

A peer-reviewed version of this preprint was published in PeerJ on 5 February 2015.

[View the peer-reviewed version](https://peerj.com/articles/755) (peerj.com/articles/755), which is the preferred citable publication unless you specifically need to cite this preprint.

Ganzinger M, Knaup P. 2015. Requirements for data integration platforms in biomedical research networks: a reference model. PeerJ 3:e755 <https://doi.org/10.7717/peerj.755>

Requirements for data integration platforms in biomedical research networks: A reference model

Matthias Ganzinger*, Petra Knaup

Institute of Medical Biometry and Informatics, Heidelberg University, Heidelberg, Germany

* Corresponding author:

Matthias Ganzinger

Heidelberg University

Institute of Medical Biometry and Informatics

Im Neuenheimer Feld 305

69120 Heidelberg

Germany

Phone: +49 6221 56-5143

E-mail: matthias.ganzinger@med.uni-heidelberg.de

Abstract

Biomedical research networks need to integrate research data among their members and with external partners. To support such data sharing activities, adequate information technology infrastructure is necessary. To facilitate the establishments of such an infrastructure, we developed a reference model for requirements. The reference model consists of five reference goals and 15 reference requirements. Using the Unified Modeling Language, the goals and requirements are set into relation to each other. In addition, all goals and requirements are described textually in tables. This reference model can be used by research networks as a basis for the resource efficient acquisition of their project specific requirements. Further, a concrete instance of the reference model is described for a research network on liver cancer. The reference model is transferred into a requirements model of the specific network. Based on this concrete requirements model, a service-oriented information technology architecture is derived and also described in this paper.

Introduction

Current biomedical research is supported by modern biotechnological methods producing vast amounts of data (Frey, Maojo & Mitchell, 2007; Baker, 2010). In order to get a comprehensive picture of the physiology and pathogenic processes of diseases, many facets of biological mechanisms need to be examined. Contemporary research, e.g. investigating cancer, is a complex endeavor that can be conducted most successfully when researchers of multiple disciplines cooperate and draw conclusions from comprehensive scientific data sets (Welsh, Jirotko & Gavaghan, 2006; Mathew et al., 2007). As a frequent measure to support cooperation, research networks sharing common resources are established.

37 To generate added value from such a network, all available scientific and clinical data should
38 be combined to facilitate a new, comprehensive perspective. This requires provision of
39 adequate information technology (IT) which is a challenge on all levels of biomedical
40 research. For example, it is inevitable for research networks to use an IT infrastructure for
41 sharing data and findings in order to leverage joint analyses. Data generated by
42 biotechnological devices can only be evaluated thoroughly by applying biostatistical methods
43 with IT tools.

44 However, data structures are often heterogeneous, resulting in the need for a data integration
45 process. This process involves the harmonization of data structures by defining appropriate
46 metadata (Cimino, 1998). Depending on the specific needs and data structures of the research
47 network, often a non-standard IT platform needs to be developed to meet the specific
48 requirements. An important requirement might be the protection of data in terms of security
49 and privacy, especially when patient data are involved.

50 In the German research network SFB/TRR77 – “Liver Cancer. From Molecular Pathogenesis
51 to Targeted Therapies” it was our task to explore the most appropriate IT-architecture for
52 supporting networked research (Woll, Manns & Schirmacher, 2013). The research network
53 consists of 22 projects sharing common resources and research data. To provide this network
54 with a data integration platform we implemented a service-oriented architecture (SOA)
55 (Taylor et al.; Papazoglou et al., 2008; Wei & Blake, 2010; Bosin, Dessì & Pes, 2011). The IT
56 system is based on the cancer Common Ontologic Representation Environment Software
57 Development Kit (caCORE SDK) components of the cancer Biomedical Informatics Grid
58 (caBIG) (Komatsoulis et al., 2008; Kunz, Lin & Frey, 2009). The resulting system is called
59 pelican (platform enabling liver cancer networked research) (Ganzinger et al., 2011).
60 Transfer of these data sharing concepts to other networks investigating different disease areas
61 is possible.

62 We consider our research network as a typical example for a whole class of biomedical
63 research networks. To support this kind of projects, we provide a framework for the
64 development of data integration platforms for such projects. Specifically, we strive for the
65 following two objectives:

66 Objective 1: Provide a reference model of requirements of biomedical research networks
67 regarding an IT platform for sharing and analyzing data.

68 Objective 2: Design a SOA of an IT platform for our research network on liver cancer. It
69 should implement the reference model for requirements. While this SOA is specific to this
70 project, parts can be reused by similar projects.

71 **Methods**

72 For the design of a data integration platform it is important to first capture the requirements of
73 the system’s intended users. To support this task, we developed a reference model for
74 requirements. A reference model is a generic model, which is valid not only for a specific
75 research network, but for a class of such organizations. For the development of the reference

76 model, we used the research network on liver cancer as primary source. These requirements
77 were consolidated and abstracted to get a generic model that can be applied to other research
78 networks.

79 In a first step, a general understanding of the network's aims and tasks was acquired by
80 analyzing written descriptions of the participating projects. In addition, questionnaires were
81 sent to the principal investigators to capture the data types and data formats used within the
82 projects. In a second step, projects were visited and their research subjects and processes were
83 captured by interviewing project members.

