

A latent space for unsupervised MR image quality control via artifact assessment

Lianrui Zuo^{a,b}, Yuan Xue^a, Blake E. Dewey^c, Yihao Liu^a,
Jerry L. Prince^a, and Aaron Carass^a

^aDepartment of Electrical and Computer Engineering,
Johns Hopkins University, Baltimore, MD 21218, USA

^bLaboratory of Behavioral Neuroscience, National Institute on Aging,
National Institutes of Health, Baltimore, MD 20892, USA

^cDepartment of Neurology, Johns Hopkins School of Medicine, Baltimore, MD 21287, USA

ABSTRACT

Image quality control (IQC) can be used in automated magnetic resonance (MR) image analysis to exclude erroneous results caused by poorly acquired or artifact-laden images. Existing IQC methods for MR imaging generally require human effort to craft meaningful features or label large datasets for supervised training. The involvement of human labor can be burdensome and biased, as labeling MR images based on their quality is a subjective task. In this paper, we propose an automatic IQC method that evaluates the extent of artifacts in MR images without supervision. In particular, we design an artifact encoding network that learns representations of artifacts based on contrastive learning. We then use a normalizing flow to estimate the density of learned representations for unsupervised classification. Our experiments on large-scale multi-cohort MR datasets show that the proposed method accurately detects images with high levels of artifacts, which can inform downstream analysis tasks about potentially flawed data.

Keywords: magnetic resonance imaging, contrastive learning, artifacts, quality assurance

1. INTRODUCTION

The recent development of deep learning (DL) has benefited various magnetic resonance (MR) image analyses, such as image synthesis,^{1,2} segmentation,^{3,4} registration,⁵ and volumetric analysis,⁶ where a large amount of images are processed without human intervention. Yet, these DL based algorithms are known to be sensitive to the quality of input images;⁷ when an image is poorly acquired or contaminated by artifacts, the DL algorithms are likely to produce erroneous or biased results. Manually inspecting DL results in large datasets is prone to errors as it is tedious and subjective. Therefore, there is demand for an automatic image quality control (IQC) method to identify potential failures cases caused by either poor quality or inappropriate data.

Various IQC methods have been developed in recent years.^{8,9} The goal of an IQC method is to provide an assessment of image quality \hat{y} based on the input image x . In general, an IQC algorithm has two parts: feature extraction and classification. Features m that capture image quality information can either be handcrafted with expert knowledge^{8,10} or learned from data.⁹ Classification is conducted based on the features m , which usually requires expert labels on a sample dataset—e.g., with $y \in \{0, 1\}$ indicating whether image x passes or fails quality inspection—from which a supervised classifier is trained. For example, MRIQC⁸ learned a binary classifier based on handcrafted features and labels generated by human experts. However, the current IQC methods face two major limitations. First, labeling datasets by experts requires domain specific knowledge, which can be subjective and time consuming. Second, because the labels y are limited in number and dataset specific, current IQC

Corresponding author: Lianrui Zuo.
Email: lr_zuo@jhu.edu

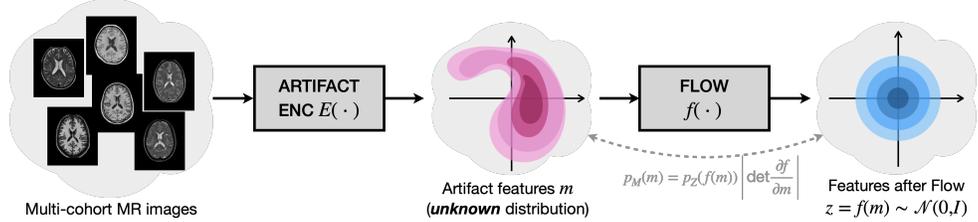


Figure 1. Schematic framework of the proposed IQC method. Artifact encoder $E(\cdot)$ extracts artifact features $m \in \mathbb{R}^2$ from multi-cohort MR images. The learned features m follow an unknown distribution $p_M(m)$. A normalizing flow $f(\cdot)$ is then applied to transform m to $z \in \mathbb{R}^2$ following a standard Gaussian distribution $\mathcal{N}(0, I)$. Due to the special property of $f(\cdot)$, the likelihood $p_M(m)$ can be evaluated using Eq. 2.

