

1  
2 **Article Type: Special Issue Paper**

3 **Genomics, Bio specimens and other Biological Data: Current status and future**  
4 **directions**

5  
6 Barry S. Rosenstein<sup>1,2</sup>, Arvind Rao<sup>3</sup>, Jean M. Moran<sup>3</sup>, Daniel E. Spratt<sup>3</sup>, Marc  
7 ■ Mendonca<sup>4,5</sup>, Bissan Al-Lazikani<sup>6</sup>, Charles S. Mayo<sup>3</sup>, Corey Speers<sup>3‡</sup>

8  
9 <sup>1</sup>Department of Radiation Oncology, Icahn School of Medicine at Mount Sinai, New  
10 York, NY, 10029, USA; <sup>2</sup>Department of Genetics and Genomic Sciences, Icahn School  
11 of Medicine at Mount Sinai, New York, NY, 10029, USA; <sup>3</sup>Department of Radiation  
12 Oncology, University of Michigan Health System, Ann Arbor, MI, 48109, USA;  
13 <sup>4</sup>Departments of Radiation Oncology, Radiation and Cancer Biology Laboratories,  
14 Indiana University School of Medicine, Indianapolis, Indiana, 46202, USA; <sup>5</sup>Department  
15 of Medical & Molecular Genetics, Indiana University School of Medicine, Indianapolis,  
16 Indiana, 46202, USA; <sup>6</sup>Cancer Research UK Cancer Therapeutics Unit, Division of  
17 Cancer Therapeutics, The Institute of Cancer Research, London, SW7 3RP, UK

18  
19 ‡ To whom correspondence should be made:

20 Corey Speers, M.D., Ph.D.  
21 Department of Radiation Oncology  
22 University of Michigan, UH B2 C490  
23 1500 E. Medical Center Dr., SPC 5010  
24 Ann Arbor, MI 48109-5010  
25 Phone: 734-936-4300  
26 E-mail: cspeers@med.umich.edu

27  
28  
29 **Keywords:** genomics, data storage, informatics, radiogenomics, biospecimens

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.12912](https://doi.org/10.1002/mp.12912)

This article is protected by copyright. All rights reserved

30

31 **Conflict of interest disclosure:** The authors have no conflicts to disclose

32

### 33 **Introduction**

34 Recent technological advances have allowed for a more sophisticated  
35 understanding of the biology of tumors and an ability to generate massive data at an  
36 unprecedented pace. These advances now routinely allow for the assessment of  
37 genomics (DNA mutations and copy number alterations), transcriptomics (RNA  
38 expression levels), methylation profiles and protein and phosphoprotein abundance to  
39 unravel the biologic underpinnings of various types of cancers. This information is now  
40 being combined with various imaging modalities, histopathology, clinical and patient  
41 characteristics, and treatment information to allow for a systems-based approach to  
42 understanding and characterizing cancer. With these advances, however, new  
43 challenges have emerged in how to acquire, store, catalog, analyze, and integrate these  
44 varying types of biologic data. This article will review examples of successful integration  
45 of genomic and biologic data, the current state of this research, issues surrounding  
46 access, extraction, collection, and curation of the genomic and biospecimens data. It will  
47 also suggest recommendations for standardizations and next steps to improve data  
48 availability.

49 With the completion of the human genome project and the subsequent inception  
50 and completion of The Cancer Genome Atlas (TCGA) project, the acquisition, storage,  
51 and subsequent availability of large-scale genomic, transcriptomic, and proteomic data  
52 has led to an accelerated pace of discovery and understanding of cancer<sup>1-4</sup>. This data  
53 has led to the development of new, effective targeted agents, and ascertainment of this  
54 genomic and biospecimen data is now making its way into routine clinical practice<sup>5</sup>.  
55 Indeed, multiple groups have recently published the findings of molecular tumor boards  
56 and this molecular data is now beginning to be used, including in the NCI-sponsored  
57 MATCH and IMPACT trials, the AACR sponsored GENIE project, and ASCO sponsored  
58 TAPUR trial, to inform clinical decision making as it relates to disease prognosis,  
59 effectiveness, therapeutic benefit, and mechanisms of treatment resistance<sup>6-8</sup>. Other  
60 examples of the successful capture and annotation of genomic and biospecimen data

61 includes the Encyclopedia of DNA elements (ENCODE) project, and the International  
62 Cancer Genome Consortium (ICGC) project.

