



Published in final edited form as:

*Neuroimage*. 2022 October 15; 260: 119474. doi:10.1016/j.neuroimage.2022.119474.

## SynthStrip: skull-stripping for any brain image

**Andrew Hoopes<sup>a</sup>, Jocelyn S. Mora<sup>a</sup>, Adrian V. Dalca<sup>a,b,c</sup>, Bruce Fischl<sup>a,b,c,d,1</sup>, Malte Hoffmann<sup>a,b,1,\*</sup>**

<sup>a</sup>Athinoula A. Martinos Center for Biomedical Imaging, Massachusetts General Hospital, 149 13<sup>th</sup> St, Charlestown, MA, USA

<sup>b</sup>Department of Radiology, Harvard Medical School, 25 Shattuck St, Boston, MA, USA

<sup>c</sup>Computer Science and Artificial Intelligence Lab, Massachusetts Institute of Technology, 32 Vassar St, Cambridge, MA, USA

<sup>d</sup>Harvard-MIT Division of Health Sciences and Technology, 77 Massachusetts Ave, Cambridge, MA, USA

### Abstract

The removal of non-brain signal from magnetic resonance imaging (MRI) data, known as skull-stripping, is an integral component of many neuroimage analysis streams. Despite their abundance, popular classical skull-stripping methods are usually tailored to images with specific acquisition properties, namely near-isotropic resolution and T1-weighted (T1w) MRI contrast, which are prevalent in research settings. As a result, existing tools tend to adapt poorly to other image types, such as stacks of thick slices acquired with fast spin-echo (FSE) MRI that are common in the clinic. While learning-based approaches for brain extraction have gained traction in recent years, these methods face a similar burden, as they are only effective for image types seen during the training procedure. To achieve robust skull-stripping across a landscape of imaging protocols, we introduce SynthStrip, a rapid, learning-based brain-extraction tool. By leveraging anatomical segmentations to generate an entirely synthetic training dataset with anatomies, intensity distributions, and artifacts that far exceed the realistic range of medical images, SynthStrip learns to successfully generalize to a variety of real acquired brain images,

---

This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>)

\*Corresponding author. [mhoffmann@mgh.harvard.edu](mailto:mhoffmann@mgh.harvard.edu).

<sup>1</sup>These authors contributed equally.

#### Declaration of competing interest

Bruce Fischl has a financial interest in CorticoMetrics, a company whose medical pursuits focus on brain imaging and measurement technologies. This interest is reviewed and managed by Massachusetts General Hospital and Mass General Brigham in accordance with their conflict-of-interest policies. The authors declare that they have no other known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

#### Credit authorship contribution statement

**Andrew Hoopes:** Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Writing – original draft, Writing – review & editing, Visualization. **Jocelyn S. Mora:** Validation, Investigation, Resources, Data curation, Writing – review & editing, Visualization. **Adrian V. Dalca:** Conceptualization, Methodology, Writing – review & editing, Supervision. **Bruce Fischl:** Conceptualization, Methodology, Software, Investigation, Writing – review & editing, Supervision, Funding acquisition. **Malte Hoffmann:** Conceptualization, Methodology, Software, Investigation, Data curation, Writing – original draft, Writing – review & editing, Supervision, Project administration.

#### Supplementary materials

Supplementary material associated with this article can be found, in the online version, at doi:10.1016/j.neuroimage.2022.119474.

removing the need for training data with target contrasts. We demonstrate the efficacy of SynthStrip for a diverse set of image acquisitions and resolutions across subject populations, ranging from newborn to adult. We show substantial improvements in accuracy over popular skull-stripping baselines – all with a single trained model. Our method and labeled evaluation data are available at <https://w3id.org/synthstrip>.

## Keywords

Skull stripping; Brain extraction; Image synthesis; MRI-contrast agnosticism; Deep learning

---

## 1. Introduction

Skull-stripping, also known as brain extraction, involves the removal of non-brain tissue signal from magnetic resonance imaging (MRI) data. This process is useful for anonymizing brain scans and a fundamental component of many neuroimage analysis pipelines, such as FreeSurfer (Fischl, 2012), FSL (Jenkinson et al., 2012), AFNI (Cox, 1996), and ANTs (Avants et al., 2011). These packages include tools that typically require brain-extracted input images and might perform inaccurately, or even fail, without removal of irrelevant and distracting tissue. One such class of algorithms that benefits from this systematic tissue extraction is image registration, a core element of atlas-based segmentation and other analyses. Nonlinear registration (Ashburner, 2007; Avants et al., 2008; Rueckert et al., 1999; Vercauteren et al., 2009) estimates local deformations between pairs of images, and these algorithms tend to produce more accurate estimates when they can focus entirely on the anatomy of interest (Klein et al., 2009; Ou et al., 2014). Similarly, skull-stripping increases the reliability of linear registration (Cox and Jesmanowicz, 1999; Friston et al., 1995; Hoffmann et al., 2015; Jenkinson and Smith, 2001; Jiang et al., 1995; Modat et al., 2014; Reuter et al., 2010) by excluding anatomy that deforms non-rigidly, such as the eyes, jaw, and tongue (Andrade et al., 2018; Fein et al., 2006; Fischmeister et al., 2013; Hoffmann et al., 2020).

Classical skull-stripping techniques are well-explored and widespread, but popular methods are often tailored to images with specific modalities or acquisition properties. Most commonly, these methods focus on three-dimensional (3D) T1-weighted (T1w) MRI scans acquired with MPRAGE sequences (van der Kouwe et al., 2008; Marques et al., 2010; Mugler and Brookeman, 1990), which are ubiquitous in neuroimaging research. While some skull-stripping tools accommodate additional contrasts, these methods are ultimately limited to a predefined set of viable image types and do not properly adapt to inputs outside this set. For example, skull-stripping tools developed for near-isotropic, adult brain images may perform poorly when applied to infant subjects or clinical scans with thick slices, such as stacks of 2D fast spin-echo (FSE) acquisitions.

When a suitable brain extraction method is not available for a particular scan type, a common workaround involves skull-stripping a compatible image of the same subject and computing a co-registration to propagate the extracted brain mask to the target image of interest (Iglesias et al., 2011). Unfortunately, an accurate intra-subject alignment can require

significant manual tuning because the target image still includes extra-cerebral matter that may impede linear registration quality (Reuter et al., 2010). Crucially, this procedure also requires the existence of an additional, strip-able image, often a high-resolution isotropic T1w or T2-weighted (T2w) scan, which is rare, for example, in clinical screening protocols, introducing a barrier to the clinical adoption of analysis tools.

While classical algorithms for skull-stripping are limited by their assumptions about the spatial features and intensity distributions in the input images, supervised deep-learning approaches, which leverage convolutional neural networks (CNNs), can, in principle, learn to extract a region of interest from any image type given sufficient anatomical contrast and resolution. In practice, these networks achieve high accuracy for data types observed during training, but their performance often deteriorates on images with characteristics unseen during training (Hendrycks et al., 2021; Hoffmann et al., 2021b; Jog et al., 2019; Karani et al., 2018). In consequence, robust, supervised learning-based approaches depend on the availability of a representative training dataset that contains accurate ground-truth annotations and exposes the network to a landscape of image types. While numerous public datasets provide access to widely used MRI acquisitions for which target brain masks can be easily derived with classical methods, curating a diverse training dataset with uncommon sequences and sufficient anatomical variability is a challenging task that requires substantial human effort. As a result, current deep-learning skull-stripping methods are trained with few different data types and deliver state-of-the-art results only for particular subsets of image characteristics (Hwang et al., 2019; Kleesiek et al., 2016; Salehi et al., 2017).

Recently, a novel learning strategy alleviates the requirement for representative acquired training data by optimizing networks with a wide array of synthetic images, each generated directly from a precomputed label map (Billot et al., 2020; Hoffmann et al., 2021b). This synthesis scheme enables networks to accurately carry out tasks on any image type at evaluation-time without ever sampling real target acquisitions during training, and it has been effectively employed for segmentation (Billot et al., 2020) and deformable image registration (Hoffmann et al., 2021b). To build on deep-learning methods for brain extraction while addressing their shortcomings, we adapt the synthesis technique and introduce SynthStrip, a flexible brain-extraction tool that can be deployed universally on a variety of brain images. By exposing a CNN to an arbitrary and deliberately unrealistic range of anatomies, contrasts, and artifacts, we obtain a model that is agnostic to acquisition specifics, as it never samples any real data during training. Consequently, this scheme enables SynthStrip to extract the brain from a wide array of neuroimaging data types, and we demonstrate its viability and improvement over popular baselines using a varied test set that spans both research scans and clinical exams (Fig. 1). The test set includes T1w, T2w, T2w fluid attenuated inversion recovery (T2-FLAIR), and proton-density (PDw) contrasts as well as clinical FSE scans with slices and high in-plane resolution, and low-resolution EPI, ranging across age and pathology. We demonstrate the ability of SynthStrip to generalize beyond structural MRI, to MR angiography (MRA), diffusion-weighted imaging (DWI), fluorodeoxyglucose positron emission tomography (FDG-PET), and even computed tomography (CT). We make our validation set publicly available to promote further development and evaluation of brain-extraction tools.

## 2. Related work

In this section, we briefly review the automated brain-extraction techniques that we use as baseline methods. We include both classical and deep-learning baselines introduced over the last two decades, focusing in particular on those with high efficacy and popularity in the research domain. For an exhaustive overview of skull-stripping methods, see Fatima et al., 2020.

### 2.1. Classical skull-stripping

Classical, or traditional, algorithms that remove non-brain image signal vary substantially in their implementation (Cox, 1996; Eskildsen et al., 2012; Iglesias et al., 2011; Roy et al., 2017; Ségonne et al., 2004; Shattuck et al., 2001; Smith, 2002). One common class of approaches leverages a deformable mesh model to reconstruct a smooth boundary of the brain matter surface. The widely-used Brain Extraction Tool (BET; Smith, 2002), distributed as part of the FSL package (Jenkinson et al., 2012), utilizes this technique by initializing a spherical mesh at the barycenter of the brain and projecting mesh vertices outwards to model the brain border. Since BET uses locally adaptive intensity thresholds to distinguish brain and non-brain voxels, it generalizes to a variety of contrasts, such as T1w, T2w, and PDw. To prevent surface leaks beyond the brain boundary, 3dSkullStrip, a component of AFNI (Cox, 1996), extends the BET strategy by considering information on the surface exterior, accounting for eyes, ventricles, and skull.

The popular hybrid approach (Ségonne et al., 2004) available in FreeSurfer also leverages a deformable surface paradigm, combing it with a watershed algorithm and statistical atlas to improve robustness. First, the watershed establishes an estimate of the white-matter mask, which is then refined to the brain boundary using a surface mesh expansion. A probabilistic atlas of intensity distributions helps prevent outliers during mesh fitting, and erroneous brain mask voxels are removed during post-processing via a graph cuts algorithm (Greig et al., 1989; Sadananthan et al., 2010) that thresholds the cerebrospinal fluid (CSF). While effective, this technique is optimized only for images with T1w contrast, since it relies on the underlying assumption that white matter is surrounded by darker gray matter and CSF. Another hybrid approach, ROBEX (Iglesias et al., 2011), exploits a joint generative-discriminative model. A Random Forest classification (Breiman, 2001) detects the brain contour, which is used to fit a point-distribution model to the brain target. The skull-stripping tool BEaST (Eskildsen et al., 2012) builds on patch-based, non-local segmentation techniques (Coupé et al., 2010; 2011; Roy et al., 2017) and assigns a label to each voxel by comparing its local neighborhood to patches in a reference set with prior labels. With the exception of BET and 3dSkullStrip, all of these tools were specifically developed for T1w images.

### 2.2. Deep-learning approaches

Innovations in deep-learning have gained popularity as methodological building blocks for an array of tasks in medical image analysis, including skull-stripping. Various learning-based extraction methods have been proposed, demonstrating accuracy and speed that often out-perform their classical counterparts. These models are optimized in a supervised

fashion, using a set of acquired training images with corresponding ground-truth brain masks, derived through classical methods or manual segmentation. An early, cross-contrast approach, Deep MRI Brain Extraction (DMBE) (Kleesiek et al., 2016), trains a 3D CNN on combinations of T1w, T2w, and FLAIR contrasts and matches the accuracy of classical baselines for several datasets, including clinical scans with brain tumors. Conversely, Auto-Net (Salehi et al., 2017) introduces two separate 2.5D architectures that skull-strip volumes by individually segmenting sagittal, coronal, and transverse views of same image and fusing the predictions with an auto-context algorithm (Tu and Bai, 2009). The first architecture leverages convolutions on single-resolution voxel-wise patches, while the second utilizes a scale-space U-Net architecture (Ronneberger et al., 2015) to predict the brain mask. Auto-Net is effective for both adult and neonatal brain scans but only trained with T1w images. CONSNNet (Lucena et al., 2019) similarly leverages a 2D U-Net, applied across image slices in each plane, to strip 3D T1w images. More recently, implementations using full 3D U-Nets (Hsu et al., 2020; Hwang et al., 2019) have robustly matched or exceeded start-of-the-art brain-extraction performance.

