# U-Net Model for Brain Extraction on Non-human Primates

Xindi Wang[1], Xin-Hui Li[2], Jae Wook Cho[2], Brian Russ[3,4,5], Nanditha Rajamani[2], Alisa Omelchenko[2], Lei Ai[2], Annachiara Korchmaros[2], Pamela Garcia-Saldivar[6], Zheng Wang[7,8], Ned H. Kalin[9], Charles E. Schroeder[3,10], R. Cameron Craddock[11], Andrew S. Fox[12], Alan C. Evans[1], Adam Messinger[13], Michael P. Milham[2,3], Ting Xu[2]

1. Montreal Neurological Institute, McGill University, Montreal, Québec, Canada.
2. The Child Mind Institute, 101 East 56th Street, New York, NY, 10022, USA.
3. Nathan Kline Institute, 140 Old Orangeburg Rd, Orangeburg, NY, USA.
4. Department of Neuroscience, Icahn School of Medicine at Mount Sinai, New York City, NY, USA
5. Department of Psychiatry, New York University School of Medicine, New York City, NY, USA
6. Instituto de Neurobiología, Universidad Nacional Autónoma de México Campus Juriquilla, Querétaro, México
7. Institute of Neuroscience, CAS Center for Excellence in Brain Science and Intelligence Technology, State Key Laboratory of Neuroscience, CAS Key Laboratory of Primate Neurobiology, Chinese Academy of Science, Shanghai, China
8. University of Chinese Academy of Science, China
9. Department of Psychiatry, University of Wisconsin School of Medicine and Public Health, 6001 Research Park Blvd, Madison, WI 53719
10. Department of Neurological Surgery, Columbia University College of Physicians and Surgeons, New York, NY 10032, USA
11. Department of Diagnostic Medicine, The University of Texas at Austin Dell Medical School
12. Department of Psychology, and the California National Primate Research Center, University of California, Davis, One Shields Ave., Davis, CA 95616
13. Laboratory of Brain and Cognition, National Institute of Mental Health, Bethesda, USA

*Corresponding Authors: sandywang.rest@gmail.com, ting.xu@childmind.org

# Abstract:

Brain extraction (a.k.a. skull stripping) is a fundamental step in the neuroimaging pipeline as it can affect the accuracy of downstream preprocess such as image registration, tissue classification, etc. Most brain extraction tools have been mainly orientated for human data and are often challenging for non-human primates (NHP). In recent attempts to improve the performance in NHP, deep learning models appear to outperform the traditional tools. However, given the minimal sample size of most NHP studies and notable variations in data quality, the deep learning models are very rarely applied in multi-site samples in NHP imaging. To overcome this challenge, we propose to use transfer-learning framework that leverages a large human imaging dataset to pretrain a convolutional neural network (i.e. U-Net Model), and then transferred to NHP data using a small NHP training sample. The resulting transfer-learning model converged faster and achieved more accurate performance than a similar U-Net Model trained exclusively on NHP samples. We improved the generalizability of the model by upgrading the transfer-learned model using additional training datasets from multiple research sites in the Primate Data-Exchange (PRIME-DE) consortium. Our final model outperformed brain extraction routines from popular MRI packages (AFNI, FSL, and FreeSurfer) across multiple heterogeneous multiple sites from PRIME-DE with less computational cost (20s~10min). Our model, code, and the skull-stripped mask repository of 136 macaque monkeys are publicly available for unrestricted use by the neuroimaging community at https://github.com/HumanBrainED/NHP-BrainExtraction.

## 1. Introduction

As the recent explosion of MRI data sharing in Nonhuman Primate (NHP) scales the amounts and diversity of data available for NHP imaging studies, researchers are having to overcome key challenges in preprocessing, which will otherwise slow the pace of progress (Autio et al., 2020a; Lepage et al., 2020; Messinger et al., 2020; Milham et al., 2018). Among them is one of the fundamental preprocessing steps - brain extraction (also referred to as skull-stripping) (Seidlitz et al., 2018; Tasserie et al., 2020; Zhao et al., 2018). In both human and NHP MRI pipelines, brain extraction is often among the first early preprocessing steps (Esteban et al., 2019; Glasser et al., 2013; Seidlitz et al., 2018; Tasserie et al., 2020; Xu et al., 2015). By removing the non-brain tissue, brain extraction dramatically improves the accuracy of later steps, such as anatomy-based brain registration, pial surface reconstruction, and cross-modality coregistration (e.g., functional MRI, diffusion MRI) (Acosta-Cabronero et al., 2008; Autio et al., 2020a; Lepage et al., 2020; Seidlitz et al., 2018). In humans, automated brain extraction tools have been developed (e.g., the Brain Extraction Tool [BET] in FSL, 3dSkullStrip in AFNI, the Hybrid Watershed Algorithm [HWA] in FreeSurfer, etc.) and easily inserted into a diversity of preprocessing pipelines (e.g. Human Connectome Project [HCP], fMRIPrep, Configurable Pipeline for the Analysis of Connectomes [C-PAC], Connectome Computational System [CCS], Data Processing & Analysis for Brain Imaging [DPABI]) (Cox, 1996; Craddock et al., 2013; Fischl, 2012; Glasser et al., 2013; Jenkinson et al., 2012; Ségonne et al., 2004; Xu et al., 2015; Yan et al., 2016). However, adaption for macaque brain extraction is significantly more challenging, as the data are often noisy due to the smaller brain and voxel sizes involved. The low signal-to-noise ratio (SNR) and strong inhomogeneity of image intensity compromise intensity-based brain extraction approaches, necessitating parameter customization to fit the macaque data (Messinger et al., 2020; Milham et al., 2018). For instance, a new option '-monkey' has been developed to customize AFNI's widely-used 3dSkullStrip function, which improves its performance for NHP data. Yet, the results are still mixed across datasets and often require further manual corrections.

In recent years, registration-based label transferring (i.e. template-driven) approaches have been proposed as a potential solution for NHP brain extraction (Jung et al., 2020; Lohmeier et al., 2019; Seidlitz et al., 2018; Tasserie et al., 2020). These approaches start by registering an individual's anatomical image to the template in order to establish the deformation between the subject-specific head and template head. Once obtained, the transform is used to bring a template-based brain mask back to the individual space, where it can be used to extract the individual brain. The performance of such approaches heavily relies on the accuracy of the transformation and whether the template is representative of the individual data. As such, factors that can compromise the appropriateness or representativeness of the template for the specific individual dataset can decrease the utility of registration-based approaches. Examples where a template may not be representative include variations in species of macaque (i.e. *M. mulatta, M. fascicularis,* etc.) or NHPs (e.g. macaque, marmoset, etc.), the field of view, age (e.g. infant, juvenile, aging), sex (e.g. thicker muscular tissue for male adult macaque), and surgical implants (e.g. with head-holder implants, or anatomical lesions). This issue might be further intensified for data from multiple study centers with different scan acquisitions and samples. In addition, non-linear registration (e.g. ANTs, 3dQwarp) for high-resolution images in brain extraction step is relatively time-consuming (e.g. over hours) and limits the computational efficiency of NHP pipelines.

1 Recognizing the continued challenges of brain extraction, researchers in human literature have begun to
2 leverage deep learning models as a potential solution. In humans, a growing number of studies have
3 demonstrated the ability of convolutional neural network (CNN) models for brain extraction, as well as
4 tissue segmentation (Henschel et al., 2020; Lyksborg et al., 2015; Rehman et al., 2018; Snehashis Roy et
5 al., 2018; Yogananda et al., 2019). Across studies, training and validation datasets in humans have included
6 hundreds, and in some cases, even thousands of datasets, to ensure accurate performance and avoid
7 overfitting. Once trained, the models have proven to be able to perform highly accurate extraction for new
8 datasets in a matter of seconds. With rare exceptions, the NHP field does not possess datasets close to the
9 multitudes used for training in humans. A recent study did, however, successfully implemented a CNN
10 model (i.e. Bayesian SegNet) for brain extraction using a relatively smaller sample in macaques collected
11 at a single site (N=50), suggesting that such training sets may not need to be as large as expected based
12 upon these preliminary human studies (Zhao et al., 2018). However, the large majority of NHP studies use
13 notably smaller sample sizes (i.e., 2-10). While combining data from multiple studies could be a solution,
14 it is important to note the NHP literature tends to have notably greater variability in imaging protocols than
15 its human counterpart.
16
17 The present work attempts to overcome the challenges at hand for NHP imaging by developing a
18 generalizable macaque brain extraction model that can handle data from previously untrained
19 protocols/sites with high accuracy. To accomplish this, we leveraged a transfer learning U-NET strategy,
20 which explicitly aims to train a model for one purpose, and then extend its utility to related problems. In
21 the present case, we trained our model on a human sample (n = 197), and then treated a nonhuman sample
22 (n = 2 for six sites) as the transfer dataset - a strategy that exploits the similarity of human and non-human
23 brain structure. Upon successful demonstration of the ability to transfer between species, we then evaluated
24 the transfer of the updated model to untrained sites in the PRIMatE Data Exchange. Finally, we improved
25 the generalizability of our model by adding a single macaque sample from each of the additional 7 sites
26 (for a total of N=20). We released our pre-trained model, code, and the brain masks outcomes via the
27 PRIMatE Resource Exchange (PRIME-RE) consortium (Messinger et al., 2020; Milham et al., 2018).
28

## 2. Methods

30
### 2.1. Human Sample
32 We made use of an open available human brain extraction sample as an initial good-standard training dataset
33 (Puccio et al., 2016). Data were collected as a part of the Enhanced Rockland Sample Neurofeedback Study
34 (N=197, 77 female, age=21-45) (McDonald et al., 2017). Anatomical images data were acquired from a 3T
35 Siemens Trio scanner using a 12 channel head matrix (T1-weighted 3D-MPRAGE sequence,
36 FOV=256x256mm$^2$, TR=2600ms, TE=3.02ms, TI=900ms, Flip angle=8°, 192 sagittal slices,
37 resolution=1x1x1 mm$^3$). T1w images were skull-stripped using a semi-automatic iterative procedure that
38 involved skull-stripping all of the data using BEaST (brain extraction based on nonlocal segmentation
39 technique) (Eskildsen et al., 2012) and manually correcting the worst results. Corrected brain masks were
40 added into the BEaST library and the procedure was repeated until the process converged. The results of
41 this procedure underwent an additional manual inspection and correction procedure to identify and fix any
42 remaining errors.
43
### 2.2. Macaque Sample

