1

3 CORRESPONDING AUTHOR EMAIL: marthamm@med.umich.edu

4

5 Performance/Outcomes Data and Physician Process Challenges for Practical Big

6 Data Efforts in Radiation Oncology

7

8 Matuszak MM[1], Fuller CD[2], Yock TI[3], Hess CB[3], McNutt T[4], Jolly S[1], Gabriel P[5], Mayo CS[1], Thor

9 M[6], Caissie A[7], Rao A[1], Owen D[1], Smith W[8], Palta J[9], Kapoor R[9], Hayman J[1], Waddle M[10],

10 Rosenstein B[11], Miller R[10], Choi S[2], Moreno A[2], Herman J[2], and Feng M[13]

11

12 1) University of Michigan, Ann Arbor, MI,  2)MD Anderson Cancer Center, Houston, TX,

13 3)Massachusetts General Hospital, Boston, MA, 4)Johns Hopkins University, Baltimore, MD , 5)

14 University of Pennsylvania, Philadelphia, PA, 6) Memorial Sloan Kettering Cancer Center, New

15 York, NY, 7)  Dalhousie University, Halifax, Nova Scotia, 8) University of Washington, Seattle,

16 WA, 9) Virginia Commonwealth University, Richmond, VA, 10) Mayo Clinic, Jacksonville, FL, 11)

17 Icahn School of Medicine at Mount Sinai, New York, NY, 12) Mayo Clinic, Phoenix, AZ,

18 13)University of California at San Francisco, San Francisco CA

19

20

21 **Abstract**

22 　　　It is an exciting time for big data efforts in radiation oncology.  The use of big data to

23 help aid both outcomes and decision making research is becoming a reality.  However, there

24 are true challenges that exist in the space of gathering and utilizing performance and outcomes

25 data.  Here, we summarize the current state of big data in radiation oncology with respect to

26 outcomes and discuss some of the efforts and challenges in radiation oncology big data.

27

## Introduction

The promise and potential of "big data" in radiation oncology cannot be overstated. There is tremendous excitement regarding the ability to learn about the efficacy of treatment, discover new interactions, and overall being able to offer our patients improved and tailored treatments based on the experience of many. There is also the hope of shared decision making between providers and patients using informed tradeoffs between cancer control and side effects. However, genuine challenges are to be faced before this can become a reality and to meet those challenges, one must first examine the nature of this "big data." There is a tendency to use the term "data mining" when thinking about informatics, when in fact, data farming is a more accurate term, reflecting the reality that the entire process, from planting the seeds of data in organized rows, watering and tending the growth of data, then harvesting it, is critical to understand and plan for (1).

Our ability to provide patients with answers about their best course of treatment relies on our a priori knowledge of how patients with similar disease, demographics, preference, and clinical characteristics were treated, and how they responded to treatment including both tumor control and treatment-induced toxicities. This data must be captured in a useable way so that it can be extracted and analyzed, with user-friendly predictive models created so that treatment can be customized for each patient.

In radiation oncology, there are two critical general issues, which must be addressed: 1.) Since radiation oncology data is different than medical/surgical oncology data, data platforms which have been designed with this in mind (many of which already exist) must be utilized. 2.) Existing standards where possible should be utilized to meet the big data needs of the multiple stakeholders (current and future patients, physicians, registries, insurance companies, the informatics community and many other groups) in radiation oncology in order to avoid duplication of work. We herein summarize the clinical aspects of big data collection in radiation oncology, and highlight the challenges and future work needed so that we can realize the potential of big data.

***Radiation Oncology Big Data is Unique***

57        An essential point that must be embraced for radiation oncology big data to reach its

58      potential is, as mentioned under 1.) above, that its format and nature is inherently different

59      from other disciplines. Fortunately, radiation oncology has recognized this leading to a number

60      of existing specialized data structures in its arsenal, including DICOM-RT structure and dose

61      files. Archiving treatment images, structures and doses in DICOM format is a relatively easy first

62      step toward ensuring that radiation oncology treatment data is captured. It also provides a

63      great step toward future quality assurance of that data. However, some features of treatment

64      are not captured in DICOM format, including, for example, motion management and use of

65      bolus (if not included in the simulation).  Recreating delivered dose requires the integration of

66      additional information (e.g. CBCT, log files from the treatment machine) in addition to the

67      treatment plan.

68        Standardizing nomenclature and definitions are crucial to our efforts to believe and

69      understand aggregated data (2). There is a recognized, but currently unmet need in radiation

70      oncology to standardize naming and delineation procedures of normal structures as well as

71      targets. Standardization includes not only naming structures, but consistency of anatomic

72      borders and instructions on the extent of normal organs to be contoured.  For example, naming

73      every esophagus "esophagus" rather than "eso" or "esoph" and contouring it from the cricoid

74      to the stomach is imperative if we hope to better understand dose-volume response-

75      relationships. If every "esophagus" in a big data set must go through independent quality

76      assurance, then the effort will not get very far. This is where planting the seeds correctly in the

77      first place pays off. Even with the best intentions, the complete OAR delineation can be

78      compromised by a treatment planning scan of limited extent, so standard nomenclature, as

79      suggested in TG263, of partial structures is recommended for clarity (2). Another often

80      overlooked element in radiation oncology big data is encoding of spatial information, especially

81      with recurrence. It is essential to know the spatial location of recurrence and its relationship to

82      the delivered dose, not just planned dose. Further, understanding why a marginal recurrence

83      occurred (e.g. variable patient positioning, inadequate GTV/IGTV delineation, poor image

84      registration, inadequate PTV margin) requires analysis of information from many steps of the

85 process. These are examples of data rarely available outside a research study, but essential to
86 determining tumor dose-response relationships.