63 While the benefits of this molecular data in areas such as targeted drug  
64 development are increasingly clear, it's utility to predict radiation treatment toxicity and  
65 therapeutic response remains uncertain. While there are many reasons for this  
66 disparity, multiple initiatives including the REQUITE, RAPPER, Gene-PARE,  
67 RadGenomics, and canSAR projects are currently underway to collect, catalog, and  
68 make available this information<sup>9-13</sup>. The success of these radiation-associated  
69 databases, and subsequent projects, however, will depend on the ability for these  
70 databases to be accessed, annotated, integrated, and updated.

71

## 72 **State of the research**

73 The acquisition and storage of genomic and biospecimen data is currently the  
74 exception, not the rule in radiation oncology clinical practice. When this information is  
75 gathered, it is usually for research purposes with variable translatability into clinical  
76 practice. While the reasons for this lack of sample collection are numerous, a major  
77 limitation to specimen collection is the requirement that it be prospectively incorporated  
78 into research and non-research protocols. The collection and analysis of patient derived  
79 biospecimens requires institutional review board (IRB) approval. This approval, in turn,  
80 is dependent on a clearly formulated rationale for collecting the information, and  
81 safeguards regarding the utilization of the information and protection of potentially  
82 identifiable information. This requires foresight, resources (monetary, staff, and space-  
83 related resources), and patient and physician buy in. Some groups have begun to  
84 address this challenge by creating "boilerplate" language that can be incorporated into  
85 the standard consenting process for any patient undergoing radiation treatment. This  
86 consent includes language that allows for the de-identified patient genomic data to be  
87 used for research purposes, and is easily included in prospective trials as well as the  
88 regular clinic workflow<sup>14</sup>. Efforts to make this language and consent template more  
89 widely available are already underway.

90 An important area of radiation oncology research utilizing big data is  
91 radiogenomics, whose goal is the identification of genomic markers that are predictive

92 for the development of outcomes resulting from cancer treatment with radiation<sup>15</sup>. Work  
93 in radiogenomics has greatly benefited from creation of the Radiogenomics Consortium  
94 (RGC). The RGC was created in 2009 and is a cancer epidemiology consortium through  
95 the Epidemiology and Genomics Research Program of the NCI of the NIH<sup>16</sup>. The RGC  
96 now has 225 member investigators located at 131 medical centers in 32 countries. The  
97 common goal of the RGC membership is to share biospecimens and data so as to  
98 achieve large scale studies with increased statistical power to enable identification of  
99 relevant genomic markers. However, in order to accomplish this work and definitively  
100 discover and validate the critical genomic markers, access to the radiotherapy treatment  
101 information and long-term longitudinal follow-up data reporting details such as outcomes  
102 must be obtained for large numbers of patients. The RGC does not maintain a  
103 centralized biorepository, but serves to facilitate the development of collaborations  
104 between investigators with similar research goals who have assembled cohorts and  
105 collected data that can be synthesized into one large study. Although the RGC has  
106 successfully assembled large cohorts to perform adequately-powered studies, data  
107 harmonization remains a challenge for studies involving multiple patient cohorts treated  
108 with a variety of radiotherapy techniques and evaluated using multiple grading systems  
109 <sup>15</sup>. Although a proposed set of reporting requirements have been promulgated for  
110 research in radiogenomics, it would be advantageous if identical, or at least similar,  
111 case report forms were utilized for all radiogenomic research<sup>17</sup>.

