

# Hierarchical brain parcellation with uncertainty

Mark S. Graham<sup>1</sup>, Carole H. Sudre<sup>1</sup>, Thomas Varsavsky<sup>1</sup>, Petru-Daniel Tudosiu<sup>1</sup>, Parashkev Nachev<sup>2</sup>, Sebastien Ourselin<sup>1</sup>, and M. Jorge Cardoso<sup>1</sup>

<sup>1</sup> Biomedical Engineering and Imaging Sciences, Kings College London, UK

<sup>2</sup> Institute of Neurology, University College London, London, United Kingdom  
`mark.graham@kcl.ac.uk`

**Abstract.** Many atlases used for brain parcellation are hierarchically organised, progressively dividing the brain into smaller sub-regions. However, state-of-the-art parcellation methods tend to ignore this structure and treat labels as if they are ‘flat’. We introduce a hierarchically-aware brain parcellation method that works by predicting the decisions at each branch in the label tree. We further show how this method can be used to model uncertainty separately for every branch in this label tree. Our method exceeds the performance of flat uncertainty methods, whilst also providing decomposed uncertainty estimates that enable us to obtain self-consistent parcellations and uncertainty maps at any level of the label hierarchy. We demonstrate a simple way these decision-specific uncertainty maps may be used to provide uncertainty-thresholded tissue maps at any level of the label tree.

## 1 Introduction

Brain parcellation seeks to partition the brain into spatially homogeneous structural and functional regions, a task fundamental for allowing us to study the brain in both function and dysfunction. The brain is hierarchically organised, with smaller subregions performing increasingly specialised functions, and the atlases classically used for parcellation typically reflect this by defining labels in a hierarchical tree structure. Manual parcellation is also typically performed hierarchically; typically semi-automated methods are used to help delineate larger structures with sufficient tissue contrast, and these are then manually sub-parcellated using anatomical or functional landmarks [1].

The state-of-the-art for brain parcellation has come to be dominated by convolutional neural networks (CNNs). These methods tend to ignore the label hierarchy, instead adopting a ‘flat’ label structure. However, methods that are aware of the label hierarchy are desirable for many reasons. Such methods could degrade their predictions gracefully, for example labelling a noisy region with the coarser label ‘cortex’ rather than trying to assign a particular cortical division. They also offer the opportunity to train on multiple datasets with differing degrees of label granularity, assuming those labels can be mapped onto a single hierarchy.

Hierarchical methods also enable uncertainty to be modelled at different levels of the label tree. There has been recent interest in using uncertainty estimates

provided by CNNs [8,9] to obtain confidence intervals for downstream biomarkers such as regional volumes [4,15], which is key if these biomarkers are to be integrated into clinical pipelines. Flat methods provide only a single uncertainty measure per voxel, which prevents attribution of the uncertainty to a specific decision. Hierarchical methods can provide uncertainty for each decision along the label hierarchy, for example enabling the network to distinguish between relatively easy decisions (e.g. cortex vs non-cortex) and more challenging decisions, such as delineating cortical sub-regions that are ill-defined on MRI. This could facilitate more specific and informative confidence bounds for derived biomarkers used in clinical decision making.

Whilst hierarchical methods have been applied to classification, [3,14,6,16], there are very few CNN-based methods that attempt hierarchical segmentation. A method proposed by Liang et al. [12] has been applied to perform hierarchical parcellation of the cerebellum [5]. A drawback of this approach is that the tree structure is directly built into the model architecture, requiring a tailored model to be built for each new label tree.

In this work we make two contributions. Firstly, we extend a method previously proposed for hierarchical classification [14] to hierarchically-aware segmentation. The method works by predicting decisions at each branch in the label tree, and has the advantage that it requires no alteration to the network architecture. Secondly, we show it is possible to use such a model to estimate uncertainty at each branch in the label tree. Our model with uncertainty matches the performance of ‘flat’ uncertainty methods, whilst providing us with decomposed uncertainty estimates that enable us to obtain consistent parcellations with corresponding uncertainty at any level of the label tree. We demonstrate how these decision-specific uncertainty maps can be used to provide uncertainty-thresholded tissue segmentations at any level of the label tree.