84 For the reference model, we use the term *goal* to describe the highest level of requirements.
85 This is in accordance with ISO/IEC/IEEE 24765 where goal is defined as “an intended
86 outcome” (ISO/IEC/IEEE, 2011). In contrast, *requirement* is defined as “a condition or
87 capability needed by a user to solve a problem or achieve an objective” (ISO/IEC/IEEE,
88 2011). In our reference model, each requirement was related to a goal, either directly or
89 indirectly. Requirements were later on mapped to concrete functions in the resulting data
90 integration system. On the other hand, goals were used to structure requirements and usually
91 do not lead to a specific function of the system.

92 To provide a more detailed characterization of the goals, we provide a standardized table for
93 each of them. It covers the reference number, name, description and weight of the goal. Table
94 1 shows the structure of such a table. The complete set of tables for all goals is available in
95 the supplementing information S1.

96 The requirements are documented in the same way as goals are. Figure 1 shows a Unified
97 Modeling Language (UML) diagram with all elements used for describing of both goals and
98 their subordinated requirements (*Object Management Group, 2012*). As for the goals, we
99 provide a set of tables with more detailed descriptions for all requirements in supplementing
100 information S1. In total, we identified 15 requirements for the reference model.

101 The instantiation of the reference model for requirements to meet the needs of a specific
102 research network provides the basis for the architecture of the desired data integration and its
103 subsequent implementation. We provide a concrete instance of a reference model as well as
104 the resulting IT-architecture in this manuscript.

105 For the research described in this paper, ethics approval was not deemed necessary. This work
106 involved no human subjects in the sense of medical research, as e.g. covered by the
107 Declaration of Helsinki (*World Medical Association, 2013*). At no time patients were included
108 for survey or interview. Data was only acquired from scientists regarding their work and data,
109 but no personal or patient related data were gathered. Participants were not required to
110 participate in this study. They consented by returning the questionnaire. No research was
111 conducted outside Germany, the authors' country of residence. However, in other countries
112 the approval of an institutional review board or other authority might be necessary to apply
113 the reference model.

114 Results

115 In this section we first describe the reference model for requirements. Then, we show how a
116 concrete model for requirements and an IT-architecture is derived from this reference model.
117 The reference model for requirements is an abstract model and thus a universally usable
118 artefact. It is mapped in several steps to the network specific system architecture.

119 Reference Model

120 The reference model for requirements covers five reference goals (RG). An overview of the
121 goals and their relations is shown in Figure 2 by means of a UML requirements diagram. The
122 reference goals are:

- 123 • Conduct research project (reference goal RG1): The ultimate goal of a research
124 network is to fulfill the intended research tasks. This usually corresponds to the project
125 specification of the funding organization.
- 126 • Answer research questions (reference goal RG2): Each research network has specific
127 research questions it pursues to answer. These questions frame the core of the network
128 and led to its establishment in the first place.
- 129 • Create, store, and retrieve data (reference goal RG3): Research networks need data to
130 conduct the project. Thus, it is necessary to generate and handle them.
- 131 • Analyze data (reference goal RG4): To generate knowledge out of the data it is
132 necessary to analyze them.
- 133 • Control data access and usage (reference goal RG5): Research networks need to
134 protect their data. This includes the prevention of unauthorized access to protected
135 data like patient data as well as aspects of intellectual property rights that need to be
136 respected by authorized users as well.

137 These goals are ordered in a hierarchical structure: Goal RG1 acts as the root node, which has
138 the two sub goals RG2 and RG 3. Goal RG4 is subordinated to Goal RG2, whereas Goal RG5
139 is a sub goal of Goal RG3.

140 Each goal has several requirements. In total, the reference model contains 15 reference
141 requirements (RR). A UML diagram with all reference goals and reference requirements is
142 shown in Figure 3. Reference requirements are associated with the reference goals as follows:

143 Goal 2 is associated with the reference requirements to *create data (RR1)* and to *retrieve*
144 *external data (RR2)*. These two requirements respect possible sources of data necessary for
145 the research network. Reference requirement RR3, *represent data*, is further defined by its
146 subordinate requirements *define syntax (RR4)*, *define data model (RR5)*, *identify data (RR6)*,
147 and *define semantics (RR7)*.

148 Goal 5, *control data access and usage*, has two aspects, which are represented by reference
149 requirements RR8 and RR9. RR9 requires the creator's contribution in the generation of data
150 for the research network to be recognized when data are used by others. As a consequence,
151 even users with legitimate access to the system have to adhere to usage regulations. These
152 regulations should be checked and enforced by the system as far as possible. In contrast, RR9
153 covers the requirement to *protect data from unauthorized access*.

154 The second group of reference requirements covers data analysis. At the highest level we
155 identify goal 2, *answer research questions*. It is associated downstream with goal RG4,
156 *analyze data*. Goal RG4 is composed of two reference requirements: RR11 *integrate data* and
157 RR13 *define analytical process*. RR11 is associated with RR3, since the technical provision of
158 data within the research network is of great relevance for the integration of data. RR13 has
159 subordinate requirement RR12, *define analytical methods*, which covers the low-level data
160 analysis methods. RR14 and RR15 cover two distinct instances of RR13: RR14 describes
161 *static workflows* with all process steps being fixed. In this case, the order of analytical steps
162 and data sources used cannot be changed by the users of the system. In contrast, RR15
163 considers *dynamic workflows*, allowing users to compose analytical steps and data sources as
164 needed. Since the type of data involved in a dynamic workflow is not known upfront, this
165 reference requirement is more demanding in terms of semantic description of data sources.
166 Precise annotation of data sources is necessary in order to perform automated transformations
167 for matching different data fields.

