

Compositional Representation Learning for Brain Tumour Segmentation

Xiao Liu^{1,2}, Antanas Kascenas¹, Hannah Watson¹, Sotirios A. Tsaftaris^{1,2,3}
and Alison Q. O’Neil^{1,2}

¹ Canon Medical Research Europe Ltd., Edinburgh, UK

² School of Engineering, University of Edinburgh, Edinburgh EH9 3FB, UK

³ The Alan Turing Institute, London, UK

`xiao.liu@mre.medical.canon`

Abstract. For brain tumour segmentation, deep learning models can achieve human expert-level performance given a large amount of data and pixel-level annotations. However, the expensive exercise of obtaining pixel-level annotations for large amounts of data is not always feasible, and performance is often heavily reduced in a low-annotated data regime. To tackle this challenge, we adapt a mixed supervision framework, vMFNet, to learn robust compositional representations using unsupervised learning and weak supervision alongside non-exhaustive pixel-level pathology labels. In particular, we use the BraTS dataset to simulate a collection of 2-point expert pathology annotations indicating the top and bottom slice of the tumour (or tumour sub-regions: peritumoural edema, GD-enhancing tumour, and the necrotic / non-enhancing tumour) in each MRI volume, from which weak image-level labels that indicate the presence or absence of the tumour (or the tumour sub-regions) in the image are constructed. Then, vMFNet models the encoded image features with von-Mises-Fisher (vMF) distributions, via learnable and compositional vMF kernels which capture information about structures in the images. We show that good tumour segmentation performance can be achieved with a large amount of weakly labelled data but only a small amount of fully-annotated data. Interestingly, emergent learning of anatomical structures occurs in the compositional representation even given only supervision relating to pathology (tumour).

Keywords: Compositionality · Representation learning · Semi-supervised · Weakly-supervised · Brain tumour segmentation.

1 Introduction

When a large amount of labelled training data is available, deep learning techniques have demonstrated remarkable accuracy in medical image segmentation [2]. However, performance drops significantly when insufficient pixel-level annotations are available [16, 17, 24]. By contrast, radiologists learn clinically relevant visual features from “weak” image-level supervision of seeing many medical scans [1]. When searching for anatomy or lesions of interest in new images, they

look for characteristic configurations of these clinically relevant features (or components). A similar compositional learning process has been shown to improve deep learning model performance in many computer vision tasks [9, 11, 25] but has received limited attention in medical applications.

In this paper, we consider a limited annotation data regime where few pixel-level annotations are available for the task of brain tumour segmentation in brain MRI scans. Alongside this, we construct slice-level labels for each MRI volume indicating the presence or absence of the tumour. These labels can be constructed from 2-point expert pathology annotations indicating the top and bottom slices of the tumour, which are fast to collect. We consider that pathology annotations are not only better suited to the task (tumour segmentation) but also to the domain (brain MRI) than the originally proposed weak supervision with anatomy annotations [15]; annotating the top and bottom slices for anatomical brain structures such as white matter, grey matter and cerebrospinal fluid (CSF) would be relatively uninformative about the configurations of structures within the image due to their whole brain distributions.

For the learning paradigm, we investigate the utility of learning compositional representations in increasing the annotation efficiency of segmentation model training. Compositional frameworks encourage identification of the visible semantic components (e.g. anatomical structures) in an image, requiring less explicit supervision (labels). We follow [11, 15, 18] in modelling compositional representations of medical imaging structures with learnable von-Mises-Fisher (vMF) kernels. The vMF kernels are learned as the cluster centres of the feature vectors of the training images, and the vMF activations determine which kernel is activated at each position. On visualising kernel activations, it can be seen that they approximately correspond to human-recognisable structures in the image, lending interpretability to the model predictions. Our contributions are summarised as:

