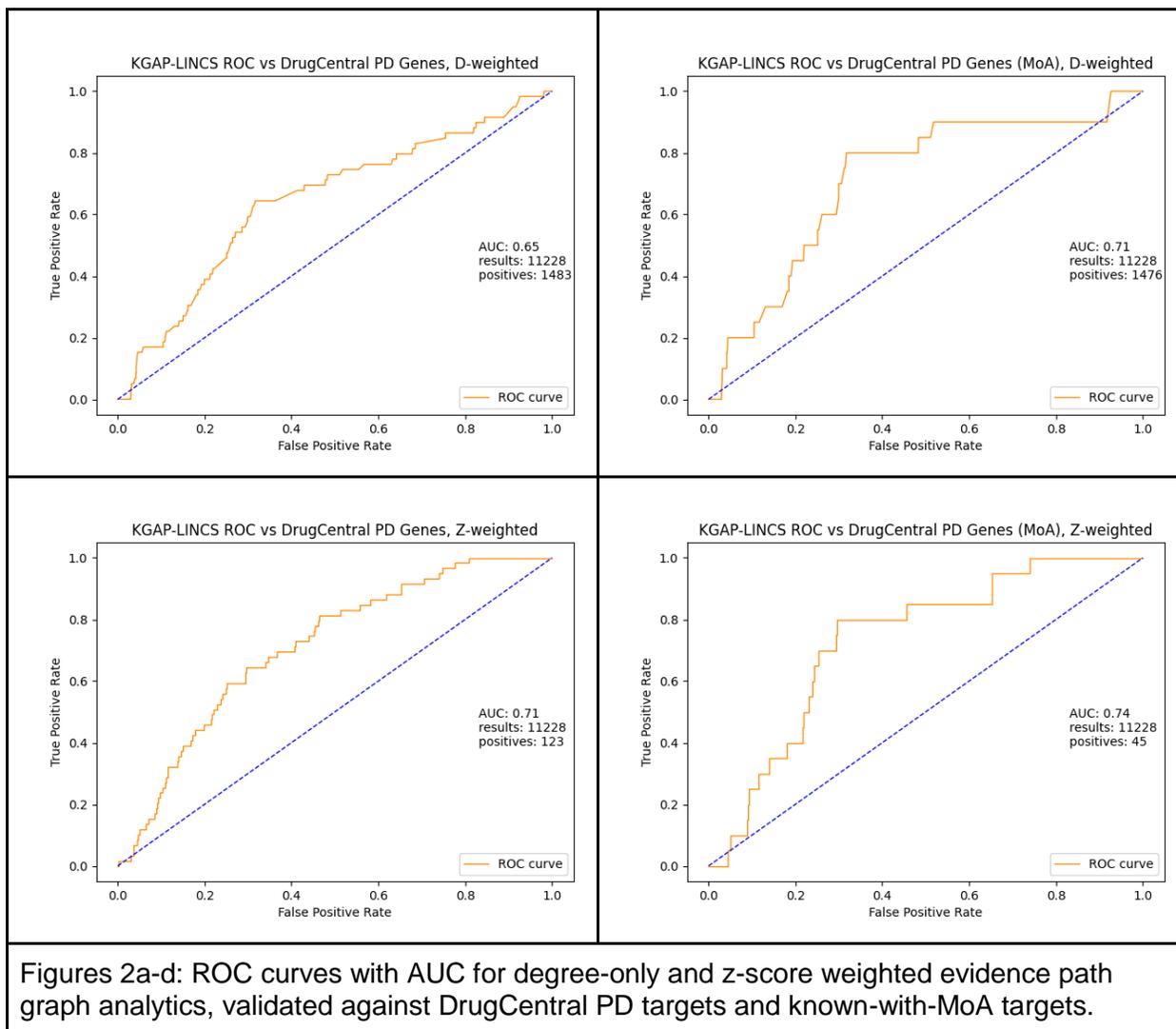
93 Graph analytics for evidence aggregation

94 We used the list of PD drugs as starting points for evidence paths in KGAP to identify likely and
95 novel PD targets. KGs in general and Neo4j in particular provide an intuitive and powerful way
96 to aggregate instances of evidence paths, yielding a score measuring the aggregated evidence.
97 As illustrated schematically in Figure 1, each evidence path connects a drug with expression
98 signature and gene. The magnitude of the expression perturbation is represented by a z-score
99 which can be used for thresholding and/or weighting. The search resulted in 641 genes, ranked
100 by KGAP score produced by the Cypher query Listing 1 in Methods.

101 Validation versus known mechanism-of-action targets

102 There is a necessity to solve the ongoing fundamental challenges of molecular biomedicine.
103 Specifically, there are no gold standard validation datasets of causal or druggable genes for
104 complex diseases, including cancers, neurodegenerative, metabolic, and cardiovascular
105 diseases. Therefore it is pertinent to elucidate and modulate the genomic basis of diseases, and
106 to a significant degree, fundamental uncertainties and difficulties concerning the definitions and
107 diagnostic criteria for complex diseases such as PD. Mindful of these difficulties, we validated
108 against a dataset of known drug targets, for the same PD drugs described above, all based on
109 peer-reviewed references, and with elucidated mechanism-of-action (MoA), as manually curated
110 in the DrugCentral database. Given the scientific and regulatory standards for efficacy met by
111 approved drugs, and the standards of evidence for peer-reviewed, manually curated MoA
112 targets, we consider this approach the most useful and informative dataset validation available.

