

HyperDense-Net: A hyper-densely connected CNN for multi-modal image segmentation

Jose Dolz, Karthik Gopinath, Jing Yuan, Herve Lombaert, Christian Desrosiers, and Ismail Ben Ayed

Abstract—Recently, dense connections have attracted substantial attention in computer vision because they facilitate gradient flow and implicit deep supervision during training. Particularly, DenseNet, which connects each layer to every other layer in a feed-forward fashion, has shown impressive performances in natural image classification tasks. We propose *HyperDenseNet*, a 3D fully convolutional neural network that extends the definition of dense connectivity to multi-modal segmentation problems. Each imaging modality has a path, and dense connections occur not only between the pairs of layers within the same path, but also between those across different paths. This contrasts with the existing multi-modal CNN approaches, in which modeling several modalities relies entirely on a single joint layer (or level of abstraction) for fusion, typically either at the input or at the output of the network. Therefore, the proposed network has total freedom to learn more complex combinations between the modalities, *within and in-between all the levels of abstraction*, which increases significantly the learning representation. We report extensive evaluations over two different and highly competitive multi-modal brain tissue segmentation challenges, iSEG 2017 and MRBrainS 2013, with the former focusing on 6-month infant data and the latter on adult images. *HyperDenseNet* yielded significant improvements over many state-of-the-art segmentation networks, ranking at the top on both benchmarks. We further provide a comprehensive experimental analysis of features re-use, which confirms the importance of hyper-dense connections in multi-modal representation learning. Our code is publicly available.

Index Terms—Deep learning, brain MRI, segmentation, 3D CNN, multi-modal imaging

1 INTRODUCTION

MULTI-MODAL imaging is of primary importance for developing comprehensive models of pathologies and increasing the statistical power of current imaging biomarkers [1]. In neuroimaging studies, different magnetic resonance imaging (MRI) modalities are often combined to overcome the limitations of independent imaging techniques. While T1-weighted images yield a good contrast between gray matter (GM) and white matter (WM) tissues, T2-weighted and proton density (PD) pulses help visualize tissue abnormalities like lesions. Likewise, fluid attenuated inversion recovery (FLAIR) images can enhance the image contrast of white matter lesions resulting from multiple sclerosis [2]. In brain segmentation, considering multiple MRI modalities is essential to obtain accurate results. This is particularly true for the segmentation of infant brains, where tissue contrast is low (Fig. 1).

Advances in multi-modal imaging, however, come at the price of an inherently large amount of data, imposing a burden on disease assessments. Visual inspections of such an enormous amount of medical images are prohibitively time-consuming, prone to errors and unsuitable for large-scale studies. Therefore, automatic and reliable multi-modal segmentation algorithms are of high interest to the clinical community.

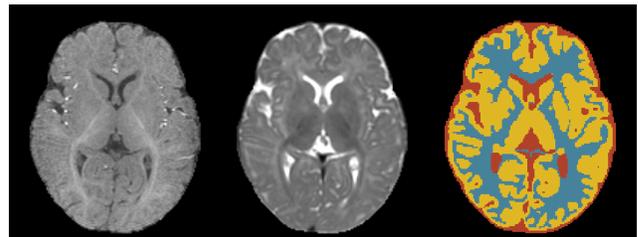


Fig. 1. Example of data from a training subject. Neonatal isointense brain images from a mid-axial T1 slice (*left*), the corresponding T2 slice (*middle*), and manual segmentation (*right*).

1.1 Prior work

Multi-modal image segmentation in brain-related applications has received a substantial research attention, for instance, brain tumors [3]–[6], brain tissues of both infant [7]–[17] and adult [18], [19], subcortical structures [20], among other problems [21]–[23]. Atlas-propagation approaches are commonly used in multi-modal scenarios [24], [25]. These methods rely on registering one or multiple atlases to the target image, followed by a propagation of manual labels. When several atlases are considered, labels from individual atlases can be combined into a final segmentation via a label fusion strategy [8], [10], [13]. When relying solely on atlas fusion, the performance of such techniques might be limited and prone to registration errors. Parametric or deformable models [11] can be used to refine prior estimates of tissue probability [14]. For example, the study in [14] investigated a patch-driven method for neonatal brain tissue segmentation, integrating the probability maps of a subject-specific atlas into a level-set framework.

• J. Dolz, K. Gopinath, H. Lombaert, C. Desrosiers and I. Ben Ayed are with the École de technologie Supérieure, Montreal, Canada. email:jose.dolz@livioa.etsmtl.ca

• J. Yuan is with the Xidian University, School of Mathematics and Statistics, Xi'an, China.

Manuscript received XXX; revised XXX.

More recently, our community has witnessed a wide adoption of deep learning techniques, particularly, convolutional neural networks (CNNs), as an effective alternative to traditional segmentation approaches. CNN architectures are supervised models, trained end-to-end, to learn a hierarchy of image features representing different levels of abstraction. In contrast to conventional classifiers based on hand-crafted features, CNNs can learn both the features and classifier simultaneously, in a data-driven manner. They achieved state-of-the-art performances in a broad range of medical image segmentation problems [26], [27], including multi-modal tasks [4]–[6], [15]–[17], [19], [22], [23], [28], [29].

1.1.1 Fusion of multi-modal CNN feature representations

Most of the existing multi-modal CNN segmentation techniques followed an *early-fusion* strategy, which integrates the multi-modality information from the original space of low-level features [5], [15], [19], [23], [28], [29]. For instance, in [15], MRI T1, T2 and fractional anisotropy (FA) images are simply merged at the input of the network. However, as argued in [30] in the context of multi-modal learning, it is difficult to discover highly non-linear relationships between the low-level features of different modalities, more so when such modalities have significantly different statistical properties. In fact, early-fusion methods implicitly assume that the relationship between different modalities are simple (e.g., linear). For instance, the early fusion in [15] learns complementary information from T1, T2 and FA images. However, the relationship between the original T1, T2 and FA image data may be much more complex than complementarity, due to significantly different image acquisition processes [16]. The work in [16] advocated *late fusion* of high-level features as a way that accounts better for the complex relationships between different modalities. They used an independent convolutional network for each modality, and fused the outputs of the different networks in higher-level layers, showing better performance than early fusion in the context infant brain segmentation. These results are in line with a recent study in the machine learning community [30], which investigated multimodal learning with deep Boltzmann machines in the context of fusing data from color images and text.

1.1.2 Dense connections in deep networks

Since the recent introduction of residual learning in [32], shortcut connections from early to late layers have become very popular in a breadth of computer vision problems [33], [34]. Unlike traditional networks, these connections back-propagate gradients directly, thereby mitigating the gradient-vanishing problem and allowing deeper networks. Furthermore, they transform a whole network into a large ensemble of shallower networks, yielding competitive performances in various applications [19], [35]–[37]. DenseNet [38] extended the concept of shortcut connections, with the input of each layer corresponding to the outputs from all previous layers. Such a dense network facilitates the gradient flow and the learning of more complex patterns, which yielded significant improvements in accuracy and efficiency for natural image classification tasks [38]. Inspired by this success, recent works have included dense connections in deep networks for medical image segmentation

[39]–[41]. However, these works have either considered a single modality [39], [40] or have simply concatenated multiple modalities in a single stream [41]. So far, the impact of dense connectivity across multiple network paths, and its application to multi-modal image segmentation, remains unexplored.

1.2 Contributions

We propose *HyperDenseNet*, a 3D fully convolutional neural network that extends the definition of dense connectivity to multi-modal segmentation problems. Each imaging modality has a path, and dense connections occur not only between the pairs of layers within the same path, but also between those across different paths; see the illustration in Fig. 2. This contrasts with the existing multi-modal CNN approaches, in which modeling several modalities relies entirely on a single joint layer (or level of abstraction) for fusion, typically either at the input (early fusion) or at the output (late fusion) of the network. Therefore, the proposed network has total freedom to learn more complex combinations between the modalities, *within and in-between all the levels of abstractions*, which increases significantly the learning representation in comparison to early/late fusion. Furthermore, hyper-dense connections facilitate the learning as they improve gradient flow and impose implicit deep supervision. We report extensive evaluations over two different¹ and highly competitive multi-modal brain tissue segmentation challenges, iSEG 2017 and MRBrainS 2013. *HyperDenseNet* yielded significant improvements over many state-of-the-art segmentation networks, ranking at the top on both benchmarks. We further provide a comprehensive experimental analysis of features re-use, which confirms the importance of hyper-dense connections in multi-modal representation learning. Our code is publicly available².

A preliminary conference version of this work appeared at ISBI 2018 [42]. This journal version is a substantial extension, including (1) a much broader, more informative/rigorous treatment of the subject in the general context of multi-modal segmentation; and (2) comprehensive experiments with additional baselines and publicly available benchmarks, as well as a thorough investigation of the practical usefulness and impact of hyper-dense connections.

2 METHODS AND MATERIALS

Convolutional neural networks (CNNs) are deep models that can learn feature representations automatically from the training data. They consist of multiple layers, each processing the imaging data at a different level of abstraction, enabling segmentation algorithms to learn from large datasets and discover complex patterns that can be further employed for predicting unseen samples. The first attempts to use CNNs in segmentation problems followed a sliding-window strategy, where the regions defined by the window are processed independently, which impedes segmentation accuracy and computational efficiency. To overcome these

1. iSEG 2017 focuses on 6-month infant data, whereas MRBrainS 2013 uses adult data. Therefore, there are significant differences between the two benchmarks in term of image data characteristics, e.g. the voxel spacing and number of available modalities.

