

## Treatment Data and Technical Process Challenges for Practical Big Data Efforts in Radiation Oncology

Mayo CS<sup>1</sup>, Phillips, M<sup>2</sup>, McNutt T<sup>3</sup>, Palta J<sup>4</sup>, Dekker A<sup>5</sup>, Miller RC<sup>6</sup>, Xiao Y<sup>7</sup>, Moran JM<sup>1</sup>, Matuszak MM<sup>1</sup>, Gabriel P<sup>8</sup>, Ayan AS<sup>9</sup>, Prisciandaro J<sup>1</sup>, Thor M<sup>9</sup>, Dixit N<sup>10</sup>, Popple R<sup>11</sup>, Killoran J<sup>12</sup>, Kaleba E<sup>1</sup>, Kantor M<sup>13</sup>, Ruan D<sup>13</sup>, Kapoor R<sup>3</sup>, Kessler M<sup>1</sup>, Lawrence T<sup>1</sup>

1)University of Michigan, Ann Arbor, MI,2) University of Washington, Seattle, WA, 3) Johns Hopkins University, Baltimore, MD , 4) Virginia Commonwealth University, Richmond VA, 5) Department of Radiation Oncology (MAASTRO), GROW School for Oncology and Developmental Biology, Maastricht University Medical Centre+, Maastricht, The Netherlands, 6) Mayo Clinic, Jacksonville, FL,6) University of Washington, Seattle, WA,7) University of Pennsylvania, Philadelphia, PA,8) Ohio State University, Columbus, OH, 9) Memorial Sloan Kettering Cancer Center, New York, NY,10) University of California at San Francisco, San Francisco CA, 11) University of Alabama at Birmingham, Birmingham, AL,12) Harvard Medical School, Boston, MA, 13) MD Anderson Cancer Center, Houston, TX, University of California at Los Angeles, Los Angeles, CA

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1002/mp.13114](https://doi.org/10.1002/mp.13114)

This article is protected by copyright. All rights reserved

**Corresponding Author** – Charles Mayo, University of Michigan, Ann Arbor, MI

cmayo@med.umich.edu

**Conflicts of Interest** – There are no conflicts of interest

**Acknowledgements**

Thanks to acknowledge Zhong, Haoyu, MS for figure 1

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

**Received Date:**

**Revised Date:**

**Accepted Date:**

**Article Type: Special Issue Paper**

**Treatment Data and Technical Process Challenges for Practical Big Data Efforts in Radiation  
Oncology**

Mayo CS<sup>1</sup>, Phillips, M<sup>2</sup>, McNutt T<sup>3</sup>, Palta J<sup>4</sup>, Dekker A<sup>5</sup>, Miller RC<sup>6</sup>, Xiao Y<sup>7</sup>, Moran JM<sup>1</sup>, Matuszak MM<sup>1</sup>, Gabriel P<sup>8</sup>, Ayan AS<sup>9</sup>, Prisciandaro J<sup>1</sup>, Thor M<sup>9</sup>, Dixit N<sup>10</sup>, Popple R<sup>11</sup>, Killoran J<sup>12</sup>, Kaleba E<sup>1</sup>, Kantor M<sup>13</sup>, Ruan D<sup>13</sup>, Kapoor R<sup>3</sup>, Kessler M<sup>1</sup>, Lawrence T<sup>1</sup>

1)University of Michigan, Ann Arbor, MI,2) University of Washington, Seattle, WA, 3) Johns Hopkins University, Baltimore, MD , 4) Virginia Commonwealth University, Richmond VA, 5) Department of Radiation Oncology (MAASTRO), GROW School for Oncology and Developmental Biology, Maastricht University Medical Centre+, Maastricht, The Netherlands, 6) Mayo Clinic, Jacksonville, FL,6) University of Washington, Seattle, WA,7) University of Pennsylvania, Philadelphia, PA,8) Ohio State University, Columbus, OH, 9) Memorial Sloan Kettering Cancer Center, New York, NY,10) University of California at San Francisco, San Francisco CA, 11) University of Alabama at Birmingham, Birmingham, AL,12) Harvard Medical School, Boston, MA, 13) MD Anderson Cancer Center, Houston, TX, University of California at Los Angeles, Los Angeles, CA

29

30

31

32 **Corresponding Author** – Charles Mayo, University of Michigan, Ann Arbor, MI

33 cmayo@med.umich.edu

34

35 **Conflicts of Interest** – There are no conflicts of interest

36

37 **Acknowledgements**

38 *Thanks to* acknowledge Zhong, Haoyu, MS for figure 1

39

40 **Abstract**

41 The term *Big Data* has come to encompass a number of concepts and uses within medicine.

42 This paper lays out the relevance and application of large collections of data in the radiation

43 oncology community. We describe the potential importance and uses in clinical practice. The

44 important concepts are then described and how they have been or could be implemented are

45 discussed. Impediments to progress in the collection and use of sufficient quantities of data are

46 also described. Finally, recommendations for how the community can move forward to achieve

47 the potential of Big Data in radiation oncology are provided.

48

49 **Introduction**

50 To the clinician, it often seems that we have too much and too little data at the same time. We

51 spend more time than we would like at computer terminals entering or reading data. Perhaps

52 it would be better stated that we would like the *data* we input to be transformed into

53 *information* that we can use. This is the aspect of *Big Data* that this manuscript addresses.  
54 Computerized data handling has been an integral part of our field since the introduction of  
55 computerized treatment planning and record and verify systems. The question is, now that  
56 there are highly successful algorithms for using computerized data to make models for  
57 predictive purposes, can the radiation oncology community harness our data for our patients'  
58 benefit?

59 Pan et al. have provided a very clear picture of the difficulties that we face in collecting and  
60 using data in the clinic [1]. The questions we must answer are: (a) is it worth making an effort  
61 to improve the situation, (b) what are the details of the clinical data environment that need to  
62 be addressed, and (c) how do we accomplish our goals? An AAPM Science Council Focused  
63 Research Meeting (FOREM) meeting, jointly sponsored with vendors, was held in Ann Arbor in  
64 May of 2017, to address these questions. In this publication we provide an overview and  
65 summary of the answers that emerged.

66

## 67 **Motivation for embracing Big Data**

### 68 **a. Need to learn from and adapt to emerging therapies such as genetics, immunotherapy**

69 It is now commonly understood that the explosion of data and knowledge that has resulted  
70 from genomics will have a great impact on all areas of cancer care, including radiation therapy.  
71 A patient's genetic profile may play an important role in how they will react to certain agents or  
72 in their ability to repair radiation damage [2]. The tumor's genetic profiles (since many tumors  
73 have a multitude of different mutations) is increasingly being used to determine the best  
74 therapy or combination of therapies [3].

75 Immunotherapy is another area of increasing importance. The ability to use different aspects  
76 of the immune system to target tumor cells is an area of great current interest [4].

77 Radiation oncology is not alone in the interest and need for better data on patients' genetic  
78 profiles. NIH has been working with a number of groups to establish a workable solution in

79 order to avoid the current problems such as laboratory-dependent formats, text-based storage,  
80 and lack of centralized storage in current electronic health records (EHR) [5].

81

## 82 **b. Cancer as chronic disease and multiple care givers**

83 As cancer therapy becomes more effective, more and more patients are living longer. As a  
84 result, the extent and complexity of information which needs to be tracked to improve  
85 understanding of outcomes is increasing. For example, for patients who are essentially cancer  
86 free, monitoring risk for treatment-related complications when their long term home location  
87 based follow up is not at the treatment center is a challenge. Parry et al. estimate that there  
88 will be 18 million cancer survivors by 2020 [6]. In addition, there are the increasing numbers of  
89 patients who survive longer than ever due to improvements in targeted therapies, better  
90 imaging and better methods for localizing dose [7]. These advances can lead to improved local  
91 control and better control of oligometastases. The upshot is that as the number of patients  
92 who suffer cancer-related health consequences increases over time, the more likely it is that  
93 they will see a wider spectrum of specialists and in a larger number of clinical settings,  
94 interacting with a large variety of recording-keeping systems.

95 Even just considering the electronic health records, there are no general standards for the  
96 selection and formatting of data to be recorded. Different vendors, different institutions,  
97 different departments and even different physicians have different methods which are often  
98 not compatible. Finally, even within well-structured organizations, much of the data exist  
99 within text documents. Lack of standards for which data elements to gather, inconsistent  
100 processes for entry and variability among commercial systems for aggregation and reporting  
101 increase the likelihood that physicians and staff will miss information or have incorrect  
102 information regarding a patient's health and/or treatments that could potentially affect  
103 decisions.

104

## 105 **c. Comparative Effectiveness Research**

106 In the last decade, comparative effectiveness research (CER) has come to be seen as an  
107 important and necessary adjunct to randomized clinical trials (RCT) [8]. In CER, two different  
108 therapies or tests that are already accepted are compared, whereas RCT's focus on comparing a  
109 new to a current therapy. The Patient Centered Outcomes Research Institute cites CER as its  
110 primary method of research. Given the relatively small numbers of cancer patients that are  
111 enrolled in RCT's (approximately 3%), the need to use the information that is available through  
112 CER is understandable.

113 Comparative effectiveness research can be tailored along a spectrum of methods ranging from  
114 essentially an RCT to a comparison of current clinical practice with an integrated practice  
115 beyond the current norm. A recent paper by Fiore et al. looked at four different trials that  
116 sought to use only data in the current EHR's [9]. Their conclusions included: "We find that EHR-  
117 based clinical trials are feasible but pose limitations on the questions that can be addressed, the  
118 processes that can be implemented, and the outcomes that can be assessed."

119 Clearly, for progress to be made using CER practical methods for the easy and accurate  
120 collection of data and for the sharing of data must be available in clinics.

#### 121 **d. Quality Improvement and Error Detection**

122 The past few years have seen an explosion in the use of data to reduce errors in radiation  
123 therapy. ASTRO and AAPM have implemented the Radiation Oncology—Incident Learning  
124 System (RO-ILS) that relies on data submitted to it to develop a shared learning platform.  
125 While this system is not "big data" in the sense that it is in text format and is a relatively small  
126 amount of data, it does count in our definition of transforming data to information. In  
127 particular, the system is set up to provide users with more knowledge about the sources of  
128 errors and how best to avoid them. Another area is in artificial intelligence applications of error  
129 detection. For example, Kalet et al. successfully mined an OIS to develop a probabilistic model  
130 of the contributing factors to errors [10].

