

# Search, Access, and Explore Life Science Nanopublications on the Web

Fabio Giachelle, Dennis Dosso, and Gianmaria Silvello

Department of Information Engineering, University of Padua, Italy

## ABSTRACT

Nanopublications are *Resource Description Framework (RDF)* graphs encoding scientific facts extracted from the literature and enriched with provenance and attribution information. There are millions of nanopublications currently available on the Web, especially in the life science domain.

Nanopublications are thought to facilitate the discovery, exploration, and re-use of scientific facts. Nevertheless, they are still not widely used by scientists outside specific circles; they are hard to find and rarely cited. We believe this is due to the lack of services to seek, find, and understand nanopublications' content.

To this end, we present the NanoWeb application to seamlessly search, access, explore, and re-use the nanopublications publicly available on the Web. For the time being, NanoWeb focuses on the life science domain where the vastest amount of nanopublications are available. It is a unified access point to the world of nanopublications enabling search over graph data, direct connections to evidence papers, and scientific curated databases, and visual and intuitive exploration of the relation network created by the encoded scientific facts.

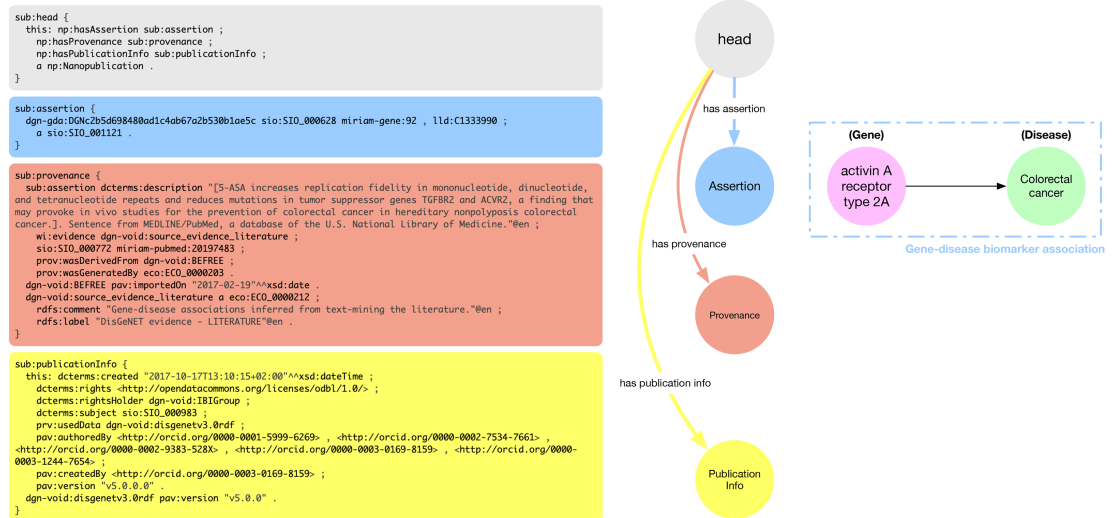
## 1 INTRODUCTION

The scientific world is swiftly becoming data-centric, embracing the principles of the so-called *fourth paradigm of science* (Hey et al., 2009). Data are at the center of scientific discovery as well as of scholarship and scholarly communication (Borgman, 2015). The growing role of data is also witnessed by the ever-increasing importance of data science and related research fields concerning the search (Chapman et al., 2020), provenance (Cheney et al., 2009), citation (Silvello, 2018), re-use (Wynholds et al., 2012), and exploration (Rahman et al., 2020) of data.

There is no “one size fits all” solution when it comes to data search, access, and re-use given the heterogeneity of data representations and models, interoperability issues, and domain-dependent requirements. In the context of scientific data, the *nanopublication model* has been proposed to target some of these issues (Groth et al., 2010). Nanopublications exploit the *Linked Open Data (LOD)* principles (Bizer et al., 2009) to represent scientific facts (*assertions* hereafter) as self-consistent, independent and machine-readable information tokens. A repository of nanopublications is to be thought of as an open and interconnected knowledge graph seamlessly integrated with the supporting scientific literature. Nanopublications can be used to support scientific claims, to explore scientific knowledge by exploiting machine intelligence and as entry points to scientific databases. Hence, this model has been embraced by several scientific fields, especially in the Life Science domain, leading to the creation of more than ten million openly available nanopublications (Kuhn et al., 2018).

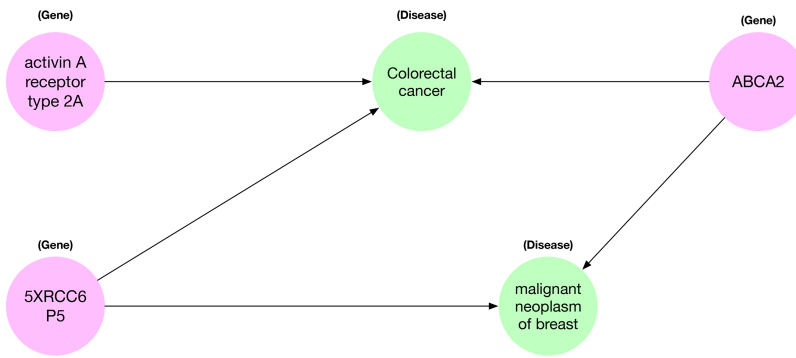
From the technical viewpoint, a nanopublication is a *Resource Description Framework (RDF)* graph built around an assertion represented as a triple (subject-predicate-object) and usually extracted, manually or automatically, from a scientific publication. The nanopublication enriches the assertion with provenance and publication information. The RDF representation format enables interoperability and thus the re-use of data, whereas provenance and publication information eases authorship recognition, credit distribution, and citation.

As an example taken from the biomedical domain, a nanopublication assertion about a gene-disease association is  $\langle \text{activin A receptor type 2A} - \text{gene-disease biomarker association} - \text{colorectal cancer} \rangle$ , where *activin A receptor type 2A* is the subject, *gene-disease biomarker association* is the predicate and *colorectal cancer* is the object of the triple. This assertion is extracted from a paper (Campregher et al., 2010), which puts in relation the *activin A receptor type 2A* gene to the *colorectal cancer* and describes a



**A** RDF serialization of a nanopublication

**B** Graphical view of the nanopublication



**C** A network of gene-disease associations

**Figure 1.** (A) RDF (trig) representation of the nanopublication encoding the assertion:  $\langle$ activin A receptor type 2A - gene-disease association - Colorectal Cancer $\rangle$ ; (B) graphical representation of the four parts of the nanopublications with a human-readable representation of the assertion graph; (C) network of gene-disease associations created by five nanopublications.

48 drug – i.e., *Mesalazine* – that reduces mutations in transforming growth factor of the gene.

49 In Figure 1.a, we can see a snippet of the RDF nanopublication serialization described above. Nanopublications are defined using the compact TriG<sup>1</sup> syntax, that enables to define *prefixes* to avoid to re-write the same IRIs multiple times. In Figure 1.a we used some prefixes within the nanopublication assertion, namely: *dgn-gda*, *sio*, *miriam-gene* and *lld*, that are specific of the life science domain. *dgn-gda* identifies a DisGeNET<sup>2</sup> gene-disease association; *sio* identifies a resource from *Semanticscience Integrated Ontology (SIO)*<sup>3</sup>, such as the type of a gene-disease association; *miriam-gene* identifies a gene in the National Center for Biotechnology Information (NCBI)<sup>4</sup> database; *lld* identifies a resource from the Linked Life Data<sup>5</sup> platform for the Biomedical domain.

57 The nanopublication is composed of four parts: (i) the *head* that acts as a connector between the other three sub-graphs; (ii) the assertion graph (blue) expressing the relationship between the two concepts of

<sup>1</sup><https://www.w3.org/TR/trig/>

<sup>2</sup><http://rdf.disgenet.org/>

<sup>3</sup><https://github.com/MaastrichtU-IDS/semanticscience>

<sup>4</sup><https://www.ncbi.nlm.nih.gov/>

<sup>5</sup><http://linkedlifedata.com/>

59 the assertion (the gene-disease association), the relationship of the concepts with external ontologies  
60 (the fact that *activin A receptor type 2A* is a gene and *colorectal cancer* is a disease), and possibly a  
61 link towards the scientific database storing related data; (iii) the provenance graph (orange) containing  
62 metadata about the assertion such as the methods used to generate the assertion and its creators; and,  
63 (iv) the publication info graph (yellow) containing the metadata about the evidence paper from which  
64 the assertion was extracted and about the nanopublication itself. In Figure 1.b, we can see a graphical  
65 representation of the four parts of the nanopublications with a human-readable representation of the  
66 gene-disease association encoded by the assertion graph.

67 A key aspect motivating the use of nanopublications is the possibility to exploit LOD features,  
68 allowing for exploring relation networks created by connecting related facts encoded in RDF. Indeed,  
69 nanopublications create a network of scientific assertions that can be explored to discover connections  
70 between facts. In the literature, there is important evidence of using nanopublications as a credible  
71 approach for expanding scientific insight, especially in the biomedical domain (Chichester et al., 2014).  
72 As a motivating example, Figure 1.c shows a small network of gene-disease associations. We can see that  
73 the genes *activin A receptor type 2A* and *5XRCC6P5* are both related to *colorectal cancer*. If we search  
74 for other connections, we find another nanopublication relating the *5XRCC6P5* gene to the *malignant*  
75 *neoplasm of breast* disease. Further expanding the relation network, we see that there exist two other  
76 nanopublications connecting the *ABCA2* gene with both *colorectal cancer* and *Malignant neoplasm of*  
77 *breast*. Figure 1.c presents a small network that shows the relationships between facts extracted from five  
78 different papers published in different venues at different times that do not cite each other. This is just a  
79 hint about how exploring the nanopublication relation network could lead to finding related concepts and  
80 assertions that might not be explicitly connected in the scientific literature and databases.

81 Nonetheless, despite these premises, nanopublications are not widely used by scientists outside  
82 specific circles (Page, 2018); they are hard to find and rarely cited. Nanopublications rarely have a  
83 human-readable accessible version and cannot be searched via keywords or natural language queries.  
84 Although nanopublications are based on LOD principles, there are still no tools that allow the user to  
85 explore their connections intuitively and discover if and how one assertion is related to others, as we have  
86 done in the example above. Leveraging on the famous *data is the new oil* metaphor (Economist, 2017),  
87 we can say that with nanopublications we have a vast oil reservoir but no active refinery, distribution net,  
88 and machines to put it into use.

89 In this work, we target these issues and present the *NanoWeb* application<sup>6</sup>, an open-source and publicly  
90 available web service enabling intuitive search, exploration, and re-use of nanopublications. The current  
91 version of *NanoWeb* is tailored for the life science domain, and it is designed to help experts of this  
92 domain in their research work. *NanoWeb* is an extensible tool to be applied to other scientific domains,  
93 even though certain customization to do so will be required. *NanoWeb* is a single entry point to the world  
94 of nanopublications enabling the seamless integration of data search, exploration, and re-use services; its  
95 central features are:

- 96 1. a crawler gathering publicly available nanopublications from the web;
- 97 2. two intuitive search functionalities, based respectively on the keyword search and boolean search  
98 paradigms;
- 99 3. a user-oriented visual interface to consult the nanopublications enriched with information gathered  
100 from external authoritative ontologies;
- 101 4. a service enabling the graph-based visualization of assertions and the exploration of their relation  
102 network;
- 103 5. data search functionalities providing entry points to external curated databases storing the scientific  
104 facts encoded by the nanopublications as well as to the scientific papers where the assertions were  
105 extracted.

106 The rest of the paper is organized as follows: Section 2 presents the background of the nanopublication  
107 model and the state of the art of systems based on it. Section 3 describes the overall architecture of the  
108 *NanoWeb* application. Section 4 reports the statistics about the nanopublications available in *NanoWeb*.

---

<sup>6</sup><https://w3id.org/nanoweb/>

109 Section 5 shows how NanoWeb works and details the functioning of the user interface. Section 6 reports  
110 the results of the expert users survey conducted on NanoWeb. Section 7 discusses the challenges to be  
111 faced with maintaining NanoWeb in the medium-long period and how it can scale up to be used in domain  
112 others than life science. Finally, Section 8 draws some final remarks and outlines future work.

## 113 2 BACKGROUND

114 **Basics of Nanopublications.** Nanopublications rely on Semantic Web technology. In particular, they  
115 are modeled via RDF (Groth et al., 2010), a widely used standard endorsed by the W3C consortium<sup>7</sup>,  
116 adopted for data publishing, accessing and sharing. RDF allows for the manipulation, enrichment, discov-  
117 ery and interoperability of data and it is at the core of the implementation of the LOD paradigm (Pound  
118 et al., 2010).

119 RDF is based on the concept of *statement*, that presents a `<subject, predicate, object>`  
120 triple-based structure. Within a triple, `subject`, `predicate` and `object` are *resources*. In particular,  
121 an RDF dataset can be represented as a graph where, given a triple, the `subject` and the `object` are  
122 the nodes representing *resources*, while the `predicate`, the direct edge connecting the two, expresses  
123 their *relationship*.

124 RDF resources can either be *IRIs* (Internationalized Resource Identifiers), *literals* or *blank nodes*. An  
125 IRI<sup>8</sup> is a more general form of URI which can also contain Unicode characters. A literal is a value which  
126 can be associated to a specific type of value, such as string, integer, date, time etc. The default value is  
127 string. Blank nodes are resources which are labeled with a URI-like string which has validity only inside  
128 the database.