### 2.3. Contribution

SynthStrip builds on a solid foundation laid by prior studies of deep-learning algorithms for brain extraction, enabling us to choose among network architectures well suited for this particular task. We emphasize that our goal is not to compare or make claims on the optimality of specific architectures – the discussed algorithms may perform equally well. Instead, our focus is on exploiting a novel training strategy using synthetic data only, to build an easy-to-use skull-stripping tool that alleviates the requirement of expanding the training set and re-optimizing network weights every time a new image type is to be supported.

## 3. Method

To predict robust brain masks for an array of real image types, we train a deep convolutional neural network on a vast landscape of images synthesized with a deliberately unrealistic range of anatomies, acquisition parameters, and artifacts. From a dataset  $\mathcal{D}$  of precomputed, whole-head segmentations with brain and non-brain tissue labels, we sample a segmentation  $s \in \mathcal{D}$  at each optimization step and use it to generate a gray-scale head scan  $x$  with randomized acquisition characteristics. In effect, this paradigm synthesizes a stream of training images used to optimize a SynthStrip network  $g_{\theta}$ , with trainable parameters  $\theta$ , in a supervised fashion:

$$\hat{\theta} = \operatorname{argmin}_{\theta} [\mathbb{E}_{\mathcal{D}} [\mathcal{L}(y, \hat{y})]], \quad (1)$$

where  $y$  is the predicted brain mask,  $\hat{y}$  is the target brain mask derived by merging the brain labels of  $s$ , and  $\mathcal{L}$  is the loss function that measures similarity between  $y$  and  $\hat{y}$ .

### 3.1. Synthesis

Building from previous work (Billot et al., 2020; Hoffmann et al., 2021b), we use a generative model to synthesize a stream of random images with substantial anatomical and intensity variation, as exhibited in Fig. 2. At each training step, parameters that dictate synthesis components are randomly sampled from predetermined ranges and probability distributions explicitly defined in Table 1. We emphasize that while the generated scans can appear implausible, these training images do not need to be realistic in order for the SynthStrip model to accurately generalize to real images at test-time.

To generate a gray-scale image  $x$  from a whole-head anatomical segmentation  $s$ , we first create spatial variability to subject the network to a landscape of possible head positions and anatomical irregularities. This is accomplished by manipulating  $s$  with a spatial transformation  $t$ , composed of an affine transform (with random translation, scaling, and rotation) and a nonlinear deformation. The deformation is generated by sampling random 3D displacement vectors from a normal distribution, with random scale, at an arbitrarily low image resolution. This random displacement field is vector-integrated, using five *scaling and squaring* steps to encourage a diffeomorphic warp (Arsigny et al., 2006; Dalca et al., 2019), and tri-linearly resampled to match the resolution of  $s$ . After applying the randomized transform, the resulting segmentation  $s_t$  serves as the basis for deriving the image  $x$  and target brain mask  $\hat{y}$ , which is obtained by merging the labels of  $s_t$  into brain and non-brain classes.

To compute  $x$ , we consider a Bayesian model of MR contrast, which assumes that the voxel intensity of each tissue type in the image can be represented by a single Gaussian distribution. Reversing this generalization, we assign a random distribution of tissue intensity to every anatomical label in  $s_t$  and use this artificial mixture model to attain an image with arbitrary contrast by replacing each label voxel in  $s_t$  with a random value drawn from its corresponding intensity distribution. Following the synthesis, we aim to simulate various artifacts and geometric properties that might exist across modality and acquisition type. First, we corrupt the image with a spatially varying intensity bias field, generated by resizing a low-resolution image sampled from a normal distribution with zero mean. The corrupted image is computed by an element-wise multiplication with the voxel-wise exponential of the bias field. Second, we perform gamma augmentation by globally exponentiating all voxels with a single value  $\exp(\gamma)$ , where  $\gamma$  is a normally sampled parameter. Lastly, to account for scans with a partial field of view (FOV) and varied resolution, we randomly crop the image content and down-sample along an indiscriminate set of axes. Before down-sampling by an arbitrary factor  $r$ , we simulate partial-volume effects by blurring the image using a Gaussian kernel with standard deviation  $\sigma = r/4$ . The image cropping and down-sampling components are applied with a 50% probability rate during synthesis.

### 3.2. Loss

We optimize  $g_\theta$  using a loss function  $\mathcal{L}$  that measures the similarity between predicted and target brain masks. Unless otherwise stated, we employ a loss  $\mathcal{L} = \mathcal{L}_{sdt}$  that encourages the network to predict a signed distance transform (SDT)  $d$  representing the minimum distance

(in *mm*) to the skull boundary at each voxel. Distances are positive within the brain and negative outside, facilitating the extraction of a binary brain mask  $y$  from  $d$  at test-time by simple thresholding. The training paradigm is outlined in Fig. 3. During training, an exact target Euclidean SDT  $\hat{d}$  is computed from the target brain mask  $\hat{y}$ , and the similarity between  $d$  and  $\hat{d}$  is measured by their mean squared difference (MSE). To concentrate optimization gradients to pertinent regions of the image during training,  $\hat{d}$  is banded such that voxel distances  $\hat{d}_i$  do not surpass a discrete threshold  $t$ , and all voxels that exceed the distance  $t$  are down-weighted in the MSE computation by a factor  $b$ . Therefore,

$$\mathcal{L}_{sdt} = \frac{\sum_{i \in \mathcal{P}} w_i (d_i - \hat{d}_i)^2}{\sum_{i \in \mathcal{P}} w_i}, \quad w_i = \begin{cases} b & \text{if } |\hat{d}_i| > t, \\ 1 & \text{otherwise} \end{cases}, \quad (2)$$

where  $i$  represents a voxel in the spatial image domain  $\mathcal{P}$ ,  $t = 5$  *mm* and  $b = 0.1$  in our experiments, optimally determined via a grid search.

As a complimentary analysis, we compare the distance-based loss  $\mathcal{L}_{sdt}$  against a soft Dice loss (Dice, 1945; Milletari et al., 2016), which is commonly used to optimize image segmentation models and quantifies volume overlap for pairs of labels. We define the loss  $\mathcal{L}_{dice}$  as

$$\mathcal{L}_{dice} = \frac{|y^j \odot \hat{y}^j|}{|y^j \oplus \hat{y}^j|} + \frac{|y^k \odot \hat{y}^k|}{|y^k \oplus \hat{y}^k|}, \quad (3)$$

where  $y^j$  and  $\hat{y}^j$  represent brain label maps,  $y^k$  and  $\hat{y}^k$  represent non-brain label maps, and  $\odot$  and  $\oplus$  represent voxel-wise multiplication and addition, respectively. While  $\mathcal{L}_{sdt}$  and  $\mathcal{L}_{dice}$  both result in effective skull-stripping networks, we favor the distance loss  $\mathcal{L}_{sdt}$  due to its smoothing effect on the outline of the predicted brain mask, as demonstrated in Experiment 4.4.

### 3.3. Implementation

We implement  $g_\theta$  using a 3D U-Net convolutional architecture, with down-sampling (encoder) and up-sampling (decoder) components that facilitate the integration of features across large spatial regions. The U-Net comprises seven resolution levels, which each include two convolutional operations with leaky ReLU activations (parameter  $\alpha = 0.2$ ) and filter numbers defined in Fig. 3. Down-sampling is achieved through max-pooling, and skip-connections are formed by concatenating the outputs of each encoder level with the inputs of the decoder level with corresponding resolution. In models using  $\mathcal{L} = \mathcal{L}_{sdt}$ , one final, single-feature convolutional layer with linear activation outputs the predicted SDT  $d$ . In models optimized with  $\mathcal{L} = \mathcal{L}_{dice}$ , the final layer is a two-feature convolution, with softmax activation, that outputs a probabilistic segmentation representing non-brain and brain regions.

We train SynthStrip using the Adam optimizer (Kingma and Ba, 2014) with a batch size of one and an initial learning rate of  $10^{-4}$ . This rate is reduced by a factor of two after every 20,000 optimization steps without a decrease in validation loss. At test-time, all inputs to the model are internally conformed to 1-*mm* isotropic voxel size using trilinear interpolation, and intensities are scaled between 0 and 1. The U-Net outputs are resampled such that the final brain mask is computed in the original input space. We implement SynthStrip in Python, using the open-source PyTorch (Paszke et al., 2019) and Neurite (Dalca et al., 2018) libraries, and make our tool and associated code available in the open-source FreeSurfer package (<https://w3id.org/synthstrip>). All experiments are conducted using Intel Xeon Silver 4214R CPUs and Nvidia RTX 8000 GPUs.

### 3.4. Data

In our experiments, we employ a small training dataset of adult and infant brain segmentations and a separate, larger dataset of acquired images for validation and testing that spans across age, health, resolution, and imaging modality. All data are 3D images, acquired either directly or as stacks of 2D MRI slices.

#### 3.4.1. Training data

**Datasets:** We compose a set of 80 training subjects, each with whole-head tissue segmentations, from the following three cohorts: 40 adult subjects from the Buckner40 dataset (Fischl et al., 2002), 30 locally scanned adult subjects from the Human Connectome Aging Project (HCP-A) (Bookheimer et al., 2019; Harms et al., 2018), and 10 infant subjects born full-term, scanned at Boston Children’s Hospital at ages between 0 and 18 months (de Macedo Rodrigues et al., 2015).

**Processing:** To compute anatomical segmentations of individual cerebral regions, adult and infant T1w scans are processed with SAM-SEG (Puonti et al., 2016) and the Infant FreeSurfer reconstruction pipeline (Zöllei et al., 2020), respectively. In order to build complete segmentation maps for robust whole-head image synthesis, we also generate six coarse labels of extra-cerebral tissue using a simple intensity-based labeling strategy with thresholds that mark label intensity boundaries. Considering only non-zero voxels without brain labels, we fit threshold values to each image by maximizing the similarity in number of voxels for each extra-cerebral label. These extra-cerebral labels do not necessarily represent or differentiate meaningful anatomical structures – their purpose is to provide intensity and spatial variability to synthesized regions outside the brain.

In total, the training segmentations contain 46 individual anatomical labels, with 40 brain-specific labels (including CSF), that we merge into the target brain mask  $\hat{y}$ . All training segmentations are fit to a  $256^3$  image shape with 1-*mm* isotropic resolution. We emphasize that this geometric preprocessing is not required at test-time.

#### 3.4.2. Evaluation data

**Datasets:** Our evaluation data comprise 620 images, split into validation and test subsets of sizes 22 and 598, respectively. We gather these images across seven public datasets, with makeup, resolution, and validation splits outlined in Table 2. The IXI<sup>1</sup> dataset features



a range of MRI contrasts and modalities, including T1w and T2w as well as PDw, MRA, and DWI. To simplify the DWI evaluation, a single diffusion direction is randomly extracted from each acquisition. The FSM subset (Greve et al., 2021) is derived from in-house data using standard acquisitions as well as quantitative T1 maps (qT1). In-house, pseudo-continuous ASL (PCASL) scans are acquired as stacks of 2D-EPI slices with low resolution and a small FOV that often crops the ventral brain region (Dai et al., 2008). The QIN (Clark et al., 2013; Mamonov and Kalpathy-Cramer, 2016; Prah et al., 2015) dataset comprises precontrast, clinical stacks of thick image slices from patients with newly diagnosed glioblastoma. We also include a subset of the infant T1w image dataset, using subjects held-out from training. Lastly, to evaluate the ability of SynthStrip to adapt to imaging modalities beyond MR, we gather a test cohort of brain CT and FDG-PET scans from the CERMEP-IDB-MRXFDG (CIM) database (Mérida et al., 2021).

**Ground-truth masks:** For each image in the evaluation dataset, we derive a reference brain mask using the following labelling strategy. Since every evaluation subject includes a corresponding T1w image, we generate brain masks for these scans using each *classical* baseline method evaluated in our analysis. Then, an “average” brain mask is computed for each subject by extracting the majority label value at every voxel. We refine the average masks manually before propagating the masks by rigidly aligning each subject’s T1w scan to the remaining image types with a robust registration approach (Reuter et al., 2010). Poor alignments are further refined by hand. We make the reference dataset available online to facilitate future development of skull-stripping techniques, including the original images if permitted by their respective licenses.

**3.4.3. Ethics**—This retrospective study re-analyzes previously published or shared datasets. The FSM and ASL studies were approved by the Mass General Brigham Internal Review Board (IRB). The HCP-A study was approved by IRBs at Washington University in St. Louis and Mass General Brigham. The infant study was approved by the Committee on Clinical Investigation at Boston Children’s Hospital and the Mass General Brigham IRB. All subjects gave written informed consent. No ethical approval was required for retrospective analysis of de-identified open-access data.

## 4. Experiments

We analyze the performance of SynthStrip on diverse whole-head images and compare its 3D skull-stripping accuracy to classical and deep-learning baseline tools.

### Baselines:

We select a group of skull-stripping baselines based on their popularity, determined by citation count, and effectiveness, as shown in prior work (Fatima et al., 2020; Iglesias et al., 2011). As classical baselines, we choose ROBEX 1.1, BET from FSL 6.0.4, 3dSkull-Strip (3DSS) from AFNI 21.0.21, BEaST 1.15, and the FreeSurfer 7.2 watershed algorithm (FSW). Unfortunately, many top-cited, learning based approaches do not make their models

---

<sup>1</sup>Acquired from <http://brain-development.org/ixi-dataset>.

available, even upon request to the authors. A notable exception is Deep MRI Brain Extraction (DMBE), which we therefore include. Default parameters are used for each method except BET, for which the `-R` option is provided for more accurate brain center estimation. All inputs to FSW and DMBE are re-sampled to 1-*mm* isotropic voxel sizes to accommodate the expected input resolution for these methods.