- We refine an existing mixed supervision compositional representation learning framework, vMFNet, for the task of brain tumour segmentation, changing the weak supervision task from anatomy presence/absence to more domain-suited pathology presence/absence and simplifying the architecture and training parameters according to the principle of parsimony (in particular reducing the number of compositional vMF kernels and removing an original training subtask of image reconstruction).
- We perform extensive experiments on the BraTS 2021 challenge dataset [4, 5, 19] with different percentages of labelled data, showing superior performance of the proposed method compared to several strong baselines, both quantitatively (better segmentation performance) and qualitatively (better compositional representations).
- We compare weak pathology supervision with *tumour* labels to richer tumour *sub-region* labels, showing that the latter increases model accuracy for the task of tumour sub-region segmentation but also reduces the generality of the compositional representation, which loses anatomical detail and increases in pathology detail, becoming more focused on the supervision task.

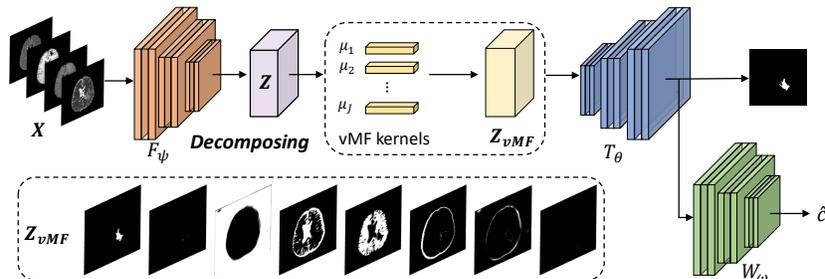


Fig. 1. Illustration of the brain tumour segmentation task using vMFBrain for compositional representation learning. We extract the weak supervision pathology labels (*presence or absence of tumour*) from 2-point brain tumour annotations; interestingly, learning of anatomical structures somewhat emerges even without supervision. Notation is specified in Section 3.

2 Related work

Compositionality is a fundamental concept in computer vision, where it refers to the ability to recognise complex objects or scenes by detecting and combining simpler components or features [13]. Leveraging this idea, compositional representation learning is an area of active research in computer vision [27]. Early approaches to compositional representation learning in computer vision include the bag-of-visual-words model [12] and part-based models [11]. Compositional representation learning has been applied to fine-grained recognition tasks in computer vision, such as recognising bird species [9,23]. In addition, compositionality has been incorporated for robust image classification [11, 25] and recently for compositional image synthesis [3, 14]. Among these works, Compositional Networks [11], originally designed for robust classification under object occlusion, are easier to extend to pixel-level tasks as they estimate spatial and interpretable vMF likelihoods. Previous work integrates vMF kernels [11] for object localisation [26] and recently for nuclei segmentation (with the bounding box as supervision) in a weakly supervised manner [28]. More recently, vMFNet [18] applies vMF kernels for cardiac image segmentation in the domain generalisation setting. Additionally, vMFNet integrated weak labels indicating the presence or absence of cardiac structures and this gave improved performance [15]. We use similar types of weak image-level annotations but apply the vMF kernels to pathology segmentation and supervise with weak labels indicating the presence or absence of pathological structures.

3 Method

We apply vMFNet [15,18], as shown in Fig. 1, a model consisting of three modules: the feature extractor F_ψ , the task network T_θ (for brain tumour segmentation in our case), and the weak supervision network W_ω , where ψ , θ and ω

denote the network parameters. Compositional components are learned as vMF kernels by decomposing the features extracted by \mathbf{F}_ψ . Then, the vMF likelihoods that contain spatial information are used to predict the tumour segmentation mask with \mathbf{T}_θ . The voxel-wise output of \mathbf{T}_θ is also input to the weak supervision network \mathbf{W}_ω to predict the presence or absence of the tumour. This framework is detailed below. We term our implementation *vMFBrain*.