## 2 Methods

We first review existing flat segmentation models with uncertainty, before describing how we apply an existing classification model to perform hierarchical parcellation. We then show how such a model can be used to provide hierarchical uncertainty estimates. We focus on modelling intrinsic uncertainty in this work, although the methods presented can be straightforwardly extended to estimating model uncertainty, too.

### 2.1 Flat parcellation

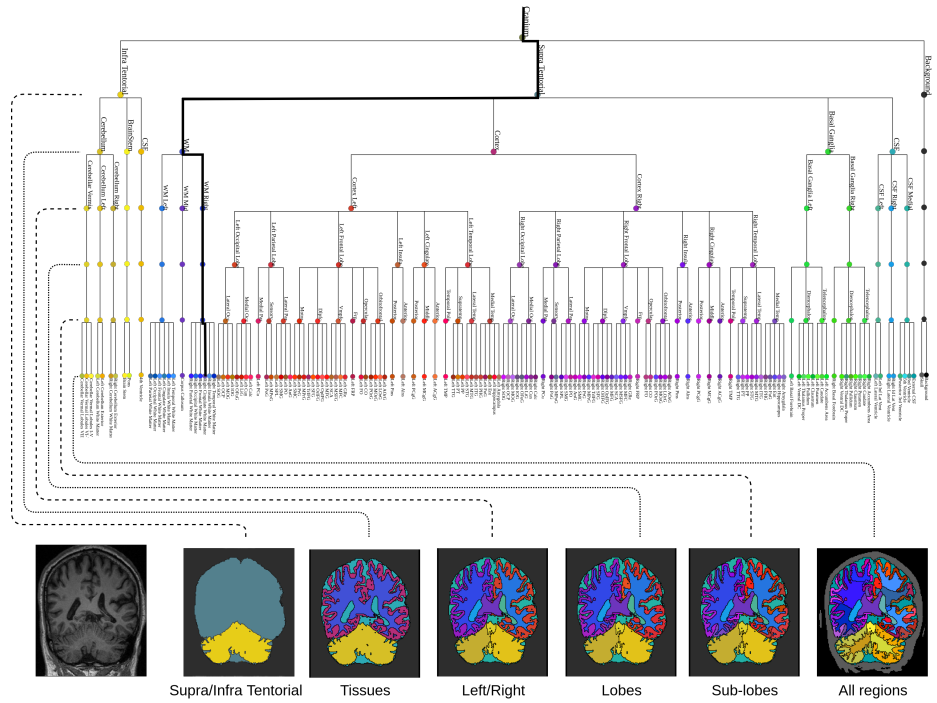
In a flat segmentation scenario, we consider the task as per-voxel classification, where the likelihood for a voxel is given by  $p(\mathbf{y}|\mathbf{W}, \mathbf{x}) = \text{Softmax}(\mathbf{f}^{\mathbf{W}}(\mathbf{x}))$  where  $\mathbf{f}^{\mathbf{W}}(\mathbf{x})$  is the output of a neural network with weights  $\mathbf{W}$ , input  $\mathbf{x}$  is a 3D image volume, and  $\mathbf{y}$  encodes the  $C$  segmentation classes. We seek the weights  $\mathbf{W}$  that minimise the negative log-likelihood, yielding the standard cross-entropy loss function,  $\text{CE}(y = c, \mathbf{f}^{\mathbf{W}}(\mathbf{x})) = -\log \text{Softmax}(f_c^{\mathbf{W}}(\mathbf{x}))$ . As in Kendall et al. [8],

