

1 **Integration among Databases and Data Sets to Support Productive** 2 **Nanotechnology: Challenges and Recommendations**

3 Egon L. Willighagen*¹, Sharon Gaheen*², Sandra Karcher*³, Christine Ogilvie Hendren⁴, Marty
4 Fritts², Dennis G. Thomas⁵, Stacey Harper⁶, Mark D. Hoover⁷, Richard L. Marchese Robinson⁸,
5 Karmann C. Mills⁹, John Rumble¹⁰, Nina Jeliazkova¹¹, Friederike Ehrhart¹, Georgia Tsiliki¹²,
6 Axel P. Mustad¹³, Nastassja Lewinski¹⁴ and Chris T. Evelo¹

7 Address: ¹Department of Bioinformatics - BiGCaT, Maastricht University, P.O. Box 616, UNS50
8 Box 19, NL-6200 MD, Maastricht, The Netherlands; ²Leidos Biomedical Research Inc., Fred-
9 erick National Laboratory for Cancer Research, Frederick, MD, 21702, USA; ³Civil and Envi-
10 ronmental Engineering, Carnegie Mellon University, Pittsburgh, PA 15213-3890, USA; Center
11 for the Environmental Implications of NanoTechnology (CEINT) Duke University, Durham, NC
12 27708-0287, USA; ⁴Center for the Environmental Implications of NanoTechnology (CEINT) Duke
13 University, Box 90287, 121 Hudson Hall, Durham, NC 27708-0287, USA; ⁵Biological Sciences
14 Division, Pacific Northwest National Laboratory, Richland, Washington, USA; ⁶Environmental
15 and Molecular Toxicology and School of Chemical, Biological and Environmental Engineering,
16 Oregon State University, Corvallis, OR 97331, USA; ⁷National Institute for Occupational Safety
17 and Health, 1095 Willowdale Road, Morgantown, WV 26505-2888, USA; ⁸School of Chemi-
18 cal and Process Engineering, University of Leeds, Leeds LS2 9JT, UK (current); School of Phar-
19 macy and Biomolecular Sciences, Liverpool John Moores University, James Parsons Building, By-
20 rom Street, Liverpool, L3 3AF, United Kingdom (previous); ⁹RTI International, 3040 Cornwallis
21 Rd., Research Triangle Park, NC 27709, USA; ¹⁰R&R Data Services, 11 Montgomery Avenue,
22 Gaithersburg MD 20877, CODATA-VAMAS Working Group on Nanomaterials, Paris, France;
23 ¹¹IdeaConsult Ltd. 4 A. Kanchev str. Sofia 1000, Bulgaria; ¹²School of Chemical Engineering, Na-
24 tional Technical University of Athens, 9 Heroon Polytechniou Street, Zografou Campus, Athens,
25 15780, Greece; ¹³Nordic Quantum Computing Group AS, Oslo Science Park, P.O. Box 1892 Vika
26 N-0124 Oslo, Norway and ¹⁴Virginia Commonwealth University, Department of Chemical and

27 Life Science Engineering, 601 West Main Street, P.O. Box 843028, Richmond, Virginia 23284-
28 3028

29 Email: Egon L. Willighagen - egon.willighagen@maastrichtuniversity.nl; Sharon Gaheen - ga-
30 heens@mail.nih.gov; Sandra Karcher - SandraKarcher44@gmail.com

31 * Corresponding author

32 **Abstract**

33 This paper is one of a series of articles by the Nanomaterial Data Curation Initiative. Other arti-
34 cles in this series discuss data curation workflows, data completeness and quality, curator respon-
35 sibilities, and metadata. Many groups within the broad nanotechnology field are already develop-
36 ing data repositories and tools driven by their individual organizational goals. Integrating these
37 data across disciplines, and with other non-nanotechnology resources, can support multiple objec-
38 tives by reusing the same information, and can serve as the impetus for novel scientific discoveries
39 through deeper data analyses. This article, framed around the results of a community-based survey
40 of organizations that maintain nanomaterial repositories, discusses current data integration prac-
41 tices in nanoinformatics and in mature fields such as genomics, as well as nanotechnology-specific
42 challenges impacting data integration. Recommendations for achieving integration of existing op-
43 erational nanotechnology resources, as based on results from the community-wide survey, are pre-
44 sented herein. Nanotechnology-specific data integration challenges, if effectively resolved, can
45 foster the application and validation of nanotechnology within and across disciplines.

46 **Keywords**

47 nanotechnology, nanoinformatics, integration, databases, web services

Introduction

Understanding and addressing complexities involved in integrating nanomaterial and non-nanomaterial data resources to further and enable scientific research is a key focus of nanoinformatics [1]. This paper is one in a series of papers focusing on different aspects of nanoinformatics produced from the Nanomaterials Data Curation Initiative (NDCI), which is part of the National Cancer Institute (NCI) Nanotechnology Working Group [2]. Other articles in this series discuss issues such as data curation workflows [3] and data completeness and quality [4]. The focus of this article is on the integration of databases and data sets across nanotechnology and non-nanotechnology resources. The conceptual integration of resources is shown in Figure 1, with databases shown in large boxes, links shown as lines, and some of the database content shown as corner boxes.

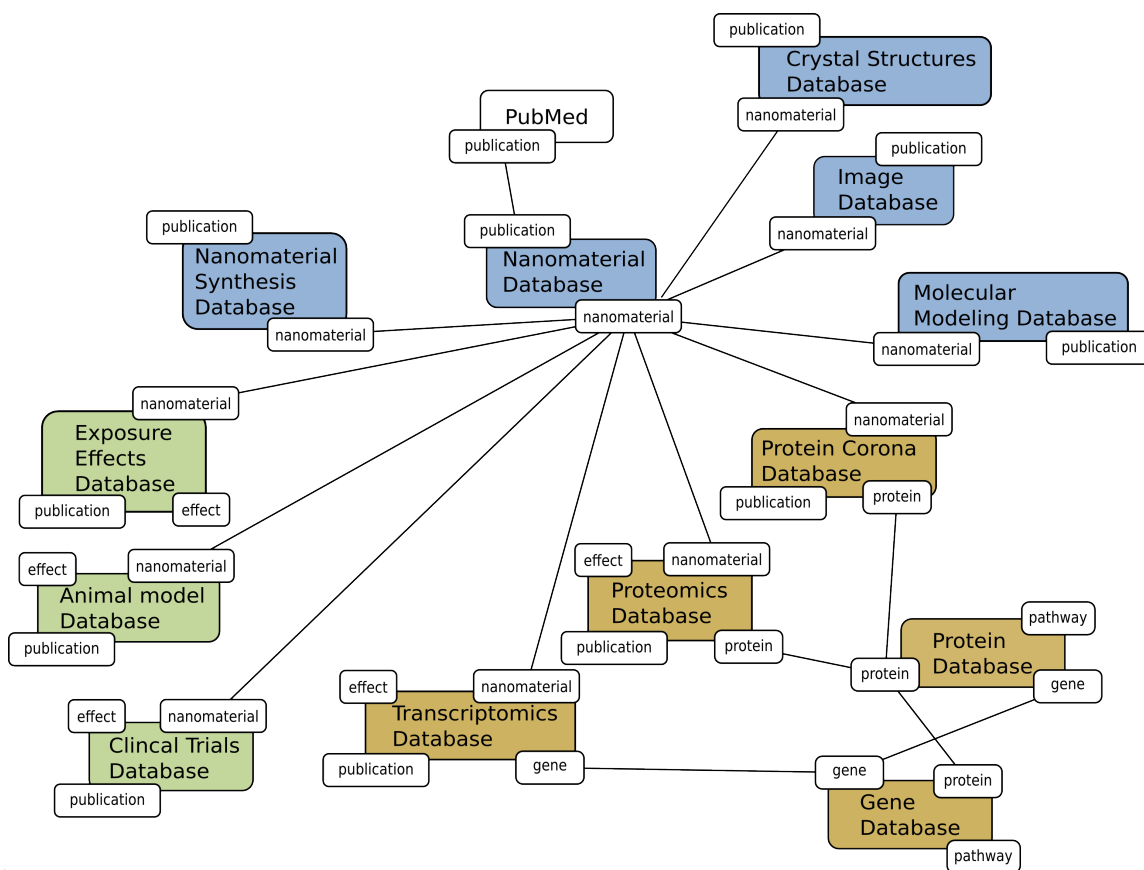


Figure 1: Conceptual Integration of Nanotechnology and Non-Nanotechnology Resources.

Figure 1 shows how nanomaterial data repositories can be integrated with other databases to enable

60 interdisciplinary decision making. All repositories containing nanomaterial information are shown
61 with a nanomaterial corner box and are linked with other repositories by nanomaterial, though they
62 often are not specific to nanomaterials. For example, Gene Expression Omnibus (GEO) and Ar-
63 rayExpress [5] are examples of gene database repositories. Sometimes the genes included in these
64 databases are the focus of studies performed using nanomaterials. The results of those studies may
65 be reported in another database, but the data can be linked using the database content relating to the
66 gene. For example, Figure 1 shows the gene database connecting with the transcriptomics database
67 through information about the gene. It should be noted that the boundaries are not always as clear-
68 cut as indicated in this conceptual diagram, and in reality, there will be many more links than are
69 shown here.

70 The NDCI is currently working to define nanoinformatics and is exploring the role of data inte-
71 gration as an essential component within the field. The following working definition (expanded
72 from the Nanoinformatics 2020 Roadmap [6]) has been proposed: "Nanoinformatics is the sci-
73 ence and practice of determining which information is relevant to meeting the objectives of the
74 nanoscale science and engineering community, and then developing and implementing effective
75 mechanisms for collecting, validating, storing, sharing, analyzing, modeling, and applying the in-
76 formation. Nanoinformatics further involves confirming that appropriate decisions were made and
77 that desired mission outcomes were achieved based on this information. Additional steps in the
78 informatics life cycle including conveying experience to the broader community, contributing to
79 generalized knowledge, and updating standards and training." [7] Successful nanoinformatics en-
80 deavors, including data integration, will apply to all of the steps in the process.

81 The structure of the rest of this paper is as follows. The article begins with an introduction (Sec-
82 tion) that discusses why data integration is important and describes common practices for achiev-
83 ing integration. Using the results of a community-wide stakeholder survey, the current practices
84 for integrating data in nanotechnology are presented (Section), followed by stakeholder identified
85 challenges to integration (Section) and a brief description of integration needs (Section). Stake-

86 holder recommendations are reviewed (Section) and the authors' recommendations presented
87 (Section). The article concludes with a few closing remarks by the authors (Section).

88 **Integration of Databases and Data Sets**

89 **A. Importance and relevance of the integration of databases and data sets to** 90 **the field of nanoinformatics**

91 Nanomaterials [8,9] are becoming ubiquitous in science and technology [10,11]. Biomedical re-
92 searchers are making multifunctional nanomaterials that can be used to diagnose, target, and treat
93 many diseases, especially cancer, looking for ways to increase nanomaterial stability and optimize
94 nanomaterial performance while minimizing potential negative effects [11]. Other researchers are
95 harnessing the same useful properties of nanoscale materials for a host of other applications rang-
96 ing from energy storage to water treatment to improved mechanical strength and flexibility of ad-
97 vanced materials [12].

98 In order to design an optimal nanomaterial and predict how the nanomaterial will behave, re-
99 searchers review numerous publications and query disparate nanomaterial repositories across the
100 biomedical, environmental, health and safety, and materials science disciplines. Where the compo-
101 sition of a nanomaterial is provided in a publication and in repositories, the nomenclature used to
102 describe the base nanomaterial formulation, the material constituents (such as core, coat, shell, and
103 any surface modifiers), and the relationships among components are not standardized and mostly
104 incompletely described. For example, the surface density of "decorator" molecules on carbon nan-
105 otubes may not be provided, resulting in the need for simplifying assumptions when preparing rep-
106 resentative structure files for computational modeling [13].

107 When storing data in a repository, the selection of the storage format and the design of the database
108 structure is often targeted to meet specific objectives. One method used to organize and store data
109 is as nodes and edges, where data are stored in the nodes and relationships are defined by the edges
110 (the lines connecting the nodes). Another method of organizing data is in tables as columns and

111 rows (fields and records), where insight regarding relationships is built into the structure of the ta-
 112 bles.

113 Some organizations have invested heavily in standardizing the content of databases used by their
 114 affiliates. For example, The National Institutes of Health has created an extensive repository of
 115 cancer data standards (caDSR) containing appropriate vocabulary and metadata content to "en-
 116 sure the longevity and agreeability of biomedical research data" (<https://wiki.nci.nih.gov/display/caDSR/caDSR+Wiki>). Other organizations, such as a group of experimental researchers who are
 117 trying to combine their data as part of an integrated study, do not have the funding or expertise to
 118 design and plan for multi-repository standardization. Their repository may be an Excel spreadsheet
 119 that holds their combined data. These two types of data repositories are represented in Figure 2.
 120

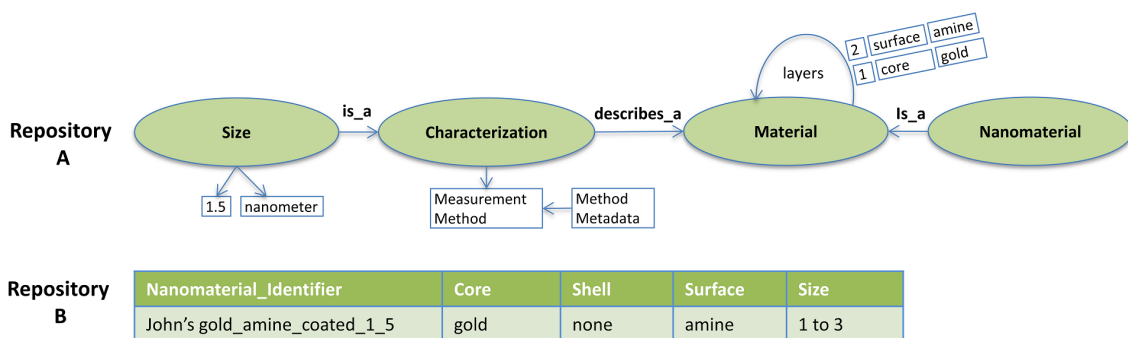


Figure 2: Showing the challenges associated with determining if a nanomaterial is the same when described in different repositories.