87

88 **Use case examples**

89      Radiation oncology has a number of early adopters of the big data paradigm that can
90 help guide the field into best practices for successful capture of patient outcomes data.  One
91 well-known example is the euroCAT infrastructure (3).   Below are several other examples that
92 were presented or discussed as part of a breakout session at the 2017 Practical Big Data
93 Workshop.   In each example, a successful workflow has been implemented to capture
94 outcomes and performance data.  The benefits and limitations of each use case are given
95 below.  It should be noted that this is a list of examples and not an exhaustive list of all of the
96 excellent big data initiatives that are ongoing in the radiotherapy community.  Table 1 attempts
97 to summarized the use cases presented here for quick reference.

98

99 *M-ROAR – University of Michigan*

100      The University of Michigan has developed the Michigan Radiation Oncology Analytics
101 Resource (M-ROAR) to aid in practice patterns and outcomes analyses in Radiation Oncology.
102 This effort involved a multi-faceted strategy of requiring entry of critical elements as discrete
103 data, building a database platform, which pulls data from the oncology information systems
104 (OIS) and electronic health records (HER), and creating a self-service interface. On the data-
105 entry size, everyone in the clinic made a commitment to entering tumor staging, diagnosis
106 code, pain scores, patient reported outcomes, and Common Terminology Criteria for Adverse
107 Events (CTCAE) scores so that this data would be available for future analysis. Also, structure
108 nomenclature was standardized. The MS SQL database aggregates data for >17,000 patients
109 treated in the department since 2002, including information from both the radiation oncology
110 and hospital information systems. The self-service interface allows users to easily create and
111 optimize reports for cohort discovery in minutes rather than waiting to get to the top of a
112 report-writer's queue with each request or iteration.

113    With implementation of this strategy, the M-ROAR database can be used to answer
114 innumerable clinical questions, such as what factors predict patient risk of hospitalizations,
115 decline in patient function, and treatment-related complications, so that patient treatment
116 protocols can be adjusted in advance. As an example, for head and neck cancer, the association
117 between radiation dose and toxicity can be stratified based on HPV status. Information to
118 optimize clinical operations can also be gathered, such as: How long does a certain treatment
119 plan take to deliver vs. another one so that therapy time slots can be scheduled properly, and
120 What patients are at risk for dehydration so that nutrition consults can be requested or
121 outpatient hydration appointments scheduled in advance?  These are only a few examples of
122 practice-changing queries, which are currently possible. This database is primarily to inform and
123 guide quality improvement, with IRB approval needed when used for research.

124    Challenges remaining in M-ROAR are consistent and standardized assessment of
125 physician and patient-reported toxicities, as well as recurrence scoring.

126

127 ***MD Anderson***
128 A vision of optimizing electronic health record (HER) utilization is currently being investigated at
129 MD Anderson Cancer Center in a multiphase process. Initiated within the Radiation Oncology
130 department, a thorough evaluation of user performance and available toolsets within EPIC was
131 performed in order to determine suboptimal practices that were limiting efficiency within the
132 clinic workflow. A general consensus of a need for standardized documentation and consistent
133 nomenclature for the purposes of improving quality and safety measures, accurate staging and
134 billing, and decreasing duplication of data entry led to the development of over 40 specialty-
135 specific templates for note generation.  These templates "pull in" discrete data elements
136 entered into EPIC by a single person (such as a nurse, midlevel, or primary referral service) so
137 that the need for dictation/manual data entry by other providers generating notes is
138 minimized. The patient's existing medical conditions, cancer stage, performance status,
139 symptoms/ROS, laboratory values, and radiologic imaging information are all structured fields
140 which are now automatically populated into specific locations within each template.
141 Furthermore, these templates utilize the Smartlist function in EPIC, which are lists of

142  customizable text that can also be retrieved at a later date as structured data. Smartlists have

143  therefore been used to define specialty-specific treatment options, protocol descriptions, and

144  structured CTCAE grading systems. Another advantage of EPIC is the ability for patient-related

145  outcome (PRO) forms to be sent to the patient electronically.  When patients fill out these

146  forms, the results are then sent back and saved in EPIC as discrete data, which is then

147  incorporated into templates and allows for more rapid documentation.

148

149  Overall, these templates offer additional advantages including increased patient screening for

150  protocol enrollment and user-friendly, electronic functionality for various research endeavors.

151  By having the variables listed above as structured, extractable data, every aspect of clinical

152  research becomes optimized. Patients can be quickly assessed and evaluated for protocol

153  eligibility, and once the patient is undergoing treatment under protocol, the collection and

154  reporting of clinical response and toxicity become more automated. Protocol-specific templates

155  have been created in order to ensure that all required data collection per individual protocol is

156  recorded in a uniform manner. Since completing phases I and II of template creation and

157  implementation within the Radiation Oncology department, there have been ongoing efforts to

158  expand standardized EHR documentation methods within other departments, beginning with

159  GI Medical Oncology and GI Surgery. So far, these services are adapting the templates to

160  maintain a similar data entry structure while tailoring sections such as the impression and plan

161  to suit their documentation needs. Our ultimate goal is to have the entire institution adopt the

162  use of standardized templates and structured data entry to 1) improve the efficiency of

163  documentation for providers and decrease the risk of provider burn-out, 2) improve patient

164  coordination within a multidisciplinary clinic setting, and 3) create an institution-wide system of

165  patient data collection for research purposes and assessment of clinical outcomes.

166

167  *Pediatric Proton Registry Consortium*

168      The Pediatric Proton Consortium Registry (PPCR) was established in 2012 to expedite

169  proton outcomes research in children and to better define the role of proton radiotherapy in

170  the pediatric cancer population (4). Approximately 1800 pediatric patients have been enrolled

171 in the PPCR across 13 participating pediatric proton centers. The PPCR is a consented registry

172 built upon the NIH supported free web-based data collection/repository platform, REDCap and

173 is currently open to any U.S. proton center that would like to participate. The PPCR collects

174 information on demographics, diagnosis and staging, baseline health status, chemotherapy and

175 surgery, radiation details, diagnostic imaging, and follow-up (5). Radiation plans are centrally

176 archived in the universal DICOM-RT format.   Due to funding issues and required manual effort,

177 there is limited participation and variable data entry.  Thus, there is an urgent need to improve

178 efficiency of data collection through automation.

