1 **Keywords:** whole body magnetic resonance imaging, machine learning, deep

2 learning, random forests, convolutional neural networks, lesion detection, cancer

3

4 **Introduction**

5 Machine learning applications are ever-present in our daily activities, whether the

6 beneficiary is aware of it or not. Medical imaging, and, more specifically, clinical

7 radiology could not have remained unaffected by these advances [1-3].

8

9 The development and application of machine learning methods in radiology, has the

10 potential to support a series of clinical tasks, such as automatic lesion detection and

11 segmentation, lesion classification, patient risk stratification or patient outcome

12 prediction and may apply to radiological images of different modalities. Recently,

13 driven by the rapid progress in computational power and speed and the availability of

14 big datasets, the use of deep learning and, more specifically, convolutional neural

15 networks has revolutionised the field of automated analysis of radiological images by

16 accomplishing some of the aforementioned tasks with remarkable accuracy [4-6].

17

18 The developed machine learning methodologies seek to improve the diagnostic and

19 predictive performance of radiological scans and generate an, 'up to the hilt', time-

20 efficient and error-proof workflow for the reporting radiologist. The role of

21 computational tools is intended to be complementary and supportive to the radiologist,

22 potentially performing time-consuming tasks such as quantitative measurements; the

23 experienced radiologists' judgement remains the reference standard, taking many

24 other factors and non-imaging information into account. However, to quote Curtis

25  Langlotz of Stanford from the Radiological Society of North America (RSNA) meeting

26  in 2017: '*radiologists who use artificial intelligence, will replace those who don't*'.

27

28  Recent technological advances in magnetic resonance imaging (MRI), have allowed

29  whole body MRI (WB-MRI) to be performed clinically with acceptable image quality

30  and within reasonable time. The addition of diffusion-weighted imaging (DWI) in whole

31  body protocols, means that WB-(DW)-MRI is now becoming an increasingly important

32  tool in oncology for cancer diagnosis, staging and treatment response monitoring [7-

33  9]. A significant challenge when reading whole body MRI scans, is the increased

34  volume of resulting imaging data, especially when multi-parametric acquisitions are

35  used. The reading process can then become rather time-consuming, with increased

36  risk of misinterpretations. Also, whole body DWI for staging cancer patients has

37  limitations with respect to its diagnostic performance [10], as it may be prone to false-

38  positives resulting from tissues with normally occurring restricted diffusivity [11].

39

40  The National Institute of Health Research (NIHR) has funded a project (EME project

41  XXXXX), which aims to develop state-of-the art machine learning algorithms for the

42  automatic detection of malignant and benign lesions in multi-centre, multi-parametric

43  whole body MRI scans [12]. The study hypothesis is that the developed machine

44  learning tools will have the potential to improve the diagnostic performance and reduce

45  the reading time of whole body MRI scans. We discuss here our experiences from this

46  study and demonstrate the methodology employed and challenges met in the pathway

47  towards translating our methods into a potentially useful clinical tool.

48

49

50    **The XXXXXX (MAchine Learning In Body Oncology) study**

51    XXXXXX is a prospective, observational study, which aims to develop machine

52    learning methods and validate them by comparing the diagnostic performance and

53    reading time of WB-(DW)-MRI, when assessed alone and when assessed in

54    conjunction with machine learning output. The study does not collect patient imaging

55    data, but relies on data collected by other NIHR and CRUK-funded trials, referred to

56    as 'contributing studies' [13, 14]. XXXXXX is funded by the NIHR, Efficacy and

57    Mechanism Evaluation programme (EME project: XXXX) and is a collaboration

58    between the XXXXX and the XXXXX. Contributing studies' data are provided by the

59    XXXXX and XXXXX.

60

61    The study is divided into three phases, whereby in Phase 1 algorithms are developed

62    and evaluated for their accuracy to identify normal structures in whole body MRI scans

63    from healthy volunteers. In Phase 2 the developed algorithms will be further trained to

64    identify benign lesions and then tested and further refined for detecting cancer lesions.

65    Finally, in Phase 3 the algorithms will be tested in a large cohort of 'unseen' whole

66    body MRI data. As far as we are aware, XXXXXX is the first study that applies machine

67    learning techniques in WB-(DW)-MRI.