121 In the case of repository A, a nanomaterial is composed of materials, in layers. The number of lay-
 122 ers is not limited and the description and order of each layer are attributes of the material. The size
 123 of the nanomaterial is measured using a specific method, and the associated method metadata is
 124 included as attributes of the measurement method. Size is reported as a single value and a unit. In
 125 the case of repository B, a nanomaterial is described by a maximum of three layers, a core, as shell,
 126 and a surface. No method information is provided in the repository and the size is described as a
 127 range. No size unit is given. When data curators of repository B are confronted with the challenge
 128 of integrating data from repository A, they face two fundamental challenges: (1) determining the

129 uniqueness and equivalency of the nanomaterials [14] and (2) assessing the value in incorporating
130 data from repository B into repository A.

131 Physical-chemical characterization information (such as the size, shape, purity, and surface prop-
132 erties) is sometimes included in repositories, but the methods and techniques used to perform the
133 characterization are not always included in sufficient detail or standardized in a way that will allow
134 for cross-study comparison of reported values [15]. Each repository collects and stores informa-
135 tion in support of their organization's needs and goals. Some repositories may include the results
136 of experimental studies focused on biomedical research, whereas others may include geospatial
137 information on the fate of nanomaterials in an environmental system. Some repositories focus ex-
138 clusively on nanomaterials, such as caNanoLab, which houses information related to biomedical nan-
139 otechnology research. Other repositories, such as the Mouse Genome Informatics (MGI) [16], Ar-
140 rayExpress, WikiPathways [17], contain information that is not specifically related to nanomaterial
141 research [18]. These disparate data, when integrated together, may provide additional insights into
142 understanding common endpoints such as nanomaterial toxicity or stability [19].

143 Consider a scenario where a species of mice have been injected with a specific nanomaterial in
144 an *in vivo* laboratory study and the same species were exposed to the same nanomaterial during a
145 mesocosm study. Integrating genomics data with the results of both studies illustrates one of many
146 real world use cases that can benefit from the integration of nanomaterial-specific resources with
147 other relevant resources (e.g. genomics data, clinical trials management systems, chemical repos-
148 itories) that may not necessarily be specific to nanomaterials. Because of the current lack of stan-
149 dardization and integration of resources, researchers must review documentation describing the
150 protocols for storing information in each repository, and sometimes retrieve and review copious
151 publications to determine what is and what is not relevant to their research. This process is time
152 consuming and redundant. The ability to fully integrate repositories across disciplines would al-
153 low rapid association of all relevant experimental results to a specific nanomaterial and could help
154 optimize allocation of resources, for example, integration could enable the prediction of adverse
155 clinical results a priori, allowing resources to go to studies that show a more promising outcome.

156 As outlined above, multidisciplinary fields are particularly demanding on data integration efforts.
157 In all domains, not just nanotechnology, data integration requires a common language (e.g. ontolo-
158 gies) , as well as standards (formal and de facto) for data exchange, communication channels, and
159 identifiers, among others. New technologies have repeatedly changed technical approaches to data
160 integration. While a paradigm based on central data platforms still predominates [20,21], the wider
161 data integration community often uses a more distributed, more-easily scalable cloud [22,23] and
162 other methods, based upon federated search approaches [24,25].

163 A key requirement for an integration effort is a shared system that crosslinks among databases.
164 This can be based on database identifiers. In disciplines close to nanotechnology, efforts such as
165 identifiers.org [26] unify how identifiers are represented, and other systems provide solutions for
166 mapping identifiers from different databases [27-29]. Identifiers, however, typically focus on en-
167 tities studied, such as chemicals, materials, genes, and proteins, but identifiers for cell lines, as-
168 says, and other key entities involved in nanosafety data are less common, though ontologies com-
169 monly provide identifiers for them [30-32]. Moreover, nanomaterials are not as well-defined as
170 small compound chemicals.

171 Linking data enables data integration; by integrating data sets, data comparisons are enabled. Link-
172 ing does not define, of course, which data need to be compared or which data can be connected.
173 Decoupling data integration into two steps, linkage and comparison, allows formalization of a
174 hypothesis into a query. For example, linking two nanomaterial data resources, one containing
175 clinical data and the other embryonic zebrafish toxicity data, by identifying records across both
176 resources as being related to the "same" nanomaterial, allows for a hypothesis (e.g. "toxicity to-
177 wards embryonic zebrafish is of clinical relevance") [33] to be converted into a query (e.g. "report
178 all nanomaterials where high toxicity with respect to embryonic zebrafish corresponds to a high
179 toxicity in a clinical setting, as a fraction of all nanomaterials with both kinds of data") which com-
180 pares data retrieved for two endpoints for the same nanomaterial. This approach becomes increas-
181 ingly powerful if links are made between entities, e.g. nanomaterials, even if they are not identical
182 (the same identifier), but show the same chemical or biological characterization for endpoints of

183 interest, i.e. are functionally equivalent (basically the difference between "the same" and "a close
 184 match"). An example of being "the same" would be two databases with data on a nanomaterial
 185 from a single paper identified with the same label (see Table 1). An example of "a close match"
 186 could be two titanium oxides from the Joint Research Centre Institute for Reference Materials and
 187 Measurements (JRC IRMM) with the same vendor identifiers. While having the same identifier,
 188 they might not be functionally equivalent, depending upon the extent to which the endpoints of in-
 189 terest were affected by aging, etc [34,35].

Table 1: Levels of Equivalence. The equivalence strengths are intended to indicate how data are intended to be combined, and does not specify why it is that it should be linked like that.

Equivalence Strength	Semantic Equivalence	Description	Example
Strong	Web Ontology Language (OWL) sameAs	Two nanomaterials that share the same properties: all properties for one are valid for the other. Moreover, if one nanomaterial is sameAs with others, the others are equally strong (transitivity).	An example would be the same nanomaterial from a journal article for which information is given in two databases.
Moderate	Simple Knowledge Organization System (SKOS) closeMatch	Two nanomaterials are said to be the same for a certain application. This match is never transitive.	An example could be two nanomaterials from the same production batch, in which the application ignores variation.
Weak	SKOS relatedMatch	Two nanomaterials are merely linked together, with an undefined similarity.	This can link two titanium oxide nanomaterials of different sizes.

190 A formalization of this approach in terms of Semantic Web technologies has been recently pro-
 191 posed through the introduction of lenses that allow users to turn on and off such equivalents
 192 based on which links they deem suited for their research question [36,37]. This approach merges
 193 the worlds of ontologies and data, by using Internationalized Resource Identifiers (IRIs), such
 194 as that found in the set of Semantic Web technologies [38,39]. The Open PHACTS project has

195 taken this approach and developed an Identifier Mapping Service (IMS) that links databases us-
196 ing IRI-based identifiers [36]. Services such as identifiers.org and the IMS itself provide routes
197 to convert between alphanumeric identifiers (e.g. CHEBI:33128) and IRI-based identifiers
198 (http://purl.obolibrary.org/obo/CHEBI_33128) as defined in the ChEBI ontology [40]. Once these
199 links are operational, allowing comparison of data for a set of similar or identical materials, the
200 cross-comparison can be used for automated data curation. During curation, automated compar-
201 isons could be enabled to automatically generate warnings that point the user towards other studies
202 reported in other data sources that contradict those being curated. Assuming the linking and sub-
203 sequent steps leading to the generation of such a warning are correct, the linking could allow re-
204 searchers of the earlier study to be automatically notified that new, related data have been added
205 to the database. Data integration for identical (or sufficiently similar) nanomaterials also enables a
206 variety of goals to be achieved that are specific to a particular organization.

207 **B. Influence of organizational purpose and goals on data integration**

208 The approaches taken by an organization or project to gathering and organizing data are governed
209 by the driving scientific questions that need to be answered in order to further its mission. Some
210 examples of use case scenarios that could benefit from multidisciplinary data integration are shown
211 in Figure 3.

212 Data that are measured, the information derived from those data, and the level of detail targeted for
213 inclusion in a resource are all informed by the purpose for which data are being collected. Such
214 purposes include building an authoritative repository of nanomaterial characterizations, parame-
215 terizing models to predict nanomaterial behavior in environmental systems, or improving perfor-
216 mance of materials, medicines, or pesticides. The goals of the individual resource also shape the
217 type of data integration of interest, with each project incentivized to link with other data sets to in-
218 crease the critical mass of data in support of its mission. The vision of the nanoinformatics field is
219 that, beyond achieving individual project goals, the potential exists for broadly-integrated data sets
220 to yield unexpected insights from deeper data mining, generating new hypotheses and knowledge

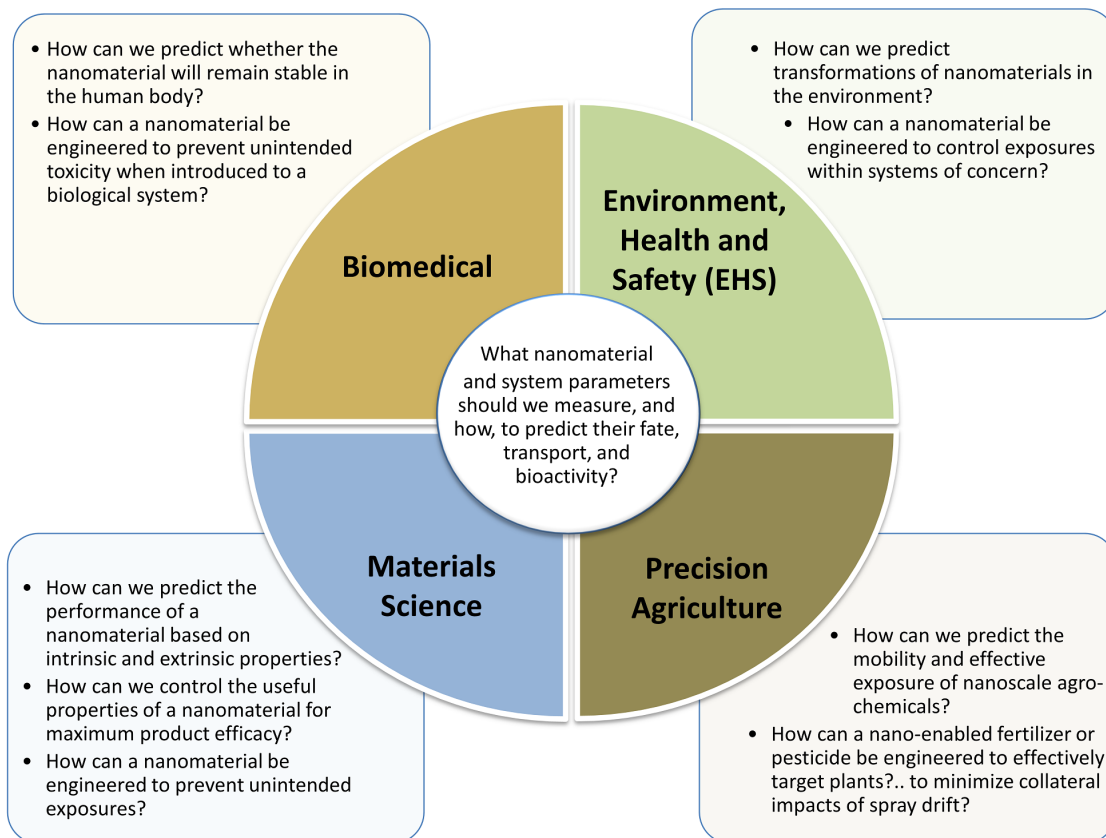


Figure 3: Examples of use cases that can be addressed, and that might mutually benefit through data integration.

221 not anticipated by the originating data resources and benefiting multiple stakeholders. To realize
 222 these secondary benefits of integration, individual projects and disciplines participating in integra-
 223 tion efforts must see improvement in their ability to meet their own objectives. Use cases for data
 224 integration efforts should therefore be selected such that the different driving forces behind their
 225 informatics interests are mutually advanced.

226 As an example, consider the overlap of interests among biomedicine, materials science, precision
 227 agriculture, and environmental, health, and safety (EHS) research as illustrated in Figure 3. Each
 228 field pursues research on its discipline-specific questions. Yet at the intersection of these fields is
 229 a common kernel of questions and answers that would advance each individual research field as
 230 well as open new vistas on a multi-disciplinary basis. Furthermore, by looking across all four dis-
 231 ciplines, data integration potentially positively affects the entire data life-cycle, from experimental
 232 design through data sharing.

233 Such use cases can guide initial pilot projects for nano-specific data integration, recognizing the
234 direct near-term value to participating projects, as well as demonstrating the benefits data sharing
235 brings to measurements outside the domain in which they were made. For example, a biomed-
236 ical nanomaterial data repository integrated with other nanomaterial data resources relevant to
237 biomedicine (e.g. toxicity) and non-nanomaterial data resources (e.g. gene expression and biomed-
238 ical images) would open interesting pathways to finding effective safe disease treatments.
239 Integrating data from different data resources for equivalent nanomaterials supports multiple goals
240 specific to diverse organizations or projects [14]. Using the example provided in Figure 3, under-
241 standing which parameters control stability of a nanomedicine in the human bloodstream could
242 provide insight when predicting nanomaterial dissolution or aggregation in a body of freshwater,
243 transport within a crop field, or efficacy in a material fabrication process. Other examples of poten-
244 tial mutually beneficial integration projects include the following:

- 245 • Calculation of a therapeutic index by integrating data from toxicology and clinical studies
246 (<http://bioportal.bioontology.org/ontologies/NCIT?p=classes&conceptid=http%3A%2F%2Fncicb.nci.nih.gov%2Fxml%2Fowl%2FEVS%2FThesaurus.owl%23C18223>)
- 248 • Development of computational models for predicting nanomaterial effects, using consistent
249 physico-chemical measurements in well-characterized media, based on integrating data from
250 physico-chemical and biological characterizations studies, where common nanomaterials
251 might be established based upon product names and/or batch identifiers [41-43]
- 252 • Predictions of nanomaterial transformations and effective exposures made through integrat-
253 ing fate and transport data across a variety of systems of interest (bloodstream, aerosolized
254 irrigation stream, polymeric matrix, water column, sediments) with implications for estab-
255 lishing treatment efficacy, product performance, and collateral toxicity [44].