113 LINCNS is experimental based, thus derived independently from the curation of MoA targets from
114 DrugCentral, or in machine learning terms, there is strict separation of training and test data.
115



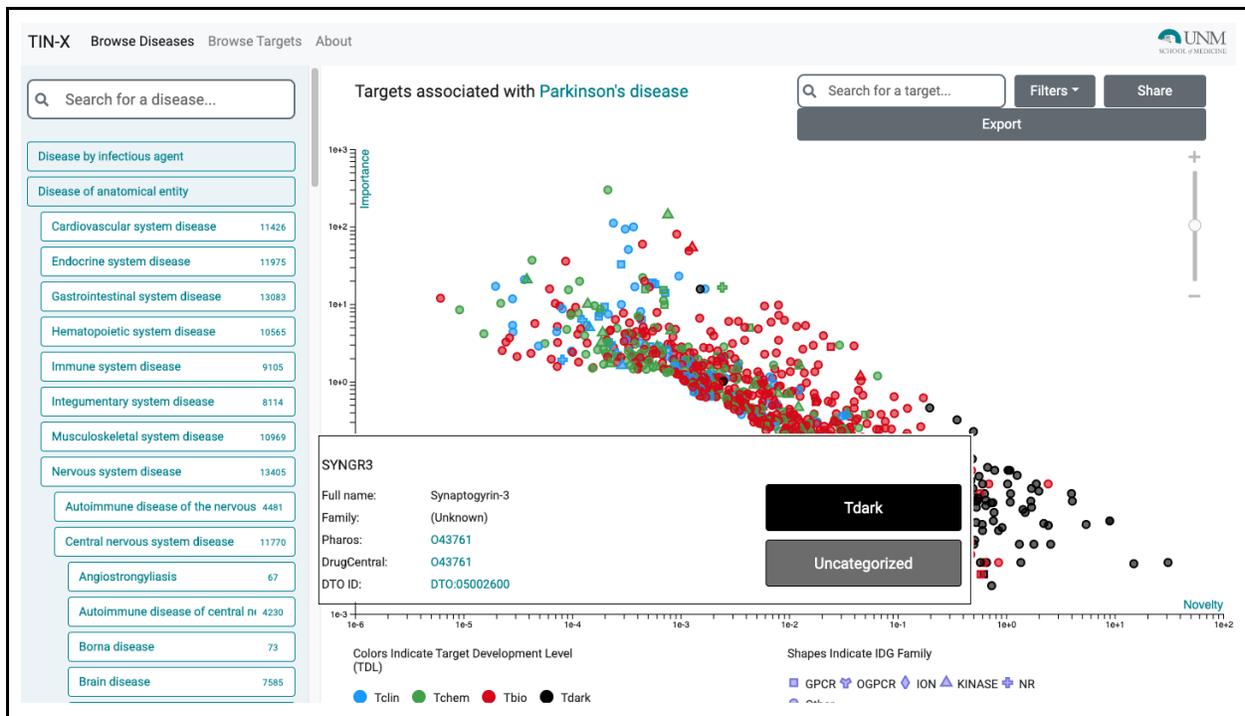
Figures 2a-d: ROC curves with AUC for degree-only and z-score weighted evidence path graph analytics, validated against DrugCentral PD targets and known-with-MoA targets.

116
117 Figures 2a-d show validation receiver operating characteristic curve (ROC) plots for two
118 variations of our method, and two validation sets. The "D-weighted" variation connotes
119 weighting of evidence paths by the sum of degree at the LINCNS expression signatures
120 associating drugs and genes. The "Z-weighted" variation combines degree with the z-score
121 expression level attribute associating signatures and genes. The two validation sets are (1)
122 DrugCentral PD targets, and (2) DrugCentral PD targets with known MoA, as described above.
123 The ROC AUC (area under the curve) values range between 0.64 and 0.74, providing
124 consistent, independent validation of the proposed method for disease to gene association
125 discovery. The validation method, for PD or other disease queries, is reproducible using the
126 source code repository referenced in the Availability of Data and Materials section. Note, the
127 validation presented does not validate KGAP as a classification method. Evaluated as such,
128 given the sparsity of known genes, the specificity is very poor. However, classification is not the

129 goal. Rather, our task is to aggregate and assess experimental evidence from LINCS which is
130 applicable to PD as proof of concept. The validation presented simply measures
131 overrepresentation, or enrichment, of known PD genes in relation to KGAP scores, indicating
132 agreement and independent confirmation between two very different sources of knowledge, one
133 purely experimental, the other expertly curated and peer-reviewed publication based.

134 Genes associated with PD, prioritized via IDG

135 Gene hitlists generated by KGAP were mapped to druggability data from IDG. Target
136 Development Level (TDL) provides a high-level classification into four groups — Tclin, Tchem,
137 Tbio and Tdark — with respect to the depth of investigation from a clinical, chemical and
138 biological standpoint. Further evaluation was explored via Target Importance and Novelty
139 Explorer (TIN-X)(8), an IDG project and bibliometric algorithm for evaluating disease-target
140 associations from scientific literature. Moreover, TIN-X defines **novelty** as a bibliometric
141 measure of occurrence rarity in the full PubMed corpus of titles and abstracts, and **importance**
142 as a bibliometric measure of co-occurrence associating a specific disease and gene. The
143 premise and motivation for IDG is that many understudied targets could offer new opportunities
144 for medicines of novel therapeutic benefit. Hence TIN-X ranks and presents targets based on
145 the principle of non-dominated solutions optimizing novelty and importance. $nds_rank=1$ is
146 assigned to all genes relative to which none are superior in both dimensions. $nds_rank=2$ is this
147 corresponding set with the first set removed, etc. In practice, typical users browse targets
148 beginning at the outer boundary, using TDL color code as an additional guide.
149