3.1 Background: learning compositional components

To learn compositional components, the image features $\mathbf{Z} \in \mathbb{R}^{H \times W \times D}$ are first extracted by \mathbf{F}_ψ . H and W are the spatial dimensions and D is the number of channels. The feature vector $\mathbf{z}_i \in \mathbb{R}^D$ is defined as the normalised vector (i.e. $\|\mathbf{z}_i\| = 1$) across channels at position i on the 2D lattice of the feature map. Then, the image features are modelled with J vMF distributions. Each distribution has a learnable mean that is defined as vMF kernel $\boldsymbol{\mu}_j \in \mathbb{R}^D$. To ensure computational tractability, a fixed variance σ is set for all distributions. The vMF likelihood for the j^{th} distribution at each position i is calculated as:

$$p(\mathbf{z}_i | \boldsymbol{\mu}_j) = \frac{e^{\sigma_j \boldsymbol{\mu}_j^T \mathbf{z}_i}}{C}, \text{ s.t. } \|\boldsymbol{\mu}_j\| = 1, \quad (1)$$

where C is a constant. This gives the vMF likelihood vector $\mathbf{z}_{i,vMF} \in \mathbb{R}^J$, a component of $\mathbf{Z}_{vMF} \in \mathbb{R}^{H \times W \times J}$, which determines which kernel is activated at each position. To update the kernels during training, the clustering loss \mathcal{L}_{clu} is defined in [11] as:

$$\mathcal{L}_{clu}(\boldsymbol{\mu}, \mathbf{Z}) = -(HW)^{-1} \sum_i \max_j \boldsymbol{\mu}_j^T \mathbf{z}_i, \quad (2)$$

where the kernel $\boldsymbol{\mu}_j$ which is maximally activated for each feature vector \mathbf{z}_i is found, and the distance between the feature vectors and their corresponding kernels is minimised by updating the kernels. Overall, feature vectors in different images corresponding to the same anatomical or pathological structure will be clustered and activate the same kernels. Hence, the vMF likelihoods \mathbf{Z}_{vMF} for the same anatomical or pathological features in different images will be aligned to follow the same distributions (with the same means).

3.2 vMFBrain for brain tumour segmentation

Taking the vMF likelihoods as input, a follow-on segmentation task module \mathbf{T}_θ , is trained to predict the tumour segmentation mask, i.e. $\hat{\mathbf{Y}} = \mathbf{T}_\theta(\mathbf{Z}_{vMF})$. Firstly, we use direct strong supervision from the available pixel-level annotations \mathbf{Y} . Secondly, we define the weak supervision label c as a scalar (or a vector \mathbf{c}) which indicates the presence or absence of the tumour (or the presence or absence of the tumour sub-regions) in the 2D image slice. We use the output of the segmentation module as the input for a weak supervision classifier \mathbf{W}_ω i.e. $\hat{c} = \mathbf{W}_\omega(\hat{\mathbf{Y}})$. We train the classifier using $L1$ distance i.e. $\mathcal{L}_{weak}(\hat{c}, c) = |\hat{c} - c|_1$.

Overall, the model contains trainable parameters ψ, θ, ω and the vMF kernel means $\boldsymbol{\mu}$. The model (including all the modules) is trained **end-to-end** with the following objective:

$$\operatorname{argmin}_{\psi, \theta, \omega, \boldsymbol{\mu}} \mathcal{L}_{clu} + \lambda_{Dice} \mathcal{L}_{Dice}(\mathbf{Y}, \hat{\mathbf{Y}})(\boldsymbol{\mu}, \mathbf{Z}) + \lambda_{weak} \mathcal{L}_{weak}(\hat{c}, c), \quad (3)$$

where \mathcal{L}_{Dice} is Dice loss [7, 20]. We set $\lambda_{Dice} = 1$ when the ground-truth mask \mathbf{Y} is available, otherwise $\lambda_{Dice} = 0$. We set λ_{weak} as 0.5 for the whole tumour segmentation task and λ_{weak} as 0.1 for the tumour sub-region segmentation task (values determined empirically).