179 The major challenges within the PPCR also present opportunities. Given that there are a

180 limited number of OIS and EHR platforms, there exists an opportunity to leverage the data

181 already contained within these platforms if appropriate programming bridges can be

182 constructed. An upfront investment of time and resources from technical personnel is needed

183 and standard interface should be created with standard basic information mapped from stable

184 locations in each OIS to minimize the need for additional customization at multiple sites.

185 Another opportunity exists with the general EHR. Given the critical mass of EPIC users in

186 the PPCR, we may be able to leverage collaboration to streamline data input and extraction. A

187 start could be the sharing and use of electronic templates and automation of population of

188 certain (standardized) fields in the database. It is key that templates must be efficient and user-

189 friendly with minimal free text so that clinicians will use them routinely and must be convinced

190 in the overall mission or be given timesaving in another area to counter-balance the extra work

191 of discrete data input.

192 The final component of PPCR is aggregation of plan information, which is eventually

193 used to help make the link between radiation dose and treatment outcomes. To facilitate this, a

194 partnership has been put in place with MiM Software (MiM Software Inc, Cleveland, OH) to

195 allow web-based archival for each participating institution. The partnership has led to the

196 development of a faster anonymization procedure and a script for automated nomenclature

197 standardization using TG263 (2).

198 In summary, the PPCR is an established and successful registry that has met some

199 hurdles along the way. As it has grown out of its funding source, it requires that we look into

200    electronic efficiencies that will help PPCR and other Radiation Oncology-related Big Data

201    efforts. Sufficient funding is critical to success of data collection. Mild funding pressure can spur

202    technological advances that can improve efficiencies, but these also need an upfront

203    investment in order to achieve them. Given the relatively few electronic radiation charts and

204    the few EHRs, we are better poised than ever to start to realize the goal of automation in data

205    entry.

206

207    ***Oncospace***

208    The Oncospace program at Johns Hopkins began with the design of a relational

209    analytical database that housed the treatment planning data in a form for fast query. The

210    database schema includes the full 3D dose for multiple radiation therapy sessions as well as the

211    3D anatomy including relevant structures (5). The system also houses features of the dose such

212    as the dose-volume histograms (DVHs) and shape relationships in the overlap volume

213    histograms (OVHs) (6). In the earlier work, the database was used for the development of

214    shape-based automated treatment planning where one could rapidly query the OVHs to

215    determine all prior treatments with critical organ that were "harder" to plan and use it to

216    predict the best achievable dose metric from DVHs (7-10). This method is in use today for both

217    plan quality evaluation and automated planning.

218    For outcomes, the Oncospace philosophy was that prospective structured data

219    collection should be integrated with the clinical workflow. Since 2007, a website enabling tablet

220    devices to be used in the clinic for data capture is available (11). Critical to the adoption is the

221    ability to generate clinical notes from the collected structured data and additional patient-

222    related information queried from the OIS. Using the same technology, electronic patient-

223    reported outcomes have been successfully captured for more than 8 years. Currently, there are

224    >5000 patients (prostate, head and neck, thoracic, breast and pancreas) in the database with

225    full treatment planning data, patient reported outcomes, clinician assessments on-treatment

226    and in follow-up, disease response as well as diagnosis, and lab data interfaced from clinical

227    systems.  Data are currently included from Johns Hopkins, the University of Washington, the

228    University of Virginia, and the University of Toronto Sunnybrook.

229　　　　The rapid access to the treatment data enables data science models to be explored (12).
230　The Oncospace group is now building predictive models for specific clinical decisions using
231　classification and regression tree models for weight loss and xerostomia prediction in head and
232　neck cancer and surgical candidacy in pancreatic cancer. The challenge in clinical prediction is to
233　focus on the decision to be made and what information truly informs it. For weight loss, the
234　decision is around the appropriate symptom management for improved nutritional support
235　such as feeding tube placement. In other cases, modifications to the treatment plan may
236　reduce risks if it does not compromise on target coverage. Additionally, the impact of the
237　spatially distributed radiation dose beyond DVHs to better understand how the patterns of
238　dose may impact the treatment related toxicities could be explored (13). The continued data
239　growth will allow continuous learning to fulfill the concept of a learning health system in the
240　future (14).

241

242　***University of Pennsylvania***

243　　　　The Penn Medicine Oncology Research and Quality Improvement Datamart (ORQID)
244　aggregates data from multiple source information systems, including Penn's enterprise EHR,
245　ROIS, TPS, Cancer Registry, and Center for Personalized Diagnostics. ORQID focuses on
246　organizing cancer patients' demographics, vital status, disease stage and prognostic indicators,
247　genomic variants, details of systemic therapy and external-beam radiotherapy, and physician-
248　reported toxicities.

249　　　　Outcomes have been among the most challenging data elements to capture. Penn
250　implemented structured, site-specific templates for documenting physician-reported toxicities
251　within the EHR in 2011. The templates are based on the CTCAE grading system, and clinical
252　teams selected the toxicities of focus for each disease site. To maximize opportunities for data
253　capture by providers at all levels, only clinically symptomatic toxicities (e.g. pain) not requiring
254　diagnostic interpretation (e.g. radiation pneumonitis) were included. Nurses have embraced
255　the effort and capture rates have been as high as 95% for on-treatment visits, which they
256　routinely staff. Physician adoption has been more challenging, and for follow-up visits (which
257　have less nursing support) capture rates have been below 50% of visits. Nevertheless, Penn has

258 amassed over 2 million toxicity observations on over 28,000 unique patients in the datamart.