256 The power of these and other use cases will hopefully pique sufficient curiosity to start significant
257 integration projects. Steps toward integration might then begin with developing an understanding
258 of the respective minimal data standards. Understanding these resource-specific data requirements

259 similarly requires an appreciation for the driving purpose of the resource. It has been suggested
260 that a method for capturing minimum data requirements by resource (e.g. MIAME for microar-
261 ray data [45]) supports the resource categorization and provides a greater understanding of data
262 requirements. One possible candidate for such a metadata resource is the BioSharing platform
263 (<https://www.biosharing.org/>) [46]. Further discussion of data and metadata requirements and their
264 explicit documentation via minimum information checklists is presented in an earlier article in the
265 NDCI series [4]

266 **C. Established methods for the integration of databases and data sets**

267 A variety of different approaches have been developed to integrate data, supported by a variety
268 of different kinds of technologies, ranging from manual integration within an Excel spreadsheet
269 (e.g. based on "VLOOKUP" matching of identifiers) to a federated search architecture based
270 on semantic web technologies (<https://www.w3.org/standards/semanticweb/>) [24,25]. The fo-
271 cus of this paper is upon approaches that best facilitate the retrieval of integrated data via au-
272 tomated queries (e.g. the data query languages SQL or SPARQL [47]); hence, these latter ap-
273 proaches will frame the following discussion. Nonetheless, it is important to note that, given the
274 preference of many scientists for data collection in Excel, tools that allow for automated inte-
275 gration of manually prepared Excel data sets into queryable databases are of considerable value
276 (<https://github.com/enanomapper/nmdataparser>) [48].

277 The extremes of the spectrum with regard to selecting an architecture that will support data integra-
278 tion through automatic querying are data warehousing and federated query [49].

- 279 • The data warehousing approach involves loading the content of different data resources into
280 the same physical database. Subsequently the "warehouse" database can be queried, which
281 involves querying all loaded data resources concurrently, with results presented to the user.
- 282 • Federated querying is implemented by sending queries to the different data resources at their
283 original locations and presenting the results to the user in one unified view as soon as they
284 are received.

285 The technology for accessing the data resources may be the same for both approaches, e.g. the data
286 warehouse approach may use extract-transform-load (ETL) procedures, connecting to external data
287 resources via web services and loading the results into the warehouse, while federated querying
288 may use wrappers for accessing several distinct databases residing on the same machine and com-
289 bine results only when presenting them to the user. Hence, a web service is a method for access-
290 ing the data, but its use does not imply anything about the data integration paradigm after data re-
291 trieval.

292 The data warehouse paradigm accomplishes the integration by transforming all the data resources
293 into a physical schema (i.e. tables and relationships for relational databases, or XML schema, etc.).

294 The federated query approach relies on a "mediated schema", i.e., a virtual schema, embedded in
295 the application, which does not store any data, but presents to the user a unified view of the do-
296 main and allows queries to be specified. The integration itself relies on how the different attributes
297 of the mediated schema match the attributes of the sources, and if the grouping of the attributes
298 corresponds to similar groupings of attributes in the data resources. This is known as "semantic
299 mapping" and is the hardest task within the integration. Regardless of the integration approach,
300 all methods require entity matching (linking associated information based on database content)
301 or mapping (virtually altering the schema of one database so that its content can be queried with
302 data from a database with a different schema). Mapping is typically performed using transfor-
303 mation procedures. There may not exist a simple one-to-one mapping between the final schema
304 and the original data resources. For example, suppose percentage cumulative mortality data were
305 required from the Nanomaterial-Biological Interactions (NBI) Knowledgebase data resource
306 (<http://nbi.oregonstate.edu/>), in order to include those data with embryonic zebrafish toxicity data
307 curated from the literature [50] in a common data warehouse. Since the NBI knowledge base [51]
308 provides mortality data in terms of the raw numbers of dead/live organisms at 24 hours post-
309 fertilization, and the additional number of zebrafish that were observed to be dead at 120 hours
310 post-fertilization, determining the total number of zebrafish observed to be dead at 120 hours post-

311 fertilization would require mathematical processing before being returned to the user in a schema
312 requiring a field "percentage cumulative mortality" to be populated.

313 Developing mapping algorithms has traditionally be done manually, however, active research is
314 producing tools for automatic schema mapping and record linkage by deterministic, probabilistic
315 and machine learning methods [52]. In the case of unstructured data resources, e.g. text, the work-
316 flow first performs data extraction and entity recognition and then proceeds with the mapping.

317 Between the two extremes of data warehousing and federated query, many hybrid architectures
318 exist combining elements of both pure data warehousing and federated querying. The choice of
319 integration architecture depends on:

320 1. How the entities can/will be matched across databases.

321 2. How the query results will be integrated.

322 Federated searching can be illustrated with an application to query several online chemical
323 databases for small molecule chemical compound properties via an Application Programming
324 Interface (API) and presenting integrated results on a single web page. Here the entities are
325 the chemical structures, and the IUPAC International Chemical Identifier (InChI, [http://www.
326 inchi-trust.org/](http://www.inchi-trust.org/)) [53] can be used as a uniform identifier across databases. The matching rule is
327 "if the search results returned include one and the same Standard InChI, then the results are for
328 the same compound". A data warehouse implementation would use the API to retrieve the results,
329 store them into a database, and then allow the user to query the database. A federated approach
330 would use the API to retrieve the results and present them in a unified format to the user. Although
331 this example may seem straightforward, there are number of complexities that must be considered
332 when matching based on an InChI. For example, small molecule chemicals which may be consid-
333 ered the same, yet correspond to rapidly interconverting structures, may still fail to match based
334 upon InChIs. Whilst InChIs are designed to be invariant to different ways of representing chemi-
335 cals based on small molecular structures, including taking into account tautomeric forms which are
336 expected to rapidly equilibrate, they cannot account for all differences in chemical structure which

337 may readily interconvert in practice - such as differences in protonation state (e.g. salicylic acid
338 will exist in dynamic equilibrium with its deprotonated form under physiological conditions) or
339 between open-and-closed ring forms, which can equilibrate for sugars in solution. If non-standard
340 InChIs are used, the situation is further complicated [53]. In spite of the challenges discussed here,
341 integration of small molecule chemical databases based on matching their Standard InChIs is cur-
342 rently viewed as best practice and may be combined with other software tools to enforce further
343 standardization of chemical structures that may facilitate desired matches [54].

344 Extending this approach to more complicated structures, e.g. proteins and genes, would require
345 expanding the queries to handle all possible synonyms used by different databases.

346 Establishing a common API for a given type of resource facilitates integration because it alleviates
347 the need of schema matching. Essentially, the API defines a common schema and if all resources
348 of the same kind are compliant with the API, the main hurdle of semantic mapping is lifted. An ex-
349 ample implementation of this approach in the genomics field is the Global Alliance for Genomics
350 and Health (GA4GH) Data Working Group (<http://ga4gh.org/#/>), which is establishing common
351 web services in support of genomic data integration and exchange. Example web services using the
352 Representational State Transfer (REST) framework [55] are provided with query requests and re-
353 sponses formatted using the JavaScript Object Notation (JSON). The common web services allow
354 the genomics community to exchange reads, variants, and reference information, provided all data
355 resources follow the API specification.

356 The implementation of a central data warehouse or repository that aggregates data from several re-
357 sources requires extract, transform, and load (ETL) processes to assist in aggregating and trans-
358 forming the data based on matching rules. Data are typically transformed into a common data
359 model (e.g. relational database or a triple store); examples of this approach are PubChem and
360 ChEMBL databases. The Open PHACTS project provides a common API to a variety of phar-
361 macological data sets. However, it does not normalize to a single data model, but addresses the
362 non-uniformity at the API level [20]. The European Bioinformatics Institute Resource Description
363 Framework (EBI-RDF) platform uses a different approach, maintaining multiple RDF repositories

364 for different resources and allowing federated searching across all of them [23]. It is mandatory for
365 all of the entities in the EBI-RDF platform to be assigned equivalent identifiers via identifiers.org
366 service, which is essentially implementing the mapping between the distributed resources.
367 The Syngenta federated search system [25] is an example of addressing the challenge of integrating
368 internal company data with public life science databases. The system has moved from data ware-
369 housing (even if that offers faster reporting) towards federated search technologies. The architec-
370 ture includes several internal relational database repositories, translated into RDF dynamically via
371 D2RQ (<http://d2rq.org/>) [56], and providing adapters in order to combine all internal and external
372 data resources into a distributed SPARQL endpoint. The implementation of this federated architec-
373 ture for data integration was found to offer clear benefits to Syngenta's chemists and biochemists.

374 **Current practice for data integration in the nanotechnology field:**

375 **perspectives of key stakeholders**

376 To understand the current practices in data integration and to identify challenges and offer recom-
377 mendations, several organizations that maintain nanomaterial repositories were asked to respond
378 to a questionnaire on data integration. The goal was to assist in defining and initiating integra-
379 tion and exchange of data resources across nanomaterial data repositories and with other non-
380 nanotechnology data resources. Questions included current and recommended functionality and
381 web services enabling data integration and exchange as well as perceived challenges associated
382 with integrating primary experimental data sets, or data sets curated from the literature, with exist-
383 ing nanomaterial and non-nanomaterial data repositories. The following sections provide details on
384 the organizations who participated in the survey along with summarized results of their feedback.
385 Information on each nanomaterial resource is provided in Table 2.

Table 2: Nanomaterial Resources Responding to Data Integration Questionnaire

Nanotechnology Resource	Resource Description	Integration Capabilities
<p>caNanoLab - caNanoLab Data Portal (https://cananolab.nci.nih.gov/) caNanoLab Wiki (https://wiki.nci.nih.gov/display/caNanoLab/caNanoLab+Wiki+Home+Page) caNanoLab Data Dictionary Resources: caNanoLab Glossary (https://wiki.nci.nih.gov/display/caNanoLab/caNanoLab+Glossary) NCI Thesaurus (https://ncit.nci.nih.gov/ncitbrowser/pages/home.jsf?version=15.05d) NCI caDSR (https://cdebrowser.nci.nih.gov/CDEBrowser/) Design Document with Domain Model caNanoLab Code Repository (https://github.com/NCIP/canolab)</p>	<p>caNanoLab is a data sharing portal designed to facilitate information sharing across the international biomedical nanotechnology research community to expedite and validate the use of nanotechnology in biomedicine. caNanoLab provides support for the annotation of nanomaterials with characterizations resulting from physico-chemical, in vitro and in vivo assays and the sharing of these characterizations and associated nanotechnology protocols in a secure fashion.</p>	<p>Provides REST-based web services supporting general sample search and retrieval of sample composition and characterizations by sample ID. Supports retrieval of samples associated with a publication. Integrates with ScienceDirect publications through an Elsevier bi-directional link and uses the PubMed and PubChem interfaces.</p>

<p>CEINT CEINT Wiki (http://www.ceint.duke.edu/)</p>	<p>CEINT is a center-wide effort focused on exploring the potential impact of exposure to nanomaterials on ecological and biological systems. The center is funded by the National Science Foundation and the US Environmental Protection Agency, and brings together researchers from several universities, NIST, the EPA, as well as other key domestic and international partners. CEINT supports fundamental research regarding the behavior of nanomaterials in laboratory studies and also in complex ecosystems. One of the goals of the center is to develop a web-based risk assessment tool that can be used to elucidate the potential risk associated with the release of nanomaterials into the environment.</p>	<p>Integration within the CEINT NIKC resource is achieved by custom API development for each collaborative project with targeted data sets.</p>
<p>Center for Safety of Substances and Products, National Institute for Public Health and the Environment (CSSP/NIHE) Netherlands Center Information (http://www.rivm.nl/en/About_RIVM/Organisation/Centres/Centre_for_Safety_of_Substances_and_Products) Software Model for Estimated Exposure from Consumer Products (http://www.rivm.nl/en/Topics/C/ConsExpo) Nanotool for Spray Products (http://www.rivm.nl/en/Topics/C/ConsExpo/Nano_tool)</p>	<p>The CSSP NIHE provides a database on ecotoxicity data focusing on nanoparticles in consumer products. The database provides a repository for modeling purposes (QSAR).</p>	<p>Does not provide any web services. In case of gathering/uploading toxicity data, the OCHEM database is commonly used. The database also allows for modeling and selection of descriptors.</p>

<p>DECHEMA http://www.dechema.de/en/ Nano-safety Wiki (http://www.nanora.eu/nano-safety)</p>	<p>DECHEMA is a network of experts in chemical engineering and biotechnology. DECHEMA supports several projects applicable to nanotechnology such as the DaNa project and the NANORA project [57]. DaNa is a Knowledge base of applied nanomaterials on health and environment. The NANORA project provides web facilities supporting the Nano Region Alliance, an alliance that facilitates market entrance for nanotechnology subject matter experts.</p>	<p>The DaNa project has been providing the web service for the NANORA project to implement the Danavis Database on the NANORA website based on JSON as data exchange format.</p>
<p>eNanoMapper Ontology http://bioportal.bioontology.org/ontologies/ENM Database https://apps.ideaconsult.net/enanomapper/ Search https://search.data.enanomapper.net Modeling http://enanomapper.net/modeling</p>	<p>eNanoMapper is a European FP7 project of eight research and industry institutes. The aim is to improve data integration and to support safe-by-design development by building up a nanosafety ontology, a database and provide tools for use of this data (e.g. modeling approaches).</p>	<p>There is a REST-based API and nanomaterials have URIs allowing a linked data approach. External databases can be indexed by uploading, for example, nanomaterial characterization or via search integration.</p>

<p>Nanomaterial Registry Websites: http://www.nanomaterialregistry.org Partner Portal at nanoHUB</p>	<p>The Nanomaterial Registry is a publicly-available database of nanomaterial characterization and biological/environmental interaction data. Data in the Registry are curated from niche databases, literature, catalogs, and reports by trained scientists. Data are curated based on a set of minimal information about nanomaterials. The data of the Registry are also available on the Portal at nanoHUB, where predictive modelers can find the data in a format that is easy for them to use.</p>	<p>Integration with the Registry is achieved on a case by case basis. Future development will include a JSON interface for analysis tools and data submission templates.</p>
<p>Nanoparticle Information Library http://nanoparticlelibrary.net/</p>	<p>The NIL is a prototype searchable database of nanoparticle properties and associated health and safety information designed to help occupational health professionals, industrial users, worker groups, and researchers organize and share information on nanomaterials, including their health and safety-associated properties.</p>	<p>Integration with the NIL is achieved on a case by case basis.</p>

386 **A. Stakeholder demographics**

387 Stakeholders who participated in the survey ranged from nanomaterial resources that have exten-
388 sive experience in integrating databases and data sets to those with limited data integration experi-
389 ence whose focus was primarily on repository development (Table 2). The diverse levels of integra-
390 tion capabilities provide insight into the challenges that need to be addressed in order to integrate
391 across nanomaterial repositories and with other non-nanotechnology resources.