<p>Harnessing the trophic and modulatory potential of statins in a dopaminergic cell line.</p> <p>Schmitt, Mathieu M, Dehay, Benjamin B, Bezard, Erwan E and Garcia-Ladona, F Javier FJ. Synapse (New York, N.Y.)</p> <p>The identification of an effective disease-modifying treatment for the neurodegenerative progression in Parkinson's disease (PD) remains a major challenge. Epidemiological studies have reported that intake of statins, cholesterol lowering drugs, could be associated to a reduced risk of developing PD. In-vivo studies suggest that statins may reduce the severity of dopaminergic neurodegeneration. The trophic potential of statins and their impact on the expression of dopaminergic synaptic markers and dopamine (DA) transport function in SH-SY5Y cells has been investigated. The findings showed that statin treatment induces neurite outgrowth involving a specific effect on the complexity of the neurite branching pattern. Statins increased the levels of presynaptic dopaminergic biomarkers such as vesicular monoamine transporter 2 (VMAT2), synaptic vesicle glycoproteins 2A</p>	<p>U18666A, an activator of sterol regulatory element binding protein pathway, modulates presynaptic dopaminergic phenotype of SH-SY5Y neuroblastoma cells.</p> <p>Schmitt, Mathieu M, Dehay, Benjamin B, Bezard, Erwan E and Garcia-Ladona, F Javier FJ. Synapse (New York, N.Y.)</p> <p>The therapeutic use of statins has been associated to a reduced risk of Parkinson's disease (PD) and may hold neuroprotective potential by counteracting the degeneration of dopaminergic neurons. Transcriptional activation of the sterol regulatory element-binding protein (SREBP) is one of the major downstream signaling pathways triggered by the cholesterol-lowering effect of statins. In a previous study in neuroblastoma cells, we have shown that statins consistently induce the upregulation of presynaptic dopaminergic proteins and changes of their function and these effects were accompanied by downstream activation of SREBP. In this study, we aimed to determine the direct role of SREBP pathway in the modulation of dopaminergic phenotype. We demonstrate that</p>
<p>Figures 3a-c: TIN-X scatterplot of genes for Parkinson's disease, DOI:14330, showing pop-up details for SYNGR3, Synaptogyrin-3, and publication details view for PD associated gene <i>SYNGR3</i>.</p>	

150

151 Case study: SYNGR3, Synaptogyrin-3

152 To illustrate a typical use case, the KGAP hitlist was browsed for drug target illumination
153 potential based on annotations from IDG. The highest ranked gene classified as Tdark is
154 *SYNGR3*, Synaptogyrin-3. The exact function of *SYNGR3* is unclear. However, recently a group
155 provides evidence in the murine brain that *SYNGR3* encodes for a synaptic vesicle protein that
156 interacts with a dopamine transporter(9). One of PD's hallmark characterizations is the loss of
157 nigrostriatal dopaminergic innervation(10). The high TIN-X rank (nds_rank=6 out of 103)
158 indicates both novelty and importance (PD-relevance), as shown in figure 3a. Additionally, in
159 figures 3b and 3c, two reference publications from the same authors, present experimental and
160 theoretical evidence for connection between statin drugs, therapeutic effectiveness for PD, and
161 the gene *SYNGR3*(11,12). Figure 4 displays evidence paths connecting *SYNGR3* with
162 associated expression signatures and drugs, matched by our method.
163