### Metrics:

We evaluate the similarity between computed and ground-truth brain masks by measuring their Dice overlap, mean and maximum (Hausdorff) surface distances, and percent difference in total volume. Baseline scores are compared to SynthStrip with a paired sample *t*-test. Sensitivity and specificity, which measure the percent of true positive and true negative brain labels, respectively, provide further insight into the properties of the computed brain masks.

**4.1. Skull-stripping accuracy**—We assess the broad skull-stripping capability of a SynthStrip model trained using images synthesized from the label maps outlined in Section 3.4.1. We compare the accuracy of our method to each of the baselines across the test set of real brain images defined in Section 3.4.2. Method runtime is compared for the FSM dataset.

The comparison demonstrates SynthStrip’s accurate and robust brain extraction, which substantially outperforms baseline methods (Tables 3, 4 and Supplementary Tables S1, S2). For every evaluation metric, brain masks predicted by SynthStrip yield significantly better scores than baseline masks ( $p < 0.05$ ) for the *vast* majority of datasets. Importantly, no baseline method significantly outperforms SynthStrip on any dataset. As shown in Fig. 4, SynthStrip achieves the highest Dice score *and* lowest mean surface distance for more than 80% of all test images, in stark contrast to the next best performing method, BET, which yields the top result for less than 10% of images. The superior performance of SynthStrip persists even when considering only T1w, near-isotropic, adult-brain images, which all of the baselines are tuned for. Across this particular subset of 127 T1w images from the IXI, FSM, and ASL datasets, SynthStrip achieves the best mean Dice, surface distance, Hausdorff distance, and volume difference (Fig. 5), and it consistently extracts the brain with high specificity and sensitivity, while other methods tend to under-perform in either of those metrics due to tendencies to substantially over- or under-label the brain. When considering the remaining non-T1w, thick-slice, and infant image types, SynthStrip’s predominance is similarly substantial (Fig. 6). For FSM T1w data, our method runs on the CPU in less than one minute (Table 5), trailing the fastest two baselines, BET and FSW, by approximately 17 seconds on average. On the GPU, SynthStrip runs substantially faster, requiring only  $1.8 \pm 0.2$  seconds.

**4.2. Qualitative brain-mask analysis**—Across the evaluation set, skull-stripping errors in SynthStrip predictions are uncommon and typically involve minimally over-segmenting the brain mask by including thin regions of extra-cerebral matter near the dorsal cortex or pockets of tissue around the eye sockets, as shown in Fig. 8. Considering only the *N* images for which SynthStrip does not achieve the best score in Fig. 4, on average,

SynthStrip lags behind the best-performing baseline by only  $-0.53 \pm 0.54$  Dice percentage points ( $N=111$ ) and  $(0.20 \pm 0.18)$  mm mean surface distance ( $N=94$ ).

The top performing baseline method is ROBEX, which yields high-quality brain extraction across many of the test datasets, with the notable exception of the qT1 cohort. ROBEX produces spatially plausible brain masks and evades drastic failure modes that exist in other base-lines, similarly to SynthStrip. However, despite its generally good performance, ROBEX has a tendency to include pockets of tissue surrounding the eyes and remove regions of cortical gray matter near the superior surface (Figs. 8 and S1).

BET and 3DSS also perform effective brain extraction across image types, but tend to fail dramatically for outlier cases. For example, BET locates the brain boundary with considerable precision when successful. However, for some image subsets, especially those with abundant non-brain matter, such as FSM, BET often includes large regions of inferior skull as well as facial and neck tissue in the brain mask. While 3DSS largely avoids such gross mislabeling, it tends to produce skull-strips that leak into neck tissue or, conversely, remove small regions of the cortical surface.

BEaST and FSW perform well for near-isotropic T1w images, such as those in the IXI, FSM, and ASL datasets. But since they are heavily optimized for the assumed spatial and intensity features of this acquisition type, they generally perform poorly or even fail completely for other contrasts. Common error modes of FSW involve the failure to remove bits of skull or inferior non-brain matter, in contrast to BEaST, which is susceptible to removing critical regions of the cortex.

The learning-based method DMBE yields suitable brain masks for near-isotropic image types with T1w contrast but frequently leaves substantial, unconnected components of non-brain matter. While DMBE extracts the brain tissue border as opposed to CSF, our analysis shows that the predominant contributor to the discrepancy between DMBE and ground-truth brain masks is the inclusion of neck and facial tissue (Figs. 8 and S1). DMBE model inference is slow, consuming more than a half hour to skull-strip a standard image.

**4.3. Variability across time-series data**—We analyze the consistency of SynthStrip brain masks across time-domain data by assessing the differences between diffusion-encoded directions acquired in the same session. For each subject in the DWI dataset, we affinely align and skull-strip all of the 16 diffusion-encoded frames in a common, average space (Reuter et al., 2010). We compute the number of discordant voxels across brain masks for a given method, defining discordant voxels (DV) as voxel locations with labels that differ in the time domain. We report the percent of DV relative to the brain mask volume, determined by the number of voxels labeled as brain in any frame. In this particular analysis, we only consider ROBEX, BET, and 3DSS as baselines since they generalize to DWI acquisitions. As shown in Fig. 7, SynthStrip demonstrates a high level of intra-subject consistency, as it predicts brain masks with substantially lower % DV across DWI directions than the baselines ( $p < 10^{-12}$ ). Since the % DV metric considers voxels labeled as brain for any direction, a single mask with gross mislabeling will substantially increase the metric

value, as is the case with ROBEX, which over-segments the brain for only a few directions per subject.

**4.4. Loss comparison**—During our experimentation, we find that training SynthStrip models using a traditional soft Dice loss yields comparable results to those trained with an SDT-based loss for *nearly* every metric. However, despite similar global accuracy, we observe that models trained with  $\mathcal{L}_{dice}$  predict brain masks characterized by relatively noisy and rough boundaries, as illustrated in Fig. 7. The high variability at the edge of the brain mask is emphasized by a  $6.4 \pm 3.2$  mm increase in maximum surface distance when using  $\mathcal{L}_{dice}$  compared to  $\mathcal{L}_{sdt}$ . We further quantify this discrepancy in brain-mask smoothness by computing the percent of exposed boundary voxels (EBV) that neighbor more non-brain labels than brain labels. Brain masks with noisier boundaries will exhibit larger EBV due to an increased mask surface area and number of sporadic border voxels. We perform this evaluation using the FSM data subset of 132 images with isotropic voxel size. Models trained with  $\mathcal{L}_{dice}$  predict masks with 4.5× higher EBV than models trained with  $\mathcal{L}_{sdt}$ . We hypothesize that as the network learns to estimate an SDT, it is encouraged to focus more on the boundary of mask, rather than the label as a whole, resulting in a smoother prediction of the brain border.

## 5. Discussion

We present SynthStrip, a learning-based, universal brain-extraction tool trained on diverse synthetic images. Subjected to training data that far exceeds the realistic range of medical images, the model learns to generalize across imaging modalities, anatomical variability, and acquisition schemes.

### 5.1. Baseline comparison

SynthStrip significantly improves upon baseline skull-stripping accuracy for nearly every image cohort tested, and the few exceptions to this improvement involve data subsets for which SynthStrip matches baseline performance. This predominance is in part due to the ability of SynthStrip to generalize across a wide variety of image types as well as its proclivity to avoid substantial mislabeling. In particular, varying specific acquisition characteristics during synthesis promotes network robustness to such characteristics across a range of protocols. For example, simulating partial-volume effects with blurring and randomizing the resolution enable SynthStrip to accurately generalize to clinical thick slice acquisitions and those with large voxel sizes. By learning robust, large-scale spatial features of representative brain masks, the model consistently predicts masks of realistic and expected shape. Baseline techniques, on the other hand, often rely on weak spatial priors and are therefore prone to over- or under-segment brain tissue when confronted with image features that are unexpected or unaccounted for (Figs. 8 and S1).

ROBEX's consistent performance across contrasts and modality is somewhat unexpected since the discriminative edge detector is trained only for T1w scans. We hypothesize that the coupled shape model is able to compensate for any intensity bias encoded in the discriminative detector. The T1w-specific approaches BEAST and FSW could be effective

for other MRI contrasts if provided known intensity priors of the brain matter. However, this work would require substantial human effort as it needs to be repeated for every new image type. The substantial, unconnected components of non-brain matter frequently left by DMBE are likely a byproduct of its convolutional architecture, which does not leverage multiple resolution levels to gather spatial features across large distances.

## 5.2. Use for brain-specific registration

Consistent brain extraction across different images from the same subject is critical for accurate analysis of time-series acquisitions. For example, diffusion (Holdsworth et al., 2012; Jones and Leemans, 2011) and functional MRI analyses (Ashburner, 2009; Jenkinson et al., 2002) depend on within-subject registration of individual frames acquired across time to undo the effect of any head motion during the scan. Unfortunately, anatomical structures that deform non-rigidly between frames, such as the neck or tongue, can hamper brain-registration accuracy and thus impinge on downstream results. While this effect can be accounted for by first removing non-brain tissue from each frame to achieve brain-specific registration (Andrade et al., 2018; Fischmeister et al., 2013), it requires consistent brain extraction across frames (Andrade et al., 2018; Fein et al., 2006; Fischmeister et al., 2013; Hoffmann et al., 2020). SynthStrip's high within-subject consistency despite substantial contrast differences across the diffusion encoding demonstrates its potential for regularizing retrospective motion correction of time-series data.

## 5.3. Model and data availability

Even as learning-based methods in neuroimaging analysis continue to grow in popularity, developers of deep-learning skull-stripping tools are sometimes disinclined to provide easy-to-use distributions of their work. Out of the three promising methods discussed in this work, only DMBE makes its models and code publicly available for use. In contrast, we make SynthStrip available as a universal, cross-platform command-line utility, distributed both as a standalone and as a built-in FreeSurfer tool. To facilitate further development and testing of robust skull-stripping tools, we also make our evaluation data and ground-truth labels available at <https://w3id.org/synthstrip>.

## 5.4. Future work

While SynthStrip facilitates state-of-the-art brain extraction, we aim to extend the tissue-extraction strategy to other applications both within and beyond neuroimaging. One such application is fetal head extraction from in-utero fetal MRI scans. Due to excessive motion, fetal MRI is limited to the acquisition of sub-second 2D slices. However, stacks of several slices are needed to cover the anatomy of interest, and while their inplane resolution is typically of the order of  $1\text{ mm} \times 1\text{ mm}$ , views across slices are hampered by slice thicknesses of 4–6 mm and between-slice motion (Hoffmann et al., 2021a). To enable full 3D views of the fetal brain, post-processing tools for super-resolution reconstruction have emerged, that aim to reconstruct a high-quality volume of isotropic resolution from a number of slice stacks acquired at different angles (Ebner et al., 2020; Iglesias et al., 2021; Kainz et al., 2015; Rousseau et al., 2006). Yet, these methods hinge on successful brain extraction which is challenging due to frequent artifacts and because the relatively small brain first needs to be localized within a wide FOV encompassing the maternal anatomy

(Gaudfernau et al., 2021). In addition, substantially fewer public fetal datasets are available for training in comparison to vast public adult brain datasets. This presents an ideal problem to be addressed with SynthStrip, as our approach synthesizes an endless stream of training data from only a handful of label maps.

## 6. Conclusion

The removal of non-brain signal from neuroimaging data is a fundamental first step for many quantitative analyses and its accuracy has a direct impact on downstream results. However, popular skull-stripping utilities are typically tailored to isotropic T1w scans and tend to fail, sometimes catastrophically, on images with other MRI contrasts or stack-of-slices acquisitions that are common in the clinic. We propose SynthStrip, a flexible tool that produces highly accurate brain masks across a landscape of imaging paradigms with widely varying contrast and resolution. We implement our method by leveraging anatomical label maps to synthesize a broad set of training images, optimizing a robust convolutional neural network that is agnostic to MRI contrasts and acquisition schemes.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

The authors thank Douglas Greve, Lilla Zöllei, and David Salat for sharing FSM, infant, and ASL data, respectively, and for providing a mechanism for distributing our reference dataset. Support for this research was provided in part by a BRAIN Initiative Cell Census Network grant [U01 MH117023], the National Institute of Biomedical Imaging and Bioengineering [P41 EB015896, R01 EB023281, R21 EB018907, R01 EB019956, P41 EB030006], the National Institute of Child Health and Human Development [K99 HD101553], the National Institute on Aging [R56 AG064027, R01 AG016495, R01 AG070988], the National Institute of Mental Health [RF1 MH121885, RF1 MH123195], the National Institute of Neurological Disorders and Stroke [R01 NS070963, R01 NS083534, R01 NS105820], and was made possible by the resources provided by Shared Instrumentation Grants [S10 RR023401, S10 RR019307, S10 RR023043]. Additional support was provided by the NIH Blueprint for Neuroscience Research [U01 MH093765], part of the multi-institutional Human Connectome Project. The research project benefitted from computational hardware generously provided by the Massachusetts Life Sciences Center (<https://www.masslifesciences.com>).