392 **B. Stakeholder experience in nanomaterial data integration**

393 The surveyed nanomaterial data resources exhibited a variety of experience in data integration in-
394 cluding integrating primary data sets and web services supporting data integration. Stakeholders
395 were asked for information on existing resource functionality supporting data integration including
396 data standards, controlled vocabulary, and common identifiers. They were also asked to identify
397 available web services supporting cross-nanomaterial resource exchange and current efforts sup-
398 porting integration with non-nanotechnology resources.

399 **Uploading / Downloading Data Sets**

400 When using a data warehousing architecture, the ability to upload and download data sets is an ini-
401 tial step towards integration as support for this feature requires the identification of data formats
402 and representation of common data elements. Federated approaches may not require the actual
403 movement of the data, but also requires identification of data formats and common data elements.
404 Stakeholders responded to questions relating to integration of primary data sets, including services
405 available in-house or services that are publicly available (Table 3). These stakeholder experiences
406 provide insights into the level of readiness the nanotechnology community has achieved with re-
407 gards to integrating databases and data sets.

408 **Web Services Supporting Data Exchange**

409 The missions of the stakeholder groups are highly diverse, with web services being of high priority
410 for some and not for others. The data exchange capabilities of each resource, as provided by each

411 stakeholder, are summarized in Table 3, and capabilities relating specifically to web services are
412 described in the following section.

Table 3: Summary of Stakeholder Responses to Upload, Download, and Mapping Questions: Does the nanomaterial data resource provide the following?

	Does the nanomaterial resource provide the following?				
Nanomaterial data resource	Uploading, downloading, or mapping	Definitions of the database fields	Controlled vocabularies, taxonomies and/or ontologies	Nanomaterial identifier uniqueness	Integration with any non-nanotechnology resources
caNanoLab	web-based forms for uploading and downloading nanomaterial composition, characterizations, publications and protocols	extensive documentation is available ^a	uses NPO and the NCI Thesaurus (http://ncit.nci.nih.gov/ncitbrowser/pages/home.jsf?version=15.05d)	Uses a pattern containing source information and a numeric identifier resulting in a unique identifier. The pattern for the sample name is: abbreviation(s) of institution names, name of the first author (without middle name), custom abbreviation of journal title, year of publication, and sample sequential number, e.g. SNL_UNM-CAshleyACSNano2012-01.	caNanoLab integrates loosely with six non-nano resources ^b .

CEINT	mapping from NBI data set	not yet	uses ontologies such as MO, NPO, UO, and ChEBI	nanomaterial associated to data source and assigned a unique identifier	not currently
CSSP/NIPHE, Netherlands	commonly use the OCHEM database for uploading toxicity data	provides a list a fields available for storing toxicity data	uses field headings as a means of controlling vocabulary	identifier assigned based on particle core composition	no
DECHEMA	no	relational model documented in Kimmig et al. [58] and Atli et al. [59]	uses the scientific wording for materials and nanomaterials, toxicology, biology ^c	not a central issue of the DECHEMA work	no
eNanoMapper	extends the OpenTox platform which has the means to download and upload data	overview of the data model documented in Hastings et al. [40]	uses the eNanoMapper ontology (composed of NPO, ChEBI, BFO, IAO, CHEMINF and others)	uses an IUC substance UUID ^d	not currently
Nanomaterial Registry	export for physico-chemical characterization	Nanomaterial Registry glossary (https://www.nanomaterialregistry.org/resources/Glossary.aspx)	uses a controlled vocabulary ^e	uses unique numeric IDs ^f	not currently
Nanoparticle Information Library	Accomplished on a case-by-case basis	Provided as drop-down lists of available fields	Uses the NPO as well as user-specified terms	Unique NIL entry numbers are assigned	The NIL integrates directly with data resources on hazardous materials [60] ^g .

413 ^aThe caNanoLab Design document (<https://github.com/NCIP/cananolab/tree/master/docs/design>) includes the object model which represents class names and at-
414 tributes associated with the data model. All class names and attributes are maintained in the NCI caDSR (<https://cdebrowser.nci.nih.gov/CDEBrowser/>). Con-
415 cepts are defined in the NCI Thesaurus (<http://ncit.nci.nih.gov/ncitbrowser/pages/home.jsf?version=15.05d>). caNanoLab also provides a user-friendly glossary
416 (<https://wiki.nci.nih.gov/display/caNanoLab/caNanoLab+Glossary>).

417 ^bcaNanoLab integrates with PubMed and ScienceDirect for access to publications, Elsevier for linking caNanoLab data to publications, PubChem for chemical
418 information, The Collaboratory for Structural Nanobiology - CSN (http://uqbar.ncifcrf.gov/Advanced_Structure_Analysis/HOME.html) for displaying 3D models
419 of specific nanomaterials, and Nanotechnology Characterization Laboratory (NCL, http://ncl.cancer.gov/working_assay-cascade.asp) assay cascade and JoVE
420 (<http://www.jove.com/>) for nanotechnology protocols.

421 ^cDECHEMA has a very diverse target group ranging from interested laymen, stakeholders to other scientists; wording is adjusted in order to tell a comprehensive
422 story without confusing the laymen on the one and hand and not losing the scientific correctness.

26

423 ^deNanoMapper is based on semantic web technologies including dereferenceable Internationalized Resource Identifiers (IRIs) and the Resource Description
424 Framework (RDF). The substance UUID does not reflect the uniqueness of the material structure, but is an identifier of the material in the database. The substances
425 (materials) are described with their composition (e.g. core, shell, and functionalization) and are linked to the chemical structures of their components. These can be
426 used to decide if the nanomaterials are the same or similar.

427 ^eThe NPO has been mapped to the Nanomaterial Registry and it was determined that a little over 80 terms used by the Registry are not yet part of the breadth of
428 the NPO.

429 ^fIt is the intent of the Nanomaterial Registry not to judge equivalence between any two nanomaterials from different data resources, as the characterization results
430 can be wildly different based on sample medium and characterization protocol.

431 ^gThe NIL integrates with the NIOSH Pocket Guide to Chemical Hazards (NPG, <http://www.cdc.gov/niosh/npg/>) and with the Registry of Toxic Effects of
432 Chemical Substances (RTECS, <http://www.cdc.gov/niosh/rtecs>). The current hosting, administration, and maintenance of the NIL web resource outside of the
433 CDC/NIOSH website is being conducted by Oregon State University in conjunction with its program to characterize nanomaterials.

434 **caNanoLab Web Services**

435 caNanoLab implements an internal and external API leveraging REST (see Table 4). The internal
436 API retrieves web forms in JSON format, while the external API retrieves web forms in HTML
437 format. caNanoLab exposes web services that retrieve publicly available information. All other
438 web services are used internally and are not exposed. caNanoLab does not publish documenta-
439 tion on web services other than The caNanoLab Design document which documents the system
440 architecture and object model. Internal web services are based on method calls on object model
441 attributes. Other NCI projects supporting genomics use Apiary for documenting web services.
442 caNanoLab uses the PubMed API to retrieve publications and interfaces with PubChem to retrieve
443 information on chemicals associated with nanomaterial composing elements.

Table 4: Web Services provided by caNanoLab (<https://cananolab.nci.nih.gov/caNanoLab/#/>)

Search Type	Possible Search Criteria	Notes and Links
protocol	protocol name	https://cananolab.nci.nih.gov/caNanoLab/#/searchProtocol
sample	specific sample, composition, and/or characterization	https://cananolab.nci.nih.gov/caNanoLab/#/advancedSampleSearch . Returns sample information by sample ID.
publication	sample name. nanomaterial characteristics	https://cananolab.nci.nih.gov/caNanoLab/#/searchPublication Retrieves publication information and associated samples by PubMed ID or DOI.

444 **CEINT Web Services**

445 CEINT does not currently provide web services for data set sharing; however, CEINT does pro-
446 vide a web-enabled service for use by CEINT members that allows them to connect with other re-
447 searchers who identify as working on the same research questions, with the same materials, and
448 with the same methods. This service facilitates Center-wide data integration through direct up-
449 stream collaboration, even in the absence of prescribed data templates that would support more
450 automated integration. CEINT uses web services provided by others, included the Nanomaterial
451 Registry, the Integrated Taxonomic Information System, Ontobee, caNanoLab, USDA Geospatial
452 Data Gateway, and the Project on Emerging Nanotechnologies.

453 **CSSP/NIPHE, Netherlands Web Services**

454 The Center for Safety of Substances and Products, National Institute for Public Health and the En-
455 vironment, Netherlands does not offer web services; however, the OCHEM database is publicly
456 available.

457 **DECHEMA Web Services**

458 DECHEMA does not provide any web services per se for the DaNa project. In the case of the
459 NANORA project, a web service was specifically created, together with an interface to imple-
460 ment the DaNaVis database on the NANORA website using JSON as the data exchange format.
461 The backend web services and customized interface for the NANORA website are not publicly
462 available but the frontend user interface is freely accessible. There is no publicly available docu-
463 mentation for the web service for the NANORA project. DECHEMA uses a content-management
464 system for the DaNa website (Joomla + several plug-ins, bootstrap framework). The DaNa website
465 is accessible for everyone without any usage restrictions. The DaNaVis database and tools use a
466 Django-framework (Python as the programming language), REST API- and JSON-based data in-
467 terchange between client and application server, client-side JavaScript widget. More details on the
468 database and tool design have been published [58,59]. DECHEMA does not use any web services
469 provided by other organizations

470 **eNanoMapper Web Services**

471 eNanoMapper provides web services based on the OpenTox API. eNanoMapper inherits and,
472 where needed, extends the machine readable API. The supported return formats include JSON ,
473 JSON-LD and RDF/XML, CSV, XLSX. Methods exist for a number of entity types, including sub-
474 stances, which is how eNanoMapper models a nanomaterial. The API is REST-like. eNanoMapper
475 separates the API design from the server implementation; AMBIT is one of the reference imple-
476 mentations of eNanoMapper services [61], and on the server-side uses Apache's Tomcat. The API
477 implements user authentication and authorization. This means that an eNanoMapper instance (it is
478 a platform rather than a single system), allows for both public data and confidential data that can
479 be shared with only a selected group of researchers. The example <http://data.enanomapper.org/>
480 instance currently hosts several public data sets, available under an Open Data license or waiver.

481 The eNanoMapper server currently does not use other web services, besides being able to re-
482 trieve chemical structures from public databases (e.g. PubChem). However, this may change when
483 eNanoMapper moves towards a more distributed platform later in the project.

484 The full details of the eNanoMapper API, including a description of the computational services
485 implementation (which uses and integrates a variety of technologies and also reads and writes
486 from/to data services) are published [62]. Interactive API documentation is available online (<http://enanomapper.github.io/API/>). A webinar using the API to visualize data in web pages is available
487 on YouTube (<https://www.youtube.com/watch?v=quy7G2mZ0gk>), and a complete list of models
488 that can be used for prediction can be found in the Swagger documentation (<http://app.jaqpot.org:8080/jaqpot/swagger/>).

491 **Nanomaterial Registry**

492 The Nanomaterial Registry does not currently have data exchange web services other than the
493 export tools described in Table 3. However, a JSON interface is in development for the connec-
494 tion with data analysis tools. The Registry website does provide a web service search tool that al-
495 lows for keyword and specific measurement values to be searched, as well as allowing the user to
496 browse nanomaterials by a variety of characteristics. Nanomaterial Registry data are also batch ex-
497 ported to a portal at nanoHUB, where users can interact with and download the data in different
498 ways.

499 **Nanoparticle Information Library**

500 The Nanoparticle Information Library website is publicly accessible to everyone with the request
501 that any use of the data be attributed to the primary source associated with the data entry. Online
502 search capabilities within the NIL are based on attributes of nanomaterial structure, elemental com-
503 position, method of synthesis, and nanomaterial size-related features including primary particle
504 diameter, agglomerate diameter, and specific surface area. Weblinks to the primary data and to the
505 principle investigators who have provided data to the NIL are included.

506 **Stakeholder identified data integration challenges**

507 Stakeholders identified several technical and operational challenges impacting current data inte-
508 gration efforts, as shown in Figure 4. These challenges, if not addressed, will continue to plague
509 the nanotechnology informatics community and greatly hinder scientific discoveries. Each are dis-
510 cussed in greater detail below.

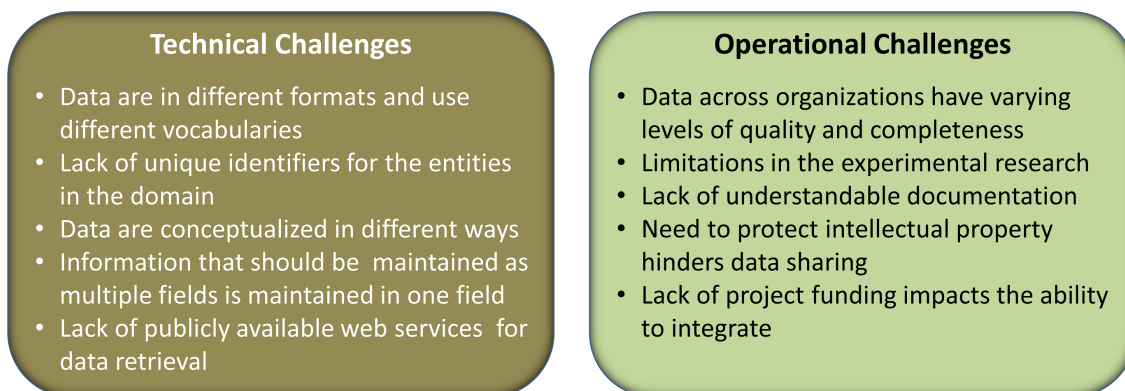


Figure 4: Technical and operational challenges impacting data integration.