68

69    The XXXXXX study relies on whole body MRI data from a range of multi-centre trials,

70    and includes a range of cancer types, and thus the setting of the study is truly

71    pragmatic in clinical terms. As a result, the imaging data is relatively heterogeneous,

72    or "messy", which poses significant challenges to applying any statistical image

73    analysis approach. Current machine learning methodology requires the data to be

74    fairly homogeneous, in the sense that the training data from which task-specific

75 features are learned should be similar to the unseen test data, on which one wishes

76 to make predictions for. Figure 1 shows a block diagram identifying the XXXXXX

77 phases, during which the most significant challenges have been encountered to date

78 and for which our methodology required adaptation.

79

80 **1. Data acquisition**

81 The use of big datasets, is a desirable feature for either clinical outcome-driven

82 imaging studies or purely machine learning outcome-driven imaging studies. A large

83 cohort of examined patients can potentially increase the statistical power of primary

84 and secondary outcomes in clinical trials and can also boost the accuracy of the

85 employed algorithms in machine learning-related imaging studies, where larger

86 datasets are more likely to sufficiently capture the natural variability of both anatomy

87 and pathology. Thus, investigators turn to the use of retrospectively-acquired imaging

88 data or look into multi-centre collaborations to maximise the amount of available data

89 for their studies. However, this means that there will be data compliance issues. In

90 studies using, for example, CT datasets, the data is likely to be fairly homogeneous,

91 although differences in slice thickness or differences in the use of contrast may pose

92 challenges. However, in the MRI setting, as encountered in XXXXXX, there may be

93 extra significant variabilities in the data, including differences in imaging sequences,

94 between manufacturers and differences in acquisition parameters posing additional

95 challenges to the training and deployment of machine learning tools, as will be

96 described below.

97

98

99

**1.1 MRI systems and acquisition protocol variabilities**

The MRI systems used in multi-centre studies, will very commonly be of different manufacturers and different field strengths, have different coil characteristics and will be quality checked to different standards, even in the context of well-designed clinical imaging studies. This implies that images of inconsistent appearance and quality will be acquired throughout different centres. These differences are of little consequence to interpretation by the flexible human reader, who is trained to readily adapt to visual differences, but pose significant challenges for current machine learning algorithms. Furthermore, the introduction of functional imaging, which can now be incorporated into whole body protocols as in XXXXXX, means that the spatial and signal intensity discrepancies between images acquired in different centres, can be of particular importance in machine learning-related imaging studies.

This protocol variability in terms of anatomical localisation and signal intensity effects is demonstrated, using XXXXXX data, in Figure 2. Methods with which a number of the variability issues mentioned above, were mitigated in XXXXXX, are described in the 'Data preparation' section.

**1.2 Image quality**

The versatility of MRI is the modality's 'blessing and curse'. It is very common that image acquisition in the body may be compromised by patient factors such as movement, bowel gas, joint prosthesis or surgical material and imaging datasets of compromised quality can be 'passed through the sieve' of the clinical workflow, often out of necessity.

124    Repeating sequences may not always be practicable, because of time constrains or

125    patient exhaustion (especially if incorporating multiple sequences including DW-MRI).

126    It should be stressed, however, that the quality of the acquired datasets might have

127    been suitable for the objectives of the clinical study, involving human readers, and not

128    all of the issues are externally-triggered (for example distortions in echo planar

129    imaging (EPI) DWI acquisitions are unavoidable [15]), but they may cause very

130    significant challenges to the machine learning algorithms and be detrimental to their

131    performance.

132

133    This, highlights the importance of having imaging data with readiness level of *'Band

134    A'*, appropriate for the task at hand, as described by Lawrence 2017 [16], for machine

135    learning studies. It is acknowledged however, that when multi-centre data are

136    collected the scenario above is unrealistic, so removal of inappropriate or

137    compromised datasets might be unavoidable for the purposes of algorithm training

138    and also at test time, when predictions are made on new, 'unseen' data. We have

139    estimated that a proportion of the datasets employed in XXXXXX, were not suited for

140    machine learning purposes and had to be discarded. Figure 3 shows some of the

141    image quality issues we encountered in XXXXXX.