4 Experiments

4.1 Dataset

We evaluate on the task of brain tumour segmentation using data from the BraTS 2021 challenge [4, 5, 19]. This data comprises native (T1), post-contrast T1-weighted (T1Gd), T2-weighted (T2), and T2 Fluid Attenuated Inversion Recovery (FLAIR) modality volumes for 1,251 patients from a variety of institutions and scanners. We split the data into train, validation and test sets containing 938, 62 and 251 subjects. The data has already been co-registered, skull-stripped and interpolated to the same resolution, each volume having 155 2D slices. Labels are provided for tumour sub-regions: the peritumoural edema (ED), the GD-enhancing tumour (ET), and the necrotic and non-enhancing tumour (NE). We additionally downscale all images to 128×128 .

4.2 Baselines

We compare to the baselines **UNet** [22], **SDNet** [6] and **vMFNet** [18]. **SDNet** [6] is a semi-supervised disentanglement model with anatomy and modality encoders to separately encode the anatomical structure information and the imaging characteristics. The anatomical features are used as the input to the segmentor for the task of segmentation; the model is also trained with unlabelled data on the task of reconstructing the image by recombining the anatomy and modality factors. We compare to **vMFNet** with the architecture and training loss as described in [18]; this setup does not use weak supervision and has an additional image reconstruction module which we found empirically not to help performance (which thus we omit from vMFBrain).

4.3 Implementation

Imaging backbone: F_ψ is a 2D UNet [22] (without the output classification layer) to extract features \mathbf{Z} . The four modalities are concatenated as the input (with 4 channels) to F_ψ . For a fair comparison, we use this same UNet implementation as the backbone for all models.

Table 1. Dice (%) and Hausdorff Distance (HD) results for the task of **whole tumour segmentation**. We report the mean and standard deviation across volumes.

Metrics	Dice (\uparrow)				HD (\downarrow)			
	0.1%	0.5%	1%	100%	0.1%	0.5%	1%	100%
Pixel labels	0.1%	0.5%	1%	100%	0.1%	0.5%	1%	100%
UNet	80.66 ₁₀	86.39 _{7.7}	87.34 _{7.0}	90.84 _{5.6}	9.18 ₁₀	6.60 _{8.1}	7.37 ₁₀	4.49_{7.2}
SDNet	79.20 ₁₁	86.38 _{7.6}	87.96 _{6.6}	90.96_{5.3}	11.85 ₁₃	7.24 _{9.3}	6.11 _{8.4}	4.87 _{8.3}
vMFNet	81.30 _{9.6}	86.14 _{7.8}	87.98 _{6.6}	90.62 _{5.8}	11.62 ₁₃	9.12 ₁₂	7.15 _{9.6}	5.20 _{8.2}
vMFBrain w/o weak	79.70 ₁₀	84.92 _{8.1}	87.26 _{6.7}	90.67 _{5.8}	13.89 ₁₄	9.80 ₁₃	7.18 _{9.4}	4.93 _{7.3}
vMFBrain	85.64_{7.8}	88.64_{6.8}	89.04_{6.7}	90.58 _{5.6}	7.75_{7.8}	6.18_{8.6}	6.14_{8.4}	4.60 _{6.5}

Table 2. Dice (%) and Hausdorff Distance (HD) results for the task of **tumour sub-region segmentation**. We report the mean and standard deviation across volumes.