112 An important example of the research projects launched by RGC investigators is  
113 the large multi-center REQUITE study (Validation of predictive models and biomarkers  
114 of radiotherapy toxicity to reduce side-effects and improve quality-of-life in cancer  
115 survivors)<sup>10</sup>. REQUITE addresses the challenge of data heterogeneity that, as for other  
116 big data projects, requires harmonization of the different outcome measures and  
117 confounding variables used in multiple cohorts. This study does not stipulate the  
118 radiotherapy protocols to be used but involves standardized case report forms across  
119 centers and countries to ensure data in identical categories are collected. The  
120 objectives of REQUITE are to: (1) Perform a multi-center, observational cohort study in  
121 which epidemiologic, treatment, longitudinal toxicity and quality-of-life data are collected  
122 from approximately 5,000 patients treated with radiotherapy for either breast, prostate or

123 lung cancer. (2) Produce a centralized biobank in which DNA is isolated from patients  
124 enrolled in the observational study and create a centralized data management system  
125 for secure collection, integration, mining, sharing and archiving of all project data. A key  
126 aspect of the centralized database is that it includes pre-treatment DICOM and DVH  
127 files. (3) Validate published SNP biomarkers of radiosensitivity and discover new  
128 variants associated with specific outcomes following radiotherapy, (4) Validate  
129 clinical/dosimetric predictors of radiotherapy toxicity and incorporate SNP biomarker  
130 data. (5) Design interventional trials to reduce long-term adverse cancer treatment  
131 effects. (6) Deliver interventional trial protocols using validated models incorporating  
132 biomarkers to identify patient sub-populations likely to benefit from interventions. (7)  
133 Serve as a resource exploitable for future studies exploring relationships between  
134 genetics and radiotherapy outcomes using developing technologies such as next  
135 generation sequencing. Those interested in becoming a member investigator of the  
136 RGC should contact Barry Rosenstein via email at [barry.rosenstein@mssm.edu](mailto:barry.rosenstein@mssm.edu).

137

### 138 **Access and Extraction**

139 With the decreasing costs and increasing availability of DNA and RNA  
140 sequencing, protein expression and metabolite assessment, the amount of data  
141 generated per patient continues to increase. Despite the increase in the availability of  
142 this data, the subsequent capture of this information for anything other than to answer a  
143 specific research question or direct a clinical treatment decision remains extremely  
144 limited. While the reasons are myriad, the amount and complexity of the data is a major  
145 factor. It is now common for germline and somatic testing of patients and tumors to  
146 include at least some of the following assessment: DNA sequencing assessing  
147 germline and/or somatic mutations and copy number variation, single nucleotide  
148 polymorphism (SNP) assessment, RNA expression (either through RNA sequencing or  
149 gene expression microarrays, epigenomic assessment, proteomic assessment (through  
150 mass spectrometry or reverse phase protein lysate arrays), metabolomic assessment,  
151 and pathologic assessment of tumor samples (through immunohistochemical staining,  
152 flow cytometry assessment, or through the creation of institutional or multi-institutional  
153 tissue microarrays (TMAs). In addition to the sheer volume of biological molecules being

154 assessed, the methods to analyze, interpret, and report the data is also varied. Issues  
155 of DNA sequencing read depth and sequence mapping for sequencing data continue to  
156 confound analyses of such data. Variation in algorithms utilized for DNA sequence  
157 mapping, variance and allelic calls also contributes to the complexity and heterogeneity  
158 in type and quality of the data. One must decide whether to collect and store pre-  
159 processed vs. raw data (i.e. normalized expression data vs. CEL files for expression  
160 microarrays or FASTQ vs. BAM files for DNA sequencing data). As an example, when  
161 raw DNA sequencing information is obtained, it usually comes from the DNA sequencer  
162 as FASTA or FASTQ files. FASTA and FASTQ format is based on simple text and  
163 contains the raw data of each sequence read. For FASTA files, each sequence starts  
164 with a ">" followed by the sequence name, a space and, optionally, the description. In  
165 addition, a separate FASTA file will include the quality information of the given read  
166 sequences. FASTQ files were developed to provide a convenient way of storing the  
167 sequence and the quality scores in the same text-based file. It bears noting that  
168 depending on the sequencing technique (Sanger vs. Illumina sequencing), different  
169 FASTQ files are generated based on the different ways in which quality is assessed  
170 between Sanger and Illumina sequencing. Because of this difference, the source of  
171 FASTQ data should be noted when storing the data as the encoding for the quality  
172 scores is different between Sanger and Illumina sequencing. In addition, paired reads  
173 are now routinely generated in which two reads are generated from the same single  
174 molecule to aid in sequence alignment. In this case (paired-reads data) two FASTQ files  
175 are created, one for the first read of the pairs and another one for the second, and the  
176 files should hold the reads exactly in the same order. Moving beyond simple read  
177 sequences FASTA or FASTQ files, alignment (SAM/BAM) and variation (VCF) files can  
178 also be created. Sequence Alignment Map (SAM) files were first created to store not  
179 only sequence and quality data (like in FASTQ files), but also mapping information for  
180 the sequences (i.e. where does each sequence align on the genome). To capture this  
181 more complex data, SAM files are tab-based and include 11-12 fields that fill one line  
182 and may include a header. SAM can express the same information as FASTQ, but also  
183 includes mapping information (see <https://samtools.github.io/hts-specs/SAMv1.pdf> for  
184 more information). SAM is rarely used as the format for data storage, instead, files are