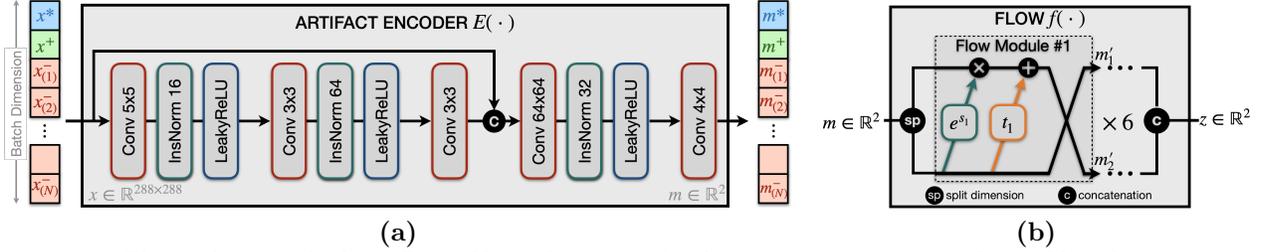


Figure 2. (a) The artifact encoder has a Dense-Net architecture. Artifact representations m are learned based on contrastive learning with both positive images x^+ and negative images $x_{(i)}^-$. x^* shares the same and different artifact levels with x^+ and $x_{(i)}^-$, respectively. (b) The flow network is based on RealNVP¹¹ with six flow modules (affine coupling layers). $\{s_i, t_i\}_{i=1}^6$ are trainable neural networks.

methods usually have limited generalizability. The feature extractor and classifier are usually generalizable to datasets similar to what they have been trained on; however, new datasets will generally require a re-training or fine tuning, critically this necessitates new labels.

To overcome the limitations of current IQC methods, we developed an unsupervised IQC method to directly assess artifact levels from MR images. Our method has two advantages. First, we propose an artifact encoder network that learns latent artifact representations in a data-driven way. Second, we use a normalizing flow¹¹ to map the learned representations m to a normal distribution, which allows us to conduct unsupervised classification without expert labels y . It is worth noting that the artifact encoder is also unsupervised, meaning that no labels are needed in our method. This means our approach is completely unsupervised, making our framework applicable to more datasets.

2. METHODS

Figure 1 shows the framework of the proposed method. MR images from multiple cohorts are first encoded into a two-dimensional latent space of artifact features m . In general, m follows an unknown distribution $p_M(m)$. We then apply a normalizing flow¹¹ $f(\cdot)$ to transform m to $z \in \mathbb{R}^2$, which follows a standard Gaussian distribution. $f(\cdot)$ also enables density estimation of $p_M(m)$ for unsupervised IQC.

2.1 Artifact encoder based on contrastive learning

Our artifact encoder $E(\cdot)$ extracts artifact representations based on contrastive learning.¹² The key concept of contrastive learning is to learn discriminative features from query, positive, and negative examples. Figure 2(a) shows the architecture of $E(\cdot)$. For each MR image $x \in \mathbb{R}^{288 \times 288}$, we assume the positive example x^+ has the same artifact level as the query example x^* . We achieve this by selecting image slices x^* and x^+ from different orientations of the same 3D volume (e.g., axial and coronal slices). Our negative examples $\{x_{(i)}^-\}_{i=1}^N$ are chosen to have different artifact levels than x^* . We prepare our negative examples by either selecting slices from a volume different from the source of x^* or augmenting x^* with simulated artifacts including noise and motion. The simulated images are used to prevent $E(\cdot)$ from learning irrelevant information such as contrast and anatomy, since slices from different volumes may differ both in their level of artifacts and in their contrasts and anatomies.

Because we also introduce real MR images as negative examples, $E(\cdot)$ after training can capture different kinds of artifacts—beyond just noise and motion—which we show in Sec. 3.