2. <https://www.github.com/josedolz/HyperDenseNet>

TABLE 1
Overview of representative works on multi-modal brain segmentation.

Work	Modality	Target	Method
Prastawa et al., 2005 [7]	T1,T2	Infant brain tissue	Multi-atlas
Weisenfeld et al., 2006 [8]	T1,T2	Infant brain tissue	Multi-atlas
Deoni et al., 2007 [20]	T1,T2	Thalamic nuclei	K-means clustering
Anbeek et al., 2008 [9]	T2,IR	Infant brain tissue	KNN
Weisenfeld and Warfield, 2009 [10]	T1,T2	Infant brain tissue	Multi-atlas
Wang et al., 2011 [11]	T1,T2,FA	Infant brain tissue	Multi-atlas + Level sets
Srhoj et al., 2012 [12]	T1,T2	Infant brain tissue	Multi-atlas + KNN
Wang et al., 2012 [13]	T1,T2	Infant brain tissue	Multi-atlas
Wang et al., 2014 [31]	T1,T2,FA	Infant brain tissue	Multi-atlas + Level sets
Kamnitsas et al., 2015 [28]	Flair, DWI, T1, T2	Brain lesion	3D FCNN + CRF
Zhang et al., 2015 [15]	T1,T2,FA	Infant brain tissue	2D CNN
Havaei et al., 2016 [4]	T1,T1c,T2,FLAIR	Multiple Sclerosis/Brain tumor	2D CNN
Nie et al., 2016 [16]	T1,T2,FA	Infant brain tissue	2D FCNN
Chen et al., 2017 [19]	T1,T1-IR,FLAIR	Brain tissue	3D FCNN
Dolz et al., 2017 [17]	T1,T2	Infant brain tissue	3D FCNN
Fidon et al., 2017 [6]	T1,T1c,T2,FLAIR	Brain tumor	CNN
Kamnitsas et al., 2017 [5]	T1,T1c,T2,FLAIR	Brain tumour/lesions	3D FCNN + CRF
Kamnitsas et al., 2017 [22]	MPRAGE,FLAIR,T2,PD	Traumatic brain injuries	3D FCNN(Adversarial Training)
Valverde et al., 2017 [23]	T1, T2,FLAIR	Multiple-sclerosis	3D FCNN

limitations, the network can be viewed as a single non-linear convolution, which is trained end-to-end, a process known as fully CNN (FCNN) [43]. The latter brings several advantages over standard CNNs. It can handle images of arbitrary sizes and avoid redundant convolution and pooling operations, enabling computationally efficient learning.

2.1 The proposed Hyper-Dense network

The concept of “the deeper the better” is considered as a key principle in deep learning [32]. Nevertheless, one obstacle when dealing with deep architectures is the problem of vanishing/exploding gradients, which hampers convergence during training. To address these limitations in very deep architectures, the study in [38] investigated densely connected networks. DenseNets are built on the idea that adding direct connections from any layer to all the subsequent layers in a feed-forward manner makes training easier and more accurate. This is motivated by three observations. First, there is an implicit deep supervision thanks to the short paths to all feature maps in the architecture. Second, direct connections between all layers help improving the flow of information and gradients throughout the entire network. Third, dense connections have a regularizing effect, which reduces the risk of over-fitting on tasks with smaller training sets.

Inspired by the recent success of densely-connected networks in medical image segmentation works [39]–[41], we propose a hyper-dense architecture for multi-modal image segmentation that extends the concept of dense connectivity to the multi-modal setting: each imaging modality has a path, and dense connections occur not only between layers within the same path, but also between layers across different paths (see Fig. 2 for an illustration).

Let \mathbf{x}_l be the output of the l^{th} layer. In CNNs, this vector is typically obtained from the output of the previous layer

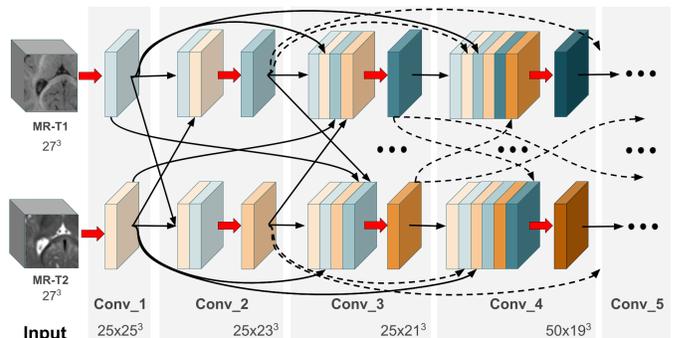


Fig. 2. A section of the proposed HyperDenseNet in the case of two image modalities. Each gray region represents a convolutional block. Red arrows correspond to convolutions and black arrows indicate dense connections between feature maps.

\mathbf{x}_{l-1} by a mapping H_l composed of a convolution followed by a non-linear activation function:

$$\mathbf{x}_l = H_l(\mathbf{x}_{l-1}). \quad (1)$$

A densely-connected network concatenates all feature outputs in a feed-forward manner,

$$\mathbf{x}_l = H_l([\mathbf{x}_{l-1}, \mathbf{x}_{l-2}, \dots, \mathbf{x}_0]), \quad (2)$$

where $[\dots]$ denotes a concatenation operation.

Pushing this idea further, HyperDenseNet introduces a more general connectivity definition, in which we link the outputs from layers in different streams, each associated with a different image modality. In the multi-modal setting, our hyper-dense connectivity yields a much more powerful feature representation than early/late fusion as the network learns the complex relationships between the modalities within and in-between all the levels of abstractions. For simplicity, let us consider the scenario of two image modalities,

although extension to N modalities is straightforward. Let \mathbf{x}_l^1 and \mathbf{x}_l^2 denote the outputs of the l^{th} layer in streams 1 and 2, respectively. In general, the output of the l^{th} layer in a stream s can then be defined as follows:

$$\mathbf{x}_l^s = H_l^s([\mathbf{x}_{l-1}^1, \mathbf{x}_{l-1}^2, \mathbf{x}_{l-2}^1, \mathbf{x}_{l-2}^2, \dots, \mathbf{x}_0^1, \mathbf{x}_0^2]). \quad (3)$$

Shuffling and interleaving feature map elements in a CNN was recently found to enhance the efficiency and performance, while serving as a strong regularizer [44]–[46]. This is motivated by the fact that intermediate CNN layers perform deterministic transformations to improve the performance, however, relevant information might be lost during these operations [47]. To overcome this issue, it is therefore beneficial for intermediate layers to offer a variety of information exchange while preserving the aforementioned deterministic functions. Motivated by this principle, we thus concatenate feature maps in a different order for each branch and layer:

$$\mathbf{x}_l^s = H_l^s(\pi_l^s([\mathbf{x}_{l-1}^1, \mathbf{x}_{l-1}^2, \mathbf{x}_{l-2}^1, \mathbf{x}_{l-2}^2, \dots, \mathbf{x}_0^1, \mathbf{x}_0^2])), \quad (4)$$

with π_l^s being a function that permutes the feature maps given as input. For instance, in the case of two image modalities, we could have:

$$\begin{aligned} \mathbf{x}_l^1 &= H_l^1([\mathbf{x}_{l-1}^1, \mathbf{x}_{l-1}^2, \mathbf{x}_{l-2}^1, \mathbf{x}_{l-2}^2, \dots, \mathbf{x}_0^1, \mathbf{x}_0^2]) \\ \mathbf{x}_l^2 &= H_l^2([\mathbf{x}_{l-1}^2, \mathbf{x}_{l-1}^1, \mathbf{x}_{l-2}^2, \mathbf{x}_{l-2}^1, \dots, \mathbf{x}_0^2, \mathbf{x}_0^1]) \end{aligned}$$

Figure 2 shows a section of the proposed architecture, where each gray region represents a convolutional block. For simplicity, we assume that the red arrows indicate convolution operations only, whereas the black arrows represent the direct connections between feature maps from different layers, within and in-between the different streams. Thus, the input of each convolutional block (maps before the red arrow) is the concatenation of the outputs (maps after the red arrow) of all the preceding layers from both paths.

2.2 Baselines

To investigate thoroughly the impact of hyper-dense connections between different streams in multi-modal image segmentation, several baselines were considered. First, we extended the semi-dense architecture proposed in [17] to a fully-dense one, by connecting the output of each convolutional layer to all subsequent layers. In this network, we follow an early-fusion strategy, in which MRI T1 and T2 are integrated at the input of the CNN and processed jointly along a single path (Fig. 3, *left*). The connectivity setting of this model corresponds to Eq. (2). Second, instead of merging both modalities at the input of the network, we considered a late-fusion strategy, where each modality is processed independently in different streams and learned features are fused before the first fully connected layer (Fig. 3, *middle*). In this model, the dense connections are included within each path, assuming the connectivity definition of Eq. (2) for each stream.

As last baseline, we used an early fusion model which combines features from different streams *after the first convolutional layer* (Fig. 3, *right*). Since this non-linear combination of features is re-used in all subsequent layers, the resulting network is similar to our hyper-dense model of Eq. (3).

