



Introducing federated queries with Wikidata

Andra Waagmeester¹, Micelio, Antwerp, Belgium | Email: andra@micel.io, Twitter: @andrawaag
Micelio, Ekeren, Antwerp, Belgium

The Gene Wiki project, circa 2008

Summarized
knowledge via
crowdsourcing

IL2-inducible T-cell kinase

Function

This gene encodes an intracellular tyrosine kinase expressed in T-cells. The protein is thought to play a role in T-cell proliferation and differentiation.^{[2][3]}

Structure

This protein contains the following domains, which are often found in intracellular kinases:^[4]

- N-terminus – PH (pleckstrin homology domain)
- BTK – Bruton's tyrosine kinase Cys-rich motif
- SH3 – (Src homology 3)
- SH2 – (Src homology 2)
- C-terminus – tyrosine kinase, catalytic domain

Interactions

ITK (gene) has been shown to interact with FYN,^{[5][6]} Wiskott-Aldrich syndrome protein,^{[7][8]} KHDRBS1,^{[8][9][10]} PLCG1,^{[10][11]} Lymphocyte cytosolic protein 2,^{[11][12]} Linker of activated T cells,^{[12][13]} Karyophormin alpha 2,^[14] Grb2^{[5][9]} and Peptidyl/prolyl isomerase A.^[15]

References

- Gibson S, Leung B, Squire JA, Hill M, Arima N, Goss P, Hogg D, Mills GB (September 1993). "Identification, cloning, and characterization of a novel human T-cell-specific tyrosine kinase located at the hematopoietin complex on chromosome 5q". *Blood* **82** (5): 1561–72. PMID 8394205.
- Kosaka Y, Felices M, Berg LJ (October 2006). "Itk and Th2 responses: action but no reaction". *Trends Immunol.* **27** (10): 453–60. doi:10.1016/j.it.2006.08.008. PMID 16931156.
- "Entrez Gene: ITK (IL2-inducible T-cell kinase)".
- Hawkins J, Marcy A (July 2001). "Characterization of Itk tyrosine kinase: contribution of noncatalytic domains to enzymatic activity". *Protein Expr. Purif.* **22** (2): 211–9. doi:10.1006/prep.2001.1447. PMID 11437596.
- Bunnell, S.C.; Diehn M; Yaffe M B; Findell P R; Cantley L C; Berg L J (Jan. 2000). "Biochemical interactions integrating Itk with the T cell receptor-initiated signaling cascade". *J. Biol. Chem. (UNITED STATES)* **275** (3): 2219–30. ISSN 0021-9258. PMID 10639929.
- Intramolecular association in a tyrosine kinase of the Tec family. *Nature (ENGLAND)* **385** (6611): 93–7. doi:10.1038/385093a0. ISSN 0028-0836. PMID 8985255.
- Perez-Villar, J. J; Kanner S B (Dec. 1999). "Regulated association between the tyrosine kinase Emh1/Tsk and phospholipase-C gamma 1 in human T lymphocytes". *J. Immunol. (UNITED STATES)* **163** (12): 6435–41. ISSN 0022-1767. PMID 10580033.
- Shim, Eun Kyung; Moon Chang Suk; Lee Gi Yeon; Ha Yun Jung; Chae Suhm-Ke; Lee Jong Ran (Sep. 2004). "Association of the Src homology 2 domain-containing leukocyte phosphoprotein of 76 kD (SLP-76) with the p85 subunit of phosphoinositide 3-kinase". *FEBS Lett. (Netherlands)* **575** (1–3): 35–40. doi:10.1016/j.febslet.2004.07.090. ISSN 0161-6759. PMID 15388330.
- Shan, X; Wang R L (Oct. 1999). "Itk/EmtSk activation in response to CD3 cross-linking in Jurkat T cells requires ZAP-70 and Lat and is independent of membrane recruitment". *J. Biol. Chem. (UNITED STATES)* **274** (41): 29323–30. ISSN 0021-9258. PMID 10501192.
- Perez-Villar, Juan J; Whitehead-Phillips, Simon

IL2-inducible T-cell kinase

Available structures

1ku, 1lu, 1um, 1un, 1sn2, 1snu, 1snx, 2eltz, 2evd

Identifiers

symbols ITK; PSCTK2; EMT; LYK; MGC128257; MGC128258

external IDs OMIM: 169873; MGI: 96621; HomoloGene: 4051

GeneCards: ITK Gene

number 2.7.10.2

Gene ontology [show]

RNA expression pattern

21109_z_at

Protein domains

Orthologs

species	Human	Mouse
entrez	3702	16428
ensembl	ENSG000000113263	ENSMUSG000000020395
UniProt	Q08881	A1A560
RefSeq	NM_005546	NM_010583
RefSeq	NP_065537	NP_034713

Data imported
from structured
databases

Reelin

From Wikipedia, the free encyclopedia

Reelin is a large secreted [extracellular matrix glycoprotein](#) that helps regulate processes of [neuronal migration](#) and positioning in the developing brain by controlling [cell–cell interactions](#).

Besides this important role in early [development](#), reelin continues to work in the adult brain. It modulates [synaptic plasticity](#) by ^{[2][3]} It also stimulates dendrite^[4] migration of [neuroblasts](#) generating [zones](#). It is found not only in the [tissues](#).

Reelin has been suggested to be implicated in pathogenesis of several brain diseases. The expression of the protein has been found to be significantly lower in schizophrenia and psychotic bipolar disorder, but the cause of this observation remains uncertain as studies show that [psychotropic medication itself affects reelin expression](#). Moreover, epigenetic hypotheses aimed at explaining the changed levels of reelin expression^[6] are controversial.^{[7][8]} Total lack of reelin causes a form of [lissencephaly](#). Reelin may also play a role in [Alzheimer's disease](#), [temporal lobe epilepsy](#) and [autism](#).

Reelin's name comes from the abnormal reeling [gait](#) of *reeler* mice,^[9] which were later found to have a deficiency of this brain [protein](#) and were [homozygous](#) for mutation of the RELN gene. The

Reelin



3D ribbon structure of the third reelin repeat domain.^[1]

Available structures

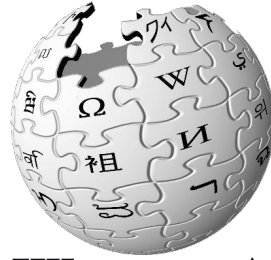
PDB Ortholog search: [PDBe](#) [RCSB](#)

List of PDB id codes [\[show\]](#)

Identifiers

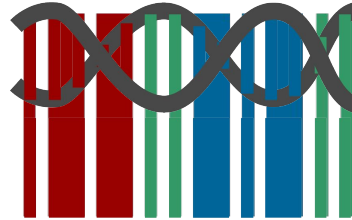
Symbols [RELN](#) ; [LIS2](#); [PRO1598](#); [RL](#)

External [OMIM: 600514](#) [MGI: 103022](#)