185 stored in binary alignment mapping (BAM) format, which is a compact binary  
186 representation of SAM. It stores the same information, just more efficiently, and in  
187 conjunction with a search index, allows fast retrieval of individual records from the  
188 middle of the file. Because of its binary nature, BAM files are also much more compact  
189 than compressed FASTQ or FASTA files. Thus, when considering storage of genomic  
190 data once must decide upon file format storage (raw data in FASTQ vs. processed and  
191 mapped data in BAM). Finally, there are data access (and limitations to access to  
192 preserve data security and patient anonymity) and extraction issues that have made the  
193 wide-spread availability of this information a challenge.

194 While there are no quick or easy solutions to these challenges, many groups  
195 have already grappled with these questions and found useful solutions. For example, in  
196 the case of the REQUITE trial, standardized case report forms were developed for data  
197 collection of epidemiological and patient characteristics. Collection of clinical/pathologic,  
198 physics and treatment data were also standardized. Of critical importance, the full  
199 radiotherapy dose volume histogram was obtained for each subject, which provides  
200 substantial detailed dosimetric data. Data collection forms were provided in the different  
201 languages of the patients located in the multiple countries where they were enrolled into  
202 the study. Paper and web-based submission methods were provided in parallel.  
203 Submitted data underwent centralized quarterly quality control and plausibility checks  
204 for quality assurance according to a standardized quality assurance protocols. The  
205 database was enhanced to enable sample tracking in conjunction with the biobank  
206 information system, and empowered with user friendly interfaces to enable flexible data  
207 mining and data downloads in various formats. The database is only accessible to  
208 authorized persons via network and database passwords.

209

210

## 211 **Collection and Curation**

212 Important lessons can be learned from the challenges facing data extraction from  
213 health system-wide electronic medical records (EMR). Natural language processing  
214 (NLP) is part of a solution to extract data that is only available in free-text fields in  
215 EMRs. We are at a junction where development of a standardized format for collection

216 and storage of genomics data (i.e all BAM format) could potentially save significant  
217 resources in connecting genomics data to patient outcomes and dosimetric data.  
218 Independent validation is essential in the path towards the robust use of genomics data  
219 in clinical practice. By standardizing our clinical data collection, we can accelerate the  
220 discovery of which data are the most beneficial for specific classes of patients. One  
221 potential opportunity for increased capture and curation is to integrate with commercial  
222 (Flatiron) or organizational (CancerLinQ) platforms. These groups are already invested  
223 in data integration from EMR systems, and the increasing amount of clinically and  
224 commercially available genomic and biospecimen testing results may be extracted  
225 using these platforms. While these commercial and organizational platforms are still in  
226 their infancy, early integration into these groups may eliminate some of the challenges  
227 with later-stage integration. In addition, initial discussions with these groups could lead  
228 to standardization of collection and storage that could lessen, if not eliminate,  
229 subsequent challenges when the data is accessed/prepared for analysis. Certainly  
230 issues to consider in these initial discussions include: which format should be used to  
231 store data and whether raw or normalized data should be collected; should the data be  
232 normalized and if so which technique will be used; can EMRs be reconfigured to host  
233 and handle this genomic and biospecimen data; should biospecimen data be built into  
234 EMRs from radiation treatment unit vendors including Varian and Elekta (with Aria and  
235 Mosaik) and how do we limit redundancy or discrepant data in these biospecimen data  
236 sets. While the answers to these questions are not immediately obvious, working group  
237 consensus and advocacy will allow for a clearer path forward as we seek to collect and  
238 curate genomic and biospecimen data.

