

1 **Title: 3D-MASNet: 3D Mixed-scale Asymmetric Convolutional Segmentation Network for**
2 **6-month-old Infant Brain MR Images**

3
4 **Author affiliations:** Zilong Zeng^{1,2,3}, Tengda Zhao^{1,2,3*}, Lianglong Sun^{1,2,3}, Yihe Zhang^{1,2,3},
5 Mingrui Xia^{1,2,3}, Xuhong Liao⁴, Jiaying Zhang^{1,2,3}, Dinggang Shen^{5,6}, Li Wang⁷, Yong He^{1,2,3,8*}

6
7 ¹ State Key Laboratory of Cognitive Neuroscience and Learning, Beijing Normal University,
8 Beijing 100875, China

9 ² Beijing Key Laboratory of Brain Imaging and Connectomics, Beijing Normal University,
10 Beijing 100875, China

11 ³ IDG/McGovern Institute for Brain Research, Beijing Normal University, Beijing 100875,
12 China

13 ⁴ School of Systems Science, Beijing Normal University, Beijing 100875, China

14 ⁵ School of Biomedical Engineering, ShanghaiTech University, Shanghai 201210, China

15 ⁶ Department of Research and Development, Shanghai United Imaging Intelligence Co., Ltd.,
16 Shanghai 200030, China

17 ⁷ Department of Radiology and BRIC, University of North Carolina at Chapel Hill, NC 27599,
18 USA

19 ⁸ Chinese Institute for Brain Research, Beijing 102206, China

20
21 * Correspondence: yong.he@bnu.edu.cn (Y.H.) or tengdazhao@bnu.edu.cn (T.Z)

22
23 **Manuscript information: 20 text pages, 6 figures, 6 tables.**

1 **Abstract**

2 Precise segmentation of infant brain MR images into gray matter (GM), white matter (WM), and
3 cerebrospinal fluid (CSF) are essential for studying neuroanatomical hallmarks of early brain
4 development. However, for 6-month-old infants, the extremely low-intensity contrast caused by
5 inherent myelination hinders accurate tissue segmentation. Existing convolutional neural
6 networks (CNNs) based segmentation models for this task generally employ single-scale
7 symmetric convolutions, which are inefficient for encoding the isointense tissue boundaries in
8 baby brain images. Here, we propose a 3D mixed-scale asymmetric convolutional segmentation
9 network (3D-MASNet) framework for brain MR images of 6-month-old infants. We replaced the
10 traditional convolutional layer of an existing to-be-trained network with a 3D mixed-scale
11 convolution block consisting of asymmetric kernels (MixACB) during the training phase and
12 then equivalently converted it into the original network. Five canonical CNN segmentation
13 models were evaluated using both T1- and T2-weighted images of 23 6-month-old infants from
14 iSeg-2019 datasets, which contained manual labels as ground truth. MixACB significantly
15 enhanced the average accuracy of all five models and obtained the most considerable
16 improvement in the fully convolutional network model (CC-3D-FCN) and the highest
17 performance in the Dense U-Net model. This approach further obtained Dice coefficient
18 accuracies of 0.931, 0.912, and 0.961 in GM, WM, and CSF, respectively, ranking first among
19 30 teams on the validation dataset of the iSeg-2019 Grand Challenge. Thus, the proposed 3D-
20 MASNet can improve the accuracy of existing CNNs-based segmentation models as a plug-and-
21 play solution that offers a promising technique for future infant brain MRI studies.

22

23 **Keywords:** infant brain segmentation, MRI, convolutional neural networks, mixed-scale
24 convolution

1 **1. Introduction**

2 The accurate tissue segmentation of infant brain magnetic resonance (MR) images into gray
3 matter (GM), white matter (WM), and cerebrospinal fluid (CSF) are essential for researchers to
4 chart the normal and abnormal early brain development of cortical regions, white matter
5 connections, and wiring topologies (Cao et al., 2017; Hazlett et al., 2017; Wang et al., 2019a;
6 Wen et al., 2019; Xu et al., 2019; Zhao et al., 2019). Notably, the tissue segmentation of 6-
7 month-old infants is the biggest challenge in baby brain segmentation tasks due to the isointense
8 phase in which the intensity distributions of GM and WM voxels become dramatically
9 overlapped in the cortical regions (Fig. 1). The effective manual annotation, which is guided by
10 longitudinal tracking of brain images with high tissue contrast in the latter children period (Wang
11 et al., 2019b), is limited by the extremely high labor costs, the requirement of specialized expert
12 knowledge (almost one week per image for an experienced neuroradiologist) and high inter- and
13 intra-rater variations (Makropoulos et al., 2018). Developing fast, automatic, and accurate brain
14 segmentation approaches is a crucial and ongoing goal for MR images of infants at 6 months of
15 age (Sun et al., 2021; Wang et al., 2019b).

17 ***1.1. Convolutional neural networks (CNNs) based methods become the mainstream***

18 In the past years, many efforts have been made for the segmentation task of 6-month-old infant
19 brain MR images. Generally, emerging CNNs-based segmentation methods that exhibiting faster
20 computational speed and higher accuracy than conventional atlas-based (Wang et al., 2014;
21 Wang et al., 2012) or machine learning based methods (Sanroma et al., 2018; Wang et al., 2015;
22 Wang et al., 2018a; Wang et al., 2014) become the mainstream. A typical example is that seven
23 of the eight top teams in iSeg-2017 challenge has utilized CNNs to segment infant brain tissues.

24 Current CNNs-based approaches for infant brain segmentation are usually variants of
25 canonical FCN (Long et al., 2015) and U-Net (Ronneberger et al., 2015) architecture. By
26 adjusting or adding specific connectional pathways within or across neural layers on classical
27 CNNs models, these approaches enhance the extraction and fusion of the semantic information in
28 multimodal features to counteract the noisy and isointense tissues boundaries in 6-month-old
29 infant brain images (Bui et al., 2019; Dolz et al., 2020; Dolz et al., 2019; Nie et al., 2019; Nie et
30 al., 2016; Wang et al., 2018b; Wang et al., 2020; Zeng and Zheng, 2018; Zhang et al., 2015).
31 Specifically, Bui et al. improved densely connected network (DenseNet) (Huang et al., 2017) by

1 concatenating fine and coarse feature maps from multiple densely connected blocks and won the
2 iSeg-2017 competition (Bui et al., 2019). Dolz et al. proposed a semi-dense network by directly
3 connecting all of the convolutional layers to the end of the network (Dolz et al., 2020) and
4 further extended it into a HyperDenseNet by adding dense connections between multimodal
5 network paths (Dolz et al., 2019). Similarly, Zeng et al. modified the classical U-Net network by
6 constructing multi-encoder paths for each modality to effectively extract targeted high-level
7 information (Zeng and Zheng, 2018). Wang et al. designed a global aggregation block in the U-
8 Net model to consider global information in the decoder path of feature maps (Wang et al.,
9 2020). Interestingly, inspired by the superiority of DenseNet and U-Net, the densely connected
10 U-Net (DU-Net) model with a combination of these two types of networks was proposed for
11 both tissue segmentation and autism diagnosis (Wang et al., 2018b).

12

13 ***1.2. Improvements from fine-grained convolution kernel designs are underestimated***

14 Although great efforts have been made, the above CNN-based segmentation models have several
15 limitations. First, the image appearance of 6-month-old infant brain MR images is quite noisy (Li
16 et al., 2019; Mostapha and Styner, 2019) which makes the effective feature extraction difficult
17 for the traditional convolution kernel design in previous works. Adopting enhanced convolution
18 kernel designs (Ding et al., 2019; Li et al., 2020; Zhang et al., 2022) that emphasizes key features
19 in the skeleton center of kernels may facilitate feature extractions throughout the network.
20 Second, the voxel-wise fuzzy tissue boundaries in infant brain images are constrained by the
21 anatomical morphology of gyrus at large spatial scales (Wang et al., 2018a). Although previous
22 infant segmentation approaches try to fuse multi-scale features by skip-connections in variants of
23 FCN and U-Net, they overlook capturing rich multi-scale features in kernel space, which
24 contains more stable and homogeneous semantic information than features between layers (Fan
25 et al., 2019). Third, all these studies focused on modifications of network layouts which need
26 seasoned expertise experience, time-consuming hyperparameter tuning, and may also bring
27 excessive graphics processing unit (GPU) burdens (Dolz et al., 2019; Wang et al., 2020) and
28 architecture incompatibility. Recent CNN studies move eyes on building architecture-
29 independent designs such as SE blocks (Hu et al., 2018), or automatically configuring models
30 such as nnU-net (Isensee et al., 2021), which requires neither rare expert knowledge nor
31 expensive manual interventions.