#### 131 **e. Modeling in Radiation Oncology**

132 Perhaps the most widespread use of data in radiation oncology is in modeling. The examples  
133 are too numerous to list, but some of the most impactful models are the QUANTEC models,  
134 outcomes, tumor control probabilities, equivalent uniform dose, and biologically effective dose  
135 [11]. As construction of Big Data Analytics Resource Systems (BDARS) aggregating a wider  
136 range of health care information (e.g. labs, medications, genomics, demographics, patient  
137 reported outcomes (PROs) etc.) expands, more comprehensive models are progressing beyond  
138 dose metrics alone [12-14]. In addition, heuristic type models have been constructed for  
139 automating the objectives of inverse planning and library-based contouring. A promising area  
140 for the more conventional use of big data is in machine learning for automated contouring. In  
141 this application, images that have been segmented by experts are fed into a machine learning  
142 algorithm and image features that predict the true contours are selected to produce anatomical  
143 contour models.

144

145

## 146 **State of the Data**

147 One of the most important concepts is that Big Data, in most cases, implies more data than may  
148 be obtained by any single institution. In order to use machine learning or any modeling  
149 techniques, there must be enough data to (a) build the model, (b) test the model, and (c)  
150 validate the model. Optimally, validation (c) can be done with data from a different institution  
151 in order to account for hidden variables that may not be appreciated [15]. In addition, as our  
152 ability to differentiate patients improves, e.g. genomics and radiomics, the number of patients  
153 suitable for any given model decreases, thereby increasing our need for more comprehensive  
154 capture of intra-institutional data as well as for multi-institutional data. This has critical  
155 implications for how organizations cooperate. Whereas success in medical research in the past  
156 has favored very large single institutions that can develop a critical mass of knowledge and  
157 resources in close physical proximity, diffuse networks of institutions able to generate and  
158 share information will have an advantage in the future.

159



160 In addition to the need for broad (many patients) data sources, we also need deep  
161 (relationships among key data elements) sources. Systems promoted as big data sources may in  
162 fact be shallow, capturing only a few data elements for a large number of patients. For  
163 example, some data sources draw upon billing records or imaging records for a large number of  
164 patients, but lack depth needed to enable linkage to diagnosis, treatment, dosimetric or  
165 outcomes details. Another impediment to obtaining the "deep" type of data is that sources  
166 often dump unstructured, "as is", data into data lakes where key data elements and  
167 relationships can in principle be extracted, but in practice carry a high overhead for extraction.  
168 Challenges for ensuring depth in aggregation of key data elements needed for radiation  
169 oncology fall into four categories

- 170 ○ Access – Staff possessing both domain knowledge of radiation oncology and of  
171 informatics need access to query data bases in source systems to construct  
172 functional big data repositories.
- 173 ○ Data Integrity – Data elements that may not require accurate entry to enable  
174 treatment but are vital for correctly identifying specific patient groups in practice  
175 quality improvement (PQI) and research efforts require changes in clinical processes  
176 to assure validity. This often implies a cultural shift to prioritize recording data in  
177 recoverable formats.
- 178 ○ Data Structure – The cost of free text is high. Lack of standardized structure for entry  
179 undermines ability to automate extraction of key data elements from text fields such  
180 as notes. To assure accurate, high volume, electronic extraction of key data  
181 elements standardized methods for encoding key data elements need to be defined  
182 and implemented in clinical processes.
- 183 ○ Lack of integration among systems – Key data elements are entered and stored in a  
184 range of commercial systems that frequently do not maintain linkages needed to  
185 identify relationships between key data elements. There is no existing standard of  
186 practice to link departmental data sets with radiation oncology-specific content with

187 large commercial and governmental datasets such as the National Cancer Database  
188 Base.

189

190

## 191 **Process and system changes**

192 In reviewing current practices, a number of obstacles stand in the way of obtaining the amount  
193 and quality of data needed to make substantial progress. The following outline provides a view  
194 that is geared towards identifying means of overcoming them.

195 *(1) Failure to collect necessary structured data*

196 *(2) Lack of data standardization*

197 *(3) Inability of different electronic data systems to communicate.*

198 Within each of these broad categories, it is useful to provide a finer-grained view of how  
199 different aspects of our clinical and electronic environments contribute to the overall difficulty  
200 in achieving the data collection and use that we seek.

### 201 (1.a) Commercial System Databases

202 Focus for development of commercial systems that store the range of data needed for clinical  
203 data repositories is often on the user interfaces rather than on the back-end databases. The  
204 situation is similar to a clinical focus on data required to treat the day's patients and support  
205 billing documentation with few resources devoted to standardizations and optimizations to  
206 increase big data extractions. Individual systems may use multiple loosely connected databases,  
207 complex compound keys, lack of indexing, poorly designed schema, lack reasonable security, or  
208 use non-standard database technologies. Vendors may also refuse to provide end-users access  
209 to extract their own data. Some commercial systems are much better than others, so end user  
210 experience is variable.

### 211 (1.b) Diagnosis and staging

212 Correct usage and quantified entry of diagnosis and staging information is central to many PQI  
213 and research efforts. For example, incorrect entry of primary disease codes (e.g. prostate 185,  
214 C61) when treating subsequent bone (C79.51), brain (C79.31) or lung (C78.00) metastasis and  
215 failure to utilize functionality in radiation oncology information systems (ROIS) to connect  
216 primary and metastatic diagnosis undermine the ability to use these codes to correctly identify  
217 patient groups by codes. Failure to utilize functionality in ROIS connecting treatment courses to  
218 these codes undermines ability to connect treatment elements (e.g. DVH metrics) to patients.  
219 The cost of not taking a few seconds to select ICD-O (International Classification of Diseases for  
220 Oncology) values linked to ICD9 (International Classification of Disease, revision 9) and ICD10  
221 (International Classification of Disease revision 10) in the ROIS means that subsequent  
222 questions about disease site location become prohibitively expensive to answer because of the  
223 manual effort required to retrospectively revisit the chart. When survival information is  
224 obtained from EHRs, failure to utilize functionality in ROIS to enter staging information  
225 undermines ability to factor staging into survival, recurrence and other factors. Typically, EHRs  
226 do not have functionality for quantifying diagnosis and staging information according to  
227 guidelines (e.g. AJCC, FIGO) or to connect primary and metastatic disease. On the other hand,  
228 ROIS generally do, but frequently this functionality is not utilized fully as part of clinical practice.

#### 229 (1.c) Outcomes

230 Patient outcomes such as toxicity and disease site status (e.g. recurrence) are frequently  
231 entered into electronic records as free text using unstandardized terminology. This renders  
232 them unavailable for automated electronic extraction. Lack of standardizations 1) for which  
233 toxicities are routinely measured, 2) how treatment site categorizations are named (e.g. breast  
234 tangents , breast tangents plus supra-clavicular field, breast tangents plus supra-clavicular field  
235 plus internal mammary node field, etc) , 3) how categorizations for disease site status are  
236 named (e.g. no-evidence-of-disease, local recurrence) or 4) in use of regular schemas for text  
237 representation of these key data elements prevent this information from being used to its full  
238 value in routine characterization of outcomes for treated patients.

#### 239 (1.d) “As-Treated Plan Sums”

240 To assess correlation of outcomes with dose volume histogram (DVH) metrics, it is necessary to  
241 first create treatment plan sums corresponding to the plans and number of fractions treated,  
242 reflecting boosts, plan revisions and incomplete treatments. When these “as treated” plan  
243 sums (ATPSs) are created as part of routine practice, then automated solutions for calculating  
244 dose-volume histograms metrics becomes possible. Unfortunately, often these are not created  
245 as part of routine practice, with the result that they must be constructed retrospectively, ad-  
246 hoc, preventing systematic, automated aggregation. Currently no major commercial system, to  
247 our knowledge, has a standard means for reporting ATPSs.

#### 248 (2.a) Prescriptions

249 Electronic prescription summaries that defined dose levels, target structures, number of  
250 treatments, fractionation groups (e.g first course, plan revision, boosts, etc) and connection to  
251 target structures, organs at risk, treated plans and DVH metrics have been developed by a few  
252 researchers [16,17]. These custom solutions were developed to fill the void left by commercial  
253 ROISs. Recently ASTRO has suggested a baseline set of guidelines for information that should be  
254 included in prescriptions to promote standardization [18]. Similar to ATPSs, commercial  
255 solutions and clinical processes often lack ability to retrospectively extract this key information.

#### 256 (2.b) Key Treatment Parameters

257 Ensuring ability to identify which patients were treated with special technologies and details of  
258 those treatments is important to being able to prove their efficacy. Examples include breath  
259 hold technologies, radio frequency or radio-opaque fiducials used for positioning,  
260 immobilization devices, etc. However, commercial systems and clinical approaches to utilizing  
261 those systems are frequently inadequate for retrospectively gathering this data.

#### 262 (3.a) Integration of Treatment Planning System (TPS) with ROIS

263 If systems do not use a common database for TPS and ROIS it is difficult to unambiguously  
264 move from the ROIS record of plans actually treated back to specific plans, plan sums and DVH  
265 curves in the TPS. Some vendors may even discard DICOM Unique Identifiers for plans from the  
266 TPS.

267 (3.b) Integration with EHR

268 ROIS and TPS systems typically do not integrate with EHR's. Connections may be made through  
269 medical record numbers and inferences around dates recorded in respective systems. This is an  
270 area where Health Level 7 (HL7) Fast Healthcare Interoperability Resources (FHIR) could  
271 significantly improve integration.

272 (3.c) Integration with specialty systems

273 Treatment devices other than conventional linear accelerators (e.g. brachytherapy, particles,  
274 specialty accelerators, MR guided linacs) may provide minimal details back to the ROIS or may  
275 use specialty tables in the ROIS that do not integrate well with tables used to manage external  
276 beam therapies. This limits the range of questions around treatment details for these specialty  
277 modalities that can be addressed at large scale for all patients treated.

278 (3.d) Integration with institutional registry data

279 Institutions with the American College of Surgeons Commission on Cancer and National  
280 Comprehensive Cancer Network (NCCN) designations are required to have medical registries  
281 that follow up on cancer patients. Registries document demographics, diagnosis, staging,  
282 survival, cause of death and other factors. Registry data is rarely linked to radiation oncology  
283 data repositories.

284 (3.e) Integration with public databases

285 Institutional registries supply data to state registries. Published state analyses are,  
286 unfortunately, many years behind current practice. Although state registries have high volumes  
287 of patients, there is no simple means to connect back to patients to check on the validity of the  
288 data or to investigate impact of cofactors on outcomes tracked in the registries.

289

## 290 **Access and Extraction Issues**

291 As radiation oncology has developed, a number of structural issues have arisen that limit  
292 clinicians', caregivers' and researchers' access to the data that we do have. Access requires

293 several key elements: knowledge of the format and schema of the stored data, software that  
294 can identify and extract the data, and permissions to view and extract the data.

295 Figure 1 illustrates the level of detail that is needed regarding the treatment of rectal carcinoma  
296 patients under three RTOG studies. To combine the data from these trials requires knowledge  
297 of how the problem is framed (which clinical data are needed, what are the key elements of  
298 those data), how the data are formatted (type of value, allowed values, units, standards if  
299 applicable), and the specific software needed to access the data (SQL, RDF triples,  
300 spreadsheets).