## Abbreviation:

**MRI**                      Magnetic resonance imaging

## References

- Andrade N, Faria FA, Cappabianco FAM, 2018. A practical review on medical image registration: From rigid to deep learning based approaches. In: 2018 31st SIBGRAPI Conference on Graphics, Patterns and Images (SIBGRAPI). IEEE, pp. 463–470.
- Arsigny V, Commowick O, Pennec X, Ayache N, 2006. A log-euclidean framework for statistics on diffeomorphisms. In: MICCAI: Medical Image Computing and Computer Assisted Interventions. Springer, pp. 924–931.
- Ashburner J, 2007. A fast diffeomorphic image registration algorithm. *Neuroimage* 38 (1), 95–113 . [PubMed: 17761438]
- Ashburner J, 2009. Preparing fMRI Data for Statistical Analysis. In: *fMRI techniques and protocols*. Springer, pp. 151–178.

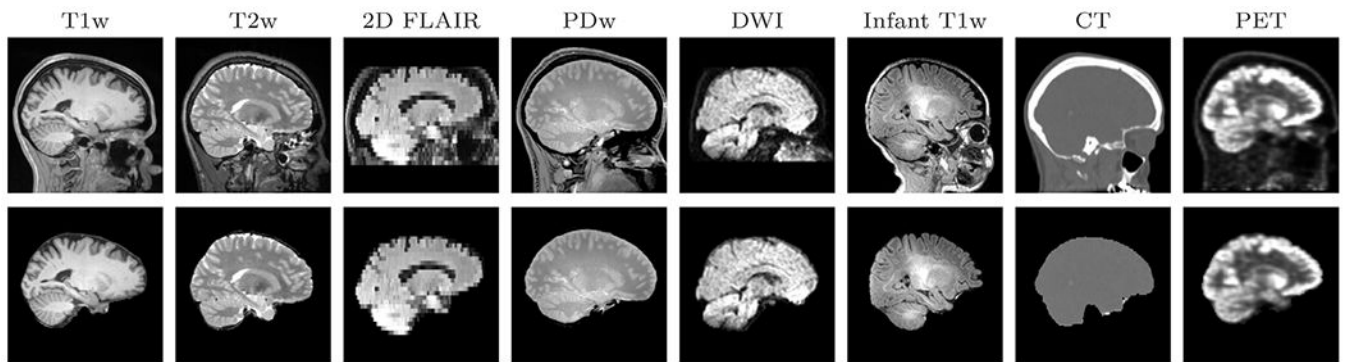
- Avants BB, Epstein CL, Grossman M, Gee JC, 2008. Symmetric diffeomorphic image registration with cross-correlation: evaluating automated labeling of elderly and neurodegenerative brain. *Med. Image Anal* 12 (1), 26–41. [PubMed: 17659998]
- Avants BB, Tustison NJ, Song G, Cook PA, Klein A, Gee JC, 2011. A reproducible evaluation of ANTs similarity metric performance in brain image registration. *Neuroimage* 54 (3), 2033–2044. [PubMed: 20851191]
- Billot B, Greve D, Van Leemput K, Fischl B, Iglesias JE, Dalca AV, 2020. A learning strategy for contrast-agnostic MRI segmentation. arXiv preprint arXiv:2003.01995.
- Bookheimer SY, Salat DH, Terpstra M, Ances BM, Barch DM, Buckner RL, Burgess GC, Curtiss SW, Diaz-Santos M, Elam JS, et al. , 2019. The lifespan human connectome project in aging: an overview. *Neuroimage* 185, 335–348. [PubMed: 30332613]
- Breiman L, 2001. Random forests. *Mach. Learn* 45 (1), 5–32.
- Clark K, Vendt B, Smith K, Freymann J, Kirby J, Koppel P, Moore S, Phillips S, Maffitt D, Pringle M, et al. , 2013. The cancer imaging archive (TCIA): maintaining and operating a public information repository. *J. Digit. Image* 26 (6), 1045–1057.
- Coupé P, Manjón JV, Fonov V, Pruessner J, Robles M, Collins DL, 2010. Non-local patch-based label fusion for hippocampus segmentation. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, pp. 129–136.
- Coupé P, Manjón JV, Fonov V, Pruessner J, Robles M, Collins DL, 2011. Patch-based segmentation using expert priors: application to hippocampus and ventricle segmentation. *Neuroimage* 54 (2), 940–954. [PubMed: 20851199]
- Cox RW, 1996. Afni: software for analysis and visualization of functional magnetic resonance neuroimages. *Comput. Biomed. Res* 29 (3), 162–173. [PubMed: 8812068]
- Cox RW, Jesmanowicz A, 1999. Real-time 3d image registration for functional MRI. *Magnet. Reson. Med* 42 (6), 1014–1018.
- Dai W, Garcia D, De Bazelaire C, Alsop DC, 2008. Continuous flow-driven inversion for arterial spin labeling using pulsed radio frequency and gradient fields. *Magnet. Reson. Med* 60 (6), 1488–1497.
- Dalca AV, Balakrishnan G, Guttag J, Sabuncu M, 2019. Unsupervised learning of probabilistic diffeomorphic registration for images and surfaces. *Med.I.A* 57, 226–236.
- Dalca AV, Guttag J, Sabuncu MR, 2018. Anatomical priors in convolutional networks for unsupervised biomedical segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9290–9299.
- Dice LR, 1945. Measures of the amount of ecologic association between species. *Ecology* 26 (3), 297–302.
- Ebner M, Wang G, Li W, Aertsen M, Patel PA, Aghwane R, Melbourne A, Doel T, Dymarkowski S, De Coppi P, et al. , 2020. An automated framework for localization, segmentation and super-resolution reconstruction of fetal brain MRI. *Neuroimage* 206, 116324. [PubMed: 31704293]
- Eskildsen SF, Coupé P, Fonov V, Manjón JV, Leung KK, Guizard N, Wassef SN, Østergaard LR, Collins DL, Initiative ADN, et al. , 2012. Beast: brain extraction based on nonlocal segmentation technique. *Neuroimage* 59 (3), 2362–2373. [PubMed: 21945694]
- Fatima A, Shahid AR, Raza B, Madni TM, Janjua UI, 2020. State-of-the-art traditional to the machine- and deep-learning-based skull stripping techniques, models, and algorithms. *J. Digit. Imaging* 33 (6), 1443–1464. [PubMed: 32666364]
- Fein G, Landman B, Tran H, Barakos J, Moon K, Di Sclafani V, Shumway R, 2006. Statistical parametric mapping of brain morphology: sensitivity is dramatically increased by using brain-extracted images as inputs. *Neuroimage* 30 (4), 1187–1195. [PubMed: 16442817]
- Fischl B, 2012. Freesurfer. *Neuroimage* 62 (2), 774–781. 20 YEARS OF fMRI [PubMed: 22248573]
- Fischl B, et al. , 2002. Whole brain segmentation: automated labeling of neuroanatomical structures in the human brain. *Neuron* 33 (3), 341–355. [PubMed: 11832223]
- Fischmeister FPS, Höllinger I, Klinger N, Geissler A, Wurnig MC, Matt E, Rath J, Robinson SD, Trattng S, Beisteiner R, 2013. The benefits of skull stripping in the normalization of clinical fMRI data. *NeuroImage: Clinical* 3, 369–380. [PubMed: 24273720]
- Friston KJ, Ashburner J, Frith CD, Poline J-B, Heather JD, Frackowiak RSJ, 1995. Spatial registration and normalization of images. *Hum. Brain Mapp* 3 (3), 165–189.

- Gaudfernau F, Blondiaux E, Allassonière S, 2021. Analysis of the Anatomical Variability of Fetal Brains with Corpus Callosum Agenesis. In: *Uncertainty for Safe Utilization of Machine Learning in Medical Imaging, and Perinatal Imaging, Placental and Preterm Image Analysis*. Springer, pp. 274–283.
- Greig DM, Porteous BT, Seheult AH, 1989. Exact maximum a posteriori estimation for binary images. *J. R. Stat. Soc.: Ser. B (Methodological)* 51 (2), 271–279.
- Greve DN, Billot B, Cordero D, Hoopes A, Hoffmann M, Dalca AV, Fischl B, Iglesias JE, Augustinack JC, 2021. A deep learning toolbox for automatic segmentation of subcortical limbic structures from MRI images. *Neuroimage* 244, 118610. [PubMed: 34571161]
- Harms MP, et al. , 2018. Extending the human connectome project across ages: imaging protocols for the lifespan development and aging projects. *Neuroimage* 183, 972–984. [PubMed: 30261308]
- Hendrycks D, Basart S, Mu N, Kadavath S, Wang F, Dorundo E, Desai R, Zhu T, Parajuli S, Guo M, Song D, Steinhardt J, Gilmer J, 2021. The many faces of robustness: A critical analysis of out-of-distribution generalization. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 8340–8349.
- Hoffmann M, Abaci Turk E, Gagoski B, Morgan L, Wighton P, Tisdall MD, Reuter M, Adalsteinsson E, Grant PE, Wald LL, et al. , 2021. Rapid head-pose detection for automated slice prescription of fetal-brain MRI. *Int. J. Imag. Syst. Technol* 31 (3), 1136–1154.
- Hoffmann M, Billot B, Iglesias JE, Fischl B, Dalca AV, 2021. Learning mri contrast-agnostic registration. In: *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)*. IEEE, pp. 899–903.
- Hoffmann M, Carpenter TA, Williams GB, Sawiak SJ, 2015. A survey of patient motion in disorders of consciousness and optimization of its retrospective correction. *Magn. Reson. Imaging* 33 (3), 346–350. [PubMed: 25485789]
- Hoffmann M, Frost R, Salat D, Tisdall MD, Polimeni J, van der Kouwe A, 2020. Real-time brain masking algorithm improves motion tracking accuracy in scans with volumetric navigators (vNavs). In: *International Society for Magnetic Resonance in Medicine. ISMRM*, p. 3367.
- Holdsworth SJ, Aksoy M, Newbould RD, Yeom K, Van AT, Ooi MB, Barnes PD, Bammer R, Skare S, 2012. Diffusion tensor imaging (DTI) with retrospective motion correction for large-scale pediatric imaging. *J. Magn. Reson. Imaging* 36 (4), 961–971. [PubMed: 22689498]
- Hsu L-M, Wang S, Ranadive P, Ban W, Chao T-HH, Song S, Cerri DH, Walton LR, Broadwater MA, Lee S-H, et al. , 2020. Automatic skull stripping of rat and mouse brain MRI data using u-net. *Front. Neurosci* 14, 568614. [PubMed: 33117118]
- Hwang H, Rehman HZU, Lee S, 2019. 3D u-net for skull stripping in brain MRI. *Appl. Sci* 9 (3), 569.
- Iglesias JE, Billot B, Balbastre Y, Tabari A, Conklin J, González RG, Alexander DC, Golland P, Edlow BL, Fischl B, et al. , 2021. Joint super-resolution and synthesis of 1 mm isotropic MP-RAGE volumes from clinical MRI exams with scans of different orientation, resolution and contrast. *Neuroimage* 237, 118206. [PubMed: 34048902]
- Iglesias JE, Liu C-Y, Thompson PM, Tu Z, 2011. Robust brain extraction across datasets and comparison with publicly available methods. *IEEE Trans. Med. Imag* 30 (9), 1617–1634.
- Jenkinson M, Bannister P, Brady M, Smith S, 2002. Improved optimization for the robust and accurate linear registration and motion correction of brain images. *Neuroimage* 17 (2), 825–841. [PubMed: 12377157]
- Jenkinson M, Beckmann CF, Behrens TEJ, Woolrich MW, Smith SM, 2012. *Fsl*. *Neuroimage* 62 (2), 782–790. [PubMed: 21979382]
- Jenkinson M, Smith S, 2001. A global optimisation method for robust affine registration of brain images. *Med. Image. Anal* 5 (2), 143–156. [PubMed: 11516708]
- Jiang A, Kennedy DN, Baker JR, Weisskoff RM, Tootell RBH, Woods RP, Benson RR, Kwong KK, Brady TJ, Rosen BR, et al. , 1995. Motion detection and correction in functional MR imaging. *Hum. Brain Mapp* 3 (3), 224–235.
- Jog A, Hoopes A, Greve DN, Van Leemput K, Fischl B, 2019. Psacnn: pulse sequence adaptive fast whole brain segmentation. *Neuroimage* 199, 553–569. [PubMed: 31129303]
- Jones DK, Leemans A, 2011. Diffusion Tensor Imaging. In: *Magnetic resonance neuroimaging*. Springer, pp. 127–144.