511 **Data are in different formats and use different (or no) common vocabularies** 512 **or ontologies**

513 The primary challenge in achieving data integration, selection, and aggregation in the nanotechnol-
514 ogy domain is the diversity of ways in which nanomaterial information is represented across data
515 repositories and the lack of standardization to a common model that represents nanomaterial enti-
516 ties, their attributes, and relationships. These issues include multiple meanings for the same word
517 (or abbreviation) and different words (or abbreviations) having the same meaning. For example,
518 cytotoxicity can have different specific meanings when different bioassays were used to measure it.
519 Similarly, examples of synonyms with the same meaning are abundant too, and include for exam-
520 ple, ZnO, zinc oxide, and nano-zinc oxide. Developing appropriate ontologies, including resolution
521 of terminology conflicts, to address the nuances of nanotechnology research are an important key
522 to achieving integration.

523 **Lack of unique identifiers for the entities in the domain**

524 Certain difficult aspects of data integration remain challenging regardless of the specific domain,
525 including deciding when entities (e.g. nanomaterials, cells, samples, people, etc.) in different
526 data contexts should be mapped as "the same" or "different" e.g. if their names have narrower or
527 broader meanings. An example would be the difference of titanium dioxide (NPO_1485) and tita-
528 nium oxide nanoparticle (NPO_1486) made by the NanoParticle Ontology. This difficulty impacts
529 the ability to perform cross material comparisons. Other fields have introduced naming conven-
530 tions for generating unique identifiers based on metadata; however the different metadata used
531 across studies made this challenging. As such, other fields such as genomics are moving forward
532 with generating a Universal Unique Identifier (UUID) for entities not based on metadata associated
533 with the UUID to support queries in support of entity comparisons. In the context of nanomaterial
534 data resource integration, metadata might include the results of physico-chemical characteriza-
535 tion required to establish whether the nanomaterials are "the same" or "sufficiently similar" to be
536 matched during data integration. However, the question of which physico-chemical properties need
537 to match [15], not to mention complexities associated with different measurement techniques and
538 experimental protocols, make uniquely identifying and matching nanomaterials a significant sci-
539 entific challenge. Further discussion of metadata (including batch identifiers) that could support
540 unique identification and matching of nanomaterial database records is provided in an earlier arti-
541 cle in the NDCI series [4].

542 **"Data" are conceptualized in different ways**

543 There has been a trend away from establishing a fixed hierarchy between database elements; a
544 trend that, in some regards, adds to the data mining challenge. Sometimes expert knowledge is
545 built into the establishment of hierarchical relationships, and that knowledge can be extracted when
546 mining a database to assure that data are appropriately aggregated when performing statistical anal-
547 yses. Often times, databases are designed to support searching, but not specifically to support min-
548 ing. In these types of database repositories, measurements are sometimes duplicated so that they

549 can be pulled in many kinds of searches. Duplicate measurements, if not handled correctly during
550 analysis, can lead to bias in statistical computations.

551 **Information that should be maintained as multiple fields is maintained in one** 552 **field**

553 Integration of data can be hampered by differences in data granularity. A common issue is that in-
554 formation in one repository may be stored in one "field", but be split into multiple "fields" in an-
555 other repository. Additionally, in some repositories, numerical data are stored without a separate
556 unit "field". For example, some repositories use a field name such as "Concentration" and expect
557 the user to know that the result should always be in a specific unit, such as "mg/l". In other cases,
558 a measured result is combined with a unit and stored together in the same field (e.g. 7 mg/l), or in-
559 clude a range of values in one field (e.g. 7-10 mg/l).

560 **Lack of publicly available web services for data retrieval**

561 Integration is often hindered by the lack of publicly available web services supporting data re-
562 trieval. Additionally, even when data services are provided, open frameworks such as REST are
563 not leveraged to ease development of integration touchpoints [55].

564 **Data across organizations has varying levels of quality and completeness**

565 Finding data that are sufficiently complete and of acceptable quality is a key challenge for nanoin-
566 formatics. At times data from external repositories are not integrated with local systems due to
567 concerns regarding the quality and completeness of those data. For example, a local knowledge
568 base can implement a screening procedure that carefully selects high quality data from the sci-
569 entific literature; data from publications not meeting the specific quality criteria are deemed un-
570 suitable and are not curated into the knowledge base. When evaluating external data for inclusion
571 in the knowledge base, if they do not come with an indicator or ranking of the reliability of those

572 data, and if the ranking is not in line with the screening procedure used by the curators, it is diffi-
573 cult to determine if and how those data should be incorporated.

574 Lack of data completeness also poses a challenge to data integration as it is often difficult to obtain
575 the necessary information to support comparison (a pre-requisite for matching and data integra-
576 tion) between material records in different databases. For example, when obtaining information on
577 physico-chemical characterization, it is important to have information on the chemical composition
578 of the particles, such as the presence/absence of coatings, and if the particle has been transformed.
579 In addition, lack of complete metadata for associated biological tests may be considered to affect
580 the clarity, hence quality, of results [63] and preclude an assessment of whether two sets of results
581 were generated under sufficiently similar conditions to allow them to be meaningfully integrated in
582 support of analysis, for example, the relationships between material characteristics and biological
583 effects. It is also critical to have information on the media properties that might affect the result of
584 (toxicity) testing as well as standardized methods of collecting the information. A lack of proper
585 particle characterization is a key problem [64], and the consequence is that often a database con-
586 tains more blank fields (no information) than actual data. This lack of high quality and complete
587 data sets discourages integration.

588 A thorough discussion of the challenges associated with assessing the completeness and quality of
589 nanomaterial data was presented in an earlier paper in the NDCI series [4].

590 **Limitations in the experimental research**

591 There are limitations in the experimental research process, such as biological variance, uniform
592 characterization, and technological and methodological constraints. One major challenge related
593 to data quality and completeness is defining the minimum data requirements for integration. The
594 continuing evolution of knowledge of the important independent variables that must be controlled
595 to make a measurement or assay accurate and reproducible can change these data requirements.
596 As is customary in science, it takes time for new scientific insights to reach every lab, and as with
597 any novel field, nanotechnology is evolving and maturing. This maturing process is evident in

598 the nanosafety field as well as in bioinformatics; the first generation of results may not be opti-
599 mal, but they must be used as a basis for improvement or the field will not progress. Another major
600 challenge in nanoinformatics is that researchers are continuing to refine measurement techniques,
601 which could change the comparability of measure results over time. These kinds of issues are re-
602 lated to the concepts of data quality and completeness, which were discussed - along with recom-
603 mendations for progress - in an earlier article in the NDCI series [4].

604 **Lack of usable documentation**

605 The available documentation for external resources often just introduces the resource and provides
606 instructions for its use, but does not convey adequate information to understand the conceptualiza-
607 tion behind the database design. A commonly accepted minimum documentation standard would
608 be helpful.

609 **Need to protect intellectual property hinders data sharing**

610 Although data sharing encourages the public to use and exploit knowledge contained in a database,
611 restrictions may be in place to protect intellectual property and investments in generating and up-
612 dating database content. Often, these restrictions have unclear statements about ownership, copy-
613 right, and licensing. Researchers are sometimes reluctant to share data until they are completely
614 done analyzing and reporting their results out of fear that someone will take their data and use it in
615 a way that limits or reduces the novelty of their work [65]. Some have even suggested that those
616 performing analysis on data they had no role in generating are "research parasites" [66]. The need
617 to maintain "unique selling points" of a database can impede data sharing. One solution to over-
618 come this challenge is to provide a web service with restricted accesses in support of data retrieval
619 while maintaining a customized interface to maintain the unique characteristics of the resource.

620 **Lack of project funding**

621 Individual projects to build data resources and repositories usually do not have funding allocated to
622 data integration. Further, it is not clear which people in the management and funding chain are the

623 correct contacts for expanding a project scope to include integration. This is also a primary con-
624 straint for driving standardization towards a common model. The funding issues extend beyond the
625 necessity to win monetary support that is shared by all research endeavors because these projects
626 can often be seen as investments in infrastructure or tools and are thus perceived to fall outside
627 the purview of basic science funding. Data projects, however, are actually significant exploratory
628 investigations into scientific questions and not just IT projects. Data resources are a major future
629 source of scientific knowledge, and integration across numerous sources expands research opportu-
630 nities.

631 **Stakeholder nanomaterial data integration needs**

632 To address key challenges, stakeholders identified the functionality and web services needed to en-
633 able data integration across nanomaterial repositories. Stakeholders also identified use case driven
634 integration needs with non-nanotechnology resources.

635 **Functionality Needed to Enable Data Integration across Nanomaterial Repos-** 636 **itories**

637 **Use of shared controlled vocabularies**

638 To integrate across resources, each resource needs either to adopt shared controlled vocabularies
639 or to be able to map to agreed-upon standards. When mapping between controlled vocabularies, it
640 is important to fully document the mappings and develop tools to assist in the mapping and trans-
641 formation of the data. Although tool development to automate mapping of terms and schemas re-
642 quires significant work, time is saved in the long run as standards evolve. Adoption of a common
643 language is important, as well as using open standards for data exchange.

644 **Data search and retrieval by ontological terms**

645 Most nanomaterial resources support basic search and retrieval by nanomaterial, characterization,
646 protocol, and publication. To facilitate search and retrieval across resources, it is necessary for re-

647 sources to support searching by ontological term. Additionally, search capability should support
648 retrieval of data (e.g. primary particle characteristics) across each nanomaterial resource and re-
649 trieval of detailed information from the same source on study endpoints applicable to the resource.
650 For example, in the case of toxicity data, it is necessary to support retrieval of particle fate char-
651 acteristics during testing as well as information on the test medium. eNanoMapper's search sys-
652 tem allows searching using ontologies, taking into account synonyms. The demonstration server
653 at <https://search.data.enanomap.net/> allows simultaneous searching over data collected by
654 eNanoMapper and by caNanoLab.

655 **User friendly web-based data submission forms**

656 Nanomaterial resources should provide user friendly tools supporting the submission of data on
657 nanomaterials, characterizations, protocols, and publications via web-based forms. These forms
658 should constrain data entry by requiring use of a controlled vocabulary.

659 **Data import and export tools**

660 Resources should provide support for the validation, import, and export of data in standard data file
661 formats such as ISA-TAB-Nano [67,68], which would allow data to be exported from one database
662 directly into another. It is understood that the development of such tools would require a significant
663 amount of work for resources not currently supporting standards like ISA-TAB-Nano.

664 **Tools to analyze and visualize data**

665 Data analysis and visualization tools within and across nanomaterial resources will facilitate cross
666 material comparisons. Visualizing nanomaterials in 3D and displaying scatter plots and distribution
667 plots across data would assist in optimizing nanomaterial design. Analytic tools need to support the
668 work of many disciplines, including chemistry, biology, toxicology, medicine, and physics.

669 **Data modeling tools**

670 Data modeling tools assist in predicting nanomaterial behavior in different biological and environ-
671 ment systems. The integration of nanomaterial resources with data modeling tools requires that
672 each resource provide access to sufficiently high quality and complete data sets [4].

673 **Facilities for rating data sets for data quality and completeness**

674 Prior to integrating with an existing nanomaterial resource, it is important to understand the data
675 quality and completeness of the resource. Facilities that rate data for completeness and/or quality
676 can assist in providing this assessment. This may include rating against minimum information as
677 well as feedback from users who try to reproduce those data. However, assessing data complete-
678 ness and quality is decidedly non-trivial. A thorough examination of this issue is presented in an
679 other article in the NDCI series [4].

680 **Data Annotations**

681 It is important that data are clearly annotated with statements such as possible provenance, includ-
682 ing ownership and licensing or rights waiving where applicable. Understandably, data can be pro-
683 prietary, and if so should be clearly marked as proprietary. The growing use of resources, such as
684 ZENODO (<http://zenodo.org/>) and FigShare (<https://figshare.com/>), which allow users to assign a
685 specific license to their research data, is arguably indicative of a growing awareness of the impor-
686 tance of clarity regarding rights to data usage within the scientific community - although these re-
687 sources do not support the application of automated data integration techniques [69]. In addition to
688 annotations on data provenance, data annotations can also be provided to further clarify the quality
689 of the data.

690 **Web Services Needed to Enable Data Integration across Nanomaterial Repositories**

691 Stakeholders supporting the use of nanotechnology in the biomedicine and the nanosafety com-
692 munity indicated that the Biomedical Community needs common web services supporting the ex-
693 change of nanomaterials, characterizations, protocols, and publications in support of cross mate-

694 rial comparison. By integrating with other nanomaterial repositories supporting biomedicine and
 695 with other repositories from environmental and health, the biomedical community hopes to better
 696 predict the bio-distribution and toxicity of nanomaterials in model organisms, including humans.
 697 Additionally, the biomedical community would like to obtain detailed information on the investi-
 698 gation, studies, and assays based on metadata identified in the ISA-TAB standard. To support data
 699 integration, ISA-TAB and ISA-TAB-Nano Application Programming Interfaces (APIs) are under
 700 development that retrieve entities based on the ISA-TAB and ISA-TAB-Nano JavaScript Object
 701 Notation (JSON) schemas (<https://github.com/ISA-tools>). The Nanosafety Community has many
 702 interests and covers many different scientific domains. But of special interest, at this moment,
 703 for linking databases, are web services oriented at two central entities in publishing: most similar
 704 nanomaterials, and anything about the same paper or experimental protocol. Common web services
 705 envisaged by these stakeholders as being needed to support integration of nanomaterial data in the
 706 biomedical nanotechnology and nanosafety domains are presented in Table 5.

