This is a post-peer-review, pre-copyedit version of an article published in Neuroinformatics. The final authenticated version is available online at: https://doi.org/10.1007/s12021-018-9387-8. The following terms of use apply: https://www.springer.com/gb/open-access/publication-policies/self-archiving-policy

AUTHOR'S PROOF

Neuroinformatics

 $\frac{1}{3}$

4

5 6

7 8

9

 $\begin{array}{c} 10 \\ 11 \end{array}$

https://doi.org/10.1007/s12021-018-9387-8

ORIGINAL ARTICLE

Using Deep Learning Algorithms to Automatically Identify the Brain MRI Contrast: Implications for Managing Large Databases

Ricardo Pizarro^{1,2} • Haz-Edine Assemlal² • Dante De Nigris² • Colm Elliott² • Samson Antel² • Douglas Arnold² • Amir Shmuel¹

© Springer Science+Business Media, LLC, part of Springer Nature 2018

12 Abstract

Neuroimaging science has seen a recent explosion in dataset size driving the need to develop database management with efficient 13processing pipelines. Multi-center neuroimaging databases consistently receive magnetic resonance imaging (MRI) data with 1415unlabeled or incorrectly labeled contrast. There is a need to automatically identify the contrast of MRI scans to save database-16managing facilities valuable resources spent by trained technicians required for visual inspection. We developed a deep learning (DL) algorithm with convolution neural network architecture to automatically infer the contrast of MRI scans based on the image 17intensity of multiple slices. For comparison, we developed a random forest (RF) algorithm to automatically infer the contrast of 18MRI scans based on acquisition parameters. The DL algorithm was able to automatically identify the MRI contrast of an unseen 1920dataset with <0.2% error rate. The RF algorithm was able to identify the MRI contrast of the same dataset with 1.74% error rate. Our analysis showed that reduced dataset sizes caused the DL algorithm to lose generalizability. Finally, we developed a 2122confidence measure, which made it possible to detect, with 100% specificity, all MRI volumes that were misclassified by the DL algorithm. This confidence measure can be used to alert the user on the need to inspect the small fraction of MRI volumes that 23are prone to misclassification. Our study introduces a practical solution for automatically identifying the MRI contrast. 24Furthermore, it demonstrates the powerful combination of convolution neural networks and DL for analyzing large MRI datasets. 25

Keywords Convolutional neural network · Deep learning · Magnetic resonance imaging · Database management · Automatic
 contrast identification

28

29 Introduction

A recent trend in the neuroimaging community has been to
increase dataset size in order to improve the power of studies
(Marcus et al. 2013; Weiner et al. 2013; Zuo et al. 2014).
Successfully managing large datasets requires multiple servers
for storage, software for efficient access and management, and
personnel, i.e., system administrators and software developers.

Ricardo Pizarro ricardo.pizarro@mcgill.ca

Amir Shmuel amir.shmuel@mcgill.ca

² NeuroRx Research, Montreal, QC, Canada

There have been great efforts to develop the framework and36software to facilitate the setup of a successful neuroimaging37database (Cheng et al. 2009).38

JrnIID 12021_ArtID 9387_Proof# 1 - 14/06/2018

Large neuroimaging databases have been setup in academ-39ic imaging centers and high-tech companies, often for accu-40 mulating data from multiple clinical trials. A generic schemat-41 ic for the hierarchy of the parties involved in acquiring brain 42magnetic resonance imaging (MRI) data is illustrated in 43Fig. 1. In the dataset we analyzed, distinct neuroimaging sites 44 acquire brain scans, then typically, a neuroimaging processing 45company serves as the MRI reading center for the trial and 46 analyzes the data to provide informative results. Clinical trials 47investigate the efficacy of a drug in one of two ways: (1) make 48 a cross sectional comparison between one patient group taking 49the drug of investigation and another patient group taking 50either placebo or the standard of care drug or (2) make a 51longitudinal comparison within a group of patients, before 52and after taking the drug. Investigators working in the clinical 53trials provide specific detailed guidelines including 54

¹ Montreal Neurological Institute, Departments of Neurology, Neurosurgery, Physiology, and Biomedical Engineering, McGill University, 3801 University, Room 786, Montreal, QC H3A 2B4, Canada

AUTIPH22 RtbS38 PRf 0 (4)0642018



Fig. 1 Schematic diagram to illustrate hierarchy of parties involved in acquiring the neuroimaging dataset we analyzed. For each trial, different neuroimaging sites were hired to scan multiple subjects. Each subject was scanned multiple times over different days. In each visit, a subject was

scanned multiple times with different MRI contrasts. The MRI datasets accounting for over 10⁵ scans are transferred for systematic and objective analysis

acquisition parameters and instructions for the patients to ac quire several different types of brain MRI contrasts acquired
 in different scans. This helps to minimize variability between
 datasets and increase precision of the results.

A typical clinical trial is associated with heterogeneity in 59scanner platforms and data export workflows across several 60 61 MRI facilities acquiring data for the study. Therefore, the MRI contrast cannot be reliably determined from the names of the 6263 corresponding data files or the description of the acquisition included within the file metadata. The MRI technicians may 64 not strictly follow the guidelines or use slightly altered acqui-6566 sition parameters. In addition, during data conversion, the DICOM header may be corrupted and not have the necessary 67 information. Finally, the scan parameters could not help to 68 69 distinguish T1-weighted images before and after contrast 70(i.e. T1P and T1C). Manual entries are required by the MRI technicians to distinguish T1P and T1C, a process that is 7172prone to error. The gold standard for identifying the contrast of a MRI volume is for a trained technician to use visual 73inspection. The number of MRI scans can reach into the thou-7475sands for a single trial, therefore human visual inspection takes up valuable time and resources (Gardner et al. 1995; 7677Pizarro et al. 2016). Thus, the ability to identify the contrast of the scan automatically would represent a significant reduction of time, labor and cost. 79

Current in-house practice for brain MRI contrast identifi-80 cation is a semi-automated process. The first step is to use a 81 decision-tree (DT) algorithm that exploits the acquisition pa-82 rameters recorded within the metadata of a MRI volume. In a 83 second step, the MRI volumes undergo interactive quality 84 control, at which time the operator can manually rename 85 MRI volumes, if the contrast has been incorrectly identified 86 by the DT algorithm. This semi-automated process is limited 87 in cases when the metadata does not contain sufficient infor-88 mation to correctly distinguish between similar contrasts e.g., 89 T1-weighted images before and after the injection of gadolin-90 ium contrast. Another limitation is the requirement for contin-91ual updates of the DT algorithm when new contrasts are intro-92duced, making the algorithm difficult to maintain. Finally, the 93 current approach fails to utilize any of the information 94 contained within the intensity values of the images, relying 95solely on 1% of the available data. 96

To our knowledge, no method currently exists to automatically identify the contrast of a MRI scan based on the image contrast. Such a fully automated algorithm can potentially overcome the problems that arise from existing solutions. 100

Neuroinform

101 Here we present a deep learning (DL) algorithm developed to automatically identify the contrast of a brain MRI volume. We 102 103 discuss the model architecture for neural networks used se-104quentially and how the DL algorithm was trained, validated, 105 and tested. The DL algorithm predicted data from an unseen dataset containing five contrasts with 0.15% error rate and 106 another dataset containing eight contrasts with 0.19% error 107 108rate. For comparison, we developed a random forest (RF) algorithm to automatically identify the contrast of a brain 109 MRI volume. The RF algorithm predicted data from an un-110seen dataset containing eight contrasts with 1.74% error rate. 111 112We demonstrate how smaller dataset sizes decrease the performance of the DL algorithm. We computed receiver operat-113ing characteristic (ROC)-based contrast-specific probability 114 thresholds, developed for the user of the DL algorithm. 115Finally, we discuss the utility of this algorithm and how it 116 117has been implemented in practice. We used the notation throughout the manuscript that boldface symbols represent 118 119vectors and matrix size is specified as $k \times l \times m$.