129 In RDF every resource and relationship is labeled. `subject` and `object` nodes can be labeled with  
130 IRIs, `object` nodes can also be labeled with literals. Relationships can only be labeled with IRIs. Blank  
131 nodes can be `subject` or `object` of a triple. A set of RDF triples can also be thought as a directed  
132 graph, where subjects and objects are nodes and predicates are the directed edges. Hence, it is also called  
133 RDF graph.

134 In recent years it has been proposed the idea to extend the basic semantic of RDF by using *quads*  
135 instead of triples, where an identifier (an IRI) is added. In this way, groups of triples may be characterized  
136 as belonging to the same subgraph, i.e. to the same *named graph* (Carroll and Stickler, 2004; Carroll  
137 et al., 2005), if they share the same extra URI.

138 Every nanopublication is made of four basic named graphs as shown in Figure 1.a:

- 139 1. *Head*: the graph composed of four triples connecting assertion, provenance and publication info  
140 graphs together and specifying that the graph at hand is a nanopublication.
- 141 2. *Assertion*: the assertion is to be thought of as the minimal unit of thought, a fact or a statement. It  
142 can be composed of one or more RDF triples and for this reason, we often call it *assertion graph*.
- 143 3. *Provenance*: the named graph made of metadata providing *context* about the assertion. The  
144 information contained in the provenance describes how the information expressed in the assertion  
145 was created (from some experiment, extrapolate from a paper or article, etc.) and the methods  
146 that were used to generate the assertion. It includes information such as authors, institutions,  
147 time-stamps, grants, links to evidence papers and other resources.
- 148 4. *Publication information*: the graph containing the information about the nanopublication itself,  
149 such as its authors, the topic of the assertion, and rights information.

150 **Nanopublication resources and datasets.** The website <http://nanopub.org/> is the most com-  
151 prehensive access point to the world of nanopublications. It collects papers and tools about nanopublica-  
152 tions. The central resource to access millions of publicly available nanopublications is the “nanomonitor”  
153 <sup>9</sup>. It provides a list of sixteen worldwide distributed servers where nanopublications can be openly ac-  
154 cessed and downloaded in several formats. The nanopublications are ordered by identifier, but no full-text  
155 or structured search service is available. The nanopublications are accessible in an RDF serialization  
156 format. Thus they are machine-readable but not human-readable (see Figure 1.a).

<sup>7</sup><https://www.w3.org/TR/rdf11-primer>

<sup>8</sup><https://tools.ietf.org/html/rfc3987>

<sup>9</sup><http://app.tkuhn.eculture.labs.vu.nl/nanopub-monitor/>

157 Kuhn et al. (2013) describes a Web-based service (i.e., *nanobrowser*) enabling access to human-  
158 readable enriched scientific statements extracted from nanopublications. The aim of *nanobrowser* is  
159 to enable easy publishing and curation of nanopublications, but unfortunately, at the time of writing,  
160 it does not work, even though the source code is publicly available.<sup>10</sup> The nanobrowser had the goal  
161 to ease the extraction of facts from scientific papers and to enable the community to curate and revise  
162 the statements; its overall objective is different from those of NanoWeb even though they share the  
163 requirement of making nanopublications human-readable and facilitate access to them. In the same  
164 direction, the *whyis* project<sup>11</sup> proposes a knowledge graph infrastructure to support domain-aware  
165 management and curation of knowledge from different sources; it leverages on the nanopublication model  
166 to represent the facts and handle their provenance in the knowledge base. *whyis* also offers some facilities  
167 to allow the users to visually explore the knowledge graph beyond a given entity by using the so-called  
168 knowledge explorer (McCusker et al., 2017, 2018); the knowledge explorer shares some similarities with  
169 the NanoWeb exploration tool. In particular, they both allow the exploration of the connections between  
170 entities in the knowledge graph. Nevertheless, *whyis* does not visualize the scientific assertions encoded by  
171 nanopublications. More specifically, the *whyis* project is oriented to the creation and user-based curation  
172 of the nanopublications rather than to the search and exploration possibilities connected to them. Hence,  
173 NanoWeb is a complementary service rather than a competitor to *whyis*.

174 Mons et al. (2011) advocated for the systematic use of nanopublications to encode scientific facts  
175 reported in published papers. They see nanopublications as the key tool to enable reasoning and fact  
176 discovery exploiting machine intelligence. Furthermore, they extracted thousands of nanopublications  
177 about valuable and hard to discover gene variations and made them publicly available. We enable the  
178 search and access to these nanopublications in NanoWeb.

179 Chichester et al. (2015) described how they created nanopublications encoding scientific facts asso-  
180 ciated with more than 38K proteins stored in the neXtProt database.<sup>12</sup> The main motivation for this  
181 work is to exploit nanopublications potential to support end-user research on human proteins enabling  
182 machine-reasoning, easy search and access to the protein-related facts. Chichester et al. (2014) showed  
183 how nanopublications as fine-grained annotations answer to complex knowledge discovery queries other-  
184 wise challenging to deal with. Also, in this case, queries are performed using the SPARQL structured  
185 language confining the use of nanopublications to technical database experts. We crawled and enable  
186 keyword-based search over all the publicly available neXtProt nanopublications.

187 Queralt-Rosinach et al. (2016) described the process that led to the publication of millions of nanop-  
188 ublications about the pathophysiology of diseases extracted from the scientific literature and backed by  
189 curated records in the DisGeNET database.<sup>13</sup> The DisGeNET nanopublications are publicly available  
190 and accessible via a SPARQL endpoint. NanoWeb collected, indexed all the available DisGeNET nanop-  
191 ublications and made them searchable and human-readable. Each nanopublication is enriched with a URL  
192 linking to the related curated record in DisGeNET.

193 Wikipathways is an online collaborative pathway resource that is made available as RDF and nanop-  
194 ublications (Waagmeester et al., 2016). The nanopublications are backed by the Wikipathways curated  
195 database and are accessible via a SPARQL endpoint (not available at the time of writing). The resource to  
196 convert the RDF triples of Wikipathways to nanopublication is publicly available.<sup>14</sup> We crawled all the  
197 Wikipathways nanopublications, that are now searchable and accessible via NanoWeb.

198 Hettne et al. (2016) extracted more than 200M assertions about gene-disease associations from  
199 the biomedical literature. 7M assertions are explicitly stated in the scientific papers and the rest is  
200 implicitly inferred. There is a publicly available dump<sup>15</sup> of the nanopublications shared as additional  
201 data for the paper. The website <https://rdf.biosemantics.org/> is intended to share all the  
202 nanopublications and to make access to the ontology required to dereference the concepts encoding the  
203 assertions. Unfortunately, at the time of writing, the nanopublications as well as the SPARQL endpoints  
204 to access them are unavailable.

205 Amith and Tao (2018) defined an ontology – VAXMO – for encoding vaccines-related information  
206 extracted from scientific literature and used nanopublications to propose a method to store misconceptions

<sup>10</sup><https://github.com/tkuhn/nanobrowser>

<sup>11</sup><http://tetherless-world.github.io/whyis/>

<sup>12</sup><https://www.nextprot.org/>

<sup>13</sup><http://rdf.disgenet.org/>

<sup>14</sup><https://github.com/wikipathways/nanopublications>

<sup>15</sup><https://datadryad.org/stash/dataset/doi:10.5061/dryad.gn219>

207 about vaccines. Unfortunately, the VAXMO ontology is not accessible as well as the associated nanopub-  
208 lications. Also, Zhang et al. (2019) recently used the nanopublication model to represent scientific facts  
209 manually extracted from the literature about cancer behavioral risk factors. They presented a prototype  
210 – AERO – to search and visualize the nanopublications; search is based on SPARQL queries and the  
211 visualization is allowed only for the results returned by the SPARQL endpoint. At the time of writing,  
212 AERO is not publicly available.

213 To the best of our knowledge, there is no available tool to visualize nanopublications and explore their  
214 connections. The tool which is closer to NanoWeb in terms of semantic search and graph visualization is  
215 BioKB (Biryukov et al., 2018). BioKB provides access to the semantic content of biomedical articles  
216 through a SPARQL endpoint and a web interface; its goal is to allow the users to search for biomedical  
217 entities and visualize their graph of relations. However, BioKB does not account for nanopublications and  
218 does not support a multi-level exploration of the graph, enabling an in-depth exploration of the entities  
219 relation network.

220 Overall, the current services for searching nanopublications are all based on sparse SPARQL end-  
221 points. To this end, NanoWeb contributes on two levels. First, it provides a unique online access point  
222 to all the publicly available nanopublications from the Life Science domain; and, second, NanoWeb  
223 provides advanced services as keyword search, visualization and human-readable access to millions of  
224 nanopublications, making them accessible to users without technical expertise in SPARQL and related  
225 technologies.

226 **Search over RDF.** RDF graphs can be interrogated through the powerful but complex SPARQL query  
227 language (Pérez et al., 2009). SPARQL is not intuitive for end-users since it presents a complex syntax,  
228 far from a natural expression of their information need (Wu, 2013). It also requires knowledge of the  
229 underlying schema of the database, and of the IRIs used in it. This knowledge is often not possessed by  
230 the average end-user.

231 A search paradigm adopted to address the issues related to the use of SPARQL is *keyword search*.  
232 Keyword-based methods have gained importance over time both in research and in industry as a paradigm  
233 to facilitate the access to structured data (Bast et al., 2016; Kopliku et al., 2014; Yu et al., 2010).

234 The main difference between SPARQL and keyword search is that, while SPARQL returns the one  
235 and only correct answer (or an empty set if there was no answer), keyword search returns a ranking of  
236 answers, ordered based on their *relevance* to the information need expressed by the user via the keyword  
237 query.

238 In the literature, keyword query search systems over structured data are mainly focused on relational  
239 databases (RDB) (Yu et al., 2010) but many are also emerging for graph-like databases such as RDF  
240 datasets (Wang and Aggarwal, 2010; Bast et al., 2016). These systems may be divided into *three*  
241 categories.

242 The first kind of systems is *schema-based*. Examples are (Balmin et al., 2003; Agrawal et al., 2002;  
243 Luo et al., 2011). These systems exploit the schema information of the database, be it relational or RDF,  
244 to formulate queries in a structured language (SQL or SPARQL depending on the type of the database)  
245 designed from the keyword query of the user.

246 The second category is *graph-based*. Originally born with relational databases (Bhalotia et al., 2002;  
247 Simitis et al., 2008), the technique at the base of these systems was based on the transformation of the  
248 relational database in a graph. These systems are relatively easily translated in the RDF scenario since  
249 these databases are already in a graph form. A core challenge of these systems is to deal with the size of  
250 big graphs, which can contain tens of millions of nodes, if not more. In several cases, it has been shown  
251 that the size makes the task unsolvable by these systems (Coffman and Weaver, 2014).

252 Stemming from this last class of systems, the last category is the one of the *virtual-document based*  
253 systems Kadilierakis et al. (2020). First described in (Lopez-Veyna et al., 2012), this approach relies  
254 on the concept of *virtual document* of a graph. Given one graph, RDF or obtained by relational tuples,  
255 its corresponding virtual document is obtained by extracting words from it in an automatic way. This  
256 produces a “flat” representation of the graph, where its syntax and topology are lost but its semantic and  
257 lexical content is somewhat maintained. The virtual document representation is convenient since systems  
258 can leverage on efficient state-of-the-art IR methods for indexing and ranking. These methods operate  
259 by first extracting subgraphs from the whole database, then converting them in their virtual document  
260 representation and ranking these documents with respect to the keyword query. The user receives at the  
261 end the ranking of graphs in the order dictated by the ranking on the corresponding documents.



- 287 to interact and control crawler activities using a Graphic User Interface (GUI).<sup>16</sup> The batch  
288 mode enables a fast and batch-based download using operating systems lacking a GUI.
- 289 – **Metadata builder** (3): the nanopublications are processed to dereference the URLs and to  
290 get additional *metadata*; for instance, the nanopublications are enriched with the label of the  
291 concepts referring to external ontologies, the names of creators and curators and the title of  
292 the evidence papers. These data are saved in a relational database (4).<sup>17</sup>
  - 293 – **Document builder** (5): The document creation phase occurs after the dereferencing and  
294 enrichment phase. The document builder creates “virtual” nanopublication documents, which  
295 are saved into a database (6), on which the keyword search system is based.
- 296 • **Search system** (Figure 2, box B): this system performs keyword search on the nanopublications  
297 and it has three components:
    - 298 – **Business logic** (11): it is the controller unit of the search system. It performs the orchestration  
299 activities such as the coordination of the crawler by feeding it with new nanopublication  
300 URLs. It takes the user keyword query as input and returns the relevant nanopublications  
301 through the Web interface as output. To perform this task, the business logic unit relies on  
302 three databases: the nanopublication documents database (6), the fields (7) and the indexes  
303 (8). The indexes database contains the inverted index extracted from the nanopublication  
304 documents required to match the query terms with the document terms. The fields database is  
305 required to provide fast access to specific nanopublication data such as the authors, curators,  
306 and evidence paper metadata.
    - 307 – **Web interface** (12): it is the front-end allowing the user to search, access, explore and cite  
308 nanopublications through an interactive interface. It communicates with the business logic  
309 unit using a REST layer that provides public API for accessing nanopublications data in  
310 JSON format.
    - 311 – **Log system** (13): it deals with the logging tasks of the search system and it relies on a  
312 specific relational database (9). It communicates with the Web interface to collect relevant  
313 user activity information and possible problems.
  - 314 • **Citation system** (10): it generates the citations text snippet for the nanopublications of interest to  
315 the user by relying on the system presented by Fabris et al. (2019). Citations are a fundamental  
316 tool to give credit to authors and curators of data and publications and help other users to recognize  
317 the value of nanopublications. When the business logic unit (11) receives the request to produce a  
318 citation for a nanopublication, it sends this request to the citation system, that in turn collects the  
319 necessary metadata from the corresponding database (4). Once produced, the citation snippet is  
320 returned to the business logic unit and then visualized in the Web interface.