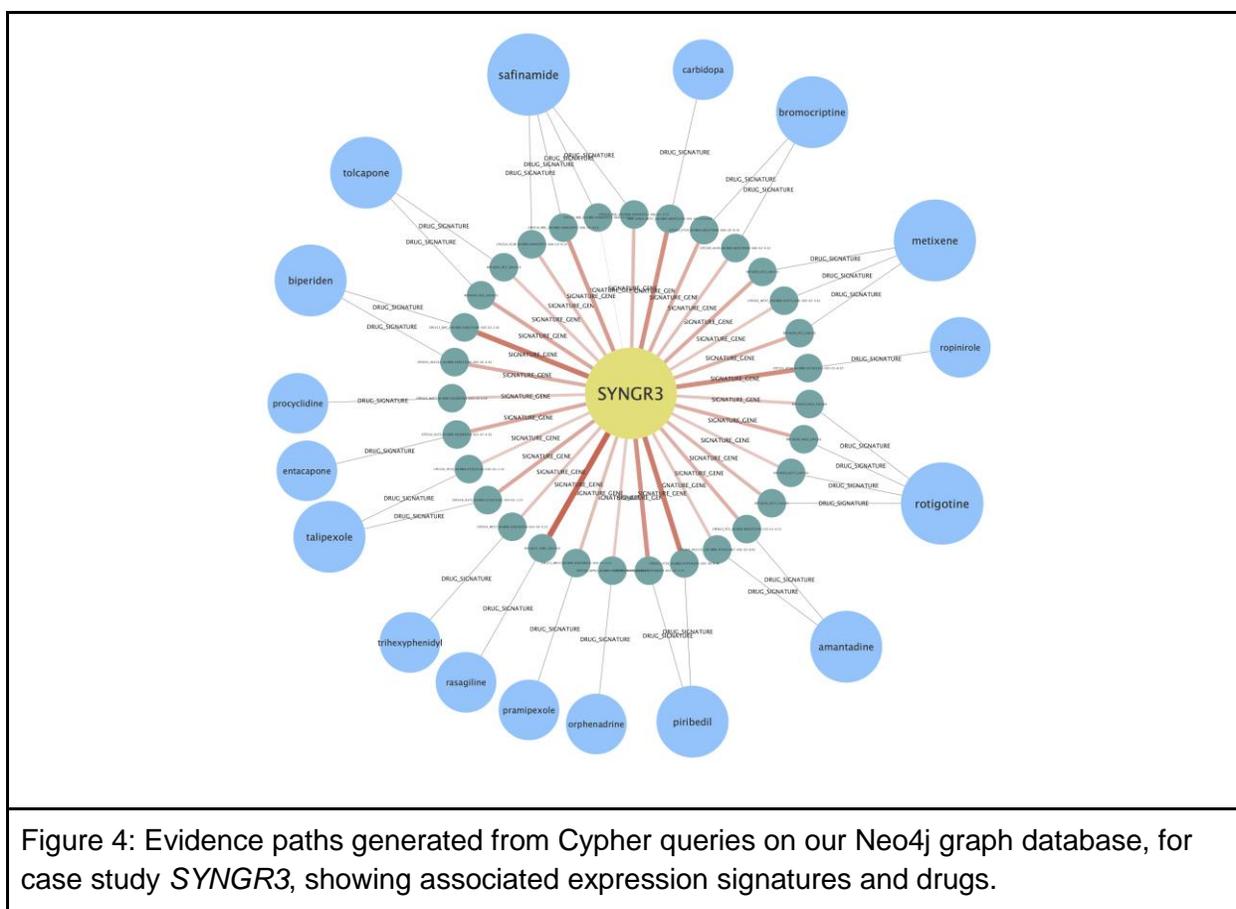


Figure 4: Evidence paths generated from Cypher queries on our Neo4j graph database, for case study *SYNGR3*, showing associated expression signatures and drugs.

164 Discussion

165 This project builds upon numerous prior efforts and resources, conceptual and technological,
166 from domain specific data semantics, to graph database advances, conferring interoperability of
167 depth and breadth(13–19). That is, beyond interoperable digital formats, scientists interacting
168 effectively and rigorously via computational tools and interfaces. Regarding the now-common
169 term *knowledge graph (KG)*, there are variations of meaning, but we emphasize: (1) Knowledge
170 with strong semantics via rigorous data modeling; (2) Graph databases are powerful and
171 intuitive for many biomedical applications and users(20); (3) Neo4j concretely implements a rich
172 suite of tools (DBMS, query language, APIs, GUI) as an exemplar among graph database
173 systems. Hence, we use the term KG unambiguously and with specific commitments.
174 Advances in KG technologies have empowered progress far beyond previous efforts. A case in
175 point is Chem2Bio2RDF(14,21), developed in the Wild Lab, enabled by RDF and Sparql but
176 limited for analytic applications by availability of high-performance triple-store technologies.
177 Another project which one of us (JY) has co-developed is CARLSBAD(22), a bioactivity
178 knowledge graph limited in analytics performance and versatility by its implementation as a
179 relational database. Additional lessons learned relate to (1) Generality versus purposefulness,
180 and (2) Volume versus veracity. Accordingly, this project has prioritized quality and simplicity to
181 mitigate the common problems of (1) unfocused, incoherent interfaces, and (2) unreliable or

182 noisy data. Interpretability also depends on understanding of the tool's purpose, logic, results,
183 and their provenance and confidence. In many cases, simple and clearly focused tools enhance
184 comprehensibility and usability. In contrast, a toolkit is multi-purpose. For this research the key
185 aim is to find gene associations based on a disease or drug query.

186 Though the Neo4j product is featured prominently in this paper, it is important to note that
187 alternative capable graph databases exist, such as DGraph(23) and Amazon Neptune(24).
188 Moreover, Neo4j and other graph databases are built on the shoulders of many diverse
189 contributors and strands of computer science, thus, Neo4j is featured as an exemplar for these
190 advances, and a powerful technology for our use case.

191 Development of this knowledge graph analytics platform (KGAP) with LINCS and IDG was
192 motivated by use cases from IDG and its core purpose to illuminate understudied genes and
193 new drug targets. Generally stated: For a given disease, what new drug targets are suggested
194 by the evidence? LINCS data provides aggregated experimental evidence to establish a cell
195 and expression signature based systems biology database, supporting rational aggregation of
196 disease-gene associations, powered by KG analytics well suited both to path-based analytics
197 and interpretation. Representing PD informatically is another key ingredient, addressed by
198 clinical indications for approved drugs, knowledge meeting high regulatory standards. The
199 definitions and nosology of PD and other complex diseases are critical issues and limitations
200 relevant to this study and many others. Any hypotheses regarding PD presumes a useful
201 definition, however biologically PD and its complexity are better described by subtypes and
202 clinical phenotypes. Yet, extant data is generally organized by diseases and clinical diagnoses
203 as historically developed.

204 Conclusions

205 A knowledge graph analytic platform (KGAP) was developed, integrating datasets from LINCS
206 and IDG, for efficient search and aggregation of evidence paths based on a disease query, to
207 identify, score and rank associated genes as drug target hypotheses. Approved indications for
208 prescription drugs were used as high confidence entry points into the KG, and published MoA
209 targets were used as a high confidence validation set. Modern graph database systems such as
210 Neo4j provide a powerful suite of tools for high performance analytics, rigorously and
211 reproducibly encoded, on semantically rich KGs, with interactivity and visualization enhancing
212 human understanding. This study has demonstrated how KGAP can generate novel and
213 plausible drug targets for Parkinson's disease. It includes a case study of the understudied
214 "Tdark" *SYNGR3* gene, scored highly by KGAP and supported independently by publications
215 identified via the IDG application TIN-X. The KGAP approach and method as implemented is
216 generalizable and applicable to drug target illumination in many other disease areas.
217 Accordingly, in future work we intend to apply KGAP more broadly.
218

219 Methods

220 Graph database

221 Graph databases are one category of NoSQL (non-SQL or non-relational) databases which are
222 designed to outperform relational databases for some specialized datasets and analyses, and
223 are well suited for representing and exploring complex biological systems(25–27). Just as graph
224 diagrams provide an intuitive way to represent and explore relationships between entities, graph
225 databases leverage this intuition to provide both a physical storage method, and approach for
226 combining diverse datasets for exploration and analysis. Neo4j was chosen as an exemplary
227 leader among graph databases with a rich set of graph based capabilities (e.g. Graph Data
228 Science library, and libraries for accessing, flat files, databases, existing open standards such
229 as SPARQL, and open standard in development GQL(28)). In addition, Neo4j offers a free,
230 open source community edition(29), powerful and user-friendly desktop client, APIs, ample
231 documentation, and a vibrant and supportive user community.