heteroscedastic intrinsic uncertainty can be modelled by considering scaling the logits by a second network output,  $\sigma^{\mathbf{W}}(\mathbf{x})$ , giving a likelihood of  $p(y|\mathbf{W}, \mathbf{x}, \sigma) = \text{Softmax}\left(\frac{1}{\sigma^2(\mathbf{x})}\mathbf{f}^{\mathbf{W}}(\mathbf{x})\right)$ .  $\sigma^{\mathbf{W}}(\mathbf{x})$  is a per-voxel estimate, so it has the same dimension as  $\mathbf{x}$ . Employing the approximation  $\frac{1}{\sigma^{\mathbf{W}}(\mathbf{x})^2} \sum_c \exp\left(\frac{1}{\sigma^{\mathbf{W}}(\mathbf{x})^2} f_c^{\mathbf{W}}(\mathbf{x})\right) \approx \left(\sum_c \exp(f_c^{\mathbf{W}}(\mathbf{x}))\right) \sigma^{\mathbf{W}}(\mathbf{x})^{-2}$  used in [9] allows us to write the negative log-likelihood as

$$\mathcal{L}(y = c, \mathbf{x}; \mathbf{W}) = \frac{\text{CE}(y = c, \mathbf{f}^{\mathbf{W}}(\mathbf{x}))}{\sigma^{\mathbf{W}}(\mathbf{x})^2} + \log \sigma^{\mathbf{W}}(\mathbf{x})$$

### 2.2 Hierarchical parcellation

Here we describe the hierarchical classification/detection model proposed by Redmon et al. [14], and discuss how it can be adapted for segmentation tasks. The methods described here are general to all label taxonomy trees, but in this work we specifically consider the tree shown in Figure 1, described in more detail in Section 3.1. The probabilities at each node obey simple rules: the probabilities



**Fig. 1.** The neuro-anatomical label hierarchy considered in this paper, with the path from the root to the right cingulate highlighted. A larger version of this tree is included in the supplementary materials.

of all a node’s children sum to the probability of the node itself, and so if we take  $p(\text{root}) = 1$  the probabilities of all leaf nodes sum to 1. Leaf node probabilities can be expressed as the product of conditional probabilities down the tree; for example using the hierarchy in Figure 1 we can express  $p(\text{Right cingulate WM})$  as

$$\begin{aligned} p(\text{Right cingulate WM}) = & p(\text{Right cingulate}|\text{Right WM})p(\text{Right WM}|\text{WM}) \dots \\ & p(\text{WM}|\text{Supra tentorial})p(\text{Supra tentorial}|\text{Cranium}) \dots \\ & p(\text{Cranium}) \end{aligned}$$

where  $p(\text{Cranium}) = 1$ . Our model predicts the conditional probabilities for each node, and is optimised using a cross-entropy loss at every level of the tree.

More formally, we label each node  $i$  at level  $l$  as  $N_{i,l}$ , where  $l = 0$  denotes the root and  $l = L$  the deepest level, giving a maximum height of  $L + 1$ . Our model  $\mathbf{f}^{\mathbf{W}}(\mathbf{x})$  produces a score for each node in the tree,  $f^{\mathbf{W}}(\mathbf{x})_{i,l}$ . We define a hierarchical softmax - essentially a softmax over the siblings for a given node - to produce the conditional probabilities at each node,

$$p_{i,l} = \frac{\exp(f^{\mathbf{W}}(\mathbf{x})_{i,l})}{\sum_{N_{j,l} \in S[N_{i,l}]} \exp(f^{\mathbf{W}}(\mathbf{x})_{j,l})}$$

where  $S[N_{i,l}]$  denotes all the sibling nodes of  $N_{i,l}$ , including itself.

In the flat case we had a single label per voxel,  $y_c$ . In the hierarchical case  $y_c$  denotes a leaf node of the tree, and we consider the label superset  $A[y_c] = \{N_{i,l}\}$  comprising all the nodes traversed from the root to the label’s leaf node, excluding the root node but including itself. The total loss is the summation of a CE loss calculated at each level of the tree,

$$\mathcal{L}(y = c, \mathbf{x}; \mathbf{W}) = - \sum_{N_{i,l} \in A[y_c]} \log p_{i,l}$$

For parcellation the network makes a prediction per voxel, that is  $\mathbf{f}^{\mathbf{W}}(\mathbf{x}) \in \mathbb{R}^{x \times y \times z \times H}$  where  $H$  is the total number of nodes, making the considerably more computationally expensive than in classification tasks. The denominator of the hierarchical softmax can be efficiently calculated as a matrix multiplication, allowing  $p_{i,l}$  to be calculated from the elementwise division of two matrices.