259 Efforts are currently underway to implement widespread patient-reported outcome collection

260 as routine standard of care to help augment and complement the physician-reported toxicities.

261       For other outcomes, progression is tracked via the institutional cancer registry, which

262 only documents the timing and nature of the first progression event after initial treatment.

263 Deaths are identified from the EHR, cancer registry, and social security death masterfile, but

264 remain a challenge, with many deaths not documented or without accurate dates.

265

266 ***US Veterans Health Administration (VHA) Radiation Oncology Practice Assessment***

267       The National Radiation Oncology Program (NROP) office of VHA, with an oversight of 40

268 radiation therapy treatment centers treating over 15,000 patients annually has launched a pilot

269 program initiative in which patient-specific radiotherapy data is collected for quality assurance

270 assessment and comparative analysis of many treatment modalities and other factors at their

271 centers (15). The NROP office collaborated with the American Society of Radiation Oncology

272 (ASTRO) disease site expert committees to define clinical measures. These clinical measures are

273 based on established clinical guidelines, patterns of care assessment done by the American

274 College of Radiology's Quality Research in Radiation Oncology program (16), and expert

275 consensus opinions. These measures have formed the basis for assessing the quality of

276 treatments and practice variations and identification of the care gaps in the VHA. Although

277 dosimetry data was automatically abstracted from treatment planning systems (TPS), clinical

278 data had to be manually abstracted from the electronic health records (EHR) for the pilot

279 project.

280       The NROP office has embarked on a project to automatically extract all data for ROPA

281 from heterogeneous data sources that include EHR, TPS and Treatment Management Systems

282 (TMS) for clinical practice assessment, outcomes, and prospective decision support analytics.

283 An integrated data curation, storage and analytics portal, titled as HINGE (Health Information

284 Gateway and Exchange), was built that can extract and aggregate data from TPS and TMS,

285 physician clinical notes and DICOM-RT files. HINGE integrates data from these disparate sources

286 coherently and standardizes it for quality assessment and predictive analytics. The HINGE

287  database is based on well-defined quality measures defined by radiation oncology disease site

288  experts. HINGE has (i) tools to aggregate data from physician note templates (ii) a built-in

289  DICOM-RT parser to extract DVH based dose constraints, (iii) a natural language processing

290  (NLP) module to extract relevant physician assessments from the clinician notes, and (iii) a

291  decision-support and genomics module to provide supplementary insight to treatment

292  predictions, treatment outcomes and research hypotheses. The HINGE application would reside

293  at each VHA radiation oncology treatment site and transmit information to a centralized

294  database server thus making big data analytics possible. HINGE is capable of seamlessly

295  connecting to local IT/medical infrastructure via network and performs data extraction and

296  aggregation. The built-in modules (TMS extraction, DICOM parser, NLP) extract defined clinical

297  data and are easily extendable. The modules of decision-support and genomics provide

298  preliminary insights into a patient's treatment and health profile. Automatic data abstraction

299  with HINGE will enable real time assessment of clinical practices and determine care gaps.

300

301  ***Mayo Clinic Florida***

302  The Mayo Clinic Florida Department of Radiation Oncology has leveraged Mayo Clinic's

303  unique cost warehouse to aggregate data on the cost of radiation therapy and other associated

304  healthcare costs in the first two years after radiotherapy on approximately 3,000 patients over

305  a five year period incurred.  The Mayo cost data warehouse is a unique resources consisting of

306  linked EMR data and administrative data from Mayo Clinic's hospital and clinics in Florida,

307  Minnesota, and Wisconsin (17). These costs were linked to other sources of institutional data,

308  such as departmental treatment records captured through its radiation oncology information

309  system, demographic, tumor specific, and outcomes data obtained through Mayo's tumor

310  registry, adverse events recorded in the EMR, and other disease specific registries containing

311  non-oncological diagnosis data, such as psychiatric comorbidities.  Waddle *et al* have used this

312  cost warehouse to demonstrate that patients with co-existing psychiatric morbidities utilize the

313  emergency department and inpatient hospitalization at rates greater than patients without

314  psychiatric co-morbidities at 6 months and two years after radiotherapy. (18)  It should be

315  noted that even with many successes, toxicity capture remains challenging.

316

### *The Radiogenomics Consortium (RGC)*

The hypothesis that genetic/genomic alterations may function as surrogate biomarkers of disease response or normal tissue toxicity represents the basis of the field of radiogenomics (19). A principal goal of research in the field of radiogenomics is to identify the genomic markers associated with the development of adverse outcomes resulting from cancer radiotherapy. However, in order to accomplish this goal and definitively discover and validate the critical genomic markers, access to the radiotherapy treatment information and long-term longitudinal follow-up data reporting details as to adverse outcomes must be obtained for large numbers of patients. In order to enable the creation of large cohorts of patients who received radiotherapy, the Radiogenomics Consortium (RGC) was created in 2009, which is a cancer epidemiology consortium through the Epidemiology and Genomics Research Program of the NCI of the NIH (20). The RGC now has 225 investigators at 132 institutions in 31 countries. Although the RGC has successfully assembled large cohorts to perform adequately-powered studies, data harmonization remains a problem when multiple cohorts involve patients treated with a variety of radiotherapy techniques and evaluated using multiple grading systems. Nevertheless, a number of large studies have been accomplished in which substantial amounts of radiotherapy data have been gathered for studies that typically comprise over a thousand patients.

Four large studies involving the use of Big Data are currently in progress whose main goal is to discover new SNPs and validate previously identified genetic biomarkers predictive of susceptibility for the development of adverse effects resulting from radiotherapy. The first project involves roughly 6,000 men treated for prostate cancer, which encompasses multiple cohorts created by RGC investigators. DNA samples from all of these men have been genotyped and detailed clinical data are available with a minimum of two-years of follow-up.