However, there are two important differences. First, each stream in our model processes its input differently, as shown by the stream-indexed function H_l^s in Eq. (3). Also, as described above, each stream performs a different shuffling of inputs, which can enhance robustness to the model and mitigate the risk of overfitting. Our experiments in Section 3 demonstrate empirically the advantages of our model compared to this baseline.

2.3 Network architecture

To have a large receptive field, FCNNs typically use full images as input. The number of parameters is then limited via pooling/unpooling layers. A problem with this approach is the loss of resolution from repeated down-sampling operations. In the proposed method, we follow the strategy in [5], where sub-volumes are used as input, avoiding pooling layers. While sub-volumes of size $27 \times 27 \times 27$ are considered for training, we used $35 \times 35 \times 35$ non-overlapping sub-volumes during inference, as in [5], [26]. This strategy offers two considerable benefits. First, it reduces the memory requirements of our network, thereby removing the need for spatial pooling. More importantly, it substantially increases the number of training examples and, therefore, does not need data augmentation.

TABLE 2

The layers used in the baselines and the proposed architecture and the corresponding values with an input of size $27 \times 27 \times 27$. In the case of multi-modal images, the convolutional layers (conv_x) are present in any network path. All the convolutional layers have a stride of one pixel.

	Conv. kernel	# kernels	Output Size	Dropout
conv_1	$3 \times 3 \times 3$	25	$25 \times 25 \times 25$	No
conv_2	$3 \times 3 \times 3$	25	$23 \times 23 \times 23$	No
conv_3	$3 \times 3 \times 3$	25	$21 \times 21 \times 21$	No
conv_4	$3 \times 3 \times 3$	50	$19 \times 19 \times 19$	No
conv_5	$3 \times 3 \times 3$	50	$17 \times 17 \times 17$	No
conv_6	$3 \times 3 \times 3$	50	$15 \times 15 \times 15$	No
conv_7	$3 \times 3 \times 3$	75	$13 \times 13 \times 13$	No
conv_8	$3 \times 3 \times 3$	75	$11 \times 11 \times 11$	No
conv_9	$3 \times 3 \times 3$	75	$9 \times 9 \times 9$	No
fully_conv_1	$1 \times 1 \times 1$	400	$9 \times 9 \times 9$	Yes
fully_conv_2	$1 \times 1 \times 1$	200	$9 \times 9 \times 9$	Yes
fully_conv_3	$1 \times 1 \times 1$	150	$9 \times 9 \times 9$	Yes
Classification	$1 \times 1 \times 1$	4	$9 \times 9 \times 9$	No

Table 2 summarizes the parameters of the baselines and the proposed HyperDenseNet. The network parameters are optimized via the RMSprop optimizer, using cross-entropy as cost function. Let θ denotes the network parameters (i.e., convolution weights, biases and a_i from the parametric rectifier units), and y_s^v the label of voxel v in the s -th image segment. We optimize the following:

$$J(\theta) = -\frac{1}{S \cdot V} \sum_{s=1}^S \sum_{v=1}^V \sum_{c=1}^C \delta(y_s^v = c) \cdot \log p_c^v(\mathbf{x}_s), \quad (5)$$

where $p_c^v(\mathbf{x}_s)$ is the softmax output of the network for voxel v and class c , when the input segment is \mathbf{x}_s .

To initialize the weights of the network, we adopted the strategy proposed in [48], which yields fast convergence for very deep architectures. In this strategy, a zero-mean Gaussian distribution of standard deviation $\sqrt{2/n_l}$ is used to initialize the weights in layer l , where n_l denotes the

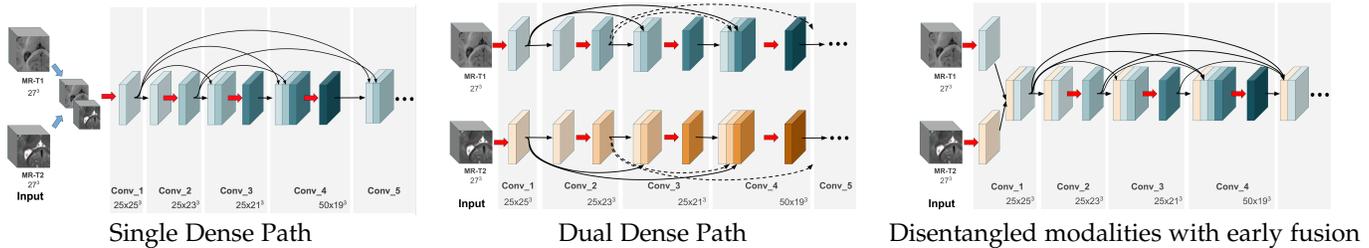


Fig. 3. Section of baseline architectures: single-path dense (*left*), dual-path dense (*middle*) with disentangled modalities and disentangled modalities with early fusion in a single path (*right*). While both modalities are concatenated at the input of the network in the first case, each modality is analyzed independently in the second architecture with the features being fused at the end of the streams. Each gray region represents a convolutional block. Red arrows correspond to convolutions and black arrows indicate dense connections between feature maps. Dense connections are propagated through the entire network.

number of connections to the units in that layer. Momentum was set to 0.6 and the initial learning rate to 0.001, being reduced by a factor of 2 after every 5 epochs (starting from epoch 10). The network was trained for 30 epochs, each composed of 20 subepochs. At each subepoch, a total of 1000 samples were randomly selected from the training images and processed in batches of size 5.

3 EXPERIMENTS AND RESULTS

The proposed HyperDenseNet architecture is evaluated on challenging multi-modal image segmentation tasks, using publicly available data from two challenges: infant brain tissue segmentation, iSEG [49], and adult brain tissue segmentation, MRBrainS³. Quantitative evaluations and comparisons with state-of-the-art methods are reported for each of these applications. First, to evaluate the impact of dense connectivity on performance, we compared the proposed HyperDenseNet to the baselines described in Section 2.2 on infant brain tissue segmentation. Then, our results, compiled by the iSEG challenge organizers on testing data, are compared to those from the other competing teams. Second, to juxtapose the performance of HyperDenseNet to other segmentation networks under the same conditions, we provide a quantitative analysis of the results of current state-of-the-art segmentation networks for adult brain tissue segmentation. This includes comparison to the participants the MRBrainS challenge. Finally, in Section 3.3, we report a comprehensive analysis of feature re-use.

3.1 iSEG Challenge

The focus of this challenge was to compare (semi-) automatic state-of-the-art algorithms for the segmentation of 6-month infant brain tissues in T1- and T2-weighted brain MRI scans. This challenge was carried out in conjunction with MICCAI 2017, with a total of 21 international teams participating in the first round [49].

3.1.1 Evaluation

The iSEG-2017 organizers used three metrics to evaluate the accuracy of the competing methods: Dice Similarity Coefficient (DSC) [50], Modified Hausdorff distance (MHD), where the 95-*th* percentile of all Euclidean distances is

employed, and Average Surface Distance (ASD). The first measures the degree of overlap between the segmentation region and ground truth, whereas the other two evaluate boundary distances.

3.1.2 Results

Validation results: Table 3 reports the performance achieved by HyperDenseNet and the baselines introduced in Section 2.2, for CSF, GM and WM brain tissues. The results were generated by splitting the 10 available iSEG-2017 volumes into training, validation and testing sets containing 6, 1 and 3 volumes, respectively. To show that improvements do not come from the higher number of learned parameters in HyperDenseNet, we also investigated a widened version of all baselines, with a similar parameter size as HyperDenseNet. The number of learned parameters of all the tested models is given in Table 4. A more detailed description of the tested architectures can be found in Table 8 of the Supplemental materials (‘Supplementary materials are available in the supplementary files /multimedia tab.’).

We observe that the late fusion of deeper-layer features in independent paths provides a clear improvement over the single-path version, with an increase on performance of nearly 5%. Fusing the feature maps from independent paths after the first convolutional layer (i.e., Dual-Single) outperformed the other two baselines by 1-2%, particularly for WM and GM, which are the most challenging structures to segment. Also, the results indicate that processing multi-modal data in separate paths, while allowing dense connectivity between all the paths, increases performance over early and late fusion, as well as over disentangled modalities with fusion performed after the first convolutional block. Another interesting finding is that increasing the number of learned parameters does not bring an important boost in performance. Indeed, in some tissues (e.g., CSF for Single path and Dual-Single path architectures), the performance slightly decreased when widening the architecture.

Figures 4 and 5 compare the training and validation accuracy between the baselines and HyperDenseNet. In these figures, the mean DSC for the three brain tissues is evaluated during training (*Top*) and validation (*Bottom*). One can see that HyperDenseNet outperforms baselines in both cases, achieving better results than architectures with a similar number of parameters. Performance improvements seen in Table 3, Fig. 4 and Fig. 5 might be due to two factors:

3. <http://mrbrains13.isi.uu.nl>

TABLE 3
Performance on the testing set, in terms of DSC, for the investigated baselines and the proposed architecture. The best performance is highlighted in bold.

	Architectures	CSF	WM	GM
No connectivity between paths	Single Path	0.9014	0.8518	0.8370
	Single Path*	0.9010	0.8532	0.8401
	Dual Path	0.9482	0.9078	0.8875
	Dual Path*	0.9503	0.9089	0.8872
Connectivity between paths	Dual-Single Path	0.9552	0.9142	0.9008
	Dual-Single Path*	0.9541	0.9159	0.9017
	HyperDenseNet	0.9580	0.9183	0.9035

* Widened version.

the high number of direct connections between different layers, which facilitates back-propagation of the gradient to shallow layers, and the freedom of the network to explore more complex patterns thanks to the combination of several image modalities at any level of abstraction.