### 2.3 Hierarchical uncertainty

We extend the model by modelling an uncertainty for every decision made along the tree. The network output  $\sigma^{\mathbf{W}}(\mathbf{x})$  is now vector-valued, and exists for every non-leaf node,  $\sigma^{\mathbf{W}}(\mathbf{x})_{i,l}$ . The loss becomes:

$$\mathcal{L}(y = c, \mathbf{x}; \mathbf{W}) = - \sum_{N_{i,l} \in A[y_c]} \frac{\log p_{i,l}}{\sigma^{\mathbf{W}}(\mathbf{x})_{i,l-1}^2} + \log \sigma^{\mathbf{W}}(\mathbf{x})_{i,l-1}$$

In this formulation the uncertainty values in a given voxel are unconstrained if they do not fall along the decision path for that voxel; for example values of  $\sigma$

relating to cortical parcellation do not enter into the loss in white matter voxels. We add a penalty term to encourage shrinking every value of  $\sigma_{i,l}$  that does not fall along the path from the true leaf node to the root node, giving a final loss of

$$\begin{aligned} \mathcal{L}(\mathbf{y} = c, \mathbf{x}; \mathbf{W}) = & - \sum_{N_{i,l} \in A[y_c]} \left( \frac{\log p_{i,l}}{\sigma^{\mathbf{W}}(\mathbf{x})_{i,l-1}^2} + \log \sigma^{\mathbf{W}}(\mathbf{x})_{i,l-1} \right) \\ & + \lambda \sum_{N_{i,l} \notin A[y_c]} \log \sigma^{\mathbf{W}}(\mathbf{x})_{i,l-1} \end{aligned} \quad (1)$$

where  $\lambda$  controls the strength of this penalty.

## 2.4 Architecture and implementation details

The network is a 3D UNet based on the implementation described in the nnUNet paper [7] and implemented in PyTorch. Our implementation contains three pooling layers and separate, identical decoder branches for the segmentation and uncertainty outputs. The parcellation branch predicts an output for each leaf node in the tree for the flat case - 151 for the tree considered in this work - and in the hierarchical case predicts an output for each node in the tree. As the hierarchical network does not make any predictions for nodes with no siblings, as  $p(\text{node}|\text{parent})=1$  always for such nodes, the hierarchical model predicts 213 outputs per voxel for the same tree. The uncertainty branch predicts a single channel for flat models, and a number of channels equal to the number of branches in the label tree for hierarchical models - 61 for the tree in this work. In practice,  $\log(\sigma^2)$  is predicted for numerical stability. We set the penalty term in the hierarchical loss  $\lambda = 0.1$ . Networks were trained on  $110^3$  patches randomly sampled from the training volume. Group normalisation was used, enabling a batch size of 1 to be coupled with gradient accumulation to produce an effective batch size of 3. Models were trained with the Adam optimiser [10] using a learning rate of  $4e^{-3}$ . Each model was trained for a maximum of 300 epochs with early stopping if the minimum validation loss did not improve for 15 epochs.

## 3 Experiments and Results

### 3.1 Data

We use the hierarchical label tree from the GIF label-fusion framework [2], which is based on the labelling from the MICCAI 2012 Grand Challenge on label fusion [11]. In total, there are 151 leaf classes and a hierarchical depth of 6, see Figure 1.

We use 593 T1-weighted MRI scans from the ADNI2 dataset [13], with an average voxel size of  $1.18 \times 1.05 \times 1.05 \text{mm}^3$  and dimension  $182 \times 244 \times 246$ . Images were bias-field corrected, oriented to a standard RAS orientation and cropped using a tight mask. Silver-standard labels were produced using GIF on multimodal input data, followed by manual quality control and editing where necessary. 543 scans were used for training and validation, and 50 were reserved for testing.