0.1% pixel labelled data	ED		ET		NE	
	Dice (\uparrow)	HD (\downarrow)	Dice (\uparrow)	HD (\downarrow)	Dice (\uparrow)	HD (\downarrow)
UNet	71.47 ₁₂	9.60 ₁₁	83.74 _{8.4}	5.19_{5.7}	79.42 ₁₀	10.24 _{7.9}
SDNet	75.87 ₁₁	10.17 ₁₁	82.45 _{8.8}	7.74 ₁₂	80.70 _{9.7}	9.89 ₁₁
vMFNet	71.11 ₁₂	10.06 _{9.8}	80.92 _{9.6}	7.99 ₁₁	78.37 ₁₁	12.97 ₁₁
vMFBrain w/o weak	70.65 ₁₃	15.65 ₁₅	79.36 ₁₂	13.13 ₁₇	79.33 _{9.8}	9.50 _{9.3}
vMFBrain w/ whole tumour weak	75.02 ₁₁	11.56 ₁₂	84.59 _{8.5}	8.17 ₁₂	79.48 _{9.6}	10.02 _{9.1}
vMFBrain w/ tumour sub-region weak	78.43_{9.8}	9.14_{8.8}	85.77_{8.2}	5.90 _{7.7}	81.31_{9.0}	8.08_{7.0}

vMFNet and vMFBrain⁴: T_θ is a shallow convolutional network. W_ω is a classifier model. Following [11], we set the variance of the vMF distributions as 30. The number of kernels is set to 8, as this number performed best empirically in our early experiments. For vMF kernel initialisation, we pre-train a 2D UNet for 10 epochs to reconstruct the input image with all the training data. After training, we extract the corresponding feature vectors and perform k-means clustering, then use the discovered cluster centres to initialise the vMF kernels.

Training: All models are implemented in PyTorch [21] and are trained using an NVIDIA 3090 GPU. Models are trained using the Adam optimiser [10] with a learning rate of $1 \times e^{-4}$ using batch size 32. In semi-supervised and weakly supervised settings, we consider the use of different percentages of fully labelled data to train the models. For this purpose, we randomly sample 2D image slices and the corresponding pixel-level labels from the whole training dataset.

4.4 Results

We compare model performance quantitatively using volume-wise Dice (%) and Hausdorff Distance (95%) (HD) [8] as the evaluation metrics, and qualitatively using the interpretability and compositionality of representations. In Table 1 and Table 2, for semi-supervised and weakly supervised approaches, the training data contains all unlabelled or weakly labelled data alongside different percentages of fully labelled data. UNet is trained with different percentages of labelled data only. Bold numbers indicate the best performance. Arrows (\uparrow, \downarrow) indicate the direction of metric improvement.

⁴ The code for vMFNet is available at <https://github.com/vios-s/vMFNet>.