142

143    It is, therefore, highly recommended that MRI acquisitions for machine learning studies

144    are standardised to the highest possible degree and are performed and monitored by

145    an experienced research radiographer or by the local MRI physicist. This issue also

146    raises the much wider question of acquisition uniformity throughout the radiology

147    community, in order to harness the potential benefits of applying machine learning

148    techniques in the future.

**2. Data preparation**

Data preparation or pre-processing is an essential step in any machine learning study, whether related to imaging or not. In XXXXXX, where whole body MRI data from multiple imaging stations were acquired, we converted all our datasets in compressed Nifti format (nii.gz), in the interest of space and machine learning pipeline efficiency, after stitching images together according to slice location to form whole body volumes. It should be noted that, in case of DICOM data conversion to other 'headerless' formats, the original data should be retained so that header information can be 'glued' back to the converted images for uploading to the reading platform, as these accommodate almost exclusively DICOM data.

**2.1 Signal intensity standardisation**

As discussed earlier, the richness of acquisition schemes in MRI, comes with a major challenge. Unlike other medical imaging modalities, the image intensities in MRI do not have a fixed interpretation, not even within the same protocol or when acquired in the same body region, using the same scanner for the same patient [17]. In XXXXXX, this even applies between imaging stations in whole body acquisitions. This lack of a fixed meaning for intensities poses problems, not only when it comes to image quantification, but also in machine learning tasks, such as image segmentation. Therefore it is essential that an MRI signal intensity standardisation step is incorporated in the preparation pipeline before extracting the features in supervised learning algorithms or feeding the images in deep learning algorithms.

In XXXXXX we designed a specific pre-processing pipeline for intensity normalisation across images. We initially experimented with simple intra-subject intensity scaling,

174 based on signal normalisation using the 4th and 94th percentiles of the intensity

175 histogram, a somewhat arbitrary choice which has been shown to work well for brain

176 imaging [18]. However, in whole body imaging there is the challenge of inconsistent

177 anatomical coverage due to protocol variability, as discussed in Section 1.1. A number

178 of whole body volumes used in XXXXXX, fully included the head and neck regions

179 down to the lower limbs, while others only covered the body from the shoulders down

180 to knees (Figure 2). This violates the assumption that statistics, such as percentiles

181 obtained from the image intensity histograms, correspond to similar anatomical

182 regions. To address this, we make use of a rigid registration technique to

183 approximately align all images to a reference image. In this way, the field of view

184 between the tested and training images is normalised and similarity between the

185 histogram statistics is ensured.

186

187 This then allows us to employ Nyul's intensity normalisation technique [19], which

188 involves two stages. In the learning stage, a standard scale is derived from the

189 intensity histograms of the training images using ten, uniformly distributed, histogram

190 landmarks ranging from the 1st to the 99th percentile. In the testing stage, any new

191 image, following rigid registration to the reference image, can then be mapped to the

192 intensity standard scale, using the learned transformation from the training stage.

193 Figure 4 shows an example of using this pipeline on a whole body T2w volume.

194

195 Other histogram-based methods to perform intra and inter-subject signal intensity

196 standardisation for the same acquisition protocol are currently explored and compared

197 to the existing pipeline [20].

198

199 **2.2 Generating training data**

200 Generating training data for machine learning algorithms is one of the most important,

201 but also laborious and time-consuming processes. Manual, volumetric segmentations

202 performed by clinical experts, should be used to ensure reliable and accurate

203 algorithmic training. These labelled data, should also be used as the reference

204 standard to compare with, when evaluating algorithmic performance. Semi-automatic

205 or fully automatic methods can also be used to alleviate part of the workload, but it is

206 suggested that these segmentations are always double-checked and finalised by a

207 clinical expert. In XXXXXX, we used ITK-SNAP [21] to manually generate annotated

208 whole body images. Labelling of heathy structures (23 anatomical structures, including

209 organs and bones) occupied a significant proportion of Phase 1 of the project, but this

210 work was of paramount importance as in Phase 2 we are using a two-stage approach,

211 to identify cancer lesions, as will be discussed in Section 3.2.

212

213 **2.3 Image registration**

214 The use of multi-modal MRI data ('multi-channel' data as commonly referred to in

215 computer science terminology) has been shown to improve algorithmic performance

216 in tasks like brain lesion segmentation [22]. However, using multi-channel inputs for

217 algorithm training requires optimally registered imaging datasets between modalities,

218 so that annotated data from a single modality are used -in the interest of time-

219 efficiency- when generating training data. Anatomically-matched datasets from

220 different modalities, is a task which can be performed efficiently enough in the brain,

221 where minimal gross motion or anatomical deformation is expected between

222 acquisitions, with using a rigid registration algorithm.