1 **1.3. Our contribution**

2 Our goal is to obtain a CNN-based building block for 6-month-old infant brain image
3 segmentation which is 1) with fine-grained kernel designs to enhance the representation and
4 abundance of features; 2) transplantable in up-to-date segmentation models in a plug-and-play
5 way; 3) without much additional hyperparameter tuning or computational burden. To this end,
6 we construct a 3D mixed-scale asymmetric segmentation network (3D-MASNet) framework by
7 embedding a well-designed 3D mixed-scale asymmetric convolution block (MixACB) into
8 existing segmentation CNNs for 6-month-old infant brain MR images (Fig. 2). The MixACB
9 design is comprised by 1) four parallel 3D convolutional layers including a symmetric kernel
10 ($d \times d \times d$) and three asymmetric 2D kernels ($1 \times d \times d$, $d \times 1 \times d$, $d \times d \times 1$) (Fig. 3A),
11 respectively; 2) multiple groups on input feature maps with different kernel sizes (Fig. 3B)
12 independently ; 3) parameter fusion for each MixACB after the training process to lower
13 inference-time computations compare to the original network. We first evaluated the
14 effectiveness of the MixACB on five canonical CNN networks using the iSeg-2019 training
15 dataset. We next compared the performance of our method with that of top-4 approaches
16 proposed in the MICCAI iSeg-2019 Grand Challenge on the iSeg-2019 validation dataset. The
17 experimental results revealed that the MixACB significantly improved the segmentation
18 accuracy of various CNNs, among which DU-Net (Wang et al., 2018b) with MixACB achieved
19 the best-enhanced average performance and obtained the highest Dice coefficients of 0.931 in
20 GM, 0.912 in WM, and 0.961 in CSF, ranking first in the iSeg-2019 Grand Challenge. All codes
21 are publicly available at <https://github.com/RicardoZiTseng/3D-MASNet>.

22

23 **2. Methods and Implementations**

24 **2.1.1. Mathematical formulation of basic 3D convolution**

25 Consider a feature map $I \in \mathbb{R}^{U \times V \times S \times C}$ with a spatial resolution of $U \times V \times S$ as input and a feature
26 map $O \in \mathbb{R}^{R \times T \times Q \times K}$ with a spatial resolution of $R \times T \times Q$ as output of a convolutional layer with
27 a kernel size of $H \times W \times D$ and K filters. Then, each filter's kernel is denoted as
28 $F \in \mathbb{R}^{H \times W \times D \times C}$, and the operation of the convolutional layer with a batch normalization (BN)
29 layer can be formulated as follows:

$$\begin{aligned}
 O_{\dots,j} &= \left(\sum_{k=1}^C I_{\dots,k} * F_{\dots,k}^{(j)} - \mu_j \right) \cdot \frac{\gamma_j}{\sigma_j} + \beta_j \\
 &= \left(\sum_{k=1}^C I_{\dots,k} * \frac{\gamma_j}{\sigma_j} F_{\dots,k}^{(j)} \right) - \frac{\mu_j \gamma_j}{\sigma_j} + \beta_j
 \end{aligned} \tag{1}$$

1 where $*$ is the 3D convolution operator, $I_{\dots,k}$ is the k^{th} channel of the input feature map I ,
 2 $F_{\dots,k}^{(j)}$ is the k^{th} channel of the j^{th} filter's kernel, μ_j and σ_j are the channel-wise mean
 3 value and standard deviation value, respectively, γ_j and β_j are the scaling factor and bias
 4 term to restore the representation ability of the network, respectively.
 5
 6

7 **2.1.2. Design of 3D asymmetric convolutions (3D-AC) during training and inference phases**

8 3D-AC was designed behaving differently during training and inference phases (Fig. 3A).
 9 Concretely, for each kernel of each layer in the network during the training phase, a 3D-AC
 10 contains 4 parallel convolutional branches, namely one standard 3D convolution layer and three
 11 orthogonal 2D asymmetric convolutional layers ($1 \times d \times d$, $d \times 1 \times d$, $d \times d \times 1$) at kernel center
 12 for the enhancement of features along axial, sagittal and coronal directions, respectively. The
 13 input feature maps are fed into these 4 branches, and the outputs of these branches are summed
 14 to fuse the knowledge learned by these 4 independent branches. During the inference phase, 3D-
 15 AC contains one standard convolutional layer with equivalently fused kernel of the training-time
 16 3D-AC (described in section 2.1.4), thus the input feature maps only need feed into this single
 17 branch which bringing low inference computations.
 18

19 **2.1.3. Constructing MixACB by multiple 3D-ACs with varying kernel scales**

20 To process the input feature map at different scales of detail, we propose the MixACB by mixing
 21 multiple 3D-ACs with different kernel sizes, as illustrated in Fig. 3B. Notably, since we used the
 22 3D-AC to strength the core skeleton part of the convolutional kernel, thus the kernel size of 3D-
 23 AC must be odd, such as 3, 5, 7. Since directly adopting multiple 3D-ACs on all feature maps
 24 then concatenating outputs will dramatically increase the models' parameters and computations,
 25 we leverage the grouped convolution approach by splitting original input feature maps into
 26 groups and apply 3D-AC independently in each input feature map's group. Assume that we split
 27 the input feature maps into g groups of tensors such that their total number of channels is equal

1 to the original feature maps' channels: $C_1 + C_2 + \dots + C_g = C$ with $C_1 > C_2 > \dots > C_g$; similarly, the
 2 output feature maps also have g groups: $K_1 + K_2 + \dots + K_g = K$ with $K_1 > K_2 > \dots > K_g$. We
 3 denote $I^{<i>} \in \mathbb{R}^{U \times V \times S \times C_i}$ as the i^{th} group of input, $\widehat{O}^{<i>} \in \mathbb{R}^{R \times T \times Q \times K_i}$ as the MixACB's i^{th}
 4 group output, and $F^{* <i>} \in \mathbb{R}^{H_i \times W_i \times D_i \times C_i}$ as the equivalent kernel of the i^{th} group of the 3D-AC
 5 whose equivalent kernel size is $H_i \times W_i \times D_i$. Thus, we have following equations:

$$6 \quad \widehat{O}_{\dots, j}^{<i>} = \left(\sum_{q=1}^{C_i} I_{\dots, q}^{<i>} * F_{s^{<i>}(j)}^{* <i>} \right) + b_j^{<i>} \quad (2)$$

7 $s.t. \quad 1 \leq i \leq g, \quad 1 \leq j \leq K_i$

7 The final output of MixACB is the concatenation of all groups' outputs:

$$8 \quad \widehat{O} = \text{concat} \left(\widehat{O}^{<1>}, \widehat{O}^{<2>}, \dots, \widehat{O}^{<g>} \right) \quad (3)$$

9 In this paper, we only split the input and output feature maps into 2 groups, and define the
 10 mix ratio as the ratio between C_1 and C_2 . For simplicity, the ratio between K_1 and K_2 is set
 11 to be equal to the mix ratio. We set a kernel size of 3 for the 1st group of 3D-AC and a kernel size
 12 of 5 for the 2nd group of 3D-AC, with the mix ratio set to 3:1. In this situation, the number of
 13 FLOPs (floating point operations) required for inference-time MixACB is

$$14 \quad 3^3 \times \frac{3}{4} C \times \frac{3}{4} K + 5^3 \times \frac{1}{4} C \times \frac{1}{4} K = 23CK, \text{ which is smaller than that } (3^3 \times C \times K = 27CK) \text{ for standard}$$

15 convolutions with kernel size of 3.

16

17 **2.1.4. Equivalently fusing kernel of each 3D-AC inside MixACB**

18 Once the training process of 3D-MASNet is completed, we equivalently fused the kernels of
 19 each 3D-AC inside the MixACB to retain the same output results as the original network. Due to
 20 the additivity of convolutional kernels, the kernels of 3D-AC's four branches can be fused to
 21 obtain an equivalent kernel in a 3D convolutional layer to produce the same output, which can be
 22 formulated as the following equation:

$$23 \quad I * F + I * \widetilde{F} + I * \widehat{F} + I * \overline{F} = I * \left(F \oplus \widetilde{F} \oplus \widehat{F} \oplus \overline{F} \right) = I * F' \quad (4)$$

1 where I is an input feature map, F , \tilde{F} , \hat{F} and \bar{F} are the 4 branches' kernels of 3D-AC.
 2 \oplus is an elementwise operator that performs parameter addition on the corresponding positions,
 3 and F' is the equivalent fused kernel of the 4 branches' kernels.

4 Here, we took a kernel size of 3 as an example. We first fused the BN parameters into the
 5 convolutional kernel term and bias term following Eq. (1). Then, we further fused the four
 6 parallel kernels by adding the asymmetric kernels onto the skeletons of the cubic kernel.