With x^* , x^+ , and $x_{(i)}^-$'s composing our input mini-batch, we expect the learned feature m^* to be similar (if not identical) to m^+ and sufficiently distinct from the $m_{(i)}^-$'s. We encourage this relationship using

$$\mathcal{L}(m^*, m^+, \{m_{(i)}^-\}_{i=1}^N) = -\log \left[\frac{\exp(m^* \cdot m^+)}{\exp(m^* \cdot m^+) + \frac{1}{N} \sum_{i=1}^N \exp(m^* \cdot m_{(i)}^-)} \right] \quad (1)$$

as our loss function for $E(\cdot)$. Since we prepare our x^+ and $x_{(i)}^-$ based on their relative extent of artifact with x^* and encourage m to preserve this relationship, we would expect m to capture the artifact information of the input image. Note that m is learned based on the *relative* extent of artifact between x^* , x^+ , and $x_{(i)}^-$, there is no assumption made about the *absolute* extent of artifact of x^* (i.e. x^* is not assumed to be free from artifacts).

2.2 Density estimation with normalizing flows

With sufficiently large datasets, one can assume that most acquired MR images have acceptable image quality with a relatively small sample of images being poorly acquired or contaminated by artifacts. This ratio is reflected by the likelihood $p_M(m)$; when an image has uncommonly high artifact level \tilde{m} , we would expect $p_M(\tilde{m})$ to be small. Unsupervised IQC can then be achieved by finding $\{x_i\}$'s with $p_M(m_i)$'s below a percentile. Evaluating $p_M(m)$ is a nontrivial task, but it can be approximated using a normalizing flow network¹¹ $f(\cdot)$. As shown in Fig. 2(b), $f(\cdot)$ is composed of six flow modules with each module affinely processing a proportion of the input variable, e.g., $m'_1 = m_2$ and $m'_2 = m_1 \cdot e^{s_1(m_2)} + t_1(m_2)$, where $m = [m_1, m_2] \in \mathbb{R}^2$ and $\{s_i, t_i\}_{i=1}^6$ are neural networks. The output variable $z = f(m)$ follows a standard Gaussian distribution $\mathcal{N}(0, I)$. It is easy to show that the Jacobian matrix of $f(m)$ is a triangular matrix with positive determinant and the density $p_M(m)$ can be calculated by

$$p_M(m) = p_Z(f(m)) \left| \det \frac{\partial f(m)}{\partial m} \right|. \quad (2)$$

During training, we use $-\log p_M(m)$ calculated with Eq. 2 as our loss function for $f(\cdot)$, where the trainable modules are $\{s_i, t_i\}_{i=1}^6$.

3. EXPERIMENTS AND RESULTS

3.1 Datasets and preprocessing

The training data for $E(\cdot)$ include 200 T₁-weighted (T₁-w) MR volumes acquired from 16 different cohorts. Detailed information about image acquisition is provided in Table 1. Our preprocessing includes inhomogeneity correction¹³ and registration to a 0.8mm³ isotropic template. For each 3D volume, we extracted and padded axial, coronal, and sagittal slices to dimension 288 × 288. $E(\cdot)$ was trained on 2D slices following Sec. 2.1. Our evaluation dataset for $E(\cdot)$ has 1,400 3D images acquired from the 16 cohorts. For each volume, we calculated the average m values of its 20 center axial slices as the artifact representation. $f(\cdot)$ was then trained and applied following Sec. 2.2 to estimate the density $p_M(m)$ based on all the 1,400 volumes.