301 The issue of framing the medical problem is difficult but rewards are high. The DICOM standard  
302 (and its radiation therapy extension) has achieved such success in large part due to its  
303 structuring of what an imaging study (radiation treatment) *is*--what are its elements and how  
304 are they related [19]. Thus, regardless of the details of the implementation of a procedure, all  
305 partners in a communication exchange agree on the essential elements. The definition of such  
306 standards in other areas of medicine is rapidly increasing. For example, a relatively commonly  
307 used standard for data exchanges between EHR's is the standard Health Level 7 (HL7). HL7  
308 version 2 standardized types of data and the allowed values and permitted organizations and  
309 vendors to develop software for the reliable interchange of certain data. However, it was  
310 considered to be quite limited, and version 3 was built around the Reference Information  
311 Model which was a much more robust view of healthcare in general [20]. Even more recently,  
312 they have started developing HL7-FHIR which instantiates an even more up-to-date view of  
313 medical practice, but also highlights the importance of appropriate technology. HL7-FHIR is  
314 built upon the REST specification that is the current industry standard for web-based  
315 applications [21]. Other data standards, such as the NCI thesaurus [22], provide additional  
316 resources that facilitate the development of software for access and extraction of data.

317 With rare exception, major vendors of ROIS, TPS and EHR systems, store information in  
318 relational databases. A few types of large volume objects (e.g. DICOM images) are stored in files  
319 that are referenced in the relational databases. Custom extractions from databases are carried  
320 out using structured query language (SQL). SQL queries may have dialectical variation among

321 relational database systems (e.g. Oracle, Microsoft SQL). Ideally, relational databases are  
322 designed with categories of data grouped into tables and views (stored SQL query results)  
323 reflecting an overall view of the procedure itself. They also use normalization strategies to  
324 prevent redundant information, reduce complexity in SQL queries and increase performance in  
325 retrieving data. Secure data retrieval requires granting read access to specific authenticated  
326 network accounts. Access may be controlled at the level of the database, table or views. Skill  
327 with SQL is essential to any staff constructing or extracting data for a data repository.

328 Application programming interfaces (APIs) are provided by vendors of many TPSs. These may be  
329 used to gather subsets of information stored in the ROIS database or elements only calculated  
330 at run time in the TPS (e.g. DVH curves for some systems). APIs allow custom software  
331 applications to be constructed by users that interface with the TPS. Access is controlled by end  
332 user system administrators, subject to constraints of the commercial system. Clinical staff  
333 members with coding skills are necessary for effective use of API's.

334 Legacy issues with vendor changes to both database and API structures are an issue for groups  
335 automating extraction from electronic records systems. Effort to re-write queries and scripts  
336 when systems are upgraded can be substantial.

337 Patient reported outcomes (PROs) are important outcome measures and their routine  
338 monitoring during cancer therapy has been demonstrated to improve survival [23]. However,  
339 use of paper based rather than electronic systems are more common. Electronic systems are  
340 significantly better for making the data accessible, but require substantial effort in setting up  
341 systems and arranging for staffing resources to assist patients with completing electronic  
342 surveys is required. In addition, lack of standardization in instruments to be used, redundant  
343 questions between surveys, excessive length diminishing patient willingness to participate, and  
344 question formats and logic that translate poorly to electronic systems already used in patient  
345 work flows are issues for generalized use of PROs.

346 Diagnostic images are stored on Picture Archive and Communication Systems (PACS) in Digital  
347 Imaging and Communication in Medicine (DICOM) format and accessed with DICOM servers.  
348 Graphic user interfaces for clinical use are not well suited to large volume, batch access of sets

349 of patient images. The objective in utilizing these resources in connection with BDARS is not  
350 creation of a parallel PACS. Instead, when large sets of images are identified for utilization in a  
351 study, e.g. developing predictive radiomics measures for a disease type, downloading a large  
352 specific set of images for batch processing is needed. Negotiating access is the primary barrier.

353 Finally, it is important to discuss the role that legal and commercial considerations play in  
354 limiting access to data. The Health Insurance Portability and Accountability Act of 1996  
355 requires certain standards to be met when exchanging private health information. The  
356 standards depend on the intended use of the data, for example, clinical decisions, insurance  
357 coverage, quality improvement and research. They also depend on the entities exchanging the  
358 information. These standards add time, effort and new procedures to any effort to obtain data  
359 access. Intra-institutional exchange, for example between a departmental data repository and  
360 the hospital EHR, is in general easier than between institutions, but even that type of  
361 transaction usually requires some level of administrative oversight and/or procedure. In  
362 addition, storing data in a clinical data repository for possible future research can be viewed as  
363 problematic under national ethics guidelines for human research [24]. Overall, it is difficult to  
364 make any broad statements or recommendations regarding these issues since they are, to  
365 some degree, institution- and use-specific. In addition, how the regulations are interpreted is  
366 evolving, particularly in response to some of the national healthcare programmatic initiatives  
367 such as the Affordable Care Act.

368 **Selecting technologies**

369 The objective is to use the treatment data, rather than to utilize a novel database technology.  
370 Selecting database technologies which minimize investment overhead and risk while  
371 maximizing productivity and interoperability for addressing particular tasks requires careful  
372 consideration [25,26].

373 At a high level, four process steps can be considered and technology choices should be made  
374 fit-for-purpose for these steps.

- 375 1. Capture of treatment data



376 The primary use for health care data is delivery of patient care. Health care database  
377 technology is often vendor dependent and under regulatory oversight. For structured data  
378 elements (e.g. record and verify, electronic health records, outcome) relational databases are  
379 the most common technologies. Images and related objects such as treatment plans and record  
380 are generally object stores (e.g. PACS) with a relational schema containing object pointers.

## 381 2. Extraction

382 Since the primary use sources have to be taken as-is, the extraction technologies providing  
383 connectors to these primary sources should be able to handle many different sources and  
384 formats including all common relational sources. They should be able to handle non-relational  
385 sources including “databases” that researchers and physicians often use (e.g. Excel, SPSS) and  
386 include JSON and XML support as these are common export format for more technical users.  
387 Ideally, the technology can be extendible to support common medical standards (HL7v2, HL7v3,  
388 HL7 FHIR and DICOM) as needed.

389 A wide range of programming languages and standard database import tools are frequently  
390 used. These have the advantage of hiding very little from the user. There are also commercial  
391 and open source software systems intended to reduce the technical skill requirements for users  
392 with the trade-off of obscuring details about the extraction, cleaning and loading processes.  
393 Since primary sources change and extraction tools generally expand and change over time, a  
394 crucial requirement is versioning. Users of technology should be able to store different versions  
395 of the extraction scripts and configurations so that subsequent users can re-use their solutions.

## 396 3. Transformation, integration and storage

397 For successful secondary use, the primary use sources need to be combined, integrated and  
398 common data elements mapped on each other. An example is the combination of ROIS/EHR  
399 data (diagnosis, comorbidities, prescriptions, treatments, follow-up), Record and Verify data  
400 (radiotherapy treatment) and DICOM data (imaging/plan). This transformation and integration  
401 is generally the most time consuming task of the process. Knowledge of the primary sources  
402 and of the secondary use data model is a requirement for staff using the tool. Again, versioning  
403 and manageability is crucial as sources change and sharing transformation scripts with others is

404 needed for work to not be duplicated. Defining distinctions between data element categories  
405 and relationships means mapping the raw values onto a schema. For example, a schema needs  
406 to be applied so that we can inform our analytics programs if an extracted value “30”  
407 corresponds to a dose, an age, a day of the month, etc. and how that value relates to other  
408 information e.g. toxicity, survival, PROs, treatment dates, etc.

409 From a technology standpoint two main approaches exist.

410 • Schema-On-Aggregate (aka schema-on-write): Upon extraction each data  
411 element from each source is considered more or less separately, transformed  
412 and mapped to the secondary use data model and then written in the secondary  
413 use data store. Schema-on-aggregate has as its main benefit that it often re-uses  
414 the knowledge contained in the primary use schema and forces one to decide up  
415 front how to map data items and think about transformation for each data  
416 element. The end-result is often a data store with a structured schema.

417 Relational databases are widely used for this approach owing to their speed,  
418 ease of integration with other systems and large pool of talent for use. Non-  
419 relational databases (e.g. object stores, graph databases and triple stores) have  
420 also been used in some research settings.

421 • Schema-On-Query (aka schema-on-read): The secondary use data model is  
422 applied when the secondary user requests, or queries, the data from the  
423 secondary source. In a schema-on-query system the data is stored from the  
424 primary source “as-is” and by necessity this is a non-relational store (e.g. a data  
425 lake). An example is Apache Hive which can be used for SQL-like schema-on-  
426 query for Apache Hadoop. NoSQL databases, such as MongoDB or CouchDB, are  
427 another example. The main benefit of this approach is that the transformation  
428 and secondary use data model can be defined fit-for-purpose, and different for  
429 different use cases. Also all primary use data can be stored immediately for later  
430 secondary use. The main drawback is that knowledge of original schema is often  
431 not available by the time the data is used and that data is stored without de-

432 identification. Variability in nomenclature for key data elements, relationships  
433 and formats among the various “as-is” sources requires creating and maintaining  
434 custom code for each to enable programmatic extraction. Care must be taken to  
435 ensure consistent meaning at the time of data entry so that contents of an  
436 element are internally consistent and stable.

437  
438 Note that many solutions allow a combination of the above approaches, with some data  
439 elements stored in a schema generation upon aggregate and some stored “as-is” for schema at  
440 a later time point. In that case, key data elements are often duplicated into the secondary use  
441 storage.

442

#### 443 *Secondary use application*

444 Secondary use of subsets of data extracted from BDARS to address specific research or clinical  
445 questions is a common use case. The secondary user usually has defined their own data model,  
446 store and the application to analyze the data. The technology choices made by secondary users  
447 vary widely and limited influence exists especially if the secondary user is external to the  
448 primary use institution. The main job of technology here is to provide the secondary end-user  
449 with a dataset and format which he or she can use (often called a data mart). Typical requested  
450 formats include SQL database dumps, Microsoft Excel, comma (or tab) separated values (CSV),  
451 DICOM, HL7 FHIR, HL7v3, HL7v2, XML and JSON. Additionally, data visualization and allowing  
452 the end-user to navigate the data store established in the previous step increase the efficiency  
453 and effectiveness of secondary use. The tools mentioned above generally allow such export to a  
454 variety of data formats. Figure 1 illustrates one such use case, a semantic triple store database  
455 (a.k.a. Resource Description Framework) was applied for the purpose of combining datasets  
456 from several clinical trials. Semantic triples can be used to define a range of relationships  
457 between objects (e.g. PTV → is a type of → target structure).