WIKIPEDIA
The Free Encyclopedia

is to text



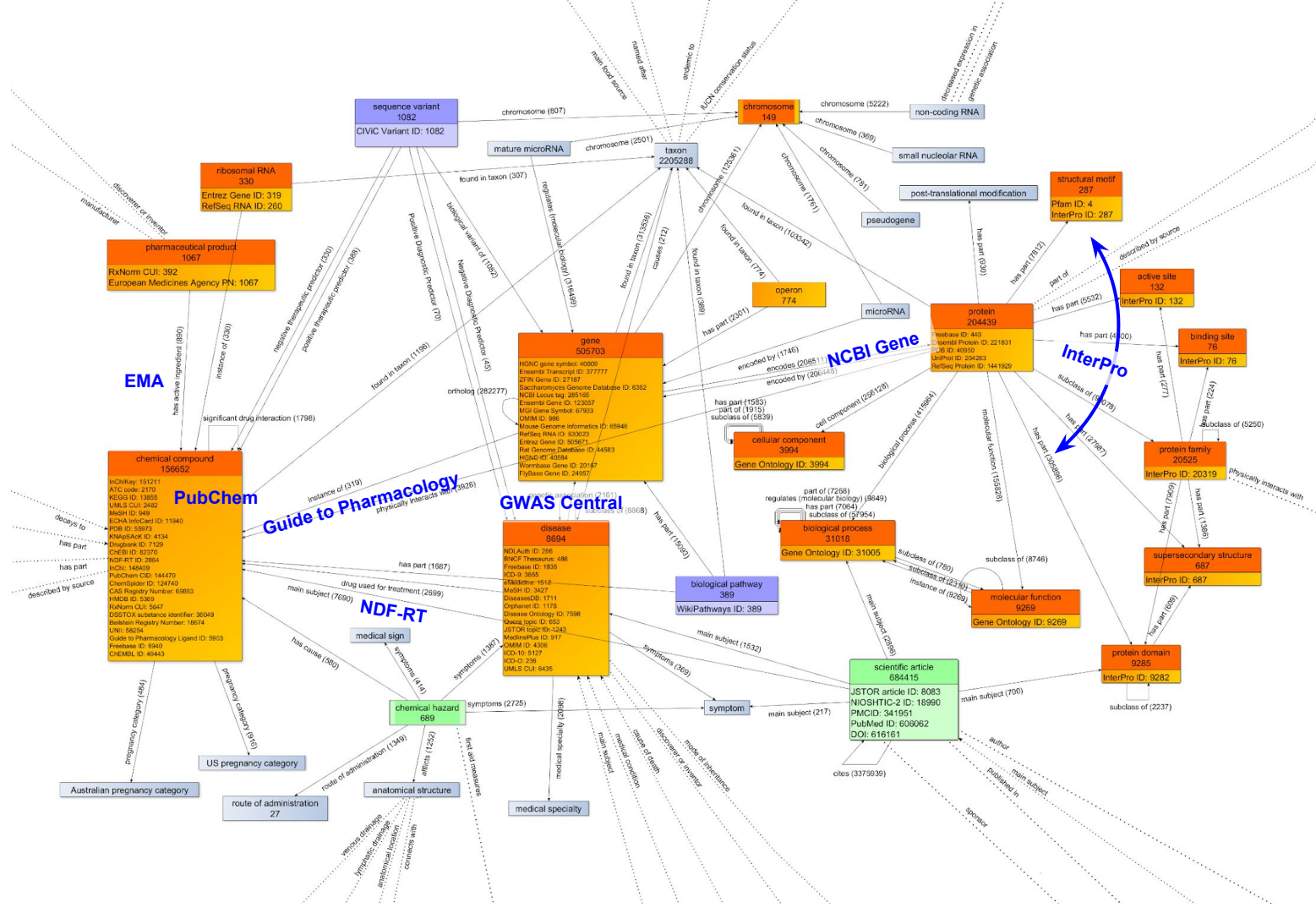
WIKIDATA

is to data

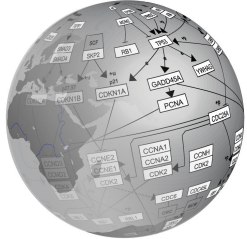
biomedi
cal

“Provide a database of the world’s
knowledge that anyone can edit

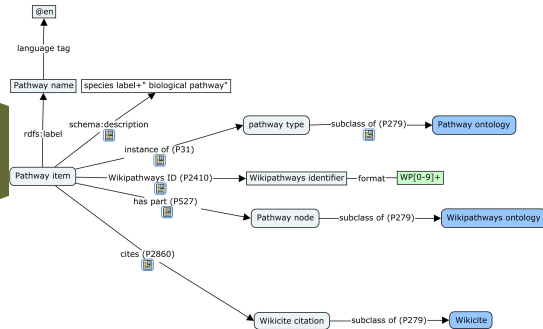
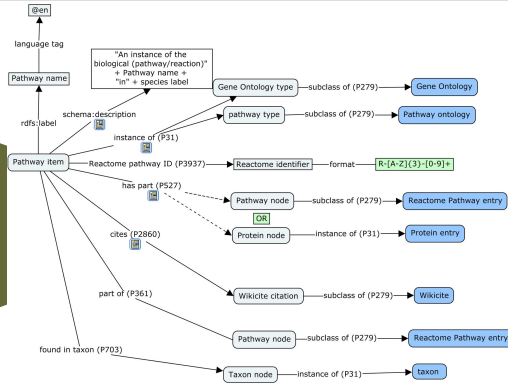
- Denny Vrandečić



Pathways in Wikidata





















WIKIPATHWAYS
Pathways for the People























Reactome Properties in Wikidata

Properties [\[edit \]](#)

V • T • E Wikidata properties related to Reactome (Q2134522)	
Required properties	instance of (P31)  • has part (P527)  •
	part of (P361)  • Reactome pathway ID (P3937)  •
	found in taxon (P703)  • exact match (P2888) 
Optional properties	cites (P2860) 
Reference properties	stated in (P248)  • retrieved (P813)  •
	Reactome pathway ID (P3937) 
Primary sources	Reactome (Q2134522)  
External Ontologies	Pathway Ontology (Q28864280)   •
	Gene Ontology (Q135085)  
Species covered	Homo sapiens (Q15978631)  

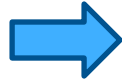
Wikipathways Properties in Wikidata

Properties [\[edit \]](#)

V • T • E Wikidata properties related to biological pathway (Q4915012)	
Required properties	instance of (P31)  • has part (P527)  •
	WikiPathways ID (P2410)  • found in taxon (P703)  •
	exact match (P2888) 
Optional properties	connects with (P2789)  • cites (P2860) 
Reference properties	stated in (P248)  • retrieved (P813)  •
	WikiPathways ID (P2410) 
Primary sources	WikiPathways (Q7999828)  
External Ontologies	Pathway Ontology (Q28864280)   •
	Cell line ontology (Q21039006)   •
	Disease Ontology (Q5282129)  
Species covered	Homo sapiens (Q15978631)  

Simple data retrieval

“Retrieve genes with
GWAS association
with asthma”



39 genes

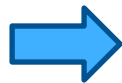
gene	geneLabel		gene	geneLabel		gene	geneLabel		gene	geneLabel
Q5013317	COL22A1		Q18027370	IGSF3		Q18053559	CDHR3		Q14903974	SMAD3
Q14912759	SLC22A5		Q18045382	HPSE2		Q18045669	ATG3		Q18033889	IL1RL1
Q14914243	PSAP		Q18048437	IL33		Q18035037	RAD50		Q17917202	ERBB4
Q14907990	SLC30A8		Q18051900	PYHIN1		Q18036984	FBXL7		Q18027836	IL6R
Q18025002	GAB1		Q17709208	ACO1		Q18033919	XPR1		Q18030185	NOTCH4
Q18035589	C6orf10		Q18027822	IL2RB		Q15326496	RORA		Q18030409	PDE4D
Q18054256	GSDMA		Q18030364	PBX2		Q18042132	GSDMB		Q18045645	IKZF4
Q18058487	C5orf56		Q18037773	ABI3BP		Q18029145	MKLN1		Q18039979	KLHL5
Q18030785	PRKG1		Q18039623	CTNNA3		Q18036729	RAP1GAP2		Q18026947	HLA-DQA1
Q18033424	IL18R1		Q18046350	ZNF665		Q14878303	IL13			

```

1 SELECT DISTINCT ?gene ?geneLabel where {
2   ?gene wdt:P2293 wd:Q35869 .    # gene has genetic association to "asthma"
3   ?gene wdt:P31 wd:Q7187 .       # gene is subclass of "gene"
4   SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
5 }
```


