



Automatic motion artefact detection in brain T1-weighted magnetic resonance images from a clinical data warehouse using synthetic data

Sophie Loizillon^a, Simona Bottani^a, Aurélien Maire^b, Sebastian Ströer^c, Didier Dormont^{c,d}, Olivier Colliot^a, Ninon Burgos^{a,*}, for the Alzheimer's Disease Neuroimaging Initiative¹, the APPRIMAGE Study Group²

^a Sorbonne Université, Institut du Cerveau – Paris Brain Institute - ICM, CNRS, Inria, Inserm, AP-HP, Hôpital de la Pitié-Salpêtrière, Paris 75013, France

^b AP-HP, Innovation & Données – Département des Services Numériques, Paris 75012, France

^c AP-HP, Hôpital de la Pitié Salpêtrière, Department of Neuroradiology, Paris 75013, France

^d Sorbonne Université, Institut du Cerveau – Paris Brain Institute - ICM, CNRS, Inria, Inserm, AP-HP, Hôpital de la Pitié-Salpêtrière, DMU DIAMANT, Paris 75013, France

ARTICLE INFO

Keywords:

Clinical data warehouse
Deep learning
Motion
MRI

ABSTRACT

Containing the medical data of millions of patients, clinical data warehouses (CDWs) represent a great opportunity to develop computational tools. Magnetic resonance images (MRIs) are particularly sensitive to patient movements during image acquisition, which will result in artefacts (blurring, ghosting and ringing) in the reconstructed image. As a result, a significant number of MRIs in CDWs are corrupted by these artefacts and may be unusable. Since their manual detection is impossible due to the large number of scans, it is necessary to develop tools to automatically exclude (or at least identify) images with motion in order to fully exploit CDWs. In this paper, we propose a novel transfer learning method from research to clinical data for the automatic detection of motion in 3D T1-weighted brain MRI. The method consists of two steps: a pre-training on research data using synthetic motion, followed by a fine-tuning step to generalise our pre-trained model to clinical data, relying on the labelling of 4045 images. The objectives were both (1) to be able to exclude images with severe motion, (2) to detect mild motion artefacts. Our approach achieved excellent accuracy for the first objective with a balanced accuracy nearly similar to that of the annotators (balanced accuracy > 80 %). However, for the second objective, the performance was weaker and substantially lower than that of human raters. Overall, our framework will be useful to take advantage of CDWs in medical imaging and highlight the importance of a clinical validation of models trained on research data.

1. Introduction

Recently, hospitals have created clinical data warehouses (CDWs) that gather medical images from thousands to millions of patients (Nordlinger et al., 2020; Karami et al., 2017; Mia et al., 2022). These resources represent an exceptional opportunity to develop computational tools (Jannot et al., 2017). In contrast to research datasets where acquisition protocols are well standardised, the quality of CDW images is highly heterogeneous. Images come from different hospitals over several decades and diverse machines are used with no homogenisation on the acquisition parameters (Mia et al., 2022). Once an image is acquired at the hospital, it will immediately be saved in the picture archiving

and communication system (PACS), meaning that a non negligible number of unusable images will be archived. For instance, if a patient moves during the acquisition, the corrupted image will still be stored in the PACS. Therefore, quality control (QC) is a fundamental first step before developing any machine learning project on a CDW. Magnetic resonance (MR) images are sensitive to motion induced by patient movement during the acquisition process. As they require a long acquisition time, subjects are likely to move during the examination, which causes artefacts in the reconstructed image. Motion artefacts primarily manifest in the phase encoding direction due to the faster readout acquisition compared to repetition time, resulting in blurring, ringing,

* Corresponding author.

E-mail address: ninon.burgos@cns.fr (N. Burgos).

¹ Data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.usc.edu). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: http://adni.loni.usc.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf.

² Members of the APPRIMAGE study group can be found at <https://www.aramislab.fr/apprimage>.

ghosting or signal loss depending on the timing and spatial changes during acquisition (Wood and Henkelman, 1985). Thus, motion can be a serious confounding factor for further neuroimaging analyses. This can have dramatic consequences when the presence of motion artefacts is correlated with a diagnosis of interest (e.g., patients with a specific disease have a tendency to move more often) since it would lead to biased models. Previously, based on a CDW gathering data from 39 different hospitals, we found that 25% of MRIs were considered totally unusable for further processing, and almost a third had a very low quality especially due to motion (Bottani et al., 2021). Another study conducted in a single hospital showed that the prevalence of repeating an MRI examination due to the presence of motion was up to 20% of all the acquisitions (Andre et al., 2015). Beyond the cost that this represents for hospitals, these studies are highlighting the fact that many images present in the PACS are simply unusable, often because corrupted by motion artefacts. Therefore, it is important to be able to automatically exclude such images before conducting any study on a CDW. What is more, several works (Hedges et al., 2022; Reuter et al., 2015; Alexander-Bloch et al., 2016) have shown the impact of motion in the use of brain imaging software packages such as Freesurfer (Fischl, 2012) or SPM (Penny et al., 2011). The presence of motion artefacts induces a constant bias in the morphometric analyses, leading to a reduced estimation of the grey matter volume (Reuter et al., 2015) as well as of the cortical thickness (Hedges et al., 2022).

Quality control is needed to fully exploit the potential of CDWs and important efforts have been made to propose automatic QC tools, including the detection of motion artefacts (Esteban et al., 2017; Sadri et al., 2020; Lei et al., 2022; Ravi et al., 2022; Shaw et al., 2021). Esteban et al. (2017) introduced MRIQC, a pipeline for the automatic QC of 3D brain T1-weighted (T1w) MRI based on image quality measures (IQMs). It enables the extraction of IQMs such as the signal-to-noise ratio, the contrast-to-noise ratio or the volume of grey and white matter. This method relies on an extensive pre-processing pipeline using neuroimaging software packages such as ANTs (Avants et al., 2008) and FSL (Jenkinson et al., 2012), which are only usable on good quality images and therefore incompatible with CDWs. Sadri et al. (2020) developed MRQy, a quantitative tool to quickly determine relative differences in MRI volumes within and between large MRI cohorts. As MRIQC, MRQy is based on the extraction of IQMs but it does not require extensive pre-processing thanks to an Otsu thresholding to distinguish the foreground, which includes the whole head, from the background. QC methods based on convolutional neural networks (CNNs) have also been proposed (Lei et al., 2022; Sujit et al., 2019; Fantini et al., 2021; Küstner et al., 2018; Iglesias et al., 2017). They have the advantage of learning features without knowing a priori which are the most adapted. Sujit et al. (2019) developed an ensemble deep learning model based on CNNs to automatically evaluate the quality of multi-centre structural brain MR images. A limitation of this work is that it relies on images acquired following a well-defined research protocol, which are not representative of the heterogeneity of clinical images. Lei et al. (2022) presented a multi-task CNN framework for artefact-based MRI quality assessment, which not only provides a quality score but interprets the cause of the poor image quality. Image rulers, which consists of several versions of the original MRI slices with one type of artefact (noise or motion), are used during inference time. Each of these MRIs will be run through the trained CNN and the different outputs will be compared with the test image. The use of a single image ruler consisting of different versions of a single scan makes this method incompatible with the high heterogeneity that characterises CDWs, where different types of artefacts can coexist in a single image (e.g., noise and motion).

Previously, we proposed the first framework for the automatic QC of T1w brain MRI in a CDW using deep learning techniques (Bottani et al., 2021). A unique set of 5500 MRIs was manually annotated with a three-level grade for three characteristics: noise, contrast and motion. According to these grades, we determined three tiers corresponding to images of good, medium and bad quality. CNNs were then trained to

rate the overall image quality. Our classifier was as efficient as manual rating for the classification of images which are not proper 3D T1w brain MRIs (e.g., images of segmented tissues or truncated images). It was also able to recognise low quality images with good accuracy. While the detection of certain features such as noise did not present any particular difficulty for our model, the detection of motion artefacts proved more problematic.

As motion quantification is a complex problem, particularly due to its sensitivity to many cofactors such as contrast, there is a lack of dataset with reliable quantitative ground truths. Some studies thus rely on synthetic motion to detect motion artefacts in a controlled way (Mohebbian et al., 2021; Pawar et al., 2022; Shaw et al., 2019). Despite the excellent results claimed in the literature, only few papers have attempted to validate their performance on data with real motion. And even when they did, their test sets were extremely limited and only composed of research data (Shaw et al., 2019; Duffy et al., 2018). It is yet unclear how they would perform on routine clinical data.

In this paper, we propose a transfer learning framework from research to clinical data for the automatic detection of motion artefacts in 3D T1w brain MRI from a CDW gathering images acquired within the 39 hospitals of the Greater Paris area. We generated synthetic motion in MR images of publicly available research databases to train a CNN classifier which was validated on synthetic and real motion artefacts. Our model was then generalised and validated on a very large clinical dataset from a CDW with an effective transfer learning technique using 4045 labelled MRIs acquired in clinical routine. Preliminary work was accepted for publication in the proceedings of the SPIE Medical Imaging 2023 conference (Loizillon et al., 2023). Contributions specific to this paper include (i) a comparison of the two main synthetic motion generation approaches for the detection of motion artefacts in a CDW: the image and k-space based approaches; (ii) the implementation and comparison of four deep learning architectures (Conv5FC3, ResNet, SE-CNN and ViT) for the detection of motion artefacts; (iii) an optimisation of the fine-tuning parameters; (iv) a comparison of the proposed framework with four state-of-the-art approaches: one based on IQMs (Sadri et al., 2020) and three relying on neural networks (Duffy et al., 2018; Oksuz, 2021; Mohebbian et al., 2021).

2. Background

In this work, we focus on the detection of motion artefacts in 3D T1w MRIs in a CDW. As most of the automatic QC tools rely on neuroimaging software packages that are only usable on good quality images, they are incompatible with CDWs. What is more, manual annotation of motion artefacts is a challenging task. When an image is degraded, it may be difficult to properly distinguish motion from noise or bad contrast. Hence the idea of simulating motion, which can be done automatically and provides reliable ground truths.