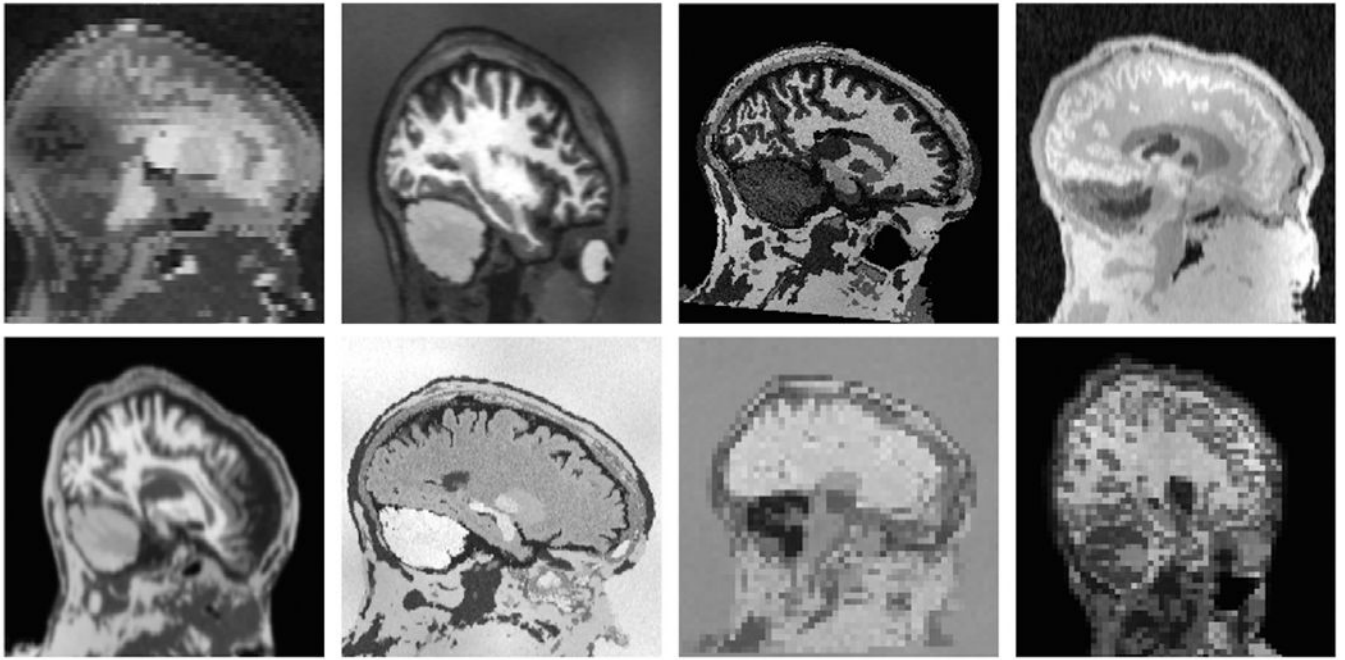


- Kainz B, Steinberger M, Wein W, Kuklisova-Murgasova M, Malamateniou C, Keraudren K, Torsney-Weir T, Rutherford M, Aljabar P, Hajnal JV, et al. , 2015. Fast volume reconstruction from motion corrupted stacks of 2d slices. *IEEE Trans. Med. Imaging* 34 (9), 1901–1913. [PubMed: 25807565]
- Karani N, Chaitanya K, Baumgartner C, Konukoglu E, 2018. A lifelong learning approach to brain MR segmentation across scanners and protocols. In: *Medical Image Computing and Computer Assisted Intervention – MICCAI 2018*. Springer International Publishing, Cham, pp. 476–484.
- Kingma DP, Ba J, 2014. Adam: a method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- 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. [PubMed: 26808333]
- Klein A, Andersson J, Ardekani BA, Ashburner J, Avants B, Chiang M-C, Christensen GE, Collins DL, Gee J, Hellier P, et al. , 2009. Evaluation of 14 nonlinear deformation algorithms applied to human brain MRI registration. *Neuroimage* 46 (3), 786–802. [PubMed: 19195496]
- van der Kouwe AJW, Benner T, Salat DH, Fischl B, 2008. Brain morphometry with multiecho MPRAGE. *Neuroimage* 40 (2), 559–569. [PubMed: 18242102]
- Lucena O, Souza R, Rittner L, Frayne R, Lotufo R, 2019. Convolutional neural networks for skull-stripping in brain MR imaging using silver standard masks. *Artif. Intell. Med* 98, 48–58. [PubMed: 31521252]
- de Macedo Rodrigues K, Ben-Avi E, Sliva DD, Choe M.-s., Drottar M, Wang R, Fischl B, Grant PE, Zöllei L, 2015. A freesurfer-compliant consistent manual segmentation of infant brains spanning the 0–2 year age range. *Front. Hum. Neurosci* 9, 21. [PubMed: 25741260]
- Mamonov AB, Kalpathy-Cramer J, 2016. Data from QIN GBM treatment response. 10.7937/k9/tcia.2016.nQF4gpn2
- Marques JP, Kober T, Krueger G, van der Zwaag W, Van de Moortele P-F, Gruetter R, 2010. MP2Rage, a self bias-field corrected sequence for improved segmentation and T1-mapping at high field. *Neuroimage* 49 (2), 1271–1281. [PubMed: 19819338]
- Mérida I, Jung J, Bouvard S, Le Bars D, Lancelot S, Lavenne F, Bouillot C, Redouté J, Hammers A, Costes N, 2021. Cermep-idb-mrxfdg: a database of 37 normal adult human brain [18f] fdg pet, t1 and flair mri, and ct images available for research. *EJNMMI Res.* 11 (1), 1–10. [PubMed: 33394212]
- Milletari F, Navab N, Ahmadi S-A, 2016. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In: *3DV*, pp. 565–571.
- Modat M, Cash DM, Daga P, Winston GP, Duncan JS, Ourselin S, 2014. Global image registration using a symmetric block-matching approach. *J. Med. Imaging* 1 (2), 024003.
- Mugler III JP, Brookeman JR, 1990. Three-dimensional magnetization-prepared rapid gradient-echo imaging (3D MP RAGE). *Magn. Reson. Med* 15 (1), 152–157. doi:10.1002/mrm.1910150117. [PubMed: 2374495]
- Ou Y, Akbari H, Bilello M, Da X, Davatzikos C, 2014. Comparative evaluation of registration algorithms in different brain databases with varying difficulty: results and insights. *IEEE Trans. Med. Imaging* 33 (10), 2039–2065. [PubMed: 24951685]
- Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, Killeen T, Lin Z, Gimelshein N, Antiga L, Desmaison A, Kopf A, Yang E, DeVito Z, Raison M, Tejani A, Chilamkurthy S, Steiner B, Fang L, Bai J, Chintala S, 2019. Pytorch: An Imperative Style, High-performance Deep Learning Library. In: *Advances in Neural Information Processing Systems* 32, pp. 8024–8035.
- Prah MA, Stuffelbeam SM, Paulson ES, Kalpathy-Cramer J, Gerstner ER, Batchelor TT, Barboriak DP, Rosen BR, Schmainda KM, 2015. Repeatability of standardized and normalized relative CBV in patients with newly diagnosed glioblastoma. *Am. J. Neuroradiol* 36 (9), 1654–1661. [PubMed: 26066626]
- Puonti O, Iglesias JE, Leemput KV, 2016. Fast and sequence-adaptive whole-brain segmentation using parametric bayesian modeling. *Neuroimage* 143, 235–249. [PubMed: 27612647]
- Reuter M, Rosas HD, Fischl B, 2010. Highly accurate inverse consistent registration: a robust approach. *Neuroimage* 53 (4), 1181–1196. [PubMed: 20637289]

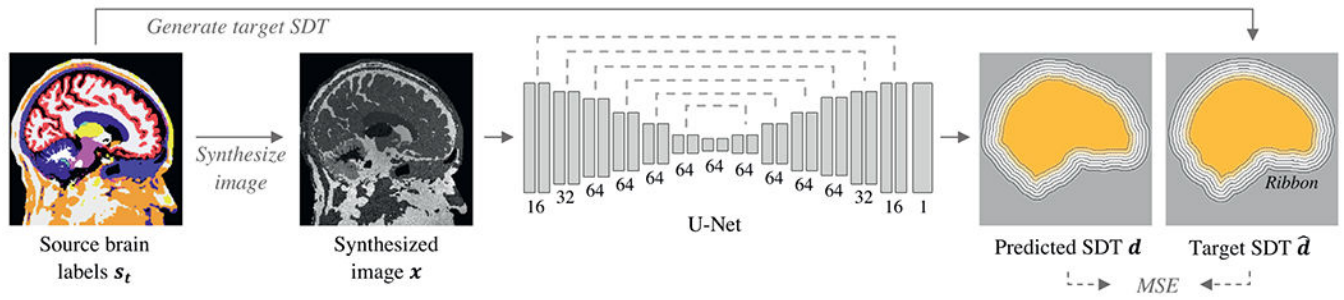
- Ronneberger O, Fischer P, Brox T, 2015. U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. Springer, pp. 234–241.
- Rousseau F, Glenn OA, Iordanova B, Rodriguez-Carranza C, Vigneron DB, Barkovich JA, Studholme C, 2006. Registration-based approach for reconstruction of high-resolution in utero fetal MR brain images. *Acad. Radiol* 13 (9), 1072–1081 . [PubMed: 16935719]
- Roy S, Butman JA, Pham DL, Initiative ADN, et al. , 2017. Robust skull stripping using multiple MR image contrasts insensitive to pathology. *Neuroimage* 146, 132–147. [PubMed: 27864083]
- Rueckert D, Sonoda LI, Hayes C, Hill DLG, Leach MO, Hawkes DJ, 1999. Non-rigid registration using free-form deformations: application to breast MR images. *IEEE Trans. Med. Image* 18 (8), 712–721.
- Sadanathan SA, Zheng W, Chee MWL, Zagorodnov V, 2010. Skull stripping using graph cuts. *Neuroimage* 49 (1), 225–239. [PubMed: 19732839]
- Salehi SSM, Erdogmus D, Gholipour A, 2017. Auto-context convolutional neural network (auto-net) for brain extraction in magnetic resonance imaging. *IEEE Trans. Med. Imaging* 36 (11), 2319–2330. [PubMed: 28678704]
- Ségonne F, Dale AM, Busa E, Glessner M, Salat D, Hahn HK, Fischl B, 2004. A hybrid approach to the skull stripping problem in MRI. *Neuroimage* 22 (3), 1060–1075. [PubMed: 15219578]
- Shattuck DW, Sandor-Leahy SR, Schaper KA, Rottenberg DA, Leahy RM, 2001. Magnetic resonance image tissue classification using a partial volume model. *Neuroimage* 13 (5), 856–876. [PubMed: 11304082]
- Smith SM, 2002. Fast robust automated brain extraction. *Hum. Brain Mapp* 17 (3), 143–155. [PubMed: 12391568]
- Tu Z, Bai X, 2009. Auto-context and its application to high-level vision tasks and 3d brain image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell* 32 (10), 1744–1757.
- Vercauteren T, Pennec X, Perchant A, Ayache N, 2009. Diffeomorphic demons: efficient non-parametric image registration. *Neuroimage* 45 (1), S61–S72. [PubMed: 19041946]
- Zöllei L, Iglesias JE, Ou Y, Grant PE, Fischl B, 2020. Infant freesurfer: an automated segmentation and surface extraction pipeline for t1-weighted neuroimaging data of infants 0–2 years. *Neuroimage* 218, 116946. [PubMed: 32442637]



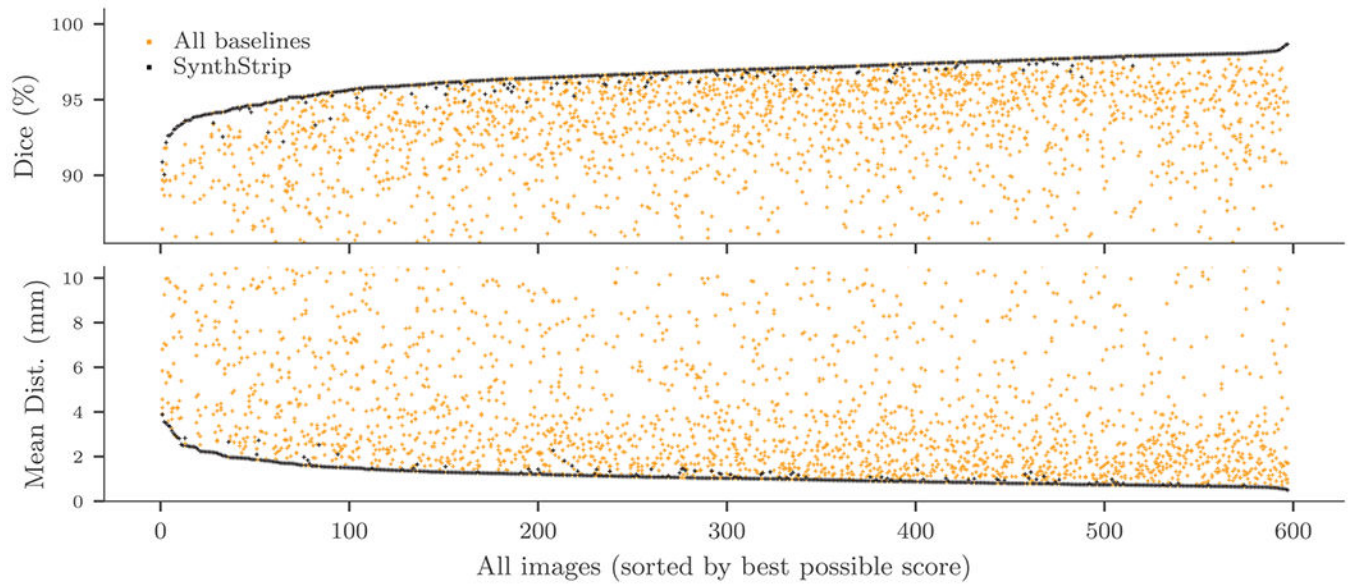
**Fig. 1.** Examples of SynthStrip brain extractions (bottom) for a wide range of image acquisitions and modalities (top). Powered by a strategy for synthesizing diverse training data, SynthStrip learns to skull-strip brain images of any type.