**Table 1.** Dice scores averaged over all classes on the test set for the flat ( $F$ ) and the proposed hierarchical ( $H$ ) model. Uncertainty-aware models are denoted with an unc subscript. Values are Median (IQR) across the 50 subjects in the test set. Bold indicates significantly better performance between model pairs ( $F$  vs  $H$ ,  $F_{\text{unc}}$  vs  $H_{\text{unc}}$ ), at  $p < 0.05$ , using p-values obtained from a Wilcoxon paired test.

Tree level	$F$	$H$ (ours)	$P$	$F_{\text{unc}}$	$H_{\text{unc}}$ (ours)	$P$
Supra/Infra	<b>0.986 (0.003)</b>	0.985 (0.002)	<0.00005	<b>0.984 (0.003)</b>	0.984 (0.002)	0.009
Tissue	<b>0.942 (0.007)</b>	0.941 (0.008)	<0.00005	0.934 (0.007)	0.934 (0.007)	0.95
Left/right	<b>0.942 (0.006)</b>	0.938 (0.008)	<0.00005	0.932 (0.005)	<b>0.933 (0.006)</b>	0.006
Lobes	<b>0.924 (0.008)</b>	0.922 (0.009)	0.00001	0.913 (0.008)	<b>0.917 (0.008)</b>	<0.00005
Sub-lobes	<b>0.891 (0.011)</b>	0.884 (0.013)	<0.00005	0.870 (0.013)	<b>0.880 (0.012)</b>	<0.00005
All regions	<b>0.861 (0.011)</b>	0.848 (0.015)	<0.00005	0.831 (0.018)	<b>0.845 (0.013)</b>	<0.00005

### 3.2 Experiments

We consider the following four models: 1) a baseline network trained on flat labels with weighted cross-entropy ( $F$ ) 2) the same as ( $F$ ) but with uncertainty estimates ( $F_{\text{unc}}$ ), 3) a network trained on hierarchical labels ( $H$ ), 4) a hierarchically-trained network with hierarchical uncertainty estimates ( $H_{\text{unc}}$ ). The following experiments were performed:

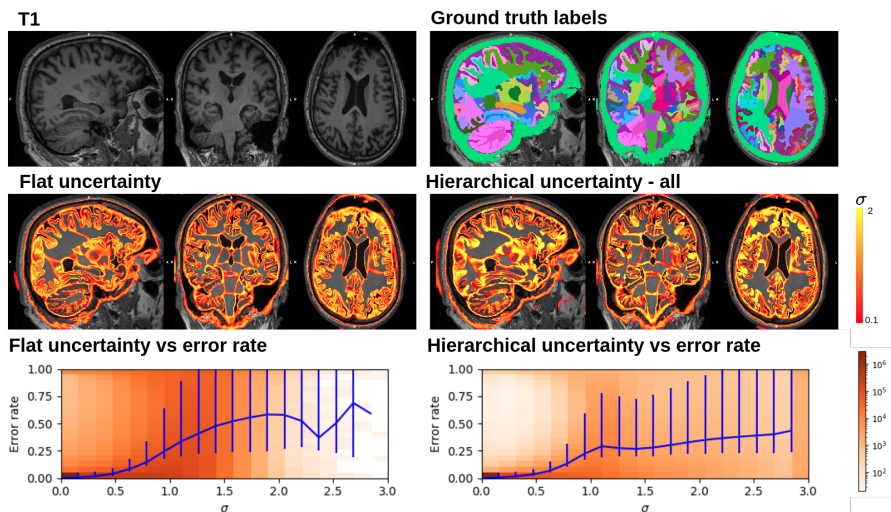
- Performance comparison using dice overlap on the withheld test data at all six levels of the tree.
- Qualitative assessment of the uncertainty maps provided by  $H_{\text{unc}}$  and  $F_{\text{unc}}$ .
- Comparison of uncertainty-thresholded segmentations from  $H_{\text{unc}}$  and  $F_{\text{unc}}$ .