3.2 Unsupervised IQC on simulated data

After training, $f(\cdot)$ transforms m to z , which follows a standard Gaussian distribution. Density $p_M(m)$ can then be evaluated using Eq. 2. Figure 3(a) shows z values of the 1,400 volumes after applying the normalizing flow $f(\cdot)$. We then applied the proposed method to a held-out simulated dataset with various kinds of artifacts that could potentially fail downstream analyses. In Fig. 3(b), m 's of eight representative images are shown on top of the density contours of the 1,400 volumes. The eight images are A) a T₁-w image that passed our manual inspection, B) a non-T₁-w image, C) a T₁-w image with a bias field, D) a T₁-w image with high noise, E) a T₁-w image with motion artifacts, F) a T₁-w image with wrap-around artifacts, G) a T₁-w image with one side of the head removed, and H) a T₁-w image with registration errors. We assume the original 1,400 volumes are fairly diverse samples of T₁-w MR images and that most of them have acceptable artifact levels with a small proportion

Table 1. Key information about image acquisition of each imaging cohort. Unavailable information is marked as “-”. Data source: C_1 and C_2 (IXI-Brain);¹⁴ C_3 thru C_6 (OASIS3);¹⁵ C_7 thru C_{10} (BLSA);¹⁶ C_{11} thru C_{16} (Private).

Cohort	C_1	C_2	C_3	C_4	C_5	C_6	C_7	C_8
Open data	✓	✓	✓	✓	✓	✓	✓	✓
Manufacturer	Philips	Philips	Siemens	Siemens	Siemens	Siemens	Philips	Philips
Field (T)	1.5	3.0	3.0	3.0	3.0	1.5	1.5	3.0
Resolution (mm)	$1.2 \times 0.9 \times 0.9$	$1.2 \times 0.9 \times 0.9$	$1.0 \times 1.0 \times 1.0$	$1.0 \times 1.0 \times 1.0$	$1.0 \times 1.0 \times 1.0$	$1.1 \times 1.1 \times 1.2$	$0.9 \times 0.9 \times 1.5$	$1.0 \times 1.0 \times 1.2$
TE/TR/TI (ms)	4.6/-/-	4.6/-/-	3.9/1900/1100	3.2/2400/1000	3.2/2400/1000	2.9/2300/900	3.3/3000/-	3.1/3000/800
Cohort	C_9	C_{10}	C_{11}	C_{12}	C_{13}	C_{14}	C_{15}	C_{16}
Open data	✓	✓	✗	✗	✗	✗	✗	✗
Manufacturer	Philips	Philips	Siemens	GE	Siemens	GE	Siemens	Siemens
Field (T)	1.5	3.0	3.0	3.0	3.0	1.5	1.5	3.0
Resolution (mm)	$1.2 \times 0.9 \times 0.9$	$1.2 \times 0.9 \times 0.9$	$1.0 \times 1.0 \times 1.0$	$1.0 \times 1.0 \times 1.0$	$1.0 \times 1.0 \times 1.0$	$1.1 \times 1.1 \times 1.2$	$0.9 \times 0.9 \times 1.5$	$1.0 \times 1.0 \times 1.2$
TE/TR/TI (ms)	3.1/3000/800	3.1/3000/800	3.0/2300/900	3.1/-/-	3.6/2500/-	2.6/-/-	3.0/2300/900	3.4/2300/900