239

#### 240 **Specific Recommendations for standardizations**

241 The utility of genomic and biospecimen data collection and utilization will depend  
242 heavily on the quality and completeness of the data collected. In an ideal world, this  
243 data would be automatically collected and seamlessly integrated into other databases of  
244 collected data (patient outcomes, comorbidities, toxicity, dosimetric and treatment  
245 related information, etc.). While this is unlikely to be reality in the near term, the



246 following recommendations will allow for the gradual transition to this new reality.

247 These recommendations include:

- 248 1. Pool genomics and bio-specimen analysis templates among centers active in  
249 genomics for clinical research so that common features are universally captured  
250 and similarly named for ease of extraction in the future. Appropriate batch effect  
251 corrections across sample acquisition and preparation sites would be necessary  
252 prior to data collation<sup>18</sup>.
- 253 2. Develop a standard nomenclature for data collection. Similar to the TG-263 task  
254 group, the formation of a similar task group to standardize genomic and  
255 biospecimen data nomenclature and reporting would significantly aid in this  
256 process.
- 257 3. Harmonize the preferred format for standard fields to store genomics data within  
258 the hospital EMR with appropriate patient privacy safeguards built in. When  
259 housing within the EMR is not practical/feasible, uploading of genomics data with  
260 clinical outcomes and de-identified patient information into cBioPortal should be  
261 done (<http://www.cbioportal.org>)
- 262 4. Publish the recommendations such that individual institutions can request the  
263 major EMR vendors implement those standard fields. Concentrated and  
264 consistent pressure by end-users is likely to be more effective in implementing  
265 change than scatter shot, disjointed requests.
- 266 5. Identify institutions that would be able to perform validation of another institution's  
267 results through a standard data query. Data standardization and completeness is  
268 a key limitation on integrating this more globally, and quality assurance measures  
269 and standardized operating procedures that are universally available and  
270 implemented will be key to subsequent data utilization and integration.
- 271 6. Following examples like TCGA, the Sharing of standardized analysis pipelines  
272 enable the communication of "best practices" for concordant, reproducible and  
273 rigorous data analysis. Methods that enable the models to learn across  
274 institutional cohorts (i.e distributed learning), rather than requiring the data to be  
275 centrally stored can create viable alternatives for effective data learning and  
276 interpretation while being cognizant of potential privacy concerns<sup>19</sup>.

277

278

**279 Recommendations for next steps**

280

281 *Develop and conduct a survey to determine the state, quantity and quality of genomics*  
282 *and bio-specimens data in hospital EMRs*

283 In order to successfully fix a problem, one must first effectively identify and define  
284 said problem. Before the integration of genomic and biospecimen data can become a  
285 reality, one must first understand the present barriers and limitations in real-world terms.  
286 A survey that includes academic and industry participants, data generators and end  
287 users is critical to further identifying and then understanding the problem. The results of  
288 this initial survey will provide the basis for subsequent task group's efforts as they seek  
289 to assist the integration of genomic and biospecimen data into the radiation oncology  
290 space. Involvement of health ethicists and geneticists as well as health policy experts  
291 will also be key in navigating issues surrounding the housing of genomic data within an  
292 EMR (including health insurance and employer privacy concerns as well as protocols for  
293 notification should actionable germline mutations be identified.