223

224   In abdominal imaging, where there might be significant organ motion and deformation

225   between acquisitions, a rigid registration might not suffice. The task proved to be even

226   more challenging with whole body MRI data. Furthermore, when we attempted to

227   register DWI volumes to anatomical volumes, we encountered the extra challenge

228   from the geometrically distorted EPI-acquired, high $b$-value DW volumes [15]. We

229   qualitatively assessed registration between DWI and anatomical volumes, when using

230   a 12 degrees-of-freedom affine registration [23], but with mixed results. A non-rigid

231   registration using free-form deformations [24] was also tested, but the time required to

232   apply on the tens of whole body datasets used in XXXXXX was unacceptably long. At

233   this stage of XXXXXX, we simply use slice-matched acquisitions, resampled to match

234   the spatial resolution of the reference (T2-weighted) volumes. This aligns the majority

235   of structures, in particular bones, very well between modalities, but ignores differences

236   due to breathing or other movements of the subjects between scans.

237

238   A block diagram of the data preparation pipeline for XXXXXX, as described in Section

239   2, is shown in Figure 5.

240

241   **3. Machine learning pipeline**

242   **3.1 Choice of algorithm and feature crafting**

243   The choice of machine learning algorithm will depend on the task at hand.

244   Unfortunately, there is no 'one-fits-all' recipe and so, the choice comes down to a

245   recursive trial-and- error process, until the desirable performance and characteristics

246   are reached. The number of supervised, state-of-the-art, algorithms suited for imaging-

247   related tasks and their variants, but also the choice for the hyper-parameters in each

248   individual method may seem infinite; previous experience, already published results

249 and the quality and quantity of available data for training should provide guidance for

250 a good starting point.

251

252 Another important consideration for algorithm selection, is whether the model

253 interpretability is of interest for the task at hand. Deep learning algorithms have

254 demonstrated great accuracy in imaging-related tasks [6], but interpreting the

255 extracted features and the complex, non-linear relationships between them, which

256 take place in the hidden layers of the network, remains an almost impossible

257 challenge. Despite the fact that there are now ways to visualise the features that

258 activate specific neurons in a layer [25], the hidden layers of a deep convolutional

259 neural network still have the traits of a 'black box'.

260

261 In XXXXXX, we mainly tested and evaluated two algorithms; one state-of-the-art

262 ensemble algorithm based on classification forests (CFs) [26, 27] and one deep

263 learning algorithm based on convolutional neural networks (CNNs) [28]. Classification

264 forests are powerful, multi-label classifiers, which facilitate the simultaneous

265 segmentation of multiple organs. They have very good generalisation properties,

266 which means they can be effectively trained using a limited number of datasets. Both

267 of these traits were desirable in XXXXXX. Our convolutional neural networks

268 implementation was based on XXXXX [28, 29], an approach which has been shown

269 to perform very well in brain lesion segmentation with multi-parametric MRI data [22].

270 The details of the hyperparameters used for the CFs and network architecture for the

271 CNNs, can be found elsewhere [30]. CNNs performed consistently better in healthy

272 organ segmentation in Phase 1 of XXXXXX, so it was the algorithm of choice for Phase

273 2 of the project (lesion detection).

274 **3.2 Pipeline adjustments for task at hand and performance evaluation**

275 Whether the task at hand is organ or lesion classification, segmentation or detection,

276 the core of the pipeline will most commonly be an accurate and robust classifier. In

277 XXXXXX Phase 2 we were interested in lesion localisation and characterisation, rather

278 than segmentation. We therefore had to employ a scheme to evaluate the

279 segmentation algorithms used in Phase 1, but now in terms of detection. A specific

280 automatic evaluation procedure was implemented to calculate detection accuracy.

281 This uses as inputs the manual reference segmentation and the detection map from

282 the segmentation algorithm and calculates the true positive rate, positive predictive

283 value and F1 score, based on a user defined distance threshold (in mm). An example

284 plot of the accuracies for a range of detected lesions and manual segmentations

285 distance is shown in Figure 6.