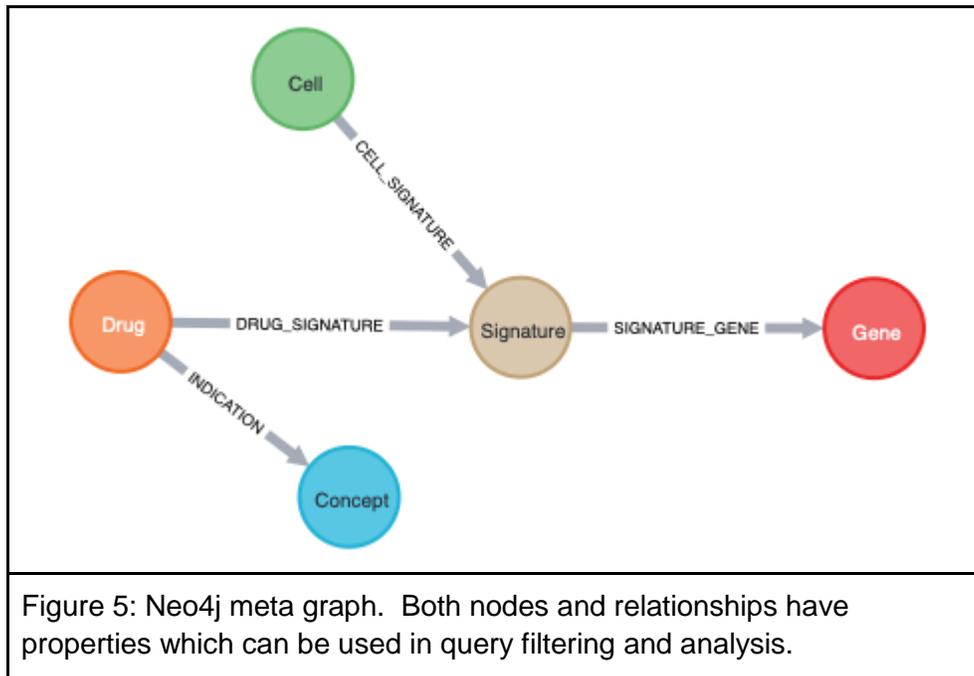
232 Data Sources

233 As shown in Figure 1 an understanding of entities and the relationships between them can be
234 quickly represented in a representative graph diagram, this establishes high level concepts and
235 methods for discovering new insights, but also provides a clear design for the required extract
236 transform and load (ETL) steps for building the capability. The datasets listed here were used
237 to construct a graph database instantiation of this design. From LINCS we used the L1000
238 Phase I and II Level 5 differential gene expression data(30). Level 5 connotes the highest level
239 of processing, normalization, evidence aggregation, and therefore confidence, interpretability,
240 and interoperability. LINCS was processed using the same pipeline used by colleagues for the
241 DrugCentral LINCS drug profiling tool, whereby differential gene expression data across 81 cell-
242 lines were mapped to 1613 unique drug active pharmaceutical ingredients(5). Data and
243 metadata files available from the LINCS Data Portal and NCBI GEO were ETL-ed via 1.5TB
244 PostgreSQL db. DrugCentral(5) (June 2020 version) was used for indications, ATC codes,
245 chemical structures and other properties. IDG's TCRD, Target Central Resource Database(31),
246 (v6.7.0) provides gene properties including target development level (TDL: Tclin, Tchem, Tbio,
247 Tdark) and gene family information.

248 Data modeling, representing the knowledge

249 Graphs in Neo4j are composed of nodes and relationships, which correspond with vertices and
250 edges in graph theory terms, respectively. Graph queries can be thought of as patterns for
251 matching paths through the graph, which consist of specified relationship types, each between
252 two nodes, a semantic triple, subject-predicate-object. Neo4j is considered a schema-less
253 database, meaning that data is loaded without constraints. But a strict schema with strong
254 semantics informed by domain knowledge can be implemented, and is essential for accurate
255 knowledge representation and interpretability. The meta-model for a graph can be reported by

256 introspective Cypher query (`CALL apoc.meta.graph`). The meta-model for this graph is
 257 shown in Figure 5.
 258



259
 260 Our graph is constructed of the following.

261
 262 **Relationships.** INDICATION is constructed from DrugCentral indications involving a Drug
 263 which is also a LINCS perturbagen. CELL_SIGNATURE indicates the Cell type for the LINCS
 264 Signature. SIGNATURE_GENE provides the differential expression z-score for the Gene.
 265 DRUG_SIGNATURE associates the Drug involved for a given Signature.

266
 267 **Nodes.** Concept consists of disease terms from DrugCentral with unique OMOP identifiers, and
 268 retaining other disease identifiers as properties for interoperability and extensibility. Drug nodes
 269 are derived by filtering for LINCS perturbagens which are in DrugCentral and have one or more
 270 indications. Genes are from LINCS, with additional properties added from TCRD (e.g.,
 271 development/druggability level). Cell and Signature nodes are from LINCS without filtering.

272
 273 Table 2 shows the complete list of relationships present in the graph and count, and Table 3 the
 274 list of nodes with counts.