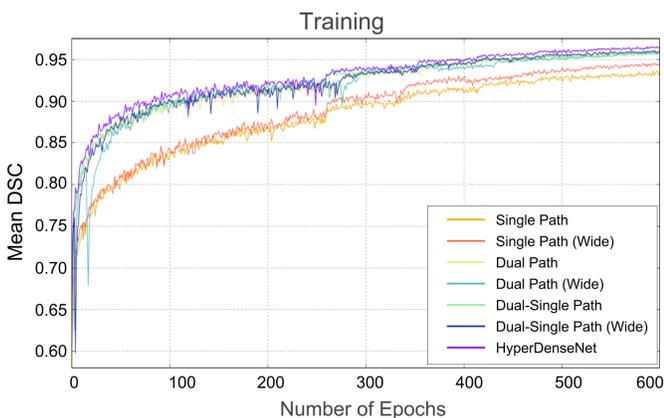


Fig. 4. Training accuracy plots for the proposed architecture and the baselines on the iSeg-2017 challenge data. The first point of each curve corresponds to the end of the first training epoch.

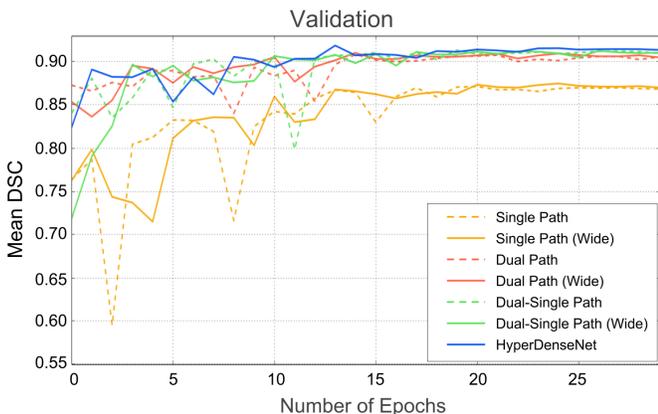


Fig. 5. Validation accuracy plots for the proposed architecture and the baselines on the iSeg-2017 challenge data. The first point of each curve corresponds to the end of the first training epoch.

The computational efficiency of HyperDenseNet and baselines is compared in Table 4. As expected, inference times are proportional to the number of model parameters. While the lightest architecture needs around 45 seconds to

segment a whole 3D brain, HyperDenseNet performs the same task in less than 2 minutes. This is acceptable from a clinical point of view.

Figure 6 depicts visual results for the subject used in validation. It can be seen that, in most cases, HyperDenseNet typically recovers thin regions better than the baselines, which can explain the improvements observed for distance-based metrics. As confirmed in Table 3, this effect is most prominent in the boundaries between the gray and white matter. Furthermore, HyperDenseNet produces fewer false positives for WM than the baselines, which tend to overestimate the segmentation in this region.

Challenge results: Table 5 compares the segmentation accuracy of HyperDenseNet to that of top-5 ranking methods in the first round of the iSEG Challenge, as well as to all the methods in the second round of submission. We observe that our network ranked among the top-3 methods in 6 out of 9 metrics, considering the results of the first and second rounds of submissions.

A noteworthy point is the general performance decrease of all the methods for the segmentation of GM and WM, with lower DSC and larger ASD values. This confirms that segmenting these tissues is more challenging due to the unclear boundaries between them.

3.2 MRBrainS Challenge

The MRBrainS challenge was initially proposed in conjunction with MICCAI 2013. It focuses on adult brain tissue segmentation in the context of aging, based on three modalities: MRI T1, MRI T1 Inversion Recovery (IR) and MR-FLAIR. To this day, a total of 47 international teams have participated in this challenge.

3.2.1 Evaluation

The organizers used three types of evaluation measures: a spatial overlap measure (DSC), a boundary distance measure (MHD) and a volumetric measure (the percentage of absolute volume difference).

3.2.2 Architectures for comparison

We compare HyperDenseNet to three state-of-the-art networks for medical image segmentation. The first architecture is a 3D fully convolutional neural network with residual connections [51], which we denote as *FCN_Res3D*. The second one, referred to as *UNet3D*, is a U-Net [52] model with residual connections in the encoder and 3D volumes as input. Finally, our comparison also includes *DeepMedic* [5], which showed an outstanding performance in brain lesion segmentation. The implementation details of these architectures are described in Supplemental materials (Supplementary materials are available in the supplementary files /multimedia tab).

3.2.3 Results

Validation results: We performed a leave-one-out-cross-validation (LOOCV) on the 5 available MRBrainS datasets, using 4 subjects for training and one for validation. We trained and tested models three times, each time using a different subject for validation, and measured the average accuracy over these three folds. For this experiment, we

TABLE 4

Number of parameters (convolution, fully-connected and total) and inference times of the baselines and the proposed architecture. Widened versions of the baselines, which we denoted using superscript *, are also included.

Architecture	Nb. of parameters			Time (sec)
	Conv.	Fully-conn.	Total	
Single Path	2,380,050	290,600	2,670,650	43.67 (± 8.37)
Single Path*	9,518,850	470,600	9,989,450	101.63 (± 12.65)
Dual Path	4,760,100	470,600	5,230,700	64.57 (± 9.45)
Dual Path*	9,381,960	614,600	9,996,560	104.31 (± 11.65)
Dual-Single Path	2,666,760	300,600	2,968,200	47.33 (± 8.74)
Dual-Single Path*	9,518,850	470,600	9,989,450	103.64 (± 13.61)
HyperDenseNet	9,518,850	830,600	10,349,450	105.67 (± 14.74)

* Widened version.

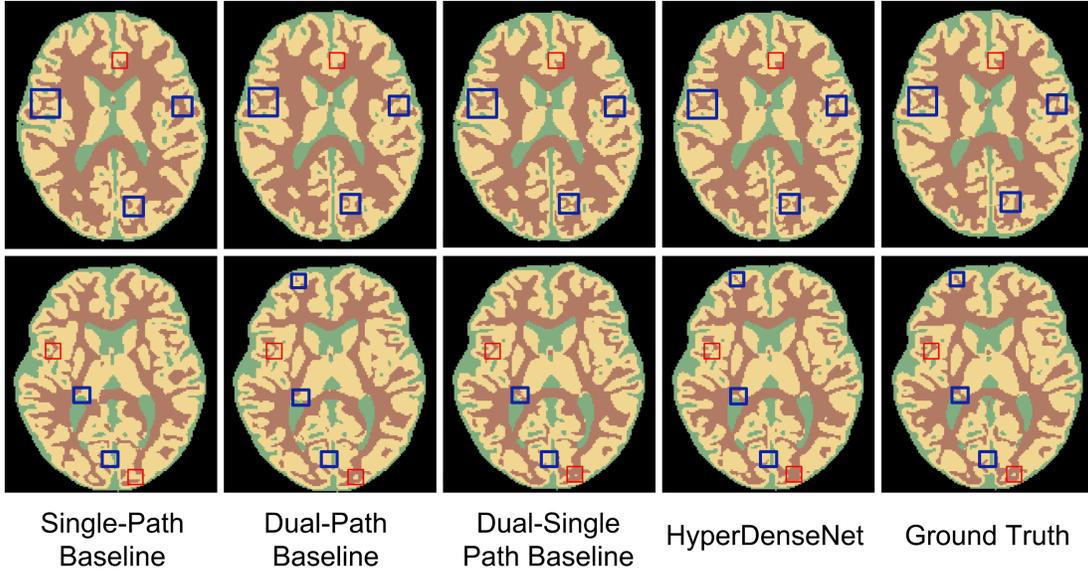


Fig. 6. Qualitative results of segmentation achieved by the baselines and HyperDenseNet on two validation subjects (each row shows a different subject). The green squares indicate some spots, where HyperDenseNet successfully reproduced the ground-truth whereas the baselines failed. Some regions where HyperDenseNet yielded incorrect segmentations are outlined in red.

TABLE 5

Results on the iSEG-2017 data for HyperDenseNet and the methods ranked in the top-5 at the first round of submissions (in alphabetical order). The bold fonts highlight the best performances. **Note:** The reported values were obtained from the challenge organizers at the time of submitting this manuscript, in February 2018. For an updated ranking, see the iSEG-2017 Challenge website for first (<http://iseg2017.web.unc.edu/rules/results/>) and second (<http://iseg2017.web.unc.edu/evaluation-on-the-second-round-submission/>) rounds of submission. The method referred to as LIVIA is a previous work from our team [17].