### 3.3 Results & Discussion

Dice scores for all the models are reported in Table 1. Despite predicting a tree-structure with > 41% more predictions per voxel than the flat model, performance for  $H$  only drops marginally when compared to  $F$ , consistent with existing performance comparisons between flat and hierarchical models in classification and object detection settings [14].  $H_{\text{unc}}$  outperforms  $F_{\text{unc}}$  for the four more fine-grained levels of the label tree. This is likely due to the empirically observed difficult in stably training  $F_{\text{unc}}$ ; we found no such problems with  $H_{\text{unc}}$ , which was easy to optimise.

Figure 2 compares the uncertainty map from  $F_{\text{unc}}$  with the total uncertainty map from  $H_{\text{unc}}$ , obtained by summing all uncertainty components at each voxel. They look visually similar, and the joint histograms demonstrate expected trade-offs between uncertainty and error rate. Ideally, we would see low counts in the top-left of the joint histograms, indicating the models do not make confidently wrong predictions with low uncertainty. We see this desired behavior for  $H_{\text{unc}}$  more strongly than  $F_{\text{unc}}$ .

Figure S1 shows uncertainty maps predicted by  $H_{\text{unc}}$  for different branches of the label tree. The model provides sensibly decomposed uncertainty maps



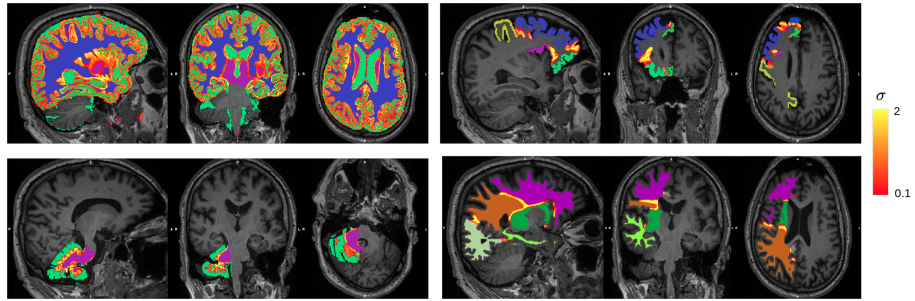
**Fig. 2.** Evaluation of the uncertainty from  $F_{\text{unc}}$  and the total uncertainty for  $H_{\text{unc}}$  obtained by summing all uncertainty components. Joint histograms show voxel counts for  $\sigma$  against  $(1 - \text{predicted probability for true class})$ , averaged across all test subjects. Blue lines represent the mean error rate and error bars are 25-75 percentiles.

for each decision along the label tree, with uncertainty strongly localised along decision boundaries. The maps reflect the uncertainty we expect for different decisions: for example there is highly localised uncertainty along the well contrasted WM-CSF boundary, but uncertainty is more spread out on boundaries between cortical regions which are poorly defined, and subject to high inter-rater variability.

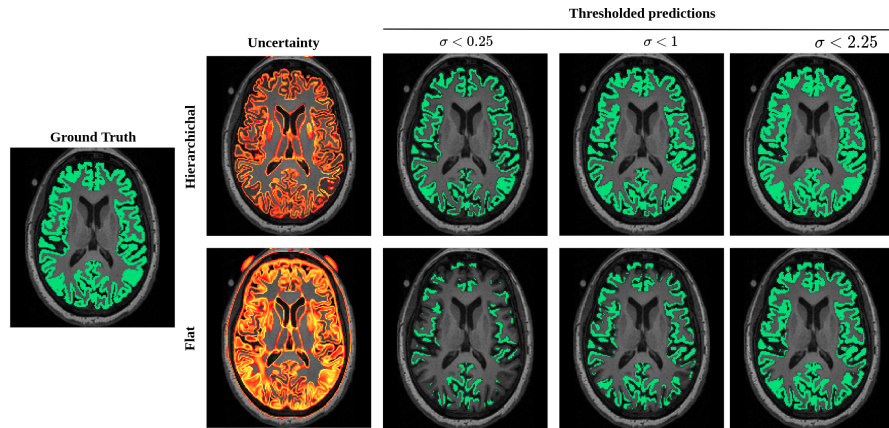
Figure 4 demonstrates a simple uncertainty-based thresholding method to obtain upper- and lower-bound cortical maps. They show that the cortical-specific uncertainty component from  $H_{\text{unc}}$  can be used to sensibly threshold predictions for non-leaf classes, in a way that is not possible for the uncertainty map from  $F_{\text{unc}}$  which lacks specificity to non-leaf nodes.