Head motion can be well approximated as rigid body motion, which requires six degrees of freedom, comprising three translations and three rotations (Duffy et al., 2018). Lee et al. (2020) described the two main approaches for motion simulation in brain MRIs: the image and the k-space based techniques.

As illustrated in Fig. 1 (left), the image-based approach assumes that the subject takes Nt different positions during the acquisition. First, Nt rigid transformations of the motion-free image are applied before computing the fast Fourier transform (FFT). A new k-space is then built by concatenating blocks for the Nt different simulated positions. Finally, in order to obtain the final synthetic image corrupted by motion, an inverse FFT is applied (Shaw et al., 2019; Pérez-García et al., 2021). While in the image-based approach, motion parameters are applied on the motion-free MRI, the k-space based method directly uses the raw k-space to perform the simulation of motion (Fig. 1, right). In their algorithm, Loktyushin et al. (2013) start by applying the rotation to the k-space grid and perform a non-uniform FFT. Then, a linear phase shift proportional to the amplitude of the translation

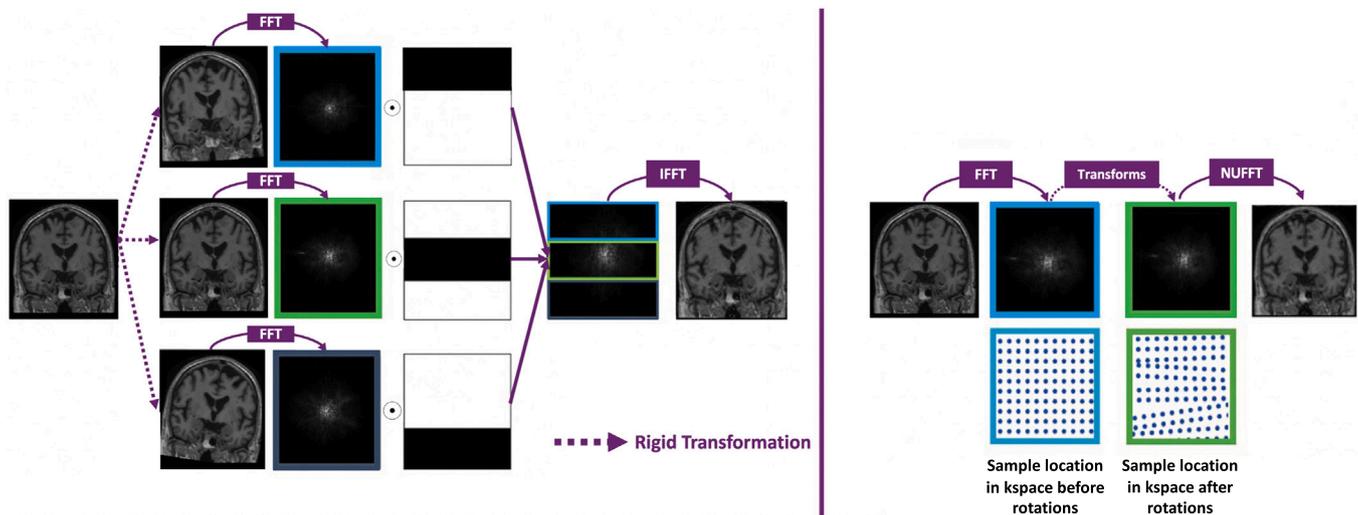


Fig. 1. Left: Image-based motion simulation. (1) N_t rigid transformations of the motion-free image are applied (here $N_t = 2$), (2) Fast Fourier transform (FFT) of the original and N_t transformed images, (3) Concatenation of the $N_t + 1$ blocks to create a new k-space, (4) Inverse FFT (IFFT) to obtain the motion corrupted image. Right: k-space based motion simulation. (1) FFT of the motion-free image. (2) Transformation (rotation + translation) for each point of the time course. (3) Non uniform FFT (NUFFT) to reconstruct the corrupted image (because of the non uniform sample spacing due to the rotation in the k-space).

is added before the final FFT. More recently, Al-masni et al. (2022) introduced a new approach by combining the image and the k-space based methods, where the translation was directly performed in the k-space domain, whereas the rotation was applied on the image. This method has the advantage of preserving the uniformity of the k-space sampling as the rotation is applied in the image domain.

Motion detection in MRI with deep learning techniques has been studied in Duffy et al. (2018), Mohebbian et al. (2021) and Oksuz (2021) using datasets of images corrupted with synthetic motion obtained from motion-free MRI. Duffy et al. (2018) proposed a modified version of the HighRes3DNet composed of 8-convolutional layers to detect and correct motion artefacts using as input MRI patches of size $80 \times 80 \times 80$ voxels. Their method was validated on real motion using images of the ABIDE dataset with promising results. Mohebbian et al. (2021) developed a stacked ensemble model to classify motion artefacts into five severity levels in brain MRIs. While their 2D model was perfectly able to predict, across different sequences (T1w and T2w), synthetic motion artefacts (balanced accuracy >90%), their approach was not validated on MRIs with real motion. Oksuz (2021) introduced a 2D dense CNN to detect motion in brain MRIs and successfully validated their binary algorithm using 28 MRIs from a research dataset (balanced accuracy: 97.8%). This method was also only validated on research MRIs corrupted with synthetic motion artefacts. Recently, Sagawa et al. (2022) presented a CNN trained using images corrupted with synthetic motion labelled with their full-reference image quality assessment (FR-IQA) metrics to predict with a high accuracy these metrics. Their approach enables a quantitative assessment of motion artefacts without any reference image. The model classified real motion artefacts from research MRIs with an area under the receiver operating characteristic (ROC) curve (AUC) of 0.928. Overall, the use of simulation for motion detection in brain MRI has only been validated on test sets with synthetic motion or on real motion with very small test sets (<40 MRIs), and always on images acquired in a research context. It still remains unclear how such methods can perform on large clinical datasets.

The requirement for accurate ground truths encourages researchers to develop solutions using synthetic data, where labels are easily available. However, the anatomical complexity and diversity of healthy and pathological brain tissues makes it difficult to generate an appropriate spectrum of synthetic MRIs, which leads to poor performance of the classifiers at the stage of inference on real data (Karani et al., 2018;

Billot et al., 2021). To benefit from the use of synthetic data, it is thus important to bridge the gap between synthetic and real images.

Transfer learning applies knowledge learned from one domain and one task to another related domain and task. In the case of motion detection using synthetic data, if we have labels for both synthetic and real images, we can resort to the use of fine-tuning (inductive transfer learning). Fine-tuning involves transferring the weights from a pre-trained network to the network to be trained. In a classification context, a common practice is to replace some of the last fully connected layers of the pre-trained CNN with new fully connected layers to target the new application. Tajbakhsh et al. (2016) demonstrated that the use of a pre-trained CNN with fine-tuning outperformed or, in the worst case, performed as well as a CNN trained from scratch for four distinct medical imaging applications. Although the distance between natural images and medical images is considerable, Ahmed et al. (2017) also showed that fine-tuning a 2D-CNN, initially trained on ImageNet, by transferring the learned feature representations to the MRI-based survival time prediction task, performed better than training from scratch. Instead of transferring knowledge from natural images to medical imaging slices, in this paper we leverage the 3D information contained in medical images and propose to transfer the knowledge learned from the detection of motion artefacts simulated from research MRIs to the detection of real motion artefacts in clinical MRIs. To our knowledge, we are the first to leverage motion simulation and research data to tackle the problem of motion detection in routine clinical data.

3. Materials and methods

We developed an approach based on the generation of synthetic motion to improve the detection of motion artefacts in clinical T1w brain MR images. We used T1w images, which were acquired with scanners from different manufacturers and different magnetic fields, from publicly available research data sources as well as from a CDW. Motion artefacts were synthetically generated by applying both image and k-space based approaches using rigid body transformations to simulate different severity degrees of artefacts. CNNs were first trained on research databases to recognise synthetic motion, and their performance was evaluated on real motion. We generalised our model to the CDW by applying an efficient transfer learning technique.

Table 1
Distribution of the sex and age over the research (ADNI, MSSEG and MNI BITE) and the clinical (AP-HP) datasets.

| | Database | N patients | N images | Age in years [range] | Sex (%F) |
|-------------------------|----------|------------|----------|------------------------|----------|
| Research databases | ADNI | 70 | 1143 | 74.31 ± 7.11 [55, 90] | 41.43% |
| | MSSEG | 53 | 53 | 45.42 ± 10.27 [24, 66] | 71.70% |
| | MNI BITE | 13 | 26 | 52.00 ± 17.70 [31, 76] | 35.71% |
| Clinical data warehouse | AP-HP | 3346 | 4045 | 55.15 ± 7.89 [18, 95] | 55.39% |

3.1. Datasets

To detect motion artefacts in routine clinical images, we first used three publicly available research datasets to pre-train a CNN on images with synthetic motion artefacts. Then, images from our CDW were exploited for transfer learning and validation.

3.1.1. Research-oriented databases

We worked with the ADNI, MSSEG and MNI BITE research-oriented databases to cover the wide spectrum of pathologies that can be found in a CDW. A special attention was paid to the search for contrast-enhanced T1w MRI as the CDW includes images acquired with and without injection of a gadolinium-based contrast agent.

The Alzheimer’s Disease Neuroimaging Initiative (ADNI) is a multi-site study of elderly individuals with normal cognition, mild cognitive impairment, or Alzheimer’s disease (Petersen et al., 2010). The ADNI-1 phase included T1w MRIs acquired on 1.5 T scanners from different manufacturers (GE, Siemens, and Philips) using a magnetisation-prepared rapid gradient echo (MPRAGE) sequence. A two-level quality control was performed, one related to the adherence to the protocol parameters and one to the series-specific quality. Part of the metadata, the IPMOTION score indicates the absence of motion (0), or the presence of mild (1), moderate (2) or severe (3) motion artefacts. This score was used to select motion-free T1w MRIs. Our selection procedure resulting in 1143 MR images for 70 subjects is detailed in supplementary material (Fig. S1). We also created a small test set with MRIs corrupted by motion artefacts based on the IPMOTION and the comments section of the corresponding metadata file.