707 **Needs for Integrating Nanotechnology Repositories with Non-nanotechnology** 708 **Resources**

709 Stakeholders identified a variety of non-nanotechnology resources that must be accessed to support
 710 use case driven data integration needs; these are summarized in Table 6.

Table 6: Non-Nanotechnology Resources needed to support use case driving data integration.

Non-nanotechnology Resource	Description
Life Sciences and Chemistry Databases	Life science and chemistry databases in general, containing information about human biology (both experimental data, as well as knowledge bases) and chemistry (functionality, chemical structure, etc.) [70,71]. Needed to inform the design new nanomaterials to avoid potential negative influences on human health.

Image Archives	Such as the National Biomedical Imaging Archive (NBIA) (https://ncia.nci.nih.gov/ncia/login.jsf), The Cancer Image Archive (TCIA) (http://www.cancerimagingarchive.net/), or other image archive to display MRIs or other image modalities of subjects in which nanomaterials are used for diagnostic and/or therapeutic purposes. A "public domain" image archive illustrating images used in articles, e.g. SEM-pictures would assist in visualizing particle characterizations (see http://www.enanomapper.net/library/image-descriptor-tutorial).
Image Contrast Agent Repository	For example, the Molecular Imaging and Contrast Agent Database (MICAD) (http://www.ncbi.nlm.nih.gov/books/NBK5330/) to obtain information on image contrast agents to compare with nanomaterials used in diagnostic imaging.
Model Organisms Repository	Such as the Mouse Genome Informatics (MGI) (http://www.informatics.jax.org/) resource to access information on animal models used in in vivo characterizations involving nanomaterials.
Publication Sources	PubMed LinkOut or publication vendors such as Elsevier (http://www.elsevier.com/books-and-journals/content-innovation/data-base-linking) to link nanomaterial data to nanomaterial publications. An example of this is the caNanoLab interface with ScienceDirect publications through Elsevier.
Clinical Trials Management Systems (CTMS)	Such as OpenClinica to access clinical data associated with the use of nanomaterials in human clinical trials.
Genomic Data / Biomarker Repositories	Such as the NCI Genomic Data Commons to maintain molecular data for transfection and targeting characterization involving nanomaterials.
Chemical and Agent Repositories	Such as PubChem, ChemSpider, ChEBI, and vendor repositories like Sigma Aldrich to obtain information on chemicals used in nanomaterial compositions. Integrate with small molecule repositories like DrugBank [72] to compare a small molecule (e.g. magnevist) with a nanomaterial formulation that associates with the small molecule (e.g. dendrimer magnevist complex).

Modeling Tools	Modeling and simulation tools as well as 3D structural modeling tools. Integrating with modeling and simulation tools will assist in modeling the effects of nanomaterial size, shape, and other properties on biodistribution and toxicity. Integrating with 3D modeling tools such as The Collaboratory for Structural Nanobiology - CSN (http://uqbar.ncifcrf.gov/Advanced_Structure_Analysis/HOME.html) facilitates the display on nanomaterial structures in 3D leveraging a Protein Data Bank (PDB) file.
Analysis and Visualization Tools	Includes various tools such as R (https://www.r-project.org/) [73], an environment for statistical computing, and Bioconductor [74], D3.js, and other tools to analyze and visualize nanomaterial data in support of nanomaterial comparisons.
Ontology / Taxonomy Resources	To obtain an up-to-date database of ontologies in a table type format so that one can easily review them. This includes resources like the NCI Thesaurus (https://ncit.nci.nih.gov/ncitbrowser/) [75], BioPortal (http://bioportal.bioontology.org/) [76], and Ontobee (http://www.ontobee.org/). This will allow databases to link to term references and accession numbers.

711 **Stakeholder recommendations for the nanotechnology community in**
712 **furthering integration**

713 To assist in providing guidance to the nanotechnology community, stakeholders provided recom-
714 mendations for furthering the integration and exchange of data sets across nanomaterial resources.
715 Guidance centered around the development of pilot projects supporting data integration and the es-
716 tablishment of a global alliance in nanotechnology for standardizing data formats and web services.

717 **Obtain commitment to integration**

718 Stakeholders expressed that the only way to achieve integration effectively is to:

- 719 1. Be committed to achieving integration,

Table 5: Common web services envisaged by these stakeholders as being needed to support integration of nanomaterial data in the biomedical nanotechnology and nanosafety domains.

Web Service Method	Description
createIdentifier	Creates a Universally Unique Identifier (UUID) for any entity such as a material, characterization, protocol, or publication
getCharacterization	Retrieves characterizations for a material by material type and characterization type (e.g. size) and returns characterization data in JSON and XML format.
getDataByDOI	Returns (pointers to) entries in the database with information about or from a specific publication.
getDataByPubMedID	Returns (pointers to) entries in the database with information about or from a specific publication.
getIdentifier	Retrieves a UUID for any entity such as a material, characterization, protocol, or publication
getIsaTabNano	Retrieves ISA-TAB-Nano files associated with a publication (DOI, PubMed)
getInvestigation	Retrieves an investigation associated with a specific disease and/or nanomaterial type and returns an investigation in JSON or XML format. The JSON and XML format would be based on metadata from ISA-TAB-Nano.
getMaterial	Retrieves materials by material type (e.g. dendrimer) or property (e.g. size) and returns a material in JSON or XML format. The JSON and XML format would represent the minimal information about a material.
getProtocol	Retrieves protocols by protocol type (e.g. in vitro) and returns a protocol document and list of materials characterized with the protocol if requested. The protocol document can be returned in a format that uses a common workflow language (e.g. CWL) and/or as a document file.
getPublication	Retrieves publications associated with a material, characterization, and/or protocol, and returns a DOI, PubMed ID, and/or URL to the publication.
getStudy	Retrieves a study associated with a specific assay type and/or nanomaterial type and returns a study in JSON or XML format. The JSON and XML format would be based on metadata from ISA-TAB-Nano.
searchByChemistry	Retrieves nanomaterials based on chemical structure or chemical similarity. Supports a function such as: "Find the most similar structure in database X".

- 720 2. Have the funding in place to complete the effort,
- 721 3. Get the right people (i.e. hands-on developers and nanomaterial experimental experts) to-
- 722 together to work through details of conceptual design and controlled vocabulary, and
- 723 4. Continue fostering a commitment to maximum possible transparency and community-wide
- 724 sharing of approaches, intentions, and techniques, despite the concurrent need of individual
- 725 teams to remain competitive for what will certainly represent limited funding opportunities.

726 This good faith collaboration is the necessary key to making enough progress to achieve the mo-

727 mentum needed for success.

728 **Initiate pilot integration projects**

729 Initiating pilot projects in data source integration efforts is critical. As it stands, individual data re-
730 sources are funded for individual purposes and collaboration and interoperability can be difficult.
731 Based on the U.S. NNI's signature initiative for a knowledge infrastructure [77], there is already
732 a documented need for collaborative resources. Now is clearly the time for funding pilot collabo-
733 rative projects focused on data integration. These should include databases, repositories, ontology
734 designers, experimental researchers, and predictive modelers for a better understanding of the data
735 life cycle and for development of meaningful plans to go forward with existing and new knowledge
736 management resources.

737 **Establish GAIN - a Global Alliance in Nanotechnology - to develop integration standards**

738 Similar to the genomics community that established a Global Alliance in Genomics and Health
739 (GA4GH), the nanotechnology community should form an organization to develop integration
740 standards. A Global Alliance in Nanotechnology (GAIN) would provide a critical mass of inter-
741 est to develop:

- 742 1. A common model for representing data and their relationships,
- 743 2. A standard data dictionary, and
- 744 3. Web service specifications enabling integration.

745 In the stakeholder survey, all stakeholders agreed to participate in a Global Alliance pending avail-
746 ability of funding and time. The eNanoMapper project already actively participates in various col-
747 laborations, including the NanoSafety Cluster (NSC) Database Working Group (along with par-
748 ticipation in other NSC working groups), the US-EU Communities of Research working group on
749 Databases and Computational Modeling for NanoEHS, the US NanoWG, the CODATA/VAMAS
750 Working Group developing the Uniform Description System for Nanomaterials [78] and applied
751 for associate partnership with the CEN/CENELEC node in Europe of the International Standards
752 Organisation (ISO). Alliances with these organizations can be strengthened to avoid unnecessary
753 duplication of effort across the broader community with the primary objective of supporting and

754 enabling concrete open source projects around ontologies, nanoinformatics tools, and data integra-
755 tion.

756 **Focus on providing high quality and complete data sets in data repositories to encourage inte-** 757 **gration**

758 Individual repositories should recognize the importance of providing high quality and sufficiently
759 complete data, rather than simply focusing on providing large amounts of data. However, assess-
760 ing data quality is a complex issue as is the related topic of data completeness. It should also be
761 recognized that the requirements for data to be considered complete and the degree of quality re-
762 quired may be contingent upon the intended purpose of the data. The extent to which different data
763 resources may have legitimately different definitions of data completeness, based upon their differ-
764 ent objectives, underscores the importance of nanoinformatics data resource developers collectively
765 recognizing the value of data integration and the need to ensure the necessary data and metadata
766 required to support integration are documented. A thorough examination of these challenges and
767 a set of recommendations to promote and extend best practice is presented in another article in the
768 NDCI series [4].

769 **Implement data stewardship**

770 Data stewardship should be central to any nanomaterial project. Good stewardship requires that
771 all researchers involved in the project actively participate throughout the process, from beginning
772 to conclusion. This effort involves experimental design, data management plans (including plan-
773 ning for data sharing and adoption of scientific methods in handling data), data citation, and more.
774 Stewardship implies setting aside resources for these tasks. Some will be monetary resources, e.g.
775 for cloud storage, data hosting, possibly commercial support in making data available in commu-
776 nity formats, but other actions should be a core part of the daily research of all the people involved
777 in the project. Postponing planning for data (handling, retrievability, and storage) inevitably jeopard-
778 izes good stewardship and increases costs substantially [69].

779 **Recommendations: A Path forward for achieving data integration**
780 **across nanomaterial resources and with non-nanotechnology reposi-**
781 **tories**

782 Taking into consideration needs of the stakeholders (Section), a multi-step path forward to achieve
783 meaningful progress in integrating nanomaterial data resources is proposed. The four phases iden-
784 tified in Figure 5 provide a roadmap for achieving data integration. Each phase is discussed in
785 greater detail below.

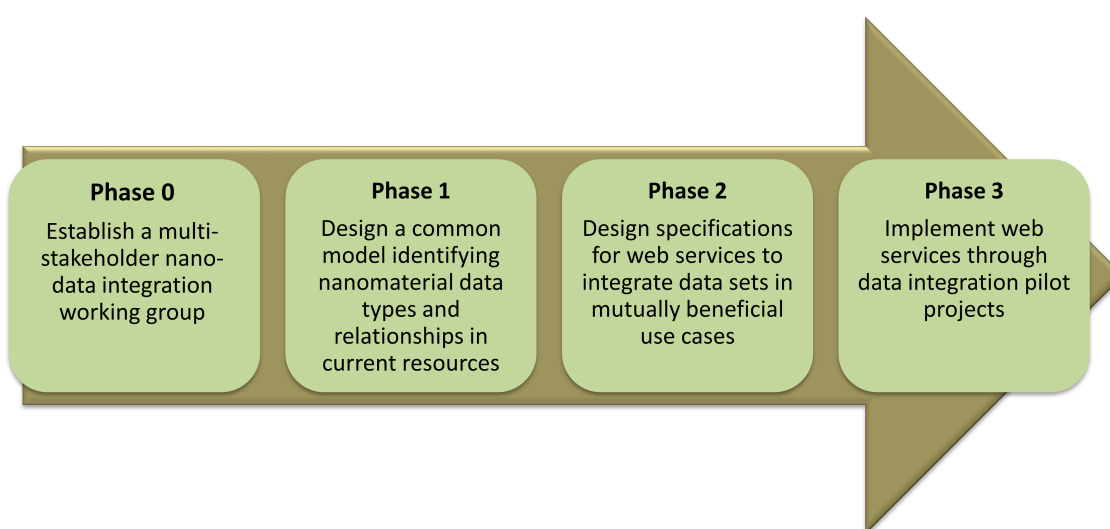


Figure 5: Roadmap of recommendations for achieving data integration across nanomaterial and non-nanomaterial repositories.

786 **Phase 0: Establishment of an organization dedicated to achieving data inte-**
787 **gration in the nanomaterial domain**

788 The time has come to establish a multi-stakeholder, multi-disciplinary, international group focused
789 on nanotechnology data integration. As described above the Global Alliance in Nanotechnology
790 (GAIN) would provide the visibility and energy to start the process towards meaningful data inte-
791 gration in nanoinformatics. GAIN could be an independent group or part of an existing working
792 group such as the Nano WG (<https://wiki.nci.nih.gov/display/ICR/Nanotechnology+Working+>

793 Group) and the NanoSafety Cluster (<http://nanosafetycluster.eu/>) [79] focused on achieving data
794 integration goals. Initial goals include development of a common model to describe the nanoma-
795 terial domain with associated web services supporting data exchange across specific nanomaterial
796 sources.

797 **Phase 1: Design of a common model that identifies nanomaterial entities and** 798 **their relationships within existing resources**

799 One of the first tasks for an organization such as GAIN would be development of a common model
800 that identifies nanomaterial entities and their relationships. It is recommended that the common
801 model be a graph model that depicts nanomaterial entities and nodes and associated relationships
802 as edges (Figure 6). A graph model can provide a flexible structure that can more readily changed
803 as the model evolves. The design of the common model can prioritize identifying the nodes and
804 edges that cross multiple fields such as nanomaterial composition and physico-chemical charac-
805 terizations [15]. Concepts from ISA-TAB-Nano and other ontologies and description systems can
806 be leveraged to represent entities associated with investigations, studies, assays, and materials. It
807 is important to note that this common model is not envisaged as a single, authoritative, federated
808 cyberinfrastructure to facilitate integration in an automated manner. Rather, this model is intended
809 to provide a centralized community-wide understanding of the nanoinformatics space, capturing an
810 overview of the data types implicated, and providing insight into where it makes sense to dedicate
811 resources toward detailed integration projects and tools.