275
 276
 277 **Table 2:** Type of relationship, source, and total counts

Vertex type	Edge type	Vertex type	Source	Count
Signature	SIGNATURE_GENE	Gene	L1000	45623637
Cell	CELL_SIGNATURE	Signature	L1000	591676

Drug	DRUG_SIGNATURE	Signature	L1000	85814
Drug	INDICATION	Concept	DrugCentral, L1000	7164

278

279

280 **Table 3:** Type of node, primary source/ source of additional properties, and total counts

Vertex type	Source	Count
Signature	L1000	591676
Gene	L1000, TCRD	12329
Concept	DrugCentral 2021	2241
Drug	L1000, DrugCentral 2021	1613
Cell	L1000	99

281

282 We filter the 45.6M signature associations at z-score threshold $|z|>3$, to identify strong evidence
 283 and patterns, which results in 10,663,228 (d:Drug)--(s:Signature)--(g:Gene) paths.

284 Parkinson's disease drug-set

285 To apply KGAP to PD knowledge discovery, the set of drugs indicated by PD and related
 286 conditions were identified using DrugCentral, which derives drug indication from FDA DailyMed
 287 (provides trustworthy information about marketed drugs in the United States), and maps
 288 SNOMED to OMOP terms(32) for interoperability. A simple substring search for "Parkinson"
 289 matches five terms: "Parkinsonism", "Parkinson's disease", "Arteriosclerotic Parkinsonism",
 290 "Dementia associated with Parkinson's Disease", and "Neuroleptic-induced Parkinsonism."
 291 Requiring Anatomical and Therapeutic Classification (ATC) "nervous system" is intended to
 292 focus on the neurological etiology of PD rather than symptoms. This query returns 25 drugs,
 293 shown in Table 4 with their PubChem CID, of which 22 / 25 are present in LINCS. This drug-set
 294 represents PD knowledge as in drug-set enrichment analysis (DSEA)(33), analogous to gene-
 295 set enrichment analysis (GSEA).

296 **Table 4:** Drug-set representing Parkinson's disease

Name	Struct_ID	PubChem_CID	In_LINCS
amantadine	144	2130	TRUE
apomorphine	228	6005	TRUE
benzatropine	333	1201549	TRUE
biperiden	374	2381	TRUE
bromocriptine	403	31101	TRUE

dexetimide	831	30843	<i>FALSE</i>
entacapone	1018	5281081	<i>TRUE</i>
levodopa	1567	6047	<i>TRUE</i>
melevodopa	1673	23497	<i>FALSE</i>
metixene	1780	4167	<i>TRUE</i>
opicapone	5143	135565903	<i>FALSE</i>
orphenadrine	1999	4601	<i>TRUE</i>
pergolide	2105	47811	<i>TRUE</i>
pimavanserin	5142	10071196	<i>TRUE</i>
piribedil	2202	4850	<i>TRUE</i>
pramipexole	2233	119570	<i>TRUE</i>
procyclidine	2276	4919	<i>TRUE</i>
rasagiline	3521	3052776	<i>TRUE</i>
rivastigmine	2392	77991	<i>TRUE</i>
ropinirole	2402	5095	<i>TRUE</i>
rotigotine	2407	59227	<i>TRUE</i>
safinamide	4921	131682	<i>TRUE</i>
selegiline	2429	26757	<i>TRUE</i>
tolcapone	2697	4659569	<i>TRUE</i>
trihexyphenidyl	2745	5572	<i>TRUE</i>

297

298 Graph analytics for evidence aggregation

299 KGAP consists of the graph database described plus code and tools used to query the
300 database for target associations from the input drug-set representing the disease of interest. In
301 the current version of the platform, the queries were implemented via the Neo4j Python
302 Driver(34) and our prototype Python3 command-line applications and notebooks. As in the data
303 modeling in the construction of the database, the Cypher graph queries must be crafted to
304 accurately express the biomedical question, and justify useful interpretation of results. However,
305 since the database was designed precisely for target association and ranking, the Cypher
306 queries involved can be relatively terse. The query used to generate PD target associations is
307 shown in Listing 1. The scoring function is responsible for aggregating evidence from signature
308 counts and z-scores, to reflect that confidence increases with additional confirmatory data, and
309 is an implementation of Stouffer's function for z-score meta-analysis, based on Fisher's function

310 using p-values(35). Though not used in KGAP for hypothesis testing, the reference indicates the
311 function is well behaved and combines z-scores as intended.

312 Listing 1: Cypher query used for PD target association.

```
MATCH p=(d:Drug)-[]-(s:Signature)-[r]-(g:Gene), p1=(s)-[]-(c:Cell)
WHERE (d.pubchem_cid in [2130, 2130, 2130, 2130, 2381, 2381, 2381, 2381, 4167, 4601, 4850, 4919, 4919, 4919, 4919, 5095, 5572, 5572, 5572, 5572, 6005, 6047, 6047, 6047, 23497, 26757, 30843, 31101, 31101, 47811, 59227, 77991, 119570, 131682, 1201549, 1201549, 1201549, 1201549, 3052776, 4659569, 5281081, 10071196, 135565903])
WITH g, sum(r.zscore)/sqrt(count(r)) AS score
RETURN g.id, g.name, score
ORDER BY score DESC
```

313

314 Exploring and visualizing the KG

315 Human interaction with the KG through an effective GUI can greatly facilitate understanding and
316 insights. In this project we used (1) Neo4j Desktop, (2) Neo4j Browser web client, used to
317 produce Figure 5, and (3) Cytoscape (v3.8.1) with the Neo4j plugin, which was used to produce
318 Figure 4.