120 Methods

The database we analyzed, courtesy of NeuroRx Research, 121 was constructed through a processing pipeline to process over 12212310⁵ MRI datasets. The contrasts of the MRI volumes were identified in the beginning of the pipeline with a semi-124automated process consisting of a DT algorithm and manual 125126intervention process. We used the results of our semi-127automated process as the ground truth for this project. We developed two algorithms to automatically identify the con-128129trast of a MRI volume based on the ground truth. The first method used a RF classifier that considered scan parameters 130and basic image statistics as input features, while the second 131132method used a DL algorithm based on convolutional neural 133network. We assessed and compared the performance of the 134two algorithms in predicting an unseen dataset. We assessed the DL algorithm in more detail to describe its efficiency in 135predicting the various contrasts and develop an ROC-based 136contrast-specific probability threshold. 137

138 Neuroimaging Dataset

139The hierarchy of the parties involved in acquiring the neuro-140imaging data is illustrated in Fig. 1. In the dataset we analyzed,141an overall hundreds of neuroimaging sites were contracted to142acquire patient brain scans. Since part of the clinical trials143required longitudinal imaging, each subject had up to N144timepoints, where N is based on the study aim; a defined set145of contrasts was acquired at each timepoint.

A pipeline was developed for efficient processing of the
data. The pipeline consisted of multiple sequential phases,
meaning successful output of a previous phase was a

prerequisite prior to processing the next phase. MRI contrast 149identification is a critical task performed at the initial phase of 150the pipeline. We evaluated the DL algorithm performance in 151two stages. In stage I, we developed the DL algorithm and 152investigated how the dataset size influenced performance by 153generating equally balanced datasets of smaller size. We used 154a balanced dataset with 40,283 MRI scans, containing the five 155most common contrasts, and referred to, hereon, as the refer-156ence dataset. After developing the DL algorithm in stage I, we 157included additional and less common contrasts in stage II to 158test whether the algorithm was able to retain high performance 159when an unbalanced sample was used. In stage II, we used a 160 dataset with 45,785 MRI scans, incorporating three additional 161contrasts, and referred to, hereon, as the extended dataset. 162

Target Contrast, the Ground Truth

We previously developed a semi-automated algorithm in 164house to identify the MRI volume contrast using the 165following two steps. In step 1, a customized DT algo-166 rithm was used on acquisition parameters contained in 167the metadata of each MRI scan to identify the contrast. 168 In step 2, trained MRI experts visually inspected the 169dataset slice-wise, as part of an interactive quality control 170process. If necessary, MRI volumes were manually 171renamed to reflect the correct contrast. We used the 172semi-automated process to generate the target contrast, 173 \mathbf{c}_i , defined as a binary class vector of size $C \times 1$ where 174C is the total number of contrasts. The target contrast 175corresponded to our ground truth, which were used to 176train, validate, and test the DL and RF algorithms eval-177uated in this manuscript. 178

Deep Learning (DL) with Neural Networks 179

We developed a DL algorithm with neutral networks to infer 180the contrast of a MRI volume based on image intensity. The 181 DL algorithm consisted of an initial convolutional neural net-182work (CNN), which inferred the contrast on a single slice of 183 the MRI, and a subsequent dense neural network (DNN), 184which relied on the CNN output to infer a contrast for the 185entire MRI volume. We made the implementation openly 186 available on GitHub (https://github.com/AS-Lab/Pizarro-et-187al-2018-DL-identifies-MRI-contrasts) and developed the 188algorithm in Python with a Theano backend (Al-Rfou et al. 1892016) and compiled on Keras (Chollet 2015). Keras is a high-190level software package that provides extensive flexibility to 191easily design and implement DL algorithms. We manually 192selected all of the parameters that defined the network archi-193 tecture, including the number and type of layers, the number 194of layer nodes, and C, the number of final possible contrasts. 195

163

AUTH1P1222Auth238PPR#0149622018

196 Convolutional Neural Network (CNN) Architecture

The CNN architecture was based on an existing 197 198convolution-based neural network, made publicly available in Keras (Chollet 2015). The architecture was pre-199viously developed to identify the content on 32×32 200 201images from CIFAR10 (Krizhevsky and Hinton 2009). In particular, it consisted of a sequential model written 202in modular form as illustrated in Fig. 2. The first three 203 204modules were convolutional modules with a final down-205sampling operation. Each convolutional module 206 consisted of the following seven operations: 2D convo-207lution (Krizhevsky et al. 2012), batch normalization (Ioffe and Szegedy 2015), rectified linear unit (ReLU) 208 activation (Dahl et al. 2013), 2D convolution, batch 209210 normalization, ReLU activation, and max pooling (Krizhevsky et al. 2012). The final CNN module was 211212essentially a fully connected network (Bengio 2009) that 213inferred the contrast of the MRI volume. It consisted of reshaping the output of the final convolutional module 214to a linear array, comprising the following operations: 215216fully connected network (i.e. dense), batch normaliza-217tion, ReLU activation, fully connected network, and a softmax activation (Dunne and Campbell 1997). 218

The cited references provide in-depth detail regarding each operation; however, a brief description and

> Fig. 2 The convolutional neural network (CNN) architecture was comprised of 27 sequential layers. There were three repeating modules (green, blue, orange) of seven layers. The first module is detailed on the left; rectified linear unit is abbreviated as ReLU. The purple module was a fully connected network that contained the bulk (90%) of the network parameters and inferred the contrast on each slice. The data size is presented above in parentheses before and after each layer. For example, the input was one slice with size $1 \times 32 \times 32$, the input to the fully connected network (in purple) was of size $128 \times 2 \times 2$, and the final output was the inference with size $C \times 1$, where C is the total number of contrasts

motivation for each layer follows. The input layer pre-221pared the data to have size $1 \times 32 \times 32$ with the image 222processing steps described in 2.5.2. A 2D convolution 223layer generated the convolution of an image and a kernel 224of size 3×3 . For instance, the first 2D convolution layer 225estimated 32 kernels of size 3×3 . An image of size $1 \times$ 226 32×32 was convolved with each kernel to generate 32 227images, making the output of size $32 \times 32 \times 32$. A 2D 228convolution "viewed" different areas of the image, and 229as depth increased, the scope widened. Batch normaliza-230tion, as the name implies, normalized each image by re-231moving the intensity mean and standard deviation. Batch 232normalization was used to accelerate the training of a deep 233network, shown to reach optimum parameters in less steps 234(Ioffe and Szegedy 2015). In a ReLU operation, any in-235tensity value <0 was set to 0, while any value ≥ 0 was 236unchanged. A ReLU operation introduced a nonlinear 237function and efficiently replaced the previously used sig-238moid operation. A max pooling operation extracted the 239element with the highest value within a window of size 240 2×2 , effectively reducing the image size by 2. Max 241pooling was used to avoid over-fitting and reduce the 242computational cost. A softmax activation operation trans-243 formed the arbitrary values to probability values between 244[0, 1]. A softmax activation was used to make a final 245selection output as a contrast probability, \mathbf{p}_{C} , of size $C \times 1$. 246



Neuroinform

247 Dense Neural Network (DNN) Architecture

The CNN generated inference, \mathbf{p}_{C} , on a per-slice basis, yet we 248 249were interested in making an inference on the entire MRI 250volume. We developed an additional dense neural network (DNN), illustrated in Fig. 3, to compute a volumetric infer-251252ence, **p**. First, the input to DNN of size $n \times C$ was generated by computing CNN-generated \mathbf{p}_C on n = 30 slices. The DNN 253architecture was made from similar types of layers as the fully 254255connected portion of the CNN algorithm (purple module in 256Fig. 2). The DNN output, **p**, of size $C \times 1$ approximated the 257probability the MRI volume belonged to one of the possible 258contrasts, C.