458

459 **Specific recommendations for work flows and standardizations**

- 460 1) Diagnosis and staging data should be entered into quantified fields in accessible, electronic  
461 systems that
- 462 ◦ have quantified fields for staging elements and overall staging, and staging guideline  
463 system used (e.g. American Joint Committee on Cancer (AJCC))
  - 464 ◦ ensure correct selection of staging from component elements
  - 465 ◦ provide explicit linkage to treatment courses and plans used to treat
  - 466 ◦ link metastatic diagnosis (e.g. C79.51, Secondary malignant neoplasm of bone) to  
467 diagnosis for originating sites (e.g. C34.1, Malignant neoplasm of upper lobe,  
468 bronchus or lung)

469

470 In the current vendor landscape, the ROIS is frequently the only system in the clinical  
471 process workflow meeting these objectives.

472

473 2) Nomenclature standardizations recommended by AAPM Task Group 263 should be adopted  
474 into routine practice. These define standardized nomenclature for structure, target and  
475 DVH metric naming to promote ability to automate aggregation [27].

476

477 3) Course cumulative as-treated plan sums should be constructed as part of routine practice.  
478 Since more than one image set may be used in the construction of the ATPS's, and relative  
479 positioning of structures may vary between sets, using the image set providing the best  
480 representation for the clinical evaluation carried out for treatment is currently the most  
481 viable approach.

482

483 4) Toxicities, recurrence and PRO outcomes need to be routinely collected as quantified fields  
484 (instead of free text fields) in accessible electronic systems. Standardizations for specific  
485 items and values are needed. This includes, for example, definition of recurrence  
486 nomenclature. Ability to automatically recover these values from the electronic record is  
487 important.

488  
489 5) Detailing of key data elements and relationships (i.e. an ontology) is needed for a broad  
490 range of practice quality improvement and translational research efforts. An initial set, drawn  
491 from experience in constructing BDARS, is presented as an appendix to this paper. Success in  
492 gathering this information requires that clinical systems should be utilized to ensure ability to  
493 accurately aggregate these elements and relationships from the electronic record (ROIS, TPS,  
494 EHR). Ideally, professional societies such as ASTRO, AAPM, ESTRO and CARO would combine  
495 efforts to eventually take the role of maintaining standardized ontologies to promote  
496 interoperability among institutions and commercial systems. Combining the ontology presented  
497 in the appendix with related ontologies would be a valuable step toward a common standard  
498 [28,29].

499  
500 6) In addition to demonstrating adherence to standard quality metrics, clinical entities will face  
501 increasing demands for demonstration of the value of the care they deliver as medicine in the  
502 transitions from fee for service to value based payments. Success in the value based payment  
503 environment will require the ability to conduct on-demand analysis of patient and tumor  
504 characteristics, all aspects of treatment delivery, outcomes, and cost of care.

505  
506 We note that the task of creating ATPSs (item 3) needs to begin as soon as possible, guided by  
507 clinical judgment, in order to replace complete lack of data with reasonable data. In addition,  
508 further refinement is needed. Collaborations between professional societies, vendors and  
509 clinical trials groups for defining standards for the end-of-treatment dose composite are

510 needed. Issues include means to quantify quality of the composite, identifying source images,  
511 identifying trade-off decisions in image registrations, uncertainties in structure dosimetric  
512 measures when multiple image sets are used, and realistic appraisal of the role of image  
513 deformation.

## 514 **Examples of Clinical Data Repositories**

515 Several groups have been actively engaged in construction of clinical data repositories (CDR),  
516 also known as data lakes and Big Data Analytic Resource Systems (BDARs). These systems  
517 become important components for both research and clinical practice efforts in their clinics.  
518 Practical recommendations from this group have been grounded in the experience of  
519 constructing, using and sharing these systems. Brief summaries of several are highlighted to  
520 convey the scope and volume of these resources.

- 521 • The University of Michigan Radiation Oncology Analytics Resource (M-ROAR) automates  
522 aggregation of electronic data from the Treatment Planning System (TPS), Radiation  
523 Oncology Information System (ROIS), Electronic Health Record (EHR) and other  
524 databases for all patients treated. Data types include demographics, treatment and  
525 dosimetric data, chemotherapy, toxicities, comorbidities, labs, medications, encounters  
526 and patient reported outcomes (PROs). The system contains records for over 20,000  
527 patients. Key data elements are extracted utilizing a combination of SQL queries, TPS  
528 application programming interface (API) based scripts and custom code to extract and  
529 process data from multiple source systems [25].
- 530 • The UCLA Clinical Informatics Management System (CIMS) consists of three major  
531 modules: a physician interaction module that interacts closely with EHR, a physics  
532 parameter module that handshakes with PACS systems, treatment planning and delivery  
533 stations for quantitative value collection and exchange, and a patient reported outcome  
534 management system (Patient Reported Outcomes Measurement Information System,  
535 PROMIS) with a web/mobile portal. The physician interaction module supports  
536 comprehensive query for collection and integration of radiotherapy relevant  
537 information from other departments. The patient reported outcome management

538 module consists of a front-end with site-specific patient-oriented Common Terminology  
539 Criteria for Adverse Events (CTCAE) questionnaires tailored to patients. As of now, the  
540 registry contains records for 1790 definitive prostate treatment, 209 post-operative  
541 prostate treatment, 1950 breast, 663 lung, 531 brain metastasis, 484 GYN, 424 glioma,  
542 409 meningioma, 209 rectum, 151 metastatic bone, 164 trigeminal, 111 pancreas, and  
543 over 3000 general cases [30].

544 • The Ohio State University Radiation Oncology Department's "Quality Database" has  
545 been designed to serve as a data aggregation platform to capture clinical, technical, and  
546 health outcome data on all patients who receive radiation treatments. All data are  
547 stored in a REDCap database. Smart texts have been implemented in EHR to enable  
548 automated capture and extraction of discrete data elements such as adverse events  
549 from provider notes. The dosimetry data for radiation therapy are extracted via TPS's  
550 API. Demographics, diagnosis, tumor biomarkers, surgery, systemic therapy, radiation  
551 therapy, and adverse events constitute the collected data and provide means for  
552 determining effectiveness of treatment modality. The Quality Database currently  
553 contains 3385 patients and is being populated prospectively with new patient data.

554 • Oncospace: Johns Hopkins University developed a comprehensive data collection and  
555 data repository system [31]. The system consists of a network of data collection  
556 systems (ROIS, clinic computer terminals, mobile devices, hospital EHR) that provides  
557 data that is transformed and loaded into a SQL database. Using a federated database  
558 approach (including University of Washington, University of Virginia, Odette Cancer  
559 Center-Sunnybrook), each institution has implemented compatible schemas so  
560 federation-wide queries will succeed. This approach has the advantages of  
561 "crowdsourcing" ideas and technology and allowing each institution to keep control of  
562 their data while still permitting individual flexibility.

563 • The Veterans Health Administration (VHA) developed a pilot Radiation Oncology  
564 Practice Assessment (ROPA) program to assess the quality of radiotherapy across the  
565 entire VHA network with 40 institutions participating [32]. Data types include quality

566 metrics targeted at workup, diagnosis, treatment planning, delivery and follow-up. The  
567 gathered quality metrics were developed by the VHA in partnership with ASTRO for  
568 locally advanced non-small cell lung cancer, limited stage small cell lung cancer, and  
569 intermediate and high-risk prostate cancer. Data extraction for the initial pilot project  
570 will be completed in 2018. At that time ROPA is anticipated to contain 45,000 scores for  
571 49 metrics aggregated from approximately 2,000 patients.

572  
573 Large data sets from sources outside of radiation oncology are now available for  
574 analysis. Waddle *et al.* recently published utilization data derived from insurance  
575 records from a commercial warehouse (Optum Labs) to examine treatment technologies  
576 used (proton, stereotactic body radiotherapy, IMRT, 3D, other) by diagnosis code used  
577 in billing records. The data base contains utilization data on a subset of 474,533  
578 radiation oncology patients from a larger database of over 100 million insured lives.  
579 However, connection of this data to clinical outcomes and other cofactors was pending  
580 at the time of that analysis [33].

581  
582  
583  
584  
585

## 586 **Recommendations for next steps needed to improve data availability.**

### 587 *Adopting national standards*

588 As discussed above, an important aspect of data exchange is employing a generally recognized  
589 view of the medical process. HL7 FHIR is an emerging standard and one that has the crucial  
590 elements of (a) flexibility, (b) state-of-the-art technologically, and (c) widespread support [34].  
591 As this standard is just not being formalized, this is an excellent time for the radiation oncology  
592 community to support efforts to develop radiation oncology-specific resources for this standard  
593 [35].



594 *Increasing multi-institutional collaborative efforts*

595 Real, effective standards emerge from being actively engaged in exchanging data with outside  
596 groups as part of more frequent collaborations. Professional and government grant support for  
597 research efforts that develop and proof these standards as by-products are important to their  
598 emergence.

599 Included in this effort is need to facilitate information exchanges that support re-treatment. As  
600 patients are able to survive longer with cancer, likelihood of visiting more than one center for  
601 subsequent treatments increases. Clinical process and data exchange standardizations needed  
602 to facilitate these exchanges also support collaborative efforts.

603

604 *Links to institutional registries*

605 Institutions which are members of the National Comprehensive Cancer Network (NCCN) are  
606 required to have access to a registry which carries out longitudinal follow-up on a few key data  
607 elements (e.g. survival, cause of death) for treated patients. EHR database records may be  
608 substantially different from registry database records. Providing electronic access registry  
609 databases provides opportunities to synchronize data sources in constructing big data analytics  
610 resource systems.

611

612 *Support for Public Data Sets*

613 The value of producing data sets that can be publicly shared (without compromising PHI) has  
614 been heralded by several authors. [36-38]. There is growing interest from funding agencies for  
615 publicly funded research to produce publically available datasets. Similarly, an increasing  
616 number of journals require publication of datasets accompanying findings. Recently Medical  
617 Physics has introduced a special publication category just for data sets. Principles for ensuring  
618 that data are findable, accessible, interoperable, and reusable (FAIR) for public access of data  
619 sets have been set out by Wilkinson *et al.* [39] and others [40].

620 The National Cancer Institute has recently begun to implement a Cancer Research Data  
621 Commons which meet the standards of FAIR. In their announcement, they echo a number of  
622 the themes that we have set forth in this article. This is clearly a propitious time for radiation  
623 oncology to join with others in the oncology fields to make these sorts of community-wide  
624 efforts more productive [41].

### 625 *Informatics Training*

626 Clinical staff bring great value to informatics efforts because of the depth of their clinical  
627 domain knowledge with respect to key data elements, their inter-relationships, clinical  
628 processes by which data is entered, end user expectations for meaning, etc. The set of clinical  
629 staff that take on expanding their informatics skills to include database, programming,  
630 statistical analysis and machine learning also improve ability to develop practical solutions  
631 bridging needs between the larger number of specialists entirely focused in either the clinical or  
632 informatics domains.