### 321 3.1 Search system

322 Let us assume that a user has an information need, and wants to retrieve the nanopublications that satisfy  
323 it. Since nanopublications are encoded in RDF, one possibility is to query the graph composed by all the  
324 nanopublications via the SPARQL query language, that, as already discusses, presents drawbacks for  
325 non-expert users.

326 We adopt two alternatives to SPARQL, i.e. keyword search and boolean search, both oriented to ease  
327 the search process for the users. Boolean search (i.e., advanced search) is adopted for domain-specific  
328 searches and it is useful to guide users in query formulation, since they often do not know in advance what  
329 they can search. We realized advanced search over the nanopublication metadata database, that allows for  
330 searching on specific fields of the indexed data (e.g. genes, diseases, proteins or authors).

331 Boolean search enables targeted search functionalities, but it does not allow for general and open  
332 full-text search over the nanopublications. To allow users to exploit natural language to search for

<sup>16</sup>A demonstration video of the crawler in action, using the graphic mode, is available at <https://bit.ly/2RV1Gz1>.

<sup>17</sup>All the relational databases are based on PostgreSQL version 10.6 allowing for the table partitioning function; this function enables efficient storage and access to the data.



333 nanopublications, we realized a keyword search system over RDF data. The system we adopt is based on  
334 the *virtual document* strategy, first presented in (Lopez-Veyna et al., 2012) and used in many other papers  
335 about *keyword search* on RDF graphs (Dosso and Silvello, 2020; Elbassuoni and Blanco, 2011; Mass and  
336 Sagiv, 2016). The underlying task of these papers is that, given an RDF graph, the user wants to query it,  
337 but for some reason, she is unable to use a SPARQL query. Keyword search is an alternative paradigm to  
338 using a structured query based on a query made of keywords.

339 The virtual document strategy is one of the many strategies deployed to face keyword search on RDF  
340 graphs. Given an RDF graph, we call its corresponding *virtual document* the textual document obtained  
341 from the concatenation of words obtained from the IRIs and Literals contained in the nodes and edges of  
342 the graph.

343 Given a collection of graphs it is therefore possible to create a corresponding collection of *virtual*  
344 *documents*. Every document is uniquely linked to the graph that generated it since they share the same  
345 identifier.

346 Then, the collection of documents is indexed and, from that moment on, this index can be used to  
347 answer keyword queries in the same way in which it is done in more classic IR scenarios, where the collec-  
348 tions are made by “real” documents. In this paper we used a probabilistic model (i.e., BM25 (Robertson  
349 et al., 1994)) as ranking function.

350 Every time a new query is issued, BM25 uses the virtual document index to create a ranking of docu-  
351 ments. The document identifiers are used to retrieve the corresponding graphs, that is, the corresponding  
352 nanopublications, from the collection. This list of nanopublications is then returned to the final user in the  
353 same order dictated by the ranking.

354 One may argue that this strategy discards information from the graphs. Since each graph is *flattened*  
355 to a document version of itself, information such as its topology and the disposition of words among  
356 nodes and edges is lost. This is certainly true, and in fact works such as (Dosso and Silvello, 2020;  
357 Elbassuoni and Blanco, 2011; Mass and Sagiv, 2016) do not limit themselves to virtual documents, but  
358 employ different kinds of heuristics to better leverage on the topology of the graphs.

359 Moreover, topology oriented heuristics often rely on the exploration of the graphs, which adds  
360 overhead to the whole computation. The more the answers returned by BM25, the bigger this overhead.  
361 Therefore, we argue that the use of topology-oriented heuristics do not guarantee a significant improvement  
362 on the effectiveness of the rankings obtained by the graphs with respect to the added overhead to the  
363 computation.

## 364 4 NANOPUBLICATION COLLECTION STATISTICS

365 In Table 1 we report the number of nanopublications per scientific platform currently available in NanoWeb.  
366 Currently, we have crawled and indexed nanopublications from the following platforms:

- 367 • **DisGeNET**:<sup>18</sup> “a discovery platform containing one of the largest publicly available collections of  
368 genes and variants associated to human diseases” (Piñero et al., 2019). DisGeNET is a knowledge  
369 management platform integrating and standardizing data about disease-associated genes and variants  
370 from multiple sources, including the scientific literature. DisGeNET covers the full spectrum of  
371 human diseases as well as normal and abnormal traits. Queralt-Rosinach et al. (2016) presented  
372 the publication of DisGeNET human Gene-Disease Associations (GDAs) as a new Linked Dataset  
373 exploiting the nanopublication approach. DisGeNET provides roughly half of the nanopublications,  
374 about 5 million, available in NanoWeb.
- 375 • **NeXtProt**:<sup>19</sup> “neXtProt is a protein knowledge platform that aims to support end-user research on  
376 human proteins” (Chichester et al., 2015). Chichester et al. (2015) converted data from neXtProt into  
377 nanopublications to show how they can be used to seamlessly query the data and gain biological  
378 insight. In particular, they converted three types of annotations of interest for the biomedical  
379 community: variation data, posttranslational modification (PTM), and tissue expression.
- 380 • **Protein Atlas**:<sup>20</sup> “A Human Pathology Atlas has been created as part of the Human Protein Atlas  
381 program to explore the prognostic role of each protein-coding gene in each cancer type by means

---

<sup>18</sup><https://www.disgenet.org/>

<sup>19</sup><https://www.nextprot.org/>

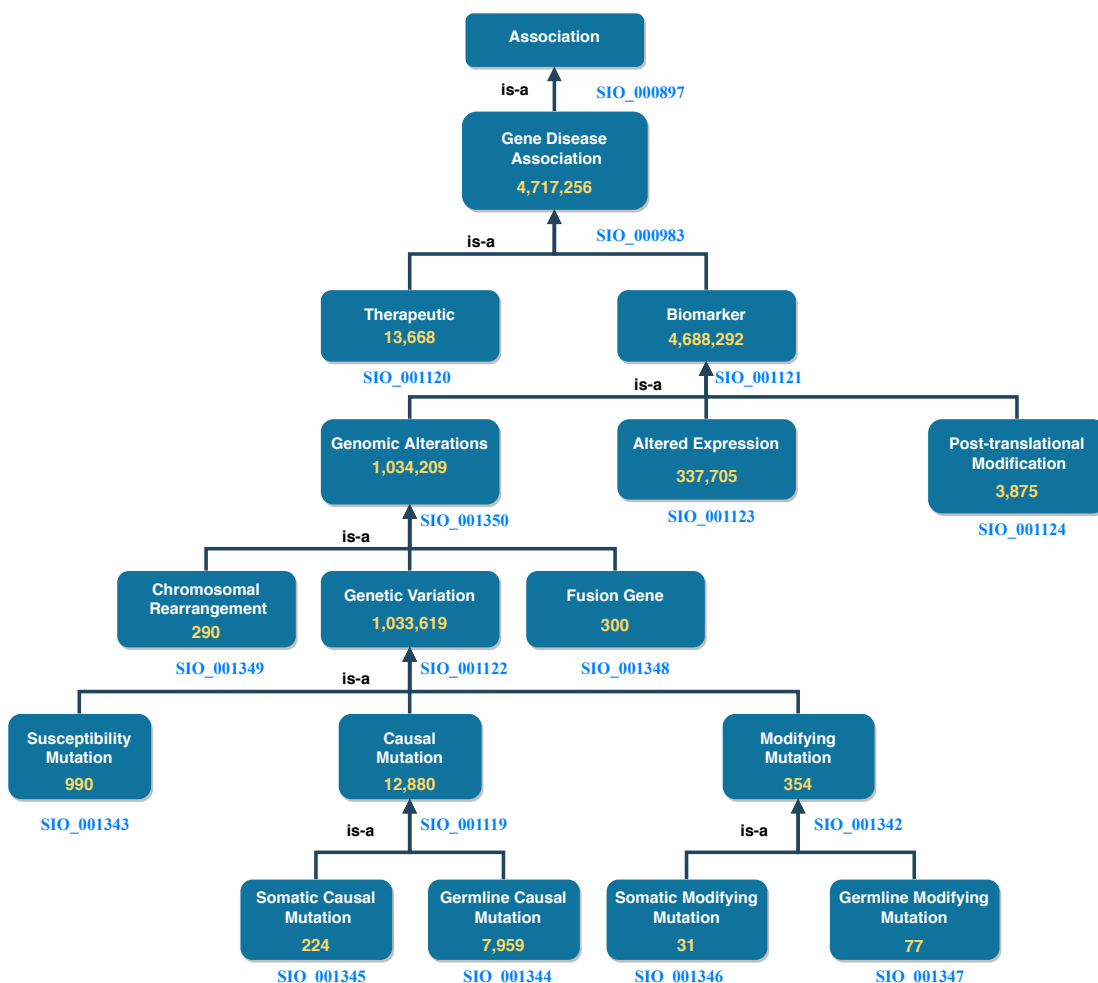
<sup>20</sup><https://www.proteinatlas.org/>

382 of transcriptomics and antibody-based profiling.” (Uhlen et al., 2017). The Human Protein Atlas is an  
 383 open-access knowledge-base providing the data to allow genome-wide exploration of the impact  
 384 of individual proteins on clinical outcomes. The Human Protein Atlas (HPA) programme aims to  
 385 “generate a comprehensive atlas of protein expression patterns in human normal and cancer tissues  
 386 as well as cell lines.” (Pontén et al., 2008).

- 387 • **WikiPathways:** <sup>21</sup> “WikiPathways is an open, collaborative platform dedicated to the curation of  
 388 biological pathways.” (Slenter et al., 2017; Waagmeester et al., 2016). WikiPathways provides rich  
 389 pathway databases with a focus on genes, proteins and metabolites. The data from WikiPathways  
 390 have been converted into a dataset of nanopublications as explained in (Kuhn et al., 2017).

Platform	Number of nanopublication
DisGeNET	4,717,256
NeXtProt	4,014,376
Protein Atlas	1,254,466
Wikipathways	26,934
<i>Total number of nanopublications</i>	<i>10,013,032</i>

**Table 1.** Number of nanopublications per platform.



**Figure 3.** DisGeNET ontology: number of assertions (yellow) for each DisGeNET association type.

<sup>21</sup><https://www.wikipathways.org/>

391 **4.1 Association analysis**

392 DisGeNET accounts for roughly half the total number of nanopublications in NanoWeb. The assertions  
 393 encoded by these nanopublications are divided into gene-disease associations of different types. In Figure  
 394 3, we report the number of assertions in NanoWeb for each association of the DisGeNET ontology. A  
 395 detailed description of the associations is available in the DisGeNET website.<sup>22</sup>

396 In the same vein, Table 2 reports the genes-tissues association types present in NeXtProt nanopublica-  
 397 tions. In particular, the *protein-coding gene expression in tissue* association describes the relationship  
 398 between a protein-coding gene in directing the production of proteins expressed in a tissue. Another type  
 399 of association regarding proteins is the *protein expression in tissue* which describes the expression level  
 400 (high, low, medium, not detected) of a protein in a tissue. Besides, the *sequence on amino-acid* associa-  
 401 tion describes the relationship between proteins and amino acids. The total number of nanopublication  
 402 assertions regarding protein associations is over 5 million.

Association	Number of assertion
protein-coding gene expression in tissue (generic)	6
protein-coding gene expression in tissue with quality high	124,261
protein-coding gene expression in tissue with quality low	184,615
protein-coding gene expression in tissue with quality medium	275,241
protein-coding gene expression in tissue with quality negative	837,144
protein-coding gene expression in tissue with quality not detected	341,062
protein-coding gene expression in tissue with quality positive	1,421,203
<i>protein-coding gene expression in tissue (total)</i>	<i>3,183,532</i>
protein expression in tissue with level high	150,366
protein expression in tissue with level low	241,325
protein expression in tissue with level medium	361,641
protein expression in tissue with level not detected	501,133
<i>protein expression in tissue (total)</i>	<i>1,254,466</i>
sequence on amino-acid	739,528
<i>protein associations (total)</i>	<i>5,177,526</i>

**Table 2.** Assertion numbers for association types: “protein-coding gene expression in tissue” and “protein expression in tissue”.