The MSSEG MICCAI challenge, which aim is to perform the segmentation of multiple sclerosis lesions, includes 53 patients across four different sites (Commowick et al., 2018). Four different MRI scanners were used: GE Discovery 3 T, Philips Ingenia 3 T, Siemens Aera 1.5 T and Siemens Verio 3 T. Each scan included four MRI sequences: 3D FLAIR, 3D T1w, 3D contrast-enhanced T1w and 2D T2w. In our study, we only considered the 3D contrast-enhanced T1w.

The Montreal Neurological Institute’s Brain Images of Tumours (MNI BITE) database made available pre and postoperative MR and ultrasound images acquired from brain tumour patients (Mercier et al., 2012). The study includes 13 patients with gliomas, who underwent a pre- and post-operative contrast-enhanced T1w MRI with the 1.5 T GE Signa EXCITE scanner using a 3D axial spoiled gradient recalled acquisition (SPGR).

If we relied on the IPMOTION score to extract MRIs without motion artefacts in the ADNI dataset, for the MSSEG and MNI BITE datasets, we conducted a manual inspection of each image. Through this manual inspection, we identified six problematic MRIs from the MNI BITE dataset and one from the MSSEG dataset that exhibited motion artefacts. Consequently, these problematic images were removed from the training set. Demographic information for each of these databases is reported in Table 1. In the following, we refer to the images coming from these three databases as the ‘research dataset’.

Acquisition parameters including repetition and echo times are presented in supplementary material for the three research-oriented datasets (Table S1).

3.1.2. Clinical data warehouse

The clinical routine data come from a large CDW containing all the T1w brain MRIs of adult patients scanned in hospitals of the Greater Paris area (Assistance Publique-Hôpitaux de Paris [AP-HP]). The large number of hospitals being part of the AP-HP consortium (39 hospitals) and the huge number of images collected every day is making this CDW a good representation of 3D T1w brain MRIs that may be acquired in other hospitals.

We used the same dataset as in our previous study, where we randomly selected 5500 images, corresponding to 4177 patients that were acquired on various scanners: Siemens Healthineers (n = 3752, 13 different scanner models), GE Healthcare (n = 1710, 12 different scanner models), Philips (n = 33, 3 different scanner models) and Toshiba (n = 5, 2 different scanner models) (Bottani et al., 2021).

Motion artefacts were manually annotated on a three-level scale by two trained annotators following an annotation protocol developed with the help of a radiologist. A score of 0 was given when no motion was seen, 1 when the structures of the brain were distinguishable despite the presence of motion and 2 when the cortical and sub-cortical structures were difficult to distinguish (Fig. S2). Some of the 5500 images did not correspond to 3D T1w brain MRIs (e.g., because of truncation) and were therefore not labelled with a motion score (SR: straight reject, n = 1455). If the annotators labelled differently a given MRI, the consensus grade was chosen as the maximum of the two grades. The weighted Cohen’s kappa was used to evaluate the inter-rater agreement between the annotators and a moderate agreement was found with a score of 0.68. Among the 4045 images that were not labelled as straight reject, 2319 had a consensus motion score of 0, 1196 a score of 1 and 530 a score of 2. Patients’ demographics are reported in Table 1.

3.2. Image pre-processing

To make the annotation process easier, MRIs were pre-processed using Clinica (Routier et al., 2021) and its `t1-linear` pipeline. First a bias field correction was applied using the N4ITK method (Tustison et al., 2010). An affine registration to the MNI space was then performed (Avants et al., 2008). Next, images were cropped to remove background resulting in images of size $169 \times 208 \times 179$, with 1 mm isotropic voxels (Wen et al., 2020). The Z-score method, which consists of subtracting the mean intensity of the entire image from each voxel value and dividing it by the corresponding standard deviation, was used to normalise the voxel intensities. Our initial aim was to obtain a rough alignment and intensity rescaling to facilitate annotation but this pre-processing is also useful when training CNNs.

3.3. Proposed approach

We developed a transfer learning approach to detect motion artefacts in clinical images based on motion simulated on research images. Our method is composed of two steps: (1) a pre-training task using synthetic motion to distinguish motion-free from motion-corrupted images, (2) a fine-tuning task to improve the generalisation ability of our pre-trained network on clinical datasets.

3.3.1. Motion generation

Because of the lack of research dataset with quantitative assessment of motion artefacts in T1w brain MRI, we adopted an approach based on synthetic motion. We compared the two main motion simulation techniques described in Section 2: the image and the k-space based approaches. We used the open-source Python library TorchIO and its function `RandomMotion` described in Pérez-García et al. (2021) for the image based approach and the `RandomMotionTimeCourseAffines` implemented in Reguig et al. (2022) for the k-space method. The `RandomMotion` function was used with a number of rigid transformations ($Nt = 4$), whereas the `RandomMotionTimeCourseAffines` was applied using 200 points of the simulated time course ($nT = 200$).

By selecting different translation and rotation range parameters, several degrees of motion severity can be generated. Different values have been tested in this study to simulate motion ranges of [2 mm, 8 mm] for translation and [2°, 8°] for rotation.

3.3.2. Network architectures

To classify motion artefacts, we used a 3D CNN composed of five convolutional blocks and of three fully connected layers (denoted as Conv5FC3) that proved successful in our previous work (Bottani et al., 2021). Each convolutional block is made of a convolutional layer, a batch normalisation layer, a ReLU activation function and a max pooling layer. The weighted binary cross-entropy was used as loss function. The learning rate of the Adam optimiser was set to $1e-4$ and the batch size to 16. The model with the lowest loss on the validation set was saved as final model. The architecture was implemented using Pytorch and is available through the ClinicaDL software on GitHub (<https://github.com/aramis-lab/clinicaDL>) (Thibeau-Sutre et al., 2022).

We compared this 3D network to more sophisticated architectures such as a 3D ResNet, a 3D Squeeze and Excitation CNN (SE-CNN) and a 3D Vision Transformer (ViT). The 3D ResNet inspired by Jonsson et al. (2019) was previously used to predict brain age from 3D T1w MRI and outperformed the Conv5FC3 (Couvry-Duchesne et al., 2020). The combination of a ResNet with Squeeze and Excitation blocks was successfully tested on brain tumour classification (Ghosal et al., 2019). SE blocks are composed of a squeeze and an excitation step. The squeeze operation is obtained through an average pooling layer and provides a global understanding of each channel. The excitation part consists of a two-layer feed-forward network that outputs a vector of n values corresponding to the weights of each channel of the feature maps. Whereas traditional CNNs weight each of their channels equally when creating feature maps, SE-CNNs weight each channel adaptively through this content-aware mechanism. Transformers, which have become the model of choice in natural language processing, have recently been applied to computer vision tasks. Even if applications in medical imaging remain limited, vision transformers have been used to perform classification (Alzheimer’s disease detection) as well as segmentation tasks (brain tumour segmentation) (Wang et al., 2021; Xing et al., 2022). The different architectures are detailed in supplementary material (Fig. S3 and S4). The same loss and hyperparameters as for the Conv5FC3 were used for the training of these networks.

3.3.3. Model generalisation using transfer learning

After pre-training a classifier on a research dataset by simulating motion and having obtained a satisfactory model capable of accurately detecting motion artefacts in research quality images, we can apply a transfer learning method based on fine tuning to generalise to routine clinical data. This approach enables closing two gaps at the same time, one between research and clinical datasets and one between real and synthetic motion.

The key idea of fine-tuning is to transfer knowledge learnt from one domain to another one. A classifier was first trained to learn features for the research dataset domain using motion simulation. The network was then optimised again for a new domain (clinical dataset with real motion) by re-training several layers of the pre-training model and freezing the weights of the other layers. Thus, we were able to generalise our model from synthetic to routine clinical data.

Table 2

Distribution of the training, validation and test sets separately for the synthetic motion detection task using the research dataset comprising images from the ADNI, MSSEG and MNI BITE databases.

| | Label | Research dataset | | |
|------------|-----------|------------------|-------|----------|
| | | ADNI | MSSEG | MNI BITE |
| Training | Motion | 400 | 19 | 9 |
| | No motion | 400 | 19 | 9 |
| Validation | Motion | 86 | 4 | 2 |
| | No motion | 86 | 4 | 2 |
| Testing | Motion | 86 | 4 | 2 |
| | No motion | 86 | 4 | 2 |

3.4. Experiments

We performed two sets of experiments that correspond to the two steps of the proposed approach. The first set focuses on the network pre-training step with research data using synthetic motion and the second set concerns the network fine-tuning step with clinical data and real motion.

3.4.1. Network pre-training on research data

At first, we aimed to test the ability of the different deep learning models to detect motion in research-quality images using only synthetic motion while training. We performed a series of experiments on the research dataset, where we corrupted motion-free MRIs with different motion severity degrees to study the influence of the translation and rotation ranges. The four different architectures as well as the two motion simulation techniques presented in Sections 3.3.2 and 3.3.1 were tested to determine the best approach for motion detection.

Before starting the experiments, we defined an independent test set by randomly selecting 184 images over the three research-oriented datasets and corrupting half of them with different motion severity degrees (rotation: [2°, 8°]; translation: [2 mm, 8 mm]). The remaining 1040 images (520 corrupted with synthetic motion and 520 with no motion) were split into training and validation using a 5-fold cross validation (CV) as shown in Table 2. The separation between training, validation and test sets was made at the subject level to avoid data leakage. Our model was also validated on a second small test set with ADNI MRIs corrupted by real motion as explained in Section 3.1.1.

3.4.2. Network fine-tuning on routine clinical data

The second set of experiments aims to evaluate the performance of our transfer learning approach. This step is required because of the quality gap that exists between research, where strict acquisition protocols are respected, and clinical data, which suffer from a lack of homogenisation of the acquisition parameters.