### 633 **Conclusions**

634 We have laid out an argument for why it is important for the radiation oncology community to  
635 improve the means by which we can collect, share and use the data that we encounter every  
636 day. However, for various reasons, much of this data remains inaccessible to us in a format  
637 that makes it easy for us to transform data to knowledge.

638 The technological challenges to implementing a community-wide system of data collection,  
639 sharing and usage are formidable but the tools have been or are currently being developed.  
640 More difficult is developing the collective will to make it happen. Such a change in our clinical  
641 behavior and workflow requires buy-in from everyone, including clinic staff, physicians, and  
642 vendors. It is our hope and expectation that this sea change has already started to occur as  
643 diffuse networks grow in size and analytic power. It is necessary to do so if we are to continue  
644 to be at the forefront of harnessing technological advances to improve the treatments that we  
645 provide our patients.

646

647 **Acknowledgements**

648 *Thanks to* acknowledge Zhong, Haoyu, MS for figure 1.

649 This work was the result of collaborative efforts from participants at a meeting supported by  
650 AAPM, Varian Medical Systems and Elekta

651

652 **Acronyms:**

653 AAPM: American Association of Physicists in Medicine

654 AJCC: American Joint Committee on Cancer

655 API: Application Programming Interface

656 ASTP: As Treated Plan Sums

657 ASTRO: American Society for Radiation Oncology

658 BDAR: Big Data Analytic Resource Systems

659 CARO: Canadian Association of Radiation Oncology

660 CDR: Clinical Data Repository

661 CER: Comparative Effectiveness Research

662 CTCAE: Common Terminology Criteria for Adverse Events

663 DB: Database

664 DICOM: Digital Imaging and Communications in Medicine

665 DVH: Dose Volume Histogram

666 ESTRO: European Society for Therapeutic Radiation Oncology

667 EHR: Electronic Health Record

668 FAIR: Findable, Accessible, Interoperable, and Reusable

669 FHIR: Fast Healthcare Interoperability Standards

- 670 FIGO: International Federation of Gynecology and Obstetrics
- 671 HIPAA: Health Insurance Portability and Accountability Act
- 672 HL7: Health Level 7
- 673 ICD-O: International Classification of Diseases for Oncology
- 674 ICD9: International Classification of Diseases, Ninth Revision
- 675 ICD10: International Classification of Diseases, Tenth Revision
- 676 JSON: JavaScript Object Notation
- 677 NCCN: National Comprehensive Cancer Network
- 678 NIH: National Institutes of Health
- 679 OIS: Oncology Information System
- 680 PACS: Picture Archive and Communication Systems
- 681 PHI: Protected Health Information
- 682 PQI: Patient Quality and Improvement
- 683 PRO: Patient Reported Outcome
- 684 PROMIS :Patient-Reported Outcomes Measurement Information System
- 685 REDCap: Research Electronic Data Capture
- 686 ROIS: Radiation Oncology Information System
- 687 RCT: Randomized Controlled Trial
- 688 ROILS: Radiation Oncology Incident Learning System
- 689 RTOG: Radiation Therapy Oncology Group
- 690 SQL: Structured Query Language
- 691 TPS: Treatment Planning System
- 692 XML: Extensible Markup Language

693 VHA: Veterans Health Administration

694 **Bibliography**

695 [1] Pan HY, Mazur LM, Martin NE, et al. Radiation Oncology Health Information Technology: Is it  
696 working for or against us? *Int J Radiat Oncol Biol Phys*, 2017; 98(2): 259-262.

697 [2] Kerns, Sarah L., Harry Ostrer, and Barry S. Rosenstein. "Radiogenomics: using genetics to  
698 identify cancer patients at risk for development of adverse effects following radiotherapy."  
699 *Cancer discovery*. 2014;4(2): 155-165.

700 [3] Stevens, Emily A., and Cristina P. Rodriguez. "Genomic medicine and targeted therapy for  
701 solid tumors." *Journal of surgical oncology*. 2015; 111(1): 38-42.

702 [4] Jaffee EM, Dang CV, Agus DB, et al. Future cancer research priorities in the USA: a *Lancet*  
703 *Oncology* Commission. 2017;18(11): e653-e706.

704 [5] Shirts BH, Salama JS, Aronson SJ, Chung WK, Gray SW, Hindorff LA, Jarvik GP, Plon SE, Stoffel  
705 EM, Tarczy-Hornoch PZ, Van Allen EM. CSER and eMERGE: current and potential state of the  
706 display of genetic information in the electronic health record. *Journal of the American Medical*  
707 *Informatics Association*. 2015 ;22(6):1231-42.

708 [6] Parry C, Kent EE, Mariotto AB, et al. Cancer survivors: a booming population. *Cancer*  
709 *Epidemiol Biomarkers Prev*. 2011;20:1996-2005

710 [7] Oeffinger KC, Argenbright KE, Levitt GA, McCabe MS, Anderson PR, Berry E, Maher J, Merrill  
711 J, Wollins DS. Models of cancer survivorship health care: moving forward. *Am Soc Clin Oncol*  
712 *Educ Book*. 2014; 4 :205-13.

713 [8] Witt CM, Herman PM, Tunis S. Comparative effectiveness research in integrative oncology.  
714 2017. *JCNI*, 52: lgx013. <https://doi.org/10.1093/jncimonographs/lgx013>

715 [9] Fiore LD, Lavori PW. Integrating randomized comparative effectiveness research with  
716 patient care. *N Engl J Med*. 2016; 74: 2152-2158.

- 717 [10] Kalet AM, Gennari JH, Ford EC, Phillips MH. Bayesian network models for error detection  
718 in radiotherapy plans. *Phys Med Biol.* 2015; 60: 2735-49.
- 719 [11] Bentzen SM, Constine LS, Deasy JO, Eisbruch A, Jackson A, Marks LB, Ten Haken RK, Yorke  
720 ED. Quantitative Analyses of Normal Tissue Effects in the Clinic (QUANTEC): an introduction to  
721 the scientific issues. *IJORBP.* 2010;76(3):S3-9.
- 722 [12] MayoCS, Matuszak MM, Schipper MJ, Jolly S, Hayman JA, Ten Haken RK: Big Data in  
723 Designing Clinical Trials: Opportunities and Challenges. *Frontiers in Oncology* 2017; 7: 187
- 724 [13] Luo Y, El Naqa I, McShan DL, et al. Unraveling biophysical interactions of radiation  
725 pneumonitis in non-small-cell lung cancer via Bayesian network analysis. *Radiother Oncol.* 2017  
726 ;123(1):85-92.
- 727 [14] El Naqa I, Kerns SL, Coates J, et al. Radiogenomics and radiotherapy response modeling.  
728 *Phys Med Biol.* 2017 ;62(16):R179-R206.
- 729 [15] Lambin P, van Stiphout RGPM, Starmans MHW, et al. Predicting outcomes in radiation  
730 oncology--multifactorial decision support systems. *Nat Rev clin Oncol*, 2013: 10;27-40.
- 731 [16] Mayo CS, Pisansky TM, Petersen IA, Elizabeth, et al. Establishment of practice standards in  
732 nomenclature and prescription to enable construction of software and databases for  
733 knowledge-based practice review, In *Practical Radiation Oncology*, Volume 2016; 6(4): e117-  
734 e126, ISSN 1879-8500, <https://doi.org/10.1016/j.prro.2015.11.001>.
- 735 [17] Matuszak M, Anderson C, Lee C et al. An Integrated Application for Radiation Therapy  
736 Treatment Plan Directives, Management, and Reporting (SU-G-TeP4-06). *Med.Phys.* 43(6):3686
- 737 [18] Evans SB, Fraass BA, Berner P, Collins KS, Nurushev T, O'Neill MJ, Zeng J, Marks LB.  
738 Standardizing dose prescriptions: An ASTRO white paper. *Pract Radiat Oncol.* 2016;6(6):e369-  
739 e381. doi: 10.1016/j.prro.2016.08.007. Epub 2016 Aug 24. PubMed PMID: 27693224
- 740 [19]<https://www.dicomstandard.org/>; accessed 1/26/18.
- 741 [20] <http://www.hl7.org>; accessed 1/26/18.
- 742 [21] <https://www.hl7.org/fhir/DSTU1/http.html>; accessed 1/26/18.

- 743 [22] <https://ncit.nci.nih.gov/ncitbrowser/>; accessed 1/26/18.
- 744 [23] Basch, E., Deal, A. M., Kris, M. G., Scher, et al. (2016). Symptom Monitoring With Patient-  
745 Reported Outcomes During Routine Cancer Treatment: A Randomized Controlled Trial. *Journal*  
746 *of Clinical Oncology*, 34(6), 557–565. <http://doi.org/10.1200/JCO.2015.63.0830>.
- 747 [24] [https://www.hhs.gov/ohrp/international/ethical-codes-and-research-](https://www.hhs.gov/ohrp/international/ethical-codes-and-research-standards/index.html)  
748 [standards/index.html](https://www.hhs.gov/ohrp/international/ethical-codes-and-research-standards/index.html); accessed 1/26/18.
- 749 [25] Mayo CS, Kessler ML, Eisbruch A, Weyburne G, Feng M, Hayman JA, et al. The big data  
750 effort in radiation oncology: data mining or data farming? *Adv Radiat Oncol*. 2016; 1(4):260–71.
- 751 [26] Bailis P, Hellerstein JM, Stonebraker M, *Readings in Database Systems* 5th edition,  
752 <http://www.redbook.io/pdf/redbook-5th-edition.pdf>
- 753 [27] Mayo CS, Moran JM, Bosch W, et al. AAPM TG-263 Standardizing Nomenclatures in  
754 *Radiation Oncology*, *IJORBP*.2018;100(4):1057-1066.
- 755 [28] <http://bioportal.bioontology.org/ontologies/ROO> , accessed 2/30/2018
- 756 [29] <https://bioportal.bioontology.org/ontologies/DLORO>, accessed 2/30/2018
- 757 [30] Wang et al, Implementation Of A Comprehensive Radiation Therapy Registry: Focus On  
758 Feasibility and reliability, *IJORBP*. 2012; 83(3): S664
- 759 [31] Robertson SP, Quon H, Kiess AP, Moore JA, Yang W, Cheng Z, Afonso S, Allen M, Richardson  
760 M, Choflet A, Sharabi A. A data-mining framework for large scale analysis of dose-outcome  
761 relationships in a database of irradiated head and neck cancer patients. *Medical*  
762 *Physics*.;42(7):4329-37.
- 763 [32] Development and implementation of quality measures for the survey based performance  
764 assessment of radiation therapy in the VA. D. Caruthers, S. Brame, J. R. Palta, et. al. *IJROBP*  
765 2017 99; E391-E392
- 766 [33] Waddle MR, Kaleem TA, Niazi S, White L, Naessens J, Rummans T, Aljabri D, Habboush JY,  
767 Miller RC. Cost of Acute and Follow up Care in patients with Pre-Existing Psychiatric Diagnoses  
768 Undergoing Radiation Therapy, *IJROBP* 2017 99(5):1231