Method	CSF			GM			WM		
	DSC	MHD	ASD	DSC	MHD	ASD	DSC	MHD	ASD
<i>First round (Top 5)</i>									
Bern_IPMI	0.954	9.616	0.127	0.916	6.455	0.341	0.896	6.782	0.398
LIVIA (ensemble)	0.957	9.029	0.138	0.919	6.415	0.338	0.897	6.975	0.376
MSL_SKKU	0.958	9.072	0.116	0.919	5.980	0.330	0.901	6.444	0.391
nic_vicorob	0.951	9.178	0.137	0.910	7.647	0.367	0.885	7.154	0.430
TU/e IMAG/e	0.947	9.426	0.150	0.904	6.856	0.375	0.890	6.908	0.433
<i>Second round (All methods)</i>									
CatholicU	0.916	10.970	0.241	0.842	7.283	0.546	0.819	8.239	0.675
MSL_SKKU	0.958	9.112	0.116	0.923	5.999	0.321	0.904	6.618	0.375
BCH_CRL_IMAGINE	0.960	8.850	0.110	0.926	9.557	0.311	0.907	7.104	0.360
HyperDenseNet (Ours)	0.956	9.421	0.120	0.920	5.752	0.329	0.901	6.660	0.382

used all three modalities (i.e., T1, T1 IR and FLAIR) for all competing methods. In a second set of experiments,

we assessed the impact of integrating multiple imaging modalities on the performance of HyperDenseNet using all possible combinations of two modalities as input.

Table 6 reports the mean DSC and standard-deviation values of tested models, with *FCN_Res3D* exhibiting the lowest mean DSC. This performance might be explained by the transpose convolutions in *FCN_Res3D*, which may cause voxel misclassification within small regions. Furthermore, the downsampling and upsampling operations in *FCN_Res3D* could make the feature maps in hidden layers sparser than the original inputs, causing a loss of image details. A strategy to avoid this problem is having skip connections as in *UNet3D*, which propagate information at different levels of abstraction between the encoding and decoding paths. This can be observed in the results, where *UNet3D* clearly outperforms *FCN_Res3D* in all the metrics.

Moreover, *DeepMedic* obtained better results than its competitors, yielding a performance close to the different two-modality configurations of HyperDenseNet. The dual multiscale path is an important feature of *DeepMedic* which gives the network a larger receptive field via two paths, one for the input image and the other processing a low-resolution version of the input. This, in addition to the removal of pooling operations in *DeepMedic*, could explain the increase in performance with respect to *FCN_Res3D* and *UNet3D*.

Comparing the different modality combinations, the two-modality versions of HyperDenseNet yielded competitive performances, although there is a significant variability between the three configurations. Using only MRI T1 and FLAIR places HyperDenseNet first for two DSC measures (GM and WM), and second for the remaining measure (CSF), even though competing methods used all three modalities. However, HyperDenseNet with three modalities yields significantly better segmentations, with the highest mean DSC values for all three tissues.

Challenge results: The MRBrainS challenge organizers compiled the results and a ranking of 47 international teams⁴. In Table 7, we report the results of the top-10 methods. We see that HyperDenseNet ranks first among competing methods, obtaining the highest DSC and HD for GM and WM. Interestingly, the BCH_CRL_IMAGINE and MSL_SKKU teams participated in both iSEG and MRBrains2013 challenges. While these two networks outperformed HyperDenseNet in the iSEG challenge, the performance of our Model was noticeably superior in the MRBrains challenge, with HyperDenseNet ranked 1st, MSL_SKKU ranked 4th and BCH_CRL_IMAGINE ranked 18th (Ranking of February 2018). Considering the fact that three modalities are employed in MRBrains, unlike the two modalities used in iSEG, these results suggest that HyperDenseNet has stronger representation-learning power as the number of modalities increases.

A typical example of segmentation results is depicted in Fig. 7. In these images, red arrows indicate regions where the two-modality versions of HyperDenseNet fail in comparison to the three-modality version. As expected, most errors of these networks occur at the boundary between the

GM and WM (see images in Fig. 1, for example). Moreover, we observe that HyperDenseNet using three modalities can handle thin regions better than its two-modality versions.

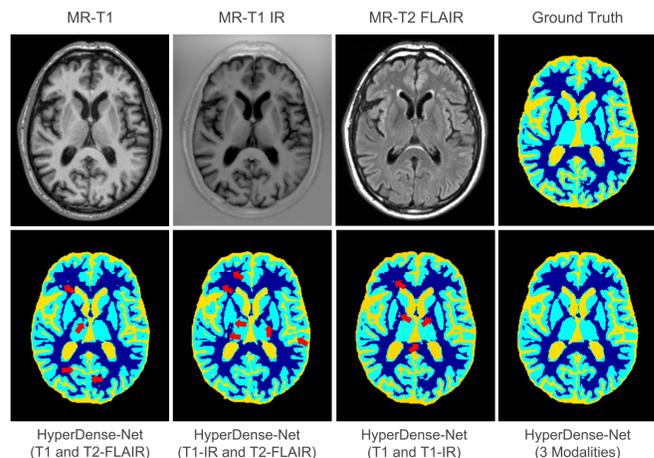


Fig. 7. A typical example of the segmentations achieved by the proposed HyperDenseNet in a validation subject (Subject 1 in the training set) for 2 and 3 modalities. The red arrows indicate some of the differences between the segmentations. For instance, one can see here that HyperDenseNet with three modalities can handle thin regions better than its two-modality versions.

3.3 Analysis of features re-use

Dense connectivity enables each network layer to access feature maps from all its preceding layers, strengthening feature propagation and encouraging feature re-use. To investigate the degree at which features are used in the trained network, we computed, for each convolutional layer, the average L_1 -norm of connection weights to previous layers in any stream. This serves as a surrogate for the dependency of a given layer on its preceding layers. We normalized the values between 0 and 1 to facilitate visualization.

Figure 8 depicts the weights of HyperDenseNet trained with two modalities, for both iSEG and MRBrainS challenges. As the MRBrainS dataset contains three modalities, we have three different two-modality configurations. The average weights for the case of three modalities are shown in Fig. 9. A dark square in these plots indicates that the target layer (on x -axis) makes a strong use of the features produced by the source layer (on y -axis). An important observation that one can make from both figures is that, in most cases, all layers spread the importance of the connections over many previous layers, not only within the same path, but also from the other streams. This indicates that shallow layer features are directly used by deeper layers from both paths, which confirms the usefulness of hyper-dense connections for information flow and learning complex relationships between modalities within different levels of abstractions.

Considering challenge datasets separately, for HyperDenseNet trained on iSEG (top row of Fig 8), immediate previous layers have typically higher impact on the connections from both paths. Furthermore, the connections having access to MRI T2 features typically have the strongest values, which may indicate that T2 is more discriminative than T1 in this particular situation. We can also observe some regions

4. <http://mrbrains13.isi.uu.nl/results.php>

TABLE 6
Comparison to several state-of-the-art 3D networks on the MRBrainS challenge.

Method	Mean DSC (std dev)		
	CSF	GM	WM
<i>FCN_Res3D</i> [53] (3-Modalities)	0.7685 (0.0161)	0.8163 (0.0222)	0.8607 (0.0178)
<i>UNet3D</i> [52] (3-Modalities)	0.8218 (0.0159)	0.8432 (0.0241)	0.8841 (0.0123)
DeepMedic [5] (3-Modalities)	0.8292 (0.0094)	0.8522 (0.0193)	0.8884 (0.0137)
HyperDenseNet (T1-FLAIR)	0.8259 (0.0133)	0.8620 (0.0260)	0.8982 (0.0138)
HyperDenseNet (T1_IR-FLAIR)	0.7991 (0.0181)	0.8226 (0.0255)	0.8654 (0.0087)
HyperDenseNet (T1-T1_IR)	0.8191 (0.0297)	0.8498 (0.0173)	0.8913 (0.0082)
HyperDenseNet (3-Modalities)	0.8485 (0.0078)	0.8663 (0.0247)	0.9016 (0.0109)

TABLE 7
Results of the MRBrainS challenge of different methods (DSC, HD (mm) and AVD). Only the top-10 methods are included in this table. **Note:** The reported values were obtained from the challenge organizers after submitting our results, in February 2018. For an updated ranking, see the MRBrainS Challenge website (<http://mrbrains13.isi.uu.nl/results.php>).

Method	GM			WM			CSF			Sum
	DSC	HD	AVD	DSC	HD	AVD	DSC	HD	AVD	
HyperDenseNet (ours)	0.8633	1.34	6.19	0.8946	1.78	6.03	0.8342	2.26	7.31	48
VoxResNet [19] + Auto-context	0.8615	1.44	6.60	0.8946	1.93	6.05	0.8425	2.19	7.69	54
VoxResNet [19]	0.8612	1.47	6.42	0.8939	1.93	5.84	0.8396	2.28	7.44	56
MSL-SKKU	0.8606	1.52	6.60	0.8900	2.11	5.54	0.8376	2.32	6.77	61
LRDE	0.8603	1.44	6.05	0.8929	1.86	5.83	0.8244	2.28	9.03	61
MDGRU	0.8540	1.54	6.09	0.8898	2.02	7.69	0.8413	2.17	7.44	80
PyraMiD-LSTM2	0.8489	1.67	6.35	0.8853	2.07	5.93	0.8305	2.30	7.17	83
3D-UNet [52]	0.8544	1.58	6.60	0.8886	1.95	6.47	0.8347	2.22	8.63	84
IDSIA [54]	0.8482	1.70	6.77	0.8833	2.08	7.05	0.8372	2.14	7.09	100
STH [55]	0.8477	1.71	6.02	0.8845	2.34	7.67	0.8277	2.31	6.73	112

with high (> 0.5) feature re-use patterns from shallow to deep layers. The same behaviour is seen for HyperDenseNet trained on two modalities from the MRBrainS challenge, where immediate previous layers have a high impact on the connections within and in-between the paths. The re-use of low-level features by deeper layers is more evident than in the previous case. For example, in HyperDenseNet trained with T1-IR and FLAIR, deep layers in the T1-IR path make a strong use of features extracted in shallower layers of the same path, as well as in the path corresponding to FLAIR. This strong re-use of early features from both paths occurred across all tested configurations. The same pattern is observed when using three modalities (Fig 9), with a strong re-use of shallow features from the network’s last layers. This reflects the importance of giving deep layers access to early-extracted features. Additionally, it suggests that learning how and where to fuse information from multiple sources is more effective than combining these sources in early or late stages.

