

1 **Knowledge Beacons: Web Service Workflow for FAIR**

2 **Data Harvesting of Distributed Biomedical Knowledge**

3 Lance M. Hannestad^{1,2}, Vlado Dančík³, Meera Godden^{1,2}, Imelda W. Suen^{1,4}, Kenneth C.
4 Huellas-Bruskiewicz^{1,5}, Benjamin M. Good⁶, Christopher J. Mungall^{6¶}, Richard M.
5 Bruskiewich^{1*¶}

6 ¹ STAR Informatics / Delphinai Corporation, Sooke, BC, Canada

7 ² Department of Computing Science, Simon Fraser University, Burnaby, BC, Canada

8 ³ Chemical Biology and Therapeutics Science Program, The Broad Institute, Cambridge,
9 MA, United States of America

10 ⁴ School of Computing Science, University of British Columbia, Vancouver, BC, Canada

11 ⁵ School of Interactive Arts and Technology, Simon Fraser University, Burnaby, BC, Canada

12 ⁶ Lawrence Berkeley National Laboratory, Berkeley, CA, United States of America

13 * Corresponding author

14 E-mail: richard.bruskiewich@delphinai.com (RMB)

15 ¶ CJM and RMB are Joint Senior Authors

16 **Abstract**

17 The continually expanding distributed global compendium of biomedical knowledge is
18 diffuse, heterogeneous and huge, posing a serious challenge for biomedical researchers in
19 knowledge harvesting: accessing, compiling, integrating and interpreting data, information and
20 knowledge. In order to accelerate research towards effective medical treatments and optimizing

21 health, it is critical that efficient and automated tools for identifying key research concepts and
22 their experimentally discovered interrelationships are developed.

23 As an activity within the feasibility phase of a project called “Translator”
24 (<https://ncats.nih.gov/translator>) funded by the National Center for Advancing Translational
25 Sciences (NCATS) to develop a biomedical science knowledge management platform, we de-
26 signed a Representational State Transfer (REST) web services Application Programming Inter-
27 face (API) specification, which we call a Knowledge Beacon. Knowledge Beacons provide a
28 standardized basic workflow for the discovery of concepts, their relationships and associated
29 supporting evidence from distributed online repositories of biomedical knowledge. This specifi-
30 cation also enforces the annotation of knowledge concepts and statements to the NCATS en-
31 dorsed the Biolink Model data model and semantic encoding standards
32 (<https://biolink.github.io/biolink-model/>). Implementation of this API on top of diverse
33 knowledge sources potentially enables their uniform integration behind client software which
34 will facilitate research access and integration of biomedical knowledge.

35 **Availability:** The API and associated software is open source and currently available for
36 access at <https://github.com/NCATS-Tangerine/translator-knowledge-beacon>.

37 **Introduction**

38 A serious challenge to impactful biomedical research is the one that biomedical research-
39 ers encounter when identifying and accessing pertinent information: the diffuse and voluminous
40 nature of such data and knowledge. The large, rapidly growing compendium of published sci-
41 entific literature is characterized by diverse data encoding standards; numerous, distinct, heteroge-
42 neous, large and often siloed public research data repositories; relatively inaccessible health rec-

43 ords; numerous clinical trial and adverse event reports, all spread across disease communities
44 and biomedical disciplines. The current distributed nature of this knowledge and associated (me-
45 ta-)data silos impedes the discovery of related concepts and the relationships between them, an
46 activity one might call “Knowledge Harvesting”. Many efforts to overcome this challenge focus
47 on data management principles to make such resources “Findable, Interoperable, Accessible and
48 Reusable” (FAIR) [1,2].

49 Web access to bioinformatics data spans many generations of web service standards
50 tagged with many acronyms, e.g. CORBA [3], SOAP/BioMOBY [4] and SADI [5], the latter an
51 exemplar of the more general paradigm of “Linked Open Data” using OWL/RDF and SPARQL
52 technology, including Linked Open Fragments [6].