**Fig. 2.** Samples of synthetic images used for SynthStrip training. To encourage the network to generalize, we synthesize images that far exceed the realistic range of whole-brain acquisitions. In this figure, each brain image is generated from the same label map. In practice, we use label maps from several different subjects.

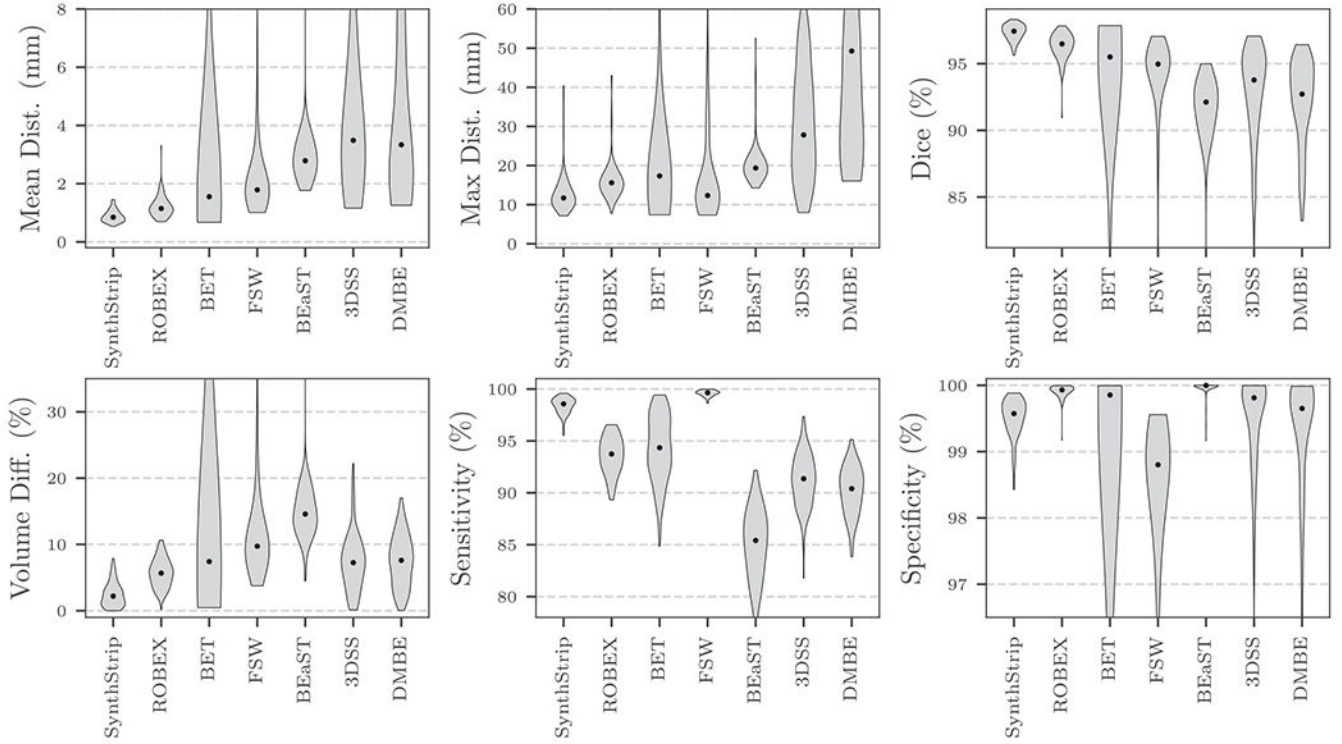


**Fig. 3.** SynthStrip training framework. At every optimization step, we sample a randomly transformed brain segmentation  $s_b$  from which we synthesize a gray-scale image  $x$  with arbitrary contrast. The skull-stripping 3D U-Net receives  $x$  as input and predicts a thresholded signed distance transform (SDT)  $d$  representing the distance of each voxel to the skull boundary. The U-Net consists of skip-connected, multi-resolution convolutional layers illustrated by gray bars, with their number of output filters indicated below. We train SynthStrip in a supervised fashion, maximizing the similarity between  $d$  and the ground-truth SDT  $\hat{a}$  within a ribbon of set distance around the brain and derived directly from the segmentation labels of  $s_t$ .

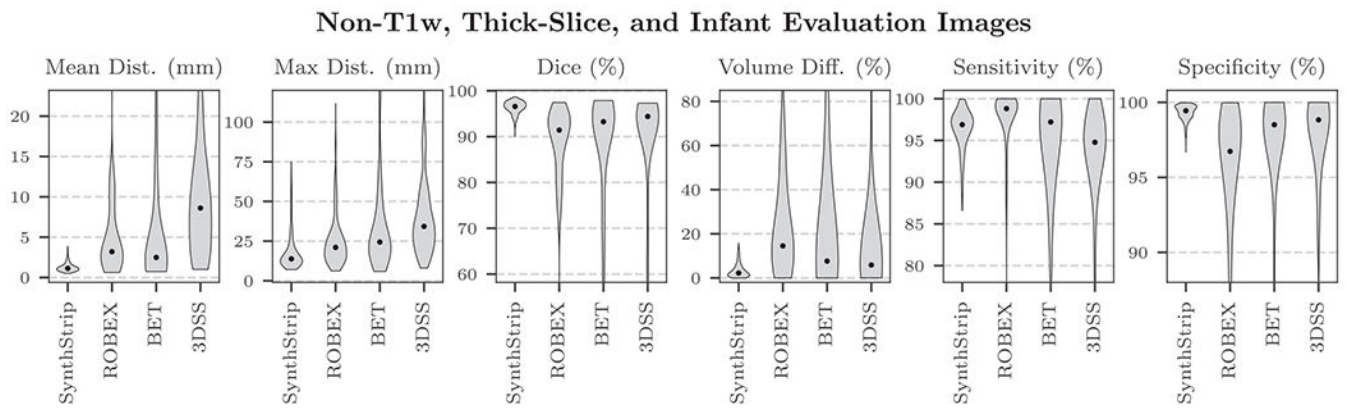


**Fig. 4.** SynthStrip accuracy compared to baseline methods, across all images in the test set. Images are sorted by the score of the top performing skull-stripping method. Each dot represents a single brain mask derived with a particular tool, and each column of dots represents the scores obtained for a single image across tools. See Supplementary Fig. S2 for a version showing each baseline in a different color.

T1w Near-Isotropic Evaluation Images

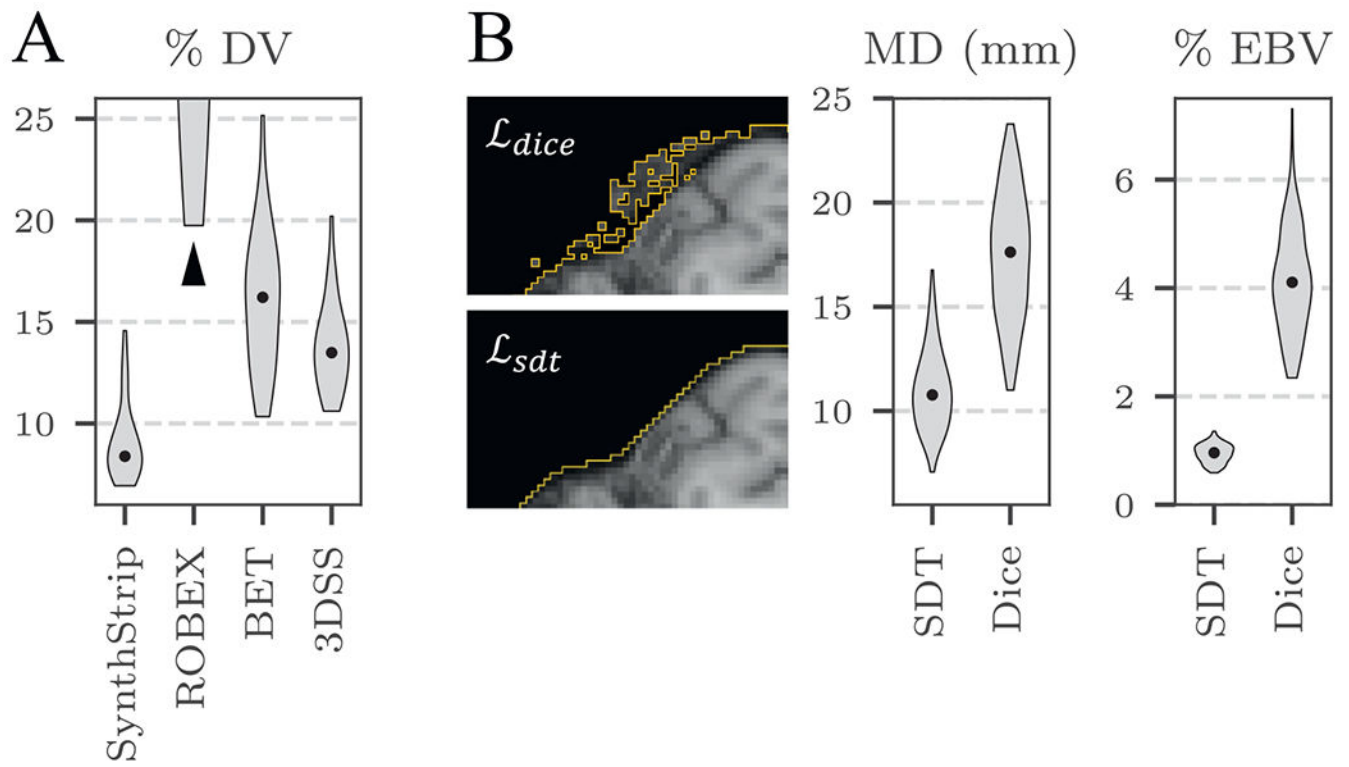


**Fig. 5.** SynthStrip and baseline skull-stripping performance for near-isotropic, T1w adult MR brain images. Median scores are represented by black dots. For all metrics except sensitivity and specificity, SynthStrip yields optimal brain masks. The high specificity achieved by ROBEX and BEaST comes at the cost of substantial under-segmentation of the brain mask, as indicated by their low sensitivity scores. The inverse is true for FSW, which tends to substantially over-segment the brain. Black dots indicate median scores.

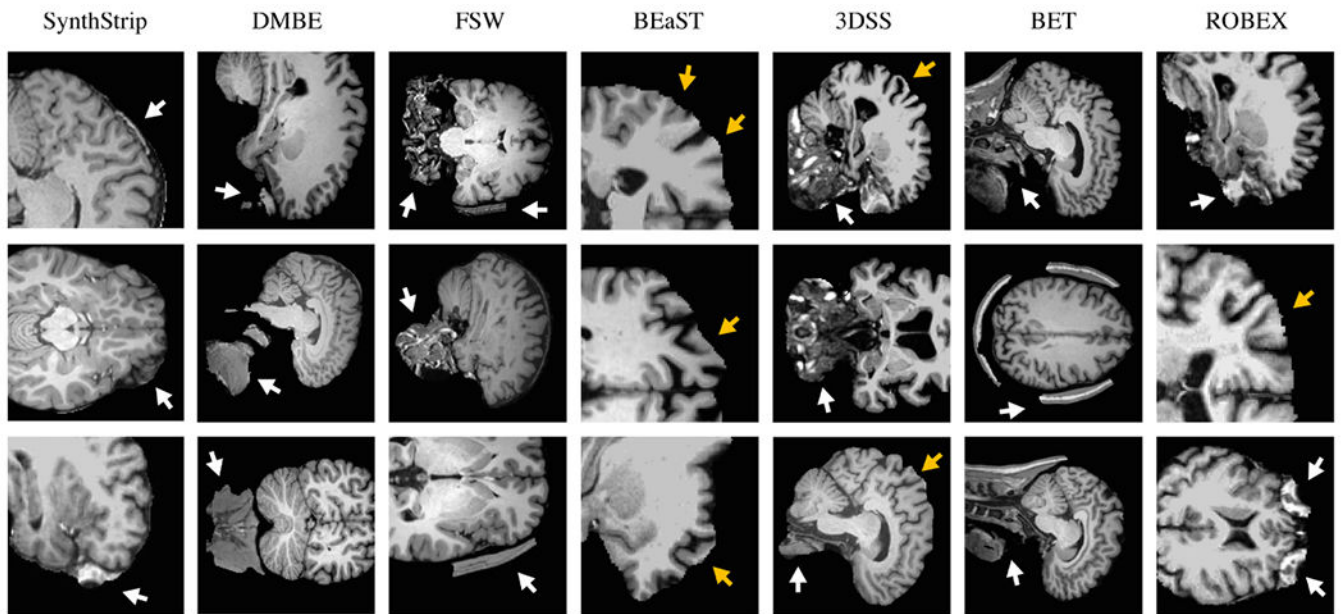


**Fig. 6.** Considering all non-T1w, thick-slice, and infant images in the evaluation set, SynthStrip surpasses baseline accuracy by a wide margin. In this figure, we include only baselines that generalize to acquisition protocols and modalities beyond the common structural T1w MRI scans. Black dots indicate median scores.