Data integration

“Retrieve genes with
GWAS association
with asthma and gene
product is localized to
membrane”



22 genes

gene	geneLabel	gene	geneLabel	gene	geneLabel	gene	geneLabel
Q14912759	SLC22A5	Q18027370	IGSF3	Q18035037	RAD50	Q18027836	IL6R
Q14914243	PSAP	Q18033424	IL18R1	Q18033919	XPR1	Q18030409	PDE4D
Q14907990	SLC30A8	Q18045382	HPSE2	Q18042132	GSDMB	Q18030185	NOTCH4
Q18035589	C6orf10	Q18027822	IL2RB	Q18036729	RAP1GAP2	Q18026947	HLA-DQA1

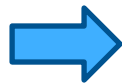
```

1 SELECT DISTINCT ?gene ?geneLabel where {
2   ?gene wdt:P2293 wd:Q35869 . # gene has genetic association to "asthma"
3
4   ?gene wdt:P31 wd:Q7187 .      # gene is subclass of "gene"
5
6   ?gene wdt:P688 ?protein .      # gene encodes a protein
7   ?protein wdt:P681 ?cc .        # protein has a cellular component
8   ?cc wdt:P279*|wdt:P361* wd:Q14349455 . # cell component is 'part of' or 'subclass of' membrane
9
10  SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
11 }
```

http://bit.ly/bosc2017_wikidata

Computing on provenance

“Retrieve genes with GWAS association with asthma and gene product is localized to membrane (non-IEA)”



15 genes

gene	geneLabel		gene	geneLabel		gene	geneLabel
Q14912759	SLC22A5		Q18045382	HPSE2		Q17917202	ERBB4
Q14914243	PSAP		Q18027822	IL2RB		Q18027836	IL6R
Q14907990	SLC30A8		Q14903974	SMAD3		Q18030409	PDE4D
Q18027370	IGSF3		Q18035037	RAD50		Q18030185	NOTCH4
Q18033424	IL18R1		Q18036729	RAP1GAP2		Q18026947	HLA-DQA1

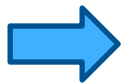
```

6  ?gene wdt:P31 wd:Q7187 ;      # gene is subclass of "gene"
7      wdt:P688 ?protein ;      # gene encodes a protein
8      rdfs:label ?geneLabel .
9  FILTER (lang(?geneLabel) = "en")
10 ?protein p:P681 ?s .          # protein's cell component statement
11     ?s ps:P681 ?cp .          # get statement value
12     FILTER NOT EXISTS { ?s pq:P459 wd:Q23190881 . } # determination method is not IEA
13     ?cp wdt:P279*|wdt:P361* wd:Q14349455 .          # statement value is 'part of' or 'subclass of' membrane
14

```

Leveraging the Disease Ontology structure

“Retrieve genes with GWAS association with any respiratory disease and gene product is localized to membrane (non-IEA)”