53 A popular web service standard currently in use is the Swagger 2.0 or OpenAPI 3.0 spec-
54 ified REST API (<https://github.com/OAI>). Many extant online biomedical data sources currently
55 provide such REST API implementations for accessing their data. API registries exist to index
56 such APIs to facilitate access (for example, the Smart API Registry; <https://smart-api.info/>) and
57 generalized tools are available to explore the space of such web services (notably, the Biothings
58 API and Explorer; <https://biothings.io/>). However, the heterogeneity of such APIs can be a barri-
59 er to efficient biomedical knowledge integration.

60 Here we present a REST-based web services specification called the Knowledge Beacon
61 API (Beacon API) that enables a basic workflow for the discovery of, and navigation through,
62 biomedical concepts, relationships and associated evidence. This work arises out of an earlier
63 effort to develop a web application called “*Knowledge.Bio*” [7] to provide enhanced navigation
64 through the knowledge base of PubMed cited concepts and relationships, captured by text mining

65 in the Semantic Medline Database [8]. The knowledge harvesting workflow underlying
66 *Knowledge.Bio* is here elaborated into a distributed web service network across diverse
67 knowledge sources hosted within the NCATS Biomedical Data Translator Consortium, a public-
68 ly funded project supporting the FAIR integration of distributed biomedical research data and
69 knowledge to accelerate the development of new disease treatments and reduce the barriers be-
70 tween basic research and clinical advances [9]. The outcome of this work was an iteratively re-
71 fined web service specification implemented in an initial set of Beacons, with validation tools
72 and client applications.

73 **Methods**

74 The Knowledge Beacon API is a Swagger 2.0 specification that defines a set of endpoint
75 paths embodying operations for accessing knowledge sources and discovering shared semantics
76 for concepts and their relationships (Fig 1).

77 **Fig 1. Knowledge Beacon Workflow.** General step-by-step flowchart illustrat-
78 ing the sequential invocation of Beacon web service endpoints, with data flows
79 as indicated. Also enumerated at the bottom left hand corner of the diagram is
80 the set of metadata endpoints that report semantic terms and namespaces used by
81 the Beacon in the annotation of results.

82 A Knowledge Beacon (hereafter abbreviated “Beacon”) initiates a workflow for
83 knowledge discovery by simple search using a concept endpoint either with a *keywords*
84 parameter (`/concepts?keywords=`) or one with a Compact Uniform Resource Identifier (CURIE;
85 <https://en.wikipedia.org/wiki/CURIE>) of the concept (`/concepts/{conceptId}`). In both cases, one
86 or more specific concepts with associated core details are retrieved.

87 Once identified, the canonical CURIE identifier of a chosen concept selected from the re-
88 trieval list is used as an input parameter to access a list of statements about the concept, docu-
89 mented as subject/predicate/object assertions (`/statements?s=...` where `s` is a subject canonical
90 concept CURIE). Additional documentation, including supporting citations, associated with re-
91 turned statements may be examined by calling the statement's endpoint again with the statement
92 identifier of one of the entries returned from the initial call (i.e. `/statements/{statementId}`).

93 The data model, concept data type (“*category*”) and relationship predicate (“*edge_label*”,
94 “*relation*”) terms in results returned by a Beacon are compliant with an emerging public Biomed-
95 ical Data Translator Consortium semantic standard and data model, the Biolink Model
96 (<https://biolink.github.io/biolink-model/>). To assist client data parsing and interpretation, a
97 Beacon supports several additional endpoints that return metadata summaries of Biolink Model
98 terms specifically employed by the Beacon to annotate concepts and statements which are re-
99 turned: concept type “categories” (`/categories`), identifier name spaces (`/namespaces`), relation-
100 ship “predicates” (`/predicates`) plus a “knowledge map” of available subject-predicate-object
101 triplet statement combinations (`/kmap`).

102 **Results**

103 **Sample workflow**

104 Knowledge Beacon workflows are implemented as a chained series of REST API end-
105 point calls that return data as JSON formatted documents, annotated using Biolink Model stand-
106 ards as noted above. Here we illustrate a basic minimal two step sequence of such calls which
107 first identifies a list of concepts with names matching a keyword, then uses the identifier of one

108 returned concept entry to retrieve *subject-predicate-object* statement assertions relating to that
109 selected concept.