286

287 We then used the CNN algorithm, developed in Phase 1 of XXXXXX, to evaluate the

288 performance of detected primary colon lesions from colorectal cancer patients,

289 scanned with whole body MRI [13]. We observed that lesion detection in whole body

290 scans was suboptimal with the CNNs, presumably due to the small fraction of lesion

291 volume occupying the scanned space, when compared to the whole body volume. The

292 complexity of intensities in background tissue and the lesion weak boundaries

293 appeared to be confusing the CNN [31].

294

295 We therefore, had to adapt our approach to become a two-stage process, whereby in

296 the first stage, the information from Phase 1 healthy organs/bones is used to identify

297 normality and in stage two the lesion is detected (Phase 2 of XXXXXX). Stage two can

298 be modular with respect to the anatomical location that the suspected lesion can be

299    found. According to this and the availability of training data, the architecture and

300    configuration of the used CNN can be modified to achieve optimal performance. This

301    work is now ongoing and the aforementioned process is depicted in Figure 7.

302

303    Finally, post-processing steps are required to prepare the machine learning output for

304    reading. In XXXXXX, the final probability maps obtained from the CNN were

305    smoothed, normalised and 'thresholded' to reduce false positives and improve visual

306    appearance for the reading process.

307

308    An integrated machine learning pipeline should also incorporate an objective

309    performance evaluation stage. The choice of performance assessment metrics will,

310    once again, depend on the examined data availability and the task at hand. In

311    XXXXXX, we evaluated segmentation tasks using cross-validation and a range of

312    overlap and distance metrics [32] and detection, using the scheme described above.

313

314    **4. Reading process**

315    **4.1 Reading platforms**

316    Traditionally, the picture archiving and communications system (PACS) is used for

317    hosting medical images and associated reader's reports. However, PACS is not

318    flexible enough to accommodate hanging protocols for machine learning outputs and

319    also, access from readers external to the hosting institution is not possible. In

320    XXXXXX, we have used a secure central imaging server (3Dnet™), provided by

321    Biotronics3D (London, UK) [33], to ensure that images and related machine learning

322    output, are hosted in an environment where customised hanging protocols can be

323    created and images are accessible by all readers via a standard internet connection.

324   A hanging protocol was created for XXXXXX readers in Biotronics3D, so that stitched

325   volumes from different imaging modalities, alongside the machine learning output, are

326   opened and browsed simultaneously, as shown in Figure 8. This setting also allows

327   for the anatomical localisation using cross-hairs and also fusion between the colour-

328   mapped machine learning output and any of the MRI modalities.

329

330   **4.2 Reading paradigm and reading process**

331   In XXXXXX, we have used a similar reading paradigm and case report forms (CRFs)

332   to the contributing studies [13, 14], with slight modifications to account for the machine

333   learning output effects in the source study's diagnostic performance and reading time.

334   Pilot testing of case report forms (CRFs) used randomised reads of anonymised scans

335   from colorectal cancer patients [13], which were performed by 6 independent readers.

336   Before the reading process, it was essential that the involved study readers met and

337   reached a consensus as to how the machine learning output will be interpreted (based

338   on suspicious lesion's size and location, detection probability value, etc.).

339

340

341   **5. Miscellaneous issues**

342   **5.1 Data and databases access**

343   In the era of machine learning in radiology, there is a need for well-organised, suitably

344   anonymised and accurately annotated database of images, annotations and metadata

345   throughout all stages of such studies. File nomenclature, which should be clearly

346   defined, needs to be available to all those involved with password-controlled access

347   to data. This may include multiple radiologists undertaking human expert

348   segmentation and standardisation of file names, which is essential for proper

349    management of the large number of files. In addition, version control is an important

350    concern, which needs attention during the iterative training process. As described in

351    Kohli 2017 [34] ideal datasets for radiology machine learning studies should be FAIR

352    (Findable, Accessible, Interoperable and Reusable). In XXXXXX, imaging data,

353    metadata and annotations were stored in a dedicated, secure workstation. Data

354    sharing and reporting was accomplished via Biotronics3D.