31 genes / 8 diseases

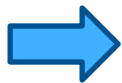
diseaseGALabel	gene_counts	geneList
asthma	15	SMAD3, RAP1GAP2, IL18R1, HPSE2, SLC30A8, SLC22A5, PSAP, ERBB4, HLA-DQA1, IGSF3, IL2RB, IL6R, NOTCH4, PDE4D, RAD50
chronic obstructive pulmonary disease	5	HLA-C, SFTPD, ANXA5, ANXA11, ATP2C2
lung cancer	3	TGM5, VTI1A, PHACTR2
interstitial lung disease	2	DSP, ATP11A
non-small-cell lung carcinoma	2	NALCN, DLST
nasopharynx carcinoma	2	ITGA9, TNFRSF19
adenocarcinoma of the lung	1	BTNL2
pulmonary emphysema	1	BICD1

```

1 SELECT ?diseaseGALabel (count (DISTINCT ?geneLabel))
2 (group_concat(DISTINCT ?geneLabel; separator=" " as geneList)
3   ?gene wdt:P2293 ?diseaseGA . # gene is subclass of "gene" and encodes protein
4   ?diseaseGA wdt:P279* wd:Q3286546 . # to get respiratory diseases
5
6   ?gene wdt:P31 wd:Q7187 ; wdt:P688 ?protein ; # gene is subclass of "gene" and encodes protein
7     rdfs:label ?geneLabel .
8   FILTER (lang(?geneLabel) = "en")
9   ?protein p:P681 ?s . # protein's cell component statement
10  ?s ps:P681 ?cp . # get statement value
11  FILTER NOT EXISTS {
12    ?cp ps:P681 ?cp2 .
13    ?cp2 ps:P681 ?cp3 .
14    ?cp3 ps:P681 ?cp4 .
15    ?cp4 ps:P681 ?cp5 .
16    ?cp5 ps:P681 ?cp6 .
17    ?cp6 ps:P681 ?cp7 .
18    ?cp7 ps:P681 ?cp8 .
19    ?cp8 ps:P681 ?cp9 .
20    ?cp9 ps:P681 ?cp10 .
21    ?cp10 ps:P681 ?cp11 .
22    ?cp11 ps:P681 ?cp12 .
23    ?cp12 ps:P681 ?cp13 .
24    ?cp13 ps:P681 ?cp14 .
25    ?cp14 ps:P681 ?cp15 .
26    ?cp15 ps:P681 ?cp16 .
27    ?cp16 ps:P681 ?cp17 .
28    ?cp17 ps:P681 ?cp18 .
29    ?cp18 ps:P681 ?cp19 .
30    ?cp19 ps:P681 ?cp20 .
31    ?cp20 ps:P681 ?cp21 .
32    ?cp21 ps:P681 ?cp22 .
33    ?cp22 ps:P681 ?cp23 .
34    ?cp23 ps:P681 ?cp24 .
35    ?cp24 ps:P681 ?cp25 .
36    ?cp25 ps:P681 ?cp26 .
37    ?cp26 ps:P681 ?cp27 .
38    ?cp27 ps:P681 ?cp28 .
39    ?cp28 ps:P681 ?cp29 .
40    ?cp29 ps:P681 ?cp30 .
41    ?cp30 ps:P681 ?cp31 .
42    ?cp31 ps:P681 ?cp32 .
43    ?cp32 ps:P681 ?cp33 .
44    ?cp33 ps:P681 ?cp34 .
45    ?cp34 ps:P681 ?cp35 .
46    ?cp35 ps:P681 ?cp36 .
47    ?cp36 ps:P681 ?cp37 .
48    ?cp37 ps:P681 ?cp38 .
49    ?cp38 ps:P681 ?cp39 .
50    ?cp39 ps:P681 ?cp40 .
51    ?cp40 ps:P681 ?cp41 .
52    ?cp41 ps:P681 ?cp42 .
53    ?cp42 ps:P681 ?cp43 .
54    ?cp43 ps:P681 ?cp44 .
55    ?cp44 ps:P681 ?cp45 .
56    ?cp45 ps:P681 ?cp46 .
57    ?cp46 ps:P681 ?cp47 .
58    ?cp47 ps:P681 ?cp48 .
59    ?cp48 ps:P681 ?cp49 .
60    ?cp49 ps:P681 ?cp50 .
61    ?cp50 ps:P681 ?cp51 .
62    ?cp51 ps:P681 ?cp52 .
63    ?cp52 ps:P681 ?cp53 .
64    ?cp53 ps:P681 ?cp54 .
65    ?cp54 ps:P681 ?cp55 .
66    ?cp55 ps:P681 ?cp56 .
67    ?cp56 ps:P681 ?cp57 .
68    ?cp57 ps:P681 ?cp58 .
69    ?cp58 ps:P681 ?cp59 .
70    ?cp59 ps:P681 ?cp60 .
71    ?cp60 ps:P681 ?cp61 .
72    ?cp61 ps:P681 ?cp62 .
73    ?cp62 ps:P681 ?cp63 .
74    ?cp63 ps:P681 ?cp64 .
75    ?cp64 ps:P681 ?cp65 .
76    ?cp65 ps:P681 ?cp66 .
77    ?cp66 ps:P681 ?cp67 .
78    ?cp67 ps:P681 ?cp68 .
79    ?cp68 ps:P681 ?cp69 .
80    ?cp69 ps:P681 ?cp70 .
81    ?cp70 ps:P681 ?cp71 .
82    ?cp71 ps:P681 ?cp72 .
83    ?cp72 ps:P681 ?cp73 .
84    ?cp73 ps:P681 ?cp74 .
85    ?cp74 ps:P681 ?cp75 .
86    ?cp75 ps:P681 ?cp76 .
87    ?cp76 ps:P681 ?cp77 .
88    ?cp77 ps:P681 ?cp78 .
89    ?cp78 ps:P681 ?cp79 .
90    ?cp79 ps:P681 ?cp80 .
91    ?cp80 ps:P681 ?cp81 .
92    ?cp81 ps:P681 ?cp82 .
93    ?cp82 ps:P681 ?cp83 .
94    ?cp83 ps:P681 ?cp84 .
95    ?cp84 ps:P681 ?cp85 .
96    ?cp85 ps:P681 ?cp86 .
97    ?cp86 ps:P681 ?cp87 .
98    ?cp87 ps:P681 ?cp88 .
99    ?cp88 ps:P681 ?cp89 .
100   ?cp89 ps:P681 ?cp90 .
101   ?cp90 ps:P681 ?cp91 .
102   ?cp91 ps:P681 ?cp92 .
103   ?cp92 ps:P681 ?cp93 .
104   ?cp93 ps:P681 ?cp94 .
105   ?cp94 ps:P681 ?cp95 .
106   ?cp95 ps:P681 ?cp96 .
107   ?cp96 ps:P681 ?cp97 .
108   ?cp97 ps:P681 ?cp98 .
109   ?cp98 ps:P681 ?cp99 .
110   ?cp99 ps:P681 ?cp100 .
111   ?cp100 ps:P681 ?cp101 .
112   ?cp101 ps:P681 ?cp102 .
113   ?cp102 ps:P681 ?cp103 .
114   ?cp103 ps:P681 ?cp104 .
115   ?cp104 ps:P681 ?cp105 .
116   ?cp105 ps:P681 ?cp106 .
117   ?cp106 ps:P681 ?cp107 .
118   ?cp107 ps:P681 ?cp108 .
119   ?cp108 ps:P681 ?cp109 .
120   ?cp109 ps:P681 ?cp110 .
121   ?cp110 ps:P681 ?cp111 .
122   ?cp111 ps:P681 ?cp112 .
123   ?cp112 ps:P681 ?cp113 .
124   ?cp113 ps:P681 ?cp114 .
125   ?cp114 ps:P681 ?cp115 .
126   ?cp115 ps:P681 ?cp116 .
127   ?cp116 ps:P681 ?cp117 .
128   ?cp117 ps:P681 ?cp118 .
129   ?cp118 ps:P681 ?cp119 .
130   ?cp119 ps:P681 ?cp120 .
131   ?cp120 ps:P681 ?cp121 .
132   ?cp121 ps:P681 ?cp122 .
133   ?cp122 ps:P681 ?cp123 .
134   ?cp123 ps:P681 ?cp124 .
135   ?cp124 ps:P681 ?cp125 .
136   ?cp125 ps:P681 ?cp126 .
137   ?cp126 ps:P681 ?cp127 .
138   ?cp127 ps:P681 ?cp128 .
139   ?cp128 ps:P681 ?cp129 .
140   ?cp129 ps:P681 ?cp130 .
141   ?cp130 ps:P681 ?cp131 .
142   ?cp131 ps:P681 ?cp132 .
143   ?cp132 ps:P681 ?cp133 .
144   ?cp133 ps:P681 ?cp134 .
145   ?cp134 ps:P681 ?cp135 .
146   ?cp135 ps:P681 ?cp136 .
147   ?cp136 ps:P681 ?cp137 .
148   ?cp137 ps:P681 ?cp138 .
149   ?cp138 ps:P681 ?cp139 .
150   ?cp139 ps:P681 ?cp140 .
151   ?cp140 ps:P681 ?cp141 .
152   ?cp141 ps:P681 ?cp142 .
153   ?cp142 ps:P681 ?cp143 .
154   ?cp143 ps:P681 ?cp144 .
155   ?cp144 ps:P681 ?cp145 .
156   ?cp145 ps:P681 ?cp146 .
157   ?cp146 ps:P681 ?cp147 .
158   ?cp147 ps:P681 ?cp148 .
159   ?cp148 ps:P681 ?cp149 .
160   ?cp149 ps:P681 ?cp150 .
161   ?cp150 ps:P681 ?cp151 .
162   ?cp151 ps:P681 ?cp152 .
163   ?cp152 ps:P681 ?cp153 .
164   ?cp153 ps:P681 ?cp154 .
165   ?cp154 ps:P681 ?cp155 .
166   ?cp155 ps:P681 ?cp156 .
167   ?cp156 ps:P681 ?cp157 .
168   ?cp157 ps:P681 ?cp158 .
169   ?cp158 ps:P681 ?cp159 .
170   ?cp159 ps:P681 ?cp160 .
171   ?cp160 ps:P681 ?cp161 .
172   ?cp161 ps:P681 ?cp162 .
173   ?cp162 ps:P681 ?cp163 .
174   ?cp163 ps:P681 ?cp164 .
175   ?cp164 ps:P681 ?cp165 .
176   ?cp165 ps:P681 ?cp166 .
177   ?cp166 ps:P681 ?cp167 .
178   ?cp167 ps:P681 ?cp168 .
179   ?cp168 ps:P681 ?cp169 .