110 **Step 1:** Query knowledge sources by keyword to identify concepts. For example, calling
111 the basic /concepts endpoint using the Fanconi Anemia complementation group C gene
112 ‘FANCC’ as a keyword, on the Monarch “Biolink API” Beacon, namely:

113 <https://kba.ncats.io/beacon/biolink/concepts?keywords=FANCC>

114 returns the following JSON result with lists of CURIE-identified concepts (one entry shown; full
115 list of concept entries truncated for conciseness):

```
116                    [  
117                    ... some JSON results  
118                    {  
119                       "categories": [  
120                          "gene",  
121                          "sequence feature"  
122                       ],  
123                          "id": "NCBIGene:102158362",  
124                          "name": "FANCC"  
125                       },  
126                       ... more JSON results  
127                    ]
```

128 **Step 2:** Using the canonical (URL-encoded) concept CURIE of a selected concept in the
129 list of concepts returned by keyword in step 1 above, e.g. NCBIGene:102158362, a search for
130 knowledge assertions (statements) is made on the same database:

131 <https://kba.ncats.io/beacon/biolink/statements?s=NCBIGene%3A102158362>

132

133 This query gives another JSON result which contains asserted “subject-predicate-object” state-
134 ments about the concept, where the predicate return defines the relationship, as follows:

```
135 [
136 {
137     "id": "biolink:125a0182-0205-44a8-a70a-c03339383177",
138     "object": {
139         "categories": [
140             "gene"
141         ],
142         "id": "NCBIGene:102158362",
143         "name": "FANCC"
144     },
145     "predicate": {
146         "edge_label": "is_about",
147         "relation": "IAO:0000136"
148     },
149     "subject": {
150         "categories": [
151             "publication"
152         ],
153         "id": "PMID:17145712",
154         "name": "PMID:17145712"
155     }
156 },
157 ... more JSON results
158 ]
159
```

160 **Beacon implementations**

161 A stable set of publicly accessible Beacons are implemented and currently hosted stably
162 online (as of February 2020) by the NCATS Biomedical Translator Consortium, as enumerated
163 in Table 1. The Java and Python software implementations of these Beacons are available in re-
164 positories of the NCATS-Tangerine (<https://github.com/NCATS-Tangerine>) GitHub organiza-
165 tion. One implementation is a generic accessor of Biolink Model compliant knowledge graph
166 databases stored in Neo4j (<https://github.com/NCATS-Tangerine/tkg-beacon>). These Beacon
167 implementations may be tested using an available validator application
168 (<https://github.com/NCATS-Tangerine/beacon-validator>). A Python command line Beacon cli-
169 ent is available (<https://github.com/NCATS-Tangerine/tkbeacon-python-client>). A Knowledge
170 Beacon Aggregator (<https://github.com/NCATS-Tangerine/beacon-aggregator-client>) was also
171 designed to manage a registered pool of Beacons, and to return consolidated knowledge using
172 “equivalent concept cliques” to merge related Beacon results.

173 **Table 1. Biomedical Translator Consortium Deployed Beacons**

Subdomain ^a	Beacon Description	Wrapped Knowledge Source
semmeddb	Semantic Medline Database [8]	https://skr3.nlm.nih.gov/SemMedDB/
biolink	Monarch Database Biolink API [10]	https://api.monarchinitiative.org/api/
hmdb	Human Metabolome Database [11]	http://www.hmdb.ca/
rhea	Rhea Annotated Biochemical Reactions database [12]	https://www.rhea-db.org/
smpdb	Small Molecular Pathway Database [13]	http://smpdb.ca/
ndex	nDex Bio Graph Archive [14]	http://www.ndexbio.org

174 ^aThe basepath of each Beacon has the form <https://kba.ncats.io/beacon/<Subdomain>>, where
175 the <Subdomain> is as listed in column 1 of the table.