4 CONCLUSION

This study investigated a hyper-densely connected 3D fully CNN, HyperDenseNet, with applications to brain tissue segmentation in multi-modal MRI. Our model leverages dense connectivity beyond recent works [39]–[41], exploiting the concept in multi-path architectures. Unlike these works, dense connections occur not only within the stream of individual modalities, but also across different streams.

This gives the network total freedom to explore complex combinations between features of different modalities, within and in-between all levels of abstraction. We reported a comprehensive evaluation using the benchmarks of two highly competitive challenges, iSEG-2017 for 6-month infant brain segmentation and MRBrainS for adult data, and showed state-of-the-art performances of HyperDenseNet on both datasets. Our experiments provided new insights on the inclusion of short-cut connections in deep neural networks for segmenting medical images, particularly in multi-modal scenarios. In summary, this work demonstrated the potential of HyperDenseNet to tackle challenging medical image segmentation problems involving multi-modal volumetric data.

ACKNOWLEDGMENTS

This work is supported by the National Science and Engineering Research Council of Canada (NSERC), discovery grant program, and by the ETS Research Chair on Artificial Intelligence in Medical Imaging. The authors would like to thank both iSEG and MRBrainS organizers for providing data benchmarks and evaluations.

REFERENCES

- [1] D. Delbeke, H. Schöder, W. H. Martin, and R. L. Wahl, “Hybrid imaging (SPECT/CT and PET/CT): improving therapeutic decisions,” in *Seminars in nuclear medicine*, vol. 39, no. 5. Elsevier, 2009, pp. 308–340.

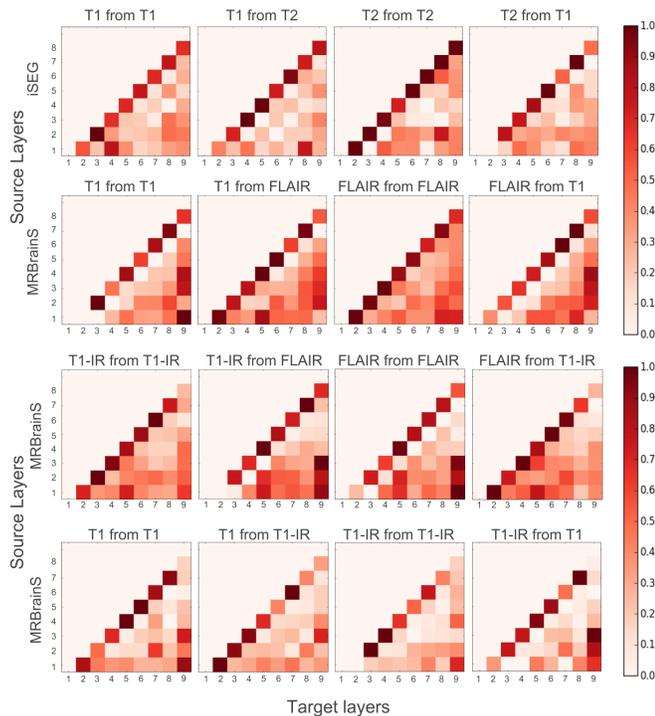


Fig. 8. Relative importance of connections in HyperDenseNet trained on the iSEG (top) and MRBrainS (from 2nd to 4th rows) challenges with two modalities. The color at each location encodes the average L1 norm of weights connecting a convolutional-layer source to a convolutional-layer target. These values were normalized between 0 and 1 by accounting for all the values within each layer.

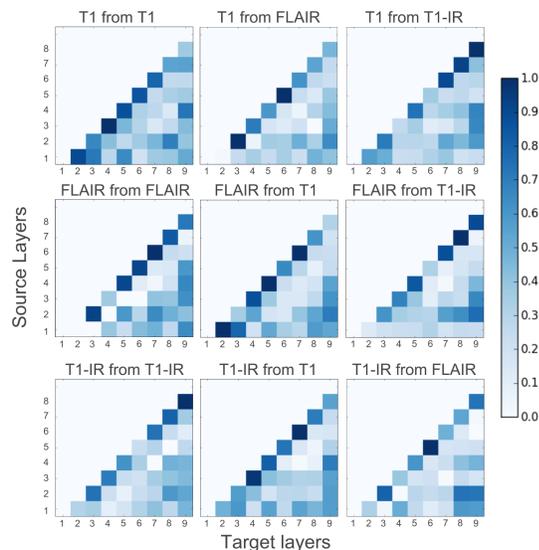


Fig. 9. Relative importance of connections in HyperDenseNet trained on the MRBrainS challenge with three modalities (MRI T1, FLAIR and T1 IR). The color at each location encodes the average L1 norm of weights connecting a convolutional-layer source to a convolutional-layer target. These values were normalized between 0 and 1 by accounting for all the values within each layer.

[2] X. Lladó, A. Oliver, M. Cabezas, J. Freixenet, J. C. Vilanova, A. Quiles, L. Valls, L. Ramió-Torrentà, and À. Rovira, "Segmentation of multiple sclerosis lesions in brain MRI: a review of automated approaches," *Information Sciences*, vol. 186, no. 1, pp. 164–185, 2012.

[3] B. H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, Y. Burren, N. Porz, J. Slotboom, R. Wiest *et al.*, "The mul-

timodal brain tumor image segmentation benchmark (BRATS)," *IEEE Transactions on Medical Imaging*, vol. 34, no. 10, pp. 1993–2024, 2015.

[4] M. Havaei, N. Guizard, N. Chapados, and Y. Bengio, "HeMIS: Hetero-modal image segmentation," in *International Conference on MICCAI*. Springer, 2016, pp. 469–477.

[5] K. Kamnitsas, C. Ledig, V. F. Newcombe, J. P. Simpson, A. D. Kane, D. K. Menon, D. Rueckert, and B. Glocker, "Efficient multi-scale 3D CNN with fully connected CRF for accurate brain lesion segmentation," *Medical image analysis*, vol. 36, pp. 61–78, 2017.

[6] L. Fidon, W. Li, L. C. Garcia-Peraza-Herrera, J. Ekanayake, N. Kitchen, S. Ourselin, and T. Vercauteren, "Scalable multimodal convolutional networks for brain tumour segmentation," in *International Conference on MICCAI*. Springer, 2017, pp. 285–293.

[7] M. Prastawa, J. H. Gilmore, W. Lin, and G. Gerig, "Automatic segmentation of MR images of the developing newborn brain," *Medical image analysis*, vol. 9, no. 5, pp. 457–466, 2005.

[8] N. I. Weisenfeld, A. Mewes, and S. K. Warfield, "Segmentation of newborn brain MRI," in *Biomedical Imaging: Nano to Macro, 2006. 3rd IEEE International Symposium on*. IEEE, 2006, pp. 766–769.

[9] P. Anbeek, K. L. Vincken, F. Groenendaal, A. Koeman, M. J. Van Osch, and J. Van der Grond, "Probabilistic brain tissue segmentation in neonatal magnetic resonance imaging," *Pediatric research*, vol. 63, no. 2, pp. 158–163, 2008.

[10] N. I. Weisenfeld and S. K. Warfield, "Automatic segmentation of newborn brain MRI," *Neuroimage*, vol. 47, no. 2, pp. 564–572, 2009.

[11] L. Wang, F. Shi, W. Lin, J. H. Gilmore, and D. Shen, "Automatic segmentation of neonatal images using convex optimization and coupled level sets," *NeuroImage*, vol. 58, no. 3, pp. 805–817, 2011.

[12] V. Srhoj-Egekher, M. Benders, K. J. Kersbergen, M. A. Viergever, and I. Isgum, "Automatic segmentation of neonatal brain MRI using atlas based segmentation and machine learning approach," *MICCAI Grand Challenge: Neonatal Brain Segmentation*, vol. 2012, 2012.

[13] S. Wang, M. Kuklisova-Murgasova, and J. A. Schnabel, "An atlas-based method for neonatal MR brain tissue segmentation," *Proceedings of the MICCAI Grand Challenge: Neonatal Brain Segmentation*, pp. 28–35, 2012.

[14] L. Wang, F. Shi, G. Li, Y. Gao, W. Lin, J. H. Gilmore, and D. Shen, "Segmentation of neonatal brain MR images using patch-driven level sets," *NeuroImage*, vol. 84, pp. 141–158, 2014.

[15] W. Zhang, R. Li, H. Deng, L. Wang, W. Lin, S. Ji, and D. Shen, "Deep convolutional neural networks for multi-modality isointense infant brain image segmentation," *NeuroImage*, vol. 108, pp. 214–224, 2015.

[16] D. Nie, L. Wang, Y. Gao, and D. Sken, "Fully convolutional networks for multi-modality isointense infant brain image segmentation," in *13th International Symposium on Biomedical Imaging (ISBI), 2016*. IEEE, 2016, pp. 1342–1345.