To generalise our pre-trained network on clinical datasets, we used a transfer learning technique that consists in fine-tuning our model on two distinct target tasks:

- the detection of severe motion (Mov01vs2): being able to detect these MRIs in CDWs is of great importance as subsequent processing steps are likely to fail on these images,
- the detection of moderate motion (Mov0vs1): MRIs labelled as motion 1 may lead to unreliable diagnostic predictions.

Before starting the experiments, we defined a test set by selecting the same MRIs as in Bottani et al. (2021) as well as the same training and validation splits with a 5-fold CV (Table 3). Because of the presence of straight reject MRIs, the different models trained in the CV were evaluated on respectively 328 and 385 images of the test set for the Mov0vs1 and Mov01vs2 tasks.

To evaluate the impact of the proposed method, we compared the results obtained with the use of fine-tuning and by training a model from scratch where we randomly initialise the neural network

Table 3

Distribution of the training validation and test sets separately for the two fine-tuning tasks on the AP-HP CDW: the detection of severe (Mov01vs2) and moderate (Mov0vs1) motion.

| | Label | CDW | |
|------------|----------|---------|----------|
| | | Mov0vs1 | Mov01vs2 |
| Training | Motion 0 | 1681 | 1681 |
| | Motion 1 | 859 | 859 |
| | Motion 2 | – | 379 |
| Validation | Motion 0 | 428 | 428 |
| | Motion 1 | 219 | 219 |
| | Motion 2 | – | 94 |
| Testing | Motion 0 | 210 | 210 |
| | Motion 1 | 118 | 118 |
| | Motion 2 | – | 57 |

parameters. We also studied the influence of the number of layers to freeze for the different architectures, as well as the impact of the number of transforms when applying the RandomMotion function.

3.4.3. Comparison with state-of-the-art QC

Several quality check tools for T1w brain MRI based on IQMs have been developed in the last years (Esteban et al., 2017; Sadri et al., 2020). Whereas MRIQC (Esteban et al., 2017) cannot be applied on CDWs because of the need for extensive image pre-processing designed for T1w brain MRI of good quality without gadolinium, MRIQY (Sadri et al., 2020) is compatible with CDW MRIs thanks to its automatic extraction and separation of the background from the foreground with an Otsu thresholding. Thus, we were able to extract 13 IQMs such as noise ratios, variation metrics, entropy, and energy criteria from our clinical images. Thanks to these IQMs, we trained a random forest (RF) classifier to detect the presence of motion artefacts in MRI. We performed a random search in order to optimise the hyper-parameters and particularly analysed: the number of decision trees, the maximum tree-depth, the minimum number of samples per split, and the minimum leaf samples. 300 different combinations were tested using a 5-fold CV.

We also compared our method with three state-of-the-art deep learning-based motion artefact detection methods (Duffy et al., 2018; Oksuz, 2021; Mohebbian et al., 2021). In order to replicate their methodology, we introduced k-space based motion artefacts into good-quality MRIs from the CDW, following the approach employed in these studies. Subsequently, we trained three different architectures: a modified 2D DenseNet (Oksuz, 2021), a 3D patches HighRes3DNet (Oksuz, 2021), and a 2D CNN architecture consisting of three convolutional layers and two dense layers (Mohebbian et al., 2021).

4. Results

4.1. Validation on research data

The ability of deep learning models to detect motion was first assessed using images from the research dataset corrupted with synthetic motion. Fig. 2 displays three corrupted images obtained with the image and the k-space approach using different translation and rotation ranges, and the original image without any motion.

We started by evaluating the performance of our Conv5FC3 model trained on synthetic motion when applied to our synthetic independent test set corrupted with different motion severity degrees (rotation: [2°, 8°]; translation: [2 mm, 8 mm]). We studied the influence of the translation and rotation ranges by performing several experiments with different motion severity degrees. We first trained a model with synthetic severe motion by applying a large rotation and translation range ([6°, 8°]; [6 mm, 8 mm]). The balanced accuracy (BA) on our independent test set is excellent with both motion simulation techniques (>98%). We also obtained very good results for smaller ranges of rotation and translation as reported in Table 4. A ROC analysis was

performed for these experiments on the research dataset and the AUC values consistently exceeded 0.98, indicating an excellent ability to distinguish between true positives and false positives, regardless of the rotation and translation parameters used (Fig. S5).

Then, we evaluated the ability of these models to detect real motion. As mentioned in Section 3.1.1, we defined a test set according to the IPMOTION score. Our models were perfectly able to detect motion on these images. No notable differences were noted in terms of performance between the two simulation techniques (Table 4).

We also compared the performance obtained by different architectures on the test set corrupted with synthetic motion. In Table 5, we report the results of the four architectures trained using k-space based motion simulation with the following parameters: rotation: [2°, 4°]; translation: [2 mm, 4 mm], as these led to the best results on synthetic and real motion for the Conv5FC3 architecture. The results of the ResNet and the SE-CNN were comparable to that of the Conv5FC3 with a BA>99%, whereas the ViT BA was lower (BA = 97.69%). Thus, more complex networks did not provide any notable improvement. The same conclusion was reached for the image-based motion simulation technique (Table S2).

4.2. Application to routine clinical data

The first set of experiments performed with the routine clinical data consisted in fine-tuning the Conv5FC3 network pre-trained on the research dataset with synthetic image-based motion to detect severe motion (mov01vs2) by unfreezing one to five layers of the Conv5FC3 architecture. We used the same training and validation split as for the mov01vs2 task (Table 3). The different models trained in the CV were evaluated on the 385 images composing the test set for the mov01vs2 task. Best results were obtained by freezing all the layers except the three fully connected ones of the model pre-trained on the research dataset with synthetic motion using four rigid transforms ($nT = 4$) (Table S3 and Table S4). We also conducted experiments involving the fine-tuning of the ResNet, SE-CNN, and ViT architectures (Table S5 and Table S6). For ResNet and SE-CNN, we froze all layers except the three fully connected ones, as we did for the Conv5FC3. As ViT architectures lack fully connected layers, we retrained all layers of this architecture. It is worth noting that ResNet demonstrated the best performance for the detection of severe motion, achieving a BA of 86.36%, while Conv5FC3 outperformed the other architectures for mild motion detection, with a BA of 62.61%. As there is no clear advantage in using a more complex architecture, the fine-tuning results presented in the following are those obtained by training the Conv5FC3, with only the three fully connected layers being retrained.

The results obtained with the proposed transfer learning framework on our independent clinical test set are presented in Table 6. For the detection of severe motion (mov01vs2), the classifier BA is almost as good as that of the annotators, which is defined as the average of the BA between each rater and the consensus (classifier: 84.52%; annotators: 86.29%). For the detection of mild motion (mov0vs1) the classifier BA is low (62.61%) and lower than that of the raters (73.21%).

We compared the results obtained with and without fine-tuning to measure the impact of our approach. When applying the network trained on the synthetic research data directly to the clinical data, we observed a large drop in performance with a particularly low specificity for both tasks. A second comparison was performed between the proposed transfer learning framework and when training with the clinical data from scratch. Our transfer learning method achieved a gain of more than 10 percent points for the detection of severe motion. A much smaller improvement was observed for the detection of mild motion in terms of BA, but specificity and sensitivity became more balanced. The ROC curves for the detection of severe (mov01vs2) and moderate (mov0vs1) motion in 3D T1w MRIs are shown in Fig. 3. The AUC of the proposed approach for severe motion detection (AUC: 0.85) outperformed that of training from scratch by 5 percent points and that

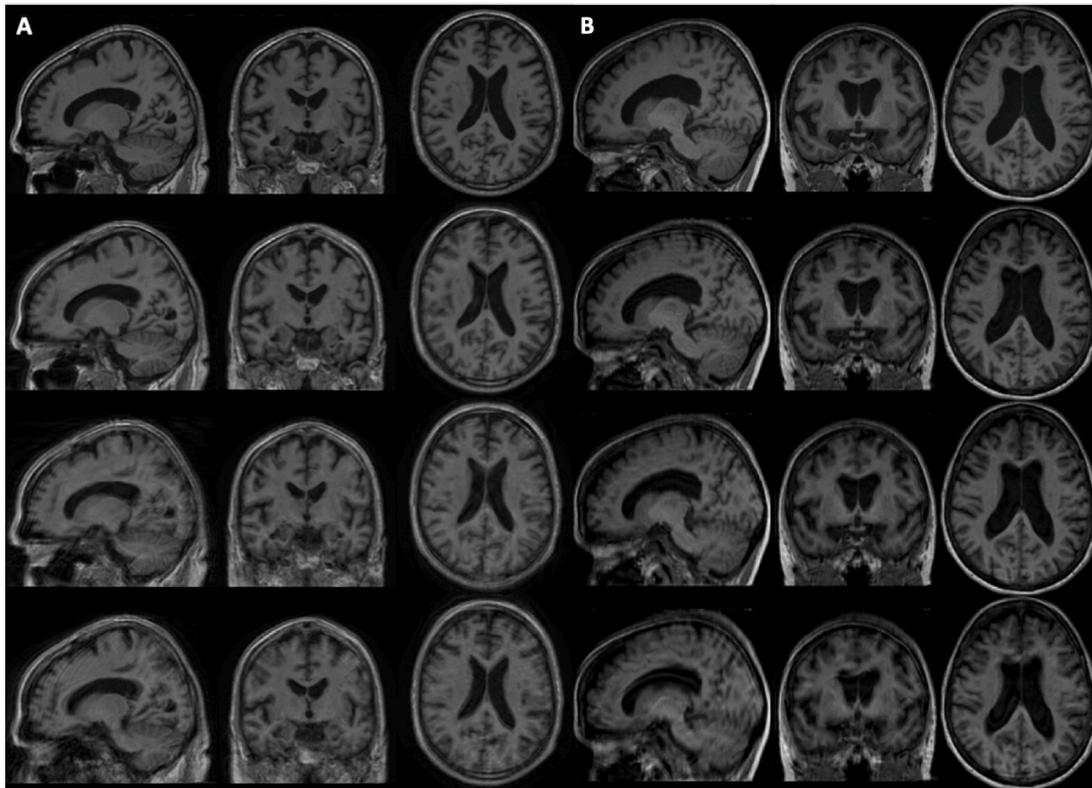


Fig. 2. Example of motion simulation on brain MRI using the k-space (A) and the image (B) based approach with different translation and rotation ranges. From top to bottom: motion-free MRI, MRIs corrupted with a rotation of 3°, 5° and 7° and a translation of 3 mm, 5 mm and 7 mm.