7 Formally, we denote $F^{(j)}$ as the j^{th} filter at the $1 \times 3 \times 3$, $3 \times 1 \times 3$ and $3 \times 3 \times 1$ layer,
 8 respectively. Hence, we obtain the following formulas:

$$9 \quad F^{(j)} = \frac{\gamma_j}{\sigma_j} F^{(j)} \oplus \frac{\tilde{\gamma}_j}{\tilde{\sigma}_j} \tilde{F}^{(j)} \oplus \frac{\hat{\gamma}_j}{\hat{\sigma}_j} \hat{F}^{(j)} \oplus \frac{\bar{\gamma}_j}{\bar{\sigma}_j} \bar{F}^{(j)} \quad (5)$$

$$10 \quad b'_j = -\frac{\mu_j \gamma_j}{\sigma_j} - \frac{\tilde{\mu}_j \tilde{\gamma}_j}{\tilde{\sigma}_j} - \frac{\hat{\mu}_j \hat{\gamma}_j}{\hat{\sigma}_j} - \frac{\bar{\mu}_j \bar{\gamma}_j}{\bar{\sigma}_j} + \beta_j + \tilde{\beta}_j + \hat{\beta}_j + \bar{\beta}_j \quad (6)$$

11 Then, we can write any output of j^{th} filter as:

$$12 \quad O_{\dots,j} + \tilde{O}_{\dots,j} + \hat{O}_{\dots,j} + \bar{O}_{\dots,j} = \sum_{k=1}^c I_{\dots,k} * F^{(j)}_{\dots,k} + b'_j \quad (7)$$

13 where $O_{\dots,j}$, $\tilde{O}_{\dots,j}$, $\hat{O}_{\dots,j}$ and $\bar{O}_{\dots,j}$ are the outputs of the original $3 \times 3 \times 3$, $1 \times 3 \times 3$,
 14 $3 \times 1 \times 3$ and $3 \times 3 \times 1$ branch, respectively.

15

16 **2.2. Candidate CNNs for the evaluation of the MixACB on 6-month-old infant brain image** 17 **segmentation**

18 We choose five representative networks to evaluate the effectiveness of the 3D-MixACB in
 19 improving the segmentation performance, including BuiNet (Bui et al., 2019), 3D-UNet (Çiçek
 20 et al., 2016), convolution and concatenate 3D fully convolutional network (CC-3D-FCN) (Nie et
 21 al., 2019), non-local U-Net (NLU-Net) (Wang et al., 2020) and DU-Net (Wang et al., 2018b).

22 Notably, these five networks are either variants of the U-type architecture (3D U-Net, NLU-Net,
 23 and DU-Net) or the FCN-type architecture (BuiNet and CC-3D-FCN) and encompass major
 24 CNN frameworks in infant brain segmentation. After replacing their original convolution layers
 25 with the 3D-MixACB design, we followed the training configurations set in the candidate CNN's
 26 release codes (Table 1) and adopted the Adam optimizer to update these models' parameters.
 27 Except for the CC-3D-FCN, which used the Xavier algorithm (Glorot and Bengio, 2010) to

1 initialize network weights, all other networks adopted the He initialization method (He et al.,
2 2015). The configuration parameters are as follows:

3 (1) BuiNet adopted four dense blocks consisting of four $3 \times 3 \times 3$ convolutional layers for
4 feature extraction. Transition blocks were applied between every two dense blocks to reduce the
5 feature map resolutions. 3D up-sampling operations were used after each dense block for feature
6 map recovery, and these upsampled features were concatenated together. (2) 3D-UNet has 4
7 levels of resolution, and each level adopts one $3 \times 3 \times 3$ convolution, which is followed by BN
8 and a rectified linear unit (ReLU). The $2 \times 2 \times 2$ max pooling and the $2 \times 2 \times 2$ transposed
9 convolution, each with a stride of 2, are employed for resolution reduction and recovery. Feature
10 maps of the same level of both paths were summed. (3) CC-3D-FCN used 6 groups of $3 \times 3 \times 3$
11 convolutional layers for feature extraction, in which the $2 \times 2 \times 2$ max pooling with a stride of 2
12 was adopted between two groups of layers. The $1 \times 1 \times 1$ convolution with a stride of 1 was
13 added between two groups with the same resolution for feature fusion. (4) DU-Net used 7 dense
14 blocks to construct the encoder-decoder structure with 4 levels of resolution and leveraged
15 transition down blocks and transition up blocks for down-sampling and up-sampling,
16 respectively. Unlike the implementations in (Wang et al., 2018b), the bottleneck layer is
17 introduced into the dense block to constrain the rapidly increasing number of feature maps, and
18 the transition down block consisted of two $3 \times 3 \times 3$ convolutions, each followed by BN and
19 ReLU. In addition, we used the $1 \times 1 \times 1$ convolution followed by a softmax activation function
20 in the last layer. (5) NLU-Net leveraged five different kinds of residual blocks to form the U-
21 type architecture with 3 levels of resolution. BN with the ReLU6 activation function was adopted
22 before each $3 \times 3 \times 3$ convolution. The global aggregation block replaced the two convolutional
23 layers of the input residual block to form the bottom residual block for the integration of global
24 information.

25 We fed the same multimodal images into these five networks and employed the same
26 inference strategy. We extracted overlapping patches of the same size as that used during the
27 training phase. The overlapping step size had to be smaller than or equal to the patch length size
28 to form the whole volume. Following the common practice in (Bui et al., 2019; Nie et al., 2019;
29 Wang et al., 2018b; Wang et al., 2020), we set the step size to 8. Since the effect of the
30 overlapping step size in the proposed framework remains unknown, we further evaluated it in
31 section 3.3. Voxels inside the overlapping regions were averaged.

1

2 **3. Experiments and Results**

3 **3.1. iSeg-2019 dataset and image preprocessing**

4 Twenty-three isointense phase infant brain MRIs, including T1w and T2w images, were offered
5 by the iSeg-2019 (<http://iseg2019.web.unc.edu/>) organizers from the pilot study of the Baby
6 Connectome Project (BCP) (Howell et al., 2019). All the infants were term-born (40 ± 1 weeks of
7 gestational age) with an average scan age of 6.0 ± 0.5 months. All experimental procedures were
8 approved by the University of North Carolina at Chapel Hill and the University of Minnesota
9 Institutional Review Boards. Detailed imaging parameters and preprocessing steps that were
10 implemented are listed in (Sun et al., 2021). Before cropping the MR images into patches, we
11 normalized the T1w and T2w images by subtracting the mean value and dividing by the standard
12 deviation value.

13 The iSeg-2019 organizers offered the ground truth labels, which were obtained by a
14 combination of initial automatic segmentation using the infant brain extraction and analysis
15 toolbox (iBEAT) (Dai et al., 2013) on follow-up 24-month scans of the same baby and manual
16 editing using ITK-SNAP (Yushkevich et al., 2006) under the guidance of an experienced
17 neuroradiologist. The MR images of 10 infants with manual labels were provided for model
18 training and validation. The images of 13 infants without labels were provided for model testing.
19 The testing results were submitted to the iSeg-2019 organizers for quantitative measurements.

20

21 **3.2. Evaluation metrics**

22 We employed the Dice coefficient (DICE), modified Hausdorff distance (MHD) and average
23 surface distance (ASD) to evaluate the model performance on segmenting 6-month-old infant
24 brain MR images.

25 **3.2.1. Dice coefficient**

26 Let A and B be the manual labels and predictive labels, respectively. The DICE can be
27 defined as:

$$28 \quad DICE(A, B) = \frac{2|A \cap B|}{|A| + |B|} \quad (8)$$

29 where $|\cdot|$ denotes the number of elements of a point set. A higher DICE indicates a larger
30 overlap between the manual and predictive segmentation areas.

1

2 **3.2.2. Modified Hausdorff distance**

3 Let C and D be the sets of voxels within the manual and predictive segmentation boundary,
4 respectively. MHD can be defined as:

$$5 \quad MHD(C, D) = \max \{h(C, D), h(D, C)\} \quad (9)$$

6 where $h(C, D) = \frac{1}{N_c} \sum_{c \in C} d(c, D)$, and $d(c, D) = \min_{d \in D} \|c - d\|$ with $\|\cdot\|$ representing the

7 Euclidean distance. We follow the calculation described in (Wang et al., 2020) by computing the
8 average MHD based on the three different vectorization directions to obtain a direction-
9 independent evaluation metric. A smaller MHD coefficient indicates greater similarity between
10 manual and predictive segmentation contours.

11

12 **3.2.3. Average surface distance**

13 The ASD is defined as:

$$14 \quad ASD(C, D) = \frac{1}{2} \cdot \left(\frac{\sum_{v_i \in S_C} \min_{v_j \in S_D} \|v_i - v_j\|}{\sum_{v_i \in S_C} 1} + \frac{\sum_{v_j \in S_D} \min_{v_i \in S_C} \|v_j - v_i\|}{\sum_{v_j \in S_D} 1} \right) \quad (10)$$

15 where S_C and S_D represent the surface meshes of C and D , respectively. A smaller ASD
16 coefficient indicates greater similarity between cortical surfaces reconstructed from manual and
17 predictive segmentation maps.

18

19 **3.3. Exploring the effectiveness of the MixACB**

20 We performed several experiments to evaluate the effectiveness of the MixACB, including 1)
21 ablation tests on five representative segmentation networks (section 2.2); 2) comparisons with
22 state-of-the-art approaches in iSeg-2019; 3) component analysis of MixACB; and 4) validation
23 of the impact of the overlapping step size.