1    The MRI macaque data used in the present study are publicly available from the recent NHP data-sharing
2    consortium – the non-human PRIMate Data-Exchange (PRIME-DE) (Milham et al. 2018), which includes
3    136 macaque monkeys from 20 laboratories. We selected one anatomical T1w image per macaque in our
4    analyses. The detailed description of the data acquisition of the magnetization-prepared rapid gradient echo
5    (MPRAGE) image for each site was described in the prior study and PRIME-DE website
6    (https://fcon_1000.projects.nitrc.org/indi/indiPRIME.html, Milham et al., 2018). As the sample size is
7    relatively small in most of the sites (N≤6 for 13 out of 20 sites), we selected six sites which have no less
8    than eight macaque monkeys collected as the first dataset pool for manual edits, model training, and testing
9    (East China Normal University Chen [ecnu-chen], Institute of Neuroscience [ion], Newcastle University
10   Medical School [newcastle], University of Oxford [oxford], Stem Cell and Brain Research Institute (sbri),
11   and the University of California, Davis [ucdavis]). The details of the sample and MRI acquisition are
12   specified in Milham et al ((Milham et al., 2018), Table 1). In total, eight macaque monkeys per site were
13   selected to create the manually-edited 'ground truth' dataset to train and evaluate our model (**Macaque**
14   **Dataset I**, N=48). To optimize the generalizability of our model to other sites, we also manually edited an
15   additional seven macaques from seven sites (one per site: East China Normal University Kwok [ecnu-k],
16   Lyon Neuroscience Research Center (lyon), Mount Sinai School of Medicine - Siemens scanner
17   [mountsinai-S], National Institute of Mental Health - Messinger [nimh1], Netherlands Institute for
18   Neuroscience [nin], and Rockefeller University [rockefeller]) to create **Macaque Dataset II** (N=7). The
19   rest of the PRIME-DE macaques across all 20 sites were used as an additional hold-out testing dataset
20   (N=81). Of note, our main model was built based on the MPRAGE data. We further made use of the
21   magnetization-prepared two rapid acquisition gradient echoes (MP2RAGE) images from site-UWO-
22   MP2RAGE (N=3) to extend our model to facilitate the brain extraction for MP2RAGE data. All animal
23   procedures were conducted in compliance with the animal care and use policies of the institution where the
24   data was collected.
25

26   **2.3. Preprocessing**
27   To improve the quality and homogeneity of input data across different sites, a minimal preprocessing was
28   carried out for all anatomical images to remove the salt-and-pepper noise and correct the intensity bias.
29   Specifically, we first re-conformed all T1w images into RPI orientation and applied a spatially adaptive
30   non-local means filtering to remove the 'salt-and-pepper' noise (DenoiseImage in ANTs) (Buades et al.,
31   2011). Next, we performed the bias field correction to normalize image intensities (N4BiasFieldCorrection
32   in ANTs) (Tustison et al., 2010). The preprocessed images were served as inputs for all the brain extraction
33   approaches.
34

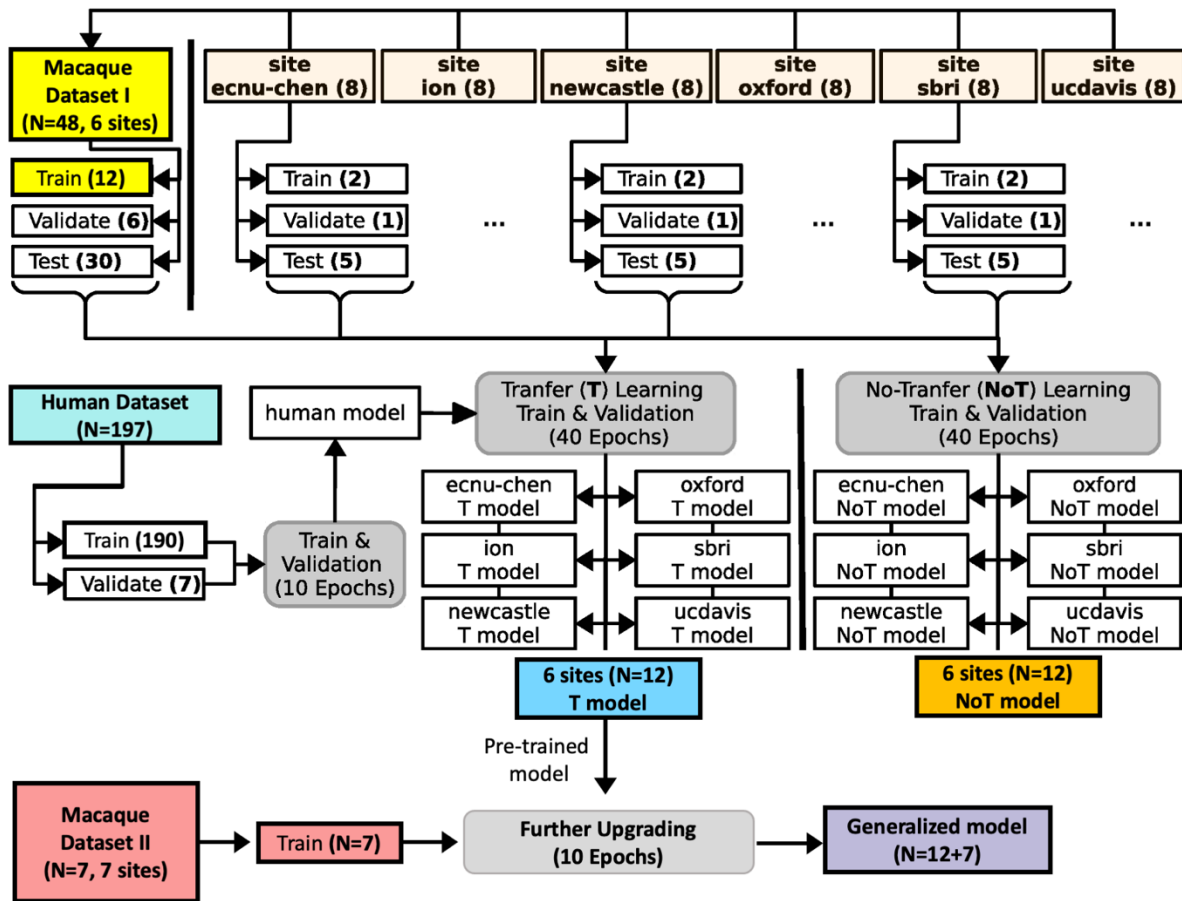35   **2.4. Traditional Methods and Manually Edited Masks**
36   To compare our deep learning models with state-of-the-art methods for brain extraction, we employed five
37   widely-used skull stripping pipelines implemented in commonly used MRI packages (AFNI, ANTs, FSL,
38   and FreeSurfer) (Avants et al., 2009; Cox, 1996; Fischl, 2012; Jenkinson et al., 2012). Specifically, we
39   tested three intensity-based approaches (FSL BET, FreeSurfer HWA, and AFNI 3dSkullStrip) and two
40   template-driven pipelines (Flirt+ANTS and AFNI @animal_warper) (Jung et al., 2020; Seidlitz et al., 2018;
41   Tustison et al., 2020). The command and parameters of intensity-based approaches were selected based on
42   the experiments and suggestions from the prior studies as follows (Xu et al., 2019; Zhao et al., 2018). 1)
43   FSL 'bet' command with a smaller fractional intensity threshold (-f 0.3) and head radius = 60 mm, 2)
44   FreeSurfer 'mri_watershed' command with default settings, and 3) AFNI '3dSkullStrip' command with

1   NHP specific option '-monkey' and shrink factor = 0.5. Template-driven approaches in both Flirt+ANTs
2   and AFNI @animal_warper pipelines were performed by first applying a linear registration to transform
3   the individual head to the template head, followed by nonlinear registration. Next, the template brain mask
4   was transformed back into the individual space to obtain the individual brain mask. Specifically, the
5   Flirt+ANTs pipeline uses 'flirt' and symmetric diffeomorphic image registration (SyN) (Avants et al.,
6   2008) for linear and nonlinear registration. Of note, we used FSL 'flirt' rather than ANTS linear registration
7   because 'flirt' is faster and performs better in our initial tests on NHP samples. AFNI @animal_warper uses
8   3dAllineate and 3dQwarp to compute affine and nonlinear alignments. The same NIMH Macaque Template
9   (NMT) was used in Flirt+ANTs and AFNI @animal_warper pipelines (Jung et al., 2020; Seidlitz et al.,
10  2018). The Macaque Dataset I and II were manually edited by well-trained experts (J.W.C, A.K, and T.X.)
11  using the best output from the traditional approaches as initial masks in ITK-SNAP
12  (http://www.itksnap.org/pmwiki/pmwiki.php)(Yushkevich et al., 2006).
13
14  **2.5. Train, Update & Evaluation Workflow for Deep Learning Models**
15  **2.5.1. Overview**
16  Figure 1 illustrates the overall analytic flow chart of the present study. First, we established a skull-stripping
17  model using the human dataset, aiming to provide an initial pre-trained model to facilitate the transfer-
18  learning from humans to macaques. The human NKI-RS dataset (N=197) was split into a training set
19  (N=190) and a validation set (N=7). We used the training set to train the U-Net model for 10 epochs (i.e.
20  full training dataset has passed through the complete neural network 10 times) and selected the best epoch
21  as the human model based on the performance (i.e. dice coefficient) in the validating set. Next, we
22  transferred the pre-trained human model to build the macaque model using Macaque Dataset I. Specifically,
23  for each of six sites, we randomly selected 2 macaque monkeys as the training set, 1 macaque as the
24  validation set and 5 macaques as the testing set. We also used all macaques across six sites in Macaque
25  Dataset I to create the merged training (N=12), validating (N=6), and testing (N=30) sets. Transfer-learning
26  from human to macaque was carried out for each site as well as the merged data. We calculated 40 transfer-
27  training epochs and selected the epoch that had the best performance in the validation set as the transfer-
28  learning model for each site and the merged samples (refer to the U-Net T model). We also created macaque
29  models that were trained only on the macaque data using the same training (N=12) and validation (N=6)
30  data for each of the six sites and the merged samples (refer to the non-transfer-learning model, i.e. U-Net
31  NoT model). We evaluated and compared the performance between the site-specific T model and the NoT
32  model in testing sets. We also compared the U-Net T and NoT model to the traditional pipelines using the
33  held out test set from Macaque Dataset I. To improve the generalizability of the U-Net model to fit more
34  macaque data from other sites, we further upgraded the U-Net transfer-learning model using both Macaque
35  Dataset I (N=12) & II (N=7) to generate the final generalized model (referred to as generalized U-Net 12+7
36  model). To evaluate the model performance, we applied the U-Net transfer-learning model, generalized
37  12+7 model, and traditional pipelines to all the T1-weighted images from PRIME-DE (136 macaques).
38  Expert ratings (details in the Model Evaluation section) were conducted to evaluate whether brain extraction
39  was successful.
40

**Figure 1. Schematic of the U-Net model training, transfer-learning, and validation on human and macaque datasets.** The human dataset was used to generate a pre-train U-Net model for transfer learning. Macaque Dataset I was adopted to train and evaluate the models with and without transfer-learning from the human dataset to small samples of macaque datasets. Further, Macaque Dataset II was added to the training set to generate a generalized model for brain extraction on macaque data.

### 2.5.2. Neural Network Model (U-Net)

We used a convolutional neural network (CNN) model, U-Net (Ronneberger et al., 2015) for brain extraction. The model was built using an open resource machine learning package (PyTorch: https://pytorch.org). Briefly, the preprocessed 3D T1w images were first resampled into slides along axial, sagittal, and coronal planes. The U-Net model predicts the brain tissue in each slice and then merges all slices to obtain a 3D brain mask. Here, we focused on the architecture of the U-Net model (Fig. S1) and illustrated the details of how the U-Net model identifies the whole brain tissue in the training, validation, and testing processes in the next section.