Table 4

Results for the detection of synthetic and real motion in T1w brain MRIs from the research dataset. For the validation on synthetic motion, we report the mean and the empirical standard deviation across the five folds for the balanced accuracy, specificity and sensitivity. For the detection of real motion, only the accuracy obtained by the best model of the 5-fold CV was reported as our independent test set contained only images with motion. Results are detailed for both simulation approaches: image and k-space based.

| Motion simulation | Rotation range | Translation range | Cross-validation on synthetic motion | | | Test on real motion |
|-------------------|----------------|-------------------|--------------------------------------|--------------|--------------|---------------------|
| | | | Balanced accuracy | Specificity | Sensitivity | Accuracy |
| Image | [6°, 8°] | [6 mm, 8 mm] | 98.22 ± 1.39 | 99.29 ± 0.86 | 97.14 ± 2.33 | 100% |
| | [4°, 6°] | [4 mm, 6 mm] | 97.06 ± 1.47 | 98.25 ± 1.92 | 95.87 ± 1.27 | 100% |
| | [2°, 4°] | [2 mm, 4 mm] | 95.51 ± 2.47 | 98.94 ± 2.11 | 92.06 ± 5.76 | 98.41% |
| k-space | [6°, 8°] | [6 mm, 8 mm] | 98.44 ± 0.05 | 97.22 ± 0.12 | 99.70 ± 0.01 | 100% |
| | [4°, 6°] | [4 mm, 6 mm] | 97.77 ± 0.03 | 95.56 ± 0.06 | 100 ± 0.00 | 100% |
| | [2°, 4°] | [2 mm, 4 mm] | 99.17 ± 0.03 | 98.33 ± 0.06 | 100 ± 0.00 | 100% |

Table 5

Results of four different CNN classifiers (Conv5 FC3, ResNet, SE-CNN and ViT) trained and tested on MRIs from the research dataset corrupted with k-space based motion simulation.

| Architectures | Balanced accuracy | Specificity | Sensitivity | Training time |
|---------------|-------------------|--------------|--------------|---------------|
| Conv5FC3 | 99.17 ± 0.03 | 98.33 ± 0.06 | 100 ± 0.00 | 3 h 52 min |
| ResNet | 99.72 ± 0.03 | 99.44 ± 0.07 | 100 ± 0.00 | 6 h 27 min |
| SE-CNN | 100 ± 0.00 | 100 ± 0.00 | 100 ± 0.00 | 6 h 18 min |
| ViT | 97.69 ± 0.07 | 98.61 ± 0.03 | 96.77 ± 0.04 | 4 h 31 min |

of training on research datasets by 24 percent points. For moderate motion detection, the AUC of the proposed approach (AUC: 0.63) is also 5 percent points higher than that of learning from scratch (AUC: 0.58).

The same set of experiments was performed using the model pre-trained on the research dataset with k-space based synthetic motion (Table S3). As with the model pre-trained on image-based motion, our transfer learning method led to a substantial performance improvement when detecting severe motion, with a gain of 6.85 percent points in BA compared with training from scratch on clinical data, but the

improvement was limited when detecting mild motion (gain of 1.7 percent point in BA).

4.3. Comparison with state-of-the-art QC

We conducted a comparative study between our proposed approach and four state-of-the-art motion artefact detection methods (Sadri et al., 2020; Duffy et al., 2018; Mohebbian et al., 2021; Oksuz, 2021).

We first compared our deep learning method with an IQM based approach (RF classifier trained on IQMs) (Sadri et al., 2020), which achieved a 61.94% BA for the detection of severe motion and 54.72% for mild motion detection (Table 7). Our proposed method outperformed the IQM approach by 8 percent points for the mild motion detection and by 23 percent points for the severe motion.

We then compared our approach with three deep learning-based methods (Duffy et al., 2018; Mohebbian et al., 2021; Oksuz, 2021). As shown in Table 7, state-of-the-art methods struggle to effectively identify real motion artefacts in 3D T1w brain routine clinical MRIs (BA < 57%). This highlights the critical importance of fine-tuning to address the gap between synthetic and real motion artefacts.

Table 6

Detection of motion artefacts within brain T1w MR images of the CDW. For both the detection of severe motion (Mov01vs2) and mild motion (Mov0vs1), we report: the agreement between human raters and the consensus (manual annotations), results of the proposed approach (pre-training on image-based synthetic motion from research data and fine-tuning on CDW), results when training on image-based synthetic motion from research datasets without fine-tuning, and results when training from scratch on CDW.

| | | BA | Specificity | Sensitivity |
|----------|--|--------|-------------|-------------|
| Mov01vs2 | Manual annotations | 86.29% | – | – |
| | Conv5FC3 fine-tuned on CDW (proposed) | 84.52% | 85.37% | 83.67% |
| | Conv5FC3 trained on research dataset | 60.26% | 33.33% | 87.19% |
| | Conv5FC3 trained from scratch on CDW | 73.75% | 49.58% | 97.93% |
| Mov0vs1 | Manual annotations | 73.21% | – | – |
| | Conv5FC3 fine-tuned on CDW (proposed) | 62.61% | 52.00% | 73.23% |
| | Conv5FC3 trained on research dataset | 53.18% | 17.96% | 88.57% |
| | Conv5FC3 trained from scratch on CDW | 58.93% | 28.81% | 89.05% |

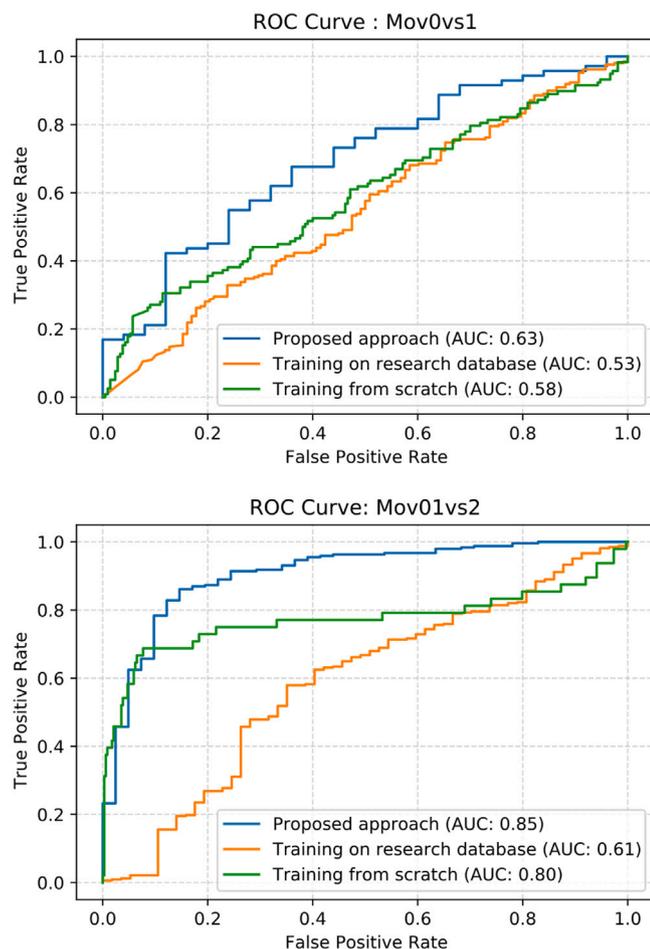


Fig. 3. Receiver operating characteristic curves (ROC) for the detection of severe (mov01vs2) and moderate (mov0vs1) motion in 3D T1w MRIs.

5. Discussion

In this work, we developed an innovative transfer learning framework from research to clinical data for the automatic detection of motion in 3D T1w brain MRI from a clinical data warehouse. After a pre-training phase using synthetic motion to distinguish images with and without motion artefacts, we enabled the generalisation of our pre-trained network on clinical data using an effective fine-tuning technique. We validated our approach on a unique set of manually-labelled MRIs acquired in clinical routine within the network of the 39 hospitals of the AP-HP that are gathered in a warehouse. To the best of our knowledge, we are the first to propose a very large-scale validation using a CDW for motion artefact detection using synthetic motion.

Table 7

Comparative study between our proposed method and state-of-the-art (SOTA) motion artefact detection methods (for the task Mov01vs2).

| Methods | BA | Specificity | Sensitivity |
|--|-------|-------------|-------------|
| Proposed | 84.52 | 85.37 | 83.67 |
| Duffy et al. (2018) | 55.20 | 14.58 | 95.82 |
| Oksuz (2021) | 52.70 | 6.94 | 98.46 |
| Mohebbian et al. (2021) | 56.42 | 16.03 | 96.05 |
| Sadri et al. (2020) (trained on research dataset) | 62.13 | 85.67 | 38.60 |
| Sadri et al. (2020) (trained on CDW) | 61.94 | 97.56 | 26.32 |

For the detection of synthetic motion on research datasets, our model achieved excellent results with a BA over 98% for the image and the k-space based simulation approaches. Trained using only synthetic motion, the model had no difficulty generalising to real motion artefacts and was able to detect every image corrupted with motion on a small independent research dataset. Despite the performance obtained on research datasets, the model was not able to generalise to the CDW (BA: 60.26%). This poor result, with a low specificity, does not come as a surprise. It highlights the critical importance of validating models trained on research datasets to clinical ones, but also the quality gap that exists between research, where strict acquisition protocols are respected, and clinical data, which suffer from a lack of homogenisation of acquisition parameters. To overcome these issues, we proposed a transfer learning framework that achieved very good results for the detection of severe motion with a BA of 84.52%, which is nearly as good as that of the annotators (86.29%) and 10 percent points higher than when training the model from scratch on clinical data (BA: 73.75%). The performance is significantly different from that of a network trained from scratch ($p = 6.60 \times 10^{-17} < 0.05$, McNemar's test). However, the result for the more difficult task of mild motion detection was less successful with a BA of 62.61% compared to the 73.21% obtained by the annotators. Our two-step approach with pre-training and fine-tuning still allowed an increase of almost 5 percent points of BA over training from scratch on this task and the difference between the two approaches was significant ($p = 1.39 \times 10^{-8} < 0.05$, McNemar's test).