168 RR13 is further associated with RR10, *show results*. RR10 covers the requirement to present
169 the results of the analysis adequately. Thus, it partially fulfils goal 2, *answer research*
170 *questions*.

171 The reference model for requirements is the basis for a network specific requirements model.
172 We present an example for creating such a model and all following steps in the next section.
173 All goals and requirements from the reference model are mapped to network specific
174 instances. In this process, elements of the reference model are checked for their applicability
175 to the specific research network. Further special requirements of the network are considered at
176 this point as well.

177 The network specific requirements model is then mapped to system properties. These are
178 qualities contributed to the system by different components. At first, we consider abstract
179 components instead of specific products. For example, in a research network, reference
180 requirement RR1 *create data* might be mapped to a system property *automated creation of*
181 *data services*. This property is then mapped to the specific component responsible for the
182 implementation of this property.

183 In a second step, the abstract components are mapped to specific components in accordance
184 with the research network's requirements. Specific components can be preexisting modules
185 with a product character, software development frameworks providing specific functionality,
186 or newly developed components.

187 In a final modelling step a distribution model of the components is created. All components
188 need to be mapped to system resources down to the hardware level. Among others, the
189 following aspects have to be considered in this step:

- 190 • Security: Components with high security requirements should be isolated against
191 other, less sensitive components and thus be run on a separate system node.
- 192 • Performance: All components must be distributed in a way that availability of
193 sufficient system resources is ensured.

- 194
- Maintainability: To ensure that the possibly complex distributed system can be managed efficiently, components should be grouped together in a sensible way.
- 195

196 **Sample Application: pelican**

197 We now describe a sample application of our reference model within the research network
198 SFB/TRR77 on liver cancer. Further, we describe two specifications for metadata we
199 developed for the research network.

200 **Specific model for requirements**

201 In this section we summarize key requirements specific to our research network. The
202 complete list of requirements is shown in supplementing information S2.

203 The first goal of the research network, an instance of reference goal RG1, is defined by its
204 research assignment of gaining a deeper understanding of the molecular basis of liver cancer
205 development. This spans research on the chronic liver disease to progression of metastatic
206 cancer. Further, the research network aims to identify novel preventive, diagnostic and
207 therapeutic approaches on liver cancer. Subordinated to goal G1 is G3, the instance of
208 reference goal R3 regarding the data necessary for the network. Since molecular processes
209 play a major role within the research network, genomic microarray data are of central
210 importance. They are complemented by imaging data like tissue microarray (TMA) data and
211 clinical data.

212 Goal G2, answering research questions, is characterized by the following two questions:

- Which generic or specific mechanisms of chronic liver diseases, especially of chronic virus infections and inflammation mediated processes predispose or initiate liver cancer?
 - Which molecular key events promoting or keeping up liver cancer could act as tumor markers or are promising targets for future therapeutic interventions?
- 217

218 Goal G5 requires making the data available for cross project analysis within the network, but
219 to protect data against unauthorized access at the same time. Especially important to the
220 members of the project is the requirement R8, subordinated to goal G5: The projects
221 contributing data to the network require to keep control over the data in order to ensure proper
222 crediting of their intellectual property. Thus, they require fine-grained rules for data access
223 control. Depending on the type of data, they should be available only to specific members of
224 the network, to all members of the project or the general public.

225 **System architecture**

226 To acknowledge the project's requirement R8 to keep the ownership over their data,
227 federation is the underlying concept of the system architecture. Technically, pelican
228 implements a SOA. All data sources of the projects are transformed into data services and
229 made available to the research network. The data services stay under the control of the
230 contributing project. This can even go as far as running the service on computer hardware on
231 the projects' premises. Data services are complemented by analytical services. All services
232 are described by standardized metadata to help finding appropriate services and allow for

233 automated access to the services' interfaces. Using a web-based user interface, researchers
234 can chain data services and analytical services to answer specific research questions.

235 **Component model**

236 The requirements are mapped to system properties first. In the next step, components are
237 identified to provide these properties as module of the new system. In Table 2 we show the
238 complete chain of mappings from requirements over system features to components. Each
239 component is realized either by a readily available product or by a newly developed module.
240 In Table 3 we give an overview of our components.

241 The portal component provides the user interface to the system. It is implemented using the
242 open source software *Liferay* (<http://www.liferay.com>, accessed: 2014-07-03) (*Sezov, Jr.,*
243 2012). Liferay provides a number of functions affecting several components of our model.
244 Thus, we provide a decomposition of the portal components in Figure 4. One important
245 subcomponent of the portal is the document management system. It is realized by the Alfresco
246 component (<http://www.alfresco.com>, accessed: 2014-07-03) (*Berman, Barnett & Mooney,*
247 2012). The user interface of Alfresco can be integrated into the Liferay portal or be accessed
248 with a separate unified resource locator (URL). The portal provides user management
249 functionality to control access to portal pages and components like portlets (*Java Community*
250 *Process*, 2008). However, the user account information including username, passwords, and
251 others is stored in a separate component using the Lightweight Directory Access Protocol
252 (LDAP). Thus, it is possible for all components of the SOA-network to commonly access the
253 users' identity information.