The second large multi-center study developed by RGC members is REQUITE (Validation of predictive models and biomarkers of radiotherapy toxicity to reduce side-effects and improve quality-of-life in cancer survivors)(21). REQUITE addresses the challenge of data heterogeneity that, as for other big data projects, requires harmonization of the different

345 outcome measures and confounding variables used in multiple cohorts. This study does not
346 stipulate the radiotherapy protocols to be used but involves standardized case report forms
347 across centers and countries to ensure data in identical categories are collected. A key aspect of
348 REQUITE is the centralized database that includes pre-treatment DICOM and DVH files.

349     A third study involves three large cohorts comprising roughly 4,500 breast cancer
350 patients treated with radiotherapy for which blood samples and detailed clinical information
351 are available. These samples and data are available from three large groups of patients: (1)
352 1,500 patients treated under a series of breast cancer clinical protocols performed at New York
353 University School of Medicine (22-25); (2) ~2,000 breast cancer patients enrolled though the
354 REQUITE study and (3) ~1,000 women who receive breast cancer treatment through
355 participation in RTOG 1005 (26).

356     The fourth effort being made is to create a biorepository with linked clinical data for
357 patients treated with charged particle therapy (CPT). With the increasing use of CPT, there is a
358 need to establish cohorts for patients treated with these advanced technology forms of
359 radiotherapy. In recognition that the formation of patient cohorts treated with CPT for
360 radiogenomic studies is a high priority, efforts are underway to establish collaborations
361 involving institutions treating cancer patients with protons and/or carbon ions as well as
362 consortia, including the Proton Collaborative Group, the Particle Therapy Cooperative Group
363 and the Pediatric Proton Consortium Registry.

364

365 **State of the data**

366     As noted by the varied workflows highlighted in the use cases, hospital-wide and
367 radiation oncology-specific EHR systems are not often designed to facilitate collection of key
368 data elements for subsequent extraction and use. Typically, when a patient is referred to
369 radiation oncology, the diagnosis for that patient has been entered to the hospital EHR system.
370 Most radiation oncology-specific EHRs can link to the hospital EHR via HL7 FHIR (27) to sync the
371 diagnosis information. However, linking the specific diagnosis relevant to a given treatment
372 plan is often a manual process requiring physician input. In addition, there is generally not a
373 mechanism to input the staging information into the radiation oncology EHR or link metastatic

374  sites to the original diagnosis, which are in general of interest for outcome analyses. Thus,

375  curation of the diagnosis and staging information that comes into radiation oncology can be

376  cumbersome. Apart from simple diagnosis information, data elements from pathology,

377  radiology, surgery, internal medicine and medical oncology that may be relevant for radiation

378  oncology outcomes are seldom entered in discrete fields or even templated free-text formats,

379  and are, therefore, often inaccessible for automatic extraction and use.

380      As the patient goes through treatment, physicians typically see the patient weekly for

381  on treatment visits. However, the documentation of these visits, including routine toxicity

382  assessments relies on each individual institution creating their own clinical practice, datasheets

383  and custom tools for reporting. While many institutions are beginning to recognize the

384  importance of standardized toxicity assessments and PROs and are putting mechanisms in place

385  to track this data, there is still inconsistency, which can lead to missing data.  Further, once

386  institutions have these tools in place, it can be challenging to share personalized templates

387  across the varying platforms and clinical workflows that exist at different institutions.  Adding

388  this to the lack of standardized key data elements and time points to track for different

389  treatment sites, multi-institutional datasets are rarely comprehensive.

390      While some existing standards can be leveraged, it is important to evaluate if these

391  standards take into account the needs of all stakeholders and if not, determine if new

392  standards or perhaps simply minor amendments can be suggested to minimize the need to

393  start at the ground up. One must recognize that efforts to standardize common data elements

394  is a complex and time-consuming endeavor, but one that is ultimately worthwhile. An excellent

395  published discussion and proposed set of standard patient-reported outcomes within oncology

396  shows the complexity of these issues (28).

397      Once collected, Big Data will perform a crucial role by providing accurate outcome data

398  in order to build clinical decision support systems (CDSS) (29).  Conversely, decision models

399  themselves can be used to guide the selection of data elements to include.  In a recent work,

400  for example, a decision cost-model in the form of an influence diagram was constructed to

401  model the choice between photons and protons for the treatment of locally advanced non-

402  small cell lung cancer (30).  By including the monetary cost of managing acute toxicities, it was

403  possible to determine the ROC characteristics of a biomarker for radiosensitivity that a

404  physician would need in order to select patients for proton radiotherapy when their total

405  expected cost for protons is below that of photons.  As this cost-model example illustrates,

406  models can guide data farming efforts by establishing outcomes that are important for clinical

407  decision making, and by placing requirements on how accurately these outcomes need to be

408  known.  In this case, the required sensitivity and specificity were established for a novel test for

409  radiosensitivity for the decision to lower treatment costs.  This use of models may be especially

410  important when resources (e.g. cost of human labor) for populating databases are limited,

411  allowing efforts to be directed towards collecting the data that is most likely to lead to

412  improved clinical decision making.

413       This in turn highlights an important issue in constructing data standards for capturing

414  outcome data, namely, the standards need to be easily expandable.  As big data results are

415  applied in the clinic, used for clinical decision support, or new interactions are discovered

416  within the data, these efforts will inevitably – and rapidly – call for the collection of different

417  types of data.  Adaptability is emerging as a feature of data and communication standards

418  throughout healthcare, as recognition grows that developing a standard which attempts to

419  include everything will fail to do so, and in the process will become unwieldly.  HL7 FHIR, for

420  example, is a communication standard which follows an 80/20 directive, whereby 80% of the

421  elements which are implemented are included in the specification itself (31).  These core

422  elements are referred to as resources, and the remaining elements, called profiles, are

423  definable by individual institutions or groups in order to alter or add properties to resources.