24

25 **3.3.1. Performance improvement on five representative CNN architectures**

26 For a given network architecture without the MixACB design, we regarded it as the baseline
27 model and further transformed it into a 3D-MASNet design. All pairs of the baseline models and
28 their corresponding 3D-MASNet followed the training strategies described in Table 1. To

1 balance the training and testing sample sizes, we adopt a 2-fold cross-validation (one fold with
2 five random selected participants for training and the left for testing) for model evaluation on the
3 iSeg-2019 training dataset. Table 2 and Table 3 and Fig. 4 show that the performance of all the
4 models was significantly improved across almost all tissue types in terms of the DICE and MHD,
5 which demonstrates the effectiveness of the MixACB on a wide range of CNN layouts.
6 Specifically, DU-Net with the MixACB achieved the highest average DICE of 0.928 and the
7 lowest average MHD value of 0.436; CC-3D-FCN with the MixACB gained the most
8 considerable DICE improvement and reached a higher average DICE than that attained by
9 BuiNet, which was a champion solution in the MICCAI iSeg-2017 grand challenge, indicating
10 that a simple network could reach excellent performance by advanced convolution designs. Fig.
11 5 further provides a visual segmentation comparison between networks with and without the
12 MixACB. The MixACB could effectively correct misclassified voxels which are indicated by red
13 squares.

14

15 **3.3.2. Comparison with state-of-the-art methods on iSeg-2019**

16 Since DU-Net, which was combined with MixACB, has achieved the highest accuracy among all
17 candidate models, we compared it with methods developed by the 29 remaining teams that
18 participated in the iSeg-2019 challenge. We employed a majority-voting strategy on 10 trained
19 networks' outputs to improve the model generalization.

20 Table 4 reports the segmentation results achieved by our proposed method and those of other
21 teams' methods that ranked in the top 4 on the validation dataset of the iSeg-2019. The mean
22 DICE, MHD value and ASD value are presented for CSF, GM, and WM, representatively.
23 Compared with other teams, our method yielded the highest DICE and lowest ASD value for the
24 three brain tissues in the validation test of iSeg-2019, with comparable MHD values. The
25 superior average value of the 3 types of brain tissues also indicates that our method has the best
26 overall performance.

27

28 **3.3.3. Component analysis of MixACB**

29 We analyzed the effect of the mix ratio on model segmentation performance. Table 5 shows that
30 segmentation accuracy reaches the highest value when the mix ratio is set to 3:1. Then, we
31 further performed an ablation test to verify the effectiveness of each part of the proposed

1 MixACB, as shown in Table 6. The segmentation accuracy was improved with large variations
2 when using different 3D-ACs alone. Moreover, when these 3D-ACs were mixed in scales for a
3 MixACB design, the model was able to achieve the best performance in both DICE and MHD
4 metrics.

5 ***3.3.4. Impact of overlapping step sizes***

6 We further performed experiments to evaluate the effectiveness of the MixACB on overlapping
7 step sizes, which controls the trade-off between accuracy and inference time. Based on 2-fold
8 cross-validation, which has been done previously, we tested the overlapping impact when the
9 step size is set to 4, 8, 16, and 32 on the DU-Net in the proposed 3D-MASNet framework. Fig.
10 6A and Fig. 6B present the changes in the segmentation performance in terms of DICE and
11 MHD, respectively, for different overlapping step sizes. Fig. 6C presents the changes in the
12 average number of inference patches for different overlapping step sizes. We found that a step
13 size of 8 is a reasonable choice for achieving fast and accurate results.

14

15 **4. Discussion**

16 Instead of designing a new network architecture to segment the brain images of 6-month-old
17 infants, we proposed a 3D-MASNet framework by replacing the standard convolutional layer
18 with MixACB on an existing mature network and reduced model parameters and computations
19 by equivalently performing fusion during the inference phase. The experimental results revealed
20 that the MixACB significantly improved the performance of several CNNs by a considerable
21 margin, in which DU-Net with MixACB showed the best average segmentation accuracy. The
22 proposed framework obtained the highest average DICE of 0.935 and lowest ASD of 0.244,
23 which ranked first among all 30 teams on the validation dataset of the iSeg-2019 Grand
24 Challenge. In addition, the CC-3D-FCN model showed the largest improvement, which indicates
25 that a simple model could achieve relatively better performance by implementing our
26 convolution design.

27

28 ***4.1. Effectiveness of the MixACB on improving segmentation accuracy***

29 The wide improvement in the segmentation accuracy of different models by the MixACB is
30 derived from several aspects. First, the mixed-scale design of MixACB enables the network to
31 collect multiscale details of local features with different receptive fields, facilitating the

1 integration of coarse-to-fine information inside the input patches at low-to-high semantic levels.
2 Second, the isointense intensity distribution and heterogeneous tissue contrasts hamper effective
3 feature extraction in baby brain images. The mix rate at 3:1 of feature maps for multi-scale
4 kernel size enriches small receptive fields with enough detail features while enabling large
5 receptive fields for capturing coarse global features. We also employed the 3D-AC inside the
6 MixACB by adding multiple orthogonal 3D asymmetric convolutional layers to emphasize
7 informative feature patterns in the central place (Fig. 3). Meanwhile, the asymmetric design is
8 also shown robustness to image rotational distortion (Ding et al., 2019), which may help the
9 network cope with the residual head motion of infants, even though these images have been
10 linearly aligned to standard space. Third, the significant improvable performance of MixACB on
11 various segmentation networks (Table 3) indicates that inter-layer architecture design may not be
12 sufficient for multi-scale information fusion. Notably, besides providing better performance than
13 the previous networks, 3D-MASNet is also more efficient than the baseline models, requiring
14 fewer model parameters once its parameters were fused in the inference phase. For example, the
15 baseline DU-Net's number of parameters is 2,492,795, while the corresponding 3D-MASNet's
16 number of parameters is reduced to 2,341,141 during the inference phase.

17

18 ***4.2. Well-designed convolution operations***

19 In recent years, researchers have begun to shift their interests from macro network layout to
20 micro neuron units by studying specific convolution operators rather than touching the overall
21 network. Previous works have proposed several advanced convolution operators by combining
22 well-designed filters, such as pyramidal convolution (PyConv), dynamic group convolution
23 (DGC), and asymmetric convolution block (ACB). PyConv employs multiple kernels in a
24 pyramidal way to capture different levels of image details (Duta et al., 2020); DGC equips a
25 feature selector for each group convolution conditioned on the input images to adaptively select
26 input features (Su et al., 2020); ACB introduces asymmetry into 2D convolution to power up the
27 representational power of the skeleton part of the kernel (Ding et al., 2019). These operators
28 implanted into existing mature networks have achieved better performance on image
29 classification or semantic segmentation tasks than in original networks. Due to the “easy-to-plug-
30 in” property, this type of design could be conveniently adopted in various advanced CNNs and
31 avoided high cost of network re-designing. However, these studies mainly concentrate on natural

1 image tasks, few were applied to infant brain segmentation tasks. Here, we design a novel
2 convolution block by combining three basic characteristics including 3D spatial convolution,
3 group convolution containing mixed-scale of kernel sizes, and asymmetry convolution. Due to
4 blurred image appearance, large individual variation of brain morphology, and limited labeled
5 sample sizes, we emphasize that effective and robust feature extraction, especially in a plug-and-
6 play form, is essential for the infant brain segmentation task. Nevertheless, exhausting the
7 combination of various convolution designs is beyond the scope of the article.

8

9 ***4.3. Limitations and future directions***

10 The current study has several limitations. First, the patching approach may cause spatial
11 consistency loss near boundaries. Although we adopted a small overlapping step size to relieve
12 this issue, it is necessary to consider further integrating guidance from global information.
13 Second, the small sample sizes of infant-specific datasets limit the generalizability of our method
14 for babies across MRI scanners and acquisition protocols. Further validation on large samples is
15 needed. Third, image indexes, such as the fractional anisotropy derived from diffusion MRI,
16 contain rich white matter information (Liu et al., 2007), which could be beneficial for
17 insufficient tissue contrast (Nie et al., 2019; Zhang et al., 2015). Importantly, determining how to
18 leverage mixed-scale asymmetric convolution to enhance specific model features needs to be
19 further explored. Fourth, we only explored the effectiveness of MixACB when input feature
20 maps are split into 2 groups. Further combination configurations of convolutional kernel sizes
21 and mix ratios are warranted.

22

23 **5. Conclusion**

24 In this paper, we proposed a 3D-MASNet framework for brain MR image segmentation of 6-
25 month-old infants, which ranked first in the iSeg-2019 Grand Challenge. We demonstrated that
26 the designed MixACB could easily migrate to various network architectures and enable
27 performance improvement without extra inference-time computations. This work shows great
28 adaptation potential for further improvement in future studies on brain segmentation.

29

1 **Acknowledgments**

2 The study was supported by the National Natural Science Foundation of China (Nos. 31830034,
3 82021004 and 81801783), Changjiang Scholar Professorship Award (T2015027), and the China
4 Postdoctoral Science Foundation (2020TQ0050 and 2022M710433).

5

6 **Conflicting Interests**

7 The authors have declared that no conflicting interests exist.

8

9 **Data and Code Availability Statement**

10 The 6-months-old infant brain MRI data was publicly offered by the iSeg-2019
11 (<http://iseg2019.web.unc.edu/>) organizers.

12 Codes developed for the proposed segmentation algorithm are released at
13 <https://github.com/RicardoZiT seng/3D-MASNet>.