176 **Discussion**

177 The Knowledge Beacon API is a basic knowledge discovery workflow (Figure 1)
178 representing a relatively high-level use case of user interaction with the biomedical knowledge
179 space, and as such, lacks the full expressive power of a general knowledge query language inter-
180 face like SPARQL. Furthermore, the beacon data model aligns with the emerging Biolink Model
181 standards of the Biomedical Translator Consortium as its template for knowledge representation.
182 As such, Beacons do not automatically express results in a generic manner as do knowledge rep-
183 resentations such as RDF, although conversion of Beacon statement results into RDF format is
184 easily accomplished. Finally, aside from some general profiling of the performance of the
185 Swagger API endpoints on various knowledge sources, we have not here conducted a rigorous
186 computing-theoretic assessment of the efficiency of this form of knowledge harvesting, although

187 early experience with Beacons point to challenges with internet latency and knowledge-source
188 specific differences in query performance. In partial response to such challenges, we prototyped
189 a “Knowledge Beacon Aggregator” to provide enhanced asynchronous query/status/retrieval
190 endpoints as a client-friendly integration layer for managing access to, and merging data from, a
191 registered catalog of multiple Beacon implementations.

192 Despite the use of some off-the-shelf API generation tools, the wrapping of knowledge
193 sources as Beacons remains a labour-intensive activity. The semantics of the knowledge source
194 being wrapped must be heuristically translated. This is somewhat easier for knowledge sources
195 which have a small number of easily resolved discrete data types (i.e. discrete Biolink Model
196 concept categories of data) and namespaces with clear mapping onto those discrete data types.

197 In contrast, some “graph” knowledge sources, for example, the NDex Bio biomedical
198 network data archive (<https://home.ndexbio.org/index/>, wrapped by this project as the *ndex* Bea-
199 con), don’t have such clear concept category and relationship predicate tagging of much of the
200 archived data. The development of useful but (so far) imprecise heuristics to tag such data on the
201 fly is required to develop a useful Beacon. In other cases, such as biomedical knowledge re-
202 sources whose data object namespace aggregates several types of concepts in a fuzzy manner
203 with limited additional concept category tagging, it may be even more challenging to semantical-
204 ly tag data entries for beacon export.

205 A few common library and reference implementations are developed for Beacons; the
206 Beacon platform would benefit from the further development of standardized tools to systemati-
207 cally assist such wrapping of native knowledge sources.

208 The availability of a shared API standard for knowledge integration doesn't, in and of it-
209 self, deal with all the challenges of FAIR data integration within the global community. Practical
210 experience with knowledge harvesting using such API implementations has revealed perfor-
211 mance issues relating to internet and service latency, bandwidth limitations. Knowledge ware-
212 housing in centralized knowledge graphs using ETL (Extract, Transform, Load) processes may
213 sometimes result in a more tractable process for biomedical knowledge integration; however,
214 such approaches are still faced with the task of merging equivalent concepts, including the elimi-
215 nation of duplicate concepts and the resolution of conflicting information, including weighting of
216 assertions differing in levels of confidence. More unique to ETL warehousing approaches is the
217 ongoing problem of keeping such resources up-to-date relative to their original knowledge
218 sources. Note that ETL warehouses and API driven distributed knowledge harvesting approaches
219 can be complementary, in that ETL data warehouses can also themselves be accessed by the ap-
220 plication of web service REST API's like the Knowledge Beacon API. In fact, some of the cur-
221 rent Beacon implementations use this approach: a back end Biolink Model compliant Neo4j
222 knowledge graph directly wrapped with the API.

223 The Linked Open Data paradigm using RDF knowledge representation and SPARQL
224 represents an alternate paradigm for distributed knowledge integration, the theoretical perfor-
225 mance of which was surveyed by Verborgh *et al* [6]. In their assessment, it was noted that down-
226 loadable RDF knowledge data sets and SPARQL endpoints to triple store knowledge bases rep-
227 resent two extremes of a continuum of RDF knowledge access, each with their characteristic ad-
228 vantages and weaknesses. They proposed that a constrained query selector specification and
229 RDF representation – with data, metadata and hypermedia controls - denoted as Linked Data
230 Fragments could be a shared design representation spanning both ends of the continuum. Fur-

231 thermore, they proposed an intermediate implementation - termed Triple Pattern Fragments -
232 partitioning RDF processing more symmetrically across client and server, thus potentially miti-
233 gate some of the challenges of both ends of the API design continuum, for more balanced client-
234 server performance and greater ease of implementation (see
235 <https://github.com/LinkedDataFragments>).