Database	Number of evidences
Bgee	5,576,047
Cancer Sanger	578
EbiQuickGo	8876
Gene Expression Omnibus (GEO)	573,648
Protein Atlas	4,125,154
UniProt	628,749
<i>Total number of evidences</i>	<i>10,913,052</i>

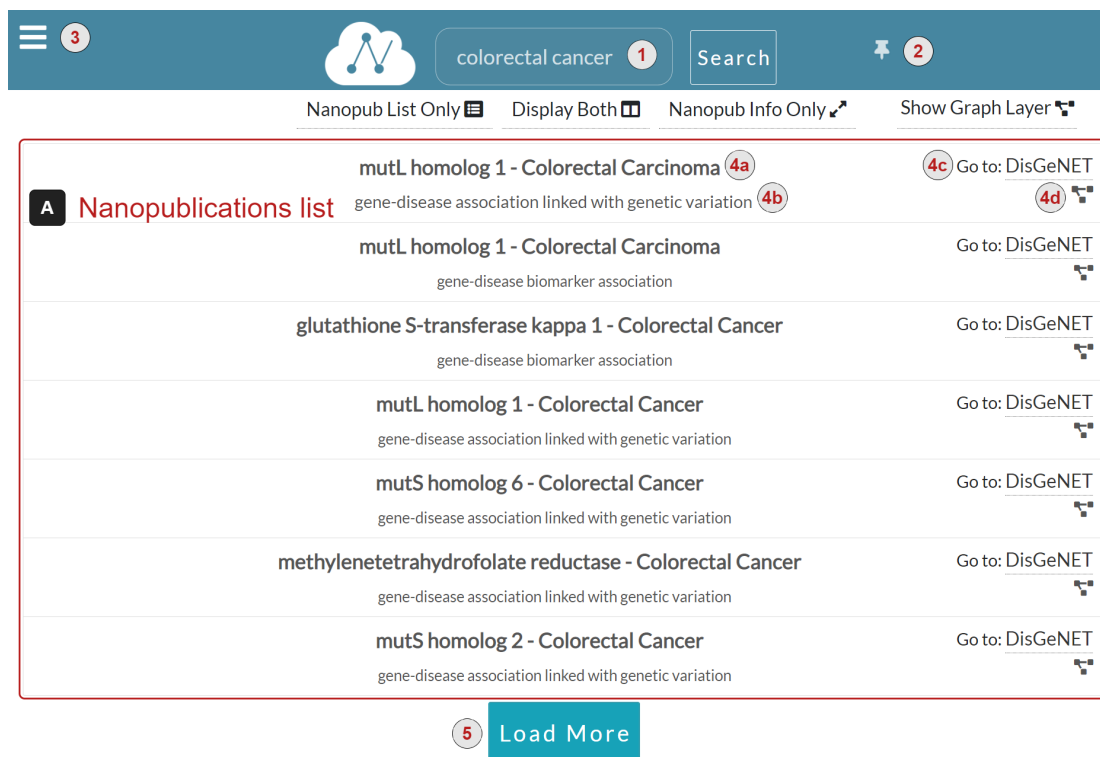
**Table 3.** Number of evidences per database.

403 **4.2 Scientific Evidences**

404 Nanopublication assertions are supported by evidences; an evidence can be a scientific publication,  
 405 a curated database record or both. The nanopublication evidences in NanoWeb come from several

<sup>22</sup><https://www.disgenet.org/dbinfo#section5>

406 institutional open-access databases such as Bgee<sup>23</sup>, Cancer Sanger<sup>24</sup>, EbiQuickGo<sup>25</sup>, Gene Expression  
 407 Omnibus (GEO)<sup>26</sup>, Protein Atlas<sup>20</sup> and UniProt<sup>27</sup>. We report the evidence databases associated to  
 408 the nanopublications available in NanoWeb in Table 3. The total number of evidences collected from  
 409 authoritative databases are about 11 million, and the evidences coming from publications are more than 6  
 410 million. All these publications are available in the PubMed<sup>28</sup> database.



**Figure 4.** NanoWeb search interface with user-provided query: *colorectal cancer*

## 5 NANOWEB GRAPHICAL USER INTERFACE

411 The NanoWeb system, available at <http://w3id.org/nanoweb/>, provides an interactive Web  
 412 interface that the user can use to search, access, explore, and cite nanopublications. A demo video  
 413 presenting NanoWeb functionalities is available at <https://bit.ly/NWURL2>.

414 Figure 4 shows the NanoWeb search interface. At the top of the page, there is the query input form (1),  
 415 where the user types the query and searches for nanopublications. There is a button (2) to pin or unpin the  
 416 query input form on the right side of the query input form. The query input form is unpinned by default;  
 417 this means that it floats at the top of the page so that it is always visible to the user even when the page  
 418 is scrolled. The user can press the button to pin the query input form, making it hidden when the page  
 419 is scrolled. On the left side of the query input form, there is the menu button (3). By clicking on it, the  
 420 sidebar appears with a list of links to the web app functionalities:  
 421

- 422 1. **Home:** takes the user to the home page.
- 423 2. **Stats:** takes the user to the Web page summarizing the NanoWeb system statistics, such as the  
 424 number of nanopublications and triples inserted in the database.

<sup>23</sup><https://bgee.org/>

<sup>24</sup><https://cancer.sanger.ac.uk/>

<sup>25</sup><https://www.ebi.ac.uk/QuickGO/>

<sup>26</sup><https://www.ncbi.nlm.nih.gov/geo/>

<sup>27</sup><https://www.uniprot.org/>

<sup>28</sup><https://pubmed.ncbi.nlm.nih.gov/>

- 425 3. **About:** takes the user to the page that briefly describes the purpose of the NanoWeb system and  
426 summarizes the provided functionalities.
- 427 4. **Contacts:** leads to a page with contact information of the authors of this project.

428 The body of the Web interface consists of three layers displayed alternatively:

- 429 • **Nanopublications list** (Figure 4.A) A list of nanopublications retrieved for the user query. Each  
430 nanopublication is represented with a row in the list, reporting the following information:

- 431 1. The title of the nanopublication (4a).
- 432 2. The assertion of the nanopublication (4b).
- 433 3. A link to the source platform of the data (4c). For instance, in Figure 4 the source platform  
434 of the data is DisGeNET.
- 435 4. The *graph button* to display the graph associated with the nanopublication (4d). When the  
436 user clicks this button, the Graph layer appears to show the nanopublication graph on the  
437 right side of the nanopublications list. If the Information layer is displayed, it is replaced  
438 with the Graph layer.

439 The *Load More* button (Figure 4.5) loads more relevant nanopublications associated with the query,  
440 if any.

441 As we can see in Figure 5, when a user clicks on a specific row, the Information layer is displayed,  
442 showing the information regarding the selected nanopublication.

- 443 • The **information layer** shows information associated with a selected nanopublication, including:
  - 444 1. **Assertion:** (Figure 5.1) This section reports the assertion of the nanopublication of interest  
445 and its title. Besides, meaningful entities, such as the disease *Colorectal Carcinoma*, are  
446 reported as links to external knowledge bases.
  - 447 2. **Publication info:** (Figure 5.2) This section reports the publication information of the clicked  
448 nanopublication. This information includes the creation date, the creators, and the source  
449 platform. Moreover, a link to the data record is provided so that the user can be redirected  
450 to the data record about the assertion; these links act as entry points to external scientific  
451 databases. For instance, Figure 6 shows the data record web page for the nanopublication  
452 with title: *mutL homolog 1 - Colorectal Carcinoma* in DisGeNET.
  - 453 3. **Provenance:** (Figure 5.3) This section shows the provenance information such as the evi-  
454 dence source and how the nanopublication was generated. It also reports the abstract of the  
455 publication, if present.
  - 456 4. **Cite:** (Figure 5.4) This section shows the citation snippet of the nanopublication. The user  
457 can copy the citation text by clicking on the *Cite this nanopub* button in the header.

458 The user can expand/collapse each section by clicking on the title or in the header section.

- 460 • **Graph layer:** Figure 7 shows the Graph layer displayed on the right side of the nanopublications  
461 list after the user click. This layer shows the graph associated with the nanopublication, leveraging  
462 on the RDF triple structure. Each graph node corresponds to the subject or the object of an assertion,  
463 while the edge represents the predicate. Each assertion is represented with a directed edge.

464 The figure shows the graph associated with the *mutL homolog 1 - Colorectal Carcinoma* nanopubli-  
465 cation. The assertion within this nanopublication has two nodes: *mutL homolog 1* as the subject and  
466 *Colorectal Carcinoma* as the object. The subject – a gene – is colored in green, while the object – a  
467 disease – is in red. The predicate connecting the two is represented as an oriented grey edge.

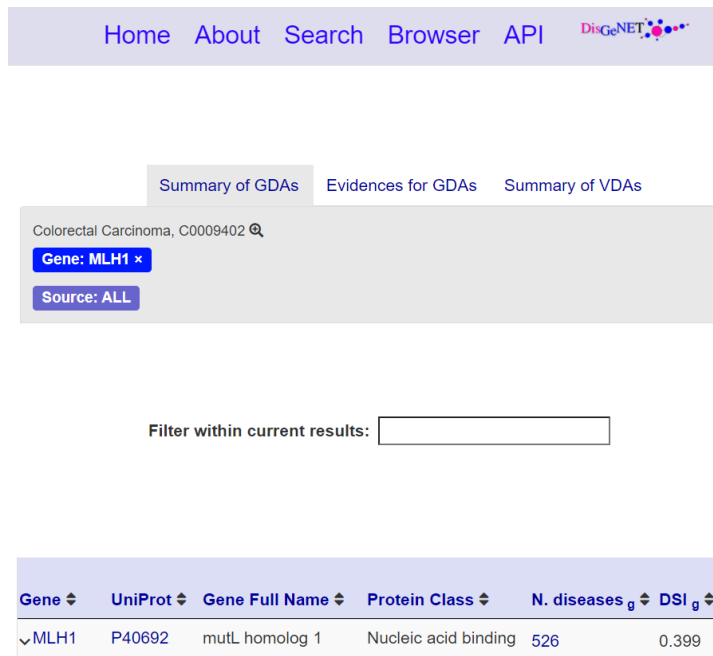
The screenshot displays the DisGeNET interface. At the top, there is a search bar with 'colorectal cancer' entered. Below the search bar, there are navigation options: 'Nanopub List Only', 'Display Both', 'Nanopub Info Only', and 'Show Graph Layer'. The main content area is divided into two parts. On the left, a list of gene-disease associations is shown, each with a 'Go to: DisGeNET' link. The first entry is 'mutL homolog 1 - Colorectal Carcinoma' with the association type 'gene-disease association linked with genetic variation'. A red box highlights the first entry in this list. On the right, the 'Information layer' is expanded for the selected entry. It contains four sections: 1. 'mutL homolog 1 - Colorectal Carcinoma' with a description of the gene-disease association and additional info. 2. 'Publication Info' with details like Nanopublication ID, creation date, and collaborators. 3. 'Provenance' with assertion generation details and evidence source. 4. 'Cite this nanopub' with a citation string and a link to the citation page.

**Figure 5.** Information layer for the nanopublication.

468 There are different ways to interact with the nanopublication graph. For instance, the user can click  
 469 on a node to expand the relation network and visualize other nodes connected to the nanopublication  
 470 of interest. The complete list of the user graphic controls available can be consulted by clicking on  
 471 the *Controls help* button indicated with number three in Figure 7. The figure shows a two-levels  
 472 expansion starting from the subject node *mutL homolog 1* and ending with the expansion of the  
 473 node associated to the *Colorectal Cancer* disease.

474 The possible actions that a user can perform on the graph are:

- 475 – **Expand/collapse graph network:** When the user left-clicks on an unexpanded node, the  
 476 graph is expanded. Thus its relation network is shown. Otherwise, if the user clicks on an  
 477 already expanded node, the graph collapses, and in turn, its relation network is hidden.
- 478 – **Show node information:** When the user right-clicks on a node, a dialog modal window  
 479 appears to show the information concerning that node. For instance, the information window  
 480 shows the type of entity node clicked, such as *gene* or *disease* in case of nodes coming from



**Figure 6.** Data record for the nanopublication with title: *mutL homolog 1 - Colorectal Carcinoma*.

- nanopublications concerning biological or medical fields.
- **Show edge information:** When the user left-clicks on edge, a dialog modal window appears to show the information regarding the nanopublication. Figure 8 shows that when the edge connecting *mutS homolog 6* and *Carcinogenesis* is clicked, the nanopublication information window appears on the right side. The modal dialog window contains the same information of the Information layer. Still, it has a smaller width and can be dragged anywhere inside the Graph layer, so it is always accessible without covering it.
  - **Drag and drop:** The user can drag and move the nanopublication graph by pressing the mouse’s left button and moving it around the graph layer. When the desired position has been chosen, the user can release the left button of the mouse to drop the graph.
  - **Zoom in/out:** Using the mouse wheel, the user can zoom in or out on the nanopublication graph.
  - **Switch between Graph and Information layers:** A button is provided to switch between Graph and Information layers. For instance, when the Graph layer is displayed to go back to the Information layer, the user can click on the *Show Nanopub Info* button (Figure 7.1). In the same way, when the Information layer is displayed, the user can switch to the Graph layer by clicking the *Show Graph Layer* button.
  - **Rearrange layers:** The Navbar menu manages layers disposition (Figure 7.2) and it is provided with the following buttons:
    1. **Nanopub List Only:** It shows a full-screen view of just the nanopublications list layer.
    2. **Display Both:** It opens a two-layers view consisting of the nanopublications list layer and the currently active layer between Graph and Information layers. For instance, Figure 7 shows the Graph layer on the right side of the nanopublications list layer.
    3. **Graph Only/Nanopub Info Only:** It shows a full-screen view of the current layer, which can be the Graph layer or the Information layer. For instance, Figure 7 shows this button with the text “Graph Only”, since the Graph layer is active.