## 4 Conclusions

We have proposed a hierarchically-aware parcellation model, and demonstrated how it may be used to produce per-decision measures of uncertainty on the label tree. Our method outperforms the flat uncertainty model in terms of dice score, and was less likely than the flat model to make wrong predictions with both high confidence and low uncertainty. Furthermore we demonstrate the decomposed uncertainty enables us to produce consistent parcellations along with uncertainty maps for classes higher up the label tree, which is not possible with flat uncertainty models.



**Fig. 3.** Demonstration of different uncertainty components for model  $H_{\text{unc}}$  at four different branches of the label hierarchy, shown alongside the tissue class options at that branch. Colours have been selected to maximise distinguishability between adjacent classes.



**Fig. 4.** Demonstration of thresholding predictions according to uncertainty. Ground truth cortical segmentation is shown on left. Using  $H_{\text{unc}}$  a cortex-specific uncertainty map can be produced, that can be sensibly thresholded to create cortical predictions at different uncertainty levels. The lack of decision specificity in the single uncertainty map provided by  $F_{\text{unc}}$  means we cannot perform cortex-specific thresholding - see in particular the map thresholded at  $\sigma < 0.25$ .



## References

1. Whole-brain segmentation protocol, <http://neuromorphometrics.com/Seg/>
2. Cardoso, M.J., Modat, M., Wolz, R., Melbourne, A., Cash, D., Rueckert, D., Ourselin, S.: Geodesic information flows: spatially-variant graphs and their application to segmentation and fusion. *IEEE transactions on medical imaging* **34**(9), 1976–1988 (2015)
3. Demyanov, S., Chakravorty, R., Ge, Z., Bozorgtabar, S., Pablo, M., Bowling, A., Garnavi, R.: Tree-loss function for training neural networks on weakly-labelled datasets. In: 2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017). pp. 287–291. IEEE (2017)
4. Eaton-Rosen, Z., Bragman, F., Bisdas, S., Ourselin, S., Cardoso, M.J.: Towards safe deep learning: accurately quantifying biomarker uncertainty in neural network predictions. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 691–699. Springer (2018)
5. Han, S., Carass, A., Prince, J.L.: Hierarchical parcellation of the cerebellum. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 484–491. Springer (2019)
6. Hu, H., Zhou, G.T., Deng, Z., Liao, Z., Mori, G.: Learning structured inference neural networks with label relations. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2960–2968 (2016)
7. Isensee, F., Petersen, J., Kohl, S.A.A., Jäger, P.F., Maier-Hein, K.H.: nnu-net: Breaking the spell on successful medical image segmentation (2019)
8. Kendall, A., Gal, Y.: What uncertainties do we need in bayesian deep learning for computer vision? In: Advances in neural information processing systems. pp. 5574–5584 (2017)
9. Kendall, A., Gal, Y., Cipolla, R.: Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7482–7491 (2018)
10. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
11. Landman, B., Warfield, S.: Miccai 2012 workshop on multi-atlas labeling. In: Medical image computing and computer assisted intervention conference (2012)
12. Liang, X., Zhou, H., Xing, E.: Dynamic-structured semantic propagation network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 752–761 (2018)
13. Petersen, R.C., Aisen, P., Beckett, L.A., Donohue, M., Gamst, A., Harvey, D.J., Jack, C., Jagust, W., Shaw, L., Toga, A., et al.: Alzheimer’s disease neuroimaging initiative (adni): clinical characterization. *Neurology* **74**(3), 201–209 (2010)
14. Redmon, J., Farhadi, A.: Yolo9000: better, faster, stronger. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7263–7271 (2017)
15. Wang, G., Li, W., Vercauteren, T., Ourselin, S.: Automatic brain tumor segmentation based on cascaded convolutional neural networks with uncertainty estimation. *Frontiers in computational neuroscience* **13**, 56 (2019)
16. Wu, C., Tygert, M., LeCun, Y.: A hierarchical loss and its problems when classifying non-hierarchically. *PloS one* **14**(12) (2019)