355

356    In another NIHR-funded study involving whole body MRI data (MAchine Learning In

357    MyelomA Response - XXXXX study, EME project XXXXX), the use of XNAT [35] for

358    the aforementioned tasks is currently being optimised. XNAT is an open-source,

359    extensible and flexible database system that allows for image, annotations and

360    metadata storage, sharing and management.

361

362    **5.2 Legal, ethical and clinical acceptance**

363    Data sharing agreements are an essential step in studies where data are being shared

364    between collaborators. Each involved party, needs to be clear and transparent

365    concerning the data to be shared and agreements with respect to background and

366    foreground intellectual property should also be in place. Local contract negotiations

367    are required prior to study commencement. Agreement for data sharing from the

368    source study funders, trial management group, trial steering committee and sponsor

369    should be obtained in writing.

370

371    Ethics considerations will vary depending on the arrangements of the primary source

372    studies. For the XXXXXX study, ethics approvals were available from each of the

373    contributing studies for use of the data and, in addition, an institutional research and

374     development approval with information governance agreement were all in place for

375     the XXXXXX protocol at the start of the study. Public and patient representation in the

376     trial management group is important to ensure that the patient's voice is heard in the

377     planning of the study and in the dissemination of the findings and public acceptance

378     of the use of machine learning support tools.

379

380     Clinical acceptance is also <span style="color:red">an important consideration in machine learning-related</span>

381     <span style="color:red">imaging studies</span>. The validation of the developed machine learning tools needs to

382     stand up to scrutiny and the methods used for testing the tools need to be clear to

383     clinical radiologists. In XXXXXX, we have devised a viewing framework that is widely

384     used by radiologists and incorporates the machine learning tools into a typical clinical

385     environment for testing.

386

387     **Discussion- Conclusion**

388     Machine learning algorithms can now perform image analysis tasks with performance

389     equal, or even superior, to the one achieved by human experts. Automatically derived

390     measurements and visual guides, obtained with machine learning techniques will

391     serve as a valuable aid in many clinical tasks and, most certainly, will transform the

392     ways we see and use medical imaging analysis tools.

393

394     We have used XXXXXX, a study that is looking into developing machine learning

395     methods for improving the diagnostic performance and reducing the reading time of

396     whole body MRI data, as a platform for identifying some of the main challenges

397     encountered in a clinical study involving machine learning. Our experiences are

398     described in this manuscript. Given the pragmatic setting of XXXXXX, we believe that

399    the methodological steps and challenges described here, can be of invaluable

400    assistance, and can serve as a guide, to groups who would like to apply similar studies

401    in the future, not only for MRI, but in radiology generally.

402

403    One of the most important considerations when designing a clinical study involving

404    machine learning, is data readiness. Acquired and used data should be assessed in

405    the context of appropriateness with quality and uniformity being the two most important

406    parameters to be considered. If these data traits cannot be assured upon design, then

407    appropriate steps towards upgrading the data level readiness should be taken or even,

408    manually identify the appropriate datasets if necessary. A robust machine learning

409    pipeline should be designed and implemented, a task which should now be

410    straightforward to accomplish, given that robust machine learning libraries, modules

411    and toolboxes are now freely available, to implement a vast amount of algorithms and

412    preparation/evaluation schemes. An important consideration for achieving the desired

413    clinical outcome is to effectively host the resulting machine learning output, along with

414    the clinical images, for reading. Once again, there are now a range of cloud-based

415    services available to facilitate this process. The reading paradigm and reading process

416    should be agreed by the readers in consensus. Finally, a range of legal, ethical and

417    clinical acceptance issues should be considered when attempting to incorporate

418    computer-assisting tools into clinical trials.

419

420    In conclusion, clinical studies involving the development and use of machine learning

421    methodology require careful design, if the study objectives are to be accomplished

422    and the employed methods to reach their full potential. The road from translating

423 computing methods into potentially useful clinical tools involves an analytical, stepwise

424 adaptation approach, as well as engagement of a multi-disciplinary team.

425

426

427 **References**

428 1. Wang S and Summers RM.Machine learning and radiology. Medical Image

429     Analysis 2012; 16(5): 933-951.

430 2. Erickson BJ, Korfiatis P, Akkus Z, and Kline TL.Machine Learning for Medical

431     Imaging. Radiographics : a review publication of the Radiological Society of

432     North America, Inc 2017; 37(2): 505-515.

433 3. Kohli M, Prevedello LM, Filice RW, and Geis JR.Implementing Machine

434     Learning in Radiology Practice and Research. American Journal of

435     Roentgenology 2017; 208(4): 754-760.

436 4. Chartrand G, Cheng PM, Vorontsov E, Drozdzal M, Turcotte S, Pal CJ, et

437     al.Deep Learning: A Primer for Radiologists. RadioGraphics 2017; 37(7): 2113-

438     2131.