**Fig. 7.**

**A:** SynthStrip variability across time-series data, measured by percent of discordant voxel locations (DV) across diffusion-encoded directions, relative to the brain mask volume. The ROBEX median % DV extends beyond the chart axis, as indicated by the black arrow. **B:** Effect of SDT- and Dice-based loss functions during training. A SynthStrip model trained using  $\mathcal{L}_{sdt}$  predicts substantially smoother brain masks (boundaries indicated in orange) than a model trained with  $\mathcal{L}_{dice}$ , resulting in considerably lower maximum surface distance (MD) to ground truth masks and percent of exposed boundary voxels (EBV).



**Fig. 8.** Representative skull-stripping errors for SynthStrip and baseline methods. White arrows indicate over-labeling of the brain mask, while orange arrows indicate removal of brain matter. SynthStrip errors are uncommon and typically involve including small regions of dura or other extracerebral tissue in the brain mask, if they occur..

**Table 1**

Uniform hyperparameter sampling ranges used for synthesizing a training image from a source segmentation map. The specific values were chosen by visual inspection of the generated images to produce a landscape of image contrasts, anatomies, and acquisition characteristics that far exceed the realistic range of medical images. We sample fields with isotropic voxels of the indicated side length, where SD abbreviates standard deviation.

| <b>Synthesis hyperparameter</b>   | <b>Uniform sampling range</b> |
|-----------------------------------|-------------------------------|
| Affine translation                | 0–50 <i>mm</i>                |
| Affine rotation                   | 0–45°                         |
| Affine scaling                    | 80–120%                       |
| Deformation voxel length          | 8–16 <i>mm</i>                |
| Deformation SD                    | 0–3 <i>mm</i>                 |
| Label intensity mean              | 0–1                           |
| Label intensity SD                | 0–0.1                         |
| Bias field voxel length           | 4–64 <i>mm</i>                |
| Bias field SD                     | 0–0.5                         |
| Exponentiation parameter $\gamma$ | –0.25–0.25                    |
| FOV cropping (any axis)           | 0–50 <i>mm</i>                |
| Down-sample factor $r$ (any axis) | 1–5                           |

**Table 2**

We employ a diverse set of acquired evaluation data, spanning across imaging modalities, MRI contrasts, and resolution (res.), where 2D indicates stacks of slice-wise acquisitions. Each individual dataset is divided into a small validation (val.) and a larger test set. For further details see Section 3.4.2.

| Dataset | Modality        | Res. ( $mm^3$ ) | Val. | Test |
|---------|-----------------|-----------------|------|------|
| IXI     | T1w MRI         | 0.9×0.9×1.2     | 0    | 48   |
|         | T2w MRI         | 0.9×0.9×1.2     | 2    | 48   |
|         | PDw MRI         | 0.9×0.9×1.2     | 2    | 48   |
|         | MRA             | 0.5×0.5×0.8     | 2    | 48   |
|         | DWI             | 1.8×1.8×2.0     | 0    | 32   |
| FSM     | T1w MPRAGE      | 1.0×1.0×1.0     | 0    | 38   |
|         | T2w 3D-SPACE    | 1.0×1.0×1.0     | 2    | 34   |
|         | PDw 3D-FLASH    | 1.0×1.0×1.0     | 2    | 30   |
|         | qT1 MP2RAGE     | 1.0×1.0×1.0     | 2    | 30   |
| ASL     | T1w MPRAGE      | 1.0×1.0×1.0     | 2    | 41   |
|         | PCASL 2D-EPI    | 3.4×3.4×5.0     | 2    | 41   |
| QIN     | T1w 2D-FLASH    | 0.4×0.4×6.0     | 2    | 52   |
|         | T2-FLAIR 2D-FSE | 0.4×0.4×6.0     | 2    | 15   |
|         | T2w 2D-FSE      | 1.0×1.0×5.0     | 2    | 37   |
| Infant  | T1w MPRAGE      | 1.0×1.0×1.0     | 0    | 16   |
| CIM     | FDG PET         | 2.0×2.0×2.0     | 0    | 20   |
|         | CT              | 0.6×0.6×1.5     | 0    | 20   |

SynthStrip and baseline method accuracy across datasets, as measured by the mean surface distance ( $\pm$  SD) between computed and ground-truth binary brain masks.  $p$ -values comparing baseline with SynthStrip results are presented below each score. SynthStrip stands out as a dominant skull-stripping technique, significantly outperforming baselines for nearly every dataset with the exception of those with  $p$ -values in bold, for which  $p > 0.05$ . FSW fails entirely for multiple subsets of non-T1w images (metrics not shown).

**Table 3**

|           | Mean surface distance (mm) |                       |  |                       |                       |                       |                       |  |
|-----------|----------------------------|-----------------------|--|-----------------------|-----------------------|-----------------------|-----------------------|--|
|           | SynthStrip                 | ROBEX                 | BET                                    | 3DSS                  | BEaST                 | FSW                   | DMBE                  |  |
| IXI T1w   | 1.0 $\pm$ 0.2              | 1.2 $\pm$ 0.3         | 1.4 $\pm$ 1.5                          | 8.9 $\pm$ 6.3         | 2.3 $\pm$ 0.3         | 2.9 $\pm$ 2.2         | 3.6 $\pm$ 3.5         |  |
|           |                            | 2.3 $\times 10^{-6}$  | <b>1.1 <math>\times 10^{-1}</math></b> | 8.6 $\times 10^{-12}$ | 7.3 $\times 10^{-33}$ | 1.2 $\times 10^{-7}$  | 5.4 $\times 10^{-6}$  |  |
| FSM T1w   | 0.7 $\pm$ 0.1              | 1.4 $\pm$ 0.3         | 18.8 $\pm$ 6.7                         | 3.0 $\pm$ 1.3         | 2.8 $\pm$ 0.2         | 2.5 $\pm$ 2.0         | 8.2 $\pm$ 5.0         |  |
|           |                            | 5.7 $\times 10^{-17}$ | 1.2 $\times 10^{-18}$                  | 3.7 $\times 10^{-12}$ | 7.2 $\times 10^{-36}$ | 6.8 $\times 10^{-6}$  | 4.6 $\times 10^{-11}$ |  |
| ASL T1w   | 0.9 $\pm$ 0.2              | 1.0 $\pm$ 0.4         | 1.7 $\pm$ 0.4                          | 3.4 $\pm$ 1.5         | 3.9 $\pm$ 2.2         | 1.6 $\pm$ 0.4         | 4.3 $\pm$ 2.3         |  |
|           |                            | 3.1 $\times 10^{-2}$  | 6.0 $\times 10^{-16}$                  | 4.4 $\times 10^{-14}$ | 9.5 $\times 10^{-12}$ | 1.6 $\times 10^{-11}$ | 2.0 $\times 10^{-12}$ |  |
| QIN T1w   | 1.1 $\pm$ 0.4              | 2.8 $\pm$ 2.4         | 2.3 $\pm$ 1.2                          | 11.7 $\pm$ 5.6        | 5.8 $\pm$ 7.8         | 9.8 $\pm$ 8.0         | 3.7 $\pm$ 2.2         |  |
|           |                            | 9.5 $\times 10^{-7}$  | 1.7 $\times 10^{-11}$                  | 2.3 $\times 10^{-19}$ | 3.9 $\times 10^{-5}$  | 1.7 $\times 10^{-10}$ | 4.0 $\times 10^{-12}$ |  |
| IXI T2w   | 1.2 $\pm$ 0.3              | 3.4 $\pm$ 1.3         | 3.2 $\pm$ 2.3                          | 9.6 $\pm$ 5.1         | 20.5 $\pm$ 12.4       | -                     | 57.2 $\pm$ 19.7       |  |
|           |                            | 1.8 $\times 10^{-17}$ | 1.6 $\times 10^{-7}$                   | 3.0 $\times 10^{-15}$ | 1.8 $\times 10^{-14}$ | -                     | 4.8 $\times 10^{-25}$ |  |
| FSM T2w   | 0.8 $\pm$ 0.1              | 2.6 $\pm$ 0.8         | 1.6 $\pm$ 0.7                          | 1.9 $\pm$ 0.9         | 14.7 $\pm$ 10.6       | -                     | 72.4 $\pm$ 24.2       |  |
|           |                            | 1.7 $\times 10^{-15}$ | 5.7 $\times 10^{-8}$                   | 8.5 $\times 10^{-9}$  | 3.6 $\times 10^{-9}$  | -                     | 6.4 $\times 10^{-19}$ |  |
| QIN T2w   | 1.6 $\pm$ 0.8              | 4.6 $\pm$ 2.5         | 3.9 $\pm$ 2.1                          | 14.9 $\pm$ 9.6        | 16.8 $\pm$ 9.3        | -                     | 11.8 $\pm$ 4.8        |  |
|           |                            | 1.9 $\times 10^{-9}$  | 9.1 $\times 10^{-8}$                   | 5.9 $\times 10^{-9}$  | 7.2 $\times 10^{-12}$ | -                     | 2.8 $\times 10^{-16}$ |  |
| QIN FLAIR | 1.0 $\pm$ 0.2              | 2.1 $\pm$ 0.7         | 1.2 $\pm$ 0.3                          | 9.9 $\pm$ 5.9         | 3.4 $\pm$ 1.9         | 7.9 $\pm$ 3.1         | 4.5 $\pm$ 1.2         |  |
|           |                            | 3.1 $\times 10^{-5}$  | 5.3 $\times 10^{-3}$                   | 1.6 $\times 10^{-5}$  | 1.2 $\times 10^{-4}$  | 1.4 $\times 10^{-7}$  | 1.5 $\times 10^{-9}$  |  |
| IXI PDw   | 1.2 $\pm$ 0.4              | 1.9 $\pm$ 0.6         | 1.6 $\pm$ 0.6                          | 9.8 $\pm$ 5.5         | 11.2 $\pm$ 11.2       | 8.5 $\pm$ 5.8         | 5.1 $\pm$ 2.4         |  |
|           |                            | 5.9 $\times 10^{-15}$ | 3.0 $\times 10^{-7}$                   | 1.4 $\times 10^{-14}$ | 1.2 $\times 10^{-7}$  | 5.7 $\times 10^{-11}$ | 2.4 $\times 10^{-15}$ |  |
| FSM PDw   | 1.0 $\pm$ 0.2              | 1.5 $\pm$ 0.4         | 1.8 $\pm$ 3.7                          | 1.6 $\pm$ 0.6         | 4.3 $\pm$ 3.3         | 20.0 $\pm$ 7.1        | 17.9 $\pm$ 3.5        |  |
|           |                            | 7.3 $\times 10^{-6}$  | <b>2.3 <math>\times 10^{-1}</math></b> | 8.9 $\times 10^{-7}$  | 4.5 $\times 10^{-6}$  | 1.2 $\times 10^{-15}$ | 5.6 $\times 10^{-23}$ |  |
| IXI MRA   | 1.3 $\pm$ 0.4              | 10.7 $\pm$ 2.8        | 16.2 $\pm$ 6.5                         | 9.1 $\pm$ 4.0         | 4.2 $\pm$ 2.4         | -                     | 34.8 $\pm$ 17.1       |  |
|           |                            | 3.0 $\times 10^{-28}$ | 7.8 $\times 10^{-21}$                  | 4.3 $\times 10^{-18}$ | 9.8 $\times 10^{-11}$ | -                     | 2.4 $\times 10^{-18}$ |  |
| FSM qT1   | 0.8 $\pm$ 0.1              | 21.8 $\pm$ 14.8       | 33.0 $\pm$ 10.2                        | 33.2 $\pm$ 3.0        | 31.9 $\pm$ 23.7       | 26.0 $\pm$ 11.9       | 54.0 $\pm$ 17.7       |  |
|           |                            | 1.2 $\times 10^{-8}$  | 1.2 $\times 10^{-17}$                  | 2.3 $\times 10^{-32}$ | 5.4 $\times 10^{-8}$  | 1.4 $\times 10^{-12}$ | 4.2 $\times 10^{-17}$ |  |