424  Single institution databases can attempt to cover a greater proportion than 80%, although the

425  principle remains.  By embedding adaptability within a database initially intended to capture,

426  for example, only traditional treatment planning data, the database may later be populated

427  with patient reported outcomes, "omics" data, or patient preferences in the form of utilities,

428  rendering it useful in significantly more applications.

429

430  **Collection and Curation**

431    In order for the promise of big data to be realized in more than just individual radiation
432    oncology departments or networks of systems, standardized key data element lists and input
433    schemas are required. For example, the connection of diagnosis information to treatment
434    courses should be automated within vended systems and reviewed for quality on an ongoing
435    basis as part of a routine workflow, such as chart rounds. In addition, the relevant staging,
436    pathology, and histology information should be automatically extracted from the EHRs into
437    appropriate fields within the radiation oncology information system. Free-text searches or
438    simple natural language processing will be necessary for scanned outside hospital reports and
439    for other information not entered in discrete fields for easy extraction, particularly for
440    information not generated in radiation oncology and thus beyond our immediate control.

441    Collection of standardized key data elements related to toxicity, disease status, and
442    patient reported outcomes requires the definition of standards, as discussed above. However,
443    even with standard elements and data entry tools, there must be a culture shift in the radiation
444    oncology community to recognize the importance of comprehensive entry of the data as part of
445    the standard care for each patient. It is our responsibility to the field and future patients to
446    make collection of key data elements related to outcomes a priority.

447

448    **Access and Extraction**

449    Accessibility and extraction of the clinical data entered by the physician and patients, in
450    the case of patient-reported outcomes, is essential. The data storage infrastructure must
451    provide a mechanism for end users to extract the key data elements and aggregate the data
452    with other related data, such as dosimetric information. The system should be designed with
453    accessible application programming interfaces enabling user data extraction in the most
454    suitable and meaningful way. However, data extraction should not be performed on a project-
455    by-project basis. Rather, institutional information technology groups, especially those housed in
456    radiation oncology, should make it a priority and be proactive in supporting the construction of
457    big data analytics resource systems (BDARS).  This may require a partnership between radiation
458    oncology users and the IT managers so that domain knowledge can be shared and the BDARS
459    designed in such a way that the information is in a complete and usable format. The

460 development and use of a radiation oncology-specific ontology will be a key development in
461 ensuring that individual BDARS can be combined into true sets of big data.

462

463 **Specific Recommendations for Standardizations**

464 While there is clear work ahead in the community to reach a point where standard key
465 data elements are recorded routinely for all patients in radiation oncology, there are first steps
466 that can be taken. Summarized in Table 2 are example standard key data elements that could
467 be collected and thus should begin to be supported by vended systems.  Note that many such
468 elements would be collected at various timepoints including baseline, during treatment, end of
469 treatment, and at follow-up.  Therefore, properly capturing dates and being consistent with
470 relative dates is essential.

471

472

473

474

475

476

477

478 While Table 1 serves as a starting point for standardization of requested data elements,
479 collection of the data requires:

480

481 1. Creation of a standardized workflow that enables collection of proper data, at the right time
482 for the right patient.
483 2. Initiation of a working group to develop standards for classifying recurrence in radiation
484 oncology that includes spatial and dose information.

485

486 **Recommendations for Next Steps Needed to Improve Data Availability**

487 The current climate is such that "big data" is becoming a known term and fills one with
488 the promise of solving mysteries of care with a lot of data and computer. There is a focus on

489     data mining, as if the data is sitting waiting to be taken and analyzed. However, it is clear that

490     the data must be created and structured in a way to make it possible to harvest and answer

491     important and relevant clinical questions. As more providers buy into the need to standardize

492     for the sake of quality and process improvement, they will become more committed to

493     inputting essential common data elements related to outcomes. Vendors must also allow the

494     data to be accessed in a variety of ways, maintaining HIPAA compliance but no longer being a

495     major barrier to quality assurance. Improved automation in both capturing and accessing data

496     within vended systems is recommended to improve efficiency and accuracy in capturing

497     outcomes data. Engagement with all stakeholders, including physicians, legislators, patients and

498     patient advocates is essential to design modern approaches to handling protected health

499     information and drafting policies and legislation regarding how health care data can be used in

500     a safe way so as to maximize healthcare value and efficiency while maintaining security.

501

502

503

504 **References**

505

506    1. Mayo CS, Kessler ML, Eisbruch A, Weyburne G, Feng M, Hayman JA, Jolly S, El Naqa I,

507       Moran JM, Matuszak MM, Anderson CJ, Holevinski LP, McShan DL, Merkel SM, Machnak

508       SL, Lawrence TS, Ten Haken RK. The big data effort in radiation oncology: Data mining or

509       data farming? Adv Radiat Oncol. 2016 Oct 13;1(4):260-271.

510    2. Mayo, C. et al.  AAPM Task Group 263: Tackling Standardization of Nomenclature for

511       Radiation Therapy.  International Journal of Radiation Oncology • Biology • Physics ,

512       Volume 93 , Issue 3 , E383 - E384

513    3. Deist, T. M., et al. "Infrastructure and distributed learning methodology for privacy-

514       preserving multi-centric rapid learning health care: euroCAT." Clin Transl Radiat Oncol 4:

515       24-31. (2017)

516    4.  Kasper HB, et al. The pediatric proton consortium registry: A multi-institutional

517        collaboration in u.S. Proton centers. International Journal of Particle Therapy

518        2014;1:323-333

519    5.  McNutt T, Wong J, Purdy J, Valicenti R, DeWeese T. OncoSpace: A New Paradigm for

520        Clinical Research and Decision Support in Radiation Oncology Proceedings of the XVIth

521        International Conference on the Use of Computers in Radiation Therapy, Amsterdam,

522        2010, Editor Jan-Jakob Sonke, Published by Het Nederlands Kanker Instituut - Antoni van