- 769 [34] [http://argonautwiki.hl7.org/index.php?title=Main\\_Page](http://argonautwiki.hl7.org/index.php?title=Main_Page); accessed 1/26/18.
- 770 [35] Phillips M, Halasz L. Radiation Oncology needs to adopt a comprehensive standard for data  
771 transfer: The case for HL7 FHIR. *Int J Radiat Oncol Biol Phys* 2017; 99: 1073-1075.
- 772 [36] Skripcak T, Belka C, Bosch W, Baumann M, et al. Creating a data exchange strategy for  
773 radiotherapy research: towards federated databases and anonymised public datasets.  
774 *Radiother Oncol.* 2014;113(3):303-9.
- 775 [37] Nyholm T, Olsson C, Montelius A et al. A national approach for automated collection of  
776 standardized and population-based radiation therapy data in Sweden. *Radiother Oncol.* 2016;  
777 119(2):344-50
- 778 [38] Roelofs E, Dekker A, Lambin P, et al. International data-sharing for radiotherapy research:  
779 an open-source based infrastructure for multicentric clinical data mining. *Radiother Oncol.* 2014  
780 ;110(2):370-4
- 781 [39] Wilkinson, M.D., Dumontier, M., Aalbersberg, I.J., et al. The FAIR guiding principles for  
782 scientific data management and stewardship. *Nature Scientific Data* 2016;3:160018
- 783 [40] <https://ncip.nci.nih.gov/blog/face-new-tragedy-commons-remedy-better-metadata/>;  
784 accessed 1/26/2018
- 785 [41] <https://ncip.nci.nih.gov/blog/towards-cancer-research-data-commons/>
- 786 [42] Gruber TR. A Translation Approach to Portable Ontology Specifications. *Knowledge*  
787 *Acquisition*, 5(2):199-220, 1993.
- 788
- 789

## 790 **Appendix**

### 791 **Key Data Elements and Relationships: A Radiation Oncology Translational Research Ontology**

792



793 We have defined of a common set of key data elements and relationships important to a broad range of  
794 patient quality improvement and translational research efforts. Ranking treatment information for  
795 effectiveness requires a broad scope of information types: Radiation Treatments, Surgery, Outcomes,  
796 etc. While it is desirable to have all the data readily available, that is not a practical starting point. Our  
797 objective here is to define a minimal set of information needed to handle frequently encountered  
798 questions as a common use starting point. With that, technical and procedural efforts attempting to  
799 automate electronic aggregation supporting Big Data efforts can use these recommendations as a guide.

800 Optimally professional organizations (e.g. AAPM, ASTRO, ESTRO, CARO) would establish an official listing  
801 of key data elements and relationships. Our intention here is to provide a practical starting point from  
802 our experience in aggregations from multiple source systems.

803 The listing of key data elements and relationships define an explicit conceptualization of a body of  
804 formally represented knowledge about Radiation Oncology, i.e. an ontology [42] The listing provided  
805 here was based on the ontology developed for M-ROAR [25] and expanded as an outgrowth of  
806 discussions at the Practical Big Data Workshop. Incorporation of the ontology into a programmatic form  
807 using Ontology Web Language (OWL) is underway.

808 Classes ( $\oplus$ ) of information, list key data elements (aka properties) denoted by one of three symbols ( $\bullet$ ,  
809  $\odot$ ,  $\circ$ ). Most elements ( $\bullet$ ) do not require special consideration for protection of patient health  
810 information (PHI). Elements that contain PHI ( $\odot$ ), are problematic for data sharing or storage in cloud  
811 based systems. Alternatives ( $\circ$ ), containing, reduced information, may be sufficient for a wide range of  
812 collaborative efforts or cloud based storage.

813 For example, dates are a type of patient health information (PHI) that institutional review boards (IRB)  
814 will not allow for many applications. For a wide range of investigations, detailing temporal relationships  
815 between events is important. Recording the patient's age at the event, rather than the date for the  
816 event is an alternative. For example, if the date of an event is 3/2/2013, and the patient's date of birth is  
817 8/17/1967, then the patient's age at the time of the event, to three decimal places (Decimal F3), is  
818 45.541. This is sufficient resolution to differentiate day on a timeline and meets requirements for  
819 protecting PHI.

820 Several key data elements typically are not present as distinct values in source data systems but have to  
821 be programmatically derived ( $\otimes$ ) from other elements. For example, the age of the patient at the time

822 of an event is derived from date of birth and date of the event. Starred (\*) items indicate particular need  
823 for recommendations of standardized values recommendations from professional societies.

824 When elements have only one instance they are indicated by the name of the class or element (e.g.  
825 DateOfBirth, Patient). When there may be more than one instance of an element, this is indicated by  
826 specifying a list of elements of this class (e.g. List<Course>).

827 Relationships among classes are categorized as Parent( $\Leftarrow$ ), Child( $\Rightarrow$ ), Sibling ( $\Leftrightarrow$ ) or Property( $\blacksquare$ ). Parent-  
828 Child are dependent relationships: a parent class object is referenced in each instance of a child class  
829 object. Sibling relationships are tracked if elements exist but do not imply dependence. Sibling  
830 relationships rather than parent-child relationships may be selected when the current state of the data  
831 will not practically support the dependent relationship. For example, Prescriptions are used in sibling  
832 relationships with respect to TreatedPlans because the current state of electronic data is inadequate to  
833 assure consistent mapping. Property relationships are used when class incorporates a set of elements  
834 grouped under a single concept.

### 835 $\oplus$ Patient -

- 836  $\odot$  PatientMRN (String) -:Medical Record Number
- 837  $\odot$  PatientGUID (String): Generalized Universal Identifier that can be used *in cloud*  
838 *based storage, when PatientMR is not.*
- 839  $\odot$  DateOfBirth (Date)
- 840  $\odot$  YearOfBirth (Int?)  $\text{\textcircled{X}}$
- 841  $\odot$  DateLastSurvivalCheck (Date?)
- 842  $\odot$  AgeAtLastSurvivalCheck (Decimal F3)  $\text{\textcircled{X}}$
- 843  $\odot$  DateOfDeath (Date?)
- 844  $\odot$  AgeAtDateOfDeath (Decimal F3)  $\text{\textcircled{X}}$
- 845  $\bullet$  IsAlive (Bool) – Status at last at Last Survival Check Date
- 846  $\bullet$  \*CauseOfDeath (String) – Need for standardized list
- 847  $\bullet$  Gender (String)
- 848  $\bullet$  Race (String)
- 849  $\bullet$  Ethnicity (String)

850

851 *Child class relationships*

- 852      ⇨ List<Radiation Therapy Course>
- 853      ⇨ List<Prescription>
- 854      ⇨ List<DiagnosisAndStaging>
- 855      ⇨ List<TreatedPlan>
- 856      ⇨ List<PatientTreatmentOutcome>
- 857      ⇨ List<PatientReportedOutcome>
- 858      ⇨ List<PlanningStructureSet>
- 859      ⇨ List<HealthInformation>
- 860      ⇨ List<Lab>
- 861      ⇨ List<Medication>
- 862      ⇨ List<Image>
- 863      ⇨ List<Chemotherapy Course>
- 864      ⇨ List<Surgical Procedure>
- 865      ⇨ List <Pathology>
- 866      ⇨ List <Charge>

867

868

869    ⊕ **RadiationTherapyCourse** ☹ – These are the treatment courses. A course Every patient has a list of  
 870    courses

- 871      ● CourseName (String)
- 872      ● NTxSessionsInCourse (Int) ☹ – Each treatment episode is a session, sessions used for  
 873      imaging only are exclude from the count
- 874      ⊙ DateFirstTreatment (Date)
- 875      ○ AgeAtFirstTreatment (Decimal F3) ☹
- 876      ⊙ DateLastTreatment (Date)
- 877      ○ AgeAtLastTreatment (Decimal F3) ☹

878

879      Sibling Class Relationships

880      ⇔ List<Prescription>

881      ⇔ List<Chemotherapy Course>

882            ⇔ List<Surgical Procedure>

883

884            Child class relationships

885            ⇨ List<TreatedPlan>

886            ⇨ List<DiagnosisAndStaging> - Typically only one per Course

887            ⇨ List<PatientTreatmentOutcome> - Typically only one per Course

888            ⇨ List<Charge>

889

890            Parent Class Relationships

891            ⇨ Patient

892

893    ⊕ **Prescription** : The prescription needs to fully convey the intent of the physician for the treatment  
894    plan. The Course contains a list of prescriptions

895            ● Name (String)

896            ● NTxSessions (Int)

897            ● NTxPerDay (Int)

898            ● DaysBetweenTxSessions (Decimal) ☹

899            ● StartOnNthDayFromCourseStart (Int) ☹

900            ● StartOnNthSessionInCourse (Int) ☹

901            ● RxDoseUnits (String) – “cGy” or “Gy” or “CGE”

902            ● IsCourseCumulativePrescription (Bool) ☹ – Only one value of True per Course

903

904

905            Class Property Relationships

906            ■ List<PrescriptionDoseLevel>

907            ■ List<PrescriptionDVHObjectives>

908

909            Sibling Class Relationships

910 ⇔ List<TreatedPlan>

911

912 Parent Class Relationships

913 ⇐ DiagnosisAndStaging

914 ⇐ Patient

915 ⇐ Course

916

917

918 ⊕ **PrescriptionDoseLevel**

919 ● RxDose (Decimal F3)

920 ● RxStructure (String) – AAPM TG263 compliant name

921 ● RxPointName (String)

922

923 Parent Class Relationships

924 ⇐ Prescription

925

926 ⊕ **PrescriptionDVHObjectives**

927 ● Structure (String) – AAPM TG263 compliant name

928 ● DVHMetric (String) – AAPM TG263 compliant name e.g. Max[Gy], V20Gy[%]

929 ● Constraint (String) - allowed values are =,<,<=,>,>=, ALARA

930 ● Value (Decimal F3) – null if constraint is ALARA

931

932 Class Property Relationships

933 ■ Prescription

934

935 ⊕ **DiagnosisAndStaging**

936 ● StagingSystem (String) - e.g. AJCC 7, FIGO

937 ● ICD9Or10 (String)

- 938 ● ICD0 (String) – Defines location of disease
- 939 ● Laterality (String) – Left, Right, Bilateral
- 940 ● Overall Staging (String): e.g. IIa, X,
- 941 ● T (String)
- 942 ● N (String)
- 943 ● M (String)
- 944 ● P (String)
- 945 ● G (String)
- 946 ● OtherStagingComponents (String)-Staging components other than T,N,M,P,G
- 947 ● PrimaryOrMetastatic (String) ☞– Either “Primary” or “Metastatic”
- 948
- 949 Child Class Relationships
- 950 ⇨ PatientTreatmentOutcome
- 951 ⇨ DiseaseSiteStatus
- 952
- 953 Parent Class relationships
- 954 ⇨ PrimaryICD9Or10? – If Metastatic, indicate Primary DiagnosisAndStaging element
- 955 ⇨ Course
- 956 ⇨ Patient
- 957
- 958
- 959 ⊕ **DiseaseSiteStatus**
- 960 ● DateOfStatus (Date)
- 961 ● AgeAtDateOfStatus (Decimal F3) ☞
- 962 ● \*Status(String) – Need standardized list e.g. (No Evidence of Disease, Local Recurrence,
- 963 Regional Recurrence, Distant Recurrence)
- 964
- 965
- 966

967 ⊕ **TreatedPlan** : Every course has a list of treated plan objects. One table for all types of plans defining  
968 key elements to track. This simplifies mixed modality tracking e.g. External + Brachy and handling of  
969 individual plans vs plan sums. Only plans actually treated are tracked. Details of actual vs number of  
970 fractions delivered are tracked.