254 Data services are generated by using caCORE SDK (*Wiley & Gagne*, 2012). With caCORE
255 SDK it is not necessary to program the software for the service in a traditional way. Instead, a
256 UML data model in Extensible Markup Language Metadata Interchange (XMI) notation has
257 to be prepared (*Object Management Group*, 2002; *Bray et al.*, 2006). From this model,
258 caCORE SDK generates several artefacts resulting in a deployment packages for Java
259 application servers like apache tomcat (*The Apache Software Foundation*, 2014). To simplify
260 this process for spreadsheet based microarray data, we developed software tool to generate the
261 XMI file as well. As a result, a service conforming to the web services specification ready for
262 deployment is generated. For the provision of network specific metadata we chose TemaTres
263 to serve our controlled vocabulary in standard formats like SKOS or Dublin Core (*Weibel,*
264 1997; *Miles & Bechhofer*, 2009; *Gonzales-Aguilar, Ramírez-Posada & Ferreyra*, 2012). Our
265 analytical services are backed by the open source language and environment for statistical
266 computing called *R* (*R Core Team*, 2014). *R* is integrated into the services using the *Rserve*
267 component (*Urbanek*, 2003).

268 **Deployment model**

269 In a final modelling step the components are distributed to the physical resources available for
270 the system. In our case we used two servers with a common virtualization layer based on
271 VMware VSphere server. Thus, all nodes in our deployment model represent virtual machines
272 (VM). Using virtual switches, routers, and firewall appliances we were able to implement our
273 Internet Protocol (IP) network infrastructure. To enhance security, we implement a network
274 zoning model comprised of an internet zone, a demilitarized zone and an internal zone. Figure

275 5 shows the deployment model in UML notation. The services shown in the model (service 1
276 to service n) are to be considered as examples, since the concrete number of services is
277 permanently changing. The deployment model also reflects the different levels of control that
278 can be executed by the owners of the data. They range from shared nodes on the common
279 servers over a dedicated VM to deployment on external hardware controlled by the respective
280 projects.

281 **Discussion**

282 In this manuscript, we describe a reference model for the requirements of research networks
283 towards an IT platform. For many funding programs including research grants of the
284 European Commission the collaboration of several research organizations at different sites is
285 mandatory. This leads to a structural similarity to our research network on liver cancer. Even
286 though other research networks will have different research aims, there are still requirements
287 that are common to most networks. Since the reference model already covers a basic set of
288 requirements, it allows future research networks to focus on defining specific requirements
289 distinguishing them from other networks.

290 Users of the reference model are responsible for assessing the reference model's applicability
291 to their project-specific needs. The reference model is based on data of a real research
292 network that were generalized. To avoid bias in the model that might hinder transferability,
293 we incorporated different views in the process of constructing the model. However, the
294 transferability of the model to another context is, as for any model, limited. As a consequence,
295 future research networks will have to derive a project specific instance of the reference model
296 to reflect the corresponding characteristics of the project. The reference model is a tool,
297 intended to help its users to create a concrete model covering the requirements of a research
298 network with a high degree of completeness. The reference model provides guidance for this
299 task. We expect, it helps reducing the effort to acquire all requirements.

300 We applied the reference model successfully to a research network on liver cancer. Some
301 specific requirements in this network led to the decision to set up a federated system allowing
302 for a maximum of control of the individual projects over their respective data. The system
303 was implemented as a service-oriented architecture using, among others, components of the
304 caBIG project. Other projects can benefit from this architecture as well, but the architecture is
305 tailored to research networks with the requirements of federating data as data services. With
306 this architecture, we try to acknowledge the data protection requirements of the participating
307 projects. Still, further research regarding the use of data and crediting creatorship of data is
308 necessary. First steps were made as part of this project (*He et al., 2013*).

309 In case the requirements regarding data control are more relaxed, an alternative would be to
310 keep the data in a central data warehouse instead of the federation. In that case, i2b2 might be
311 a suitable component to provide the data warehouse component (*Murphy et al., 2007*). Such a
312 centralized approach also affects, how and when data are harmonized: In a central research
313 data warehouse data are harmonized at the time of loading the database which ideally leads to
314 a completely and consistently harmonized data base. In a service-oriented approach data are

315 provided by means of data services as they are. All services are described by corresponding
316 metadata enabling automated transformation of the data at time of access.

317 Our sample research network concentrates more on basic research than clinical application. In
318 the future, we plan to apply our reference model to further projects with a stronger
319 translational component. By doing so, we will be able to reevaluate the framework in a more
320 clinical context.