439 5. Erickson BJ, Korfiatis P, Kline TL, Akkus Z, Philbrick K, and Weston AD.Deep

440     Learning in Radiology: Does One Size Fit All? Journal of the American College

441     of Radiology : JACR 2018; 15(3 Pt B): 521-526.

442 6. Mazurowski M, Buda M, Saha A, and R. Bashir M, *Deep learning in radiology:*

443     *an overview of the concepts and a survey of the state of the art*. 2018.

444 7. Takahara T, Imai Y, Yamashita T, Yasuda S, Nasu S, and Van Cauteren

445     M.Diffusion weighted whole body imaging with background body signal

446     suppression (DWIBS): technical improvement using free breathing, STIR and

447     high resolution 3D display. Radiation Medicine 2004; 22(4): 275-282.

448    8.    Koh DM and Collins DJ.Diffusion-Weighted MRI in the Body: Applications and

449          Challenges in Oncology. American Journal of Roentgenology 2007; 188(6):

450          1622-1635.

451    9.    Schmidt GP, Reiser MF, and Baur-Melnyk A.Whole-body MRI for the staging

452          and follow-up of patients with metastasis. European Journal of Radiology 2009;

453          70(3): 393-400.

454    10.    Wu L-M, Gu H-Y, Zheng J, Xu X, Lin L-H, Deng X, et al.Diagnostic value of

455          whole-body magnetic resonance imaging for bone metastases: a systematic

456          review and meta-analysis. Journal of Magnetic Resonance Imaging 2011;

457          34(1): 128-135.

458    11.    Padhani AR, Koh D-M, and Collins DJ.Whole-Body Diffusion-weighted MR

459          Imaging in Cancer: Current Status and Research Directions. Radiology 2011;

460          261(3): 700-718.

461    12.    XXXXX

462    13.    XXXXX

463    14.    XXXXX

464    15.    Le Bihan D, Poupon C, Amadon A, and Lethimonnier F.Artifacts and pitfalls in

465          diffusion MRI. Journal of Magnetic Resonance Imaging 2006; 24(3): 478-488.

466    16.    D. Lawrence N, *Data Readiness Levels*. 2017.

467    17.    Nyul LG, Udupa JK, and Xuan Z.New variants of a method of MRI scale

468          standardization. IEEE Transactions on Medical Imaging 2000; 19(2): 143-150.

469    18.    Sun X, Shi L, Luo Y, Yang W, Li H, Liang P, et al.Histogram-based

470          normalization technique on human brain magnetic resonance images from

471          different acquisitions. BioMedical Engineering OnLine 2015; 14(1): 73.

472    19.    Nyúl LG and Udupa JK.On standardizing the MR image intensity scale.
473           Magnetic Resonance in Medicine 1999; 42(6): 1072-1081.

474    20.    Madabhushi A and Udupa JK.New methods of MR image intensity
475           standardization via generalized scale. Medical Physics 2006; 33(9): 3426-3434.

476    21.    Yushkevich PA, Piven J, Hazlett HC, Smith RG, Ho S, Gee JC, et al.User-
477           guided 3D active contour segmentation of anatomical structures: Significantly
478           improved efficiency and reliability. NeuroImage 2006; 31(3): 1116-1128.

479    22.    Geremia E, Zikic D, Clatz O, Menze BH, Glocker B, Konukoglu E, et al.,
480           *Classification Forests for Semantic Segmentation of Brain Lesions in Multi-*
481           *channel MRI*, in *Decision Forests for Computer Vision and Medical Image*
482           *Analysis*, A. Criminisi and J. Shotton, Editors. 2013, Springer London: London.
483           p. 245-260.

484    23.    Studholme C, Hill DLG, and Hawkes DJ.An overlap invariant entropy measure
485           of 3D medical image alignment. Pattern Recognition 1999; 32(1): 71-86.