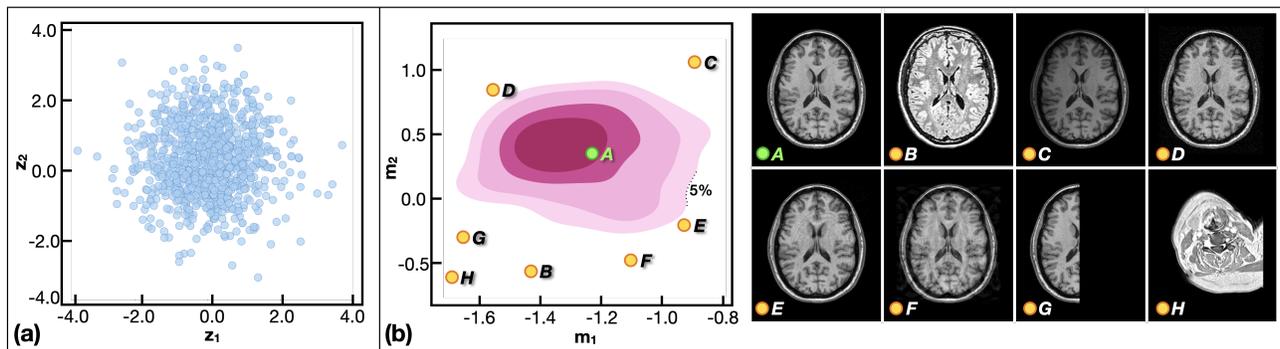


Figure 3. (a) scatter plot of z from 1,400 T₁-w MR volumes after a normalizing flow $z = f(m)$. (b) m values of held-out MR images are shown on top of the density contours fitted on m . Eight example MR images are shown on the right. Indexes A to H represent an MR image A) a T₁-w image that passed our manual inspection, B) a non-T₁-w image, C) a T₁-w image with a bias field, D) a T₁-w image with high noise, E) a T₁-w image with motion artifacts, F) a T₁-w image with wrap-around artifacts, G) a T₁-w image with one side of the head removed, and H) a T₁-w image with registration errors.

being poorly acquired. Unsupervised IQC is achieved by thresholding $p_M(m)$ with a predefined threshold τ . We found $\tau = 5\%$ achieved satisfactory results on our simulated dataset. As shown in Fig. 3(b), the image that passed our manual inspection has an m located in the high density region, while the remaining seven images have $p_M(m)$ below τ . Our unsupervised IQC method has two advantages over existing works. First, we do not require knowledge of the absolute artifact levels of training images, so that our method can be trained on very large datasets. In fact, we only assume that most our training data have acceptable image quality; this is likely true in many application scenarios. Based on contrastive learning, our artifact encoder $E(\cdot)$ during training only needs to know if a sample has the same (for positive examples) or a different (for negative examples) artifact level as the query image x^* . Second, since we construct our negative examples with both real data and simulated artifacts (i.e., motion and noise), $E(\cdot)$ after training can capture various kinds of artifacts, many of which have not been simulated in training. This makes our model more generalizable.

3.3 Quantitative evaluation on real data

To quantitatively evaluate the proposed method on real MR datasets, we manually inspected and rated 569 T₁-w MR images acquired from cohorts C_{11} to C_{16} (see Table 1 for more details). After manual inspection, each image was assigned a label from one of the three labels low, medium, or high based on the level of artifacts present in the volume. We assume images with low levels of artifacts passed our manual quality check, and assume images with either medium or high artifact levels as failed cases. Figure 4(a) shows the learned m values of the images with manual ratings. Green, orange, and red represent low, medium, and high levels of artifacts, respectively. Two

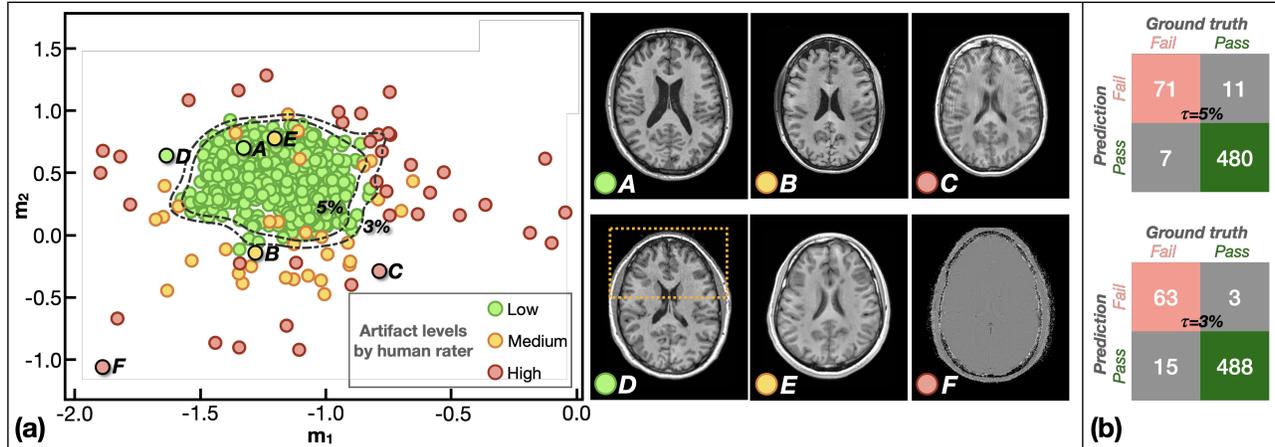


Figure 4. (a) Scatter plot of m values on a dataset with manual ratings. Example images are shown on the right. Dashed lines show 5% and 3% likelihood contours of $p_M(m)$. (b) Contingency tables of the proposed method based on $\tau = 5\%$ and $\tau = 3\%$. Images with either medium or high artifacts levels are categorized as failed cases.