321

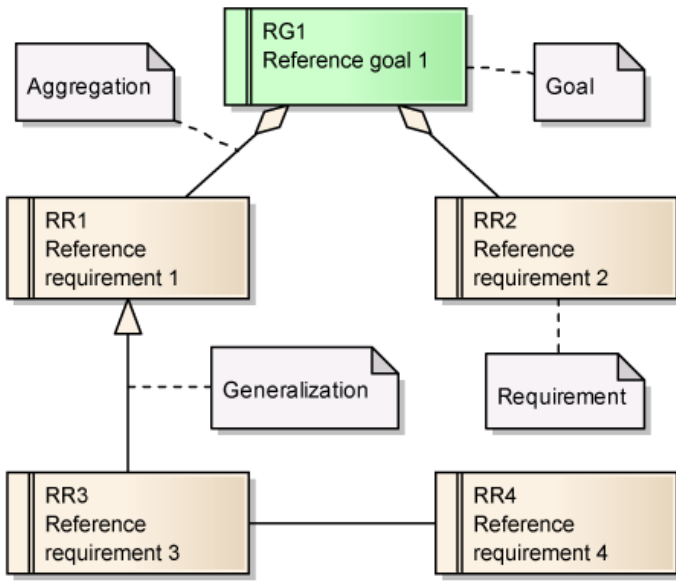


Figure 1 Overview of the UML elements used in requirements diagrams.

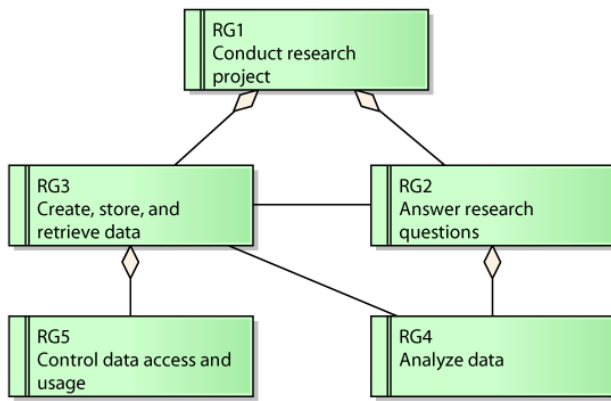


Figure 2 Reference model for goals of a research network.

324

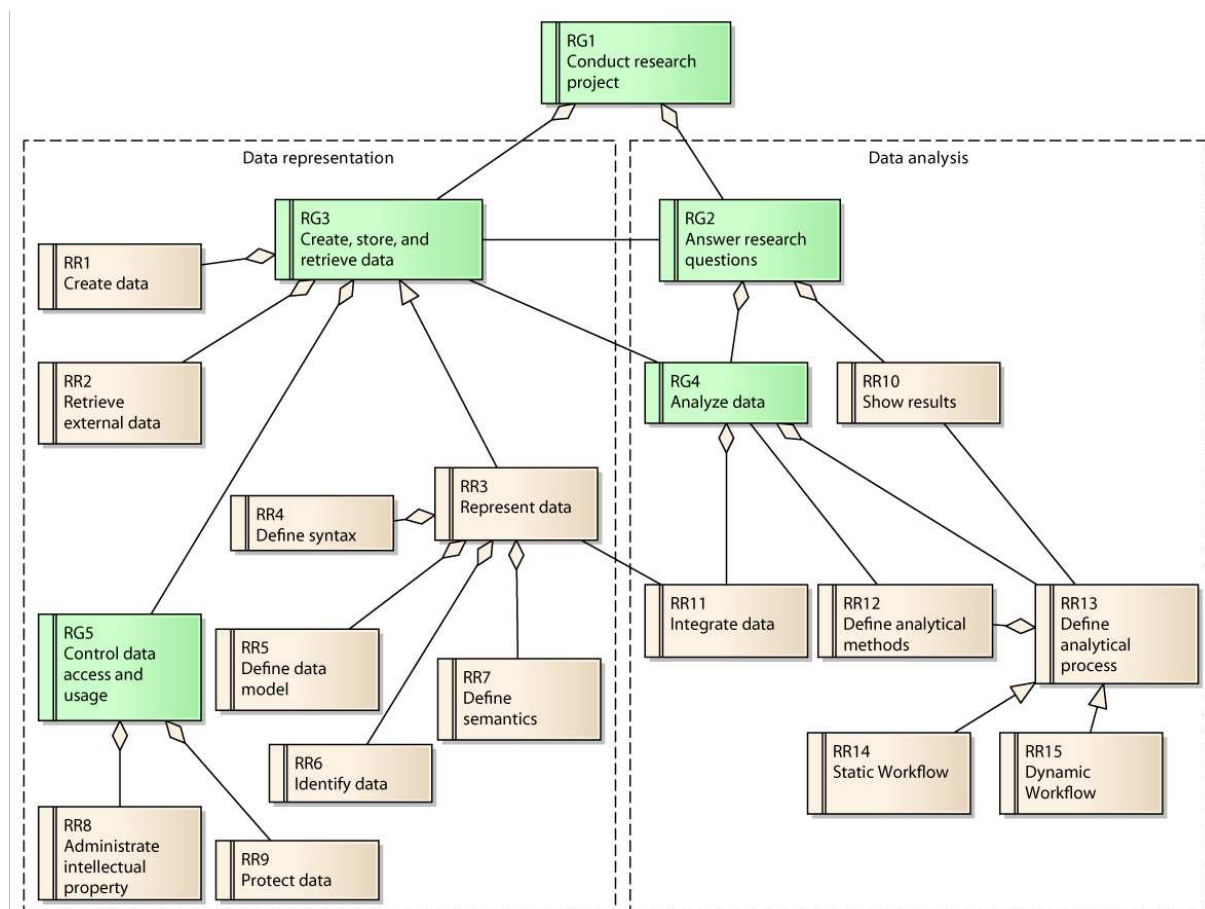


Figure 3 Reference model for goals and requirements of a research network.