For practical application, our primary focus is to exclude MRIs with severe motion artefacts (mov2), since they can lead to failures in the preprocessing pipeline. However, it is important to acknowledge that identifying mild artefacts can still be of interest, as they may still be a source of bias in morphometric analyses (Reuter et al., 2015; Hedges et al., 2022). Regarding the detection of mild motion (mov₁), our model achieved a sensitivity of 73%. We acknowledge that this sensitivity value remains relatively low, and that we should seek ways to improve the detection of mild artefacts.

The detection of mild motion remains challenging even for manual annotators as highlighted by the moderate graders BA (73.21%). We ensured the reliability of the 5500 manual annotations by calculating the inter-rater agreement for motion by computing the weighted Cohen's kappa coefficient between the two annotators. The agreement was

found to be moderate, with a value of 0.68 for the motion characteristics. We feel confident that the manual annotations are reliable enough for the application targeted.

In the scope of our work, we have been comparing the two main approaches of motion simulation: the image and the k-space based methods. While the image-based motion simulation is sufficient for our application – detecting motion artefacts within a clinical data warehouse – it is crucial to acknowledge its limitations. Firstly, unlike the k-space based implementation, the image based TorchIO transforms do not account for the actual readout direction of the image. Additionally, the `RandomMotion` function lacks the flexibility to manipulate the order of the k-space filling, which is known to influence the direction of motion artefacts. Another limitation comes from the motion model, which is randomly sampled from probability density functions, whereas the patient's motion during acquisition has some redundancy, for example due to the swallowing mechanism, making it a process that is not entirely random. In contrast, the k-space based approach using motion time courses simulates more realistic motion by considering slow, sudden and swallowing motion with respectively a Perlin noise, a step displacement and a transient motion. Despite these considerations, our proposed transfer learning framework obtained the best performance with the image-based approach. As we only used motion simulation to pre-train our model, it appears that we can limit ourselves to a very simplified motion simulation with the image-based technique considering four positions. What is more, the latter allows corrupting in a simplified way an MRI in 3 s where the k-space based approach, which generates a much more complex motion, takes 20 s.

One of the strong assumptions made when simulating motion with the image-based approach is that the subject takes nT different positions during the acquisition. We evaluated the optimal number of positions to simulate based on a comparative study where we tested our proposed method with different nT values (2, 4, 6 and 8). A maximum of eight rigid transformations was considered for computational reasons. Based on the performance obtained in these experiments, we fixed $nT = 4$. One limitation of this approach is that it did not consider the echo train length of the acquisition, i.e., the number of echoes acquired in a given TR interval, which is strongly linked to the number of transforms (nT), as we can assume that minimum motion occurs within acquisitions.

We also compared the proposed 3D architecture, Conv5FC3, with the SE-CNN, ResNet and ViT architectures. The BA obtained with the first four networks on the detection of synthetic motion in the research dataset is comparable: the SE-CNN performance (100.00 ± 0.00) was slightly higher than that of the ResNet (99.72 ± 0.03), the Conv5FC3 (99.17 ± 0.03) and the ViT (97.69 ± 0.07) architectures. The performance of the different classifiers was not statistically significantly different ($p > 0.05$, McNemar test). The same comparison of performance was made between the four fine-tuned models for the detection of severe and mild motion in the CDW MRIs. The ResNet architecture achieved the best results for the detection of severe motion with a BA of 86.36%. On the other hand, the Conv5FC3 model performed best in detecting mild motion, with a BA of 62.61%. Due to the absence of a clear advantage in utilising a more complex architecture, we considered the training time of the four networks to select the best model for our transfer learning framework. Consequently, we selected the Conv5FC3, which exhibited satisfying performance while requiring less than 4 h of training on our research databases for the pre-training task (Table 5).

Finally, we compared our proposed fine-tuning framework with four state-of-the-art approaches. Our method outperformed the one proposed by MRQy (Sadri et al., 2020) that consists in training a random forest classifier with IQMs. The proposed approach reached a BA 23 percent points higher for the detection of severe motion detection. This result highlights the limitations of MRQY that computes IQMs between the head and the background thanks to an Otsu thresholding. IQMs based on noise measurements such as the coefficient of joint variation or the contrast-to-noise ratio are much more relevant when computed

between brain tissues to evaluate how separated their distributions are and thus conclude on the overall image quality. This is why most of the IQMs in MRIQC (Esteban et al., 2017) are evaluated between grey and white matter thanks to a substantial pre-processing requiring good quality data that is incompatible with a CDW. Our approach exhibited superior performance compared to the methods proposed by Duffy et al. (2018), Oksuz (2021), and Mohebbian et al. (2021), achieving an improvement of approximately 30 percent points. These outcomes point out the limitations of current state-of-the-art techniques, which often rely on validation restricted to synthetic data or small research test sets. By bridging the gaps between research and clinical data and between synthetic and real artefacts, our method effectively detects motion artefacts in a CDW.

One of the limitations of our study is the pre-processing needed before applying our model, which might prevent its direct application to new data. Our motion detection model was trained using T1w MRIs pre-processed with the `t1-linear` pipeline of Clinica (Routier et al., 2021). This pre-processing step includes an affine registration to the MNI space that facilitates the manual annotation for the graders. Such spatial normalisation could also be beneficial when training neural networks as it reduces the variability between the images. Another limitation is the annotation process of the CDW images. As the IT environment is extremely limited and data cannot be downloaded locally, the annotations performed in our previous study had to rely on only three slices (a central slice in each plane). Thus the annotators may have missed some artefacts (Bottani et al., 2021). Finally, it would be interesting to study in the future the potential generalisation of such QC models to other MRI sequences available in CDWs, such as FLAIR.

6. Conclusion

In this study, we proposed a transfer learning framework from research to clinical data for the automatic detection of motion artefacts of 3D T1w brain MRI which was validated on a large clinical data warehouse. We trained and validated different CNNs on a research dataset comprising images from three publicly available databases using motion simulation and we successfully tested them on an independent test set with synthetic and real motion. We were able to generalise our pre-trained model to clinical images thanks to the motion labelling of 4045 MRIs. Our deep learning classifier was almost as reliable as manual rating for the detection of severe motion artefacts. Our work demonstrated the usefulness of synthetic motion to improve the detection of motion artefacts in MRI, as well as the crucial need of transfer learning to generalise models trained on research to routine clinical data.

CRedit authorship contribution statement

Sophie Loizillon: Conceptualization, Data curation, Formal analysis, Methodology, Software, Visualization, Writing – original draft. **Simona Bottani:** Data curation, Software, Writing – review & editing. **Aurélien Maire:** Data curation, Writing – review & editing. **Sebastian Ströer:** Data curation, Writing – review & editing. **Didier Dormont:** Conceptualization, Funding acquisition, Project administration, Supervision, Writing – review & editing. **Olivier Colliot:** Conceptualization, Funding acquisition, Methodology, Project administration, Supervision, Writing – review & editing. **Ninon Burgos:** Conceptualization, Data curation, Methodology, Project administration, Supervision, Writing – review & editing.

Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Competing financial interests related to the present article: none to disclose for all authors. Competing financial interests unrelated to the present article: OC reports having received consulting fees from AskBio and Therapanacea and having received fees for writing a lay audience short paper from Expression Santé. O.C. holds a patent registered at

the International Bureau of the World Intellectual Property Organization (PCT/IB2016/0526993, Schiratti J-B, Allassonniere S, Colliot O, Durrleman S, A method for determining the temporal progression of a biological phenomenon and associated methods and devices). OC is in the editorial board of Medical Image Analysis.

Data availability

The authors do not have permission to share data.

Acknowledgements

The research was done using the Clinical Data Warehouse of the Greater Paris University Hospitals. The authors are grateful to the members of the AP-HP WIND and URC teams, and in particular Yannick Jacob, Julien Dubiel, Antoine Rozès, Cyrina Saussol, Rafael Gozlan, Stéphane Bréant, Florence Tubach, Jacques Ropers, Christel Daniel, and Martin Hilka. They would also like to thank the “Collégiale de Radiologie of AP-HP” as well as, more generally, all the radiology departments from AP-HP hospitals. The authors would also like to thank Romain Valabregue and Ghiles Reguig for their help implementing the k-space motion simulation and their feedback.

The research leading to these results has received funding from the Abeona Foundation, France (project Brain@Scale), from the French government under management of Agence Nationale de la Recherche, France as part of the “Investissements d’avenir” program, reference ANR-19-P3IA-0001 (PRAIRIE 3IA Institute) and reference ANR-10-IAIHU-06 (Agence Nationale de la Recherche-10-IA Institut Hospitalo-Universitaire-6).

Data collection and sharing for this project was funded by the Alzheimer’s Disease Neuroimaging Initiative (ADNI), USA (National Institutes of Health Grant U01 AG024904) and DOD ADNI, France (Department of Defense award number W81XWH-12-2-0012). ADNI is funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: AbbVie, Alzheimer’s Association; Alzheimer’s Drug Discovery Foundation; Araclon Biotech; BioClinica, Inc.; Biogen; Bristol-Myers Squibb Company; CereSpir, Inc.; Cogstate; Eisai Inc.; Elan Pharmaceuticals, Inc.; Eli Lilly and Company; EuroImmun; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; Fujirebio; GE Healthcare; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.; Lumosity; Lundbeck; Merck & Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Neurotrack Technologies; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging; Servier; Takeda Pharmaceutical Company; and Transition Therapeutics. The Canadian Institutes of Health Research is providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health (www.fnih.org). The grantee organisation is the Northern California Institute for Research and Education, and the study is coordinated by the Alzheimer’s Therapeutic Research Institute at the University of Southern California. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of Southern California.