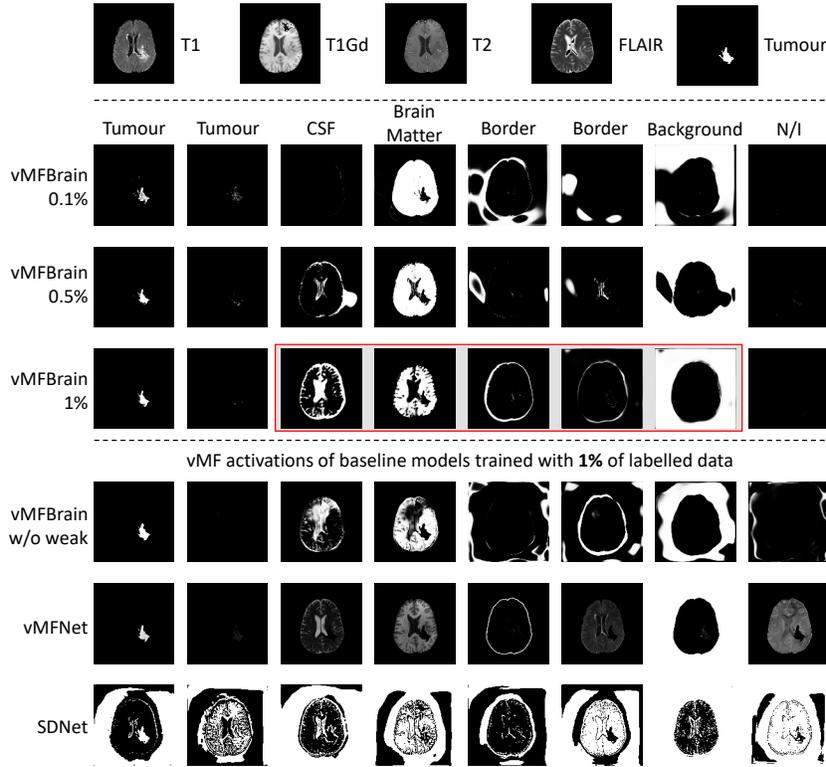


Fig. 2. Visualisation of vMF compositional representations (whole tumour supervision). We show the 4 input image modalities, the ground truth tumour segmentation mask, and all 8 vMF channels for the vMFBrain and baseline models trained with different percentages of labelled data. In the red boxes, the other interpretable vMF activations (excluding the tumour kernels) are highlighted. The vMF channels are ordered manually. For the vMFBrain channels, we label with a clinician’s visual interpretation of which image features activated each kernel. N/I denotes non-interpretable.

Brain tumour segmentation with weak labels: Overall, as reported in Table 1, the proposed vMFBrain model achieves best performance for most of the cases, particularly when very few annotations are available, i.e. the 0.1% case. When dropping the weak supervision (vMFBrain w/o weak), we observe reduced performance, which confirms the effectiveness of weak supervision. We also observe that the reconstruction of the original image (in vMFNet) does not help. It is possible that reconstruction of the tumour does not help here because the tumour has inconsistent appearance and location between different scans. With more annotated data, all models gradually achieve better performance, as expected. Notably, with only 1% labelled data vMFBrain achieves comparable performance (89.04 on Dice and 6.14 on HD) to the fully supervised UNet trained with all labelled data (90.84 on Dice and 4.49 on HD).

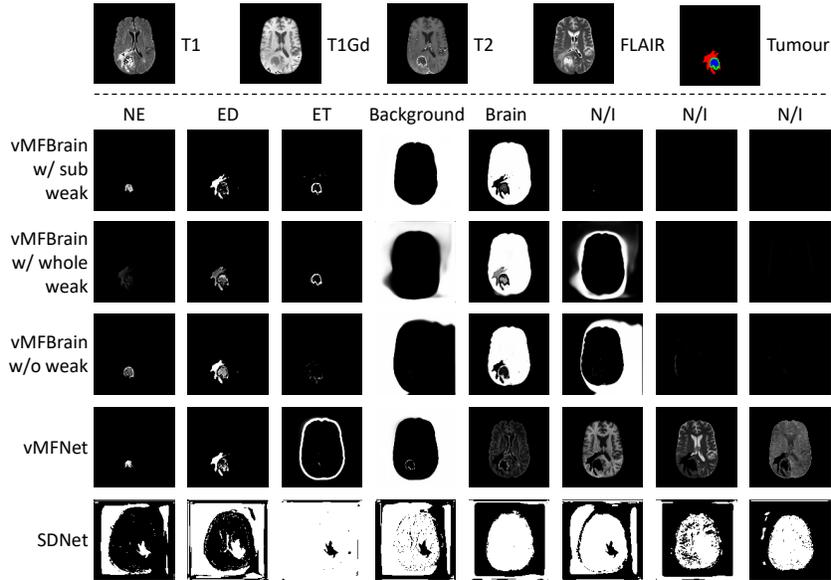


Fig. 3. Visualisation of vMF compositional representations (tumour sub-region supervision). We show the 4 input image modalities, the corresponding ground truth tumour sub-region segmentation mask and all 8 channels of the representations for the models trained with 1% of labelled data. The channels are ordered manually. For the vMFBrain channels, we label with a clinician’s visual interpretation of which image features activated each kernel. N/I denotes non-interpretable.

Tumour sub-region segmentation: We also report the results of tumour sub-region segmentation task in Table 2. For this task, we perform experiments using different weak labels: a) the weak label indicating the presence of the whole tumour i.e. vMFBrain w/ whole weak and b) the weak label indicating the presence of the tumour sub-regions i.e. vMFBrain w/ sub weak. It can be seen that our proposed vMFBrain performed best with both types of weak labels. Predictably, the best performance occurs when more task-specific weak labels (i.e. weak supervision on the tumour sub-regions) are provided.