As shown in Fig. S1, the U-Net model consists of a contraction (i.e. encoding) and an expansion (i.e. decoding) path; each includes five convolution units (Ronneberger et al., 2015). Of note, for a given slice resampled from the T1 images, the input of the U-Net model also included its neighboring slices (i.e. 3 slices in total) as an input (dimension: 3x256x256 blocks). Next, for each convolution unit, two 3×3 convolution layers are built and each is followed by a batch normalization and a leaky rectified linear

1  (ReLU) operation (Fig. S1, blue arrow). In the encoding step, a 2×2 max pooling with stride 2 was adopted
2  for down-sampling data from the upper unit to the lower unit. We used 16 feature channels in the initial
3  unit, and doubled every unit. The expansive path consists of four up-convolution units. For each up-
4  convolution unit, a 2×2 up-convolution and ReLU operation were applied to the lower unit to yield the
5  feature map for the current unit. This feature map was then concatenated with the feature map at the same
6  level in the contracting path to generate the combined feature map. Similar 3×3 convolution layers with
7  batch normalization and ReLU operation were then performed on the feature maps at each up-convolution
8  unit. Next, we employed a 1×1 convolution layer at the upper un-convolution unit to map the final feature
9  maps to a two-classes map. Finally, a SoftMax layer was used to obtain the probability map for brain tissue.
10 Of note, the initial weights of the convolution and up-convolution layers were randomly selected using a
11 Gaussian distribution $N(0, 0.2)$, and the initial bias of layers was set to 0. For the transfer-learning and
12 model-upgrading model, the parameters from the pre-trained model were used in the initial setting.
13

14 **2.5.3. Model Training and Validating Procedure**
15 In Figure S2, we illustrated the training procedure of how 3D T1w data was processed for the U-Net model.
16 First, we normalized the intensity of the preprocessed T1w image in the range 0 to 1 across all voxels. We
17 then resampled the T1w volume to a 3D intensity matrix where the highest sampled dimension of T1w
18 volume was forced to be rescaled to 256. In the example shown in Fig. S2, T1w (voxel size: 176x176x96)
19 were rescaled to a 256x256x140 matrix. Next, for each slice along the axial, sagittal, and coronal direction,
20 we generated a 3-slice block; the slice and its two neighboring slices. As a result, we obtained 254, 254,
21 and 138 blocks for axial, coronal, and sagittal directions respectively. Next, we conformed each of the 3-
22 slice blocks into a 3x256x256 matrix. When the dimension of the slice plane was less than 256, we filled
23 the matrix with 0. In total, 646 conformed blocks were generated for the T1w image. Similarly, we
24 processed the manually edited brain mask of the T1w image and generated 646 corresponding blocks. After
25 that, we used the above U-Net model to estimate the probability of brain tissue for each T1w block and
26 calculated the cross-entropy (PyTorch function: CrossEntropyLoss) between the probability map and the
27 'ground truth' map (Ketkar, 2017) as the model cost. A stochastic optimization (learning rate=0.0001, batch
28 size=20) was then used for backpropagation (Kingma and Ba, 2014).
29

30 To evaluate and select the model from the training epochs, we used the probability map generated from the
31 U-Net model to create the final predicted brain mask for each epoch, and then examined the Dice coefficient
32 between the predicted mask and 'ground truth' mask (Fig. S3). Specifically, we processed the validation
33 T1w images into 3-slices blocks following the above procedure. Next, we used the U-Net model at each
34 training epoch to estimate the probability map for each block. All the probability maps (646 blocks) were
35 then combined along the axial, sagittal, and coronal direction and yielded an averaged 3D probability matrix
36 (256x256x256). After that, we rescaled and cropped the matrix back to the original voxel dimension to
37 create a probability volume for the given T1w image. Finally, we thresholded (>0.5) this probability volume
38 to obtain the predicted brain mask. The Dice coefficient between the predicted mask and the 'ground truth'
39 mask was computed for each epoch during training. The epoch which showed the highest Dice coefficient
40 was then selected as our final model.
41

42 **2.6. Model Evaluation**
43 We carried out a quantitative examination in the testing set of Macaque Dataset I and evaluated the degree
44 to which methods provided more similar brain masks as compared to the manually edited 'ground truth'.

1 Specifically, we calculated Dice coefficients (Dice) (Sørensen, 1948) between the predicted mask and the
2 manually edited 'ground truth' using the following equation:

$$Dice = \frac{2|P \cap T|}{|P+T|},$$

4 Where |...| represents the total number of voxels in a mask; $P$ and $T$ are predicted and 'ground truth' masks,
5 respectively. Dice equals 1 when the predicted mask and 'ground truth' mask are perfectly overlapped. In
6 addition, we also calculated the voxel-wise false-positive (FP), and false-negative (FN) rate to examine
7 where the predicted mask falsely includes the non-brain tissue (higher FP) or misses any brain tissue (higher
8 FN). For each voxel in a given macaque, we tested whether the voxel is falsely assigned as brain tissue
9 (FP=1) or falsely assigned as non-brain tissue (FN=1) in individual space. Next, we transferred the
10 individual FP and FN map to the NMT space by applying affine and warp transforms. Linear affine here
11 was created by aligning the manually skull-stripped brain to the NMT brain ('flirt') and nonlinear warp was
12 generated by registering individual head to the NMT head ('fnirt'). The FP and FN maps were averaged
13 across macaques for each brain extraction approach in the NMT space.

14

15 We also employed a qualitative evaluation to compare the success rate across different brain extraction
16 approaches for PRIME-DE data without the 'ground truth'. Three experts with rich experience in imaging
17 quality control visually rated the brain masks (J.W.C., X.L., and T.X.). First, each expert independently
18 reviewed all the images and rated them with four grades, i.e. 'good', 'fair' (the brain tissue was identified
19 with a slightly inaccurate prediction at the edge), 'poor' (i.e. most of the brain regions are identified but
20 with significant errors, in particular missing brain tissue at the edge), and 'bad'. 'Poor' and 'bad' scores
21 were recognized as failed. The brain masks with inconsistent ratings among experts were reviewed and
22 discussed for a final consensus rating.

23

24 ## 3. Results

25

26 ### 3.1. The Convergence of Human U-Net Model
27 To acquire a pre-trained model for macaque samples, we first trained a U-Net model using the human
28 samples. Fig. S4 demonstrates the sum of loss on the training set and the mean Dice coefficients across
29 macaques on the validation set for each epoch. After the first epoch, the loss decreased steeply and the mean
30 Dice coefficient reached above 0.985. After that, the mean Dice coefficient gradually improved, showing
31 its highest value (0.9916±0.0012) after the 9th epoch. To avoid over-fitting the human samples for the
32 subsequent macaque training, we only carried out 10 epochs and selected the model after the 9th epoch as
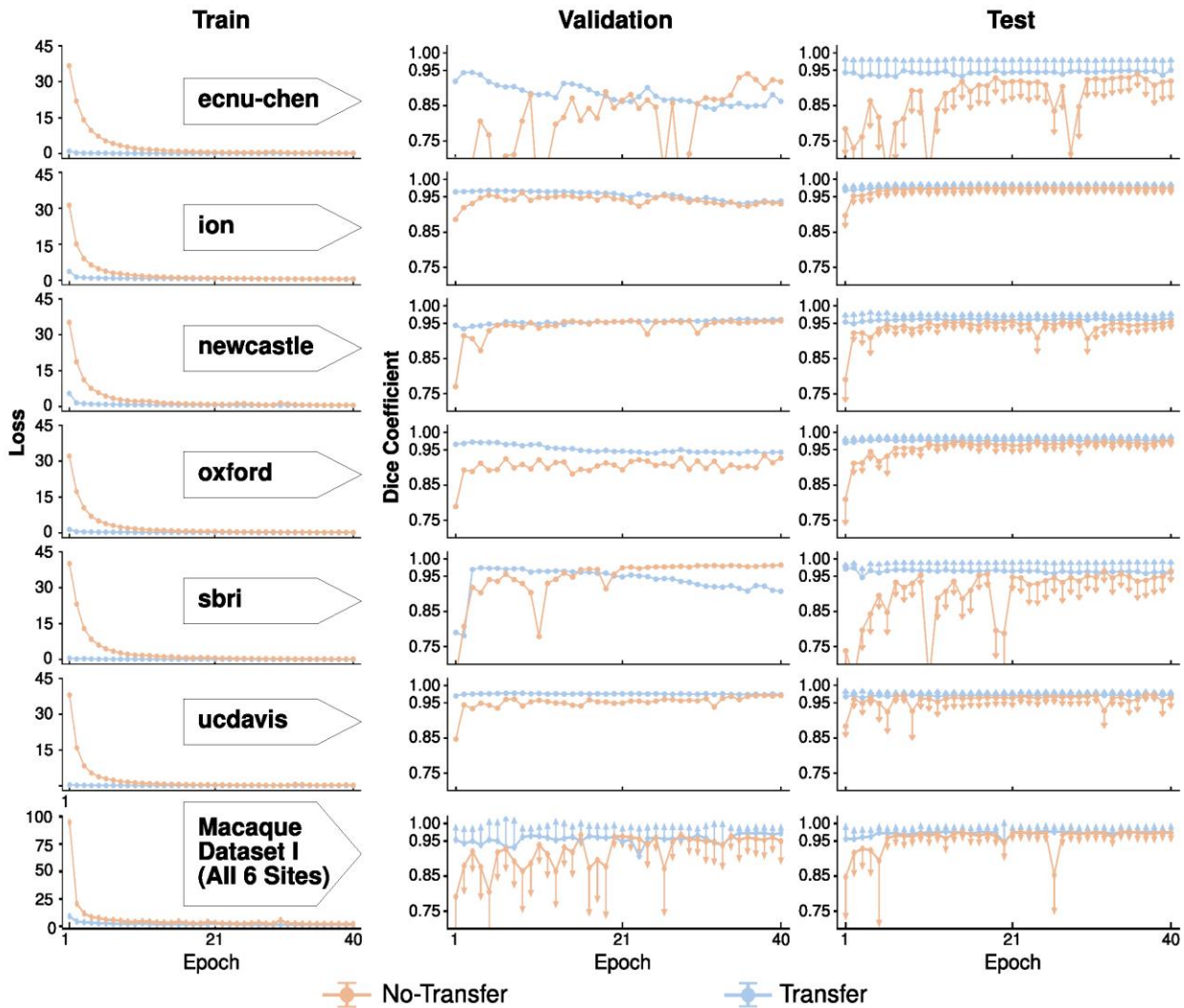33 the pre-trained model for transfer-learning to the macaque samples.