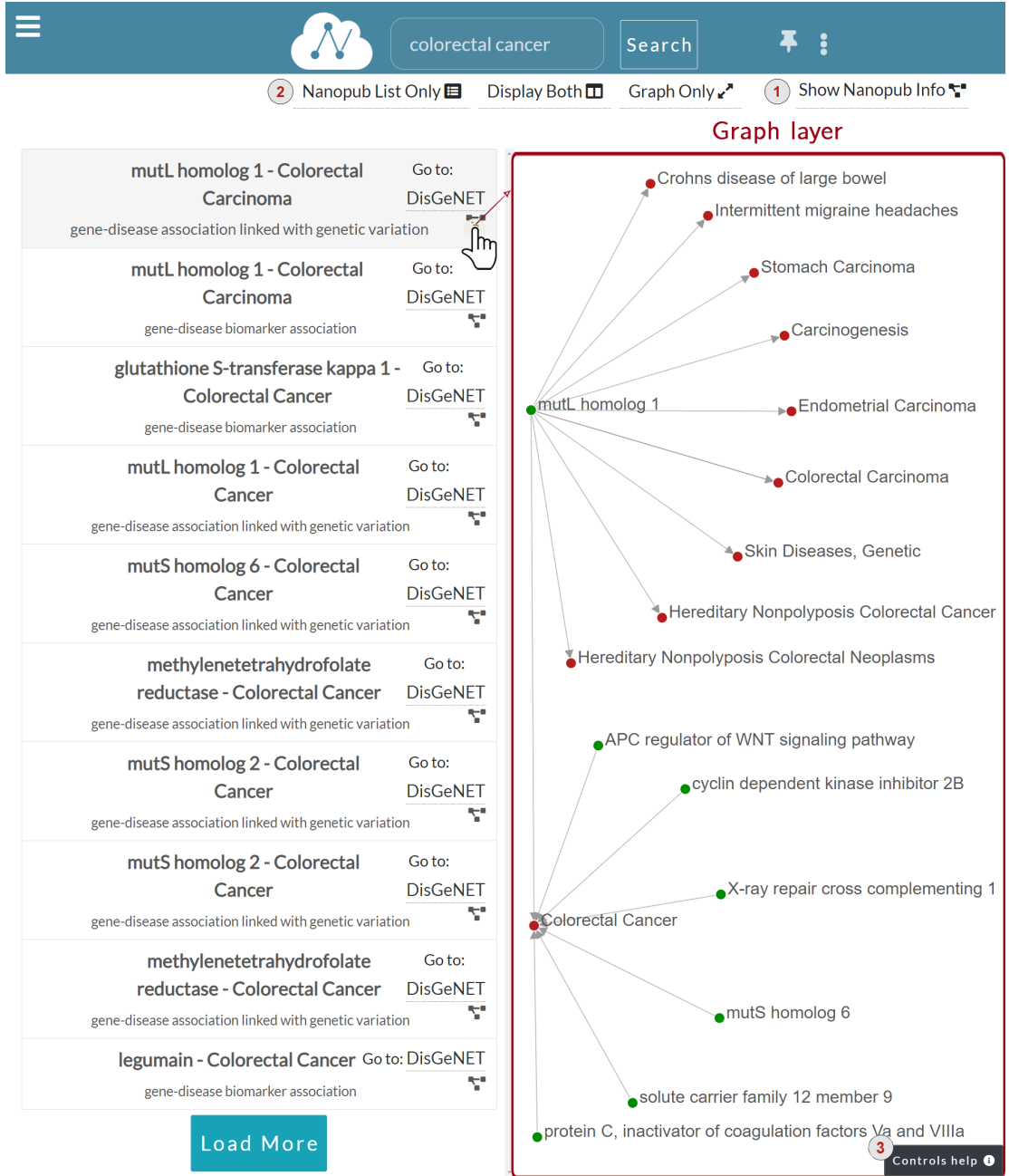
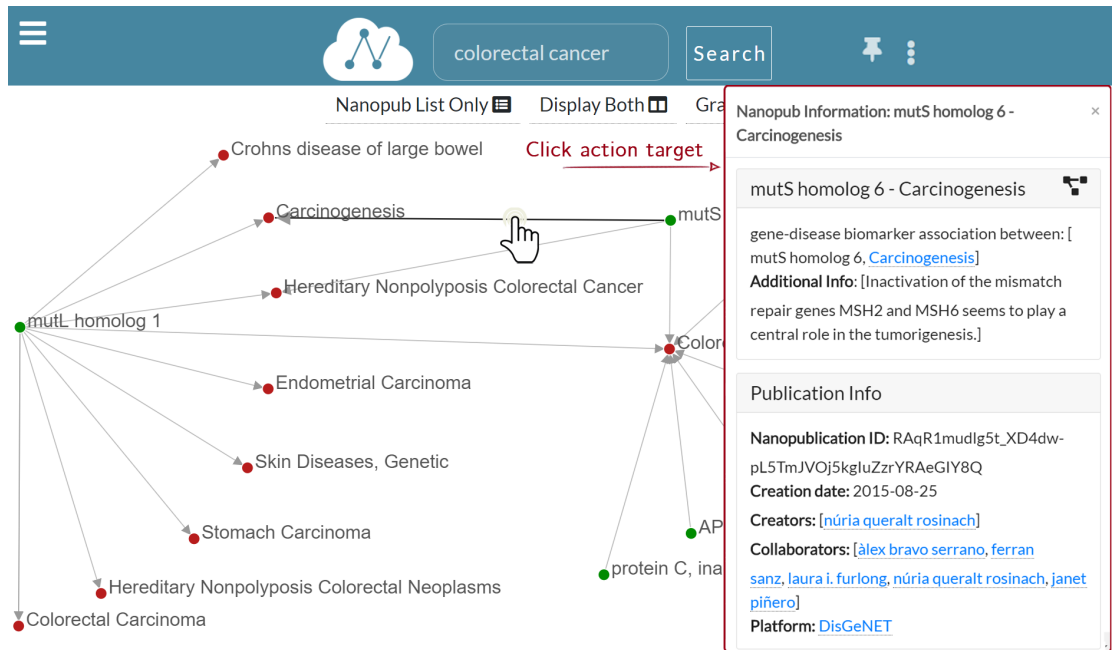
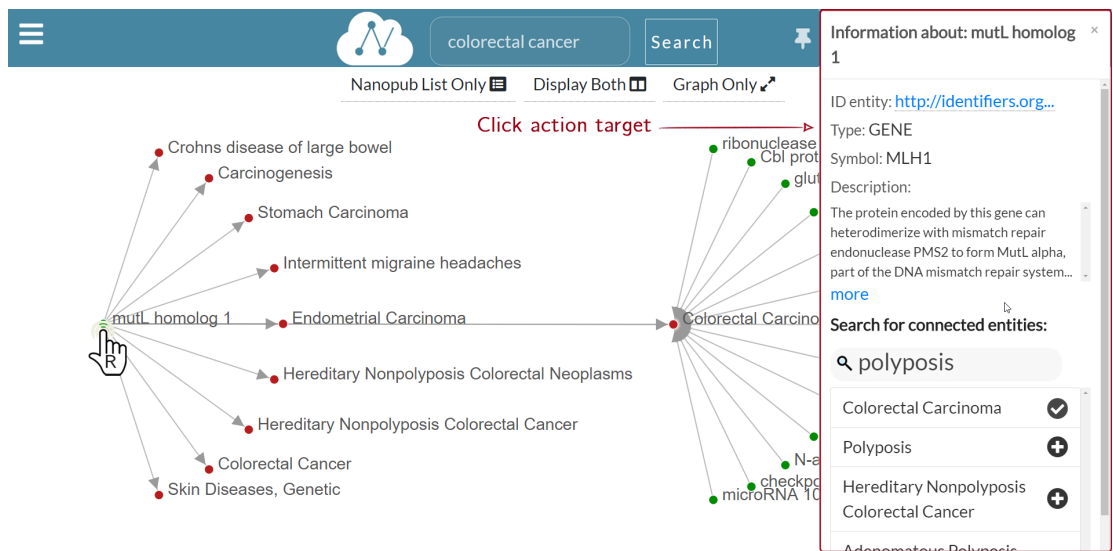


Figure 7. Graph layer for the nanopublication clicked by the user.





**Figure 8.** Graph exploration: the information window for *mutS homolog 6 - Carcinogenesis* is displayed as a result for the user click on the edge.



**Figure 9.** Graph exploration: search for *mutL homolog 1 (MLH1)* connected entities.

## 507 GRAPH EXPLORATION

508 Figure 8 shows a multi-level graph exploration for the nanopublication with the title *mutL homolog 1 -*  
509 *Colorectal Carcinoma*, which describes a gene-disease association. This functionality allows the user to  
510 explore the relation network of the considered nanopublications. Besides, the graph exploration allows the  
511 user to understand how and why different nanopublications are connected. There is no limit to the depth  
512 of the exploration, i.e., to the graph's dimension visualized. The user can potentially expand the graph at  
513 will until all the nodes connected in the relation network are displayed. In this way, the synthesis power  
514 of nanopublications is enhanced by the value of the relation network; it provides a greater information  
515 contribution than the sum of the single nanopublications taken separately. Since the graph can have a high  
516 density of connections, only a portion of the connected nodes is shown for a new graph expansion request.  
517 However, the user could be interested in a specific connection between two nodes, which may not be  
518 shown by default. Hence, it is possible to search for specific connections directly on the nanopublication  
519 semantic network – we call this functionality “connected entities search”. Figure 9 shows the connected  
520 entities search in action. In particular, we see the entities connected to the *mutL homolog 1* gene. When  
521 the user right-clicks on the node associated with the *mutL homolog 1* gene, the information window is  
522 shown on the right side. Inside the information window, there is the “connected entities” input field,  
523 where the user can specify the entity name s/he is looking for. For instance, when the user types *polyposis*,  
524 a list of matching entities appear, and the user can choose which entities to add to the graph by clicking  
525 on the plus button. Using the connected entities search, users can quickly verify whether a direct link  
526 between two nodes exists. The “connected entities search” is provided with auto-completion to ease the  
527 work of the user.

### 528 Implementation specifications

529 NanoWeb back-end is developed using Django,<sup>29</sup> which is a Python-based free and open-source Web  
530 framework. The Web app front-end is developed using HTML5, CSS3, Bootstrap framework,<sup>30</sup> JavaScript,  
531 jQuery,<sup>31</sup> and the library D3.js.<sup>32</sup> In particular, to draw the nanopublication graphs, we used the *D3*  
532 *Force Layout*,<sup>33</sup> which is specifically designed to implement force-directed graphs. A force-directed  
533 graph is a graph where nodes are subjected to forces of two types: attractive and repulsive. These kinds  
534 of forces try to simulate physics scenarios where particles attract or repel each other. Here, the particles  
535 are the nodes of the graph, and the edges represent the presence of forces between nodes. When a new  
536 instance of a force-directed layout is created, a new D3 simulation starts, and the nodes become subjected  
537 to forces. The force-directed layout can be used both for cyclic and acyclic graphs, which can be either  
538 directed or not.

539 To implement the graph exploration, we developed a custom, collapsible force-directed layout where  
540 nodes can be expanded or collapsed at will. This layout enables a user-friendly exploration of graphs  
541 leveraging on a functional disposition of children nodes around the parents.

542 In particular, Figure 8 shows that children nodes are displayed around parents at evenly spaced angles  
543 of an arc. This disposition is designed to facilitate the horizontal expansion of the graph and prevent  
544 nodes from overlapping in a multi-level expansion. The custom force-directed layout developed and the  
545 NanoWeb code are publicly available<sup>34</sup>.