[17] J. Dolz, C. Desrosiers, L. Wang, J. Yuan, D. Shen, and I. Ben Ayed, "Deep CNN ensembles and suggestive annotations for infant brain MRI segmentation," *arXiv preprint arXiv:1712.05319*, 2017.

[18] A. M. Mendrik, K. L. Vincken, H. J. Kuijf, M. Breeuwer, W. H. Bouvy, J. De Bresser, A. Alansary, M. De Bruijne, A. Carass, A. El-Baz *et al.*, "MRBrainS challenge: online evaluation framework for brain image segmentation in 3T MRI scans," *Computational intelligence and neuroscience*, vol. 2015, p. 1, 2015.

[19] H. Chen, Q. Dou, L. Yu, J. Qin, and P.-A. Heng, "VoxResNet: Deep voxelwise residual networks for brain segmentation from 3D MR images," *NeuroImage*, 2017.

[20] S. C. Deoni, B. K. Rutt, A. G. Parrent, and T. M. Peters, "Segmentation of thalamic nuclei using a modified k-means clustering algorithm and high-resolution quantitative magnetic resonance imaging at 1.5T," *Neuroimage*, vol. 34, no. 1, pp. 117–126, 2007.

[21] O. Commowick, F. Cervenansky, and R. Ameli, "MSSEG Challenge proceedings: Multiple Sclerosis Lesions Segmentation Challenge using a data management and processing infrastructure," in *MICCAI*, 2016.

[22] K. Kamnitsas, C. Baumgartner, C. Ledig, V. Newcombe, J. Simpson, A. Kane, D. Menon, A. Nori, A. Criminisi, D. Rueckert *et al.*, "Unsupervised domain adaptation in brain lesion segmentation with adversarial networks," in *International Conference on IPMI*. Springer, 2017, pp. 597–609.

[23] S. Valverde, M. Cabezas, E. Roura, S. González-Villà, D. Pareto, J. C. Vilanova, L. Ramió-Torrentà, À. Rovira, A. Oliver, and X. Lladó, "Improving automated multiple sclerosis lesion seg-

- mentation with a cascaded 3D convolutional neural network approach," *NeuroImage*, vol. 155, pp. 159–168, 2017.
- [24] S. González-Villà, A. Oliver, S. Valverde, L. Wang, R. Zwigelaar, and X. Lladó, "A review on brain structures segmentation in magnetic resonance imaging," *Artificial intelligence in medicine*, vol. 73, pp. 45–69, 2016.
- [25] A. Makropoulos, S. J. Counsell, and D. Rueckert, "A review on automatic fetal and neonatal brain MRI segmentation," *NeuroImage*, 2017.
- [26] J. Dolz, C. Desrosiers, and I. Ben Ayed, "3D fully convolutional networks for subcortical segmentation in MRI: A large-scale study," *NeuroImage*, 2017.
- [27] T. Fechter, S. Adebahr, D. Baltas, I. Ben Ayed, C. Desrosiers, and J. Dolz, "Esophagus segmentation in CT via 3D fully convolutional neural network and random walk," *Medical Physics*, 2017.
- [28] K. Kamnitsas, L. Chen, C. Ledig, D. Rueckert, and B. Glocker, "Multi-scale 3D convolutional neural networks for lesion segmentation in brain MRI," *Ischemic Stroke Lesion Segmentation*, vol. 13, 2015.
- [29] P. Moeskops, M. A. Viergever, A. M. Mendrik, L. S. de Vries, M. J. Benders, and I. Išgum, "Automatic segmentation of MR brain images with a convolutional neural network," *IEEE Transactions on Medical Imaging*, vol. 35, no. 5, pp. 1252–1261, 2016.
- [30] N. Srivastava and R. Salakhutdinov, "Multimodal learning with deep boltzmann machines," *Journal of Machine Learning Research*, vol. 15, pp. 2949–2980, 2014.
- [31] L. Wang, F. Shi, Y. Gao, G. Li, J. H. Gilmore, W. Lin, and D. Shen, "Integration of sparse multi-modality representation and anatomical constraint for isointense infant brain MR image segmentation," *NeuroImage*, vol. 89, pp. 152–164, 2014.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on CVPR*, 2016, pp. 770–778.
- [33] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger, "Deep networks with stochastic depth," in *ECCV*. Springer, 2016, pp. 646–661.
- [34] G. Larsson, M. Maire, and G. Shakhnarovich, "Fractalnet: Ultra-deep neural networks without residuals," *arXiv preprint arXiv:1605.07648*, 2016.
- [35] S. Zagoruyko and N. Komodakis, "Wide residual networks," *arXiv preprint arXiv:1605.07146*, 2016.
- [36] R. Ranjan, V. M. Patel, and R. Chellappa, "Hyperface: A deep multi-task learning framework for face detection, landmark localization, pose estimation, and gender recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [37] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *AAAI*, 2017, pp. 4278–4284.
- [38] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten, "Densely connected convolutional networks," in *Proceedings of the IEEE CVPR*, 2017.
- [39] X. Li, H. Chen, X. Qi, Q. Dou, C.-W. Fu, and P. A. Heng, "H-DenseUNet: Hybrid densely connected UNet for liver and liver tumor segmentation from CT volumes," *arXiv:1709.07330*, 2017.
- [40] L. Yu, J.-Z. Cheng, Q. Dou, X. Yang, H. Chen, J. Qin, and P.-A. Heng, "Automatic 3D cardiovascular MR segmentation with densely-connected volumetric convnets," in *International Conference on MICCAI*. Springer, 2017, pp. 287–295.
- [41] L. Chen, Y. Wu, A. M. DSouza, A. Z. Abidin, A. Wismuller, and C. Xu, "MRI tumor segmentation with densely connected 3D CNN," *arXiv preprint arXiv:1802.02427*, 2018.
- [42] J. Dolz, I. Ben Ayed, J. Yuan, and C. Desrosiers, "Isointense infant brain segmentation with a Hyper-dense connected convolutional neural network," in *Biomedical Imaging (ISBI), 2018 IEEE 15th International Symposium on*. IEEE, 2018, pp. 616–620.
- [43] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE CVPR*, 2015, pp. 3431–3440.
- [44] T. Zhang, G.-J. Qi, B. Xiao, and J. Wang, "Interleaved group convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4373–4382.
- [45] Y. Chen, H. Wang, and Y. Long, "Regularization of convolutional neural networks using shufflenode," in *Multimedia and Expo (ICME), 2017 IEEE International Conference on*. IEEE, 2017, pp. 355–360.
- [46] X. Zhang, X. Zhou, M. Lin, and J. Sun, "Shufflenet: An extremely efficient convolutional neural network for mobile devices," *arXiv preprint arXiv:1707.01083*, 2017.
- [47] G. Alain and Y. Bengio, "Understanding intermediate layers using linear classifier probes," *arXiv preprint arXiv:1610.01644*, 2016.
- [48] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE ICCV*, 2015, pp. 1026–1034.
- [49] L. Wang, D. Nie, G. Li, . Puybureau, J. Dolz, Q. Zhang, F. Wang, J. Xia, Z. Wu, J. Chen, K. Thung, T. D. Bui, J. Shin, G. Zeng, G. Zheng, V. S. Fonov, A. Doyle, Y. Xu, P. Moeskops, J. P. W. Pluim, C. Desrosiers, I. B. Ayed, G. Sanroma, O. M. Benkarim, A. Casamitjana, V. Vilaplana, W. Lin, G. Li, and D. Shen, "Benchmark on automatic 6-month-old infant brain segmentation algorithms: The iseg-2017 challenge," *IEEE Transactions on Medical Imaging*, pp. 1–1, 2019.
- [50] L. R. Dice, "Measures of the amount of ecologic association between species," *Ecology*, vol. 26, no. 3, pp. 297–302, 1945.
- [51] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *ECCV*. Springer, 2016, pp. 630–645.
- [52] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3D U-Net: learning dense volumetric segmentation from sparse annotation," in *International Conference on MICCAI*. Springer, 2016, pp. 424–432.
- [53] N. Pawłowski, S. I. Ktena, M. C. Lee, B. Kainz, D. Rueckert, B. Glocker, and M. Rajchl, "DLTK: State of the art reference implementations for deep learning on medical images," *arXiv preprint arXiv:1711.06853*, 2017.
- [54] M. F. Stollenga, W. Byeon, M. Liwicki, and J. Schmidhuber, "Parallel multi-dimensional LSTM, with application to fast biomedical volumetric image segmentation," in *NIPS*, 2015, pp. 2998–3006.
- [55] A. Mahbod, M. Chowdhury, Ö. Smedby, and C. Wang, "Automatic brain segmentation using artificial neural networks with shape context," *Pattern Recognition Letters*, vol. 101, pp. 74–79, 2018.
- [56] Y. D. Reijmer, A. Leemans, M. Brundel, L. J. Kappelle, G. J. Biessels *et al.*, "Disruption of the cerebral white matter network is related to slowing of information processing speed in patients with type 2 diabetes," *Diabetes*, vol. 62, no. 6, pp. 2112–2115, 2013.
- [57] S. Klein, M. Staring, K. Murphy, M. A. Viergever, and J. P. Pluim, "Elastix: a toolbox for intensity-based medical image registration," *IEEE Transactions on Medical Imaging*, vol. 29, no. 1, pp. 196–205, 2010.
- [58] W. D. Penny, K. J. Friston, J. T. Ashburner, S. J. Kiebel, and T. E. Nichols, *Statistical parametric mapping: the analysis of functional brain images*. Academic press, 2011.