34

35 ### 3.2. Comparison Between Models With and Without Transfer-learning
36 We first evaluated the site-specific models with and without transfer-learning on Macaque Dataset I. For
37 each site and for the merged sample across the six sites, the loss converged faster for the transfer-learning
38 model than the model without transfer-learning (Fig. 2, left); the loss had nearly reached its minimum after
39 2 epochs. In addition, the Dice coefficients were more stable for transfer-learning models in the validation
40 sets (Fig. 2, center). Of note, the mean Dice coefficients slowly dropped after 20 epochs in the validation
41 set across three sites (i.e. ecnu-chen, ion, and sbri), which may reflect overfitting of the model selection on
42 small training samples (Fig. 2, center). Nevertheless, the mean Dice coefficients still remained stable in the
43 testing set after the first epoch (Fig. 2, right). In addition, compared to the model trained solely on macaque

1    samples, the transfer-learning model also showed lower variation across testing macaques, suggesting its
2    potential generalizability across macaque samples (Fig. 2, right).
3
4    We also computed and evaluated the model with and without transfer-learning based on the merged training
5    samples across six sites (N=12). Similarly, the transfer-learning model (referred to as U-Net T12 model)
6    showed lower loss and higher Dice coefficients across 40 epochs than the model without transfer-learning
7    (Fig. 2, last row); its best performance epoch (i.e. the 37th epoch) was selected as our macaque transfer-
8    learning model (i.e. U-Net T12).
9



10
11   **Figure 2. Comparison of skull-stripping performance of the U-Net models with and without transfer-**
12   **learning on Macaque Dataset I.** The models with transfer-learning (blue curve) outperform the one
13   without (orange curve) across each of six sites (row 1-6), as well as the merged sample (the last row). The
14   loss (i.e. the sum of cross-entropy, first column) between the predicted mask and ground truth mask on the
15   train set converges faster for the model with transfer-learning. Similarly, models with transfer-learning
16   show higher Dice coefficients with lower variation in the validation and testing sets than models without
17   transfer-learning. Of note, one-side error bars (standard deviation) were used to avoid dense overlaps
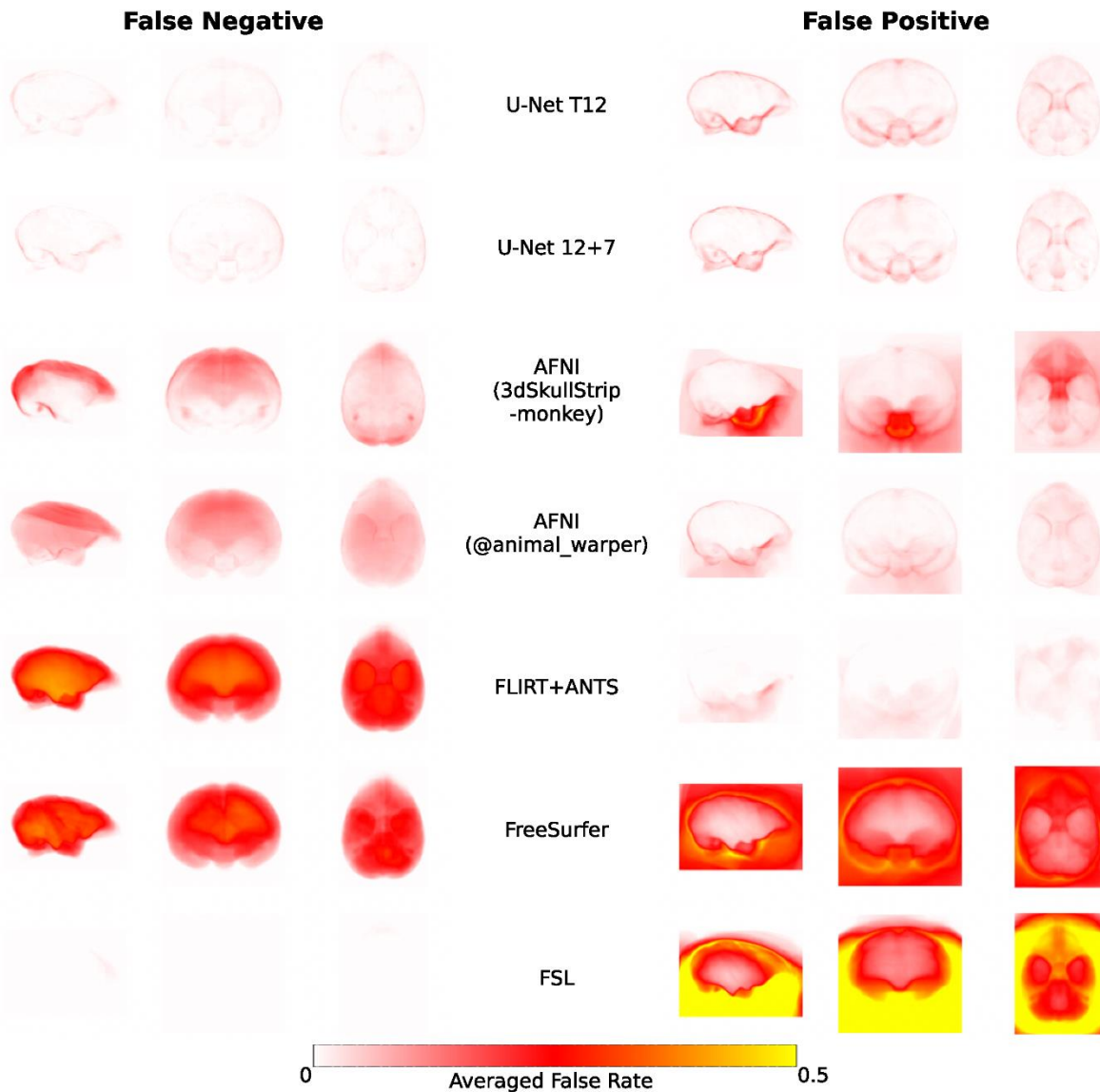18   between two models.
19

1  **3.3. Comparison Between the U-Net Model and Traditional Approaches**
2  Recognizing that the model with transfer-learning is superior to the one without, we then updated the
3  transfer-learning model (i.e. U-Net T12) using Macaque Dataset I (N=12) & II (N=7) to generate our final
4  generalized model (U-Net 12+7 model). Here, we evaluated the performance of the U-Net T12 model and
5  the U-Net 12+7 model to other traditional brain extraction approaches. Brain masks from the two U-Net
6  models showed significantly higher Dice coefficients than those from the traditional pipelines (Fig. 3,
7  $F$=38.531, $p<10^{-30}$ repeated ANOVA, all post-hoc p<0.05). Skull-stripping using the U-Net models was
8  successful (Dice>0.95) for all of the testing macaques (N=30) across six sites. Notably, the U-Net 12+7
9  model showed relatively higher Dice coefficients than the U-Net T12 model (Fig. 3, *paired-t*=3.62,
10 *p*=0.001), though the additional training samples used in U-Net 12+7 model were not included in the sites
11 where the testing samples were selected from. This indicated the generalizability of model-upgrading across
12 sites. At the voxel level, both U-Net T12 and '12+7' models exhibited fewer false negatives and false
13 positives than traditional pipelines (Fig. 4). The UNet T12 model showed slightly more false positives than
14 false negatives, which indicated that the model tends to include a few non-brain voxels on the edge of the
15 brain mask rather than miss the brain tissue. Overall, these results demonstrated the feasibility of transfer-
16 learning and model-upgrading using small training samples.
17
18 Comparing amongst the traditional approaches, the template-driven approaches (i.e. AFNI
19 @animal_warper and Flirt+ANTS), and AFNI 3dSkullStrip with parameters customized for NHP data
20 (3dSkullStrip -monkey) showed better performance than FSL and FreeSurfer (Fig. 3-4). Both template-
21 driven approaches appear to be more conservative and have missed the brain tissue (low false positives and
22 high false negatives). AFNI 3dSkullStrip missed identifying the brain in the superior regions (higher false
23 negatives on the top) and falsely included non-brain voxels in the bottom (Fig. 4).
24



25

1 **Figure 3. Performance of the U-Net models and traditional approaches.** The boxplot shows the Dice
2 coefficients in the testing datasets (N=30) of Macaque Dataset I across brain extraction approaches
3 including 1) the transfer-learning model (i.e. U-Net T12), 2) generalized model (i.e. U-Net 12+7), 3) AFNI
4 3dSkullStrip command with '-monkey' option, 4) AFNI @animal_warper pipeline, 5) template-driven
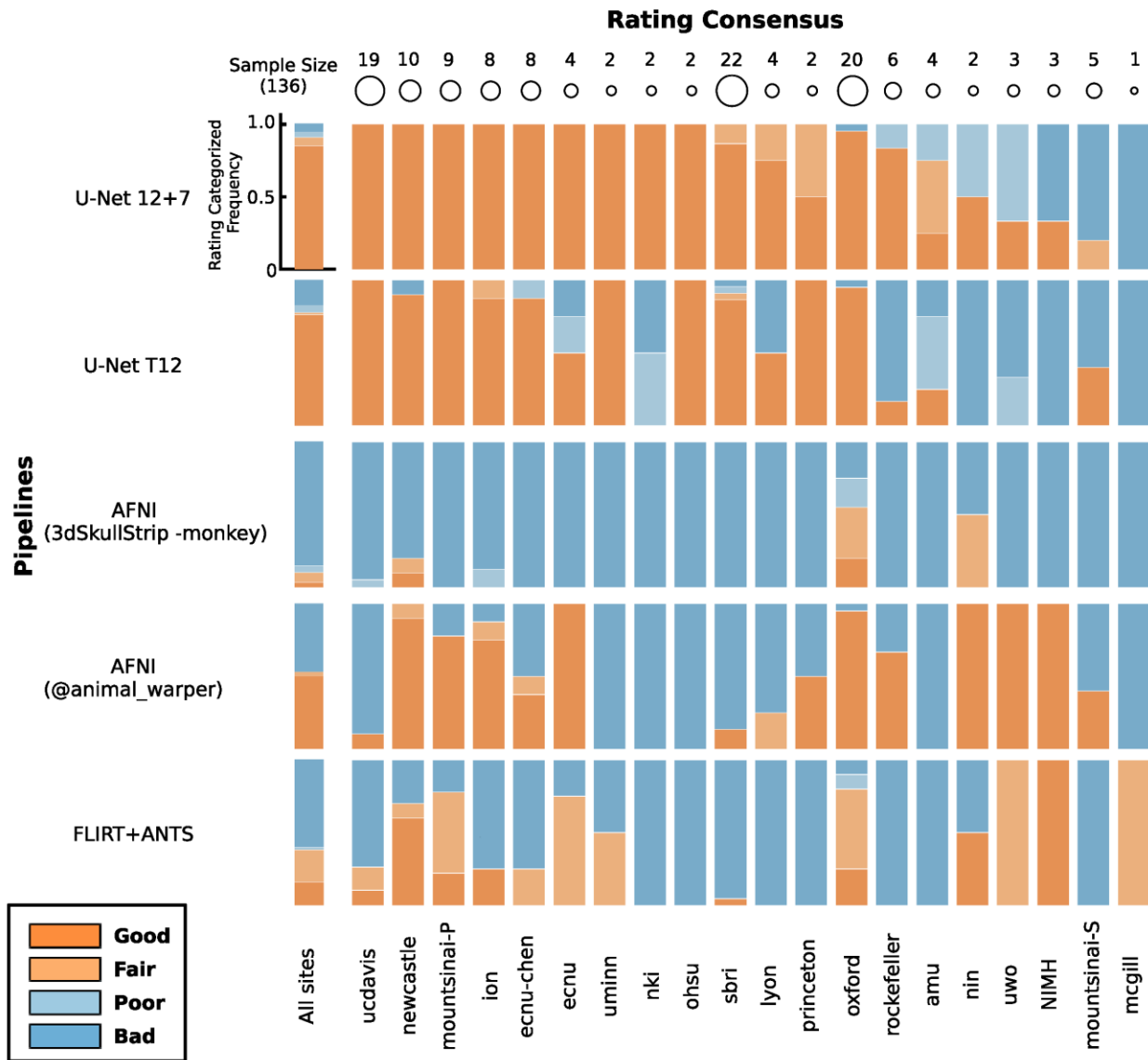5 FLIRT+ANTS pipeline, 6) FreeSurfer HWA approach, and 7) FSL BET approach.
6