Interpretability of compositional representations: We are particularly interested in the compositionality of the representations when pixel labels are not sufficient. In Fig. 2, we show the kernel activations. Note that the channels are ordered manually. For different runs, the learning is emergent such that kernels randomly learn to represent different components. Clearly, one of the kernels corresponds to the tumour in all cases. Using this kernel, we can detect and locate the tumours. For vMFBrain, training with more labelled data improves the compositionality of the kernels and the activations i.e. different kernels correspond to different anatomical or pathological structures, which are labelled by a clinician performing visual inspection of which image features activated each

channel. The most interpretable and compositional representation is vMFBrain trained with 1% labelled data. As highlighted in the red boxes, the kernels relate to CSF, brain matter, and the border of the brain even without any information about these structures given during training. Qualitatively, vMFBrain decomposes this information better into each kernel i.e. learns better compositional representations compared to other baseline models. Notably, weak supervision improves compositionality. We also show in Fig. 3 the representations for sub-region segmentation. Overall, we observe that with the more task-specific weak labels, the kernels learn to be more aligned with the sub-region segmentation task, where less information on other clinically relevant features is learnt.

5 Conclusion

In this paper, we have presented vMFBrain, a compositional representation learning framework. In particular, we constructed weak labels indicating the presence or absence of the brain tumour and tumour sub-regions in the image. Training with weak labels, better compositional representations can be learnt that produce better brain tumour segmentation performance when the availability of pixel-level annotations is limited. Additionally, our experiments show the interpretability of the compositional representations, where each kernel corresponds to specific anatomical or pathological structures. Importantly, according to our experiments and the results reported in previous studies [15, 18], the vMF-based compositional representation learning framework is robust and applicable to different medical datasets and tasks. In future work, we might consider transferring vMFBrain to 3D in order to process wider spatial context for each structure.

Acknowledgements S.A. Tsafaris acknowledges the support of Canon Medical and the Royal Academy of Engineering and the Research Chairs and Senior Research Fellowships scheme (grant RCSR1819\8\25). Many thanks to Patrick Schrempf and Joseph Boyle for their helpful review comments.

References

1. Alexander, R.G., Waite, S., Macknik, S.L., Martinez-Conde, S.: What do radiologists look for? Advances and limitations of perceptual learning in radiologic search. *Journal of Vision* **20**(10), 17–17 (2020)
2. Antonelli, M., Reinke, A., Bakas, S., Farahani, K., Kopp-Schneider, A., Landman, B.A., Litjens, G., Menze, B., Ronneberger, O., Summers, R.M., et al.: The medical segmentation decathlon. *Nature communications* **13**(1), 4128 (2022)
3. Arad Hudson, D., Zitnick, L.: Compositional transformers for scene generation. In: *Proc. Advances in Neural Information Processing Systems (NeurIPS)* (2021)
4. Bakas, S., Akbari, H., Sotiras, A., Bilello, M., Rozycki, M., Kirby, J.S., Freymann, J.B., Farahani, K., Davatzikos, C.: Advancing the cancer genome atlas glioma MRI collections with expert segmentation labels and radiomic features. *Scientific data* **4**(1), 1–13 (2017)