486    24.    Rueckert D, Sonoda LI, Hayes C, Hill DLG, Leach MO, and Hawkes
487           DJ.Nonrigid registration using free-form deformations: application to breast MR
488           images. IEEE Transactions on Medical Imaging 1999; 18(8): 712-721.

489    25.    *Convolutional Neural Networks for Visual Recognition.  Available via:*
490           *http://cs231n.github.io/understanding-cnn/. . Accessed in September 2018.*

491    26.    Breiman L.Random Forests. Machine Learning 2001; 45(1): 5-32.

492    27.    Glocker B, Konukoglu E, and Haynor DR, *Random Forests for Localization of*
493           *Spinal Anatomy*, in *Medical Recognition, Segmentation and Parsing*, S. Zhou,
494           Editor. 2015, Academic Press, Elsevier: London. p. 94-109.

495    28.    XXXXX

496    29.    XXXXX

497   30.   XXXXX

498   31.   Valindria V, Lavdas I, Cerrolaza J, O. Aboagye E, G. Rockall A, Rueckert D, et

499         al., *Small Organ Segmentation in Whole-body MRI using a Two-stage FCN and*

500         *Weighting Schemes.* 2018.

501   32.   Heimann T, van Ginneken B, Styner MA, Arzhaeva Y, Aurich V, Bauer C, et

502         al.Comparison and Evaluation of Methods for Liver Segmentation From CT

503         Datasets. Medical Imaging, IEEE Transactions on 2009; 28(8): 1251-1265.

504   33.   Biotronics3D. *Biotronics3D, Analyze-Collaborate-Discover. Available via*

505         *https://www.biotronics3d.com/public/. Accessed in September 2018.* 2018.

506   34.   Kohli MD, Summers RM, and Geis JR.Medical Image Data and Datasets in the

507         Era of Machine Learning—Whitepaper from the 2016 C-MIMI Meeting Dataset

508         Session. Journal of Digital Imaging 2017; 30(4): 392-399.

509   35.   XNAT. *XNAT, the most widely-used informatics platform for imaging research.*

510         *Available via https://www.xnat.org/. Accessed in September 2018.* 2018.

511

512

513

514   **Figure 1.** Block diagram depicting the methodological components that were

515   considered in XXXXXX study.

516

517   **Figure 2.** Different variants of a T2-weighted whole body MRI protocol. (a): Non-fat-

518   suppressed T2w images covering the body from the neck to mid-thighs (b): Non-fat-

519   suppressed T2w images covering the body from the top of the head to mid-calves and

520   (c): Fat-suppressed T2w images covering the body from the middle of the head to the

521 pelvis. Note the anatomical and signal intensity variability, which is of particular

522 importance in machine learning imaging studies.

523

524 **Figure 3.** Demonstrating some of the data quality challenges (artefacts) we

525 encountered in the datasets used in XXXXXX. Missing slices (a), RF interference (b)

526 and motion artefacts (c) on T2w images. RF field inhomogeneities leading to dielectric

527 shading (d) and RF noise in DW images.

528

529 **Figure 4.** Using intensity normalisation pipeline on a test image. (a): Original T2w

530 volume. (b): Same image, but scale-matched using Nyul's histogram-based method

531 described in the text, following rigid registration. The two volumes are displayed using

532 the same window/level settings. Employing Nyul's histogram-based method improved

533 healthy organ detection on previously unseen T2w images (c), when compared to

534 using the simple signal normalisation based on the $4^{th}$ and $94^{th}$ percentiles of the

535 intensity histogram (d).

536

537 **Figure 5.** Block diagram of the XXXXXX data preparation pipeline.

538

539 **Figure 6.** Primary colon lesion detection accuracies (true positive rate-TPR, positive

540 predictive value-PPV and F1 score) for different ground truth-detection distances,

541 when using the CF algorithm.

542

543 **Figure 7.** Two-stage lesion detection process, employed in XXXXXX Phase 2. During

544 stage one, the normal organs/bones are identified, based on Phase 1 training. During

545      stage two, lesion detection takes place. Stage two can be modular, with each module

546      algorithm training depending on anatomical position.

547

548      **Figure 8.** Biotronics3D view of the whole body volumes from different modalities and

549      the algorithm output, fused with the diffusion-weighted image from a colon lesion.

550