523        Leeuwenhoek Ziekenhuis ISBN: 978-90-75575-29-3

524    6.  Michael Kazhdan, Patricio Simari, Todd McNutt, Binbin Wu, Robert Jacques, Ming

525        Chuang, and Russell Taylor, "A Shape Relationship Descriptor for RadiationTherapy

526        Planning" Medical Image Computing and Computer-Assisted Intervention

527        5762/2009(12), 100–108 (2009)

528    7.  Binbin Wu, Francesco Ricchetti, Giuseppe Sanguineti, Misha Kazhdan, Patricio Simari,

529        Ming Chuang, Russell Taylor, Robert Jacques, Todd McNutt, "Patient Geometry-Driven

530        Information Retrieval for IMRT Treatment Plan Quality Control", Medical Physics, 2009

531        Dec;36(12):5497-505

532    8.  Wu, B., Ricchetti, F., Sanguineti, G., Kazhdan, M., Simari, P., Jacques, R., Taylor, R.,

533        McNutt, T.: "Data-driven approach to generating achievable dose-volume histogram

534        objectives in intensity modulated radiotherapy planning". International Journal of

535        Radiation Oncology, Biology, Physics 2011 Mar 15;79(4):1241-7. Epub 2010 Aug

536    9.  Steven F. Petit, Binbin Wu, Michael Kazhdan , André Dekker, Patricio Simari, Rachit

537        Kumar, Russel Taylor, Joseph M. Herman, Todd McNutt," Increased organ sparing using

538        shape-based treatment plan optimization for intensity modulated radiation therapy of

539        pancreatic adenocarcinoma", Radiotherapy and Oncology, 102 (2012) 38–44.

540    10. Binbin Wu, Todd McNutt, Marianna Zahurak, Patricio Simari, Dalong Pang, Russell

541        Taylor, Giuseppe Sanguineti, "Fully Automated Simultaneous Integrated Boosted-

542        Intensity Modulated Radiation Therapy Treatment Planning Is Feasible for Head-and-

543        Neck Cancer: A Prospective Clinical Study". International journal of radiation oncology,

544        biology, physics 2012 Dec 1;84(5):e647-53

545  11. W Y Yang, J Moore, H Quon, K Evans, A Sharabi, J Herman, A Hacker-Prietz and T
546      McNutt, "Browser Based Platform in Maintaining Clinical Activities – Use of The iPads in
547      Head and Neck Clinics," Phys.: Conf. Ser. 489 012095 doi:10.1088/1742-
548      6596/489/1/012095 Int't Confernce on Computers in Radiotherapy, Melbourne AUS
549      2013

550  12. Robertson S, Quon H, Kiess, A., Moore J, Yang W.,Cheng Z, Afonso S., Allen M.,
551      Richardson M., Choflet A., Sharabi A., McNutt T, "A Data-Mining Framework for Large
552      Scale Analysis of Dose-Outcome Relationships in a Database of Irradiated Head and
553      Neck (HN) Cancer Patients", Medical Physics 42(7), 4329 (7/2015)

554  13. Marungo F, Robertson S, Quon H, Rhee J, Paisley H, Taylor R,  McNutt T, "Creating a
555      Data Science Platform for Developing Complication Risk Models for Personalized
556      Treatment Planning in Radiation Oncology."  48th Hawaii International Conference on
557      System Sciences (HICSS). Kauai, HI USA: IEEE; 2015

558  14. McNutt T., Moore K., Quon H. "Needs and Challenges for Big Data in Radiation
559      Oncology," Int'l J. of Radiation Oncology, Biology, Physics.  Published online: November
560      27 2015.  Print .  July 1, 2016 Volume 95, Issue 3, Pages 909–915

561  15. D Caruthers, S Brame, JR Palta, MP Hagan, E Wilson, C Cowan, L Yun, S Brown, LM
562      DeBerry, S Mutic, et al. Development and implementation of quality measures for the
563      survey based performance assessment of radiation therapy in the VHA. International
564      Journal of Radiation Oncology, Biology, and Physics, 99(2): E391-E392, 2017.

565  16. Owen, J., White, J., Zelefsky, M. and Wilson, J. (2009). Using QRRO™ Survey Data to
566      Assess Compliance with Quality Indicators for Breast and Prostate Cancer. Journal of the
567      American College of Radiology, 6(6), pp.442-447

568  17. Visscher SL, Naessens JM, Yawn BP, Reinalda MS, Anderson SS, Borah BJ. Developing a
569      standardized healthcare cost data warehouse BMC Health Serv Res. 2017 Jun
570      12;17(1):396. doi: 10.1186/s12913-017-2327-8.

571  18. M.R. Waddle, T. Kaleem, S.K. Niazi, L.J. White, J.M. Naessens, T.A. Rummans, D.I. Aljabri,
572      J.Y. Habboush, R.C. Miller. Cost of Acute and Follow-Up Care in Patients with Pre-
573      Existing Psychiatric Diagnoses Undergoing Radiation Therapy. IJROBP 2017, 99(5): 1321.

574      19. Rosenstein, B.S. (2017) Radiogenomics: Identification of Genomic Predictors for
575          Radiation Toxicity. Seminars in Radiation Oncology 2017 Oct;27(4):300-309.

576      20. West, C. and B.S. Rosenstein (2010) Establishment of a radiogenomics consortium.
577          Radiotherapy and Oncology 94(1):117-118.

578      21. West, C., D. Azria D et al. (2014) The REQUITE project: validating predictive models and
579          biomarkers of radiotherapy toxicity to reduce side-effects and improve quality of life in
580          cancer survivors. Clin. Oncol. (R. Coll. Radiol.) 26(12):739-742.

581      22. Constantine C., P. Parhar et al. (2008) Feasibility of accelerated whole-breast radiation in
582          the treatment of patients with ductal carcinoma in situ of the breast. Clinical Breast
583          Cancer 8(3):269-274.