325

326

327

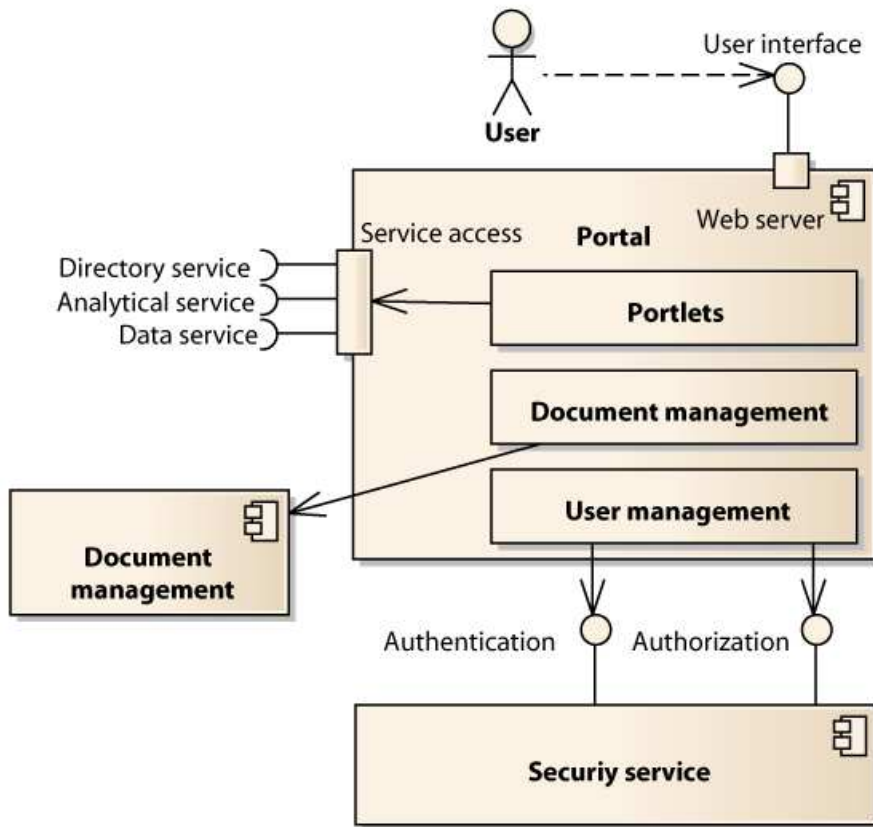


Figure 4 Structure of the component *portal*

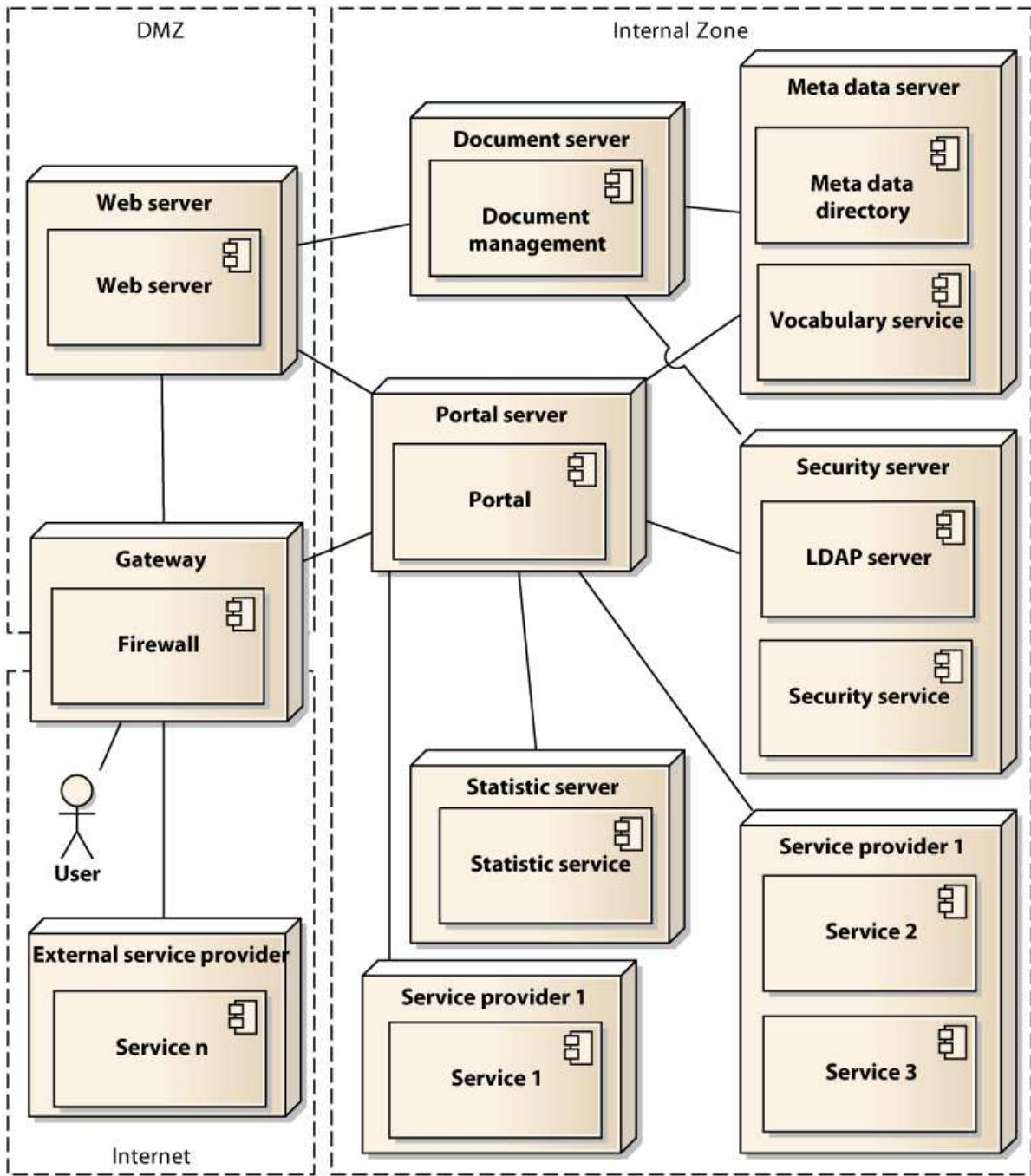


Figure 5 Deployment diagram of the components of the architecture in UML notation.

332 **Tables**

333 **Table 1. Schema for documenting reference requirements and goals.**

Feature	Explanation
Number	Number for uniquely identifying requirements
Name	Name of the requirement
Description	Verbal description of the requirement and its properties
Weighting	Importance of the requirement for fulfilling the goals of the project (low, medium, high)

334

335 **Table 2 Mapping of requirements to corresponding system features and components.**

Requirement	System feature	Component
Requirements associated with G2		
R1 Create data	Automated creation of data services	Data service framework
R2 Retrieve external data	Integration of external data services	Data service framework, portal
R3 Represent data	Data service, document service	Data service framework, document management system
R4 Define syntax	Service description	Data service framework
R5 Define data model	Defined data model	Data service framework
R6 Identify data	Provisioning of information on service location	Meta data directory
R7 Define semantics	Definition of controlled vocabulary and ontologies	Terminology server
Requirements associated with G4		
R8 Administrate intellectual property	Log data usage	Portal, data service
R9 Protect data	User authentication User authorization	Portal, security service Portal, data service, security service
Requirements associated with G5		
R10 Show results	Data specific portlets	Portal
R11 Integrate data	Analytical services	Portal, statistics service, Data service framework
R12 Define analytical methods	Statistical methods	Statistics service
R13 Define analytical process	Documentation service	Document management system
R14 Static Workflow	Workflow in portal application	Portal
R15 Dynamic Workflow	Flexible pipeline	Pipeline management

336

Abstract component	Implementing component
Portal	Liferay
Data service framework	caCORE SDK
Meta data directory	Internal development (based on caCORE SDK)
Terminology server	TemaTres
Security service	caCORE SDK, LDAP