---

<sup>29</sup><https://www.djangoproject.com/>

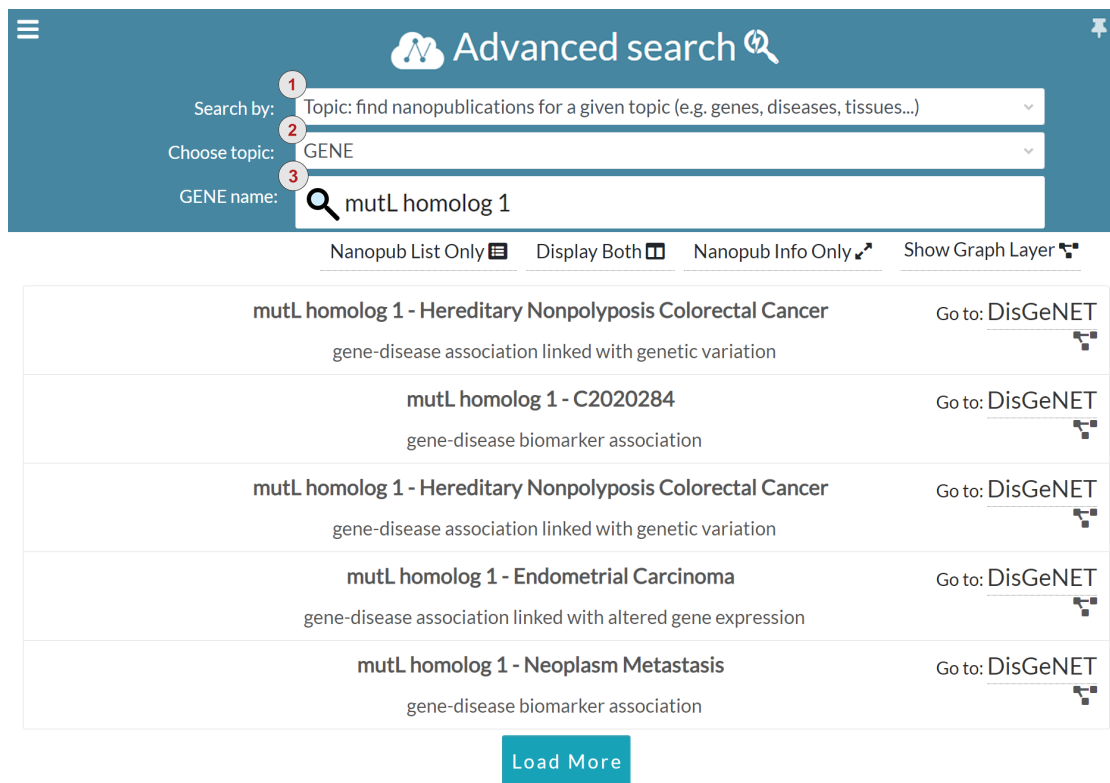
<sup>30</sup><https://getbootstrap.com/>

<sup>31</sup><https://jquery.com/>

<sup>32</sup><https://d3js.org/>

<sup>33</sup>[https://d3-wiki.readthedocs.io/zh\\_CN/master/Force-Layout/](https://d3-wiki.readthedocs.io/zh_CN/master/Force-Layout/)

<sup>34</sup><https://github.com/giachell/nanoweb>



**Figure 10.** Advanced search: search for nanopublications regarding the *mutL homolog 1* gene.

## 546 ADVANCED SEARCH

547 In addition to keyword search, we introduced the advanced search to guide users in query formulation.  
 548 The advanced search is based on structured terms that can be general purpose (e.g. nanopublication  
 549 URLs, author ORCID and scientific evidence identifiers) or domain-specific (e.g. genes, diseases, proteins  
 550 and tissues). Figure 10 shows one of the configurations available in the advanced search interface. The  
 551 interface is based on filters enabling the users to perform boolean search and restrict the search results.  
 552 Users can choose the search modality in the *Search by* drop-down menu, marked with number one in  
 553 Figure 10. The interface provides four different search modalities:

- 554 1. **Topic:** topic-based search is domain-specific, and it allows the user to find nanopublications for a  
 555 specific topic. Currently, the available topics are genes, diseases, proteins, and tissues. The user can  
 556 specify the chosen topic in the *Choose topic* drop-down menu, indicated with number two in Figure  
 557 10. The user can also specify the name of the entity that s/he is looking for in the *Entity name* input  
 558 field, marked with number three in Figure 10. For instance, in Figure 10 the chosen topic is *GENE*  
 559 and the gene name is *mutL homolog 1*. Since gene and protein names could be quite complex to  
 560 remember, the *Entity name* input field is provided with and auto-completion functionality. Once the  
 561 user specifies the details about the topic, the list of related nanopublications is returned, so that the  
 562 user can visualize and explore them as described for the keyword search interface.
- 563 2. **Author:** allows the user to find all the nanopublications related to a nanopublication/evidence  
 564 author. The provided author could be a nanopublication author or the author of the scientific  
 565 publications containing the evidence of nanopublication assertions. Users can search for a specific  
 566 author by providing the author's name or her/his ORCID identifier. The author input field is  
 567 provided with auto-completion for both author names and ORCID identifiers.
- 568 3. **Nanopublication ID:** using this mode, users can search for a specific nanopublication via its  
 569 identifier/URL. The users can take advantage of the auto-completion feature to search for all the  
 570 nanopublications.

571 4. **Evidence:** this mode allows the users to get all the nanopublications extracted from a given  
572 scientific publication (i.e., evidence) starting from the publication DOI or PubMed URL (e.g.,  
573 <http://identifiers.org/pubmed/29970664>).

574 To define the advanced search interface filters we used structured terms (entities) collected from several  
575 public ontologies, databases and terminology resources concerning both life science and medical domains.  
576 For instance, we consider *genes*, *diseases*, *proteins* and *tissues* categories that users can use as filters. The  
577 machine-readable versions of the entities are contained in the nanopublications indexed by NanoWeb. To  
578 obtain their human-readable version, we leverage on public ontologies and databases. From these resources  
579 the associated labels are extracted, stored into the NanoWeb database and then linked to the respective  
580 machine-readable entities. To do so, we used some ontologies: *Basic Formal Ontology (BFO)*<sup>35</sup>, *Chemical*  
581 *Entities of Biological Interest Ontology (CHEBI)*<sup>36</sup>, *Evidence and Conclusion Ontology (ECO)*<sup>37</sup>, *Open*  
582 *Biological and Biomedical Ontology (OBO)*<sup>38</sup>, *Pathway Ontology (PW)*<sup>39</sup> *Semanticscience Integrated*  
583 *Ontology (SIO)*<sup>40</sup>, *Sequence Ontology (SO)*<sup>41</sup>. Additionally, as terminology resources we employed the  
584 *National Center for Biotechnology Information (NCBI)*<sup>42</sup>, *National Cancer Institute Thesaurus (NCIT)*<sup>43</sup>  
585 and the *Unified Medical Language System (UMLS)*<sup>44</sup>.

586 The entities extracted from the resources mentioned above are also used for the mapping of nanopub-  
587 lication assertions – originally modeled as machine-readable RDF statements – into a human-readable  
588 form. To do so, NanoWeb exploits the entity types to determine the proper visual representation of  
589 nanopublication assertions. For instance, in the case of a DisGeNET gene-disease association (dgn-gda),  
590 the entity types are *gene* or *disease*. The entities are represented as nodes labeled with the human-readable  
591 versions of the corresponding URI used in the RDF serialization of the nanopublication. The nodes are  
592 connected together by an oriented edge from *gene* to *disease*. As an example let us consider the assertion  
593 of the nanopublication with identifier: *RA3WLHsGFZrDU4kULrSa\_pTa0gk8-mwadaj-LZ7kAqog*:

```
594 miriam-gene:351 a nci:C16612 .  
595 lld:C0002395 a nci:C7057 .  
596 dgn-gda: DGNa4c88520d1a84e659043089fff632d78 sio:SIO_000628 miriam-gene:351 , lld:  
597 C0002395 ;  
598 a sio:SIO_001121 .
```

599 The assertion describes a *gene-disease association* (dgn-gda) between the NCBI gene *amyloid beta*  
600 *precursor protein* (miriam-gene:351) and the *Alzheimer's disease* (lld:C0002395). The association type is  
601 more specifically a *gene-disease biomarker association* (SIO:001121). NanoWeb enriches the entities  
602 with additional information that can be inferred from the RDF graph of the nanopublication. For instance,  
603 additional information are the types of the entities – e.g. the fact that first entity (miriam-gene:351) is a  
604 gene (nci:C16612) and that (lld:C0002395) is a disease (nci:C7057). All these additional information  
605 are treated as entity properties that the user can access via the interactive visual representation of the  
606 nanopublication. The entity labels *amyloid beta precursor protein* and *Alzheimer's disease* are taken  
607 respectively from the NCBI and Linked Life Data platforms. The entity labels are resolved from entity  
608 identifiers by relying on public API endpoints such as the *Entrez Programming Utilities (E-utilities)*<sup>45</sup>  
609 provided by NCBI. Nanopublications from the same platform (e.g. DisGeNET, NeXtProt, Protein Atlas,  
610 and Wikipathways) use the same authorities to identify entities (e.g. genes, diseases, proteins and tissues).  
611 However, when nanopublications from different platforms are visualized, it is sometimes necessary to  
612 reconcile different resource identifiers across authorities to link the same entities to others using different  
613 identifiers. In the visual representation only one valid identifier is presented for each entity to keep the  
614 interface as clean as possible.

<sup>35</sup><https://basic-formal-ontology.org/>

<sup>36</sup><https://www.ebi.ac.uk/chebi/>

<sup>37</sup><https://www.evidenceontology.org/>

<sup>38</sup><http://www.obofoundry.org/>

<sup>39</sup><https://rgd.mcw.edu/rgdweb/ontology/search.html>

<sup>40</sup><https://github.com/MaastrichtU-IDS/semanticscience>

<sup>41</sup><http://www.sequenceontology.org/>

<sup>42</sup><https://www.ncbi.nlm.nih.gov/>

<sup>43</sup><https://ncit.nci.nih.gov/ncitbrowser/>

<sup>44</sup><https://www.nlm.nih.gov/research/umls/index.html>

<sup>45</sup><https://www.ncbi.nlm.nih.gov/books/NBK25501/>

## 6 EXPERT USERS SURVEY

To better understand the needs of the nanopublication community and improve the critical functionalities of NanoWeb, we conducted an expert users survey to collect feedback from nanopublication and domain experts. We advertised NanoWeb on the nanopublication public mailing lists, on social media targeting the potentially interested communities and private emails to the authors of papers about nanopublications. We asked the nanopublication experts involved in the survey to use NanoWeb, and then to answer a questionnaire. It should be noticed that we did not provide any tutorial to inform the users about NanoWeb functions because we also wanted to investigate how intuitive the system is for first-time users and how steep its learning curve is.

The survey was composed of sixteen questions (Q[1-16]) divided in four sections. The majority of the questions is answered through the Likert five-point scale, ranging from 1 to 5 points, meaning different things depending on the question.

1. **Personal information.** This section is composed of four questions and collects basic information about the participants and their experience with nanopublications:

- **Q1:** *Do you have any experience with nanopublications?*

In this case the answer with 1 point in the Likert scale means: “Not at all” (i.e., I heard someone mentioning nanopublications once), while the 5 points one means: “Quite a lot” (i.e., I created some nanopublications myself)

- **Q2:** *Current Position?*

Single choice between: Academic, Industry, Master Student, PhD Student, PostDoc.

- **Q3:** *Primary domain of expertise?*

Multiple choices between: Art and architecture, Biology, Chemistry, Communication Science, Computers and the humanities, Computer Science, Economics, Life Science, Linguistics, Mathematics, Medicine, Physics, Psychology, Sociology.

The survey considered fourteen participants in total, counting seven highly-experienced users (5 on the Likert scale) and nine experienced users (4 on the Likert scale). According to the data collected, the majority of the participants (85.7%) are from Academia. Also, according to Q3, the main domains of expertise of the participants are: Computer Science (57.1%), Chemistry (35.7%), Life Science (35.7%), Biology (28.6%), Medicine (14.3%). Computer Science indicates experts in the creation of nanopublications from the technical viewpoint, whereas the others are domain experts who might curate or use nanopublications in their daily work.

2. **The relevance of the addressed problem.** This section explores the existence and quality of other services enabling search, access, exploration, and re-use of nanopublications (all questions are answered according to a 1 (not at all) to 5 (quite a lot) Likert scale):

- **Q4:** *Is searching, accessing, and consulting nanopublications relevant for the stakeholders (e.g., researchers, developers, domain experts)?*

- **Q5:** *To the best of your knowledge, are the currently available tools and services adequate for searching and accessing nanopublications?*

- **Q6:** *To the best of your knowledge, do other tools and services offer interactive visualizations to interact with nanopublications?*

- **Q7:** *To the best of your knowledge, do other available tools and services offer visual exploration possibilities of the nanopublication relation network?*

According to the data collected for questions Q[4-7], the majority of the participants (57%) considers the problem addressed by NanoWeb relevant or very relevant, pointing out the lack of other tools and services for the interactive visualization and exploration of nanopublications and their relation network. About Q5, 50% of the participants consider the currently available tools and services for searching and accessing nanopublications inadequate (1 or 2 points on the Likert scale) and 42%

662 are not enthusiastic about them (3 points on the Likert scale). 71% of the participants answered  
663 that there are no other available tools offering interactive visualizations of nanopublications and  
664 57% say there are no alternative tools to visually explore the nanopublication network. From these  
665 answers, we can see that the participants confirm our analysis highlighting the lack of intuitive and  
666 visual tools for the access and exploration of the nanopublications despite the confirmed utility of  
667 searching and accessing nanopublications for the stakeholders.

668 3. **NanoWeb - Search Engine and Interface.** The questions of this section are designed to evaluate  
669 the search capabilities of NanoWeb and the usability of its interface. This section was answered by  
670 twelve participants over fourteen.

- 671 • **Q8:** *Is NanoWeb search interface intuitive and easy-to-use?*
- 672 • **Q9:** *Is NanoWeb capable of retrieving relevant nanopublications for a given query?*
- 673 • **Q10:** *In your opinion, is a search based on keywords an effective way to seek for nanopubli-*  
674 *cations?*
- 675 • **Q[11-12]:** *In your opinion, for the not technologically savvy, what is the most effective way*  
676 *to search nanopublications? Q11 and Q12 are the same, but the answers are different since*  
677 *for Q11 the range of answers is from 1: SPARQL end-point to 5: Keyword-based search;*  
678 *whereas, for Q12 the range is from 1: Faceted search to 5: Keyword search.*
- 679 • **Q13:** *“Will NanoWeb enhance the productivity of involved stakeholders (researchers, devel-*  
680 *opers, nanopublication experts)?*

681 About question Q8, the majority of the participants consider NanoWeb search interface intuitive  
682 and easy-to-use (75% answered 4 or above and none answered below 3). There is no accordance  
683 instead for Q9 (median = 3, mean = 3.08, STD = 1.04), 42% of the participants answered 3 which  
684 means “not sure” and the rest of them is divided into the two other classes “not really” ( $\leq 2$ : 33%)  
685 and “quite a lot” ( $\geq 4$ : 25%). One reason that could motivate this kind of distribution might be that  
686 participants did not know what they could search in advance, thus many user queries might have not  
687 produced the expected results. To address this issue, after the survey we introduced the advanced  
688 search which guides users on NanoWeb search capabilities. Participants are well-distributed for  
689 Q10 (median = 3, mean = 3.25, STD = 1.16), there is not a preferred opinion about keyword  
690 search; nevertheless, 46% of the participants consider the search based on keywords quite an  
691 effective or highly effective (answer 4 or above) way to seek for nanopublications. About Q[11-12],  
692 the majority of the participants (58%) consider that keyword-based search is more effective than  
693 SPARQL end-point but less effective than faceted search (67%) for the non-technologically savvy.  
694 This answer shows how domain experts are more accustomed to use faceted search rather than  
695 keyword search for searching structured data as nanopublications are. Keyword search is considered  
696 useful, but it should not substitute faceted search as a means to access RDF scientific data. Finally,  
697 all the participants believe NanoWeb can moderately (58%) or substantially (42%) enhance the  
698 productivity of researchers and nanopublication experts.