7
8 **Figure 4. The averaged voxel-wise false-negative and false-positive rates of the U-Net models and**
9 **traditional approaches in the testing datasets of Macaque Dataset I.** The false-negative rate examines
10 where the predicted mask falsely misses the brain tissue in the brain (left column) while the false-positive
11 rate examines where the predicted mask falsely includes the non-brain tissue. Of note, FSL BET tends to
12 falsely include large amounts of non-brain tissue in the brain mask, thus false-negative rates are close to
13 zeros.
14

15 **3.4. Generalizability of U-Net model and Skull-Stripping PRIME-DE Samples**
16 To evaluate the generalizability of the U-Net model, we applied the U-Net T12 and 12+7 models to all the
17 other macaque samples (i.e. T1w images) contributed to PRIME-DE from differing research sites. The

1  results of brain masks were visually reviewed by three experts and rated into four grades (good, fair, poor,
2  and bad). Good and fair ratings were considered to be successes, while poor and bad ratings were failures.
3  We also used the traditional pipelines and similarly rated their brain mask results. The performance of
4  different       pipelines       for       each       macaque       is       shown       on       github       repository:
5  https://github.com/HumanBrainED/NHP-BrainExtraction. Figure 5 shows the proportion of macaques with
6  good, fair, poor, and bad skull-stripping masks for each of the 20 sites, for each of the five approaches.
7  Again, the U-Net models outperformed the traditional approaches for most sites. In particular, the final U-
8  Net 12+7 model showed the best performance across macaque samples (successful rate=90.4%). All
9  macaque samples (N=123) were successfully skull-stripped across the twelve sites with thirteen exceptions.
10  This result demonstrated the generalizability of the U-Net 12+7 model across sites.
11



**Figure 5. Expert rating consensus of brain extraction performance across the U-Net models and traditional approaches for the PRIME-DE datasets.** The stacked bar plots show the proportion of macaques with good, fair, poor, and bad skull-stripping masks on T1w images for each of 20 sites in PRIME-DE. Of note, FreeSurfer HWA and FSL BET default pipelines failed to obtain fair/good brain

masks, and thus their ratings were not displayed. The U-Net 12+7 model (the first row) shows the highest rating across pipelines.

We also noticed that template-driven approaches performed well in some cases, but failed in others (successful rate: AFNI @animal_warper=52.9%, Flirt+ANTS=39.7%). Visual inspection of the intermediate outputs from AFNI @animal_warper and Flirt+ANTS pipelines showed that all successful cases had the first linear alignment of the individual's head to the template head somewhat close. The failures mostly occurred in the first linear registration step; 79.4% of 63 skull-stripping failures from @animal_warper and 89.3% of 84 failures from Flirt+ANTS pipeline had failed in the first linear registration step.

For the failed 13 macaques in the U-Net 12+7 model, when we used the U-Net output mask to perform an initial brain extraction and estimated the affine transformation from the initial skull-stripped brian to the template brain, the template-driven approaches (AFNI @animal_warper and ANTS) turned seven failed cases into successes. We have selected the best successful mask for each macaque and shared all brain masks on PRIME-RE (Messinger et al., 2020).

**3.5. Application of U-Net Model-Upgrading for an External Dataset and MP2RAGE Images**
Here we used a large external dataset (N=454, UW-Madison dataset with manually drawn whole-brain masks, details described in previous studies) to evaluate the utility of our U-Net -based approach (Fox et al., 2018, 2015; Oler et al., 2010). First, we directly applied the model-prediction module using our U-Net 12+7 model to extract the brain masks for 40 macaques and found relatively good results (Dice=0.923+/-0.025). We further used the U-Net 12+7 model as the pre-trained model and upgraded the U-Net model on these 40 macaques for 10 epochs. This procedure took about 90 min. The upgraded model showed substantial improvement and successfully skull-stripped all remaining macaques (N=414, All Dice>0.95, Mean=0.977±0.005). We further challenged our U-Net model-upgrading module for MP2RAGE images collected from UWO site (N=8) in PRIME-DE, which have different intensity profiles with opposite contrast in gray matter and white matter (Fig. S5). By upgrading the pre-trained U-Net 12+7 model with three MP2RAGE brain masks, the U-Net model enabled skull-stripping on the rest of the hold out MP2RAGE images (N=5).

# 4. Discussion
The present work demonstrated the feasibility of developing a generalizable brain extraction model for NHP imaging through transfer learning. Central to the success of our effort was our leveraging of human data as a base training set, upon which additional learning for NHP and site-related variations could be readily achieved. We employed the heterogenous, multi-site PRIME-DE resource to evaluate the effectiveness of our framework, finding that the transfer-learning U-Net model identified macaque brain tissue with higher accuracy than traditional methods and also proved more generalizable across data collections. Notably, the transfer-learning U-Net model provides a notably faster solution (approximately 1-10 min) than the next best performing algorithms, which tend to rely on template-based strategies and require over an hour. We have released our model and code, including the utilities for skull-stripping, model-training, transfer-learning, as well as model-updating modules (https://github.com/HumanBrainED/NHP-BrainExtraction). Additionally, we created and shared the skull-

1   stripped repository of 136 macaque monkeys to facilitate the large-scale macaque MRI imaging for the
2   PRIME-DE and NHP community (Messinger et al., 2020; Milham et al., 2018).
3
4   *A priori*, the major roadblocks that one would anticipate for the application of deep learning in NHP imaging
5   are the small sample sizes and variations in imaging protocols (Autio et al., 2020b; Milham et al., 2018). In
6   part, this is a reflection of the experiences of the human imaging community, where large sample sizes have
7   been required to accurately segment the brain using deep learning (Henschel et al., 2020; Kleesiek et al.,
8   2016). The superiority of the transfer learning models in the present work emphasizes the unique advantages
9   of transfer learning in overcoming such challenges for nonhuman primate imaging and offers a model that
10  may be considered in future efforts to overcome similar obstacles in other imaging populations (e.g.,
11  macaques with surgical implants and/or in-dwelling electrodes, brain extraction in other species, pathologic
12  models, early development, aging) (Pontes-Filho et al., 2019; S. Roy et al., 2018; Salehi et al., 2018). The
13  success of the transfer-learning model emphasizes the similarity of the general structure of brain tissue in
14  both species (e.g. gray matter, white matter) - despite the anatomical differences between human and
15  macaque heads (e.g. head size, muscular tissue surrounding the skull, skull thickness, etc.) (Yosinski et al.,
16  2014). Aside from accuracy, the transfer-learning model also converged faster (leveling off after 2 epochs)
17  and yielded more stable results across epochs, even though the pre-trained model was established using a
18  different species.
19
20  The success of transfer learning in the present work may signal the ability to use smaller samples than
21  previously employed for human imaging studies (Ghafoorian et al., 2017). However, this is not necessarily
22  the case, as this may instead reflect that the folded surface of the macaque is much less complex than that
23  of humans. The macaque central sulcus is less meandering on the lateral parietal lobe (Hopkins et al., 2014).
24  There is only one superior temporal sulcus in the temporal lobe and two less curved sulci (i.e. rectus,
25  arcuate) in the frontal lobe, which results in relatively smooth surface edges for brain extraction (Bogart et
26  al., 2012; Hopkins, 2018). In addition, such folded and meandering brain morphology is substantially more
27  similar across individual macaques than humans (Croxson et al., 2018). As such, in comparison to human
28  brain extraction, the macaque model is feasible using small training samples. Future work in the human
29  imaging community would benefit from a systematic examination of minimal sample sizes needed for
30  successful training and generalization.
31
32  Beyond transfer across species, a key finding of the present work is the ability to improve model
33  generalizability across independent imaging sites relative to traditional methods. By further upgrading the
34  pre-trained transfer-learning model based on the secondary training samples across multiple datasets, the
35  upgraded U-Net model has improved the brain extraction performance and showed a higher successful rate
36  than the traditional methods across multiple sites. Of note, the upgraded model enabled successfully skull-
37  stripping datasets acquired from three additional sites that were not included in any training sites. This
38  demonstrates the out-of-site generalizability of upgrading the pre-trained U-Net model across sites. More
39  impressively, we found that the model-upgrading module offers a solution of generalizing the pre-trained
40  model to other modalities (i.e. MP2RAGE), which is usually difficult to achieve with traditional methods.
41  These findings highlight the important role of pre-trained models in brain extraction for small samples,
42  regardless of site (across or within sites), modality (MPRAGE or MP2RAGE), and species (across or within
43  species). Further improvement for the user-specific dataset can be achieved by adding small samples (N≥2)
44  using our U-Net result as the pre-trained model.

It is worth noting that the U-Net model we used in the current study is a 2D convolutional neural network. Although the 3D model usually tends to have higher accuracy (Hwang et al., 2019), the 2D model has a smaller network size, much less memory cost, and requires less computational time. More importantly, the 2D model is more accessible in general computational platforms for users without large amounts of video memory, which can be costly. In addition, unlike the single slice 2D model, our 2D model leveraged the local interslice features by using each slice and its neighboring slices (i.e. 3 slices in total as opposed to just one) as the input image in the first layer. We also resampled the slices in axial, sagittal, and coronal planes and combined the predictions from all three planes into a final probability map (Lyksborg et al., 2015). By doing so, we effectively tripled the number of training slices, which is especially useful for optimizing prediction when training sample sizes are limited (i.e. data augmentation technique). In future work, a conditional random field can be considered in the prediction layer for further refining the weights in the tissue prediction (Chen et al., 2019; Zhao et al., 2018). Beyond improving brain extraction, future efforts may place a greater focus on tissue classification (e.g. GM, WM, subcortical structures). Central to such efforts will be the sharing and amassing of manually segmented brain images, to which 3D CNNs can be applied.

To promote pipeline development in the NHP field, we have released the skull-stripped brain masks, our generalized model, and code via the PRIMate Resource Exchange platform (PRIME-RE: https://prime-re.github.io) (Messinger et al., 2020). Researchers can access the code and perform the brain extraction on their own macaque datasets. The model-prediction for a dataset takes about 20 seconds on a GTX1070 GPU with 700MB GPU memory or 2-10 min on a single CPU core. We also included the model-building and model-upgrading modules in the code, which has been implemented in a recent version of a Configurable Pipeline for the Analysis of Connectomes (C-PAC v1.6: https://github.com/FCP-INDI/C-PAC/releases/tag/v1.6.0). When the results from the current model need improvement, users can upgrade the U-Net using our current generalized model as a pre-trained model and expect a stabilized solution after several training epochs (N<10 epochs) without validation datasets. The model-upgrading module takes about 1-5 hours on a single CPU, or 20 min on a GPU. In addition, we released our manually skull-stripped masks (40 macaques across 7 sites) which can be used as 'gold standards' for other deep-learning algorithms. Additionally, we included the successful brain mask outputs to facilitate the further preprocessing analysis of PRIME-DE data.