294

295 *Share institutional best practices in data collection and storage and identify institutions,*  
296 *organizations, and companies who are willing to share current data templates*

297 While an effort that begins by trying to capture all data at all institutions on all  
298 patients is unlikely to be successful in the near term, an effort that includes multi-  
299 institutional and multi-tiered (i.e. academic, organizational, and industry) collaborations  
300 is likely to help move the field towards this greater goal. Critical to the successful  
301 advocacy for the integration of this data is receiving the support of large organizations  
302 already operating within this space. This includes dialogue with and endorsement from  
303 the American Society for Radiation Oncology (ASTRO), the American Society of Clinical  
304 Oncology (ASCO), the NRG, the National Institutes of Health (NIH), the National Cancer  
305 Institute (NCI) and the Global Alliance for Genomics in Health. As the sources for  
306 genomic and biospecimen data becomes increasingly available and complex, the  
307 inclusion of all parties associated with the data generation and usage in these

308 collaborations will be important. By identifying those groups that both have an interest in  
309 the integration of this data and who “touch” the data on a daily basis, potential pitfalls  
310 will be more readily identified, and avoided, as this process continues.

311  
312 *Develop and publish the harmonized template to standardize data collection, generation*  
313 *and analysis to facilitate connection to patient outcome and dosimetric data*

314 As was noted earlier in the article, the formation of a task group to address  
315 issues pertaining to standardized collection and nomenclature will be crucial to the  
316 successful integration of genomic and biospecimen data into radiation oncology  
317 treatment paradigms. One of the mandates of this working group will be the publication  
318 of recommended templates and nomenclature standardization that will allow for the data  
319 to be universally accessed and utilized. The publication of uniform access requirements  
320 and sharing of “Best Practices for Data Collection” will be critical to the success of this  
321 project. Similar efforts for establishing standardized analysis templates (for variant  
322 interpretation, gene expression analysis etc), will be essential to create datasets  
323 amenable to sharing and joint mining in the context of corresponding imaging and  
324 outcome data.

325  
326 *Final considerations*

327 In addition to the previously noted “next steps”, integration “discovery” and  
328 “validation” pipelines into the workflow will enable the more effective utilization of this  
329 data in the future. By carefully considering the collection and partitioning of these  
330 “discovery” and “validation” cohorts, subsequent integration of findings utilizing  
331 genomics and biospecimen data will be expedited. Critical to this collection, curation,  
332 and storage is the need for data housing standardizations that are HIPAA compliant and  
333 removes patient identifiable information. This also includes the need to incorporate our  
334 medical ethicist colleagues to consider the ethical issues surrounding the acquisition,  
335 storage, and reporting of genomic and biospecimen data.

336  
337 **Conclusions**

338 As we continue to translate the use of genomics data to guide treatment  
339 decisions for individual patients, we have an opportunity to accelerate this translation by  
340 developing and applying standard templates for data collection. By standardizing, they  
341 can be used to more robustly connect genomics and bio-specimen data directly to  
342 patient outcomes and dosimetric data. There is a lot of enthusiasm for how  
343 standardized nomenclature for organs-at-risk and targets will accelerate the analysis of  
344 dose and patient outcomes (AAPM TG-263)<sup>20</sup>. Similar potential exists within the  
345 collection, annotation, and storage of genomic and biospecimen space. Initial steps  
346 should include: standardizing nomenclature for data collection and harmonizing format  
347 of data collection and entry, pressuring EMR-vendors to build genomic and biospecimen  
348 data collection into the EMR platform, and establishing a task-group to generate specific  
349 guidelines governing the collection, analysis and reporting of genomic data. Successful  
350 completion of these steps will allow genomic and biospecimen data to be integrated into  
351 future data analysis as we seek to improve treatment efficacy and limit normal tissue  
352 toxicity.