14

1 **References**

- 2 Bui, T.D., Shin, J., Moon, T., 2019. Skip-connected 3D DenseNet for volumetric infant brain
3 MRI segmentation. *Biomedical Signal Processing and Control* 54, 101613.
- 4 Cao, M., Huang, H., He, Y., 2017. Developmental connectomics from infancy through early
5 childhood. *Trends in neurosciences* 40, 494-506.
- 6 Çiçek, Ö., Abdulkadir, A., Lienkamp, S.S., Brox, T., Ronneberger, O., 2016. 3D U-Net: learning
7 dense volumetric segmentation from sparse annotation. *International conference on medical
8 image computing and computer-assisted intervention*. Springer, pp. 424-432.
- 9 Dai, Y., Shi, F., Wang, L., Wu, G., Shen, D., 2013. iBEAT: A Toolbox for Infant Brain
10 Magnetic Resonance Image Processing. *Neuroinformatics* 11, 211-225.
- 11 Ding, X., Guo, Y., Ding, G., Han, J., 2019. Acnet: Strengthening the kernel skeletons for
12 powerful cnn via asymmetric convolution blocks. *Proceedings of the IEEE International
13 Conference on Computer Vision*, pp. 1911-1920.
- 14 Dolz, J., Desrosiers, C., Wang, L., Yuan, J., Shen, D., Ayed, I.B., 2020. Deep CNN ensembles
15 and suggestive annotations for infant brain MRI segmentation. *Computerized Medical Imaging
16 and Graphics* 79, 101660.
- 17 Dolz, J., Gopinath, K., Yuan, J., Lombaert, H., Desrosiers, C., Ayed, I.B., 2019. HyperDense-
18 Net: A Hyper-Densely Connected CNN for Multi-Modal Image Segmentation. *IEEE
19 Transactions on Medical Imaging* 38, 1116-1126.
- 20 Duta, I.C., Liu, L., Zhu, F., Shao, L., 2020. Pyramidal convolution: Rethinking convolutional
21 neural networks for visual recognition. *arXiv preprint arXiv:2006.11538*.
- 22 Fan, J., Cao, X., Yap, P.-T., Shen, D., 2019. BIRNet: Brain image registration using dual-
23 supervised fully convolutional networks. *Medical Image Analysis* 54, 193-206.
- 24 Glorot, X., Bengio, Y., 2010. Understanding the difficulty of training deep feedforward neural
25 networks. *Proceedings of the thirteenth international conference on artificial intelligence and
26 statistics. JMLR Workshop and Conference Proceedings*, pp. 249-256.
- 27 Hazlett, H.C., Gu, H., Munsell, B.C., Kim, S.H., Styner, M., Wolff, J.J., Elison, J.T., Swanson,
28 M.R., Zhu, H., Botteron, K.N., 2017. Early brain development in infants at high risk for autism
29 spectrum disorder. *Nature* 542, 348-351.

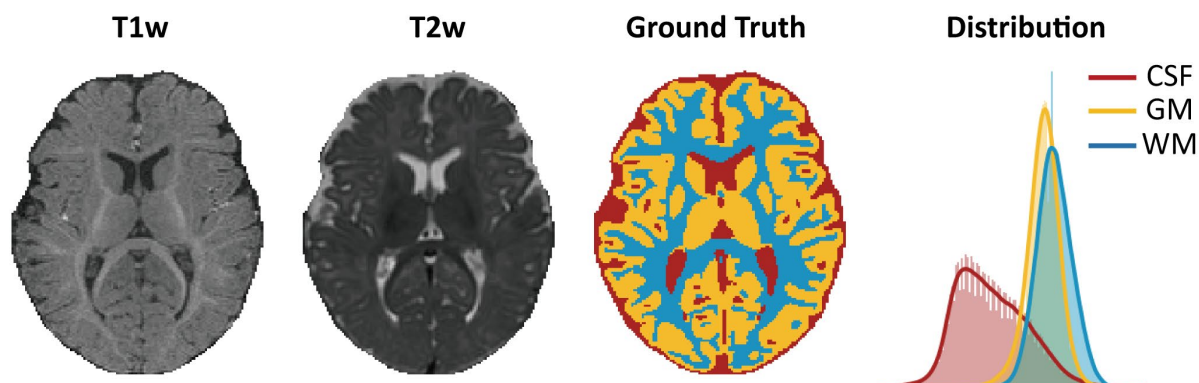
- 1 He, K., Zhang, X., Ren, S., Sun, J., 2015. Delving Deep into Rectifiers: Surpassing Human-
2 Level Performance on ImageNet Classification. 2015 IEEE International Conference on
3 Computer Vision (ICCV), 1026-1034.
- 4 Howell, B.R., Styner, M.A., Gao, W., Yap, P.-T., Wang, L., Baluyot, K., Yacoub, E., Chen, G.,
5 Potts, T., Salzwedel, A., 2019. The UNC/UMN baby connectome project (BCP): an overview of
6 the study design and protocol development. *NeuroImage* 185, 891-905.
- 7 Hu, J., Shen, L., Sun, G., 2018. Squeeze-and-excitation networks. *Proceedings of the IEEE*
8 *conference on computer vision and pattern recognition*, pp. 7132-7141.
- 9 Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q., 2017. Densely connected
10 convolutional networks. *Proceedings of the IEEE conference on computer vision and pattern*
11 *recognition*, pp. 4700-4708.
- 12 Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H., 2021. nnU-Net: a self-
13 configuring method for deep learning-based biomedical image segmentation. *Nature methods* 18,
14 203-211.
- 15 Li, G., Wang, L., Yap, P.T., Wang, F., Wu, Z., Meng, Y., Dong, P., Kim, J., Shi, F., Rekić, I.,
16 Lin, W., Shen, D., 2019. Computational neuroanatomy of baby brains: A review. *Neuroimage*
17 185, 906-925.
- 18 Li, R., Duan, C., Zheng, S., 2020. MACU-Net Semantic Segmentation from High-Resolution
19 Remote Sensing Images. *arXiv preprint arXiv:2007.13083*.
- 20 Liu, T., Li, H., Wong, K., Tarokh, A., Guo, L., Wong, S.T., 2007. Brain tissue segmentation
21 based on DTI data. *Neuroimage* 38, 114-123.
- 22 Long, J., Shelhamer, E., Darrell, T., 2015. Fully Convolutional Networks for Semantic
23 Segmentation. *IEEE transactions on pattern analysis and machine intelligence* 39, 640-651.
- 24 Makropoulos, A., Counsell, S.J., Rueckert, D., 2018. A review on automatic fetal and neonatal
25 brain MRI segmentation. *Neuroimage* 170, 231-248.
- 26 Mostapha, M., Styner, M., 2019. Role of deep learning in infant brain MRI analysis. *Magn*
27 *Reson Imaging* 64, 171-189.
- 28 Nie, D., Wang, L., Adeli, E., Lao, C., Lin, W., Shen, D., 2019. 3-D Fully Convolutional
29 Networks for Multimodal Isointense Infant Brain Image Segmentation. *IEEE Trans Cybern* 49,
30 1123-1136.

- 1 Nie, D., Wang, L., Gao, Y., Shen, D., 2016. Fully Convolutional Networks for Multi-Modality
2 Isointense Infant Brain Image Segmentation. Proc IEEE Int Symp Biomed Imaging 2016, 1342-
3 1345.
- 4 Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical
5 image segmentation. International Conference on Medical image computing and computer-
6 assisted intervention. Springer, pp. 234-241.
- 7 Sanroma, G., Benkarim, O.M., Piella, G., Lekadir, K., Hahner, N., Eixarch, E., Ballester,
8 M.A.G., 2018. Learning to combine complementary segmentation methods for fetal and 6-month
9 infant brain MRI segmentation. Computerized Medical Imaging and Graphics 69, 52-59.
- 10 Su, Z., Fang, L., Kang, W., Hu, D., Pietikäinen, M., Liu, L., 2020. Dynamic group convolution
11 for accelerating convolutional neural networks. European Conference on Computer Vision.
12 Springer, pp. 138-155.
- 13 Sun, Y., Gao, K., Wu, Z., Li, G., Zong, X., Lei, Z., Wei, Y., Ma, J., Yang, X., Feng, X., 2021.
14 Multi-site infant brain segmentation algorithms: The iSeg-2019 Challenge. IEEE Transactions on
15 Medical Imaging 40, 1363-1376.
- 16 Wang, F., Lian, C., Wu, Z., Zhang, H., Li, T., Meng, Y., Wang, L., Lin, W., Shen, D., Li, G.,
17 2019a. Developmental topography of cortical thickness during infancy. Proceedings of the
18 National Academy of Sciences 116, 15855-15860.
- 19 Wang, L., Gao, Y., Shi, F., Li, G., Gilmore, J.H., Lin, W., Shen, D., 2015. LINKS: Learning-
20 based multi-source IntegratioN framework for Segmentation of infant brain images.
21 NeuroImage 108, 160-172.
- 22 Wang, L., Li, G., Adeli, E., Liu, M., Wu, Z., Meng, Y., Lin, W., Shen, D., 2018a. Anatomy-
23 guided joint tissue segmentation and topological correction for 6-month infant brain MRI with
24 risk of autism. Hum Brain Mapp 39, 2609-2623.
- 25 Wang, L., Li, G., Shi, F., Cao, X., Lian, C., Nie, D., Liu, M., Zhang, H., Li, G., Wu, Z., Lin, W.,
26 Shen, D., 2018b. Volume-Based Analysis of 6-Month-Old Infant Brain MRI for Autism
27 Biomarker Identification and Early Diagnosis. Med Image Comput Comput Assist Interv 11072,
28 411-419.
- 29 Wang, L., Nie, D., Li, G., Puybureau, E., Dolz, J., Zhang, Q., Wang, F., Xia, J., Wu, Z., Chen, J.,
30 Thung, K.H., Bui, T.D., Shin, J., Zeng, G., Zheng, G., Fonov, V.S., Doyle, A., Xu, Y.,
31 Moeskops, P., Pluim, J.P.W., Desrosiers, C., Ayed, I.B., Sanroma, G., Benkarim, O.M.,