812 **Phase 2: Design specifications for web services that implement the common** 813 **model**

814 Once the common model is established, specifications for common web services can be developed,
815 including defining service endpoints based on entities in the common model. Web service speci-
816 fication should be prioritized to focus on a basic query to retrieve nanomaterials by nanomaterial
817 characteristics and other properties. Web services can be further expanded to accommodate use-

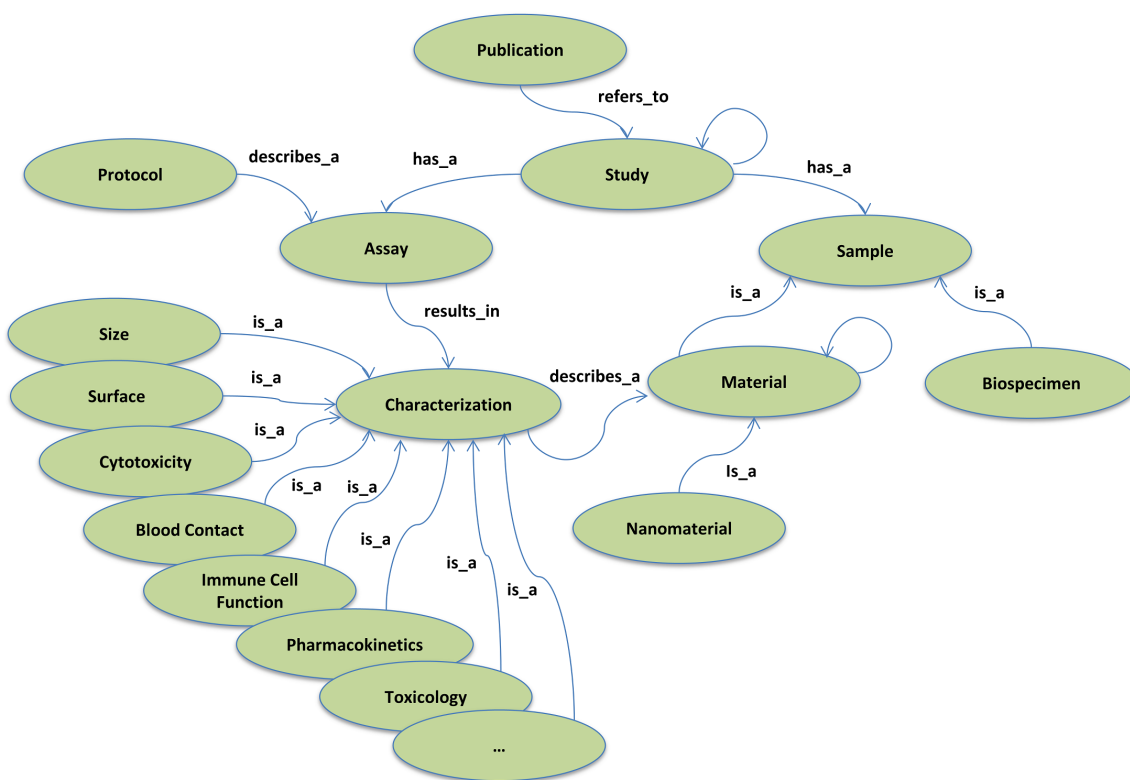


Figure 6: Example graph model depicting nodes (e.g. nanomaterial) and edges (describes_a).

818 case-dependent data exchange with non-nanotechnology sources. In support of data exchanges
 819 with non-nanotechnology sources, established interfaces could be published and organizations
 820 could collaborate with resource providers on developing a common interface to facilitate re-use.

821 **Phase 3: Implementation of web services through pilot projects**

822 Once an initial web service is designed, pilot projects should be started as soon as possible to
 823 implement the web service with an ultimate goal of querying across nanomaterial resources. Pi-
 824 lot projects should focus on developing re-usable software that can be extended in support of
 825 other pilot efforts. Software should be made available as open source and published as a GitHub
 826 (<https://github.com/>) repository. Lessons learned from pilot efforts should result in improvements
 827 to the common model and web services design specifications.

828 **Closing remarks**

829 The various challenges recognized by members of the nanoinformatics community are hampering
830 efforts to integrate across nanomaterial and other non-nanotechnology resources in a meaningful
831 way. The technical and operational challenges summarized in Figure 4 are significant barriers to
832 scientific progress in designing new and higher impact nanomaterials and in understanding how
833 nanomaterials interact with biological, environmental, and other systems. The tools to take advan-
834 tage of high quality nanotechnology data exist but cannot be exploited unless true data sharing and
835 integration is possible. This paper analyzes these challenges and outlines a path forward to real
836 progress.

837 The authors encourage readers to share feedback or join the National Cancer Informatics Program
838 (NCIP) Nanotechnology Working Group (<https://nciphub.org/groups/nanowg/overview>) and learn
839 more about the Nanomaterial Data Curation Initiative, in particular, by visiting [https://nciphub.org/
840 groups/nanotechnologydatacurationinterestgroup/wiki/MainPage](https://nciphub.org/groups/nanotechnologydatacurationinterestgroup/wiki/MainPage).

841 **Acknowledgements**

842 Authors would like to gratefully acknowledge several organizations who participated in the data
843 integration survey including the: Center for the Environmental Implications of NanoTechnol-
844 ogy (CEINT) which is funded by the National Science Foundation (NSF) and the Environmen-
845 tal Protection Agency (EPA) under NSF Cooperative Agreement DBI-1266252 and EF-0830093;
846 Nanomaterial Registry, which is funded by the National Institutes of Health (NIH) under contract
847 HHSN268201000022C; caNanoLab which is funded in whole or in part with Federal funds from
848 the National Cancer Institute (NCI), NIH, under Contract No. HHSN261200800001E; Center
849 for Safety of Substances, and Products, National Institute of Public Health and the Environment,
850 Netherlands; DECHEMA; and eNanoMapper, funded by the European Union's Seventh Frame-
851 work Programme for research, technological development and demonstration (FP7-NMP-2013-
852 SMALL-7) under grant agreement no. 604134. RLMR is grateful for funding from the European
853 Union Seventh Framework Programme (FP7/2007-2013) under grant agreement #309837 (NanoP-
854 UZZLES project).

855 Authors would also like to acknowledge the Nano WG leaders for their time and expertise in pro-
856 viding the necessary tools that supported collaboration on this article. Authors would like to give
857 special thanks to Dr. Mervi Heiskanen, the Nano WG lead from the NCI Center for Biomedical
858 Informatics and Information Technical (CBIIT) and Dr. Stephanie Morris, from the NCI Office of
859 Cancer Nanotechnology Research (OCNR) for their leadership and support for Nano WG initia-
860 tives including the development of this article.

861 The views, opinions, and content in this article are those of the authors and do not necessarily rep-
862 resent the views, opinions, or policies of their respective employers or organizations. Mention of
863 trade names, commercial products, or organizations does not imply endorsement by the U.S. gov-
864 ernment.

865 **References**

- 866 1. Thomas, D. G.; Klaessig, F.; Harper, S. L.; Fritts, M.; Hoover, M. D.; Gaheen, S.;
867 Stokes, T. H.; Reznik-Zellen, R.; Freund, E. T.; Klemm, J. D.; Paik, D. S.; Baker, N. A. *Wi-*
868 *ley Interdisciplinary Reviews: Nanomedicine and Nanobiotechnology* **2011**, 3 (5), 511–532.
- 869 2. Hendren, C. O.; Powers, C. M.; Hoover, M. D.; Harper, S. L. *Beilstein Journal of Nanotech-*
870 *nology* **2015**, 6, 1752–1762.
- 871 3. Powers, C. M.; Mills, K. A.; Morris, S. A.; Klaessig, F.; Gaheen, S.; Lewinski, N.;
872 Ogilvie Hendren, C. *Beilstein Journal of Nanotechnology* **2015**, 6, 1860–1871.
- 873 4. Marchese Robinson, R.; Lynch, I.; Peijnenburg, W.; Rumble, J.; Klaessig, F.; Hendren, C.;
874 Marquardt, C.; Rauscher, H.; Puzyn, T.; Purian, R.; ÅËberg, C.; Karcher, S.; Vriens, H.;
875 Hoet, P.; Hoover, M.; Harper, S. *Nanoscale* **2016**, 8, 9919–9943.
- 876 5. Rustici, G.; Kolesnikov, N.; Brandizi, M.; Burdett, T.; Dylag, M.; Emam, I.; Farne, A.; Hast-
877 ings, E.; Ison, J.; Keays, M.; Kurbatova, N.; Malone, J.; Mani, R.; Mupo, A.; Pereira, R. P.;
878 Pilicheva, E.; Rung, J.; Sharma, A.; Tang, Y. A.; Ternent, T.; Tikhonov, A.; Welter, D.;

- 879 Williams, E.; Brazma, A.; Parkinson, H.; Sarkans, U. *Nucleic Acids Research* **2013**, *41* (D1),
880 D987–D990.
- 881 6. de la Iglesia, D.; Harper, S.; Hoover, M. D.; Klaessig, F.; Lippell, P.; Maddux, B.; Morse, J.;
882 Nel, A.; Rajan, K.; Reznik-Zellen, R. et al. *Nanoinformatics 2020 Roadmap*; National
883 Nanomanufacturing Network, 2011.
- 884 7. Hoover, M. D.; Myers, D. S.; Cash, L. J.; Guilmette, R. A.; Kreyling, W. G.; Oberdörster, G.;
885 Smith, R.; Cassata, J. R.; Boecker, B. B.; Grissom, M. P. *Health physics* **2015**, *108* (2),
886 179–194.
- 887 8. Boholm, M.; Arvidsson, R. *NanoEthics* **2016**, *10* (1), 25–40.
- 888 9. Rauscher, H.; Sokull-Klüttgen, B.; Stamm, H. *Nanotoxicology* **2012**, *7* (7), 1195–1197.
- 889 10. Vance, M. E.; Kuiken, T.; Vejerano, E. P.; McGinnis, S. P.; Hochella Jr, M. F.; Rejeski, D.;
890 Hull, M. S. *Beilstein journal of nanotechnology* **2015**, *6* (1), 1769–1780.
- 891 11. Xia, Y. *Angewandte Chemie International Edition* **2014**, *53* (46), 12268–12271.
- 892 12. Roco, M. C.; Mirkin, C. A.; Hersam, M. C. *Nanotechnology research directions for societal*
893 *needs in 2020 retrospective and outlook*; World Technology Evaluation Center ; Springer,
894 2011.
- 895 13. Shao, C.-Y.; Chen, S.-Z.; Su, B.-H.; Tseng, Y. J.; Esposito, E. X.; Hopfinger, A. J. *J. Chem.*
896 *Inf. Model.* **2013**, *53* (1), 142–158.
- 897 14. Rumble, J.; Freiman, S. *Data Science Journal* **2012**, *11* (0), ASMD1–ASMD6.
- 898 15. Stefaniak, A. B.; Hackley, V. A.; Roebben, G.; Ehara, K.; Hankin, S.; Postek, M. T.; Lynch, I.;
899 Fu, W.-E.; Linsinger, T. P. J.; Thünemann, A. F. *Nanotoxicology* **2013**, *7* (8), 1325–1337.
- 900 16. Bult, C. J.; Eppig, J. T.; Kadin, J. A.; Richardson, J. E.; Blake, J. A.; Group, M. G. D. et al.
901 *Nucleic acids research* **2008**, *36* (suppl 1), D724–D728.

- 902 17. Kutmon, M.; Riutta, A.; Nunes, N.; Hanspers, K.; Willighagen, E. L.; Bohler, A.; Mélius, J.;
903 Waagmeester, A.; Sinha, S. R.; Miller, R.; Coort, S. L.; Cirillo, E.; Smeets, B.; Evelo, C. T.;
904 Pico, A. R. *Nucleic Acids Research* **2016**, *44* (D1), D488–D494.
- 905 18. Jeliaskova, N.; Doganis, P.; Fadeel, B.; Grafstrom, R.; Hastings, J.; Jeliaskov, V.; Kohonen, P.;
906 Munteanu, C. R.; Sarimveis, H.; Smeets, B.; Tsiliki, G.; Vorgrimmler, D.; Willighagen, E. The
907 first eNanoMapper prototype: A substance database to support safe-by-design. In *Bioinform-*
908 *matics and Biomedicine (BIBM), 2014 IEEE International Conference on*; IEEE, 2014; pp
909 1–9.
- 910 19. Izak-Nau, E.; Huk, A.; Reidy, B.; Uggerud, H.; Vadset, M.; Eiden, S.; Voetz, M.; Himly, M.;
911 Duschl, A.; Dusinska, M.; Lynch, I. *RSC Adv.* **2015**, *5* (102), 84172–84185.
- 912 20. Williams, A. J.; Harland, L.; Groth, P.; Pettifer, S.; Chichester, C.; Willighagen, E. L.;
913 Evelo, C. T.; Blomberg, N.; Ecker, G.; Goble, C.; Mons, B. *Drug Discovery Today* **2012**, *17*
914 (21-22), 1188–1198.
- 915 21. Maglott, D.; Ostell, J.; Pruitt, K. D.; Tatusova, T. *Nucleic acids research* **2005**, *33* (suppl 1),
916 D54–D58.
- 917 22. Samwald, M.; Jentzsch, A.; Bouton, C.; Kallesoe, C.; Willighagen, E.; Hajagos, J.; Mar-
918 shall, M.; Prud'hommeaux, E.; Hassanzadeh, O.; Pichler, E.; Stephens, S. *Journal of Chem-*
919 *informatics* **2011**, *3* (1), 19+.
- 920 23. Jupp, S.; Malone, J.; Bolleman, J.; Brandizi, M.; Davies, M.; Garcia, L.; Gaulton, A.;
921 Gehant, S.; Laibe, C.; Redaschi, N.; Wimalaratne, S. M.; Martin, M.; Le Novère, N.; Parkin-
922 son, H.; Birney, E.; Jenkinson, A. M. *Bioinformatics* **2014**, *30* (9), 1338–1339.
- 923 24. Cheung, K.-H.; Frost, H. R.; Marshall, M. S.; Prud'hommeaux, E.; Samwald, M.; Zhao, J.;
924 Paschke, A. *BMC Bioinformatics* **2009**, *10* (Suppl 10), S10+.
- 925 25. Eyres, T. P. *EMBnet.journal* **2013**, *19*, 36–39.