180   ?cp169 ps:P681 ?cp170 .
181   ?cp170 ps:P681 ?cp171 .
182   ?cp171 ps:P681 ?cp172 .
183   ?cp172 ps:P681 ?cp173 .
184   ?cp173 ps:P681 ?cp174 .
185   ?cp174 ps:P681 ?cp175 .
186   ?cp175 ps:P681 ?cp176 .
187   ?cp176 ps:P681 ?cp177 .
188   ?cp177 ps:P681 ?cp178 .
189   ?cp178 ps:P681 ?cp179 .
190   ?cp179 ps:P681 ?cp180 .
191   ?cp180 ps:P681 ?cp181 .
192   ?cp181 ps:P681 ?cp182 .
193   ?cp182 ps:P681 ?cp183 .
194   ?cp183 ps:P681 ?cp184 .
195   ?cp184 ps:P681 ?cp185 .
196   ?cp185 ps:P681 ?cp186 .
197   ?cp186 ps:P681 ?cp187 .
198   ?cp187 ps:P681 ?cp188 .
199   ?cp188 ps:P681 ?cp189 .
200   ?cp189 ps:P681 ?cp190 .
201   ?cp190 ps:P681 ?cp191 .
202   ?cp191 ps:P681 ?cp192 .
203   ?cp192 ps:P681 ?cp193 .
204   ?cp193 ps:P681 ?cp194 .
205   ?cp194 ps:P681 ?cp195 .
206   ?cp195 ps:P681 ?cp196 .
207   ?cp196 ps:P681 ?cp197 .
208   ?cp197 ps:P681 ?cp198 .
209   ?cp198 ps:P681 ?cp199 .
210   ?cp199 ps:P681 ?cp200 .
211   ?cp200 ps:P681 ?cp201 .
212   ?cp201 ps:P681 ?cp202 .
213   ?cp202 ps:P681 ?cp203 .
214   ?cp203 ps:P681 ?cp204 .
215   ?cp204 ps:P681 ?cp205 .
216   ?cp205 ps:P681 ?cp206 .
217   ?cp206 ps:P681 ?cp207 .
218   ?cp207 ps:P681 ?cp208 .
219   ?cp208 ps:P681 ?cp209 .
220   ?cp209 ps:P681 ?cp210 .
221   ?cp210 ps:P681 ?cp211 .
222   ?cp211 ps:P681 ?cp212 .
223   ?cp212 ps:P681 ?cp213 .
224   ?cp213 ps:P681 ?cp214 .
225   ?cp214 ps:P681 ?cp215 .
226   ?cp215 ps:P681 ?cp216 .
227   ?cp216 ps:P681 ?cp217 .
228   ?cp217 ps:P681 ?cp218 .
229   ?cp218 ps:P681 ?cp219 .
230   ?cp219 ps:P681 ?cp220 .
231   ?cp220 ps:P681 ?cp221 .
232   ?cp221 ps:P681 ?cp222 .
233   ?cp222 ps:P681 ?cp223 .
234   ?cp223 ps:P681 ?cp224 .
235   ?cp224 ps:P681 ?cp225 .
236   ?cp225 ps:P681 ?cp226 .
237   ?cp226 ps:P681 ?cp227 .
238   ?cp227 ps:P681 ?cp228 .
239   ?cp228 ps:P681 ?cp229 .
240   ?cp229 ps:P681 ?cp230 .
241   ?cp230 ps:P681 ?cp231 .
242   ?cp231 ps:P681 ?cp232 .
243   ?cp232 ps:P681 ?cp233 .
244   ?cp233 ps:P681 ?cp234 .
245   ?cp234 ps:P681 ?cp235 .
246   ?cp235 ps:P681 ?cp236 .
247   ?cp236 ps:P681 ?cp237 .
248   ?cp237 ps:P681 ?cp238 .
249   ?cp238 ps:P681 ?cp239 .
250   ?cp239 ps:P681 ?cp240 .
251   ?cp240 ps:P681 ?cp241 .
252   ?cp241 ps:P681 ?cp242 .
253   ?cp242 ps:P681 ?cp243 .
254   ?cp243 ps:P681 ?cp244 .
255   ?cp244 ps:P681 ?cp245 .
256   ?cp245 ps:P681 ?cp246 .
257   ?cp246 ps:P681 ?cp247 .
258   ?cp247 ps:P681 ?cp248 .
259   ?cp248 ps:P681 ?cp249 .
260   ?cp249 ps:P681 ?cp250 .
261   ?cp250 ps:P681 ?cp251 .
262   ?cp251 ps:P681 ?cp252 .
263   ?cp252 ps:P681 ?cp253 .
264   ?cp253 ps:P681 ?cp254 .
265   ?cp254 ps:P681 ?cp255 .
266   ?cp255 ps:P681 ?cp256 .
267   ?cp256 ps:P681 ?cp257 .
268   ?cp257 ps:P681 ?cp258 .
269   ?cp258 ps:P681 ?cp259 .
270   ?cp259 ps:P681 ?cp260 .
271   ?cp260 ps:P681 ?cp261 .
272   ?cp261 ps:P681 ?cp262 .
273   ?cp262 ps:P681 ?cp263 .
274   ?cp263 ps:P681 ?cp264 .
275   ?cp264 ps:P681 ?cp265 .
276   ?cp265 ps:P681 ?cp266 .
277   ?cp266 ps:P681 ?cp267 .
278   ?cp267 ps:P681 ?cp268 .
279   ?cp268 ps:P681 ?cp269 .
280   ?cp269 ps:P681 ?cp270 .
281   ?cp270 ps:P681 ?cp271 .
282   ?cp271 ps:P681 ?cp272 .
283   ?cp272 ps:P681 ?cp273 .
284   ?cp273 ps:P681 ?cp274 .
285   ?cp274 ps:P681 ?cp275 .
286   ?cp275 ps:P681 ?cp276 .
287   ?cp276 ps:P681 ?cp277 .
288   ?cp277 ps:P681 ?cp278 .
289   ?cp278 ps:P681 ?cp279 .
290   ?cp279 ps:P681 ?cp280 .
291   ?cp280 ps:P681 ?cp281 .
292   ?cp281 ps:P681 ?cp282 .
293   ?cp282 ps:P681 ?cp283 .
294   ?cp283 ps:P681 ?cp284 .
295   ?cp284 ps:P681 ?cp285 .
296   ?cp285 ps:P681 ?cp286 .
297   ?cp286 ps:P681 ?cp287 .
298   ?cp287 ps:P681 ?cp288 .
299   ?cp288 ps:P681 ?cp289 .
300   ?cp289 ps:P681 ?cp290 .
301   ?cp290 ps:P681 ?cp291 .
302   ?cp291 ps:P681 ?cp292 .
303   ?cp292 ps:P681 ?cp293 .
304   ?cp293 ps:P681 ?cp294 .
305   ?cp294 ps:P681 ?cp295 .
306   ?cp295 ps:P681 ?cp296 .
307   ?cp296 ps:P681 ?cp297 .
308   ?cp297 ps:P681 ?cp298 .
309   ?cp298 ps:P681 ?cp299 .
310   ?cp299 ps:P681 ?cp300 .
311   ?cp300 ps:P681 ?cp301 .
312   ?cp301 ps:P681 ?cp302 .
313   ?cp302 ps:P681 ?cp303 .
314   ?cp303 ps:P681 ?cp304 .
315   ?cp304 ps:P681 ?cp305 .
316   ?cp305 ps:P681 ?cp306 .
317   ?cp306 ps:P681 ?cp307 .
318   ?cp307 ps:P681 ?cp308 .
319   ?cp308 ps:P681 ?cp309 .
320   ?cp309 ps:P681 ?cp310 .
321   ?cp310 ps:P681 ?cp311 .
322   ?cp311 ps:P681 ?cp312 .
323   ?cp312 ps:P681 ?cp313 .
324   ?cp313 ps:P681 ?cp314 .
325   ?cp314 ps:P681 ?cp315 .
326   ?cp315 ps:P681 ?cp316 .
327   ?cp316 ps:P681 ?cp317 .
328   ?cp317 ps:P681 ?cp318 .
329   ?cp318 ps:P681 ?cp319 .
330   ?cp319 ps:P681 ?cp320 .
331   ?cp320 ps:P681 ?cp321 .
332   ?cp321 ps:P681 ?cp322 .
333   ?cp322 ps:P681 ?cp323 .
334   ?cp323 ps:P681 ?cp324 .
335   ?cp324 ps:P681 ?cp325 .
336   ?cp325 ps:P681 ?cp326 .
337   ?cp326 ps:P681 ?cp327 .
338   ?cp327 ps:P681 ?cp328 .
339   ?cp328 ps:P681 ?cp329 .
340   ?cp329 ps:P681 ?cp330 .
341   ?cp330 ps:P681 ?cp331 .
342   ?cp331 ps:P681 ?cp332 .
343   ?cp332 ps:P681 ?cp333 .
344   ?cp333 ps:P681 ?cp334 .
345   ?cp334 ps:P681 ?cp335 .
346   ?cp335 ps:P681 ?cp336 .
347   ?cp336 ps:P681 ?cp337 .
348   ?cp337 ps:P681 ?cp338 .
349   ?cp338 ps:P681 ?cp339 .
350   ?cp339 ps:P681 ?cp340 .
351   ?cp340 ps:P681 ?cp341 .
352   ?cp341 ps:P681 ?cp342 .
353   ?cp342 ps:P681 ?cp343 .
354   ?cp343 ps:P681 ?cp344 .
355   ?cp344 ps:P681 ?cp345 .
356   ?cp345 ps:P681 ?cp346 .
357   ?cp346 ps:P681 ?cp347 .
358   ?cp347 ps:P681 ?cp348 .
359   ?cp348 ps:P681 ?cp349 .
360   ?cp349 ps:P681 ?cp350 .
361   ?cp350 ps:P681 ?cp351 .
362   ?cp351 ps:P681 ?cp352 .
363   ?cp352 ps:P681 ?cp353 .
364   ?cp353 ps:P681 ?cp354 .
365   ?cp354 ps:P681 ?cp355 .
366   ?cp355 ps:P681 ?cp356 .
367   ?cp356 ps:P681 ?cp357 .
368   ?cp357 ps:P681 ?cp358 .
369   ?cp358 ps:P681 ?cp359 .
370   ?cp359 ps:P681 ?cp360 .
371   ?cp360 ps:P681 ?cp361 .
372   ?cp361 ps:P681 ?cp362 .
373   ?cp362 ps:P681 ?cp363 .
374   ?cp363 ps:P681 ?cp364 .
375   ?cp364 ps:P681 ?cp365 .
376   ?cp365 ps:P681 ?cp366 .