236 Generally, API approaches to knowledge harvesting may work best with use cases in-
237 volving smaller batches of knowledge retrieval based on a focused navigation of the knowledge
238 space from larger open-ended data sources which would be refractory to import into centralized
239 knowledge graphs.

240 Finally, there are two other API standards of the Biomedical Data Translator Consortium:
241 the “NCATS Reasoner API” (<https://github.com/NCATS-Tangerine/NCATS-ReasonerStdAPI>)
242 and the Biothings API (<https://biothings.io/>). Although there are parallels between them, the
243 Beacon API is a simpler lower level interface to knowledge resources than the Reasoner API,
244 and is somewhat more constrained to the Biolink Model than the Biothings API. But the utility
245 of the Reasoner API is inspiring efforts to publish a Reasoner API interface on top of an imple-
246 mentation of Knowledge Beacons (<https://github.com/NCATS-Tangerine/kba-reasoner>). We did
247 also implement prototype Beacon wrapper for the Biothings API (<https://github.com/NCATS-Tangerine/biothings-explorer-beacon>).

249 Availability

250 Knowledge Beacon software is open source licensed and available for access in GitHub.
251 A suitable introduction to the API, containing references to related software components, can be
252 found at <https://github.com/NCATS-Tangerine/translator-knowledge-beacon>.

253 Acknowledgments

254 BMG and RMB collaborated on the predecessor “Knowledge.Bio” application embody-
255 ing the workflow captured by the Knowledge Beacon API. CJM and BMG coined the name
256 “Knowledge Beacon” to express the architectural vision of uniformly wrapped knowledge
257 sources for distributed knowledge discovery and harvesting. RMB and his team elaborated the
258 original software design of web service endpoints and initial code implementations embodying
259 Beacon functionality, then guided further iterations of the API based on feedback from col-
260 leagues within the Biomedical Translator Consortium, with special mention to co-author VD
261 who proposed insightful revisions to the API, based on his direct experience implementing a
262 Beacon to wrap HMDB. Under overall supervision by RMB, the heavy lifting of iterative soft-
263 ware development of Beacon implementations - including several beacons, client (including ag-
264 gregator) and validation applications - was undertaken by LMH while he was a member of the
265 STAR Informatics team, assisted by the valuable software programming contributions of several
266 computing science cooperative education students: MG, WTS and KCHB.

267 The authors would like to sincerely thank Nomi Harris and Marcin Joachimiak of LBNL
268 for their very helpful editorial feedback on, and suggested revisions to the draft manuscript.

269 The authors would also like to thank the various members of the Biomedical Data Trans-
270 lator Consortium who gave helpful user needs feedback and support of the Knowledge Beacon
271 API during its development, in particular, Chris Bizon and Stephen Ramsay. We also
272 acknowledge here Greg Stupp who, while employed at TSRI, implemented an earlier version of
273 a Beacon wrapper for biomedical knowledge in Wikidata.

274 **References**

275 1. Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, et al. The
276 FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* 2016
277 Mar 16;3:160018. doi:10.1038/sdata.2016.18. PubMed PMID: 26978244; PubMed Cen-
278 tral PMCID: PMC4792175

279 2. Wilkinson MD, Dumontier M, Jan Aalbersberg I, Appleton G, Axton M, Baak A, et al.
280 Addendum: The FAIR Guiding Principles for scientific data management and steward-
281 ship. *Sci Data* 2019 Mar 19;6(1):6. doi: 10.1038/s41597-019-0009-6. PubMed PMID:
282 30890711; PubMed Central PMCID: PMC6427092.

283 3. Wilkinson MD, McCarthy L, Vandervalk B, Withers D, Kawas E, Samadian S. SADI,
284 SHARE, and the *in silico* scientific method *BMC Bioinformatics*. 2010 Dec 21;11 Suppl
285 12:S7. doi: 10.1186/1471-2105-11-S12-S7. PubMed PMID: 21210986; PubMed Central
286 PMCID: PMC3040533.