density contours (5% and 3%) of $p_M(m)$ are also shown in Fig. 4(a). It is encouraging to see that most images that passed our manual inspection (green) have m values with $p_M(m) > 5\%$, while most images that failed our manual inspection (with medium and high artifact levels) have $p_M(m) < 5\%$. Furthermore, images with high levels of artifacts (red) usually have even lower $p_M(m)$ than images with medium levels of artifacts (orange). Figure 4(a) also shows six example images with different levels of artifacts, where *A*) has passed our manual quality check and it has $p_M(m) > 5\%$. *B*) has a medium level of artifacts due to the intensity inhomogeneity, and *C*) has strong motion artifacts. Interestingly, image *D*) has passed our manual quality check, but the proposed method identified it as a low density example ($p_M(m) < 5\%$). We hypothesize the reason for this is because the uncommon noise pattern of the image; the noise level is only high inside the orange box. *E*) shows an example with medium artifact level according to our manual inspections, but our method failed to identify it as a poor quality image. *F*) is an extreme case where a non-T₁-w image was processed and identified by our algorithm as potential bad data.

In Fig. 4(b), we show the contingency tables of the proposed method based on two thresholds: $\tau = 5\%$ and $\tau = 3\%$. Here, we assume any images with $p_M(m) < \tau$ at test time should be highlighted as potential artifact-laden images (potential failed cases). $\tau = 5\%$, which we used on simulated data in Sec. 3.2, achieves a sensitivity of 91.0% and a specificity of 97.8%. $\tau = 3\%$ achieves a sensitivity of 80.0% and a specificity of 99.4%.

4. DISCUSSION AND CONCLUSION

In this paper, we present a novel unsupervised IQC approach by assessing the levels of artifacts from MR images. Our approach learns representations of image artifacts without domain knowledge. This unsupervised nature enables our approach to be trained on a large variety of datasets with improved applicability over existing IQC methods. We showcase using normalizing flow that after artifact representations are learned, classification can be achieved with a simple thresholding on feature densities. The fact that the threshold τ needs to be determined at test time is a limitation of our work, as it may vary from dataset to dataset. We regard this as a direction for future improvements. We believe introducing a very small amount of labels during training (for semi-supervised training) would benefit the proposed method to learn more robust feature extractors and classifiers.

Experiments on both simulated and real MR datasets show that the proposed method achieves both high sensitivity and specificity. Our approach can be used to inform downstream analyses about potential bad quality data by accurately highlighting different kinds of artifact cases as low likelihood examples.

ACKNOWLEDGMENTS

This work was supported in part by the Intramural Research Program of the NIH, National Institute on Aging and in part by the TREAT-MS study funded by the Patient-Centered Outcomes Research Institute (PCORI/MS-