- 971 ● PlanName (String): Corresponds to PlanID in ARIA
- 972 ● \*TreatmentAreaClassifier (String) : e.g. Head and Neck, Lung\_L, Breast\_R+SC
- 973 ● TPSSourceSystem (String) ☒
- 974 ● IsCourseCumulativePlan (Bool): The plan or plan sum(ATPS) represents all plans treated in  
975 the course
- 976 ● IsPlanSum (Bool): The dose associated with the plan is created by summing dose from other  
977 plans

978

979 ⊙ DateOfFirstPlanTreatment (DateTime)

- 980 ○ AgeAtFirstPlanTreatment ☒

981 ⊙ DateOfLastPlanTreatment (DateTime)

- 982 ○ AgeAtLastPlanTreatment ☒

983

984 ● PrimaryTxDeliveryFacility (String) – Facility where most of plan fractions were delivered

985 ● PrimaryTxDeliveryMachine (String) – Machine on which most of the plan fractions were  
986 delivered

987 ● NFractions\_Planned (Int)

988 ● NFractions\_Delivered (Int)

989 ● TotalDose\_Planned (Decimal) – Dose planned for highest dose structure e.g. PTV\_High

990 ● TotalDose\_Delivered (Decimal) – Dose delivered for highest dose structure e.g. PTV\_High

991 ● TotalDose\_Units (String) – Gy, cGy, CGE

992

993 ● UsedFiducials (Bool) ☒

994 ● FiducialType (String) – Gold, Calypso, Carbon

995 ● UsedBreathMotionControl (Bool) ☒

996 ● BreathMotionControlType (String): SDX, ABC, Compression

997

- 998 ● MeanSessionTimeMinutes(Int) ☒
- 999 ● MeanSessionBeamOnTimeMinutes (Int) ☒
- 1000 ● MeanSessionImagingTimeMinutes (Int) ☒
- 1001
- 1002 ● NImages\_MV (Int) ☒ - Total number of MV images for all sessions treating this plan
- 1003 ● NImages\_kV (Int) ☒ - Total number of kV images for all sessions treating this plan
- 1004 ● NImages\_CBCT (Int) ☒ :Total number of CBCT for all sessions treating this plan
- 1005 ● NImages\_MR (Int) ☒: Total number of MR images for all sessions treating this plan
- 1006
- 1007 ● List<SupplementalTreatmentDetail>
- 1008
- 1009 Sibling Class Relationships
- 1010 ⇔ Prescription
- 1011 ⇔ List<Images> - Image Class Objects related to the TreatedPlan e.g. CBCT, kV
- 1012
- 1013 Child Class Relationships
- 1014 ⇨ PlanningStructureSet
- 1015 ⇨ List<DVHCurve>
- 1016 ⇨ List<DVHMetric>
- 1017 ⇨ List<PatientPositioningDevice>
- 1018 ⇨ TreatmentPlanDetails\_XRT
- 1019 ⇨ TreatmentPlanDetails\_Brachy
- 1020 ⇨ TreatmentPlanDetails\_Particles
- 1021 ⇨ PlanningStructureSet
- 1022
- 1023 Parent Class Relationships
- 1024 ⇨ Patient
- 1025 ⇨ Course
- 1026 ⇨ ComponentOfATPS (TreatedPlan) - Plans that are components of ATPS link back to the ATPS



1027

1028

1029 ⊕ **PlanningStructureSet**

1030 ● StructureSetName (String)

1031 ● ImageModality (String)⌘ : e.g. CT, MR

1032 ⊙ DateOfImageAcquisition (Date)

1033 ○ AgeAtImageAcquisition (Decimal F3)⌘

1034 ⊙ DICOMImage\_UID (String) DICOM\_UID of image use for the plan. In the Image list attached to the patient.

1035 ⊙ DICOMPlan\_UID (String)

1036 ⊙ DICOMStructure\_UID (String)

1037 ⊙ DICOMDose\_UID (String)

1038 ● PatientPosition (String)

1040

1041 Parent Class Relationships

1042 ↩ Patient

1043 ↩ TreatedPlan

1044

1045 ⊕ **PatientPositioningDevice**

1046 ● \*DeviceCategory (String) – Need standardized list

1047 ● DeviceName(String)

1048 ● SetupDetails (String)

1049

1050

1051 ⊕ **TreatmentPlanDetails\_XRT**

1052 ● List<EnergyModality>

1053 ● TotalPlanMU (Decimal)

1054 ● UsedIMRT (Bool) ⌘

1055 ● UsedVMAT (Bool) ⌘

1056	● UsedFIF (Bool) ☒
1057	● UsedWedges (Bool) ☒
1058	● UsedBolus (Bool) ☒
1059	● UsedNonCoplanarBeams (Bool) ☒
1060	● NBeams (Int) ☒
1061	● NFractionsPlanned (Int)
1062	● NFractionsDelivered (Int)
1063	● List<SupplementalTreatmentDetail>
1064	Parent Class Relationship
1065	↳ TreatedPlan
1066	
1067	
1068	⊕ <b>TreatmentPlanDetails_Brachy</b>
1069	● List<EnergyModality>
1070	● NSourcesTotal (Int)
1071	● TotalActivity (Decimal)
1072	● *TotalActivityUnits (String)- Need standardized list e.g. MBq, Ci, mCi, GBq
1073	● UsedRadiopharm (Bool)
1074	● UsedApplicator (Bool)
1075	● TotalHDRDwellTimeMin (Decimal)
1076	● TotalPDRDwellTimeMin (Decimal)
1077	● TotalLDRImplantTimeMin (Decimal)
1078	● List<SupplementalTreatmentDetail>
1079	
1080	Child Class Relationships
1081	↳ List<Applicator>
1082	
1083	Parent Class Relationship
1084	↳ TreatedPlan

1085

1086

1087 ⊕ **Applicator**

1088 ● \*ApplicatorType (String) Need standardized list e.g. Needle, BrachyCath, TandemAndOvoid,  
1089 Cylinder, Mamosite, Savi

1090 ● NApplicatorsInserted (Int) ⌘

1091 ● NApplicatorsUsedInTx (Int) ⌘

1092

1093 Parent Class Relationships

1094 ⇐ TreatmentPlanDetails\_Brachy

1095 ⊕ **TreatmentPlanDetails\_Particles**

1096 ● List<EnergyModality>

1097 ● UsedPassiveScattering (Bool)

1098 ● UsedSpotScanning (Bool)

1099 ● UsedEndOfRangeToSpareCriticalOAR (Bool)

1100 ● List<SupplementalTreatmentDetail?>

1101 Parent Class Relationships

1102 ⇐ TreatedPlan

1103

1104 ⊕ **EnergyModality**

1105 ● Energy (String) – Need standardized list e.g. X06, X06FFF, X10, X10FFF, E06, E09, E12, E16,  
1106 E20, Ir192, I125, P70, C250

1107 ● \*Modality (String) – Need standardized list e.g. XRT, HDR, LDR, Proton, CyberKnife,  
1108 GammaKnife

1109

1110 Parent Class Relationship

1111 ⇐ TreatedPlanDetails\_XRT

1112 ⇐ TreatedPlanDetails\_Brachy

1113      ⇐ TreatedPlanDetails\_Particles

1114

1115

1116    ⊕ **SupplementalTreatmentDetail**

1117      ● Name (String)

1118      ● Value (String)

1119      ● ValueType (String)

1120

1121      Parent Class Relationships

1122      ⇐ TreatedPlanDetails\_XRT

1123      ⇐ TreatedPlanDetails\_Brachy

1124      ⇐ TreatedPlanDetails\_Particles

1125      ⇐ TreatedPlan

1126

1127

1128    ⊕ **Image** : Information about image objects relevant to patient's treatment

1129      ● ImageName (String)

1130      ● DICOM\_UID (String)

1131      ● ImageModality (String) e.g. CT, kV, CBCT, MR-T1w, MR-T2w,PET,etc

1132      ● SourceSystem (String) ⌘ Where to find the image and how to get it e.g. ARIA, Velocity,

1133      Hospital PACS, etc

1134      ⊙ AccessionNumber (String)

1135      ● StudySeries (String)

1136      ● BodySite (String)

1137      ⊙ DateOfImageAcquisition (Date)

1138      ○ AgeAtImageAcquisition (Decimal F3) ⌘

1139      ● RelevanceComment (String?) e.g. TumorResponse

1140

1141      Sibling Class Relationships

1142 ⇔ List<ImageDataFeature>

1143 ⇔ TreatedPlan

1144 ⇔ Course

1145

1146 Parent Class Relationships

1147 ⇨ Patient

1148

1149 ⊕ **DVHCurve** : Store the DVH curve for as treated (i.e. number of fractions delivered) plans and plan  
1150 sums. Every Treated Plan has a list of DVH curves

1151 ● StructureName (String) – Use TG263 Standardization

1152 ● Volume[cc] (Decimal)

1153 ● Min[Gy] (Decimal)

1154 ● Max[Gy] (Decimal)

1155 ● Mean[Gy] (Decimal)

1156 ● Median[Gy] (Decimal)

1157 ● Stdev[Gy] (Decimal)

1158 ● DVHCurve (String) ☿ – Dose, Volume tuples separated by semi colons. Dose is in units of Gy,  
1159 Volume is in units of percent of structure volume e.g. 0,100; 50,100;50.5,99.5;....