377   ?cp366 ps:P681 ?cp367 .
378   ?cp367 ps:P681 ?cp368 .
379   ?cp368 ps:P681 ?cp369 .
380   ?cp369 ps:P681 ?cp370 .
381   ?cp370 ps:P681 ?cp371 .
382   ?cp371 ps:P681 ?cp372 .
383   ?cp372 ps:P681 ?cp373 .
384   ?cp373 ps:P681 ?cp374 .
385   ?cp374 ps:P681 ?cp375 .
386   ?cp375 ps:P681 ?cp376 .
387   ?cp376 ps:P681 ?cp377 .
388   ?cp377 ps:P681 ?cp378 .
389   ?cp378 ps:P681 ?cp379 .
390   ?cp379 ps:P681 ?cp380 .
391   ?cp380 ps:P681 ?cp381 .
392   ?cp381 ps:P681 ?cp382 .
393   ?cp382 ps:P681 ?cp383 .
394   ?cp383 ps:P681 ?cp384 .
395   ?cp384 ps:P681 ?cp385 .
396   ?cp385 ps:P681 ?cp386 .
397   ?cp386 ps:P681 ?cp387 .
398   ?cp387 ps:P681 ?cp388 .
399   ?cp388 ps:P681 ?cp389 .
400   ?cp389 ps:P681 ?cp390 .
401   ?cp390 ps:P681 ?cp391 .
402   ?cp391 ps:P681 ?cp392 .
403   ?cp392 ps:P681 ?cp393 .
404   ?cp393 ps:P681 ?cp394 .
405   ?cp394 ps:P681 ?cp395 .
406   ?cp395 ps:P681 ?cp396 .
407   ?cp396 ps:P681 ?cp397 .
408   ?cp397 ps:P681 ?cp398 .
409   ?cp398 ps:P681 ?cp399 .
410   ?cp399 ps:P681 ?cp400 .
411   ?cp400 ps:P681 ?cp401 .
412   ?cp401 ps:P681 ?cp402 .
413   ?cp402 ps:P681 ?cp403 .
414   ?cp403 ps:P681 ?cp404 .
415   ?cp404 ps:P681 ?cp405 .
416   ?cp405 ps:P681 ?cp406 .
417   ?cp406 ps:P681 ?cp407 .
418   ?cp407 ps:P681 ?cp408 .
419   ?cp408 ps:P681 ?cp409 .
420   ?cp409 ps:P681 ?cp410 .
421   ?cp410 ps:P681 ?cp411 .
422   ?cp411 ps:P681 ?cp412 .
423   ?cp412 ps:P681 ?cp413 .
424   ?cp413 ps:P681 ?cp414 .
425   ?cp414 ps:P681 ?cp415 .
426   ?cp415 ps:P681 ?cp416 .
427   ?cp416 ps:P681 ?cp417 .
428   ?cp417 ps:P681 ?cp418 .
429   ?cp418 ps:P681 ?cp419 .
430   ?cp419 ps:P681 ?cp420 .
431   ?cp420 ps:P681 ?cp421 .
432   ?cp421 ps:P681 ?cp422 .
433   ?cp422 ps:P681 ?cp423 .
434   ?cp423 ps:P681 ?cp424 .
435   ?cp424 ps:P681 ?cp425 .
436   ?cp425 ps:P681 ?cp426 .
437   ?cp426 ps:P681 ?cp427 .
438   ?cp427 ps:P681 ?cp428 .
439   ?cp428 ps:P681 ?cp429 .
440   ?cp429 ps:P681 ?cp430 .
441   ?cp430 ps:P681 ?cp431 .
442   ?cp431 ps:P681 ?cp432 .
443   ?cp432 ps:P681 ?cp433 .
444   ?cp433 ps:P681 ?cp434 .
445   ?cp434 ps:P681 ?cp435 .
446   ?cp435 ps:P681 ?cp436 .
447   ?cp436 ps:P681 ?cp437 .
448   ?cp437 ps:P681 ?cp438 .
449   ?cp438 ps:P681 ?cp439 .
450   ?cp439 ps:P681 ?cp440 .
451   ?cp440 ps:P681 ?cp441 .
452   ?cp441 ps:P681 ?cp442 .
453   ?cp442 ps:P681 ?cp443 .
454   ?cp443 ps:P681 ?cp444 .
455   ?cp444 ps:P681 ?cp445 .
456   ?cp445 ps:P681 ?cp446 .
457   ?cp446 ps:P681 ?cp447 .
458   ?cp447 ps:P681 ?cp448 .
459   ?cp448 ps:P681 ?cp449 .
460   ?cp449 ps:P681 ?cp450 .
461   ?cp450 ps:P681 ?cp451 .
462   ?cp451 ps:P681 ?cp452 .
463   ?cp452 ps:P681 ?cp453 .
464   ?cp453 ps:P681 ?cp454 .
465   ?cp454 ps:P681 ?cp455 .
466   ?cp455 ps:P681 ?cp456 .
467   ?cp456 ps:P681 ?cp457 .
468   ?cp457 ps:P681 ?cp458 .
469   ?cp458 ps:P681 ?cp459 .
470   ?cp459 ps:P681 ?cp460 .
471   ?cp460 ps:P681 ?cp461 .
472   ?cp461 ps:P681 ?cp462 .
473   ?cp462 ps:P681 ?cp463 .
474   ?cp463 ps:P681 ?cp464 .
475   ?cp464 ps:P681 ?cp465 .
476   ?cp465 ps:P681 ?cp466 .
477   ?cp466 ps:P681 ?cp467 .
478   ?cp467 ps:P681 ?cp468 .
479   ?cp468 ps:P681 ?cp469 .
480   ?cp469 ps:P681 ?cp470 .
481   ?cp470 ps:P681 ?cp471 .
482   ?cp471 ps:P681 ?cp472 .
483   ?cp472 ps:P681 ?cp473 .
484   ?cp473 ps:P681 ?cp474 .
485   ?cp474 ps:P681 ?cp475 .
486   ?cp475 ps:P681 ?cp476 .
487   ?cp476 ps:P681 ?cp477 .
488   ?cp477 ps:P681 ?cp478 .
489   ?cp478 ps:P681 ?cp479 .
490   ?cp479 ps:P681 ?cp480 .
491   ?cp480 ps:P681 ?cp481 .
492   ?cp481 ps:P681 ?cp482 .
493   ?cp482 ps:P681 ?cp483 .
494   ?cp483 ps:P681 ?cp484 .
495   ?cp484 ps:P681 ?cp485 .
496   ?cp485 ps:P681 ?cp486 .
497   ?cp486 ps:P681 ?cp487 .
498   ?cp487 ps:P681 ?cp488 .
499   ?cp488 ps:P681 ?cp489 .
500   ?cp489 ps:P681 ?cp490 .
501   ?cp490 ps:P681 ?cp491 .
502   ?cp491 ps:P681 ?cp492 .
503   ?cp492 ps:P681 ?cp493 .
504   ?cp493 ps:P681 ?cp494 .
505   ?cp494 ps:P681 ?cp495 .
506   ?cp495 ps:P681 ?cp496 .
507   ?cp496 ps:P681 ?cp497 .
508   ?cp497 ps:P681 ?cp498 .
509   ?cp498 ps:P681 ?cp499 .
510   ?cp499 ps:P681 ?cp500 .
511   ?cp500 ps:P681 ?cp501 .
512   ?cp501 ps:P681 ?cp502 .
513   ?cp502 ps:P681 ?cp503 .
514   ?cp503 ps:P681 ?cp504 .
515   ?cp504 ps:P681 ?cp505 .
516   ?cp505 ps:P681 ?cp506 .
517   ?cp506 ps:P681 ?cp507 .
518   ?cp507 ps:P681 ?cp508 .
519   ?cp508 ps:P681 ?cp509 .
520   ?cp509 ps:P681 ?cp510 .
521   ?cp510 ps:P681 ?cp511 .
522   ?cp511 ps:P681 ?cp512 .
523   ?cp512 ps:P681 ?cp513 .
524   ?cp513 ps:P681 ?cp514 .
525   ?cp514 ps:P681 ?cp515 .
526   ?cp515 ps:P681 ?cp516 .
527   ?cp516 ps:P681 ?cp517 .
528   ?cp517 ps:P681 ?cp518 .
529   ?cp518 ps:P681 ?cp519 .
530   ?cp519 ps:P681 ?cp520 .
531   ?cp520 ps:P681 ?cp521 .
532   ?cp521 ps:P681 ?cp522 .
533   ?cp522 ps:P681 ?cp523 .
534   ?cp523 ps:P681 ?cp524 .
535   ?cp524 ps:P681 ?cp525 .
536   ?cp525 ps:P681 ?cp526 .
537   ?cp526 ps:P681 ?cp527 .
538   ?cp527 ps:P681 ?cp528 .
539   ?cp528 ps:P681 ?cp529 .
540   ?cp529 ps:P681 ?cp530 .
541   ?cp530 ps:P681 ?cp531 .
542   ?cp531 ps:P681 ?cp532 .
543   ?cp532 ps:P681 ?cp533 .
544   ?cp533 ps:P681 ?cp534 .
545   ?cp534 ps:P681 ?cp535 .
546   ?cp535 ps:P681 ?cp536 .
547   ?cp536 ps:P681 ?cp537 .
548   ?cp537 ps:P681 ?cp538 .
549   ?cp538 ps:P681 ?cp539 .
550   ?cp539 ps:P681 ?cp540 .
551   ?cp540 ps:P681 ?cp541 .
552   ?cp541 ps:P681 ?cp542 .
553   ?cp542 ps:P681 ?cp543 .
554   ?cp543 ps:P681 ?cp544 .
555   ?cp544 ps:P681 ?cp545 .
556   ?cp545 ps:P681 ?cp546 .
557   ?cp546 ps:P681 ?cp547 .
558   ?cp547 ps:P681 ?cp548 .
559   ?cp548 ps:P681 ?cp549 .
560   ?cp549 ps:P681 ?cp550 .
561   ?cp550 ps:P681 ?cp551 .
562   ?cp551 ps:P681 ?cp552 .
563   ?cp552 ps:P681 ?cp553 .
564   ?cp553 ps:P681 ?cp554 .
565   ?cp554 ps:P681 ?cp555 .
566   ?cp555 ps:P681 ?cp556 .
567   ?cp556 ps:P681 ?cp557 .
568   ?cp557 ps:P681 ?cp558 .
569   ?cp558 ps:P681 ?cp559 .
570   ?cp559 ps:P681 ?cp560 .
571   ?cp560 ps:P681 ?cp561 .
572   ?cp561 ps:P681 ?cp562 .
573   ?cp562 ps:P681 ?cp563 .
57
```