259 Neural Network Parameters Estimation

The neural network parameters were estimated to minimize
the loss, defined in Eq. (2) as the categorical cross-entropy
(Murphy 2012). The parameter space was trained using
Adam RMSprop with Nesterov momentum (Dozat 2015).
The DL algorithm was compiled to run on a Nvidia Quadro



Fig. 3 A dense neural network (DNN) was used to make a volumetric inference from n slices. The DNN was comprised of five sequential layers and made a final inference on the MRI volume; rectified linear is abbreviated as ReLU. The parentheses represent the data size as input and output to each layer

K2200. The GPU ran about 20 times faster when compared to265a CPU. This increase in speed provided an efficient way to266iteratively explore and improve the DL algorithm. The full267training process took approximately 20 h to complete.268

We used a cross-validation scheme described in 2.5.4 to 269generate training, validation, and testing subsets. We used 270the training subset to estimate the neural network parameters. 271The training subset was loaded in batches, with sizes empiri-272cally determined by the limit of the GPU memory. The vali-273dation subset was used to estimate performance at the end of 274each estimation epoch. The algorithm ran for a total of 1000 275estimation epochs, where each epoch consisted of 20 estima-276tion steps. During each step, the algorithm used 600 slices, 277consisting of n = 30 slices taken from 20 randomly selected 278MRI volumes. The estimation procedure used 400 volumes 279per epoch for 1000 epochs, resulting in the algorithm going 280through the entire training subset approximately 10 times. 281This redundancy increased the probability that the algorithm 282used data from each MRI volume in the training subset at least 283once. After the estimation procedure completed 1000 epochs, 284the neural network parameters were saved to predict the data 285from the testing subset. 286

andom Forest (RF) Approach, Developed	287
or Comparison with DL	288

We developed a RF classifier to characterize the baseline per-289formance for a comparable algorithm that can automate the 290inference of a MRI contrast. A RF algorithm is a discrimina-291tive classifier that consists of an ensemble of DT classifiers, 292where the final classification is determined by summing the 293votes cast by each individual tree (Breiman 2001). RFs have 294been shown to be a powerful automatic classification ap-295proach in a wide range of classification tasks. The RF input 296 features used in this work consisted of the acquisition param-297eters, including the echo time (TE), repetition time (TR), and 298flip angle, which were extracted from the MRI scan metadata. 299Additional input features included basic image intensities sta-300 tistics: percentiles and mean. The complete list of features 301used for the RF algorithm can be found in section "Feature 302 extraction for the RF algorithm". 303

Automated Algorithms Evaluation

F

304

We evaluated the DL and RF automated algorithms in the 305 following way. First, we defined the metrics used to estimate 306 performance of the DL and RF algorithms at different devel-307 opmental stages. Second, we processed all the MRI volumes 308 to prepare the slices for input into the CNN. Third, we extract-309 ed features for the RF algorithm from the DICOM header and 310 MRI intensity profile. Fourth, we developed a cross-validation 311 scheme to generate uncorrelated subsets and evaluate the two 312algorithms with unseen data. Fifth, we assessed the algorithms 313

J Junio 1202 Arti D 9387 Proof# 1-14/06/2018

314in two stages, with datasets comprising of five contrasts and eight contrasts, respectively. Sixth, we studied how the dataset 315size affected performance for five contrasts. Finally, we plot-316

317 ted the distribution of the inferences made to describe how to

318 compute a contrast-specific probability threshold.

Performance Evaluation 319

We evaluated the DL and RF algorithms at different develop-320 mental stages by computing metrics that estimated perfor-321mance. We computed the error rate, ε , of each algorithm, 322 323 based on accuracy, to focus our attention on the incorrectly identified MRI contrasts. We estimated performance of an 324algorithm by measuring ε , defined to be: 325

$$\boldsymbol{\epsilon} = 1 - accuracy = 1 - \frac{1}{N} \sum_{i=1}^{N} c_i \cdot \mathbf{p}_i \tag{1}$$

326 where \mathbf{c}_i is the target contrast and \mathbf{p}_i is the algorithm-generated 328 contrast probability. The error rate was estimated by averaging the total number of MRI volumes used over N. 329

The DL algorithm was trained to minimize the loss defined 330 to be the categorical cross entropy, H. We defined the cross 331332 entropy to be a measure of the distance of the algorithmgenerated contrast probability from the target contrast. For 333 334 one MRI volume, *i*, we defined H_i as follows:

$$H_i(\mathbf{c}_i, \mathbf{p}_i) = -\sum_{j=1}^C c_i(j) \log p_i(j)$$
(2)

336 The cross entropy was estimated over all possible modalities, C. For multiple MRI volumes, the categorical cross en-339tropy was averaged over N, as in Eq. (1). 340

We generated confusion matrices to visualize the DL 341 and RF algorithms' classification performance per con-342 trast. We tracked the number of MRI volumes per con-343 trast where each algorithm-generated prediction agreed 344 or disagreed with the ground truth. The vertical axis 345of the confusion matrix is the contrast as determined 346 by the DL or RF algorithm, while the horizontal axis 347 of the confusion matrix is the ground truth. The num-348 bers along the diagonal of the confusion matrix reflect 349 the number of MRI volumes when the algorithm-350generated prediction and the ground truth agreed. The 351352numbers off of the diagonal of the confusion matrix reflect the number of MRI volumes when the 353algorithm-generated prediction and the ground truth 354disagreed. 355

We characterized the performance of the DL and RF 356 algorithms by computing the sensitivity and specificity 357 to estimate the ability of detecting each contrast. 358 359Sensitivity estimates the algorithm's capacity to correctly identify that a MRI volume is a particular contrast 360 while specificity estimates the algorithm's capacity to 361

337

correctly identify that a MRI volume is *not* a particular 362 contrast. The two metrics were defined as follows: 363

$$sensitivity = \frac{TP}{TP + FN}$$

$$specificity = \frac{TN}{TN + FP}$$
(3)

where TP is the true positive count, FN is the false 364 negative count, TN is the true negative count, and FP 366 is the false positive count. 367

We developed contrast-specific probability thresholds 368 for the user of the algorithm, which minimized the er-369 rors made by the algorithm, reflected by maximizing 370 sensitivity and specificity. A contrast-specific probability 371threshold, as opposed to taking the maximum value of 372 the probability vector, would increase the confidence in 373 the algorithm's ability of making the prediction. To that 374end, we computed ROC curves to find the operating 375 point that equally maximizes the algorithm's sensitivity 376 and specificity. We considered any probability value in 377 the range [0, 1] to be a candidate threshold. For each 378 candidate value, we computed the true positive rate 379 (TPR) and false positive rate (FPR), as follows: 380

$$TPR = sensitivity = \frac{TP}{TP + FN}$$

$$FPR = 1 - specificity = \frac{FP}{TN + FP}$$
(4)

383

We then computed the operating point by weighing TPR 384 and FPR equally and maximizing Youden's index (Youden 385 1950), defined as: 386

Youden's index =
$$TPR$$
- FPR (5)

388

390

389

Image Processing

We processed all MRI volumes to prepare the slices for input 391 into the CNN. All MRI volumes were masked, down-sam-392 pled, and normalized, as follows. (1) We defined a mask to 393 select n = 30 slices centered on the central slice of the volume. 394(2) Each slice was down-sampled to a 32×32 resolution, cho-395sen empirically to provide sufficient data for distinguishing 396 contrasts. Higher resolution slices did not improve results 397 and caused memory issues. (3) Each down-sampled slice 398 was then normalized over the intensity by subtracting the in-399 tensity mean and dividing by the intensity standard deviation. 400 The processed images were then recombined and used as in-401 put for the DL algorithm. In summary, each MRI volume, *i*, 402generated *n* images with 32×32 resolution, each labeled with 403 \mathbf{c}_i that specified the target contrast. 404

t1.1

t1.2

t1.3

t1.4 t1.5 t1.6 t1.7 t1.8 t1.9 t1.10 t1.11t1.12 t1.13 t1.14 t1.15

Table 1 DICOM acquisition parameters used as input features	Parameter Name	Dicom Field	Units
for automatic MRI contrast identification using a Random	Repetition Time	0018 × 0023	ms
Forest classifier	Echo Time	0018×0081	ms
	Echo Train Length	0018×0091	ms
	Inversion Time	0018×0082	ms
	Slice Spacing	0018 imes 0088	mm
	Percent Sampling	0018 imes 0093	Percent of acquisition matrix lines acquired
	Percent Phase Field of View	0018 imes 0094	Percent
	Pixel Bandwidth	0018 imes 0095	Hz
	Flip Angle	0018×1314	Degrees
	SAR	0018×1316	Watts per kilogram
	Contrast Media	0018 imes 0010	Name of contrast agent, if present ^{a,d}
	Sequence Variant	0018 imes 0021	Name of Sequence Variant ^{b,d}
	Scan Options	0018×0022	Name of Scan Options ^{c,d}