|            | Mean surface distance (mm) |                              |                         |                              |                         |                         |                         |  |
|------------|----------------------------|------------------------------|-------------------------|------------------------------|-------------------------|-------------------------|-------------------------|--|
|            | SynthStrip                 | ROBEX                        | BET                     | 3DSS                         | BEaST                   | FSW                     | DMBE                    |  |
| ASL EPI    | 1.6 ± 0.6                  | 4.8 ± 2.3                    | 2.0 ± 0.7               | 9.6 ± 4.0                    | 14.8 ± 11.8             | 16.1 ± 5.4              | 2.6 ± 0.6               |  |
|            |                            | 1.6 × 10 <sup>-10</sup>      | 9.9 × 10 <sup>-5</sup>  | 9.8 × 10 <sup>-16</sup>      | 9.8 × 10 <sup>-9</sup>  | 5.0 × 10 <sup>-20</sup> | 2.5 × 10 <sup>-15</sup> |  |
| Infant T1w | 1.0 ± 0.3                  | 4.1 ± 6.1                    | 14.3 ± 10.7             | 22.2 ± 12.1                  | 19.0 ± 21.7             | 17.6 ± 15.5             | 6.6 ± 3.4               |  |
|            |                            | <b>5.9 × 10<sup>-2</sup></b> | 2.1 × 10 <sup>-4</sup>  | 5.3 × 10 <sup>-6</sup>       | 5.5 × 10 <sup>-3</sup>  | 1.3 × 10 <sup>-3</sup>  | 1.6 × 10 <sup>-5</sup>  |  |
| IXI DWI    | 1.6 ± 0.6                  | 6.2 ± 3.0                    | 2.4 ± 0.9               | 8.5 ± 2.7                    | 11.1 ± 9.3              | 11.2 ± 4.7              | 6.8 ± 1.2               |  |
|            |                            | 2.5 × 10 <sup>-9</sup>       | 2.7 × 10 <sup>-11</sup> | 9.8 × 10 <sup>-15</sup>      | 2.5 × 10 <sup>-6</sup>  | 2.3 × 10 <sup>-12</sup> | 5.1 × 10 <sup>-27</sup> |  |
| CIM PET    | 1.5 ± 0.4                  | 3.9 ± 2.2                    | 9.3 ± 4.0               | 2.2 ± 2.9                    | 69.0 ± 21.2             | 16.2 ± 3.5              | 17.6 ± 6.4              |  |
|            |                            | 7.7 × 10 <sup>-5</sup>       | 8.0 × 10 <sup>-8</sup>  | <b>2.4 × 10<sup>-1</sup></b> | 2.6 × 10 <sup>-10</sup> | 5.8 × 10 <sup>-13</sup> | 1.0 × 10 <sup>-9</sup>  |  |
| CIM CT     | 2.0 ± 0.4                  | 11.4 ± 1.2                   | 34.1 ± 3.7              | 20.6 ± 2.0                   | 74.8 ± 18.4             | 29.1 ± 5.9              | 34.7 ± 8.2              |  |
|            |                            | 4.1 × 10 <sup>-20</sup>      | 3.6 × 10 <sup>-19</sup> | 4.1 × 10 <sup>-20</sup>      | 7.7 × 10 <sup>-12</sup> | 2.9 × 10 <sup>-12</sup> | 5.5 × 10 <sup>-13</sup> |  |

Skull-stripping accuracy across datasets, as measured by the mean Dice overlap ( $\pm$  SD) between computed and ground-truth binary brain masks.  $p$ -values comparing baseline with SynthStrip results are presented below each score. Across each dataset, SynthStrip significantly outperforms most baselines except those with  $p$ -values in bold, for which  $p > 0.05$ .

**Table 4**

|           | SynthStrip     | ROBEX                 | BET                                    | 3DSS                  | BE <sub>a</sub> ST    | FSW                   | DMBE                  |
|-----------|----------------|-----------------------|--|-----------------------|-----------------------|-----------------------|-----------------------|
| IXI T1w   | 97.0 $\pm$ 0.5 | 96.2 $\pm$ 0.8        | 96.1 $\pm$ 3.1                         | 95.4 $\pm$ 1.4        | 93.4 $\pm$ 0.8        | 92.6 $\pm$ 4.4        | 93.7 $\pm$ 3.0        |
|           |                | 4.5 $\times 10^{-8}$  | <b>5.1 <math>\times 10^{-2}</math></b> | 1.3 $\times 10^{-9}$  | 1.1 $\times 10^{-32}$ | 5.8 $\times 10^{-9}$  | 8.0 $\times 10^{-10}$ |
| FSM T1w   | 97.8 $\pm$ 0.3 | 95.9 $\pm$ 0.7        | 65.8 $\pm$ 11.2                        | 92.0 $\pm$ 3.0        | 92.1 $\pm$ 0.7        | 93.8 $\pm$ 3.7        | 89.5 $\pm$ 3.3        |
|           |                | 3.3 $\times 10^{-18}$ | 1.7 $\times 10^{-19}$                  | 1.4 $\times 10^{-13}$ | 1.9 $\times 10^{-35}$ | 2.2 $\times 10^{-7}$  | 5.6 $\times 10^{-18}$ |
| ASL T1w   | 97.3 $\pm$ 0.5 | 96.8 $\pm$ 1.1        | 94.9 $\pm$ 1.1                         | 90.8 $\pm$ 3.5        | 89.0 $\pm$ 5.5        | 95.5 $\pm$ 0.9        | 92.7 $\pm$ 1.7        |
|           |                | 3.0 $\times 10^{-3}$  | 5.7 $\times 10^{-17}$                  | 1.5 $\times 10^{-15}$ | 7.7 $\times 10^{-13}$ | 1.1 $\times 10^{-11}$ | 1.1 $\times 10^{-19}$ |
| QIN T1w   | 96.3 $\pm$ 0.8 | 92.8 $\pm$ 4.6        | 93.8 $\pm$ 2.7                         | 92.4 $\pm$ 3.1        | 85.2 $\pm$ 17.1       | 79.9 $\pm$ 13.3       | 89.0 $\pm$ 7.3        |
|           |                | 3.7 $\times 10^{-7}$  | 5.2 $\times 10^{-11}$                  | 2.3 $\times 10^{-15}$ | 1.5 $\times 10^{-5}$  | 5.2 $\times 10^{-12}$ | 7.4 $\times 10^{-10}$ |
| IXI T2w   | 96.4 $\pm$ 0.7 | 91.3 $\pm$ 2.7        | 91.0 $\pm$ 6.2                         | 94.9 $\pm$ 1.7        | 63.0 $\pm$ 12.4       | -                     | 6.7 $\pm$ 8.0         |
|           |                | 2.9 $\times 10^{-19}$ | 2.0 $\times 10^{-7}$                   | 1.4 $\times 10^{-9}$  | 7.8 $\times 10^{-24}$ | -                     | 1.1 $\times 10^{-52}$ |
| FSM T2w   | 97.7 $\pm$ 0.3 | 93.2 $\pm$ 1.7        | 95.7 $\pm$ 1.6                         | 94.9 $\pm$ 2.0        | 69.8 $\pm$ 15.0       | -                     | 3.4 $\pm$ 7.2         |
|           |                | 4.9 $\times 10^{-18}$ | 9.8 $\times 10^{-9}$                   | 4.1 $\times 10^{-10}$ | 7.3 $\times 10^{-13}$ | -                     | 1.1 $\times 10^{-40}$ |
| QIN T2w   | 95.2 $\pm$ 1.1 | 87.3 $\pm$ 5.3        | 89.6 $\pm$ 4.0                         | 71.4 $\pm$ 21.2       | 57.7 $\pm$ 19.1       | -                     | 61.5 $\pm$ 16.5       |
|           |                | 1.1 $\times 10^{-10}$ | 1.7 $\times 10^{-9}$                   | 6.3 $\times 10^{-7}$  | 9.1 $\times 10^{-14}$ | -                     | 2.9 $\times 10^{-15}$ |
| QIN FLAIR | 96.4 $\pm$ 0.5 | 93.7 $\pm$ 1.3        | 95.9 $\pm$ 0.9                         | 93.8 $\pm$ 1.1        | 90.3 $\pm$ 4.7        | 83.4 $\pm$ 6.2        | 87.9 $\pm$ 3.0        |
|           |                | 5.7 $\times 10^{-7}$  | 1.3 $\times 10^{-2}$                   | 4.2 $\times 10^{-9}$  | 8.7 $\times 10^{-5}$  | 2.6 $\times 10^{-7}$  | 1.5 $\times 10^{-9}$  |
| IXI PDw   | 96.4 $\pm$ 1.0 | 94.6 $\pm$ 1.3        | 95.5 $\pm$ 1.4                         | 95.1 $\pm$ 1.6        | 78.3 $\pm$ 15.8       | 81.3 $\pm$ 10.9       | 90.0 $\pm$ 3.8        |
|           |                | 5.3 $\times 10^{-17}$ | 1.9 $\times 10^{-7}$                   | 1.4 $\times 10^{-8}$  | 1.8 $\times 10^{-10}$ | 2.5 $\times 10^{-12}$ | 3.8 $\times 10^{-16}$ |
| FSM PDw   | 97.2 $\pm$ 0.5 | 95.8 $\pm$ 1.0        | 95.7 $\pm$ 6.5                         | 95.5 $\pm$ 1.4        | 88.8 $\pm$ 8.4        | 67.4 $\pm$ 9.8        | 79.1 $\pm$ 4.4        |
|           |                | 1.1 $\times 10^{-6}$  | <b>2.3 <math>\times 10^{-1}</math></b> | 4.7 $\times 10^{-8}$  | 5.1 $\times 10^{-6}$  | 4.9 $\times 10^{-17}$ | 3.4 $\times 10^{-21}$ |
| IXI MRA   | 97.7 $\pm$ 0.5 | 82.0 $\pm$ 4.3        | 62.3 $\pm$ 14.9                        | 95.1 $\pm$ 1.0        | 91.5 $\pm$ 5.3        | -                     | 18.1 $\pm$ 19.5       |
|           |                | 4.5 $\times 10^{-30}$ | 1.3 $\times 10^{-21}$                  | 3.0 $\times 10^{-24}$ | 6.7 $\times 10^{-10}$ | -                     | 4.7 $\times 10^{-32}$ |
| FSM qT1   | 97.7 $\pm$ 0.3 | 63.1 $\pm$ 19.6       | 47.2 $\pm$ 13.5                        | 44.6 $\pm$ 3.8        | 36.3 $\pm$ 15.6       | 49.5 $\pm$ 10.7       | 3.1 $\pm$ 7.1         |
|           |                | 1.1 $\times 10^{-10}$ | 8.8 $\times 10^{-20}$                  | 4.9 $\times 10^{-36}$ | 9.4 $\times 10^{-20}$ | 1.7 $\times 10^{-21}$ | 1.6 $\times 10^{-36}$ |
| ASL EPI   | 95.2 $\pm$ 1.4 | 88.3 $\pm$ 4.6        | 94.2 $\pm$ 1.8                         | 95.2 $\pm$ 1.4        | 67.4 $\pm$ 19.6       | 69.4 $\pm$ 11.9       | 92.8 $\pm$ 1.5        |

|            | SynthStrip | ROBEX                 | BET                   | 3DSS                  | BEaST                 | FSW                   | DMBE                  |
|------------|------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
|            |            | $3.6 \times 10^{-11}$ | $5.2 \times 10^{-4}$  | $5.8 \times 10^{-1}$  | $2.0 \times 10^{-11}$ | $8.9 \times 10^{-17}$ | $2.1 \times 10^{-15}$ |
| Infant T1w | 96.1 ± 1.4 | 87.4 ± 15.4           | 66.8 ± 20.3           | 71.6 ± 17.7           | 63.2 ± 38.4           | 61.3 ± 31.3           | 84.4 ± 7.3            |
|            |            | $3.5 \times 10^{-2}$  | $3.0 \times 10^{-5}$  | $4.4 \times 10^{-5}$  | $3.9 \times 10^{-3}$  | $7.9 \times 10^{-4}$  | $4.9 \times 10^{-6}$  |
| IXIDWI     | 95.5 ± 1.3 | 85.3 ± 5.7            | 93.4 ± 2.1            | 90.4 ± 2.3            | 79.5 ± 11.4           | 75.1 ± 8.8            | 80.5 ± 3.1            |
|            |            | $2.4 \times 10^{-10}$ | $1.2 \times 10^{-11}$ | $1.6 \times 10^{-17}$ | $6.8 \times 10^{-9}$  | $1.5 \times 10^{-13}$ | $1.4 \times 10^{-27}$ |
| CIM PET    | 95.4 ± 1.1 | 89.5 ± 4.8            | 78.3 ± 8.3            | 93.7 ± 6.1            | 4.6 ± 8.8             | 51.6 ± 10.1           | 42.3 ± 14.3           |
|            |            | $1.8 \times 10^{-5}$  | $2.9 \times 10^{-8}$  | $2.1 \times 10^{-1}$  | $1.4 \times 10^{-18}$ | $4.1 \times 10^{-13}$ | $1.6 \times 10^{-12}$ |
| CIM CT     | 94.3 ± 0.9 | 73.7 ± 2.0            | 45.4 ± 4.0            | 60.4 ± 2.9            | 2.3 ± 3.3             | 48.8 ± 5.1            | 58.2 ± 6.5            |
|            |            | $3.8 \times 10^{-23}$ | $4.8 \times 10^{-22}$ | $4.0 \times 10^{-22}$ | $3.5 \times 10^{-26}$ | $1.4 \times 10^{-16}$ | $2.1 \times 10^{-15}$ |



Average single-threaded CPU runtime ( $\pm$  SD) for T1w images from the FSM dataset. In addition to BET and FSW, SynthStrip is one of only three skull-stripping methods that consistently runs in under one minute.

**Table 5**

| CPU runtime (minutes) |                 |                 |                 |                 |                 |                  |
|-----------------------|-----------------|-----------------|-----------------|-----------------|-----------------|------------------|
| SynthStrip            | ROBEX           | BET             | 3DSS            | BEaST           | FSW             | DMBE             |
| 0.48 $\pm$ 0.01       | 2.45 $\pm$ 0.11 | 0.20 $\pm$ 0.14 | 2.27 $\pm$ 1.12 | 4.14 $\pm$ 0.24 | 0.19 $\pm$ 0.01 | 48.89 $\pm$ 4.72 |