Appendix A. APPRIMAGE study group

Olivier Colliot, Ninon Burgos, Simona Bottani, Sophie Loizillon¹
Didier Dormont^{1,2}, Stéphane Lehericy^{2,21,22}, Samia Si Smail Belkacem,
Sebastian Ströer²
Nathalie Boddaert³
Farida Benoudiba, Ghaida Nasser, Claire Ancelet, Laurent Spelle⁴
Aurélien Maire, Stéphane Bréant, Christel Daniel, Martin Hilka, Yan-
nick Jacob, Julien Dubiel, Cyrina Saussol, Rafael Gozlan¹⁹
Florence Tubach, Jacques Ropers, Antoine Rozès, Camille Nevoret²⁰

Hubert Ducou-Le-Pointe⁵, Catherine Adamsbaum⁶, Marianne Alison⁷,
Emmanuel Houdart⁸, Robert Carlier^{9,17}, Myriam Edjlali⁹, Betty
Marro^{10,11}, Lionel Arrive¹⁰, Alain Luciani¹², Antoine Khalil¹³, Elisabeth
Dion¹⁴, Laurence Rocher¹⁵, Pierre-Yves Brillet¹⁶, Paul Legmann, Jean-
Luc Drape¹⁸

¹ Sorbonne Université, Institut du Cerveau - Paris Brain Institute, Inserm, CNRS, AP-HP, Hôpital de la Pitié Salpêtrière, Inria, Aramis project-team, F-75013, Paris, France

² AP-HP, Hôpital de la Pitié Salpêtrière, Department of Neuroradiology, F-75013, Paris, France

³ AP-HP, Hôpital Necker, Department of Radiology, F-75015, Paris, France

⁴ AP-HP, Hôpital Bicêtre, Department of Radiology, F-94270, Le Kremlin-Bicêtre, France

⁵ AP-HP, Hôpital Armand-Trousseau, Department of Radiology, F-75012, Paris, France

⁶ AP-HP, Hôpital Bicêtre, Department of Pediatric Radiology, F-94270, Le Kremlin-Bicêtre, France

⁷ AP-HP, Hôpital Robert-Debré, Department of Radiology, F-75019, Paris, France

⁸ AP-HP, Hôpital Lariboisière, Department of Neuroradiology, F-75010, Paris, France

⁹ AP-HP, Hôpital Raymond-Poincaré, Department of Radiology, F-92380, Garches, France

¹⁰ AP-HP, Hôpital Saint-Antoine, Department of Radiology, F-75012, Paris, France

¹¹ AP-HP, Hôpital Tenon, Department of Radiology, F-75020, Paris, France

¹² AP-HP, Hôpital Henri-Mondor, Department of Radiology, F-94000, Créteil, France

¹³ AP-HP, Hôpital Bichat, Department of Radiology, F-75018, Paris, France

¹⁴ AP-HP, Hôpital Hôtel-Dieu, Department of Radiology, F-75004, Paris, France

¹⁵ AP-HP, Hôpital Antoine-Béclère, Department of Radiology, F-92140, Clamart, France

¹⁶ AP-HP, Hôpital Avicenne, Department of Radiology, F-93000, Bobigny, France

¹⁷ AP-HP, Hôpital Ambroise Paré, Department of Radiology, F-92100 104, Boulogne-Billancourt, France

¹⁸ AP-HP, Hôpital Cochin, Department of Radiology, F-75014, Paris, France

¹⁹ AP-HP, WIND department, F-75012, Paris, France

²⁰ AP-HP, Unité de Recherche Clinique, Hôpital de la Pitié Salpêtrière, Department of Neuroradiology, F-75013, Paris, France

²¹ ICM, Centre de NeuroImagerie de Recherche – CENIR, Paris, France

²² Sorbonne Université, Institut du Cerveau – Paris Brain Institute, Inserm, CNRS, AP-HP, Hôpital de la Pitié Salpêtrière, F-75013, Paris, France

Appendix B. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.media.2023.103073>.

References

- Ahmed, K.B., Hall, L.O., Goldgof, D.B., Liu, R., Gatenby, R.A., 2017. Fine-tuning convolutional deep features for MRI based brain tumor classification. In: SPIE Medical Imaging 2017, Vol. 10134. pp. 613–619. <http://dx.doi.org/10.1117/12.2253982>.
- Al-masni, M.A., Lee, S., Yi, J., Kim, S., Gho, S.-M., Choi, Y.H., Kim, D.-H., 2022. Stacked U-Nets with self-assisted priors towards robust correction of rigid motion artifact in brain MRI. *NeuroImage* 259, 119411. <http://dx.doi.org/10.1016/j.neuroimage.2022.119411>.
- Alexander-Bloch, A., Clasen, L., Stockman, M., Ronan, L., Lalonde, F., Giedd, J., Raznahan, A., 2016. Subtle in-scanner motion biases automated measurement of brain anatomy from in vivo MRI. *Hum. Brain Mapp.* 37 (7), 2385–2397. <http://dx.doi.org/10.1002/hbm.23180>.
- Andre, J.B., Bresnahan, B.W., Mossa-Basha, M., Hoff, M.N., Smith, C.P., Anzai, Y., Cohen, W.A., 2015. Toward quantifying the prevalence, severity, and cost associated with patient motion during clinical MR examinations. *J. Am. Coll. Radiol.* 12 (7), 689–695. <http://dx.doi.org/10.1016/j.jacr.2015.03.007>.