584      23. Formenti S.C., D. Gidea-Addeo et al. (2007) Phase I-II trial of prone accelerated intensity
585          modulated radiation therapy to the breast to optimally spare normal tissue. Journal of
586          Clinical Oncology 25(16):2236-2242

587      24. Freedman, G.M., J.R. White et al. (2013) Accelerated fractionation with a concurrent
588          boost for early stage breast cancer. Radiotherapy and Oncology 106(1):15-20.

589      25. Raza S., S.C. Lymberis et al. (2012) Comparison of acute and late toxicity of two
590          regimens of 3- and 5-Week concomitant boost prone IMRT to standard 6-week breast
591          radiotherapy. Frontiers in Oncology 2:44.

592      26. Cooper, B.T., S.C. Formenti et al.  (2016) Prospective randomized trial of prone
593          accelerated intensity modulated breast radiation therapy with a daily versus weekly
594          boost to the tumor bed. International journal of radiation oncology, biology, physics.
595          95(2):571-578.

596      27. https://www.hl7.org/fhir/

597      28. Reeve BB, Mitchell SA, Dueck AC, et al. Recommended patient-reported core set of
598          symptoms to measure in adult cancer treatment trials. J Natl Cancer Inst 2014;106

599      29. Gaebel, Jan, Mario A. Cypko, and Heinz U. Lemke. "Accessing patient information for
600          probabilistic patient models using existing standards." Proceedings of the 10th
601          eHealth2016 Conference. Vol. 223. 2016.

602    30. Smith, W.P., Richard, P.J., Zeng, J., Apisarnthanarax, S., Rengan, R., Phillips, M.P.,

603        "Decision Analytic Modeling for the Economic Analysis of Proton Radiotherapy for

604        NSCLC."  Translational Lung Cancer Research, in press.

605    31. https://www.hl7.org/fhir/overview-arch.html.

## *Table 1. Examples of Big Data Use Cases in Radiation Oncology*

| Institution/Entity | Type of Database/Project | Source of Data/Tools | Magnitude | Key Features | Key Challenges |
|---|---|---|---|---|---|
| M-ROAR/University of Michigan | tumor staging, diagnosis code, pain scores, patient reported outcomes, and CTCAE scores | Oncology Information Systems, Treatment Planning System, and Electronic Health Record | >17,000 Patients since 2002 | Microsoft SQL Database; Self-service report building interface | Consistent/standardized physician and patient reported toxicities and reccurrence scoring |
| MD Anderson | Creation of Radiation Oncology Site Specific Templates for Data Input | Electronic Health Record (EPIC) | >40 specialty specific templates in Radiation Oncology with expansion into other departments | Specialty specific templates for standardized note generation | High level of customization in each site and department limits standardization in some elements |
| Pediatric Proton Registry Consortium | demographics, diagnosis and staging, baseline health status, chemotherapy and surgery, radiation details, diagnostic imaging, and follow-up | Oncology Information Systems, Treatment Planning Systems, and Electronic Health Record | >1800 patients from at least 13 centers | RedCap Tools; Collection of DICOM plan data | Funding; Data input efficiency |
| Oncospace | treatment planning data, patient reported outcomes, clinician assessments, disease response, diagnosis, and lab data | Oncology Information Systems, Treatment Planning System, and Electronic Health Record | >5000 patients from 4 centers | Tablet and web based data capture; Generation of notes from structured data entry; | Multi-institutional data standardization; Funding for maintenance and expansion |
| University of Pennsylvania | demographics, vital status, disease stage and prognostic indicators, genomic variants, details of systemic therapy and external-beam radiotherapy, and physician-reported toxicities | Oncology Information Systems, Electronic Health Record, Treatment Planning System, Cancer Registry, and Center for Personalized Diagnostics | >28,000 patients | Structure, site-specific templates; Only capture clinically symptomatic toxicities; Strong adoption by nurses | Physician adoption; Gathering of detailed progression information; Accurate identification of death events |

| US Veterans Health Administration (VHA) Radiation Oncology Practice Assessment | clinical measures, treatment planning information | Oncology Information Systems, Electronic Health Record, Treatment Planning System | Development is being finalized | novel tools to extra data including note processing; secure environment where data is housed locally | Development of custom tools to minimize manual data entry and support heterogeneous data sources |
|---|---|---|---|---|---|
| Mayo Clinic Florida | institutional data, demographics, tumor specific data, outcomes data, adverse events recorded in the EMR, and non-oncological diagnosis data | Electronic health record, administrative data, oncology information system,tumor registry, other disease specific registries | >3,000 patients | Includes administrative component with healthcare cost data capture | Toxicity reporting and data capture |
| The Radiogenomics Consortium | genomic data, treatment data, toxicity and outcomes data | Electronic health record, treatment planning systems | 132 institutions; > 6000 prostate patients and >4500 breast patients in specific projects | combined captured of genomic and treatment data | Data harmonization across different techniques and reporting methods |

**Table 2. Example Key Data Elements for Radiation Oncology**

| Key Data Element Category | Diagnosis = breast cancer | Diagnosis = lung cancer | Diagnosis = bone met |
|---|---|---|---|
| ICD-10 code | All, including laterality info | All, including laterality info | All, including location(s) |
| TNM staging | TNM staging | TNM staging | N/A |
| Performance Status | KPS | KPS | KPS |
| Toxicity Data Elements with CTCAE grade | Dermatitis | Dermatitis | Dermatitis |
| | Pain | Pain | Pain |
| | | Esophagitis | |
| | | Pneumonitis | |

| Recurrence Data Elements | Local recurrence | Local recurrence | Local recurrence |
|---|---|---|---|
|  | Regional recurrence | Regional recurrence |  |
|  | Distant recurrence | Distant recurrence | Distant recurrence |
| **Generic Data Element** **{name=___, description=___}** | Custom | Custom | Custom |