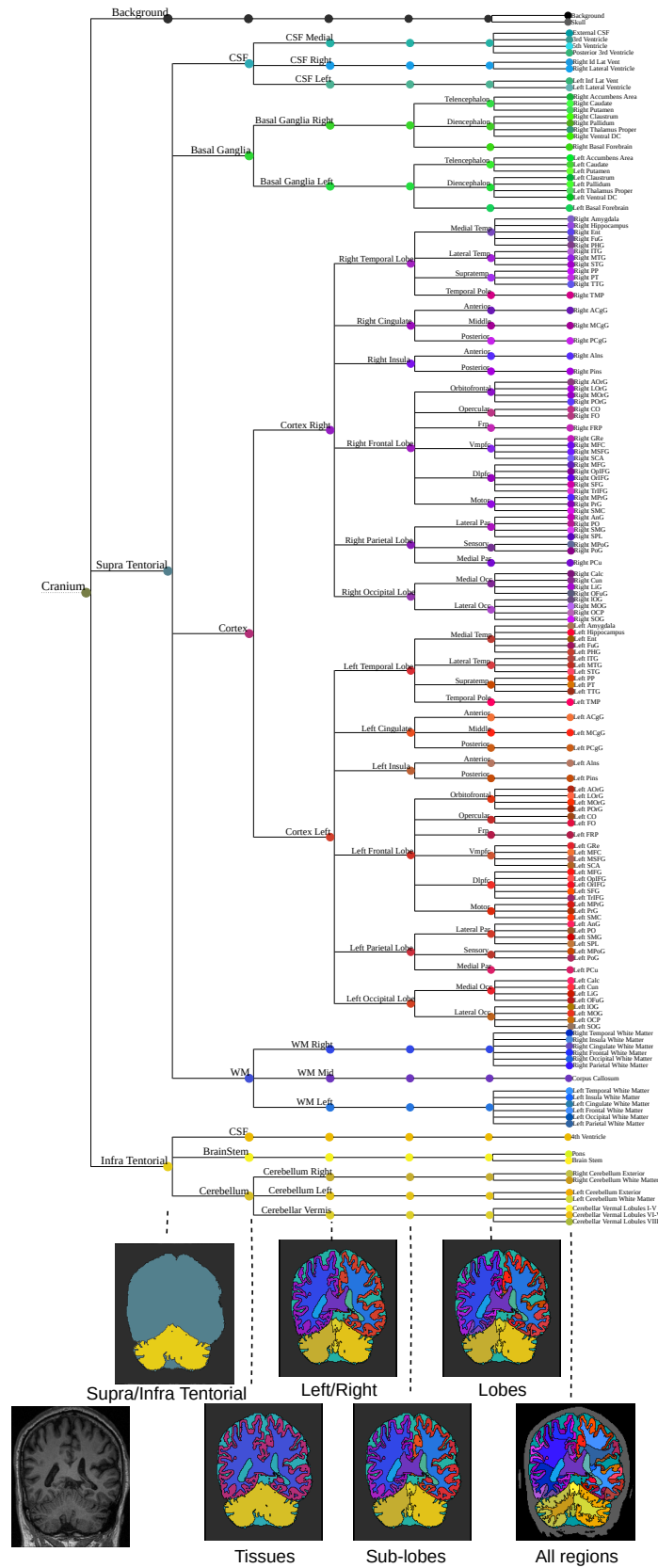


Fig. S1. Larger version of the label hierarchy considered in this paper.