1160

1161 Sibling Class Relationships

1162 ⇔ List<DVHMetric>

1163

1164 Parent Class Relationships

1165 ⇔ TreatedPlan

1166

1167 ⊕ **DVHMetric** : Metrics provide quick look up of most important values. Sibling relationship to DVH  
1168 curves is maintained so that they can be reported separately if needed.

1169 ● StructureName (String) - Use standard nomenclature from TG263

1170 ● MetricName (String) - Use standard nomenclature from TG263

- 1171 ● Value
- 1172
- 1173 Sibling Class Relationships
- 1174 ⇔ List<DVHCurve>
- 1175
- 1176 Parent Class Relationships
- 1177 ⇐ TreatedPlan
- 1178
- 1179 ⊕ **ImageDataFeature** : specific values associated with the image that e.g Radiomics values.
- 1180 Every Image has a list of image data features
- 1181 ● \*FeatureName(String) – Need for a standardized list of defined feature names and
- 1182 acceptable values
- 1183 ● Data Type (String): text, number, datetime, bool
- 1184 ● Value (String)⌘
- 1185 ○ DateOfImageDataFeature (Date)
- 1186 ○ AgeAtImageDataFeature (Decimal F3) ⌘
- 1187
- 1188
- 1189 Parent Class Relationships
- 1190 ⇐ Image
- 1191 ⇐ Patient
- 1192 ⊕ **PatientTreatmentOutcome**
- 1193 ● \*DiseaseStatus (String) – Need standardized list e.g. Local Recurrence, NED,
- 1194 BiochemicalFailure
- 1195 ○ DateOfStatus (Date)
- 1196 ○ AgeAtStatus (Decimal F3) ⌘
- 1197
- 1197 Class Property Relationship

1198 ■ DiagnosisAndStaging

1199

1200 Parent Class Relationships

1201 ⇐ Patient

1202 ⇐ Course

1203

1204 ⊕ **PatientReportedOutcome**

1205 ● \*SurveyInstrumentName (String) – Need for standardized list

1206 ● \*ElementName (String) – Need for standardized list

1207 ⊙ DateOfPRO (Date)

1208 ○ AgeAtPRO (Decimal F3) ⌘

1209 ● Value (String)

1210 ● ValueType (String) – e.g. Bool, Date, Number

1211 Sibling Class Relationship

1212 ⇐ Course

1213

1214 Parent Class Relationship

1215 ⇐ Patient

1216

1217

1218 ⊕ **ProviderReportedToxicity**

1219 ● \*ToxicityName – Use standard names from CTCAE or other standards

1220 ● ToxicityStandard (String) e.g. CTCAE

1221 ⊙ DateOfReportedToxicity (Date)

1222 ○ AgeAtReportedToxicity(Decimal F3) ⌘

1223 ● Value (String)

1224 ● ValueType (String) – e.g. Bool, Date, Number

1225 ● Attribution (String)

- 1226 Sibling Class Relationship
- 1227 ⇔ Course
- 1228
- 1229 Parent Class Relationship
- 1230 ⇐ Patient
- 1231
- 1232
- 1233 ⊕ **HealthInformation:** Used to record data elements relevant to patient status e.g. smoker, rock
- 1234 climber, diabetes, etc.
- 1235 ● \*HealthInformationItemName (String) –Need for standardized list e.g. HasDiabetes,
- 1236 IsCurrentSmoker, SmokingPackYears
- 1237 ⊙ Date (Date)
- 1238 ○ AgeDate (Decimal F3) ⌘
- 1239 ● Value (String) – e.g. True, 20
- 1240 ● ValueType (String) – Decimal, Bool, Date, String
- 1241
- 1242 Sibling Class Relationships
- 1243 ⇔ List<Course>
- 1244
- 1245 Parent Class Relationships
- 1246 ⇐ Patient
- 1247
- 1248 ⊕ **Lab**
- 1249 ● LabName (String)
- 1250 ● LOINCShortName (String)
- 1251 ● LOINCCodeName (String)
- 1252 ⊙ Date (Date)
- 1253 ○ AgeAtDate (Decimal F3) ⌘



- 1254 ● Value (String)
- 1255 ● Units (String)
- 1256 ● ValueType (String) – Decimal, Bool, Date, String

1257

1258 Sibling Class Relationships

1259 ⇔ Course

1260

1261 Parent Class Relationships

1262 ⇐ Patient

1263

1264 ⊕ **Medication**

- 1265 ● MedicationType (String)
- 1266 ● MedicationName (String)
- 1267 ● DosageValue (Decimal)
- 1268 ● DosageUnit (String)
- 1269 ● Frequency (String)
- 1270 ⊙ DateOfMedicationRecord
- 1271 ○ AgeAtMedicationRecord (Decimal F3) ⌘

1272

1273 Sibling Class Relationships

1274 ⇔ Course

1275

1276 Parent Class Relationships

1277 ⇐ Patient

1278

1279 ⊕ **ChemotherapyCourse**: Set of Chemotherapy administrations

- 1280 ● \*Protocol (String) – Need standardized list
- 1281 ● Agent (String)

- 1282 ● Facility (String)
- 1283 ● IsNeoAdjuvant (Bool)
- 1284 ● IsConcurrent (Bool)
- 1285 ● IsAdjuvant (Bool)
- 1286 ○ DateFirstTreatment (Date)
- 1287 ○ AgeAtFirstTreatment (Decimal F3) ⌘
- 1288 ○ DateLastTreatment (Date)
- 1289 ○ AgeAtLastTreatment (Decimal F3) ⌘
- 1290
- 1291 Sibling Class Relationships
- 1292 ⇔ Radiation Therapy Course
- 1293 ⇔ Surgical Procedure
- 1294
- 1295 Child Class Relationships
- 1296 ⇨ List<Chemotherapy Administration>
- 1297
- 1298 Parent Class Relationships
- 1299 ⇐ Patient
- 1300 ⇐ DiagnosisAndStaging
- 1301
- 1302 ⊕ **ChemotherapyAdministration**
- 1303 ● Agent (String)
- 1304 ● Dosage (String)
- 1305 ○ DateOfAdministration (Date)
- 1306 ○ AgeAtAdministration (Decimal F3) ⌘
- 1307
- 1308
- 1309 ⊕ **SurgicalProcedure**

- 1310 ● Facility (String)
- 1311 ● \*Purpose (String) – Need for standardized list
- 1312 ● \*Margins (String) – Need for standardized values
- 1313 ● \*BiopsyStatus (String) – Need for standardized values
- 1314 ● IsPrelrradiation (Bool)
- 1315 ⊙ DateOfSurgery (Date)
- 1316 ■ ○ AgeAtSurgery (Decimal F3) ⌘
- 1317
- 1318 Sibling Class Relationships
- 1319 ⇔ RadiationTherapyCourse
- 1320 ⇔ ChemoTherapyCourse
- 1321
- 1322 Parent Class Relationships
- 1323 ⇐ Patient
- 1324 ⇐ DiagnosisAndStaging
- 1325
- 1326 ⊕ **Pathology**
- 1327 ● \*ElementName(String) – Need standardized list
- 1328 ● \*ElementValue (String)
- 1329 ● \*ElementType (String)
- 1330 ⊙ DateOfPathology (Date)
- 1331 ○ AgeAtPathology (Decimal F3) ⌘
- 1332
- 1333 Sibling Class Relationships
- 1334 ⇔ DiagnosisAndStaging
- 1335
- 1336 Parent Class Relationships
- 1337 ⇐ Patient

1338

1339

1340 ⊕ **Charge**

1341 ● CPTCode (String)

1342 ● NCodeInstances(Int)

1343 ● DateStartRange (Date)

1344 ○ AgeAtStartRange (Decimal F3) ⌘

1345 ● DateEndRange (Date)

1346 ○ AgeAtEndRange (Decimal F3) ⌘

1347

1348 Parent Class Relationships

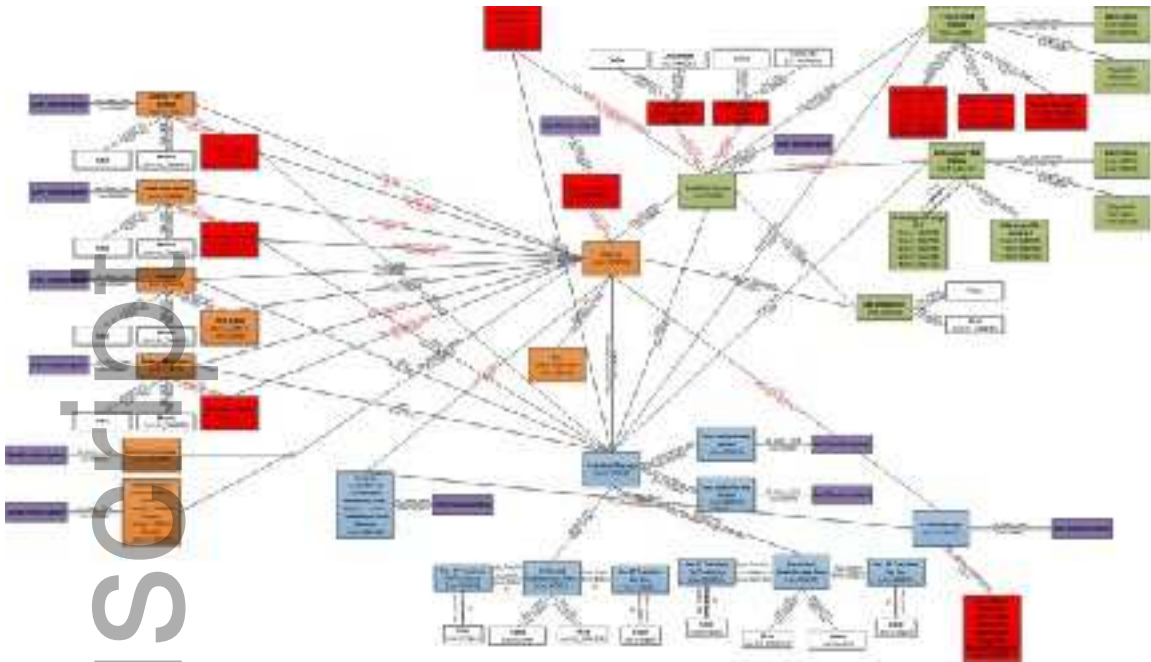
1349 ↩ Patient

1350 ↩ Course

1351

1352 **Figure Legend**

1353 Figure 1: The data from RTOG 0012, RTOG 0247, and RTOG 0822 were converted into Resource  
1354 Description Framework (RDF) specifications and were uploaded onto the NRG/IROC/ACR node  
1355 of the Varian learning portal. The mapping was performed according to the diagram shown  
1356 above. Distributed learning is enabled for contracted institutions. The distributed learning  
1357 between this node and another node on the Varian learning portal (MAASTRO Clinic,  
1358 Netherlands) was tested successfully.



mp\_13114\_f1.tif