There are some limitations in our studies. First, although the final U-Net model showed better generalizability across research centers than traditional approaches, it is unable to accurately skull strip all macaque datasets (success rate: 90.4%). This is probably caused by notable variation of macaque samples, where the majority of failed skull strips came from. Of note, even in failed cases, it is possible to use the output of the U-Net model as the initial brain mask to create a prior linear transformation for traditional template-driven approaches - a process that can turn most of the failed cases into success. Second, the U-Net model requires denoising and bias correction using traditional approaches prior to the model prediction. Future work may consider leaving this image noise and the bias field information in the image during training of the network to simplify the processing steps and possibly improve performance.

## 5. Conclusion

In the present work, we proposed and evaluated a fast and stable U-Net based pipeline for brain extraction that exhibited performance superior to traditional approaches in a heterogenous, multisite NHP dataset. We have released the code for brain mask prediction, model-building, and model-updating, as well as macaque brain masks of PRIME-DE data. We hope this open repository of code and brain masks can promote pipeline development in the NHP imaging field and accelerate the pace of investigations.

## Acknowledgement

# Reference:

Acosta-Cabronero, J., Williams, G.B., Pereira, J.M.S., Pengas, G., Nestor, P.J., 2008. The impact of skull-stripping and radio-frequency bias correction on grey-matter segmentation for voxel-based morphometry. Neuroimage 39, 1654–1665.

Autio, J.A., Glasser, M.F., Ose, T., Donahue, C.J., Bastiani, M., Ohno, M., Kawabata, Y., Urushibata, Y., Murata, K., Nishigori, K., Yamaguchi, M., Hori, Y., Yoshida, A., Go, Y., Coalson, T.S., Jbabdi, S., Sotiropoulos, S.N., Kennedy, H., Smith, S., Van Essen, D.C., Hayashi, T., 2020a. Towards HCP-Style macaque connectomes: 24-Channel 3T multi-array coil, MRI sequences and preprocessing. Neuroimage 215, 116800.

Autio, J.A., Zhu, Q., Li, X., Glasser, M.F., Schwiedrzik, C.M., Fair, D.A., Zimmermann, J., Yacoub, E., Menon, R.S., Van Essen, D.C., Hayashi, T., Russ, B., Vanduffel, W., 2020b. Minimal Specifications for Non-Human Primate MRI: Challenges in Standardizing and Harmonizing Data Collection. arXiv [q-bio.NC].

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, 26–41.

Avants, B.B., Tustison, N., Song, G., 2009. Advanced normalization tools (ANTS). Insight J. 2, 1–35.

Bogart, S.L., Mangin, J.-F., Schapiro, S.J., Reamer, L., Bennett, A.J., Pierre, P.J., Hopkins, W.D., 2012. Cortical sulci asymmetries in chimpanzees and macaques: a new look at an old idea. Neuroimage 61, 533–541.

Buades, A., Coll, B., Morel, J.-M., 2011. Non-local means denoising. Image Processing On Line 1, 208–212.

Chen, W., Liu, B., Peng, S., Sun, J., Qiao, X., 2019. S3D-UNet: Separable 3D U-Net for Brain Tumor Segmentation. Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries. https://doi.org/10.1007/978-3-030-11726-9_32

Cox, R.W., 1996. AFNI: software for analysis and visualization of functional magnetic resonance neuroimages. Comput. Biomed. Res. 29, 162–173.

Craddock, C., Sikka, S., Cheung, B., Khanuja, R., Ghosh, S.S., Yan, C., Li, Q., Lurie, D., Vogelstein, J., Burns, R., Others, 2013. Towards automated analysis of connectomes: The configurable pipeline for the analysis of connectomes (c-pac). Front. Neuroinform. 42.

Croxson, P.L., Forkel, S.J., Cerliani, L., Thiebaut de Schotten, M., 2018. Structural Variability Across the Primate Brain: A Cross-Species Comparison. Cereb. Cortex 28, 3829–3841.

Eskildsen, S.F., Coupé, P., Fonov, V., Manjón, J.V., Leung, K.K., Guizard, N., Wassef, S.N., Østergaard, L.R., Collins, D.L., Alzheimer's Disease Neuroimaging Initiative, 2012. BEaST: brain extraction based on nonlocal segmentation technique. Neuroimage 59, 2362–2373.

Esteban, O., Markiewicz, C.J., Blair, R.W., Moodie, C.A., Isik, A.I., Erramuzpe, A., Kent, J.D., Goncalves, M., DuPre, E., Snyder, M., Oya, H., Ghosh, S.S., Wright, J., Durnez, J., Poldrack, R.A., Gorgolewski, K.J., 2019. fMRIPrep: a robust preprocessing pipeline for functional MRI. Nat. Methods 16, 111–116.

Fischl, B., 2012. FreeSurfer. Neuroimage 62, 774–781.

Fox, A.S., Oler, J.A., Birn, R.M., Shackman, A.J., Alexander, A.L., Kalin, N.H., 2018. Functional Connectivity within the Primate Extended Amygdala Is Heritable and Associated with Early-Life Anxious Temperament. J. Neurosci. 38, 7611–7621.

Fox, A.S., Oler, J.A., Shackman, A.J., Shelton, S.E., Raveendran, M., McKay, D.R., Converse, A.K., Alexander, A., Davidson, R.J., Blangero, J., Rogers, J., Kalin, N.H., 2015. Intergenerational neural mediators of early-life anxious temperament. Proc. Natl. Acad. Sci. U. S. A. 112, 9118–9122.

Ghafoorian, M., Mehrtash, A., Kapur, T., Karssemeijer, N., Marchiori, E., Pesteie, M., Guttmann, C.R.G., de Leeuw, F.-E., Tempany, C.M., van Ginneken, B., Fedorov, A., Abolmaesumi, P., Platel, B.,
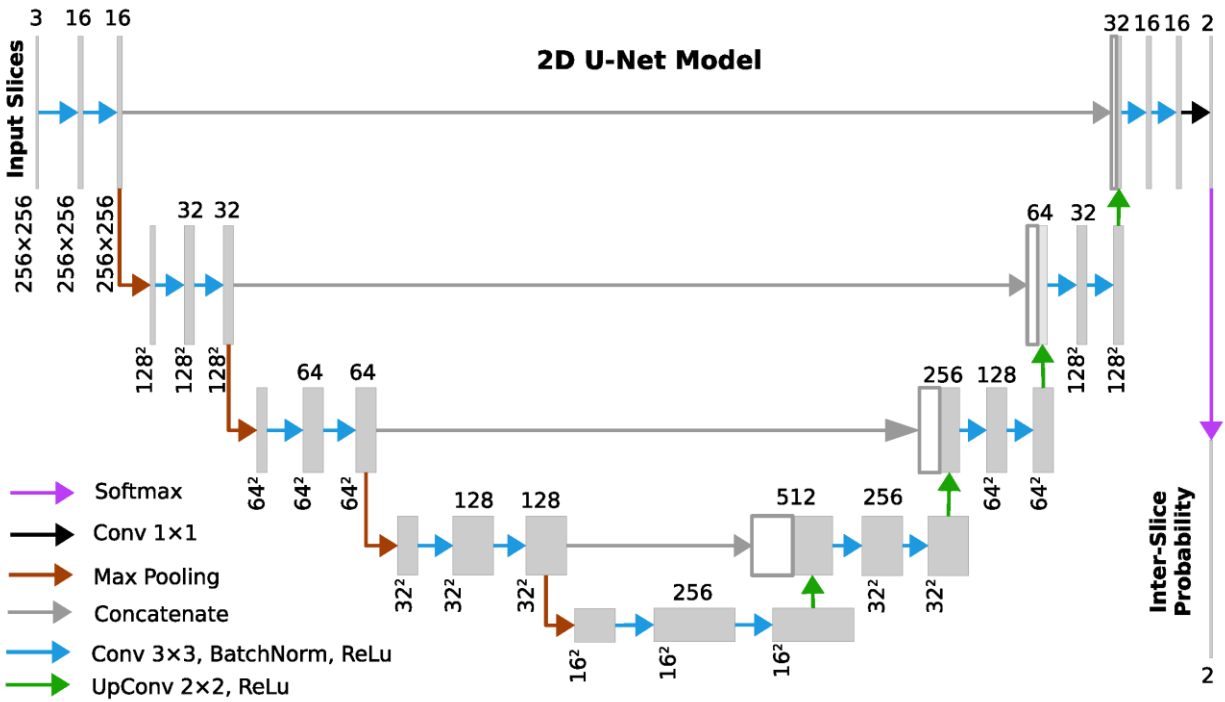
Wells, W.M., 2017. Transfer Learning for Domain Adaptation in MRI: Application in Brain Lesion Segmentation. Medical Image Computing and Computer Assisted Intervention − MICCAI 2017. https://doi.org/10.1007/978-3-319-66179-7_59

Glasser, M.F., Sotiropoulos, S.N., Wilson, J.A., Coalson, T.S., Fischl, B., Andersson, J.L., Xu, J., Jbabdi, S., Webster, M., Polimeni, J.R., Van Essen, D.C., Jenkinson, M., WU-Minn HCP Consortium, 2013. The minimal preprocessing pipelines for the Human Connectome Project. Neuroimage 80, 105–124.

Henschel, L., Conjeti, S., Estrada, S., Diers, K., Fischl, B., Reuter, M., 2020. FastSurfer - A fast and accurate deep learning based neuroimaging pipeline. Neuroimage 219, 117012.

Hopkins, W.D., 2018. Motor and Communicative Correlates of the Inferior Frontal Gyrus (Broca's Area) in Chimpanzees. Origins of Human Language: Continuities and Discontinuities with Nonhuman Primates 153.

Hopkins, W.D., Meguerditchian, A., Coulon, O., Bogart, S., Mangin, J.-F., Sherwood, C.C., Grabowski, M.W., Bennett, A.J., Pierre, P.J., Fears, S., Woods, R., Hof, P.R., Vauclair, J., 2014. Evolution of the central sulcus morphology in primates. Brain Behav. Evol. 84, 19–30.

Hwang, H., Rehman, H.Z.U., Lee, S., 2019. 3D U-Net for Skull Stripping in Brain MRI. Applied Sciences. https://doi.org/10.3390/app9030569

Jenkinson, M., Beckmann, C.F., Behrens, T.E.J., Woolrich, M.W., Smith, S.M., 2012. FSL. Neuroimage 62, 782–790.

Jung, B., Taylor, P.A., Seidlitz, J., Sponheim, C., Perkins, P., 2020. A comprehensive macaque fMRI pipeline and hierarchical atlas. BioRxiv.

Ketkar, N., 2017. Introduction to PyTorch. Deep Learning with Python. https://doi.org/10.1007/978-1-4842-2766-4_12

Kingma, D.P., Ba, J., 2014. Adam: A Method for Stochastic Optimization. arXiv [cs.LG].