- 926 26. Juty, N.; Le Novère, N.; Laibe, C. *Nucleic Acids Research* **2012**, *40* (D1), D580–D586.
- 927 27. van Iersel, M. P.; Pico, A. R.; Kelder, T.; Gao, J.; Ho, I.; Hanspers, K.; Conklin, B. R.;
928 Evelo, C. T. *BMC Bioinformatics* **2010**, *11* (1), 5+.
- 929 28. Chambers, J.; Davies, M.; Gaulton, A.; Papadatos, G.; Hersey, A.; Overington, J. *Journal of*
930 *Cheminformatics* **2014**, *6* (1), 43+.
- 931 29. Wohlgemuth, G.; Haldiya, P. K.; Willighagen, E.; Kind, T.; Fiehn, O. *Bioinformatics* **2010**, *26*
932 (20), 2647–2648.
- 933 30. Hastings, J.; Chepelev, L.; Willighagen, E.; Adams, N.; Steinbeck, C.; Dumontier, M. *PLoS*
934 *ONE* **2011**, *6* (10), e25513+.
- 935 31. Thomas, D. G.; Pappu, R. V.; Baker, N. A. *Journal of Biomedical Informatics* **2011**, *44* (1),
936 59–74.
- 937 32. Hastings, J.; Owen, G.; Dekker, A.; Ennis, M.; Kale, N.; Muthukrishnan, V.; Turner, S.;
938 Swainston, N.; Mendes, P.; Steinbeck, C. *Nucleic Acids Research* **2016**, *44* (D1),
939 D1214–D1219.
- 940 33. Harper, S. L.; Dahl, J. A.; Maddux, B. L. S.; Tanguay, R. L.; Hutchison, J. E. *International*
941 *Journal of Nanotechnology* **2008**, *5* (1), 124+.
- 942 34.
- 943 35. Izak-Nau, E.; Huk, A.; Reidy, B.; Uggerud, H.; Vadset, M.; Eiden, S.; Voetz, M.; Himly, M.;
944 Duschl, A.; Dusinska, M. et al. *RSC Advances* **2015**, *5* (102), 84172–84185.
- 945 36. Batchelor, C.; Brenninkmeijer, C.; Chichester, C.; Davies, M.; Digles, D.; Dunlop, I.;
946 Evelo, C.; Gaulton, A.; Goble, C.; Gray, A.; Groth, P.; Harland, L.; Karapetyan, K.;
947 Loizou, A.; Overington, J.; Pettifer, S.; Steele, J.; Stevens, R.; Tkachenko, V.; Waag-
948 meester, A.; Williams, A.; Willighagen, E. Scientific Lenses to Support Multiple Views over

- 949 Linked Chemistry Data. In *The Semantic Web - ISWC 2014*; Mika, P., Tudorache, T., Bern-
950 stein, A., Welty, C., Knoblock, C., Vrandečić, D., Groth, P., Noy, N., Janowicz, K., Goble, C.,
951 Eds.; Springer International Publishing, 2014; Vol. 8796, pp 98–113.
- 952 37. Brenninkmeijer, C.; Evelo, C.; Goble, C.; Gray, A. J. G.; Groth, P.; Pettifer, S.; Stevens, R.;
953 William, A. J.; Willighagen, E. L. Scientific Lenses over Linked Data: An Approach to Sup-
954 port Task Specific Views of the Data. A Vision. In *Linked Science 2012 - Tackling Big Data*;
955 CEUR-WS.org, 2012; Chapter 5.
- 956 38. Berners-Lee, T.; Hendler, J.; Lassila, O. *Scientific American* **2001**, 284 (5), 34–43.
- 957 39. Marshall, M. S.; Boyce, R.; Deus, H. F.; Zhao, J.; Willighagen, E. L.; Samwald, M.; Pich-
958 ler, E.; Hajagos, J.; Prud'hommeaux, E.; Stephens, S. *Web Semantics: Science, Services and*
959 *Agents on the World Wide Web* **2012**, 14, 2–13.
- 960 40. Hastings, J.; Jeliaskova, N.; Owen, G.; Tsiliki, G.; Munteanu, C. R.; Steinbeck, C.; Willigha-
961 gen, E. *Journal of Biomedical Semantics* **2015**, 6 (1), 10+.
- 962 41. Oksel, C.; Ma, C. Y.; Wang, X. Z. *SAR and QSAR in Environmental Research* **2015**, 26 (2),
963 79–94.
- 964 42. Crist, R. M.; Grossman, J. H.; Patri, A. K.; Stern, S. T.; Dobrovolskaia, M. A.; Adishesha-
965 iah, P. P.; Clogston, J. D.; McNeil, S. E. *Integrative Biology* **2013**, 5 (1), 66–73.
- 966 43.
- 967 44. Hendren, C. O.; Lowry, G. V.; Unrine, J. M.; Wiesner, M. R. *Science of The Total Environ-*
968 *ment* **2015**, 536, 1029–1037.
- 969 45. Brazma, A.; Hingamp, P.; Quackenbush, J.; Sherlock, G.; Spellman, P.; Stoeckert, C.;
970 Aach, J.; Ansorge, W.; Ball, C. A.; Causton, H. C. et al. *Nature Genetics* **2001**, 29 (4),
971 365–371.

- 972 46. Field, D.; Sansone, S.-A. A.; Collis, A.; Booth, T.; Dukes, P.; Gregurick, S. K.; Kennedy, K.;
973 Kolar, P.; Kolker, E.; Maxon, M.; Millard, S.; Mugabushaka, A.-M. M.; Perrin, N.;
974 Remacle, J. E.; Remington, K.; Rocca-Serra, P.; Taylor, C. F.; Thorley, M.; Tiwari, B.;
975 Wilbanks, J. *Science (New York, N.Y.)* **2009**, *326* (5950), 234–236.
- 976 47. Hartig, O.; Langegger, A. A Database Perspective on Consuming Linked Data on the Web. In
977 *Datenbank-Spektrum*; Springer-Verlag, 2010; Vol. 10, pp 57–66.
- 978 48.
- 979 49. Doan, A.; Halevy, A.; Ives, Z. G. *Principles of data integration*; Morgan Kaufmann, 2012.
- 980 50. Kovřížnych, J. A.; Sotníková, R.; Zeljenková, D.; Rollerová, E.; Szabová, E.; Wimmerová, S.
981 *Interdisciplinary Toxicology* **2013**, *6* (2), 67–73.
- 982 51. Truong, L.; Harper, S. L.; Tanguay, R. L. Evaluation of Embryotoxicity Using the Zebrafish
983 Model. In *Drug Safety Evaluation*; Gautier, J.-C., Ed.; Humana Press, 2011; Vol. 691, pp
984 271–279.
- 985 52. Christen, P. *Data Matching*; Springer Berlin Heidelberg: Berlin, Heidelberg, 2012.
- 986 53. Heller, S. R.; McNaught, A.; Pletnev, I.; Stein, S.; Tchekhovskoi, D. *Journal of Cheminform-*
987 *matics* **2015**, *7* (1), 23+.
- 988 54. Hersey, A.; Chambers, J.; Bellis, L.; Bento, A. P.; Gaulton, A.; Overington, J. P. *Drug Discov-*
989 *ery Today: Technologies* **2015**, *14*, 17–24.
- 990 55. Fielding, R. T. Architectural styles and the design of network-based software architectures.
991 Ph. D. Thesis, University of California, Irvine, 2000.
- 992 56.
- 993 57. Kühnel, D.; Marquardt, C.; Nau, K.; Krug, H. F.; Mathes, B.; Steinbach, C. *Environmental*
994 *Sciences Europe* **2014**, *26* (1), 21.

- 995 58. Kimmig, D.; Marquardt, C.; Nau, K.; Schmidt, A.; Dickerhof, M. *Computational Science &*
996 *Discovery* **2014**, 7 (1), 014001.
- 997 59. Atli, A.; Nau, K.; Schmidt, A. Navigation along Database Relationships-An Adaptive Frame-
998 work for Presenting Database Contents as Object Graphs. In *WEBIST*; Institute for Systems
999 and Technologies of Information, Control and Communication, 2011; pp 372–379.
- 1000 60. Miller, A. L.; Hoover, M. D.; Mitchell, D. M.; Stapleton, B. P. *Journal of occupational and*
1001 *environmental hygiene* **2007**, 4 (12), D131–D134.
- 1002 61. Jeliaskova, N.; Jeliaskov, V. *J. Cheminformatics* **2011**, 3, 18.
- 1003 62. Jeliaskova, N.; Chomenidis, C.; Doganis, P.; Fadeel, B.; Grafström, R.; Hardy, B.; Hast-
1004 ings, J.; Hegi, M.; Jeliaskov, V.; Kochev, N. et al. *Beilstein journal of nanotechnology* **2015**,
1005 6 (1), 1609–1634.
- 1006 63. Klimisch, H.-J.; Andreae, M.; Tillmann, U. *Regulatory toxicology and pharmacology* **1997**,
1007 25 (1), 1–5.
- 1008 64. Krug, H. F. *Angewandte Chemie International Edition* **2014**, 53 (46), 12304–12319.
- 1009 65. Reichman, O.; Jones, M. B.; Schildhauer, M. P. *Science* **2011**, 331 (6018), 703–705.
- 1010 66. Longo, D. L.; Drazen, J. M. *New England Journal of Medicine* **2016**, 374 (3), 276–277.
- 1011 67. Thomas, D. G.; Gaheen, S.; Harper, S. L.; Fritts, M.; Klaessig, F.; Hahn-Dantona, E.; Paik, D.;
1012 Pan, S.; Stafford, G. A.; Freund, E. T. et al. *BMC biotechnology* **2013**, 13 (1), 1.
- 1013 68. Marchese Robinson, R. L.; Cronin, M. T. D.; Richarz, A.-N.; Rallo, R. *Beilstein Journal of*
1014 *Nanotechnology* **2015**, 6, 1978–1999. doi:10.3762/bjnano.6.202.
- 1015 69. Wilkinson, M. D.; Dumontier, M.; Aalbersberg, I. J.; Appleton, G.; Axton, M.; Baak, A.;
1016 Blomberg, N.; Boiten, J.-W.; da Silva Santos, L. B.; Bourne, P. E.; Bouwman, J.;
1017 Brookes, A. J.; Clark, T.; Crosas, M. A.; Dillo, I.; Dumon, O.; Edmunds, S.; Evelo, C. T.;

- 1018 Finkers, R.; Gonzalez-Beltran, A.; Gray, A. J. G.; Groth, P.; Goble, C.; Grethe, J. S.;
- 1019 Heringa, J.; Hoen, P. A. C.; Hooft, R.; Kuhn, T.; Kok, R.; Kok, J.; Lusher, S. J.; Mar-
- 1020 tone, M. E.; Mons, A.; Packer, A. L.; Persson, B.; Rocca-Serra, P.; Roos, M.; van Schaik, R.;
- 1021 Sansone, S.-A.; Schultes, E.; Sengstag, T.; Slater, T.; Strawn, G.; Swertz, M. A.; Thomp-
- 1022 son, M.; van der Lei, J.; van Mulligen, E.; Velterop, J.; Waagmeester, A.; Wittenburg, P.; Wol-
- 1023 stencroft, K.; Zhao, J.; Mons, B. *Scientific Data* **2016**, *3*, 160018+.
- 1024 70. Kim, S.; Thiessen, P. A.; Bolton, E. E.; Chen, J.; Fu, G.; Gindulyte, A.; Han, L.; He, J.; He, S.;
- 1025 Shoemaker, B. A.; Wang, J.; Yu, B.; Zhang, J.; Bryant, S. H. *Nucleic Acids Research* **2016**, *44*
- 1026 (D1), D1202–D1213. doi:10.1093/nar/gkv951.
- 1027 71. Wang, Y.; Suzek, T.; Zhang, J.; Wang, J.; He, S.; Cheng, T.; Shoemaker, B. A.; Gindulyte, A.;
- 1028 Bryant, S. H. *Nucleic Acids Research* **2014**, *42* (D1), D1075–D1082. doi:10.1093/nar/gkt978.
- 1029 72. Knox, C.; Law, V.; Jewison, T.; Liu, P.; Ly, S.; Frolkis, A.; Pon, A.; Banco, K.; Mak, C.;
- 1030 Neveu, V. et al. *Nucleic acids research* **2011**, *39* (suppl 1), D1035–D1041.
- 1031 73. *R: A Language and Environment for Statistical Computing*; 2016.
- 1032 74. Gentleman, R. C.; Carey, V. J.; Bates, D. M.; Bolstad, B.; Dettling, M.; Dudoit, S.; Ellis, B.;
- 1033 Gautier, L.; Ge, Y.; Gentry, J. et al. *Genome biology* **2004**, *5* (10), R80.
- 1034 75. Sioutos, N.; de Coronado, S.; Haber, M. W.; Hartel, F. W.; Shaiu, W.-L.; Wright, L. W. *Jour-*
- 1035 *nal of biomedical informatics* **2007**, *40* (1), 30–43.
- 1036 76. Noy, N. F.; Shah, N. H.; Whetzel, P. L.; Dai, B.; Dorf, M.; Griffith, N.; Jonquet, C.; Ru-
- 1037 bin, D. L.; Storey, M.-A.; Chute, C. G. et al. *Nucleic acids research* **2009**, gkp440.
- 1038 77. Roco, M. C. *Journal of Nanoparticle Research* **2011**, *13* (2), 427–445.
- 1039 78. Rumble, J.; Freiman, S.; Teague, C. *Chemistry International* **2015**, *37* (4), 3–7.

- 1040 79. Savolainen, K.; Backman, U.; Brouwer, D.; Fadeel, B.; Fernandes, T.; Kuhlbusch, T.; Land-
1041 siedel, R.; Lynch, I.; Pylkkänen, L. *Nanosafety in Europe 2015-2025: Towards Safe and Sus-*
1042 *tainable Nanomaterials and Nanotechnology Innovations*; Finnish Institute of Occupational
1043 Health: Helsinki, 2013.