- 1 Casamitjana, A., Vilaplana, V., Lin, W., Li, G., Shen, D., 2019b. Benchmark on Automatic 6-
2 month-old Infant Brain Segmentation Algorithms: The iSeg-2017 Challenge. *IEEE Trans Med*
3 *Imaging*.
- 4 Wang, L., Shi, F., Gao, Y., Li, G., Gilmore, J.H., Lin, W., Shen, D., 2014. Integration of sparse
5 multi-modality representation and anatomical constraint for iso-intense infant brain MR image
6 segmentation. *NeuroImage* 89, 152-164.
- 7 Wang, L., Shi, F., Yap, P.-T., Gilmore, J.H., Lin, W., Shen, D., 2012. 4D multi-modality tissue
8 segmentation of serial infant images. *PLoS One* 7, e44596.
- 9 Wang, Z., Zou, N., Shen, D., Ji, S., 2020. Non-Local U-Nets for Biomedical Image
10 Segmentation. *Proceedings of the AAAI Conference on Artificial Intelligence* 34, 6315-6322.
- 11 Wen, X., Zhang, H., Li, G., Liu, M., Yin, W., Lin, W., Zhang, J., Shen, D., 2019. First-year
12 development of modules and hubs in infant brain functional networks. *NeuroImage* 185, 222-
13 235.
- 14 Xu, Y., Cao, M., Liao, X., Xia, M., Wang, X., Jeon, T., Ouyang, M., Chalak, L., Rollins, N.,
15 Huang, H., 2019. Development and emergence of individual variability in the functional
16 connectivity architecture of the preterm human brain. *Cerebral Cortex* 29, 4208-4222.
- 17 Yushkevich, P.A., Piven, J., Hazlett, H.C., Smith, R.G., Ho, S., Gee, J.C., Gerig, G., 2006. User-
18 guided 3D active contour segmentation of anatomical structures: significantly improved
19 efficiency and reliability. *Neuroimage* 31, 1116-1128.
- 20 Zeng, G., Zheng, G., 2018. Multi-stream 3D FCN with multi-scale deep supervision for multi-
21 modality iso-intense infant brain MR image segmentation. *international symposium on*
22 *biomedical imaging*, pp. 136-140.
- 23 Zhang, J., Jiang, Z., Liu, D., Sun, Q., Hou, Y., Liu, B., 2022. 3D asymmetric expectation-
24 maximization attention network for brain tumor segmentation. *NMR in Biomedicine* 35, e4657.
- 25 Zhang, W., Li, R., Deng, H., Wang, L., Lin, W., Ji, S., Shen, D., 2015. Deep convolutional
26 neural networks for multi-modality iso-intense infant brain image segmentation. *Neuroimage* 108,
27 214-224.
- 28 Zhao, T., Xu, Y., He, Y., 2019. Graph theoretical modeling of baby brain networks. *NeuroImage*
29 185, 711-727.

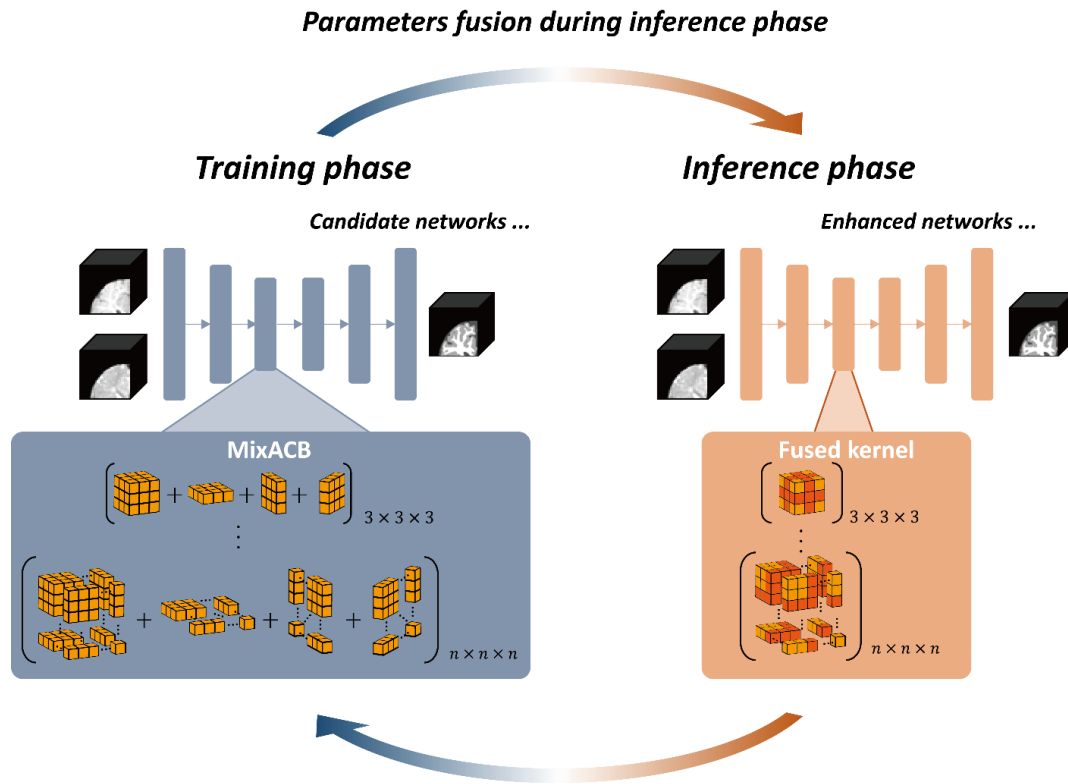
30

1 Figures and Tables

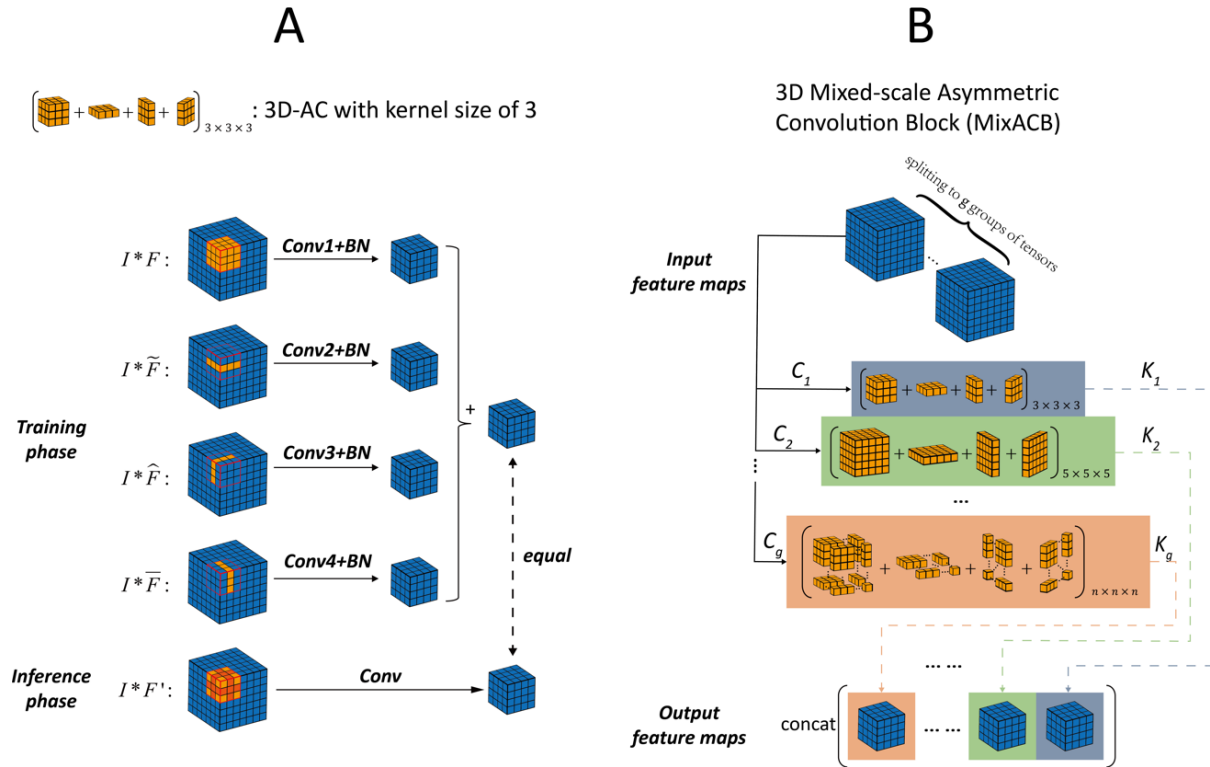


2

3 **Figure 1.** Data of a 6-month-old infant from the training set in iSeg-2019. The isointense brain
4 appearance of an axial slice in T1-weighted (T1w) and T2-weighted (T2w) images. An axial
5 view of the manual segmentation label (ground truth) and the corresponding brain tissue
6 intensity distribution of the T1w image (distribution).