^a Encoded as 0 if no contrast agent present, 1 if any contrast agent present

^b Encoded as magnetization transfer (MT) = 1, inversion recovery (IR) = 2, SAT (saturation band) = 3, VB (variable bandwidth) = 4, other/empty = 0

^c Encoded as MTC/crSP=1, SK/crSP/crMP/crOSP=2, SK/crSP/crOSP=3, SS/crSP=4, SK/crSP=5, SP=6, SK = 7, other/empty = 0, where MTC = magnetization transfer contrast, SP = spatial presaturation, MP = magnetizationtization prepared, OSP = oversampling phase, SK = segmented k-space, SS = steady state

^d Manually entered by operator

405Feature Extraction for the RF Algorithm

We extracted features for the purpose of automatic con-406 407 trast identification using a RF algorithm. Table 1 summarizes the acquisition parameters extracted from the 408 MRI DICOM header file. Non-numerical parameters 409410 were mapped to integer values prior to being used as 411 an input feature.

In addition to the acquisition parameters described in Table 412 1, we extracted percentile intensities, as follows. We linearly 413normalized each MRI volume to the range 0-100, then we 414 included the 70th, 80th, 90th, 99th, and 99.5th percentile in-415416 tensities for each volume as features for input to the RF classifier. Lower percentile intensities were not considered, as in 417 general approximately 66.7% of the MRI volume is back-418 419 ground. These features were included as a coarse representation of intensity histogram shape, primarily in an effort to help 420 distinguish T1-weighted volumes with and without a contrast 421422 agent, i.e., T1C and T1P in Table 2.

Cross-Validation Scheme 423

We developed a cross-validation scheme to divide the dataset 424 into uncorrelated subsets: training, validation, and testing 425(Ripley 2007). The dataset contained MRI scans that were 426 427 highly correlated. Subjects were typically scanned in the same site and scanner, causing the MRI volumes in a trial to be 428 correlated not only in terms of underlying anatomical 429

structures, but also in terms of the image formation model. 430There was a need to provide a training subset which effective-431ly characterized the variability of the scans. Our cross-432validation scheme consisted of using a training subset with 433 scans from all clinical trials, which proved to be the key to 434getting generalizable results. To that end, we incorporated the 435following two steps into splitting the data. First, we used a 436single timepoint for each subject, even if multiple timepoints 437 were acquired. The first step ensured the subjects did not re-438 peat, thus reducing correlation across subsets. Second, we 439constructed the training, validation, and testing subsets, with 440randomly selected MRI volumes from all clinical trials and 441 imaging centers. The second step ensured we characterized 442 scanner variability. 443

The training subset corresponded to 60% of the MRI vol-444 umes from the dataset and was used to estimate the DL and RF 445 algorithms' parameters. The validation subset corresponded to 44620% of the MRI volumes from the dataset and was used ex-447 clusively to track the DL algorithm performance after each 448 estimation epoch completed, without feedback to the training. 449The testing subset corresponded to 20% of the MRI volumes 450from the dataset and was used to assess performance of the 451algorithms after training completed. 452

We used the cross-validation scheme described above to 453divide the reference dataset containing five contrasts used in 454stage I. The stage I testing subset was used to generate the 455 confusion matrix in Fig. 4. The stage I training and validation 456subsets were used to track the performance by training epoch 457

AU THIP 122 Rub S38 PR 10 1406 2018

Q1 t2.1 Table 2 Contrast distribution of MRI volumes used for cross-

validation is shown below

 $\begin{array}{c} t2.4 \\ t2.5 \\ t2.6 \\ t2.7 \\ t2.8 \\ t2.9 \\ t2.10 \\ t2.11 \\ t2.12 \\ t2.13 \end{array}$

Contrast	Abbreviation	Cross-validati	on subsets	
		Training	Validation	Testing
fluid-attenuated inversion recovery	FLR	4897	1615	1616
proton-density weighted	PDW	4854	1606	1606
T1-weighted post-contrast	T1C	4800	1590	1582
T1-weighted pre-contrast	T1P	4825	1593	1593
T2-weighted	T2W	4880	1612	1615
Reference dataset (stage I)		24,256	8016	8011
high-resolution T1-weighted	T1G	545	181	170
magnetic transfer ON	MTON	1399	446	450
magnetic transfer OFF	MTOFF	1408	449	453
Extended dataset (stage II)		27,608	9092	9085

The contrast abbreviation and number of MRI volumes in cross-validation subsets for training, validation and testing

plotted in Fig. 5. We repeated the cross-validation scheme
described above to divide the extended dataset used in stage
II containing eight contrasts. The stage II testing subset was
used to generate the confusion matrix in Fig. 6.

462 Stage I – Five Contrasts

463 In stage I, we used the reference dataset to develop the DL algorithm and assessed how the size of the dataset affected 464465 performance. As shown in Table 2, we divided the reference dataset as follows: 24,256 in the training subset, 8016 in the 466 467 validation subset, and 8011 in the testing subset. The MRI scans were labeled to be one of C = 5 contrasts, whose name 468 we abbreviated as follows: fluid-attenuated inversion recovery 469 (FLR), proton-density weighted (PDW), T1-weighted post-470



Fig. 4 The confusion matrix is illustrated above for five contrasts. The target contrast was taken from the ground truth identification. The inferred contrast was determined using the proposed deep learning algorithm



Neuroinform

Fig. 5 (a) Categorical loss and (b) accuracy for the deep learning as a function of estimation epoch for three dataset sizes. The solid line is training accuracy and the dashed line is validation accuracy. The shaded region is an estimate for the standard error (SEM). The legend specifies the size of the training set and validation set. We repeated size 1 five times. We tracked 13 randomly selected parts of size 2 and 25 randomly selected parts of size 3

Neuroinform



Fig. 6 Confusion matrices are illustrated for eight contrasts generated using (**a**) the random forest algorithm and (**b**) the deep learning algorithm. The target contrast was defined as the ground truth identification. The inferred contrast was determined using the two different algorithms

471 contrast agent (T1C), T1-weighted pre-contrast agent (T1P),
472 and T2-weighted (T2W). The DL algorithm was developed
473 using the reference dataset and results were summarized for
474 the testing subset as a confusion matrix in Fig. 4.

475 We investigated how the size of the dataset affected the performance of the algorithm, which allowed us to estimate the 476 number of samples needed to maintain performance. In other 477 words, how much data would we need to replicate an algorithm 478performing with similar accuracy? We tracked performance as 479the value of the accuracy, $1-\varepsilon$, and loss function, H, for both the 480training and validation subsets after each epoch completed. We 481 482explored three sizes: 40,283 MRI scans, 1610 MRI scans, and 75 MRI scans. First, we estimated the network parameters for the 483reference dataset of size 40,283 MRI scans to compute the per-484 formance mean and variance by epoch. We repeated the estima-485tion procedure for the reference dataset a total of five runs. Each 486run involved a random partition of the reference dataset to the 487488 training, validation, and testing subsets. Second, we divided the 489reference dataset into equal parts, each consisting of 1610 MRI scans. Each part was then divided into three subsets: 970 in 490 491 training subset, 320 in validation subset, and 320 in testing subset. We tracked the performance by estimation epoch in 13 492randomly selected parts. Using only 13 parts proved sufficient to 493estimate performance based on the statistical stationarity of the 494 mean and standard error of the mean (SEM) of the accuracy and 495loss function. Third, we divided the reference dataset into equal 496 parts, each consisting of 75 MRI scans. Each part was then 497 divided into three subsets: 45 in training subset, 15 in validation 498subset, and 15 in testing subset. We tracked the performance by 499estimation epoch in 25 randomly selected parts. Using only 25 500parts proved sufficient to estimate performance based on the 501statistical stationarity of the mean and SEM of the accuracy 502and loss function. To summarize the performance, we computed 503 the arithmetic mean over the iterations for each size. We then fit a 504decaying exponential function to the mean and SEM with pa-505rameters that minimized the error between fit and true value. The 506resulting performance by epoch for all three sizes is plotted in 507 Fig. 5. 508