699 4. **NanoWeb - Visual Exploration.** This section of the questionnaire evaluates the experience with  
700 the NanoWeb user interface for visual exploration of nanopublications. We designed the questions  
701 of this section to investigate whether the visual exploration of nanopublication graphs could lead  
702 to the discovery of meaningful relationships and information potentially unknown to the experts.  
703 Moreover, we asked the participants to compare NanoWeb with the currently available alternative  
704 tools. This section consists of three questions:

- 705 • **Q14:** *Do you feel comfortable with the interface for the visual exploration?*
- 706 • **Q15:** *Could the visual exploration of the nanopublication graphs lead to the discovery of*  
707 *meaningful relationships and information not known in advance?*
- 708 • **Q16:** *Is NanoWeb visual exploration innovative with respect to the currently available*  
709 *alternative tools and techniques?*

710 With reference to Q14, the majority (64%) of the participants felt very comfortable with the interface  
711 for the visual exploration and only 14% gave a score below three points. Moreover, 57% of the  
712 participants believe the visual exploration of the nanopublication graphs could lead to the discovery  
713 of meaningful relationships and information not known in advance. Finally, half of the participants  
714 think that NanoWeb is highly innovative (four or five points) with respect to the state of the art,  
715 while only 21% thinks it is only marginally innovative.

## 716 6.1 User feedback

717 Finally, we asked the participants to provide some feedback and suggestions to improve NanoWeb. The  
718 feedback collected shows that users have appreciated the system:

- 719 • *“I very much appreciate the tool, and I think it can be a great push for better accessing and using  
720 nanopublications by everyone!”*
- 721 • *“I consider the NanoWeb proposal a smart insight for searching nanopublications.”*

722 We also received useful suggestions to improve the system:

- 723 • *“I found the visual exploration innovative, but I think it could be improved by a better UI/UX.”*
- 724 • *“Good work! I would suggest that you enable URL-based searching.”*
- 725 • *“Consider replacement of keyword search with a concept-based search. This can also be used to  
726 enable auto-suggest functionality based on the resources (genes, diseases, etc)”*
- 727 • *“I really like the application, but at the end of the day it is dependent on the indexed data. It would  
728 be great if there were a possibility to suggest datasets to be included or even better, to be able to  
729 add them myself!”*
- 730 • *“Downloading of the results as a dataset of nanopublications would be most welcome too. Even  
731 better, a Cytoscape plugin that allows me to pull in the full network. I’m looking forward to seeing  
732 where you are taking this. Success!”*

733 We consider the user feedback of great value, so we decided to improve NanoWeb according to the  
734 received suggestions. Firstly, we improved both the user interface and experience (UI/UX), providing  
735 a responsive mobile device layout. Then, we improved the search system so that a user can perform  
736 URL-based searching. Currently, NanoWeb allows the users to find the authors from the ORCID ids; a  
737 specific nanopublication from its URL/identifier; and, all the nanopublications related to one particular  
738 evidence paper provided its DOI.

739 The prominent feature we added to NanoWeb, thanks to the user feedback, is the advanced search, as  
740 described in Section 5. The Advanced search interface is based on structured terms extracted from the  
741 life science domain, it enables users to search for nanopublications based on topics (e.g., genes, diseases,  
742 proteins, etc.), scientific evidence, and authors. Finally, based on the collected feedback, we planned  
743 several further improvements to the system that we discuss as future work in Section 8.

## 744 7 DISCUSSION ON MAINTAINING ASPECTS

745 NanoWeb aims to provide users unified access to nanopublications and to search and explore them through  
746 a human-readable interface. Since NanoWeb is tailored for both the life science and medical domains, it  
747 is designed to help the experts of these domains in their research work. It also allows users that do not  
748 have a prior knowledge about nanopublication to easily interpret and understand the returned content.

749 Several challenges need to be addressed to maintain a stable, citable system like what NanoWeb aims  
750 to be. The major system maintaining challenges are:

- 751 1. **Ensure persistent access and re-use of data:** to guarantee persistent and reliable access to data  
752 and avoid broken URLs, NanoWeb uses persistent URLs and identifiers to refer to resources. All  
753 the indexed nanopublications are directly accessible through a persistent URL provided by the  
754 *W3C Permanent Identifier Community Group*<sup>46</sup>. The nanopublication’s persistent URL format

---

<sup>46</sup><https://w3id.org/>

755 is: <http://w3id.org/nanoweb/landingpage/<ID>>, where the *ID* in brackets is the  
756 nanopublication identifier and satisfies the regular expression:  $\text{^RA[A-Za-z0-9_\-]{43}\$}$ .  
757 Nanopublications use persistent identifiers, that allow to access them across different providers.  
758 Even if one of the several nanopublication providers is unreachable in a given moment, the others  
759 can provide access by using the same identifiers. As for nanopublications, NanoWeb itself is  
760 reachable through the persistent URL: <http://w3id.org/nanoweb/>.

- 761 2. **Long-term preservation of resources:** every information concerning nanopublications is saved  
762 in NanoWeb databases, that are stored in network hard drives using redundancy policies such as  
763 *Redundant Array of Independent Disks (RAID)*. The redundancy policies adopted and daily back-up  
764 routines are designed to prevent loss of data and ensure long-term preservation.
- 765 3. **Ongoing hosting:** NanoWeb is hosted within the cloud architecture of the University of Padova.  
766 The institutional cloud architecture and network infrastructure provide a reliable connection service  
767 as well as a protection layer from external attacks. A team of system administrators actively control  
768 the cloud/network infrastructure and support NanoWeb. NanoWeb is developed in the context of  
769 the European project ExaMode<sup>47</sup> which guarantees financial support until 2023. Within the project  
770 there are sustainability policies that should guarantee the maintenance of the developed tools well  
771 beyond the termination of the project.

## 772 8 CONCLUSIONS

773 Scientific and scholarly communications are growing at an incredible speed, and it is hardly possible  
774 to keep track of the discoveries and statements presented in the literature, even considering only a  
775 specific domain. Moreover, the “redundancy of statements in multiple fora makes it difficult to identify  
776 attribution, quality, and provenance” (Groth et al., 2010). Hence, the nanopublication model has been  
777 proposed to quickly identify, search, and access scientific facts extracted from papers. Nanopublications  
778 are represented as graphs centered on a scientific statement (i.e., the assertion) that makes provenance,  
779 attribution, and scientific information machine-readable.

780 Nanopublications are concise noise-free resources characterized by high information density. Lever-  
781 aging on the semantic-oriented RDF structure, nanopublications efficiently convey information and  
782 concepts. Hence, these features make nanopublications particularly suitable for enabling data search,  
783 information extraction, and automatic reasoning over scientific facts. Despite the promising features of  
784 nanopublications, their use is still restricted to highly-specialized scientific circles.

785 The central limit to the full exploitation of nanopublications is the lack of services enabling their  
786 search, access, exploration, and re-use. Search is limited to the use of structured query languages as  
787 SPARQL, and a service to search over all the publicly available nanopublications at once is not available.  
788 Nanopublications are machine-readable, but no human-readable counterpart is generated and open to the  
789 public. Nanopublications create a vast relation network of scientific facts that could lead to discoveries,  
790 but up to now, there are no automatic or manual services enabling graph exploration.

791 The goal of this work is to provide unified access to Life Science nanopublications in order to allow  
792 users to search, access, explore, and re-use them on the Web. To this end, we have designed and developed  
793 a Web application called *NanoWeb*, that allows the users to (i) search for domain-specific nanopublications  
794 using keywords (as they are accustomed to do with Web search engines); (ii) explore their relation  
795 network to discover new nanopublications and meaningful connections; (iii) access and understand their  
796 content; (iv) connect to the evidence paper and access the related data record in external curated scientific  
797 databases; and, (v) easily cite nanopublications when they are re-used in new scientific contexts.

798 We also presented the benefits of the serendipity-oriented perspective enabled by NanoWeb in the  
799 Life Science domain. We showed how the exploration of nanopublication graphs could enrich domain  
800 knowledge and point out interesting gene-disease connections.

801 As future work, we plan to extend the system by providing the user with the capability of exploring a  
802 new graph generated from an arbitrary set of Life Science nanopublications selected by the user. This  
803 functionality represents a significant improvement for the graph exploration since the initial relation  
804 network already considers different nanopublications, instead of starting the graph exploration from a

---

<sup>47</sup>European Union Horizon 2020 program under Grant Agreement no. 825292



805 single one. In this way it is possible to highlight, for instance, the set of common diseases due to a  
806 selection of genes or, conversely, the set of common genes that cause the disease of interest. Moreover,  
807 we plan to crawl and index the Life Science nanopublications that are not currently available on the Web,  
808 if not downloading large archive files which are hardly usable.

809 As future work, we plan to further improve NanoWeb according to the expert users survey's feedback.  
810 We will allow the users to add datasets or other domain-specific nanopublication sources to be crawled  
811 and indexed by the system. We will add the possibility to select and download custom-made sets of  
812 nanopublications. We will propose a customized user experience to save lists of favorite nanopublications,  
813 entities, and associations and notify when something new is published.

814 We will dedicate a fair amount of work to the extension of search functionalities to improve keyword  
815 search and to include faceted search which is required by the stakeholders. Indeed, faceted search is  
816 commonly adopted solution (Arenas et al., 2016) to search RDF data. A faceted search is particularly  
817 useful when it is applied to domain-specific data. For instance, in gene-disease associations, the faceted  
818 search can be used to search for specific genes or specific diseases, filtering out all the entities not relevant  
819 to the search. Faceted search can be associated with auto-completion functionalities to ease the users'  
820 work. Finally, we plan to improve keyword-based searches with ontology and database ID lookups.

## 821 ACKNOWLEDGMENTS

822 This work is supported by the Computational Data Citation (CDC-STARs) project of the University of  
823 Padua and by the ExaMode project, as part of the European Union Horizon 2020 program under Grant  
824 Agreement no. 825292. There was no additional external funding received for this study. The authors  
825 wish to thank Erika Fabris for the work on the citation of nanopublications and the development of some  
826 APIs used in this work.

## 827 REFERENCES

- 828 Agrawal, S., Chaudhuri, S., and Das, G. (2002). Dbxplorer: A system for keyword-based search over  
829 relational databases. In *Proceedings of the 18th International Conference on Data Engineering, ICDE*  
830 *2002*, pages 5–16. IEEE Computer Society.
- 831 Amith, M. and Tao, C. (2018). Representing Vaccine Misinformation using Ontologies. *Journal of*  
832 *Biomedical Semantics*, 9(1):22.
- 833 Arenas, M., Cuenca Grau, B., Kharlamov, E., Marciuška, S., and Zheleznyakov, D. (2016). Faceted search  
834 over RDF-based knowledge graphs. *Journal of Web Semantics*, 37-38:55 – 74.
- 835 Balmin, A., Hristidis, V., Koudas, N., Papakonstantinou, Y., Srivastava, D., and Wang, T. (2003). A system  
836 for keyword proximity search on XML databases. In *Proceedings of 29th International Conference on*  
837 *Very Large Data Bases, VLDB*, pages 1069–1072. Morgan Kaufmann.
- 838 Bast, H., Buchhold, B., and Haussmann, H. (2016). Semantic search on text and knowledge bases.  
839 *Foundations and Trends in Information Retrieval (FnTIR)*, 10(2-3):119–271.
- 840 Bhalotia, G., Hulgeri, A., Nakhe, C., Chakrabarti, S., and Sudarshan, S. (2002). Keyword searching and  
841 browsing in databases using BANKS. In *Proceedings of the 18th International Conference on Data*  
842 *Engineering*, pages 431–440. IEEE Computer Society.
- 843 Biryukov, M., Groues, V., Satagopam, V., and Schneider, R. (2018). Biokb-text mining and semantic  
844 technologies for biomedical content discovery.
- 845 Bizer, C., Heath, T., and Berners-Lee, T. (2009). Linked Data – The Story So Far. *Int. J. Semantic Web*  
846 *Inf. Syst.*, 5(3):1–22.
- 847 Borgman, C. L. (2015). *Big Data, Little Data, No Data*. MIT Press.
- 848 Campregher, C., Honeder, C., Chung, H., Carethers, J. M., and Gasche, C. (2010). Mesalazine reduces  
849 mutations in transforming growth factor  $\beta$  receptor ii and activin type ii receptor by improvement of  
850 replication fidelity in mononucleotide repeats. *Clin Cancer Res*, 16(6):1950–1956.
- 851 Carroll, J. J., Bizer, C., Hayes, P., and Stickler, P. (2005). Named graphs, provenance and trust. In  
852 *Proceedings of the 14th international conference on World Wide Web*, pages 613–622. ACM Press.
- 853 Carroll, J. J. and Stickler, P. (2004). RDF triples in XML. In *Proc. of the WWW 2004 Conference*  
854 *(Alternate Track Papers & Posters)*, pages 412–413. ACM Press.
- 855 Chapman, A., Simperl, E., Koesten, L., Konstantinidis, G., Ibáñez, L. D., Kacprzak, E., and Groth, P.  
856 (2020). Dataset search: a survey. *VLDB J.*, 29(1):251–272.