- Avants, B.B., Epstein, C.L., Grossman, M., Gee, J.C., 2008. Symmetric diffeomorphic image registration with cross-correlation: Evaluating automated labeling of elderly and neurodegenerative brain. *Med. Image Anal.* 12 (1), 26–41. <http://dx.doi.org/10.1016/j.media.2007.06.004>.
- Billot, B., Greve, D.N., Puonti, O., Thielscher, A., Van Leemput, K., Fischl, B., Dalca, A.V., Iglesias, J.E., 2021. SynthSeg: Domain randomisation for segmentation of brain MRI scans of any contrast and resolution. [arXiv:2107.09559](https://arxiv.org/abs/2107.09559).
- Bottani, S., Burgos, N., Maire, A., Wild, A., Strer, S., Dormont, D., Colliot, O., 2021. Automatic quality control of brain T1-weighted magnetic resonance images for a clinical data warehouse. *Med. Image Anal.* 102219. <http://dx.doi.org/10.1016/j.media.2021.102219>.
- Commowick, O., Istace, A., Kain, M., Laurent, B., Leray, F., Simon, M., Pop, S.C., Girard, P., Ameli, R., Ferré, J.-C., et al., 2018. Objective evaluation of multiple sclerosis lesion segmentation using a data management and processing infrastructure. *Sci. Rep.* 8 (1), 13650.
- Couvy-Duchesne, B., Faouzi, J., Martin, B., Thibeau-Sutre, E., Wild, A., Ansart, M., Durrleman, S., Dormont, D., Burgos, N., Colliot, O., 2020. Ensemble learning of convolutional neural network, support vector machine, and best linear unbiased predictor for brain age prediction: ARAMIS contribution to the predictive analytics competition 2019 challenge. *Front. Psychiatry* 11.
- Duffy, B.A., Zhang, W., Tang, H., Zhao, L., Law, M., Toga, A.W., Kim, H., 2018. Retrospective correction of motion artifact affected structural MRI images using deep learning of simulated motion. In: *Medical Imaging with Deep Learning*, p. 8.
- Esteban, O., Birman, D., Schaer, M., Koyejo, O.O., Poldrack, R.A., Gorgolewski, K.J., 2017. MRIQC: Advancing the automatic prediction of image quality in MRI from unseen sites. *PLoS One* 12 (9), e0184661. <http://dx.doi.org/10.1371/journal.pone.0184661>.
- Fantini, I., Yasuda, C., Bento, M., Rittner, L., Cendes, F., Lotufo, R., 2021. Automatic MR image quality evaluation using a deep CNN: A reference-free method to rate motion artifacts in neuroimaging. *Comput. Med. Imaging Graph.* (ISSN: 0895-6111) 90, 101897. <http://dx.doi.org/10.1016/j.compmedimag.2021.101897>.
- Fischl, B., 2012. FreeSurfer. *NeuroImage* 62 (2), 774–781.
- Ghosal, P., Nandanwar, L., Kanchan, S., Bhadra, A., Chakraborty, J., Nandi, D., 2019. Brain tumor classification using ResNet-101 Based Squeeze and excitation deep neural network. In: *International Conference on Advanced Computational and Communication Paradigms*. pp. 1–6. <http://dx.doi.org/10.1109/ICACCP.2019.8882973>.
- Hedges, E.P., Dimitrov, M., Zahid, U., Brito Vega, B., Si, S., Dickson, H., McGuire, P., Williams, S., Barker, G.J., Kempton, M.J., 2022. Reliability of structural MRI measurements: The effects of scan session, head tilt, inter-scan interval, acquisition sequence, FreeSurfer version and processing stream. *NeuroImage* 246, 118751. <http://dx.doi.org/10.1016/j.neuroimage.2021.118751>.
- Iglesias, J.E., Lerma-Usabiaga, G., Garcia-Peraza-Herrera, L.C., Martinez, S., Paz-Alonso, P.M., 2017. Retrospective head motion estimation in structural brain MRI with 3D CNNs. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, pp. 314–322.
- Jannot, A.-S., Zapletal, E., Avillach, P., Mamzer, M.-F., Burgun, A., Degoulet, P., 2017. The Georges Pompidou university hospital clinical data warehouse: A 8-years follow-up experience. *Int. J. Med. Inform.* 102, 21–28. <http://dx.doi.org/10.1016/j.ijmedinf.2017.02.006>.
- Jenkinson, M., Beckmann, C.F., Behrens, T.E., Woolrich, M.W., Smith, S.M., 2012. FSL. *NeuroImage* 62 (2), 782–790.
- Jonsson, B.A., Björnsdóttir, G., Thorgeirsson, T.E., Ellingsen, L.M., Walters, G.B., Gudbjartsson, D.F., Stefansson, H., Stefansson, K., Ulfarsson, M.O., 2019. Brain age prediction using deep learning uncovers associated sequence variants. *Nature Commun.* 10 (1), 5409. <http://dx.doi.org/10.1038/s41467-019-13163-9>.
- Karami, M., Rahimi, A., Shahmirzadi, A.H., 2017. Clinical data warehouse: an effective tool to create intelligence in disease management. *Health Care Manager* 36 (4), 380–384.
- Karani, N., Chaitanya, K., Baumgartner, C., Konukoglu, E., 2018. A lifelong learning approach to brain MR segmentation across scanners and protocols. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, pp. 476–484.
- Küstner, T., Liebgott, A., Mauch, L., Martirosian, P., Bamberg, F., Nikolaou, K., Yang, B., Schick, F., Gatidis, S., 2018. Automated reference-free detection of motion artifacts in magnetic resonance images. *Magn. Reson. Mater. Phys.* 31 (2), 243–256.
- Lee, S., Jung, S., Jung, K.-J., Kim, D.-H., 2020. Deep learning in MR motion correction: A brief review and a new motion simulation tool (view2Dmotion). *Invest. Magn. Reson. Imaging* 24 (4), 196. <http://dx.doi.org/10.13104/imri.2020.24.4.196>.
- Lei, K., Syed, A.B., Zhu, X., Pauly, J.M., Vasanawala, S.S., 2022. Artifact- and content-specific quality assessment for MRI with image rulers. *Med. Image Anal.* (ISSN: 1361-8415) 77, 102344. <http://dx.doi.org/10.1016/j.media.2021.102344>.
- Loizillon, S., Bottani, S., Maire, A., Ströer, S., Dormont, D., Colliot, O., Burgos, N., 2023. Transfer learning from synthetic to routine clinical data for motion artefact detection in brain T1-weighted MRI. In: *SPIE Medical Imaging 2023*, p. 6.
- Loktyushin, A., Nickisch, H., Pohmann, R., Schölkopf, B., 2013. Blind retrospective motion correction of MR images. *J. Magn. Reson. Imaging* 70 (6), 1608–1618.
- Mercier, L., Del Maestro, R.F., Petrecca, K., Araujo, D., Haegelen, C., Collins, D.L., 2012. Online database of clinical MR and ultrasound images of brain tumors. *Med. Phys.* 39 (6Part1), 3253–3261. <http://dx.doi.org/10.1118/1.4709600>.
- Mia, M.R., Hoque, A.S.M.L., Khan, S.I., Ahamed, S.I., 2022. A privacy-preserving national clinical data warehouse: Architecture and analysis. *Smart Health* 23, 100238. <http://dx.doi.org/10.1016/j.smhl.2021.100238>.
- Mohebbian, M., Wallia, E., Habibullah, M., Stapleton, S., Wahid, K.A., 2021. Classifying MRI motion severity using a stacked ensemble approach. *J. Magn. Reson. Imaging* 75, 107–115. <http://dx.doi.org/10.1016/j.mri.2020.10.007>.
- Nordlinger, B., Villani, C., Rus, D., 2020. *Healthcare and Artificial Intelligence*. Springer International Publishing.
- Oksuz, I., 2021. Brain MRI artefact detection and correction using convolutional neural networks. *Comput. Methods Prog. Biol.* 199, 105909. <http://dx.doi.org/10.1016/j.cmpb.2020.105909>.
- Pawar, K., Chen, Z., Shah, N.J., Egan, G.F., 2022. Suppressing motion artefacts in MRI using an inception-ResNet network with motion simulation augmentation. *NMR Biomed.* 35 (4), e4225. <http://dx.doi.org/10.1002/nbm.4225>.
- Penny, W.D., Friston, K.J., Ashburner, J.T., Kiebel, S.J., Nichols, T.E., 2011. *Statistical Parametric Mapping: the Analysis of Functional Brain Images*. Elsevier.
- Pérez-García, F., Sparks, R., Ourselin, S., 2021. TorchIO: A Python library for efficient loading, preprocessing, augmentation and patch-based sampling of medical images in deep learning. *Comput. Methods Prog. Biol.* 208, 106236. <http://dx.doi.org/10.1016/j.cmpb.2021.106236>.
- Petersen, R.C., Aisen, P.S., Beckett, L.A., Donohue, M.C., Gamst, A.C., Harvey, D.J., Jack, C.R., Jagust, W.J., Shaw, L.M., Toga, A.W., Trojanowski, J.Q., Weiner, M.W., 2010. Alzheimer's disease neuroimaging initiative (ADNI). *Neurology* 74 (3), 201–209. <http://dx.doi.org/10.1212/WNL.0b013e3181cb3e25>.
- Ravi, D., Barkhof, F., Alexander, D.C., Parker, G.J., Eshghi, A., 2022. An efficient semi-supervised quality control system trained using physics-based MRI-artefact generators and adversarial training. [arXiv:2206.03359](https://arxiv.org/abs/2206.03359).
- Reguig, G., Lapert, M., Lehericy, S., Valabregue, R., 2022. Global displacement induced by rigid motion simulation during MRI acquisition. [arXiv:2204.03522](https://arxiv.org/abs/2204.03522).
- Reuter, M., Tisdall, M.D., Qureshi, A., Buckner, R.L., van der Kouwe, A.J.W., Fischl, B., 2015. Head motion during MRI acquisition reduces gray matter volume and thickness estimates. *NeuroImage* 107, 107–115. <http://dx.doi.org/10.1016/j.neuroimage.2014.12.006>.
- Routier, A., Burgos, N., Díaz, M., Bacci, M., Bottani, S., El-Rifai, O., Fontanella, S., Gori, P., Guillon, J., Guyot, A., Hassanaly, R., Jacquemont, T., Lu, P., Marcoux, A., Moreau, T., Samper-González, J., Teichmann, M., Thibeau-Sutre, E., Vaillant, G., Wen, J., Wild, A., Habert, M.-O., Durrleman, S., Colliot, O., 2021. Clinica: An open-source software platform for reproducible clinical neuroscience studies. *Front. Neuroinform.* 15, 689675. <http://dx.doi.org/10.3389/fninf.2021.689675>.
- Sadri, A.R., Janowczyk, A., Zou, R., Verma, R., Beig, N., Antunes, J., Madabhushi, A., Tiwari, P., Viswanath, S.E., 2020. MRQy: An open-source tool for quality control of MR imaging data. *Med. Phys.* 47 (12), 6029–6038. <http://dx.doi.org/10.1002/mp.14593>.
- Sagawa, H., Itagaki, K., Matsushita, T., Miyati, T., 2022. Evaluation of motion artifacts in brain magnetic resonance images using convolutional neural network-based prediction of full-reference image quality assessment metrics. *J. Med. Imaging* 9 (1), 015502. <http://dx.doi.org/10.1117/1.JMI.9.1.015502>.
- Shaw, R., Sudre, C., Ourselin, S., Cardoso, M.J., 2019. MRI k-space motion artefact augmentation: model robustness and task-specific uncertainty. In: *Medical Imaging with Deep Learning*, p. 10.
- Shaw, R., Sudre, C.H., Ourselin, S., Cardoso, M.J., Pemberton, H.G., 2021. A decoupled uncertainty model for MRI segmentation quality estimation. [arXiv:2109.02413](https://arxiv.org/abs/2109.02413).
- Sujit, S.J., Coronado, I., Kamali, A., Narayana, P.A., Gabr, R.E., 2019. Automated image quality evaluation of structural brain MRI using an ensemble of deep learning networks. *J. Magn. Reson. Imaging* 50 (4), 1260–1267.
- Tajbakhsh, N., Shin, J.Y., Gurudu, S.R., Hurst, R.T., Kendall, C.B., Gotway, M.B., Liang, J., 2016. Convolutional neural networks for med. Image anal.: Full training or fine tuning? *IEEE Trans. Med. Imaging* 35 (5), 1299–1312.
- Thibeau-Sutre, E., Díaz, M., Hassanaly, R., Routier, A., Dormont, D., Colliot, O., Burgos, N., 2022. ClinicaDL: An open-source deep learning software for reproducible neuroimaging processing. *Comput. Methods Prog. Biol.* 220, 106818. <http://dx.doi.org/10.1016/j.cmpb.2022.106818>.
- Tustison, N.J., Avants, B.B., Cook, P.A., Zheng, Y., Egan, A., Yushkevich, P.A., Gee, J.C., 2010. N4ITK: Improved N3 bias correction. *IEEE Trans. Med. Imaging* 29 (6), 1310–1320. <http://dx.doi.org/10.1109/TMI.2010.2046908>.
- Wang, W., Chen, C., Ding, M., Yu, H., Zha, S., Li, J., 2021. TransBTS: Multimodal brain tumor segmentation using transformer. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, pp. 109–119.
- Wen, J., Thibeau-Sutre, E., Diaz-Melo, M., Samper-González, J., Routier, A., Bottani, S., Dormont, D., Durrleman, S., Burgos, N., Colliot, O., 2020. Convolutional neural networks for classification of Alzheimer's disease: Overview and reproducible evaluation. *Med. Image Anal.* 101694. <http://dx.doi.org/10.1016/j.media.2020.101694>.
- Wood, M.L., Henkelman, R.M., 1985. Truncation artifacts in magnetic resonance imaging. *Magn. Reson. Med.* <http://dx.doi.org/10.1002/MRM.1910020602>.
- Xing, X., Liang, G., Zhang, Y., Khanal, S., Lin, A.-L., Jacobs, N., 2022. ADViT: Vision transformer on multi-modality PET images for alzheimer disease diagnosis. In: *IEEE International Symposium on Biomedical Imaging*, pp. 1–4. <http://dx.doi.org/10.1109/ISBI52829.2022.9761584>.