REFERENCES

- [1] Zuo, L., Dewey, B. E., Carass, A., He, Y., Shao, M., Reinhold, J. C., and Prince, J. L., “Synthesizing realistic brain MR images with noise control,” in [*International Workshop on Simulation and Synthesis in Medical Imaging*], *Lecture Notes in Computer Science* **12417**, 21–31, Springer (2020).
- [2] Zuo, L., Liu, Y., Xue, Y., Han, S., Bilgel, M., Resnick, S. M., Prince, J. L., and Carass, A., “Disentangling a Single MR Modality,” in [*Data Augmentation, Labelling, and Imperfections*], 54–63, Springer Nature Switzerland, Cham (2022).
- [3] Huo, Y., Xu, Z., Xiong, Y., Aboud, K., Parvathaneni, P., Bao, S., Bermudez, C., Resnick, S. M., Cutting, L. E., and Landman, B. A., “3D whole brain segmentation using spatially localized atlas network tiles,” *NeuroImage* **194**, 105–119 (2019).
- [4] Liu, Y., Carass, A., Zuo, L., He, Y., Han, S., Gregori, L., Murray, S., Mishra, R., Lei, J., Calabresi, P. A., Saidha, S., and Prince, J. L., “Disentangled representation learning for OCTA vessel segmentation with limited training data,” *IEEE Trans. Med. Imag.* **41**(12), 3686–3698 (2022).
- [5] Liu, Y., Zuo, L., Han, S., Xue, Y., Prince, J. L., and Carass, A., “Coordinate Translator for Learning Deformable Medical Image Registration,” in [*International Workshop on Multiscale Multimodal Medical Imaging (MMMI 2022)*], **13594**, 98–109 (2022).
- [6] Duan, P., Han, S., Zuo, L., An, Y., Liu, Y., Alshareef, A., Lee, J., Carass, A., Resnick, S. M., and Prince, J. L., “Cranial meninges reconstruction based on convolutional networks and deformable models: Applications to longitudinal study of normal aging,” in [*Medical Imaging 2022: Image Processing*], **12032**, 299–305, SPIE (2022).
- [7] Kügler, D., Distergoft, A., Kuijper, A., and Mukhopadhyay, M., “Exploring Adversarial Examples,” in [*Understanding and Interpreting Machine Learning in Medical Image Computing Applications*], *Lecture Notes in Computer Science* **11038**, 70–78, Springer International Publishing (2018).
- [8] Esteban, O., Birman, D., Schaer, M., Koyejo, O. O., Poldrack, R. A., and Gorgolewski, K. J., “MRIQC: Advancing the automatic prediction of image quality in MRI from unseen sites,” *PloS one* **12**(9), e0184661 (2017).
- [9] Kang, L., Ye, P., Li, Y., and Doermann, D., “Convolutional neural networks for no-reference image quality assessment,” in [*Proceedings of the IEEE conference on computer vision and pattern recognition*], 1733–1740 (2014).
- [10] Zuo, L., Carass, A., Han, S., and Prince, J. L., “Automatic outlier detection using hidden Markov model for cerebellar lobule segmentation,” in [*Medical Imaging 2018: Biomedical Applications in Molecular, Structural, and Functional Imaging*], **10578**, 84–90, SPIE (2018).
- [11] Dinh, L., Sohl-Dickstein, J., and Bengio, S., “Density estimation using Real NVP,” *arXiv preprint arXiv:1605.08803* (2016).
- [12] Park, T., Efros, A. A., Zhang, R., and Zhu, J.-Y., “Contrastive learning for unpaired image-to-image translation,” in [*European Conference on Computer Vision*], 319–345, Springer (2020).
- [13] Tustison, N. J., Avants, B. B., Cook, P. A., Zheng, Y., Egan, A., Yushkevich, P. A., and Gee, J. C., “N4ITK: improved N3 bias correction,” *IEEE Trans. Med. Imag.* **29**(6), 1310–1320 (2010).
- [14] “IXI Brain Development Dataset.” <https://brain-development.org/ixi-dataset/>.
- [15] LaMontagne, P. J., Benzinger, T. L., Morris, J. C., Keefe, S., Hornbeck, R., Xiong, C., Grant, E., Hassenstab, J., Moulder, K., Vlassenko, A. G., Raichle, M. E., Cruchaga, C., and Marcus, D., “OASIS-3: Longitudinal neuroimaging, clinical, and cognitive dataset for normal aging and Alzheimer disease,” *medRxiv* (2019).
- [16] Resnick, S. M., Goldszal, A. F., Davatzikos, C., Golski, S., Kraut, M. A., Metter, E. J., Bryan, R. N., and Zonderman, A. B., “One-year age changes in mri brain volumes in older adults,” *Cerebral Cortex* **10**(5), 464–472 (2000).