Opportunistic integration

“Retrieve genes with GWAS association with any respiratory disease and gene product is localized to membrane (non-IEA) **and show causative chemical hazards**”



4 diseases / 6 chemical hazards

diseaseGALabel	exposureLabel
lung cancer	arsenic pentoxide exposure
lung cancer	HN1 exposure
lung cancer	mechlorethamine exposure
lung cancer	HN3 exposure
asthma	Phenacyl chloride exposure
pulmonary emphysema	phosgene exposure

```

11 .cp wdt:P279 | wdt:P501 wd:Q1167512 . # statement value is part of or is
12
13 ?exposure wdt:P1542 ?diseaseGA . # something causes disease
14 ?exposure wdt:P279 wd:Q21167512 . # and that something is a chemical hazard
15
16 SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
17 }

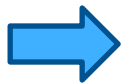
```

http://bit.ly/bosc2017_wikidata

... and show associated pathways

16 genes / 59 pathways

“Retrieve genes with GWAS association with any respiratory disease and gene product is localized to membrane (non-IEA), show causative chemical hazards and **show pathways where they have a role.**”



gene	pathway
SMAD3	Androgen receptor signaling pathway
SMAD3	TGF-beta Receptor Signaling
SMAD3	mechlorethamine exposure
HLA-C	Allograft Rejection
SFTPD	Regulation of toll-like receptor signaling pathway
....