Stage II – Eight Contrasts

In stage II we compared the performance between the DL 510algorithm and the RF algorithm. We then explored the ability 511for the DL algorithm to work as an identifying tool for the 512different contrasts and generated ROC curves to develop a 513contrast-specific probability threshold. We used the extended 514dataset which comprised of 45,785 MRI scans, which we 515divided into three subsets: 27,608 in training subset, 9092 in 516validation subset, and 9085 in testing subset. The MRI scans 517were labeled to be one of C = 8 contrasts, including the five 518stage I contrasts and the three additional contrasts: high-519resolution T1-weighted (T1G), magnetic transfer ON 520(MTON), and magnetic transfer OFF (MTOFF). The contrast 521acronym and number of MRI volumes in each of the three 522cross-validation subsets are summarized in Table 2. 523

Results

We developed a DL algorithm to automatically identify the 525contrast of a brain MRI. We designed the DL architecture to 526make inferences on unseen MRI data. We generated a confu-527 sion matrix to summarize the results obtained in stage I for 528five contrasts. We plotted training and validation performance 529metrics as a function of estimation epoch for three dataset 530sizes from stage I. For comparison to the DL algorithm, we 531designed a RF algorithm to make inferences based on features 532extracted from the MRI metadata and image statistics. In stage 533II, we generated confusion matrices to summarize the results 534obtained for eight contrasts for the DL and RF algorithms. We 535characterized the DL and RF algorithms' capacity as an iden-536tification tool for each contrast. Finally, we developed a meth-537 od to maximize performance by selecting a contrast-specific 538probability threshold accessible to the users of the algorithm. 539

509

524

540 Stage I – Five Contrasts

In stage I, we developed the DL algorithm for the five MRI 541 542contrasts that were acquired most frequently in the dataset we 543 analyzed. First, we evaluated the relationship between slice orientation and resolution. Most volumes in the database have 544545approximately 60 axial slices. Therefore, testing the performance with high resolution images was limited to the axial 546orientation. We compared the following three cases using the **Q2** 547 proposed algorithm for five contrasts (Table 3): 548

549 In all cases, we selected the slices centered around the 550 corresponding midline. Using the high-resolution images with 551 axial orientation did not improve the results. Therefore, we 552 continued developing and exploring sagittal orientation and 553 32×32 dimensions. It was encouraging that the algorithm 554 performed with low error rate; however, we were surprised 555 the higher resolution image did not improve the results.

We summarized the classification results from stage I with 556557a confusion matrix detailing the number of MRI volumes where the DL algorithm and the ground truth agreed and 558disagreed. Fig. 4 illustrates the confusion matrix for five con-559trasts generated with the DL algorithm on the testing subset. 560561The DL algorithm and the ground truth identification inferred the same contrast on nearly all MRI volumes quantified along 562the diagonal. There were some volumes where the DL algo-563564rithm and the ground truth identification did not agree, as quantified by the numbers off of the diagonal. The error rate 565of the DL algorithm for five contrasts was $\varepsilon = 0.15\%$. This 566 result raised the question of when the DL algorithm breaks 567down due to smaller dataset sizes and including additional 568contrasts, as we did in stage II. 569

570Next, we investigated how the size of the dataset affected the performance of the algorithm in order to estimate the 571dataset size required to preserve performance obtained when 572using the reference dataset. We tracked accuracy and categor-573ical loss of the training and validation subsets of three different 574sizes after each estimation epoch. Fig. 5 presents the accuracy 575576and categorical loss as a function of estimation epoch for the training and validation subsets plotted for different dataset 577578sizes. We plotted performance for estimation epochs [0, 150] 579 as the trends remained constant above 150 epochs. The mean value is plotted as dots by estimation epochs and the fits are 580plotted for the training and validation subsets in solid and 581582dashed lines, respectively. The performance generated from

t3.1 Table 3 Orientation and dimension explored with the DL algorithm

t3.2	Orientation	# Slices per volume	Dimension	Error rate (%)
t3.3	sagittal	30	32 × 32	0.15
t3.4	axial	30	32 × 32	0.35
t3.5	axial	10*	128 × 128	0.49

*- 10 slices were used because of memory issues

592

Table 4 The capacity of the DL and RF algorithms to detect eacht4.1contrast, as summarized by measures of sensitivity and specificity.Sensitivity estimates the capacity of the algorithm to correctly identifythat a MRI volume is a particular contrast and specificity estimates itscapacity to correctly identify that a MRI volume is not the particularcontrast

Contrast	deep learning (DL)		random forest (RF)	
	Sensitivity	Specificity	Sensitivity	Specificity
FLR	100.00%	100.00%	99.94%	100.00%
PDW	99.94%	100.00%	99.75%	99.95%
T1C	99.62%	99.95%	95.25%	99.29%
T1P	99.50%	99.91%	96.35%	99.01%
T2W	99.94%	99.99%	99.69%	99.96%
T1G	99.41%	99.98%	100.00%	99.97%
MTON	100.00%	100.00%	97.54%	99.93%
MTOFF	100.00%	99.97%	99.11%	99.83%

the data of size 1, in red, did not significantly differ from the 583performance generated from the data of size 2, in green. The 584DL algorithm reached peak performance in less estimation 585epochs when the dataset size was reduced to size 3, as 586reflected by the solid blue line compared to the solid green 587 and red lines. However, the DL algorithm was not able to 588generalize as well when the dataset size was reduced to size 5893, as reflected by the dashed blue line compared to the dashed 590green and red lines. 591

Stage II – Eight Contrasts

In stage II, we extended the application of the DL algorithm to 593eight MRI contrasts, and compared the DL algorithm perfor-594mance to that of the RF algorithm. We summarized the classifi-595cation results with a confusion matrix detailing the number of 596MRI volumes where the two algorithms and the ground truth 597agreed and disagreed. Fig. 6 illustrates the confusion matrices 598for eight contrasts generated with the RF and DL algorithms on 599the testing subset. The metadata was not readable for 21 MRI 600 volumes, resulting in 9066 total number of volumes tested with 601 RF instead of the possible 9087. The RF algorithm misclassified 602 several MRI contrasts that were correctly classified by the DL 603 algorithm. The DL algorithm outperformed the RF algorithm 604 across all contrasts, except for T1G. The error rate of the DL 605 algorithm for eight contrasts was $\varepsilon = 0.19\%$. The error rate of 606 the RF algorithm for eight contrasts was $\varepsilon = 1.74\%$. A lower 607 error rate generated by the DL algorithm indicates that there is 608 relevant information in the image intensity that is not captured by 609 the features used in the RF algorithm. 610

We characterized the performance of the DL and RF algorithms as a contrast-identifying tool to describe the results generated from the testing subset obtained in stage II. We computed the sensitivity and specificity describing the 614

663

Neuroinform

615 capacity for the DL algorithm to identify whether a MRI volume is or is not a particular contrast. We computed sensitivity 616 and specificity based on each contrast and summarized the 617 618 two metrics in Table 4. The resultant DL values were all 619 >99.41% and there were multiple contrasts with 100.00% sensitivity and specificity. These metrics break down the ac-620 curacy by highlighting that the DL algorithm could improve 621 622 overall accuracy by improving the sensitivity in detecting T1 contrasts: T1P, T1C, and T1G. The DL algorithm 623 outperformed the RF algorithm across all contrasts except 624 for the sensitivity generated from the T1G contrast. This indi-625 626 cates that the RF algorithm is using a feature extracted from the DICOM header or image intensity profile that is helping its 627 classification of the T1G contrast. 628

We visually inspected the MRI volumes that were
misclassified by the DL algorithm to better understand the
reason for misclassification. Our visual inspection identified
three groups of misclassification.