287 4. BioMoby Consortium, Wilkinson MD, Senger M, Kawas E, Bruskiewich R, Gouzy J, et
288 al. Interoperability with Moby 1.0--it's better than sharing your toothbrush! *Brief
289 Bioinform.* 2008 Jan 31;9(3):220-31. doi:10.1093/bib/bbn003. PubMed PMID:
290 18238804.

291 5. Stevens R, Miller C. (2000) Wrapping and interoperating bioinformatics resources using
292 CORBA. *Brief Bioinform.* 2000 Feb;1(1):9-21. PubMed PMID: 11466976.

293 6. Verborgh R, Vander Sande M, Hartig O, Van Herwegen J, De Vocht L, De Meester B, et
294 al. Triple Pattern Fragments: a Low-cost Knowledge Graph Interface for the Web. Journal
295 of Web Semantics. 2016 Mar; 37:184-206. doi:10.1016/j.websem.2016.03.003

296 7. Bruskiewich RM, Huellas-Bruskiewicz KC, Ahmed F, Kaliyaperumal R, Thompson M,
297 Erik Schultes E, et al. Knowledge.Bio: A web application for exploring, building and
298 sharing webs of biomedical relationships mined from PubMed. 2016.
299 doi: <http://dx.doi.org/10.1101/055525>

300 8. Kilicoglu H, Shin D, Fiszman M, Rosemblat G, Rindflesch TC. SemMedDB: a PubMed-
301 scale repository of biomedical semantic predications. Bioinformatics. 2012 Dec
302 1;28(23):3158-60. doi: 10.1093/bioinformatics/bts591. PubMed PMID: 23044550; Pub-
303 Med Central PMCID: PMC3509487.

304 9. The Biomedical Data Translator Consortium. Toward A Universal Biomedical Data
305 Translator. Clin Transl Sci. 2019 Mar;12(2):86-90. doi: 10.1111/cts.12591. . PubMed
306 PMID: 30412337; PubMed Central PMCID: PMC6440568.

307 10. Shefchek KA, Harris NL, Gargano M, Matentzoglu N, Unni D, Brush M, et al. The Mon-
308 arch Initiative in 2019: an integrative data and analytic platform connecting phenotypes
309 to genotypes across species. Nucleic Acids Res. 2020 Jan 8;48(D1):D704-D715. doi:
310 10.1093/nar/gkz997. PubMed PMID: 31701156.

311 11. Wishart DS, Feunang YD, Marcu A, Guo AC, Liang K, Vázquez-Fresno R, et al. HMDB
312 4.0: the human metabolome database for 2018. Nucleic Acids Res. 2018 Jan
313 4;46(D1):D608-D617. doi: 10.1093/nar/gkx1089. PubMed PMID: 29140435; PubMed
314 Central PMCID: PMC5753273.

315 12. Morgat A, Lombardot T, Axelsen KB, Aimo L, Niknejad A, Hyka-Nouspikel N, et al.
316 Updates in Rhea - an expert curated resource of biochemical reactions. Nucleic Acids
317 Res. 2017 Apr 20;45(7):4279. doi: 10.1093/nar/gkw1299. PubMed PMID: 27980062;
318 PubMed Central PMCID: PMC5397180.

319 13. Jewison T, Su Y, Disfany FM, Liang Y, Knox C, Maciejewski A, et al. SMPDB 2.0: big
320 improvements to the Small Molecule Pathway Database. Nucleic Acids Res. 2014
321 Jan;42(Database issue):D478-84. doi: 10.1093/nar/gkt1067. Epub 2013 Nov 6. PubMed
322 PMID: 24203708; PubMed Central PMCID: PMC3965088.

323 14. Pratt D, Chen J, Pillich R, Rynkov V, Gary A, Demchak B, Ideker T. NDEx 2.0: A Clear-
324 inghouse for Research on Cancer Pathways. Cancer Res. 2017 Nov 1;77(21):e58-e61.
325 doi: 10.1158/0008-5472.CAN-17-0606. PubMed PMID: 29092941; PubMed Central
326 PMCID: PMC5679399.