Supplemental Materials

Datasets

iSEG

Images were acquired at the UNC-Chapel Hill on a Siemens head-only 3T scanner with a circular polarized head coil, and were randomly chosen from the pilot study of the Baby Connectome Project (BCP)⁵. During scan, infants were asleep, unседated and fitted with ear protection, with the head secured in a vacuum-fixation device. T1-weighted images were acquired with 144 sagittal slices using the following parameters: TR/TE = 1900/4.38 ms, flip angle = 7° and resolution = 1×1×1 mm³. Likewise, T2-weighted images were obtained with 64 axial slices, TR/TE = 7380/119 ms, flip angle = 150° and resolution = 1.25×1.25×1.95 mm³. T2 images were linearly aligned onto their corresponding T1 images. All the images were resampled into an isotropic 1×1×1 mm³ resolution. Standard image pre-processing steps were then applied using in-house tools, including skull stripping, intensity inhomogeneity correction, and removal of the cerebellum and brain stem. For this application, 9 subjects were employed for training and 1 for validation. To obtain manual annotations, the organizers used 24-month follow-up scans to generate an initial automatic segmentation for 6-month subjects by employing a publicly available software iBEAT⁶. Then, based on the initial automatic contours, an experienced neuroradiologist corrected manually the segmentation errors (based on both T1 and T2 images) and geometric defects via ITK-SNAP, with surface rendering.

MRBrainS

20 subjects with a mean age of 71 ± 4 years (10 male, 10 female) were selected from an ongoing cohort study of older (65 – 80 years of age), functionally-independent individuals without a history of invalidating stroke or other brain diseases [56]. To test the robustness of the segmentation algorithms in the context of aging-related pathology, the subjects were selected to have varying degrees of atrophy and white-matter lesions, and the scans with major artifacts were excluded. The following sequences were acquired and used for the evaluation framework: 3D T1 (TR: 7.9 ms, TE: 4.5 ms), T1-IR (TR: 4416 ms, TE: 15 ms, and TI: 400 ms) and T2- FLAIR (TR: 11000 ms, TE: 125 ms, and TI: 2800 ms). The sequences were aligned by rigid registration using Elastix [57], along with a bias correction performed using SPM8 [58]. After the registration, the voxel size within all the provided sequences (T1, T1 IR, and T2 FLAIR) was 0.96×0.96×3.00 mm³. Five subjects that were representative for the overall data (2 male, 3 female and varying degrees of atrophy and white-matter lesions) were selected for training. The remaining fifteen subjects were provided as testing data. While ground truth was provided for the 5 training subjects, manual segmentations were unknown for the testing data set. The following structures were segmented and were available for training: (a) cortical gray matter, (b) basal ganglia, (c) white matter, (d) white matter

lesions, (e) peripheral cerebrospinal fluid, (f) lateral ventricles, (g) cerebellum and (h) brainstem. These structures can be merged into gray matter (a-b), white matter (c-d), and cerebrospinal fluid (e-f). The cerebellum and brainstem were excluded from the evaluation. Manual segmentations were drawn on the 3mm slice thickness scans by employing an in-house manual segmentation tool based on the contour segmentation objects tool in Mevislab⁷, starting with the inner most structures. While the outer border of the CSF was segmented using both T1 and T1 IR scans, the other regions were segmented on the T1 scan.

Performance metrics

Dice similarity coefficient (DSC)

Let V_{ref} and V_{auto} be, respectively, the reference and automatic segmentations of a given tissue class and for a given subject. The DSC for this subject is defined as

$$\text{DSC}(V_{\text{ref}}, V_{\text{auto}}) = \frac{2 | V_{\text{ref}} \cap V_{\text{auto}} |}{| V_{\text{ref}} | + | V_{\text{auto}} |} \quad (6)$$

DSC values are within a $[0, 1]$ range, 1 indicating perfect overlap and 0 corresponding to a total mismatch.

Average volume distance (AVD)

Using the same definitions for V_{auto} and V_{ref} , AVD corresponds to

$$\text{AVD}(V_{\text{ref}}, V_{\text{auto}}) = \frac{| V_{\text{ref}} - V_{\text{auto}} |}{V_{\text{ref}}} \cdot 100 \quad (7)$$

Modified Hausdorff distance (MHD)

Let P_{ref} and P_{auto} denote the sets of voxels within the reference and automatic segmentation boundary, respectively. MHD is given by

$$\text{MHD}(P_{\text{ref}}, P_{\text{auto}}) = \max \left\{ \max_{q \in P_{\text{ref}}} d(q, P_{\text{auto}}), \max_{q \in P_{\text{auto}}} d(q, P_{\text{ref}}) \right\}, \quad (8)$$

where $d(q, P)$ is the point-to-set distance defined by: $d(q, P) = \min_{p \in P} \|q - p\|$, with $\|\cdot\|$ denoting the Euclidean distance. Low MHD values indicate high boundary similarity.

Average surface distance (ASD)

Using the same notation as the Hausdorff distance above, the ASD corresponds to

$$\text{ASD}(P_{\text{ref}}, P_{\text{auto}}) = \frac{1}{|P_{\text{ref}}|} \sum_{p \in P_{\text{ref}}} d(p, P_{\text{auto}}), \quad (9)$$

where $|\cdot|$ denotes the cardinality of a set. In distance-based metrics, smaller values indicate higher proximity between two point sets and, thus, a better segmentation.

5. <http://babyconnectomeproject.org>

6. <http://www.nitrc.org/projects/ibeat/>

7. <https://www.mevislab.de/>

Implementation details

We extended our 3D FCNN architecture proposed in [26], which is based on Theano. The source code of this architecture is publicly available⁸. Training and testing was performed on a server equipped with a NVIDIA Tesla P100 GPU with 16 GB of RAM memory. Training HyperDenseNet took around 70 min per epoch, and around 35 hours in total for the two-modality version. With three image modalities, training each epoch took nearly 3 hours. Inference on a whole 3D MR scan took on average from 70-80 to 250-270 seconds, for the two- and three-modality versions, respectively.

The number of kernels per layer in each of the baselines and the proposed network are detailed in Table 8.

TABLE 8

Number of kernels (in convolutional and fully-connected layers) of the baselines and the proposed architecture. The architecture with two paths have the same number of kernels in both paths for the same convolutional block.

Architecture	Conv. kernels	Fully-conn. kernels
Single Path	[25,25,25,50,50,50,75,75,75]	[400,200,150]
Single Path*	[50,50,50,75,75,75,150,150,150]	[400,200,150]
Dual Path	[25,25,25,50,50,50,75,75,75]	[400,200,150]
Dual Path*	[40,40,40,70,70,70,100,100,100]	[400,200,150]
Dual-Single Path	[25,25,25,50,50,50,75,75,75]	[400,200,150]
Dual-Single Path*	[25,50,50,100,100,100,150,150,150]	[400,200,150]
HyperDenseNet	[25,25,25,50,50,50,75,75,75]	[400,200,150]

* Widened version.

FCN_Res3D

The architecture of *FCN_Res3D* consists on 5 convolutional blocks with residual units on the encoder path, with 16, 64, 128, 256 and 512 kernels. The decoding path contains 4 convolutional upsampling blocks, each composed of 4 kernels, one per class. At each residual block, batch normalization and a Leaky ReLU with a leakage value of 0.1 are employed before the convolution. Instead of including max-pooling operations to re-size the images, stride values of $2 \times 2 \times 2$ are used in layers 2, 3 and 4. Volume size at the input of the network is $64 \times 64 \times 24$. The implementation of this network is provided in [53]⁹.

UNet3D

Although quite similar to *FCN_Res3D*, *UNet3D* presents some differences, particularly in the decoding path. It contains 9 convolutional blocks in total, 4 in the encoding and 5 in the decoding path. The number of kernels in the encoding path are 32, 64, 128 and 256, with strides of $2 \times 2 \times 2$ at layers 2, 3 and 4. In the decoding path, the number of kernels are 256, 128, 64, 32 and 4, from the first to the last layer. Furthermore, skip connections are added at the convolutional blocks of the same scale between the encoding and decoding paths. As in *FCN_Res3D*, batch normalization and a Leaky ReLU with a leakage value of 0.1 are employed before the convolution at each block. Volume size at the input of the network is also $64 \times 64 \times 24$. The implementation is provided in [53].

DeepMedic

We used the default architecture of *DeepMedic* in our experiments. This architecture includes two paths with 8 convolutional blocks: 30, 30, 40, 40, 40, 40, 50, 50 kernels of size $3 \times 3 \times 3$. At the end of both paths, two fully connected convolutional layers with 150 $1 \times 1 \times 1$ filters each are added, before the last classification layer. The second path is used with a low-resolution version of the input at the first path, for a larger receptive field. The input patch size is $27 \times 27 \times 27$ and $35 \times 35 \times 35$ for training and segmentation, respectively. The official code¹⁰ is employed to evaluate this architecture.

8. <https://github.com/josedolz/SemiDenseNet>

9. <https://github.com/DLTK/DLTK>

10. <https://github.com/Kamnitsask/deepmedic>