- (i) The following list identified six cases where the
 ground truth identification process failed and the DL
 algorithm succeeded in correctly inferring the
 contrast:
- Two T1C volumes were mislabeled as T1P and correctly
 identified by the DL algorithm.
- 639 Three mtOFF volumes were provided by the clinic as T1P.
 640 These cases were actually mtOFF and identified as such
 641 by the DL algorithm.
- One T1C volume did not have sufficient contrast failing
 the quality control process. The DL algorithm identified
 the volume as T1P.
- (ii) The DL algorithm wrongly inferred the contrast in thefollowing seven cases, as a result of an acquisition error:
- In two cases, T1P volumes confused by the algorithm as
 T1C because they contained high-intensity voxels at the
 extreme right edge of the image due to a ghosting artifact.
- In three cases, the delay between gadolinium injection and acquisition of the T1C volume was too short or too long, resulting in minimal gadolinium enhancement, thereby confusing the algorithm to infer the T1C volumes as T1P volumes
- In one case, the wrong parameter for TR was used, caus ing the PDW to appear similar to a T2W volume.
- In one case, a wrap-around artifact at the top of the head
 caused the algorithm to infer a T1G volume as a T1P
 volume.
- (iii) There were three cases where the DL algorithm failed toinfer the correct contrast, without a clear explanation:
- One T1P was wrongly inferred as a T1G.

- One T1C was wrongly inferred as a T1P.
- One T1C showed the effect of the contrast agent, but the effect was not as bright as in the regular case. This MRI volume was wrongly inferred as a T1P.

In order to develop a contrast-specific probability thresh-667 old, we computed the DL-generated probability for each MRI 668 volume in the testing subset to belong to a specific contrast, c_s . 669 Fig. 7 presents the distributions of the probabilities in the 670 corresponding subplot for each c_s . Let c_t be the target contrast 671 of a MRI volume. When $c_s = c_t$ the distribution of the proba-672 bilities was identified in green; conversely, when $c_s \neq c_t$ the 673 distribution of the probabilities was identified in red. This 674 presentation indicates a successful identification if data points 675 in green are close to 1 and data points in red are close to 0. The 676 DL algorithm was designed to generate a probability vector 677 $C \times 1$ with high probability for c_t , and low probability for all 678 other entries of the probability vector. The corresponding 679 global distribution with the range [0, 9000] is shown in the 680 inset of each subplot where only the MRI volumes whose 681 probabilities were 0.0 or 1.0 can be clearly seen. The zoomed 682 perspective distribution with the range [0, 20] is shown in the 683 subplot, where the MRI volumes whose probability was any-684 where between [0, 1] can be seen as well. The distribution in 685 Fig. 7 is reflective of the results in the confusion matrix in Fig. 686 6(b). The algorithm selected the contrast whose value was the 687 highest across the generated probability vector. It can be seen 688 that the specified contrasts that resulted in more errors in Fig. 689 6(b), such as T1P and T1C, generated probability distributions 690 with higher entropy in Fig. 7. Conversely, the contrasts with 691 fewer errors in Fig. 6(b), generated a nearly binary probabili-692 ties distribution in Fig. 7. 693

Next, we developed a contrast-specific probability thresh-694 old for the DL algorithm to minimize the errors reflected by 695 maximizing sensitivity and specificity. For each c_s and for 696 each candidate threshold, we computed TPR and FPR. We 697 generated ROC curves by plotting TPR versus FPR in Fig. 8 698 for each c_s . The inset shows a global perspective of the ROC 699 curve with the computed operating point as a red circle in the 700 upper left corner. For all contrasts, the red circle is proximal to 701the ideal operating point located in the upper left corner. In the 702 magnified corner of the ROC curve, we included the ideal 703 operating point with a green star and the operating point with 704a red circle. The red lines oriented at 45° reflect that we com-705 puted the operating point by weighing TPR and FPR equally 706 to maximize Youden's index (Youden 1950). Compared to 707 other contrasts, the T1G and T1P contrasts were associated 708 with a larger gap between the green star and the red circle. 709

In addition, we tracked each candidate threshold whose 710 Youden's index exceeded 0.98 to characterize how the threshold influences the performance. We plotted Youden's index as 712 a function of candidate threshold in Fig. 9. By comparing to 713 the results in Fig. 7, it is visually possible to identify how the 714

AU Jmi P 202 Att S38 P Rtf 1 (40642018



DL-generated probability

Fig. 7 Distribution for the deep learning (DL) generated probabilities by target MRI contrast. The subtitle specifies the target contrast. The inset is the same plot with the range from [0, 9000] to provide a global

perspective. The green bars reflect the probability for the MRI volumes targeted by the specified contrast, and the red bars reflect the probability when the target was not the specified contrast

715DL-generated probability distribution determined Youden's index. The plots in Fig. 9 are equivalent to Fig. 8, but they 716 illustrate that the threshold corresponding to the final operat-717 ing point was selected as the highest threshold to maximize 718 Youden's index. In the general case, there are multiple thresh-719 olds that equally maximize Youden's index. For some con-720721trasts, a particular threshold is critical to maximize Youden's index. However, for other contrasts there exists a range of 722 723 candidate thresholds that provide identical results. We selected 724the highest threshold in this range because it is the most con-725 servative threshold, thus providing a safety zone that mini-726 mizes classification errors. Importantly, the algorithm clas-727 sifies each MRI volume, and outputs not only the selected class but also the estimated probability that the classification 728 729 is correct. For probabilities lower than the contrast-specific 730 probability threshold, the user can inspect the volume in question, and make the final decision. 731

732 Discussion

An overview of the results is presented in this paragraph to frame the discussion. We developed a DL algorithm to automatically identify the contrast of a MRI volume. For five contrasts, the DL algorithm identified the contrast in unseen testing data with $\varepsilon = 0.15\%$. The algorithm converged to optimum parameters in fewer estimation epochs when training on smaller subsets. However, the DL algorithm did not generalize when the size was reduced to size 3, reflected by 740 a decrease in performance on the validation subset. For eight 741 contrasts, the RF algorithm identified the contrast from the 742 testing subset with $\varepsilon = 1.74\%$ and the DL algorithm with $\varepsilon =$ 743 0.19%. We characterized the DL algorithm's capacity to detect 744each contrast with sensitivity that ranged between [99.41%, 745100.00%] and specificity that ranged between [99.91%, 746 100.00%]. We characterized the RF algorithm's capacity to 747 detect each contrast with sensitivity that ranged between 748 [95.25%, 100.00%] and specificity that ranged between 749 [99.01%, 100.00%]. A contrast-specific probability threshold 750was computed for the DL algorithm with ROC analysis to 751indicate the user when to double-check particular contrasts. 752

We modified an existing algorithm to develop a new 753neural network to perform DL and this application has 754proven useful. The algorithm has been successfully im-755plemented into the processing pipeline. The tool we 756have developed can save database management teams 757 valuable resources, including hours spent by technicians 758 doing the trivial task of identifying the contrast of the 759MRI volume. In addition, this tool provides the methods 760 for inputting MRI into DL for new applications. 761

A convolutional neural network is analogous in function to the visual system. The first few layers extract low level features such as edges and background. As the layers get deeper, the low level features get combined to increase abstraction. We traced the output of all the layers of the CNN when using different MRI contrasts as input to reveal that: 767





Fig. 8 Receiver operating characteristic (ROC) curves are plotted in blue as true positive rate (TPR) versus false positive rate (FPR) for each contrast. We computed (FPR, TPR) points for each threshold based on the DL-generated probabilities distribution in Fig. 7. The green star indicates the ideal point (0, 1) with the corresponding TPR and FPR equality line. The red point indicates the Youden's index, defined in Eq.