Statistics service	R
Document management system	Alfresco
Pipeline Management	Galaxy (planned)

338

References

339

340 **Baker M. 2010.** Next-generation sequencing: adjusting to data overload. *Nature Methods*
 341 **7 (7):**495–499.

342 **Berman AE, Barnett WK, Mooney SD. 2012.** Collaborative software for traditional and
 343 translational research. *Human genomics* **6:**21.

344 **Bosin A, Dessì N, Pes B. 2011.** Extending the SOA paradigm to e-Science environments.
 345 *Future Generation Computer Systems* **27 (1):**20–31.

346 **Bray T, Paoli J, Sperberg-McQueen CM, Maler E, Yergeau F, Cowan J. 2006.**
 347 Extensible Markup Language (XML) 1.1 (Second Edition). World Wide Web
 348 Consortium. Available at <http://www.w3.org/TR/xml11/> (accessed 9 December 2014).

349 **Cimino JJ. 1998.** Desiderata for controlled medical vocabularies in the twenty-first
 350 century. *Methods of information in medicine* **37 (4-5):**394–403.

351 **Frey LJ, Maojo V, Mitchell JA. 2007.** Bioinformatics linkage of heterogeneous clinical
 352 and genomic information in support of personalized medicine. *Yearbook of medical*
 353 *informatics:*98–105.

354 **Ganzinger M, Noack T, Diederichs S, Longerich T, Knaup P. 2011.** Service oriented
 355 data integration for a biomedical research network. *Studies in health technology and*
 356 *informatics* **169:**867–871.

357 **Gonzales-Aguilar A, Ramírez-Posada M, Ferreyra D. 2012.** TemaTres: software para
 358 gestionar tesoros. *El Profesional de la Información* **21 (3):**319–325.

359 **He S, Ganzinger M, Hurdle JF, Knaup P. 2013.** Proposal for a data publication and
 360 citation framework when sharing biomedical research resources. *Studies in health*
 361 *technology and informatics* **192:**1201.

362 **ISO/IEC/IEEE. 2011.** *Systems and software engineering--vocabulary: Ingénierie des*
 363 *systèmes et du logiciel--vocabulaire.* ISO/IEC/IEEE 24765. Geneva, Switzerland,
 364 New York: ISO/IEC; Institute of Electrical and Electronics Engineers.

365 **Java Community Process. 2008.** Java™ Portlet Specification Version 2.0. Available at
 366 <http://www.jcp.org/en/jsr/detail?id=286> (accessed 9 December 2014).