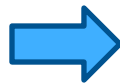
```

11  .cp wdt:P279 | wdt:P301 wd:Q4915012 .
12
13  ?pathway wdt:P31 wd:Q4915012 ;           # instance of a biological pathway
14      wdt:P527 ?gene .
15
16  SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
17  }
```

http://bit.ly/bosc2017_wikidata

... and show remote pathway annotations

“Retrieve genes with GWAS association with any respiratory disease and gene product is localized to membrane (non-IEA), show causative chemical hazards, show pathways where they have a role and **show pathway annotations in WikiPathways**”



16 genes / 59 pathways

Remote pathway annotations	
Pathway ontology	55
Cell line ontology	12
Disease ontology	11
....

```

11 .cp wdt:P2888 | wdt:P501 | wdt:P14349493 .
12
13 wdt:P2888 ?source_pathway .
14 SERVICE <http://sparql.wikipathways.org/> {
15   ?wp_pathway dc:identifier ?source_pathway ;
16   wp:ontologyTag ?pwTag .
17 }

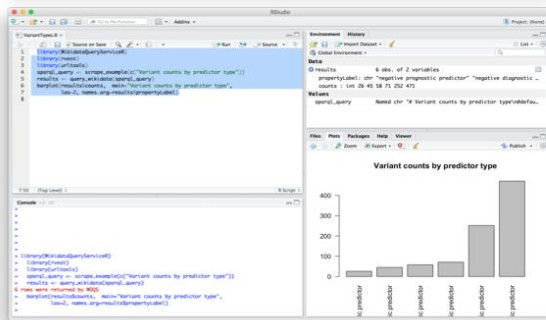
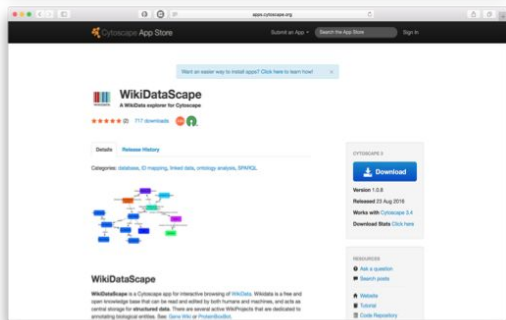
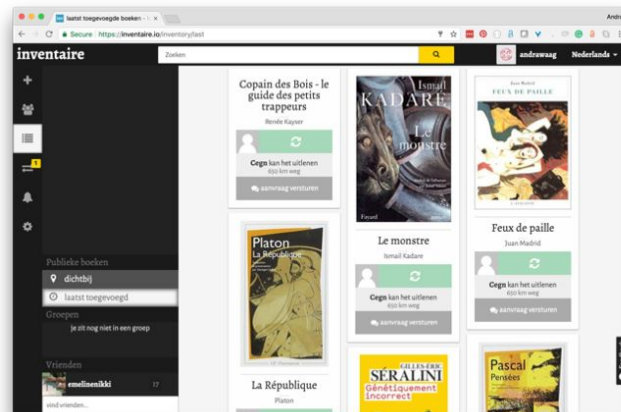
```


Tools using Wikidata

<http://www.wikigenomes.org> Wikipedia and Wikidata
google plugin



<http://inventaire.io>



Cytoscape

R plugins

Availability

www.wikidata.org/wiki/User:Pathwaybot

- www.wikidata.org/wiki/User:Pathwaybot/data_model - Semantic models
- www.wikidata.org/wiki/User:Pathwaybot/query_examples – Pathway examples
- www.wikidata.org/wiki/User:ProteinBoxBot/SPARQL_Examples– Genewiki example

github.com/SuLab/GeneWikiCentral

- github.com/SuLab/wikidataintegrator – python module for Wikidata
- github.com/SuLab/scheduled-bots – bot automation framework
- github.com/SuLab/Genewiki-ShEx – data models

Structure of the federated query

```
SELECT * WHERE {
```

```
?localitem ?localproperty ?shareditem .
```

Local pattern

```
SERVICE <http://example.remote.sparql> {  
  ?shareditem ?remoteproperty ?remoteitem .  
}
```

Remote pattern

}
A federated query contains query patterns for both the local endpoint (green box) and a remote endpoint (blue box). The address of the remote SPARQL endpoint is expressed with the SERVICE keyword

Why federate

- Incorporate more fine-grained data not captured in Wikidata
- Link to other (non-public licenses)
- Link to volatile data

Wikidata Query Service support federated queries

- From Wikidata to an external SPARQL endpoint (Wikipathways)
- From a remote SPARQL endpoint to Wikidata
- From a local SPARQL endpoint to Wikidata
- From a local Wikibase (with WDQS) to Wikidata

PREFIX wp: <http://vocabularies.wikipathways.org/wp#>
PREFIX dcterms: <http://purl.org/dc/terms/>
PREFIX dc: <http://purl.org/dc/elements/1.1/>
SELECT DISTINCT ?metabolite1Label ?metabolite2Label ?mass1 ?mass2 WITH {

SELECT ?metabolite1 ?metabolite2 WHERE {
?pathwayItem wdt:P2410 "WP706";
wdt:P2888 ?pwIri.

Wikidata

SERVICE <http://sparql.wikipathways.org/> {
?pathway dc:identifier ?pwIri.
?interaction rdf:type wp:Interaction;
wp:participants ?wpmb1, ?wpmb2;
dcterms:isPartOf ?pathway.

FILTER (?wpmb1 != ?wpmb2)
?wpmb1 wp:bdbWikidata ?metabolite1.
?wpmb2 wp:bdbWikidata ?metabolite2.
}

Wikipathways

} AS %metabolites WHERE {

INCLUDE %metabolites.
?metabolite1 wdt:P2067 ?mass1.
?metabolite2 wdt:P2067 ?mass2.
SERVICE wikibase:label { bd:serviceParam wikibase:language "[AUTO_LANGUAGE],en". }

Wikidata

}

From a remote SPARQL endpoint to Wikidata



SPARQL Downloads

Documentation/Help

Your query

Add common prefixes

```
20 SELECT DISTINCT ?wd_item ?physically_interacts_with ?interactswithLabel ?type ?iri ?uniprot ?text WHERE {
21   {SELECT * WHERE { ?iri a up:Protein ;
22     up:organism taxon:9606 ;
23     up:annotation ?annotation .
24     ?annotation a up:Natural_Variant_Annotation ;
25     rdfs:comment ?text .
26     FILTER (CONTAINS(?text, 'loss of function'))
27   }}
28   SERVICE <https://query.wikidata.org/bigdata/namespace/wdq/sparql> {
29     VALUES ?use {wd:Q427492}
30     ?wd_item wdt:P352 ?uniprot ;
31     wdt:P129 ?physically_interacts_with ;
32     wdt:P2888 ?iri ;
33     wdt:P703 wd:Q15978631 .
34     ?wd_item p:P129 ?phys_interacts_with_node .
35     ?phys_interacts_with_node ps:P129 ?physically_interacts_with ;
36     pg:P366 ?use .
37     ?physically_interacts_with wdt:P31 ?type ;
38     rdfs:label ?interactswithLabel .
39     FILTER (lang(?interactswithLabel) = "en")
40   }}
```

UniProt

Wikidata

Submit Query

Cancel

We use ShEx to describe contents of a remote endpoint (<http://shex.io>)

- Draft Shape Expression for the Europeana Sparql endpoint:

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX edm: <http://www.europeana.eu/schemas/edm/>
PREFIX rdvocab: <http://rdvocab.info/ElementsGr2/>
PREFIX dce: <http://purl.org/dc/elements/1.1/>
PREFIX dct: <http://purl.org/dc/terms/>
PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
PREFIX dbp: <http://dbpedia.org/resource/>
PREFIX oaterms: <http://www.openarchives.org/ore/terms/>
PREFIX bibgalicia: <http://www.galiciana.bibliotecadegalicia.xunta.es/aut/>
PREFIX gtaa: <http://data.beeldengeluid.nl/gtaa/>
PREFIX owl: <http://www.w3.org/2002/07/owl#>
```

```
<aggregation> {
  edm:dataProvider LITERAL* ;
  edm:rights <licenses>? ;
  edm:isShownBy IRI* ;
}
```

```
<concept> {
  rdf:type skos:Concept ;
  skos:prefLabel LITERAL* ;
  skos:broader [dbp:~]* ;
  skos:exactMatch IRI* ;
  skos:closeMatch IRI* ;
  skos:related [dbp:~ bibgalicia:~]* ;
  skos:notation LITERAL? ;
}
```