319 List of abbreviations

320

IDG	Illuminating the druggable genome (NIH Common Fund)
KG	Knowledge graph
KGAP	Knowledge graph analytics platform
LINCS	Library of integrated network-based cellular signatures (NIH Common Fund)
MoA	Mechanism of action, describing biomolecular details of drug effect
PD	Parkinson's disease
ROC AUC	Receiver operator characteristic area under curve
TCRD	Target central resource database (IDG)
TDL	Target development level (IDG)

321

322 **Declarations**

323 **Ethics approval and consent to participate**

324 Not applicable.

325 **Consent for publication**

326 Not applicable.

327 **Availability of data and materials**

328 Source code for processing and analysis is available publicly under Creative Commons Zero
329 v1.0 Universal license at https://github.com/IUIDSL/kgap_lincs-idg. A Knime workflow is
330 employed to extract, transform and load the graph database from sources, and generate TSV
331 files of nodes and relationships. The database is queried, and results reported and visualized
332 using Python and Cypher. A full dump of the KGAP_LINCS-IDG dataset is available at
333 <http://cheminfov.informatics.indiana.edu/projects/kgap/data/dclneodb.dump>. Source datasets
334 used to build Neo4j database are: (1) DrugCentral-LINCS database dump:
335 http://cheminfov.informatics.indiana.edu/projects/kgap/data/drugcentral_lincs.pgdump; (2)
336 TCRD targets: http://cheminfov.informatics.indiana.edu/projects/kgap/data/tcrd_targets.tsv.gz.
337 Neo4j Community Edition, Python, Scikit-learn, Cytoscape and other packages used in this work
338 are freely and easily findable and accessible.

339 **Competing interests**

340 JY, JD, BF, DB, KS, YD, and DW are founders, employees or contractors of Data2Discovery, a
341 private company spun off from Indiana University to develop and commercialize knowledge
342 graph technologies.

343 **Funding**

344 None.

345 **Authors' contributions**

346 JY and DW conceived the project. JD designed and built the graph database and ETL
347 workflows. CG and JD developed graph analytic algorithms. JB led the SYNGR3 case study.
348 JY, JD, DB, CG and MO developed client code. JY, JD, CG, DB, and JB co-authored the
349 manuscript. All reviewed and approved the manuscript.

350 Acknowledgements

351 We are grateful to the LINCS and IDG projects and investigators responsible for the high value
352 datasets shared with the community and used in this research.

353 References

- 354 1. National Institutes of Health, U.S. Department of Health and Human Services. National
355 Institutes of Health Common Fund [Internet]. NIH Office of Strategic Coordination - The
356 Common Fund. [cited 2020 Dec 21]. Available from: <https://commonfund.nih.gov/>
- 357 2. Keenan AB, Jenkins SL, Jagodnik KM, Koplev S, He E, Torre D, et al. The Library of
358 Integrated Network-Based Cellular Signatures NIH Program: System-Level Cataloging of
359 Human Cells Response to Perturbations. *Cell Syst*. 2018 Jan 24;6(1):13–24.
- 360 3. Oprea TI, Bologa CG, Brunak S, Campbell A, Gan GN, Gaulton A, et al. Unexplored
361 therapeutic opportunities in the human genome. *Nat Rev Drug Discov*. 2018
362 May;17(5):377.
- 363 4. Sheils T, Mathias SL, Siramshetty VB, Bocci G, Bologa CG, Yang JJ, et al. How to
364 illuminate the Druggable Genome Using Pharos. *Curr Protoc Bioinformatics*. 2020
365 Mar;69(1):e92.
- 366 5. Avram S, Bologa CG, Holmes J, Bocci G, Wilson TB, Nguyen D-T, et al. DrugCentral 2021
367 supports drug discovery and repositioning. *Nucleic Acids Res* [Internet]. 2020 Nov 5;
368 Available from: <http://dx.doi.org/10.1093/nar/gkaa997>
- 369 6. Chen B, Ma L, Paik H, Sirota M, Wei W, Chua M-S, et al. Reversal of cancer gene
370 expression correlates with drug efficacy and reveals therapeutic targets. *Nat Commun*.
371 2017 Jul 12;8:16022.
- 372 7. Ursu O, Holmes J, Bologa CG, Yang JJ, Mathias SL, Stathias V, et al. DrugCentral 2018:
373 an update. *Nucleic Acids Res* [Internet]. 2018 Oct 29; Available from:
374 <http://dx.doi.org/10.1093/nar/gky963>
- 375 8. Cannon DC, Yang JJ, Mathias SL, Ursu O, Mani S, Waller A, et al. TIN-X: target
376 importance and novelty explorer. *Bioinformatics*. 2017 Aug 15;33(16):2601–3.
- 377 9. Egaña LA, Cuevas RA, Baust TB, Parra LA, Leak RK, Hochendoner S, et al. Physical and
378 functional interaction between the dopamine transporter and the synaptic vesicle protein
379 synaptogyrin-3. *J Neurosci*. 2009 Apr 8;29(14):4592–604.
- 380 10. Stoker TB, Greenland JC, editors. *Parkinson's Disease: Pathogenesis and Clinical Aspects*.
381 Brisbane (AU): Codon Publications; 2019.
- 382 11. Schmitt M, Dehay B, Bezard E, Javier Garcia-Ladona F. Harnessing the trophic and
383 modulatory potential of statins in a dopaminergic cell line [Internet]. Vol. 70, *Synapse*. 2016.
384 p. 71–86. Available from: <http://dx.doi.org/10.1002/syn.21881>
- 385 12. Schmitt M, Dehay B, Bezard E, Garcia-Ladona FJ. U18666A, an activator of sterol