Kleesiek, J., Urban, G., Hubert, A., Schwarz, D., Maier-Hein, K., Bendszus, M., Biller, A., 2016. Deep MRI brain extraction: A 3D convolutional neural network for skull stripping. Neuroimage 129, 460–469.

Lepage, C., Wagstyl, K., Jung, B., Seidlitz, J., Sponheim, C., 2020. CIVET-Macaque: an automated pipeline for MRI-based cortical surface generation and cortical thickness in macaques. bioRxiv.

Lohmeier, J., Kaneko, T., Hamm, B., Makowski, M.R., Okano, H., 2019. atlasBREX: Automated template-derived brain extraction in animal MRI. Sci. Rep. 9, 12219.

Lyksborg, M., Puonti, O., Agn, M., Larsen, R., 2015. An Ensemble of 2D Convolutional Neural Networks for Tumor Segmentation. Image Analysis. https://doi.org/10.1007/978-3-319-19665-7_17

McDonald, A.R., Muraskin, J., Van Dam, N.T., Froehlich, C., Puccio, B., Pellman, J., Bauer, C.C.C., Akeyson, A., Breland, M.M., Calhoun, V.D., Carter, S., Chang, T.P., Gessner, C., Gianonne, A., Giavasis, S., Glass, J., Homann, S., King, M., Kramer, M., Landis, D., Lieval, A., Lisinski, J., Mackay-Brandt, A., Miller, B., Panek, L., Reed, H., Santiago, C., Schoell, E., Sinnig, R., Sital, M., Taverna, E., Tobe, R., Trautman, K., Varghese, B., Walden, L., Wang, R., Waters, A.B., Wood, D.C., Castellanos, F.X., Leventhal, B., Colcombe, S.J., LaConte, S., Milham, M.P., Craddock, R.C., 2017. The real-time fMRI neurofeedback based stratification of Default Network Regulation Neuroimaging data repository. Neuroimage 146, 157–170.

Messinger, A., Sirmpilatze, N., Heuer, K., Loh, K.K., Mars, R., Sein, J., Xu, T., Glen, D., Jung, B., Seidlitz, J., Taylor, P., Toro, R., Garza-Villareal, E., Sponheim, C., Wang, X., Benn, A., Cagna, B., Dadarwal, R., Evrard, H., Garcia-Saldivar, P., Giavasis, S., Hartig, R., Lepage, C., Liu, C., Majka, P., Merchant, H., Milham, M., Rosa, M., Tasserie, J., Uhrig, L., Margulies, D., Klink, P.C., this issue. A collaborative resource platform for non-human primate neuroimaging. NeuroImage. https://doi.org/10.1101/2020.07.31.230185

Milham, M.P., Ai, L., Koo, B., Xu, T., Amiez, C., Balezeau, F., Baxter, M.G., Blezer, E.L.A., Brochier, T., Chen, A., Croxson, P.L., Damatac, C.G., Dehaene, S., Everling, S., Fair, D.A., Fleysher, L., Freiwald, W., Froudist-Walsh, S., Griffiths, T.D., Guedj, C., Hadj-Bouziane, F., Ben Hamed, S., Harel, N., Hiba, B., Jarraya, B., Jung, B., Kastner, S., Klink, P.C., Kwok, S.C., Laland, K.N., Leopold, D.A., Lindenfors, P., Mars, R.B., Menon, R.S., Messinger, A., Meunier, M., Mok, K.,

Morrison, J.H., Nacef, J., Nagy, J., Rios, M.O., Petkov, C.I., Pinsk, M., Poirier, C., Procyk, E., Rajimehr, R., Reader, S.M., Roelfsema, P.R., Rudko, D.A., Rushworth, M.F.S., Russ, B.E., Sallet, J., Schmid, M.C., Schwiedrzik, C.M., Seidlitz, J., Sein, J., Shmuel, A., Sullivan, E.L., Ungerleider, L., Thiele, A., Todorov, O.S., Tsao, D., Wang, Z., Wilson, C.R.E., Yacoub, E., Ye, F.Q., Zarco, W., Zhou, Y.-D., Margulies, D.S., Schroeder, C.E., 2018. An Open Resource for Non-human Primate Imaging. Neuron 100, 61–74.e2.

Oler, J.A., Fox, A.S., Shelton, S.E., Rogers, J., Dyer, T.D., Davidson, R.J., Shelledy, W., Oakes, T.R., Blangero, J., Kalin, N.H., 2010. Amygdalar and hippocampal substrates of anxious temperament differ in their heritability. Nature 466, 864–868.

Pontes-Filho, S., Dahl, A.G., Nichele, S., Gustavo Borges Moreno, 2019. A deep learning based tool for automatic brain extraction from functional magnetic resonance images in rodents. arXiv [eess.IV].

Puccio, B., Pooley, J.P., Pellman, J.S., Taverna, E.C., Craddock, R.C., 2016. The preprocessed connectomes project repository of manually corrected skull-stripped T1-weighted anatomical MRI data. Gigascience 5, 45.

Rehman, S., Ajmal, H., Farooq, U., Ain, Q.U., Riaz, F., Hassan, A., 2018. Convolutional neural network based image segmentation: a review. Pattern Recognition and Tracking XXIX. https://doi.org/10.1117/12.2304711

Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation. Lecture Notes in Computer Science. https://doi.org/10.1007/978-3-319-24574-4_28

Roy, S., Butman, J.A., Reich, D.S., Calabresi, P.A., Pham, D.L., 2018. Multiple Sclerosis Lesion Segmentation from Brain MRI via Fully Convolutional Neural Networks. arXiv [cs.CV].

Roy, S., Knutsen, A., Korotcov, A., Bosomtwi, A., Dardzinski, B., Butman, J.A., Pham, D.L., 2018. A deep learning framework for brain extraction in humans and animals with traumatic brain injury, in: 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018). pp. 687–691.

Salehi, S.S.M., Hashemi, S.R., Velasco-Annis, C., Ouaalam, A., Estroff, J.A., Erdogmus, D., Warfield, S.K., Gholipour, A., 2018. Real-time automatic fetal brain extraction in fetal MRI by deep learning, in: 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018). pp. 720–724.

Ségonne, F., Dale, A.M., Busa, E., Glessner, M., Salat, D., Hahn, H.K., Fischl, B., 2004. A hybrid approach to the skull stripping problem in MRI. Neuroimage 22, 1060–1075.

Seidlitz, J., Sponheim, C., Glen, D., Ye, F.Q., Saleem, K.S., Leopold, D.A., Ungerleider, L., Messinger, A., 2018. A population MRI brain template and analysis tools for the macaque. Neuroimage 170, 121–131.

Sørensen, T., 1948. A Method of Establishing Groups of Equal Amplitude in Plant Sociology Based on Similarity of Species Content and Its Application to Analyses of the Vegetation on Danish Commons.

Tasserie, J., Grigis, A., Uhrig, L., Dupont, M., Amadon, A., Jarraya, B., 2020. Pypreclin: An automatic pipeline for macaque functional MRI preprocessing. Neuroimage 207, 116353.

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, 1310–1320.

Tustison, N.J., Cook, P.A., Holbrook, A.J., Johnson, H.J., Muschelli, J., Devanyi, G.A., Duda, J.T., Das, S.R., Cullen, N.C., Gillen, D.L., Others, 2020. ANTsX: A dynamic ecosystem for quantitative biological and medical imaging. medRxiv.

Xu, T., Sturgeon, D., Ramirez, J.S.B., Froudist-Walsh, S., Margulies, D.S., Schroeder, C.E., Fair, D.A., Milham, M.P., 2019. Interindividual Variability of Functional Connectivity in Awake and Anesthetized Rhesus Macaque Monkeys. Biol Psychiatry Cogn Neurosci Neuroimaging 4, 543–553.

Xu, T., Yang, Z., Jiang, L., Xing, X.-X., Zuo, X.-N., 2015. A Connectome Computation System for discovery science of brain. Sci Bull. Fac. Agric. Kyushu Univ. 60, 86–95.

Yan, C.-G., Wang, X.-D., Zuo, X.-N., Zang, Y.-F., 2016. DPABI: Data Processing & Analysis for (Resting-State) Brain Imaging. Neuroinformatics 14, 339–351.

Yogananda, C.G.B., Wagner, B.C., Murugesan, G.K., Madhuranthakam, A., Maldjian, J.A., 2019. A Deep Learning Pipeline for Automatic Skull Stripping and Brain Segmentation, in: 2019 IEEE 16th

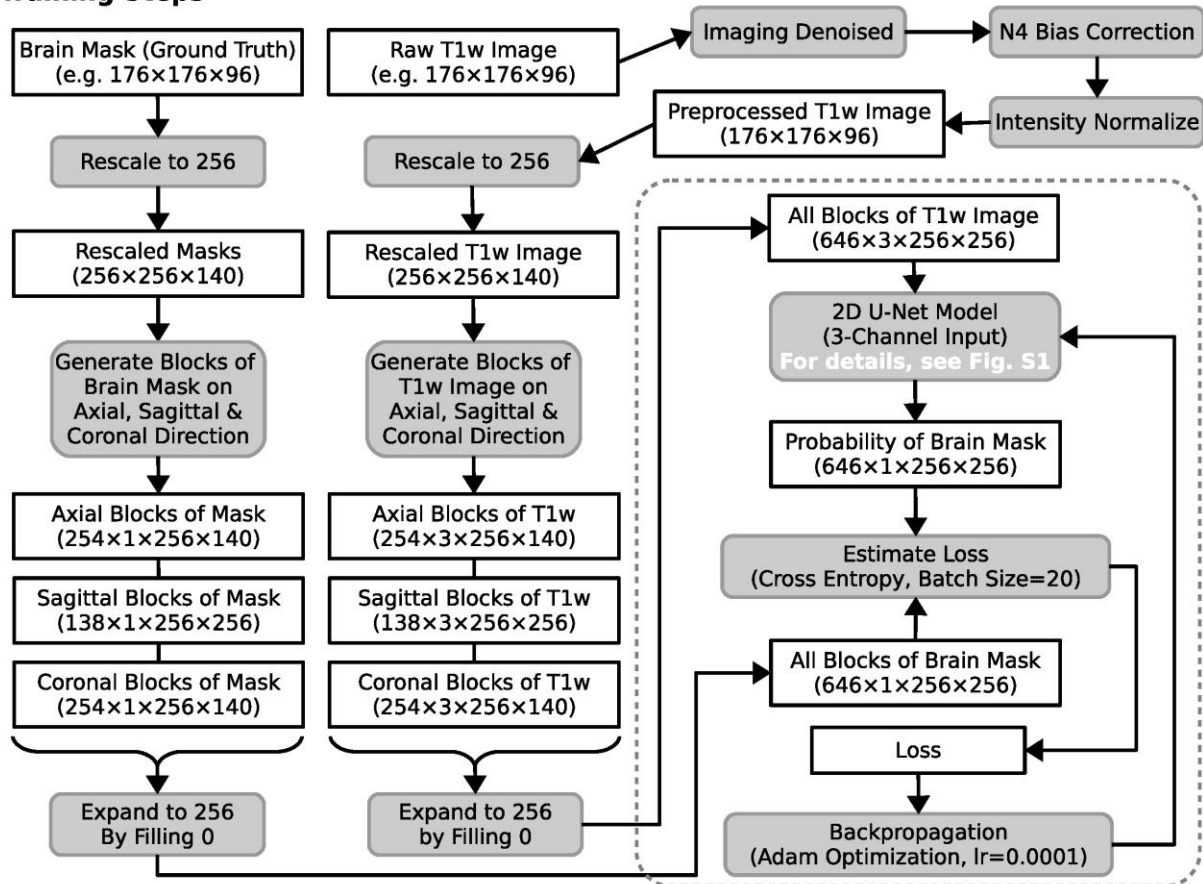International Symposium on Biomedical Imaging (ISBI 2019). pp. 727–731.