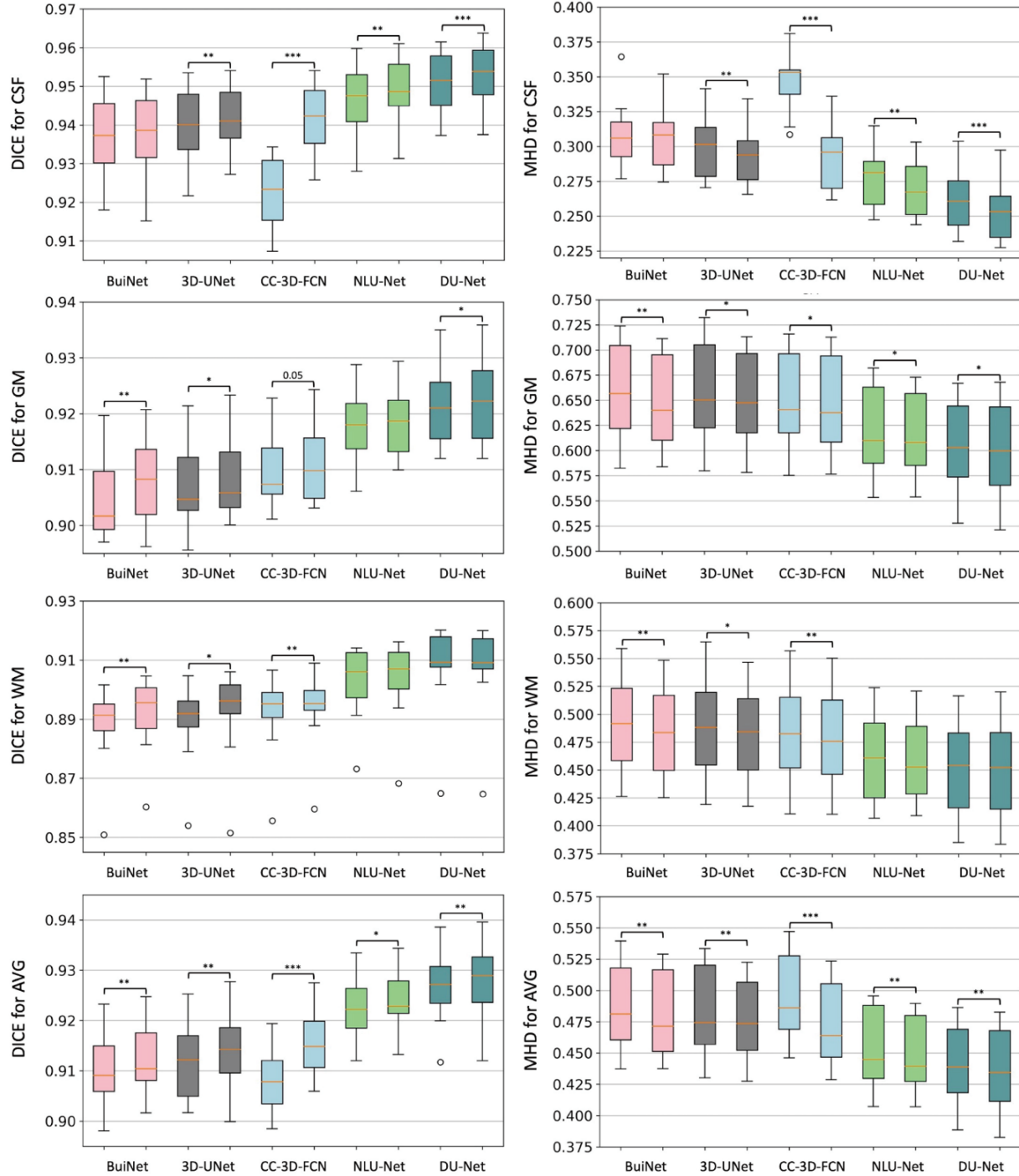


- 1 **Replace conv with MixACB during training phase**
- 2 **Figure 2.** Overview of the 3D-MASNet framework. For a candidate network, we replace its
- 3 traditional convolutional layers with MixACB during the training phase. Once the training
- 4 process is complete, we fuse the parameters of MixACB to obtain an enhanced model containing
- 5 fewer parameters after equivalent fusion.



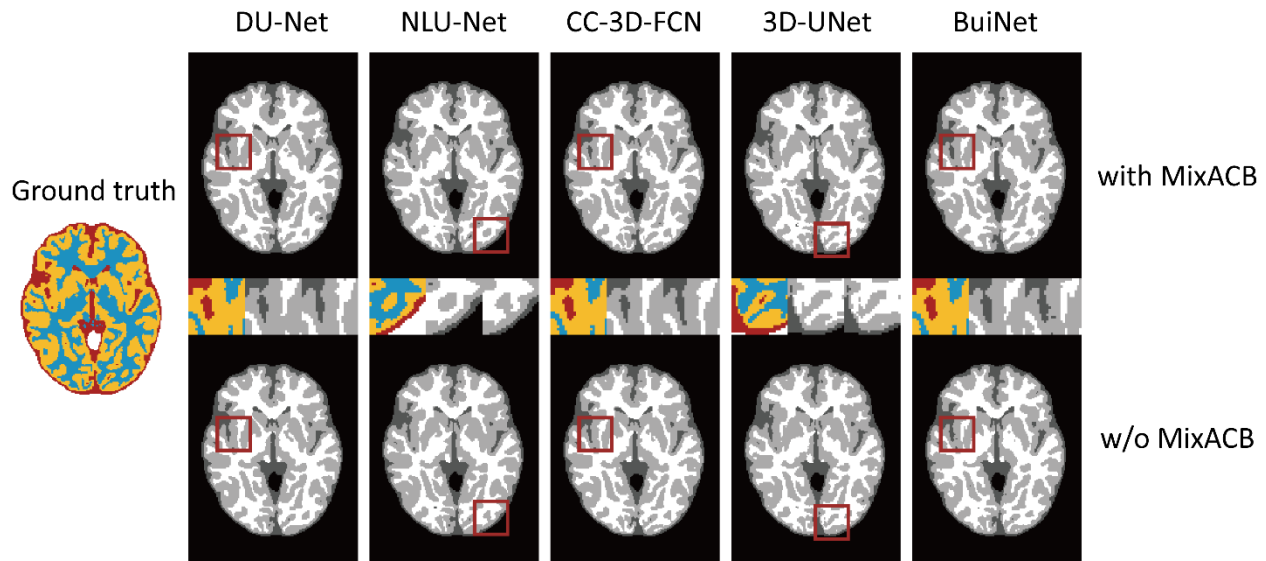
1

2 **Figure 3.** (A) Diagram of 3D-AC (taking a kernel size of 3 as an example), which has four
 3 convolutional layers during the training phase and one convolutional layer once kernel
 4 parameters have been fused during the inference phase. (B) Diagram of MixACB, which is
 5 composed of multiple 3D-ACs with different kernel sizes. MixACB splits input feature maps
 6 into several groups, applies asymmetric convolution on each group of feature maps, and then
 7 concatenates each group's output as the output feature maps.

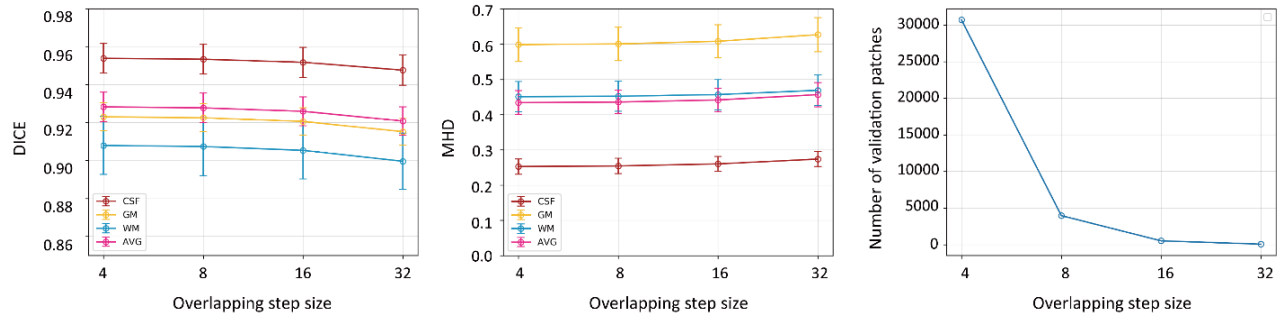


1

2 **Figure 4.** Box plot of the segmentation performance improvement on five candidate CNN
3 architectures in the 3D-MASNet framework. The first column shows the measurement of DICE
4 to represent the segmentation accuracy for each tissue type. The second column shows the results
5 of MHD. In each subgraph, we use two neighbor box plots to represent a candidate model (first
6 bar) and its corresponding 3D-MASNet (second bar). The significance of model comparison is
7 evaluated by 2-fold cross-validation. “**” denotes that $0.01 \leq p < 0.05$, “***” denotes that
8 $0.001 \leq p < 0.01$, and “0.05” denotes that $p < 0.05$.