857 Cheney, J., Chiticariu, L., and Tan, W. (2009). Provenance in databases: Why, how, and where. *Foundations and Trends in Databases*, 1(4):379–474.

858 Chichester, C., Gaudet, P., Karch, O., Groth, P. T., Lane, L., Bairoch, A., Mons, B., and Loizou, A. (2014). Querying neXtProt nanopublications and their value for insights on sequence variants and tissue expression. *J. Web Semant.*, 29:3–11.

862 Chichester, C., Karch, O., Gaudet, P., Lane, L., Mons, B., and Bairoch, A. (2015). Converting neXtProt into Linked Data and nanopublications. *Semantic Web*, 6(2):147–153.

864 Coffman, J. and Weaver, A. C. (2014). An Empirical Performance Evaluation of Relational Keyword Search Techniques. *IEEE Trans. Knowl. Data Eng.*, 26(1):30–42.

865 Dosso, D. and Silvello, G. (2020). Search text to retrieve graphs: A scalable RDF keyword-based search system. *IEEE Access*, 8:14089–14111.

866 Elbassuoni, S. and Blanco, R. (2011). Keyword Search over RDF Graphs. In *Proc. of the 20th ACM Conference on Information and Knowledge Management, CIKM 2011*, pages 237–242. ACM Press, New York, USA.

870 Fabris, E., Kuhn, T., and Silvello, G. (2019). A framework for citing nanopublications. In Doucet, A., Isaac, A., Golub, K., Aalberg, T., and Jatowt, A., editors, *Proc. of the 23rd International Conference on Theory and Practice of Digital Libraries, TPD 2019*, volume 11799 of *Lecture Notes in Computer Science*, pages 70–83. Springer.

874 Fafalios, P. and Tzitzikas, Y. (2013). X-ENS: semantic enrichment of web search results at real-time. In *The 36th International ACM SIGIR conference on research and development in Information Retrieval, SIGIR '13*, pages 1089–1090. ACM Press, New York, USA.

876 Groth, P., Gibson, A., and Velterop, J. (2010). The Anatomy of a Nanopublication. *Inf. Serv. Use*, 30(1-2):51–56.

878 Heim, P., Ziegler, J., and Lohmann, S. (2008). gFacet: A Browser for the Web of Data. In *Proceedings of the International Workshop on Interacting with Multimedia Content in the Social Semantic Web (IMC-SSW'08)*, volume 417 of *CEUR Workshop Proceedings*. CEUR-WS.org.

880 Hettne, K. M., Thompson, M., van Haagen, H., van der Horst, E., Kaliyaperumal, R., Mina, E., Tatum, Z., Laros, J. F. J., van Mulligen, E. M., Schuemie, M., Aten, E., Li, T. S., Bruskiewich, R., Good, B. M., Su, A. I., Kors, J. A., den Dunnen, J., van Ommen, G.-J. B., Roos, M., 't Hoen, P. A., Mons, B., and Schultes, E. A. (2016). The Implicitome: A Resource for Rationalizing Gene-Disease Associations. *PLOS ONE*, 11(2):1–21.

882 Hey, T., Tansley, S., and Tolle, K., editors (2009). *The Fourth Paradigm: Data-Intensive Scientific Discovery*. Microsoft Research, USA.

884 Kadilierakis, G., Fafalios, P., Papadakos, P., and Tzitzikas, Y. (2020). Keyword Search over RDF Using Document-Centric Information Retrieval Systems. In *The Semantic Web*, pages 121–137, Cham. Springer International Publishing.

888 Koplíku, A., Pinel-Sauvagnat, K., and Boughanem, M. (2014). Aggregated search: A new information retrieval paradigm. *ACM Comput. Surv.*, 46(3):41:1–41:31.

892 Kuhn, T., Barbano, P. E., Nagy, M. L., and Krauthammer, M. (2013). Broadening the Scope of Nanopublications. In *Proc. of the Semantic Web: Semantics and Big Data, 10th International Conference, ESWC 2013*, volume 7882 of *LNCS*, pages 487–501. Springer.

896 Kuhn, T., Meroño-Peñuela, A., Malic, A., Poelen, J. H., Hurlbert, A. H., Ortiz, E. C., Furlong, L. I., Queralt-Rosinach, N., Chichester, C., Banda, J. M., Willighagen, E. L., Ehrhart, F., Evelo, C. T. A., Malas, T. B., and Dumontier, M. (2018). Nanopublications: A Growing Resource of Provenance-Centric Scientific Linked Data. In *14th IEEE International Conference on e-Science, e-Science 2018*, pages 83–92. IEEE Computer Society.

900 Kuhn, T., Willighagen, E., Evelo, C., Queralt-Rosinach, N., Centeno, E., and Furlong, L. I. (2017). Reliable granular references to changing linked data. In d'Amato, C., Fernandez, M., Tamma, V., Lecue, F., Cudré-Mauroux, P., Sequeda, J., Lange, C., and Heflin, J., editors, *The Semantic Web – ISWC 2017*, pages 436–451, Cham. Springer International Publishing.

904 Lopez-Veyna, J. I., Sosa, V. J. S., and López-Arévalo, I. (2012). A Virtual Document Approach for Keyword Search in Databases. In *DATA*, pages 39–48. SciTePress.

908 Luo, Y., Wang, W., Lin, X., Zhou, X., Wang, J., and Li, K. (2011). SPARK2: top-k keyword query in relational databases. *IEEE Trans. Knowl. Data Eng.*, 23(12):1763–1780.

910 Mass, Y. and Sagiv, Y. (2016). Virtual Documents and Answer Priors in Keyword Search over Data

912 Graphs. In *Proc. of the Workshops of the EDBT/ICDT 2016 Joint Conference*, volume 1558 of *CEUR*  
913 *Workshop Proceedings*. CEUR-WS.org.

914 McCusker, J., Rashid, S. M., Agu, N., Bennett, K. P., and McGuinness, D. L. (2018). The Whyis  
915 Knowledge Graph Framework in Action. In *Proc. of the ISWC 2018 Posters & Demonstrations,*  
916 *Industry and Blue Sky Ideas Tracks co-located with 17th International Semantic Web Conference (ISWC*  
917 *2018)*, volume 2180 of *CEUR Workshop Proceedings*. CEUR-WS.org.

918 McCusker, J. P., Dumontier, M., Yan, R., He, S., Dordick, J. S., and McGuinness, D. L. (2017). Finding  
919 melanoma drugs through a probabilistic knowledge graph. *PeerJ Computer Science*, 3:e106.

920 Mons, B., van Haagen, H., Chichester, C., Hoen, P.-B., den Dunnen, J. T., van Ommen, G., van Mulligen,  
921 E., Singh, B., Hooft, R. and Roos, M., Hammond, J., Kiesel, B., Giardine, B., Velterop, J., Groth, P.,  
922 and Schultes, E. (2011). The value of data. *Nature Genetics*, 43(4):281–283.

923 Page, R. (2018). Liberating links between datasets using lightweight data publishing: an example using  
924 plant names and the taxonomic literature. *Biodiversity Data Journal*, 6:e27539.

925 Pérez, J., Arenas, M., and Gutierrez, C. (2009). Semantics and Complexity of SPARQL. *ACM Trans.*  
926 *Database Syst.*, 34(3):1–45.

927 Piñero, J., Ramírez-Anguita, J. M., Saüch-Pitarch, J., Ronzano, F., Centeno, E., Sanz, F., and Furlong,  
928 L. I. (2019). The DisGeNET knowledge platform for disease genomics: 2019 update. *Nucleic Acids*  
929 *Research*, 48(D1):D845–D855.

930 Pontén, F., Jirström, K., and Uhlen, M. (2008). The human protein atlas—a tool for pathology. *The*  
931 *Journal of Pathology*, 216(4):387–393.

932 Pound, J., Mika, P., and Zaragoza, H. (2010). Ad-hoc object retrieval in the web of data. In *Proc. of the*  
933 *19th International Conference on World Wide Web, WWW 2010*, pages 771–780. ACM Press, New  
934 York, USA.

935 Queralt-Rosinach, N., Kuhn, T., Chichester, C., Dumontier, M., Sanz, F., and Furlong, L. I. (2016).  
936 Publishing DisGeNET as Nanopublications. *Semantic Web*, 7(5):519–528.

937 Rahman, P., Jiang, L., and Nandi, A. (2020). Evaluating Interactive Data Systems. *VLDB J.*, 29(1):119–  
938 146.

939 Robertson, S. E., Walker, S., Jones, S., Hancock-Beaulieu, M. M., and Gatford, M. (1994). Okapi at  
940 TREC–3. In Harman, D. K., editor, *Overview of the Third Text REtrieval Conference (TREC-3)* ,  
941 pages 109–126. National Institute of Standards and Technology (NIST), Special Publication 500-225,  
942 Washington, USA. [http://trec.nist.gov/pubs/trec3/t3\\_proceedings.html](http://trec.nist.gov/pubs/trec3/t3_proceedings.html) [last  
943 visited 2007, March 23].

944 Silvello, G. (2018). Theory and Practice of Data Citation. *Journal of the American Society for Information*  
945 *Science and Technology (JASIST)*, 69(1):6–20.

946 Simitsis, A., Koutrika, G., and Ioannidis, Y. E. (2008). Précis: from unstructured keywords as queries to  
947 structured databases as answers. *VLDB J.*, 17(1):117–149.

948 Slenter, D. N., Kutmon, M., Hanspers, K., Riutta, A., Windsor, J., Nunes, N., Mélius, J., Cirillo, E.,  
949 Coort, S. L., Digles, D., Ehrhart, F., Giesbertz, P., Kalafati, M., Martens, M., Miller, R., Nishida,  
950 K., Rieswijk, L., Waagmeester, A., Eijssen, L. M. T., Evelo, C. T., Pico, A. R., and Willighagen,  
951 E. L. (2017). WikiPathways: a multifaceted pathway database bridging metabolomics to other omics  
952 research. *Nucleic Acids Research*, 46(D1):D661–D667.

953 The Economist, (2017). The world’s most valuable resource is no longer oil, but data. The  
954 Economist: New York, USA. [https://www.economist.com/leaders/2017/05/06/](https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data)  
955 [the-worlds-most-valuable-resource-is-no-longer-oil-but-data](https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data) [online Oc-  
956 tober 2020]

957 Uhlen, M., Zhang, C., Lee, S., Sjöstedt, E., Fagerberg, L., Bidkhor, G., Benfiteas, R., Arif, M., Liu, Z.,  
958 Edfors, F., Sanli, K., von Feilitzen, K., Oksvold, P., Lundberg, E., Hober, S., Nilsson, P., Mattsson, J.,  
959 Schwenk, J. M., Brunnström, H., Glimelius, B., Sjöblom, T., Edqvist, P.-H., Djureinovic, D., Micke,  
960 P., Lindskog, C., Mardinoglu, A., and Ponten, F. (2017). A pathology atlas of the human cancer  
961 transcriptome. *Science*, 357(6352).

962 Waagmeester, A., Kutmon, M., Riutta, A., Miller, R., Willighagen, E. L., Evelo, C. T., and Pico, A. R.  
963 (2016). Using the Semantic Web for Rapid Integration of WikiPathways with Other Biological Online  
964 Data Resources. *PLOS Computational Biology*, 12(6):1–11.

965 Wang, H. and Aggarwal, C. C. (2010). A survey of algorithms for keyword search on graph data. In  
966 *Managing and Mining Graph Data*, pages 249–273. Springer.

- 967 Wu, W. (2013). Proactive Natural Language Search Engine: Tapping into Structured Data on the Web. In  
968 *Proc. of the Joint 2013 EDBT/ICDT Conferences*, pages 143–148. ACM Press, New York, USA.
- 969 Wynholds, L. A., Wallis, J. C., Borgman, C. L., Sands, A., and Traweek, S. (2012). Data, Data Use, and  
970 Scientific Inquiry: Two Case Studies of Data Practices. In *Proc. 12th ACM/IEEE-CS Joint Conference*  
971 *on Digital Libraries (JCDL 2012)*, pages 19–22. ACM Press, New York, USA.
- 972 Yu, J. X., Qin, L., and Chang, L. (2010). Keyword Search in Relational Databases: A Survey. *IEEE Data*  
973 *Eng. Bull.*, 33(1):67–78.
- 974 Zhang, H., He, X., Harrison, T., and Bian, J. (2019). Aero: An Evidence-based Semantic Web Knowledge  
975 Base of Cancer Behavioral Risk Factors. In *Proc. of the 4th International Workshop on Semantics-*  
976 *Powered Data Mining and Analytics co-located with the 18th International Semantic Web Conference*  
977 *(ISWC 2019)*, volume 2427 of *CEUR Workshop Proceedings*, pages 7–11. CEUR-WS.org.