367 **Komatsoulis GA, Warzel DB, Hartel FW, Shanbhag K, Chilukuri R, Fragoso G,**
 368 **Coronado S de, Reeves DM, Hadfield JB, Ludet C, Covitz PA. 2008.** caCORE
 369 version 3: Implementation of a model driven, service-oriented architecture for
 370 semantic interoperability. *Journal of biomedical informatics* **41 (1):**106–123.

371 **Kunz I, Lin M, Frey L. 2009.** Metadata mapping and reuse in caBIG™. *BMC*
 372 *bioinformatics* **10 (Suppl 2):**S4.

373 **Mathew JP, Taylor BS, Bader GD, Pyarajan S, Antoniotti M, Chinnaiyan AM,**
 374 **Sander C, Burakoff SJ, Mishra B. 2007.** From Bytes to Bedside: Data Integration

- 375 and Computational Biology for Translational Cancer Research. *PLoS Computational*
376 *Biology* **3** (2):e12.
- 377 **Miles A, Bechhofer S. 2009.** SKOS simple knowledge organization system reference.
378 World Wide Web Consortium. Available at <http://www.w3.org/TR/skos-reference/>
379 (accessed 9 December 2014).
- 380 **Murphy SN, Mendis M, Hackett K, Kuttan R, Pan W, Phillips LC, Gainer V,**
381 **Berkowicz D, Glaser JP, Kohane I, Chueh HC. 2007.** Architecture of the open-
382 source clinical research chart from Informatics for Integrating Biology and the
383 Bedside. *AMIA ... Annual Symposium proceedings / AMIA Symposium. AMIA*
384 *Symposium*:548–552.
- 385 **Object Management Group. 2002.** OMG XML Metadata Interchange (XMI)
386 Specification: Version 1.2. Available at [www.omg.org/cgi-bin/doc?formal/02-01-](http://www.omg.org/cgi-bin/doc?formal/02-01-01.pdf)
387 [01.pdf](http://www.omg.org/cgi-bin/doc?formal/02-01-01.pdf) (accessed 9 December 2014).
- 388 **Object Management Group. 2012.** Information technology - Object Management Group
389 Unified Modeling Language (OMG UML), Infrastructure. Version 2.4.1. Available at
390 <http://www.omg.org/spec/UML/2.4.1/Infrastructure/PDF> (accessed 9 December
391 2014).
- 392 **Papazoglou MP, Traverso P, Dustdar S, Leymann F. 2008.** Service-Oriented
393 Computing: a Research Roadmap. *Int J Coop Info Syst (International Journal of*
394 *Cooperative Information Systems)* **17** (02):223–255.
- 395 **R Core Team. 2014.** R: A Language and Environment for Statistical Computing.
- 396 **Sezov R, Jr. 2012.** *Liferay in action: The official guide to Liferay portal development.*
397 Shelter Island: Manning.
- 398 **Taylor KL, O'Keefe CM, Colton J, Baxter R, Sparks R, Srinivasan U, Cameron MA,**
399 **Lefort L.** A service oriented architecture for a health research data network. In:
400 *Proceedings. 16th International Conference on Scientific and Statistical Database*
401 *Management, 2004*, 443–444.
- 402 **The Apache Software Foundation. 24.11.2014.** Apache Tomcat - Welcome! Available
403 at <http://tomcat.apache.org/> (accessed 2 December 2014).
- 404 **Urbanek S. 2003.** Rserve: A Fast Way to Provide R Functionality to Applications. In:
405 Hornik K, Leisch F, Zeileis A, eds. *Proceedings of the 3rd International Workshop on*
406 *Distributed Statistical Computing (DSC 2003)*. Wien.
- 407 **Wei Y, Blake MB. 2010.** Service-Oriented Computing and Cloud Computing:
408 Challenges and Opportunities. *IEEE Internet Computing* **14** (6):72–75.
- 409 **Weibel S. 1997.** The Dublin Core: A Simple Content Description Model for Electronic
410 Resources. *Bulletin of the American Society for Information Science and Technology*
411 **24** (1):9–11.
- 412 **Welsh E, Jirotko M, Gavaghan D. 2006.** Post-genomic science: cross-disciplinary and
413 large-scale collaborative research and its organizational and technological challenges
414 for the scientific research process. *Philosophical Transactions of the Royal Society A:*
415 *Mathematical, Physical and Engineering Sciences* **364** (1843):1533–1549.
- 416 **Wiley A, Gagne B. 2012.** caCORE SDK Version 4.3 Object Relational Mapping Guide.
417 Available at
418 [https://wiki.nci.nih.gov/display/caCORE/caCORE+SDK+Version+4.3+Object+Relati](https://wiki.nci.nih.gov/display/caCORE/caCORE+SDK+Version+4.3+Object+Relational+Mapping+Guide)
419 [onal+Mapping+Guide](https://wiki.nci.nih.gov/display/caCORE/caCORE+SDK+Version+4.3+Object+Relational+Mapping+Guide) (accessed 9 December 2013).
- 420 **Woll K, Manns M, Schirmacher P. 2013.** Sonderforschungsbereich SFB/TRR77:
421 Leberkrebs. *Der Pathologe* **34** (S2):232–234.
- 422 **World Medical Association. 2013.** World Medical Association Declaration of Helsinki:
423 ethical principles for medical research involving human subjects. *JAMA* **310**
424 (20):2191–2194.