1
2 **Figure 5.** Visualization of the segmentation results on different models with (w) and without
3 (w/o) the MixACB. The ground truth map is shown in color, and CNNs-based segmentation
4 maps are shown in the gray scale. The regions in the red square are magnified in the middle row
5 following an order from with MixACB to without MixACB.



1
2 **Figure 6.** Changes in segmentation performance in terms of DICE (A) and MHD (B) with
3 respect to different overlapping step sizes on 10 subjects during inference, where 2-fold cross-
4 validation is used. (C) Changes of the average number of the 10 subjects' patches with respect to
5 different overlapping step sizes during inference.

1 **Table 1.** Training strategy of each candidate network

Candidate Network	Training Batch Size	Training/Inference Patch Size	Learning Rate Schedule
BuiNet	4	64	Train for 20,000 iterations. The initial learning rate is set to $2e-4$ and is decreased by a factor of 0.1 every 5000 iterations.
3D-UNet	10	32	Train for 80 epochs for a total of 5000 patches that are randomly extracted per epoch. The learning rate is decreased every 20 epochs and is set to $3e-4$, $1e-4$, $1e-5$ and $1e-6$. Train for 80 epochs for a total of 5000 patches that are randomly extracted per epoch. The learning rate is decreased every 20 epochs and is set to $3e-4$, $1e-4$, $1e-5$ and $1e-6$.
CC-3D-FCN	10	32	The same as 3D-UNet.
NLU-Net	5	32	Train for 80 epochs for a total of 5000 patches that are randomly extracted per epoch. The learning rate is set to $1e-3$.
DU-Net	16	32	The cosine annealing strategy with a maximum learning rate of $3e-4$ and a minimum learning rate of $1e-6$ is adopted. The model is trained for 500 epochs and a total of 1000 patches are randomly extracted at each epoch.

2

1 **Table 2.** Ablation study performed by comparing the segmentation accuracy between different
 2 models and their corresponding 3D-MASNet in terms of DICE by 2-fold cross validation.

Netw ork	CSF		GM		WM		Avg	
	Baselin e	MixACB	Baselin e	MixACB	Baselin e	MixACB	Baselin e	MixACB
BuiN et	0.938±0 .010	0.938±0. 011	0.905±0 .007	0.908*± 0.007	0.888±0 .014	0.892*± 0.013	0.910±0 .007	0.912*± 0.007
3D- UNet	0.940±0 .010	0.942*± 0.008	0.907±0 .007	0.909*± 0.007	0.889±0 .014	0.892*± 0.015	0.912±0 .007	0.914*± 0.008
CC- 3D- FCN	0.923±0 .010	0.942*± 0.008	0.910±0 .006	0.911±0. 007	0.892±0 .013	0.894*± 0.013	0.908±0 .006	0.915*± 0.006
NLU -Net	0.947±0 .009	0.949*± 0.008	0.918±0 .007	0.919±0. 006	0.903±0 .012	0.904±0. 014	0.922±0 .006	0.924*± 0.006
DU- Net	0.951±0 .008	0.953*± 0.008	0.922±0 .007	0.923*± 0.007	0.907±0 .015	0.907±0. 015	0.927±0 .007	0.928*± 0.008

3 Note that the best values are highlighted in bold font. “Baseline” denotes that the corresponding
 4 model adopted the standard convolutional operation; “MixACB” denotes that the corresponding
 5 model was transformed into 3D-MASNet; “*” denotes that the difference between baseline and
 6 3D-MASNet is statistically significant (p<0.05).

1 **Table 3.** Ablation study performed by comparing the segmentation accuracy between different
 2 models and their corresponding 3D-MASNet in terms of MHD by 2-fold cross validation.

Netw ork	CSF		GM		WM		Avg	
	Baselin e	MixACB	Baselin e	MixACB	Baselin e	MixACB	Baselin e	MixACB
BuiN et	0.308±0 .024	0.307±0. 023	0.659±0 .048	0.649*± 0.045	0.493±0 .043	0.485*± 0.042	0.487±0 .035	0.480*± 0.034
3D- UNet	0.299±0 .022	0.293*± 0.020	0.658±0 .050	0.651*± 0.046	0.490±0 .046	0.485*± 0.042	0.483±0 .036	0.476*± 0.033
CC- 3D- FCN	0.348±0 .022	0.292*± 0.023	0.649±0 .047	0.645*± 0.048	0.485±0 .046	0.480*± 0.046	0.494±0 .034	0.473*± 0.034
NLU -Net	0.278±0 .022	0.270*± 0.020	0.619±0 .043	0.615*± 0.040	0.461±0 .040	0.460±0. 037	0.453±0 .032	0.448*± 0.030
DU- Net	0.261±0 .021	0.254*± 0.022	0.605±0 .046	0.601*± 0.047	0.452±0 .041	0.452±0. 043	0.439±0 .032	0.436*± 0.034

3 Note that the best values are highlighted in bold font. “Baseline” denotes that the corresponding
 4 model adopted the standard convolutional operation; “MixACB” denotes that the corresponding
 5 model was transformed into 3D-MASNet; “*” denotes that the difference between baseline and
 6 3D-MASNet is statistically significant ($p < 0.05$).
 7

1 **Table 4.** Comparison of segmentation performance of the proposed method and the methods of
 2 the top-4 ranked teams on the 13 validation infant MRI images of iSeg-2019.

Method (Top 5)	CSF			GM			WM			AVG		
	DICE	MHD	ASD	DICE	MHD	ASD	DICE	MHD	ASD	DICE	MHD	ASD
Brain_Tech	0.961	8.873	0.108	0.928	5.724	0.300	0.911	7.114	0.347	0.933	7.237	0.252
FightAutism	0.960	9.233	0.110	0.929	5.678	0.300	0.911	6.678	0.341	0.933	7.196	0.250
OxfordIBME	0.960	8.560	0.112	0.927	5.495	0.307	0.907	6.759	0.353	0.931	6.938	0.257
QL111111	0.959	9.484	0.114	0.926	5.601	0.307	0.908	7.028	0.353	0.931	7.371	0.258
Proposed	0.961	9.293	0.107	0.931	5.741	0.292	0.912	7.111	0.332	0.935	7.382	0.244

3 Note that the best values are highlighted in bold font.

1 **Table 5.** Ablation study performed by comparing the segmentation accuracy in different mix
 2 ratios by 2-fold cross validation.

Mix Ratio	CSF		GM		WM		AVG	
	DICE	MHD	DICE	MHD	DICE	MHD	DICE	MHD
1:0	0.952±0 .010	0.261±0 .024	0.922±0 .008	0.604±0 .045	0.906±0 .014	0.453±0 .041	0.927±0 .008	0.440±0 .033
1:1	0.953±0 .008	0.258±0 .023	0.921±0 .007	0.605±0 .045	0.905±0 .013	0.455±0 .040	0.926±0 .007	0.439±0 .033
3:1 (propose)	0.953±0 .008	0.254±0 .022	0.923±0 .007	0.601±0 .047	0.907±0 .015	0.452±0 .043	0.928±0 .008	0.436±0 .034
5:1	0.953±0 .009	0.257±0 .025	0.922±0 .008	0.601±0 .047	0.907±0 .015	0.452±0 .042	0.926±0 .008	0.437±0 .034

3 Note that the best values are highlighted in bold font.

1 **Table 6.** Component analysis of MixACB by 2-fold cross validation.

	CSF		GM		WM		AVG	
	DICE	MHD	DICE	MHD	DICE	MHD	DICE	MHD
CON	0.951±0	0.261±0	0.922±0	0.605±0	0.907±0	0.452±0	0.927±0	0.439±0
V_3	.008	.021	.007	.046	.015	.041	.007	.032
AC_3	0.952±0	0.261±0	0.922±0.	0.604±0	0.906±0	0.453±0	0.927±0	0.440±0
	.010	.024	.008	.045	.014	.041	.008	.033
CON	0.947±0	0.276±0	0.918±0.	0.619±0	0.903±0	0.463±0	0.922±0	0.453±0
V_5	.012	.023	.008	.047	.016	.043	.008	.034
AC_5	0.952±0	0.261±0	0.920±0.	0.610±0	0.904±0	0.460±0	0.925±0	0.443±0
	.008	.022	.008	.046	.016	.043	.008	.033
MixA	0.953±0	0.254±0	0.923±0.	0.601±0	0.907±0	0.452±0	0.928±0	0.436±0
CB	.008	.022	.007	.047	.015	.043	.008	.034

2 Note that the best values are highlighted in bold font. “CONV_3” denotes that the 3D
3 convolution with a kernel size of 3; “AC_3” denotes that the 3D-AC with a kernel size of 3;
4 “CONV_5” denotes that the 3D convolution with a kernel size of 5; “AC_5” denotes that the 3D-
5 AC with a kernel size of 5.