- 386 regulatory element binding protein pathway, modulates presynaptic dopaminergic
387 phenotype of SH-SY5Y neuroblastoma cells. *Synapse* [Internet]. 2017 Sep;71(9). Available
388 from: <http://dx.doi.org/10.1002/syn.21980>
- 389 13. Callahan A, Cruz-Toledo J, Ansell P, Dumontier M. Bio2RDF Release 2: Improved
390 Coverage, Interoperability and Provenance of Life Science Linked Data [Internet]. *The*
391 *Semantic Web: Semantics and Big Data*. 2013. p. 200–12. Available from:
392 http://dx.doi.org/10.1007/978-3-642-38288-8_14
- 393 14. Chen B, Dong X, Jiao D, Wang H, Zhu Q, Ding Y, et al. Chem2Bio2RDF: a semantic
394 framework for linking and data mining chemogenomic and systems chemical biology data.
395 *BMC Bioinformatics*. 2010 May 17;11:255.
- 396 15. Digles D, Caracoti A, Jacoby E. Accessing the Open PHACTS Discovery Platform with
397 Workflow Tools. *Methods Mol Biol*. 2018;1787:183–93.
- 398 16. Szklarczyk D, Gable AL, Lyon D, Junge A, Wyder S, Huerta-Cepas J, et al. STRING v11:
399 protein-protein association networks with increased coverage, supporting functional
400 discovery in genome-wide experimental datasets. *Nucleic Acids Res*. 2019 Jan
401 8;47(D1):D607–13.
- 402 17. Greene CS, Krishnan A, Wong AK, Ricciotti E, Zelaya RA, Himmelstein DS, et al.
403 Understanding multicellular function and disease with human tissue-specific networks. *Nat*
404 *Genet*. 2015 Jun;47(6):569–76.
- 405 18. Koscielny G, An P, Carvalho-Silva D, Cham JA, Fumis L, Gasparyan R, et al. Open
406 Targets: a platform for therapeutic target identification and validation. *Nucleic Acids Res*.
407 2017 Jan 4;45(D1):D985–94.
- 408 19. Morton K, Wang P, Bizon C, Cox S, Balhoff J, Kebede Y, et al. ROBOKOP: an abstraction
409 layer and user interface for knowledge graphs to support question answering.
410 *Bioinformatics*. 2019 Dec 15;35(24):5382–4.
- 411 20. Lysenko A, Roznovat IA, Saqi M, Mazein A, Rawlings CJ, Auffray C. Representing and
412 querying disease networks using graph databases. *BioData Min*. 2016 Jul 25;9:23.
- 413 21. Chen B, Ding Y, Wild DJ. Assessing Drug Target Association Using Semantic Linked Data
414 [Internet]. Vol. 8, *PLoS Computational Biology*. 2012. p. e1002574. Available from:
415 <http://dx.doi.org/10.1371/journal.pcbi.1002574>
- 416 22. Mathias SL, Hines-Kay J, Yang JJ, Zahoransky-Kohalmi G, Bologna CG, Ursu O, et al. The
417 CARLSBAD database: a confederated database of chemical bioactivities. *Database* . 2013
418 Jun 21;2013:bat044.
- 419 23. Dgraph Labs, Inc. Dgraph [Internet]. Dgraph. [cited 2020 Dec 21]. Available from:
420 <https://dgraph.io/>
- 421 24. Amazon Web Services, Inc. Amazon Neptune [Internet]. Amazon Neptune. [cited 2020 Dec
422 21]. Available from: <https://aws.amazon.com/neptune/>
- 423 25. Have CT, Jensen LJ. Are graph databases ready for bioinformatics? *Bioinformatics*. 2013
424 Dec 15;29(24):3107–8.

- 425 26. Himmelstein DS, Baranzini SE. Heterogeneous Network Edge Prediction: A Data
426 Integration Approach to Prioritize Disease-Associated Genes. *PLoS Comput Biol*. 2015
427 Jul;11(7):e1004259.
- 428 27. Yoon B-H, Kim S-K, Kim S-Y. Use of Graph Database for the Integration of Heterogeneous
429 Biological Data. *Genomics Inform*. 2017 Mar;15(1):19–27.
- 430 28. JCC Consulting, Inc. , acting on behalf of an unincorporated association of ISO Graph
431 Query Language Proponents, and licensed under the Apache License, Version 2. GQL
432 Standard [Internet]. Graph Query Language GQL. [cited 2020 Dec 21]. Available from:
433 <https://www.gqlstandards.org/>
- 434 29. Neo4j, Inc. Neo4j Licensing Overview [Internet]. Neo4j. [cited 2020 Dec 21]. Available from:
435 <https://neo4j.com/licensing/>
- 436 30. Subramanian A, Narayan R, Corsello SM, Peck DD, Natoli TE, Lu X, et al. A Next
437 Generation Connectivity Map: L1000 Platform and the First 1,000,000 Profiles. *Cell*. 2017
438 Nov 30;171(6):1437–52.e17.
- 439 31. IDG-KMC (Illuminating the Druggable Genome Knowledge Management Center). IDG-KMC
440 Target Central Resource Database [Internet]. Target Central Resource Database. [cited
441 2020 Nov 30]. Available from: <http://juniper.health.unm.edu/tcrd/>
- 442 32. Hripcsak G, Duke JD, Shah NH, Reich CG, Huser V, Schuemie MJ, et al. Observational
443 Health Data Sciences and Informatics (OHDSI): Opportunities for Observational
444 Researchers. *Stud Health Technol Inform*. 2015;216:574–8.
- 445 33. Napolitano F, Sirici F, Carrella D, di Bernardo D. Drug-set enrichment analysis: a novel tool
446 to investigate drug mode of action. *Bioinformatics*. 2016 Jan 15;32(2):235–41.
- 447 34. Neo4j, Inc. Neo4j Python Driver [Internet]. Neo4j. [cited 2020 Dec 21]. Available from:
448 <https://neo4j.com/docs/api/python-driver/current/>
- 449 35. Rosenthal R. Combining results of independent studies [Internet]. Vol. 85, *Psychological*
450 *Bulletin*. 1978. p. 185–93. Available from: <http://dx.doi.org/10.1037/0033-2909.85.1.185>

451
452
453
454
455