(5), with the red line corresponding to TPR and FPR equality. The corresponding threshold is included in the subtitle with the contrast name. The inset provides the global perspective with the range [0.0, 1.0] and the black-dashed line represents the random chance of correctly selecting the contrast, generated with a naïve contrast identifier

Fig. 9 Youden's index is plotted for the candidate thresholds. The red point indicates the selected threshold corresponding to the operating point. The green line indicates the optimal operating point with no errors. Note: the red point was selected to be the maximum threshold to preserve the maximum value for Youden's index. The T2W contrast curve subtly increases at the selected threshold 1.0×10^{-5}



AUTIPH22 Rtb \$38 PRt 0 140 2018

Much to our surprise, a 32 × 32 sample of the entire image
 was sufficient to generate such high performance. A
 higher density image was not necessary to reach high ac curacy levels.

776 The flexible nature of the neural network architecture of the DL algorithm allows for any alterations that may be necessary 777 778 for future studies. Within contrast identification problems, the developer can change the number of contrasts by editing the 779 number of possible outputs, C. The interested investigator 780 could then either: (1) estimate the network parameters for 781the entire architecture of the DL algorithm from the beginning 782 783 or (2) use the stored parameters and only estimate the parameters for the part of the architecture that pertains to the new 784 785 number of contrasts, i.e., the last two layers of the CNN. In addition, the investigator could explore the CNN architecture 786depths by adding or removing convolutional modules. Finally, 787 it is possible to restrict the choices of outputs to those that are 788 789 relevant for a given clinical trial. One possible way is to incorporate a bias before the final layer of the CNN to restrict 790 specific choices that are not indicated within a study. 791

792 Our performance versus dataset size analysis, summarized in Fig. 5, brought forth two concepts regarding data sampling. First, 793 reducing the dataset to size 2 did not significantly affect the 794 795 performance of the algorithm, as illustrated by the red and green 796 curves. The DL algorithm did not generalize across imaging sites during our initial attempts. The key to getting robust generaliz-797 798 able results stemmed from taking representative samples from each site when we implemented cross-validation. Second, reduc-799 ing the dataset to size 3 significantly reduced the performance of 800 the DL algorithm, as illustrated in Fig. 5 by the red and blue 801 curves. A dataset size of 45 samples was unable to represent 802 the reference dataset, resulting in worse performance metrics. 803 804 This result emphasized the point that more data improves the power of the study and that a minimal number of samples is 805 required to represent the reference dataset such that performance 806 807 is not compromised.

We computed a contrast-specific probability threshold with 808 ROC curves from the generated probability distribution. As 809 810 illustrated in Fig. 9, six of the eight contrasts produced a range of values that generated identical maximum values for 811 Youden's index. The wider this range is, the more likely the 812 algorithm makes correct classification decisions. We selected 813 the most conservative value. In other words, although there is 814 a range of thresholds that give the same Youden's index, we 815 selected the highest threshold in the range. This translates to a 816 817 'safety zone', assuring that whenever the reported probability for correct classification is higher than the threshold, the pos-818 sibility of an error is reduced. 819

It should be stated that the DL algorithm is orders of mag-820 nitude slower to implement than the RF algorithm. The pa-821 rameters for the DL algorithm were estimated on a GPU. The 822 DL estimation procedure ran for 1000 epochs lasting approx-823 imately 6 h. However, the GPU ran about 20 times faster than 824 the CPU; therefore, a similar estimation run on a CPU would 825 take approximately 120 h. After training, the DL algorithm 826 was used for testing, a process that lasted approximately 827 78 min on the GPU, equivalent to 24 h on a CPU. In compar-828 ison, the RF algorithm ran for 6 min during the parameter 829 estimation process and 3 min during the testing phase, both 830 realized on a CPU. 831

One possible improvement to the DL algorithm is to in-832 clude the acquisition parameters, which we used as features 833 for the RF algorithm, as additional input to some layer deeper 834 than the convolution process. We demonstrated that RF was 835 able to classify all T1G volumes with 100% accuracy. 836 However, as a result of the feedback from users regarding 837 the single T1G volume misclassified by the DL, we attributed 838 the misclassification to technical acquisition errors. 839 Acquisition parameters could improve classification for one 840 contrast, but it may confuse the DL algorithm and cause ad-841 ditional error, as in our experience with the RF algorithm. In 842 our experience, relying on the header file is problematic for 843 the following reasons: (i) the header file is corrupted or miss-844 ing altogether, (ii) the operator manually enters parameters 845 incorrectly, and (iii) inconsistent parameters across different 846 neuroimaging sites. We work around these problems by 847 exploiting the advantage of the CNN approach in that it 848 doesn't rely on any headers, thus reducing the errors that arise 849 from manual intervention. 850

The ground truth identification process failed in the six 851 cases when MRI technicians made mistakes. Human work is 852 prone to error which is one motivation behind developing an 853 automated algorithm. However, using supervised learning to 854 develop a deep learning algorithm requires "ground truth" 855 data. It is the goal of an automated deep learning algorithm 856 to learn and supersede the current process. The six cases high-857 light the value in using a DL algorithm to avoid mistakes 858 made by technicians. Note that the ground truth identification 859 process does not suffer a systematic bias, since there were 860 multiple annotators, and a scan can be assigned to any given 861annotator. Thus, the noisy labels are uncorrelated and have 862

Table 5Comparison of error rate results generated from deep learningt5.1by using with 2D and 3D convolution

Algorithm	Error rate (%)		
	stage I – five contrasts	stage II – eight contrasts	
deep learning 2D	0.15	0.19	
deep learning 3D	0.49	0.87	

Neuroinform

863 little effect on the overall analysis, considering the size of the dataset. 864 Another potential improvement could be to implement 3D 865 866 convolution rather than 2D convolution. We edited the archi-867 tecture of the DL algorithm to incorporate 3D convolution instead of 2D convolution. We tested the trained algorithm **Q3**868

869 in stages I and II and summarize the results here (Table 5): Incorporating 3D convolution did not reduce the error rate 870 in detecting the contrast of the MRI volume. One possibility is 871 872 that when using 2D convolution, each MRI volume generates multiple samples at once. With 3D convolution, there are less 873 874 samples as the entire volume is loaded. In addition, 3D convolution requires excessive computational power, potentially 875 causing memory issues. Meanwhile, 2D convolution results in 876 a simpler algorithm requiring less memory and time to exe-877 cute, two key concepts when dealing with a large neuroimag-878 879 ing database.

880 The DL algorithm on its own did not accurately predict the 881 entire testing subset with 100% accuracy. This indicates that the algorithm cannot be used alone. One workaround is to 882 alert the user once the probability of correct classification 883 reported by the algorithm is lower than the operating thresh-884 885 old. The user can then inspect the MRI volume in question. When employing this approach, our proposed algorithm com-886 bined with the user alert process resulted in 100% success rate. 887 888 Yet another measure that can be added alongside our proposed algorithm is the DT portion of the semi-automated current 889 approach. In the rare case that the two methods disagree, the 890 user can go back and take a second look at the volume to 891 resolve which of the two methods is incorrect. 892

The algorithm was implemented into the processing pipe-893 894 line and is currently being used by technicians to validate contrasts of unknown MRI scans. The algorithm makes a pre-895 diction on the contrast within 4-5 s. The results thus far show 896 897 that, considering the user's inspection in cases of probabilities lower than the contrast-specific probability thresholds, the 898 overall success rate is 100%. This implementation will gener-899 900 ate feedback from a user perspective to allow for improve-901 ments in the future.