Yosinski, J., Clune, J., Bengio, Y., Lipson, H., 2014. How transferable are features in deep neural networks?, in: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N., Weinberger, K.Q. (Eds.), Advances in Neural Information Processing Systems. Curran Associates, Inc., pp. 3320–3328.

Yushkevich, P.A., Piven, J., Hazlett, H.C., Smith, R.G., Ho, S., Gee, J.C., Gerig, G., 2006. User-guided 3D active contour segmentation of anatomical structures: significantly improved efficiency and reliability. Neuroimage 31, 1116–1128.

Zhao, G., Liu, F., Oler, J.A., Meyerand, M.E., Kalin, N.H., Birn, R.M., 2018. Bayesian convolutional neural network based MRI brain extraction on nonhuman primates. Neuroimage 175, 32–44.
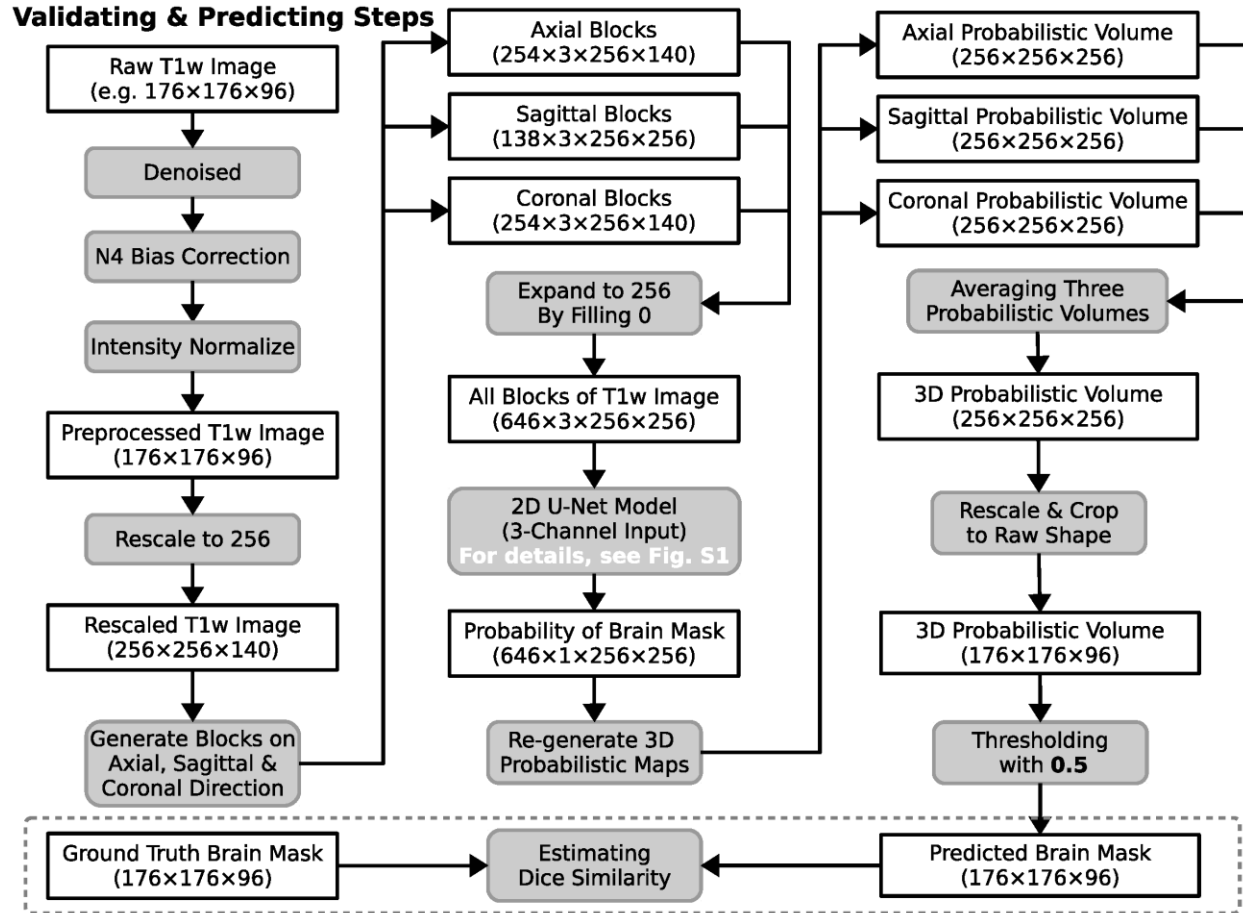
## Supplementary Information:



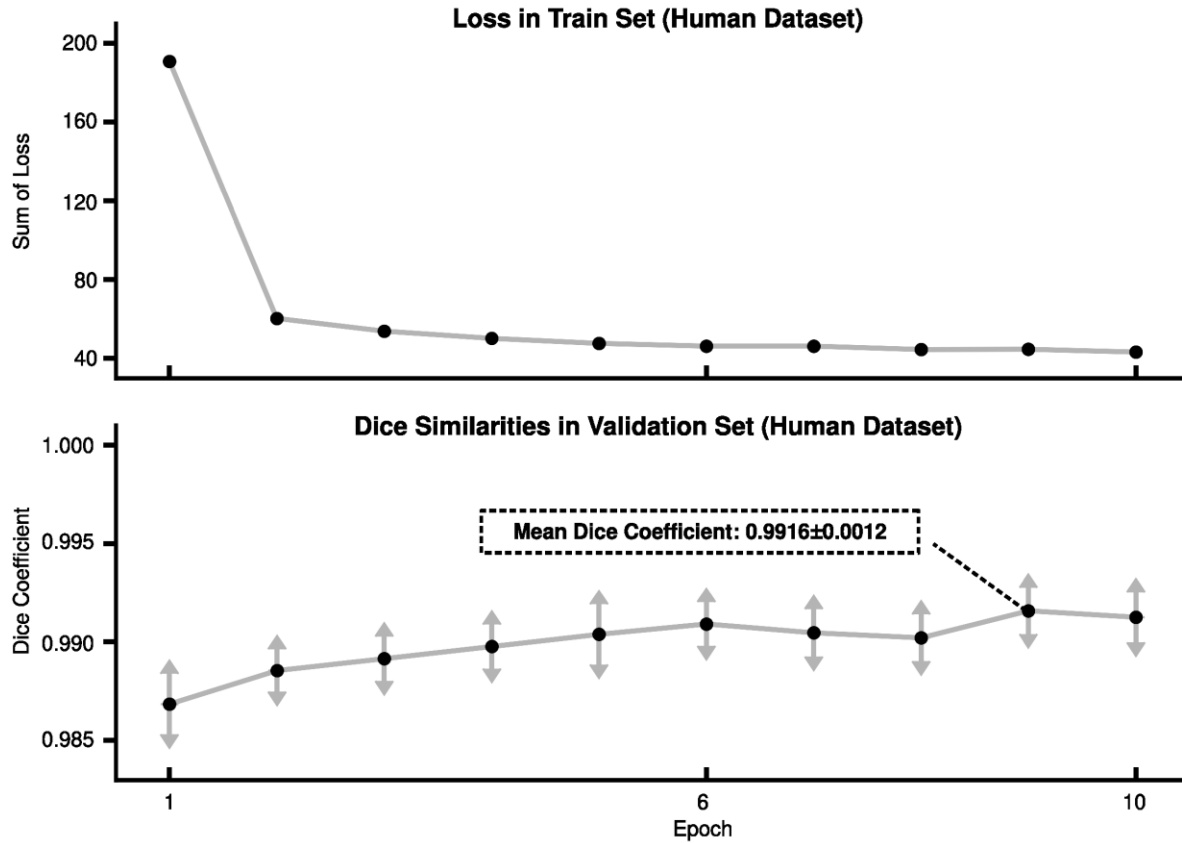**Figure S1.** Outline of the architecture of 2D U-Net.

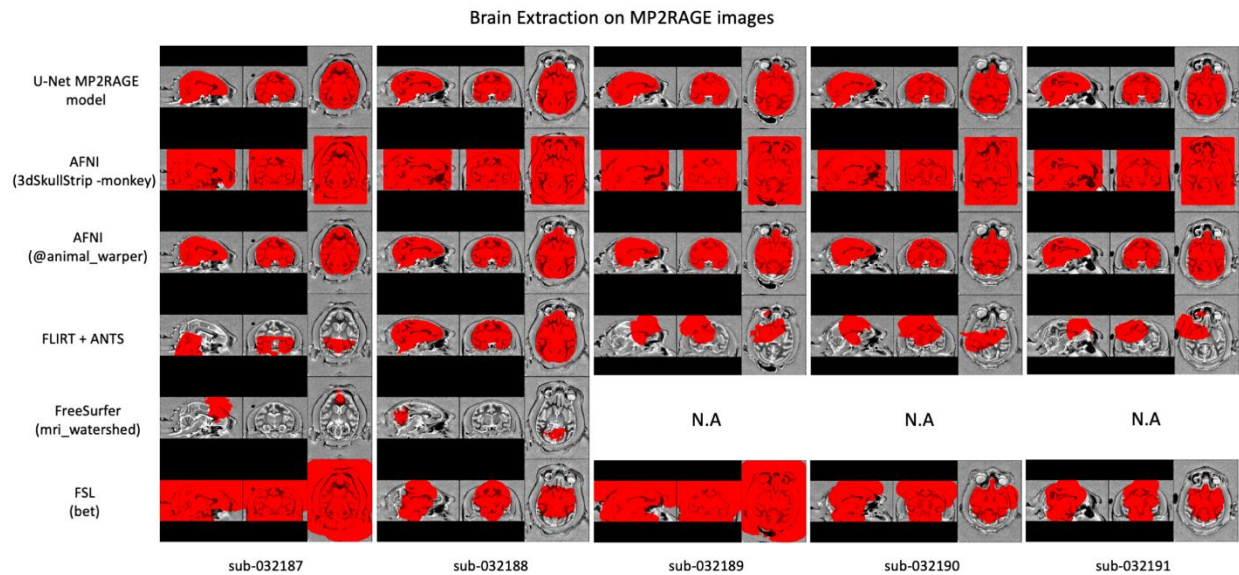**Figure S2.** U-Net training procedure for each 3D image.

**Figure S3.** U-Net validating and prediction procedure for each 3D image.

**Figure S4.** The loss and Dice coefficient of the U-Net model for each training epoch on the human dataset. The mean and standard deviation (arrows) across validation samples was calculated for each epoch.

**Figure S5.** Skullstripping results of U-Net MP2RAGE model and traditional pipelines for MP2RAGE dataset from PRIME-DE (site-uwo).