353

354

355 References:

356

- 357 1. Lander ES, Linton LM, Birren B, et al. Initial sequencing and analysis of the human  
358 genome. *Nature*. 2001;409(6822):860-921.
- 359 2. Comprehensive genomic characterization defines human glioblastoma genes and core  
360 pathways. *Nature*. 2008;455(7216):1061-1068.
- 361 3. Comprehensive molecular portraits of human breast tumours. *Nature*.  
362 2012;490(7418):61-70.
- 363 4. An integrated encyclopedia of DNA elements in the human genome. *Nature*.  
364 2012;489(7414):57-74.
- 365 5. Bombard Y, Bach PB, Offit K. Translating genomics in cancer care. *J Natl Compr Canc*  
366 *Netw*. 2013;11(11):1343-1353.
- 367 6. TAPUR: Testing the Use of Food and Drug Administration (FDA) Approved Drugs That  
368 Target a Specific Abnormality in a Tumor Gene in People With Advanced Stage Cancer

- 369 (TAPUR). 2017; <https://clinicaltrials.gov/ct2/show/NCT02693535>. Accessed 9-25-2017,  
370 2017.
- 371 7. Coyne GO, Takebe N, Chen AP. Defining precision: The precision medicine initiative  
372 trials NCI-MPACT and NCI-MATCH. *Current problems in cancer*. 2017;41(3):182-193.
- 373 8. AACR Project GENIE: Powering Precision Medicine through an International  
374 Consortium. *Cancer discovery*. 2017;7(8):818-831.
- 375 9. Burnet NG, Barnett GC, Elliott RM, et al. RAPPER: the radiogenomics of radiation  
376 toxicity. *Clinical oncology*. 2013;25(7):431-434.
- 377 10. West C, Azria D, Chang-Claude J, et al. The REQUITE project: validating predictive  
378 models and biomarkers of radiotherapy toxicity to reduce side-effects and improve  
379 quality of life in cancer survivors. *Clinical oncology*. 2014;26(12):739-742.
- 380 11. Ho AY, Atencio DP, Peters S, et al. Genetic predictors of adverse radiotherapy effects:  
381 the Gene-PARE project. *Int J Radiat Oncol Biol Phys*. 2006;65(3):646-655.
- 382 12. Iwakawa M, Imai T, Harada Y, et al. [RadGenomics project]. *Nihon Igaku Hoshasen*  
383 *Gakkai zasshi Nippon acta radiologica*. 2002;62(9):484-489.
- 384 13. Tym JE, Mitsopoulos C, Coker EA, et al. canSAR: an updated cancer research and drug  
385 discovery knowledgebase. *Nucleic Acids Research*. 2016;44(D1):D938-D943.
- 386 14. Roychowdhury S, Iyer MK, Robinson DR, et al. Personalized Oncology Through  
387 Integrative High-Throughput Sequencing: A Pilot Study. *Science translational medicine*.  
388 2011;3(111):111ra121-111ra121.
- 389 15. Rosenstein BS. Radiogenomics: Identification of Genomic Predictors for Radiation  
390 Toxicity. *Semin Radiat Oncol*. 2017;27(4):300-309.
- 391 16. West C, Rosenstein BS. Establishment of a radiogenomics consortium. *Radiotherapy*  
392 *and oncology : journal of the European Society for Therapeutic Radiology and Oncology*.  
393 2010;94(1):117-118.
- 394 17. Kerns SL, de Ruyscher D, Andreassen CN, et al. STROGAR - STrengthening the  
395 Reporting Of Genetic Association studies in Radiogenomics. *Radiotherapy and oncology*  
396 *: journal of the European Society for Therapeutic Radiology and Oncology*.  
397 2014;110(1):182-188.
- 398 18. Leek JT, Scharpf RB, Bravo HC, et al. Tackling the widespread and critical impact of  
399 batch effects in high-throughput data. *Nature reviews Genetics*. 2010;11(10):733-739.
- 400 19. Jochems A, Deist TM, van Soest J, et al. Distributed learning: Developing a predictive  
401 model based on data from multiple hospitals without data leaving the hospital - A real life

- 402 proof of concept. *Radiotherapy and oncology : journal of the European Society for*  
403 *Therapeutic Radiology and Oncology*. 2016;121(3):459-467.
- 404 20. Mayo C, Moran JM, Xiao Y, et al. AAPM Task Group 263: Tackling Standardization of  
405 Nomenclature for Radiation Therapy. *International Journal of Radiation Oncology •*  
406 *Biology • Physics*.93(3):E383-E384.
- 407

Author Manuscript