Conclusions 902

903 We developed an automated algorithm to identify the contrast of a MRI volume using DL with CNN architecture. The CNN 904 inferred contrast over n sagittal slices, followed by realizing 905 the volumetric inference using a DNN. The DL algorithm 906 identified between five and eight contrasts with a < 0.2% error 907 rate. We developed a RF algorithm for comparison and ob-908 tained a higher, 1.74% error rate for identification amongst 909 910 eight contrasts. We demonstrate that reducing the number of MRI volumes used for training to size 2 did not affect the 911 performance of the DL for five contrasts. Reducing the 912

number of MRI volumes used for training to size 3 signifi-913 cantly reduced the algorithm's capacity to generalize. We 914 characterized the DL algorithm for eight contrasts as a 915 contrast-specific identifying tool and computed contrast-916 specific probability thresholds as a reference for the end-user. 917

Information Sharing Statement

918

We made the implementation of our software openly available 919 on GitHub (see link below). We developed the algorithm in 920 Python with a Theano backend and compiled on Keras. Keras 921 is a high-level software package that provides extensive flex-922 ibility to easily design and implement deep learning algo-923 rithms. We created a Python virtual environment named 924 "deep env" with the requirements found in the following 925 GitHub link: https://github.com/AS-Lab/Pizarro-et-al-2018-926 DL-identifies-MRI-contrasts 927

Acknowledgements This work was supported by the Mathematics of 928 Information Technology and Complex Systems (Mitacs) Canada through 929 the Mitacs Elevate grant. This research was undertaken thanks in part to 930 funding from the Canada First Research Excellence Fund, awarded to 931 McGill University for the Healthy Brains for Healthy Lives initiative. 932933 We thank Laura Diamond and Micah Watts for English editing.

References

- 936 Al-Rfou, R., Alain, G., Almahairi, A., Angermueller, C., Bahdanau, D., Ballas, N., et al. (2016). Theano: A Python framework for fast com-937 putation of mathematical expressions arXiv preprint 938 939 arXiv:1605.02688.
- Bengio, Y. (2009). Learning deep architectures for AI. Foundations and 940 trends® in Machine Learning, 2(1), 1-127 %@ 1935-8237. 941
- Breiman, L. (2001). Random forests. Mach Learn, 45(1), 5-32 %@ 942 0885-6125. 943
- Cheng, X., Pizarro, R., Tong, Y., Zoltick, B., Luo, Q., Weinberger, D. R., 944 et al. (2009). Bio-swarm-pipeline: A light-weight, extensible batch 945 processing system for efficient biomedical data processing. Front 946 Neuroinform, 3, 35. https://doi.org/10.3389/neuro.11.035.2009. 947948
- Chollet, F. (2015). Keras.
- Dahl, G. E., Sainath, T. N., & Hinton, G. E. (2013). Improving deep 949 neural networks for LVCSR using rectified linear units and dropout 950 (acoustics, speech and signal processing (ICASSP), 2013 IEEE in-951ternational conference on): IEEE. 952
- Dozat, T. (2015). Incorporating Nesterov momentum into Adam. 953 Stanford University, Tech Rep, 2015.[Online]. Available: http:// 954cs229.stanford.edu/proj2015/054 report.pdf. 955**Q4**
- Dunne, R. A., & Campbell, N. A. (1997). On the pairing of the softmax 956 activation and cross-entropy penalty functions and the derivation of 957 the softmax activation function (Vol. 185, proc. 8th Aust. Conf. On 958959 the neural networks, Melbourne, 181).
- Gardner, E. A., Ellis, J. H., Hyde, R. J., Aisen, A. M., Quint, D. J., & 960 961 Carson, P. L. (1995). Detection of degradation of magnetic resonance (MR) images: Comparison of an automated MR image-962 quality analysis system with trained human observers. Acad 963 Radiol, 2(4), 277-281. 964

934

Deringer

AUIniP122 RtbS38 PR # 0 (4)06/2018

Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep
 network training by reducing internal covariate shift *arXiv preprint arXiv*:1502.03167.

- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). *Imagenet classifi- cation with deep convolutional neural networks* (advances in neural
 information processing systems).
- Marcus, D. S., Harms, M. P., Snyder, A. Z., Jenkinson, M., Wilson, J. A.,
 Glasser, M. F., Barch, D. M., Archie, K. A., Burgess, G. C.,
 Ramaratnam, M., Hodge, M., Horton, W., Herrick, R., Olsen, T.,
- McKay, M., House, M., Hileman, M., Reid, E., Harwell, J., Coalson,
 T., Schindler, J., Elam, J. S., Curtiss, S. W., van Essen, D., & WUMinn HCP Consortium. (2013). Human connectome project informatics: Quality control, database services, and data visualization.
- 980
 Neuroimage, 80, 202–219. https://doi.org/10.1016/j.neuroimage.

 981
 2013.05.077.

 982
 Mumbur K. P. (2012). Machine learning is a purchabilistic parametrization.
- Murphy, K. P. (2012). *Machine learning : a probabilistic perspective*(adaptive computation and machine learning). Cambridge, mass.:
 MIT Press.
- Pizarro, R. A., Cheng, X., Barnett, A., Lemaitre, H., Verchinski, B. A., Goldman, A. L., Xiao, E., Luo, Q., Berman, K. F., Callicott, J. H., Weinberger, D. R., & Mattay, V. S. (2016). Automated quality assessment of structural magnetic resonance brain images based on a supervised machine learning algorithm. *Front Neuroinform, 10*, 52. https://doi.org/10.3389/fninf.2016.00052.
- Ripley, B. D. (2007). Pattern recognition and neural networks:
 Cambridge university press.

UNCORDECT

1021

- Weiner, M. W., Veitch, D. P., Aisen, P. S., Beckett, L. A., Cairns, N. J., 993 Green, R. C., Harvey, D., Jack, C. R., Jagust, W., Liu, E., Morris, J. 994 C., Petersen, R. C., Saykin, A. J., Schmidt, M. E., Shaw, L., Shen, 995 996 L., Siuciak, J. A., Soares, H., Toga, A. W., Trojanowski, J. Q., & Alzheimer's Disease Neuroimaging Initiative. (2013). The 997 Alzheimer's disease neuroimaging initiative: A review of papers 998 published since its inception. Alzheimers Dement, 9(5), e111-999 e194. https://doi.org/10.1016/j.jalz.2013.05.1769. 1000
- Youden, W. J. (1950). Index for rating diagnostic tests. *Cancer*, *3*(1), 32– 1001 35.
- Zuo, X. N., Anderson, J. S., Bellec, P., Birn, R. M., Biswal, B. B., 1003Blautzik, J., Breitner, J. C. S., Buckner, R. L., Calhoun, V. D., 10041005 Castellanos, F. X., Chen, A., Chen, B., Chen, J., Chen, X., 1006 Colcombe, S. J., Courtney, W., Craddock, R. C., di Martino, A., Dong, H. M., Fu, X., Gong, Q., Gorgolewski, K. J., Han, Y., He, 1007 Y., He, Y., Ho, E., Holmes, A., Hou, X. H., Huckins, J., Jiang, T., 1008 Jiang, Y., Kelley, W., Kelly, C., King, M., LaConte, S. M., Lainhart, 1009J. E., Lei, X., Li, H. J., Li, K., Li, K., Lin, Q., Liu, D., Liu, J., Liu, X., 1010 Liu, Y., Lu, G., Lu, J., Luna, B., Luo, J., Lurie, D., Mao, Y., 1011 Margulies, D. S., Mayer, A. R., Meindl, T., Meyerand, M. E., 1012 Nan, W., Nielsen, J. A., O'Connor, D., Paulsen, D., Prabhakaran, 1013 V., Qi, Z., Qiu, J., Shao, C., Shehzad, Z., Tang, W., Villringer, A., 1014 Wang, H., Wang, K., Wei, D., Wei, G. X., Weng, X. C., Wu, X., Xu, 1015 T., Yang, N., Yang, Z., Zang, Y. F., Zhang, L., Zhang, Q., Zhang, Z., 1016 Zhang, Z., Zhao, K., Zhen, Z., Zhou, Y., Zhu, X. T., & Milham, M. 1017 P. (2014). An open science resource for establishing reliability and 1018 reproducibility in functional connectomics. Sci Data, 1, 140049. 1019 https://doi.org/10.1038/sdata.2014.49. 1020

🖄 Springer

Krizhevsky, A., & Hinton, G. (2009). Learning multiple layers of featuresfrom tiny images.

AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES.

- Q1. Please specify the significance of the data with bold emphasis in Table 2.
- O2. Missing citation for Table 3 was inserted here. Please check if appropriate. Otherwise, please provide citation for Table 3.
- Q3. Missing citation for Table 5 was inserted here. Please check if appropriate. Otherwise, please provide citation for Table 5.
- Q4. Please check if the URL captured correctly.

p