5. Bakas, S., Reyes, M., Jakab, A., Bauer, S., Rempfler, M., Crimi, A., Shinohara, R.T., Berger, C., Ha, S.M., Rozycki, M., et al.: Identifying the best machine learning algorithms for brain tumor segmentation, progression assessment, and overall survival prediction in the BraTS challenge. *arXiv preprint arXiv:1811.02629* (2018)
6. Chartsias, A., Joyce, T., Papanastasiou, G., Semple, S., Williams, M., Newby, D.E., Dharmakumar, R., Tsaftaris, S.A.: Disentangled representation learning in cardiac image analysis. *Medical image analysis* **58**, 101535 (2019)
7. Dice, L.R.: Measures of the amount of ecologic association between species. *Ecology* **26**(3), 297–302 (1945)
8. Dubuisson, M.P., Jain, A.K.: A modified Hausdorff distance for object matching. In: *Proc. International Conference on Pattern Recognition (ICPR)*. vol. 1, pp. 566–568. IEEE (1994)
9. Huynh, D., Elhamifar, E.: Compositional zero-shot learning via fine-grained dense feature composition. In: *Proc. Advances in Neural Information Processing Systems (NeurIPS)*. vol. 33, pp. 19849–19860 (2020)
10. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: *Proc. International Conference on Learning Representations (ICLR)* (2015)
11. Kortylewski, A., He, J., Liu, Q., Yuille, A.L.: Compositional convolutional neural networks: A deep architecture with innate robustness to partial occlusion. In: *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 8940–8949 (2020)
12. Kortylewski, A., Liu, Q., Wang, H., Zhang, Z., Yuille, A.: Combining compositional models and deep networks for robust object classification under occlusion. In: *Proc. IEEE/CVF winter conference on applications of computer vision (CVPR)*. pp. 1333–1341 (2020)
13. Lake, B.M., Salakhutdinov, R., Tenenbaum, J.B.: Human-level concept learning through probabilistic program induction. *Science* **350**(6266), 1332–1338 (2015)
14. Liu, N., Li, S., Du, Y., Tenenbaum, J., Torralba, A.: Learning to compose visual relations. In: *Proc. Advances in Neural Information Processing Systems (NeurIPS)*. vol. 34 (2021)
15. Liu, X., Sanchez, P., Thermos, S., O’Neil, A.Q., Tsaftaris, S.A.: Compositionally equivariant representation learning. *arXiv preprint arXiv:2306.07783* (2023)
16. Liu, X., Thermos, S., Chartsias, A., O’Neil, A., Tsaftaris, S.A.: Disentangled representations for domain-generalized cardiac segmentation. In: *Proc. International Workshop on Statistical Atlases and Computational Models of the Heart (STA-COM)*. pp. 187–195 (2020)
17. Liu, X., Thermos, S., O’Neil, A., Tsaftaris, S.A.: Semi-supervised meta-learning with disentanglement for domain-generalised medical image segmentation. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. pp. 307–317. Springer (2021)
18. Liu, X., Thermos, S., Sanchez, P., O’Neil, A.Q., Tsaftaris, S.A.: vmfnet: Compositionality meets domain-generalised segmentation. In: *Proc. International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. pp. 704–714. Springer (2022)
19. Menze, B.H., Jakab, A., Bauer, S., Kalpathy-Cramer, J., Farahani, K., Kirby, J., Burren, Y., Porz, N., Slotboom, J., Wiest, R., et al.: The multimodal brain tumor image segmentation benchmark (BraTS). *IEEE transactions on medical imaging* **34**(10), 1993–2024 (2014)
20. Milletari, F., Navab, N., Ahmadi, S.A.: VNet: Fully convolutional neural networks for volumetric medical image segmentation. In: *3DV*. pp. 565–571. IEEE (2016)

21. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. In: Proc. Advances in Neural Information Processing Systems (NeurIPS). vol. 32 (2019)
22. Ronneberger, O., Fischer, P., Brox, T.: UNet: Convolutional networks for biomedical image segmentation. In: Proc. International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI). pp. 234–241. Springer (2015)
23. Singh, K.K., Ojha, U., Lee, Y.J.: Finegan: Unsupervised hierarchical disentanglement for fine-grained object generation and discovery. In: Proc. IEEE/CVF conference on computer vision and pattern recognition (CVPR). pp. 6490–6499 (2019)
24. Thermos, S., Liu, X., O’Neil, A., Tsiftaris, S.A.: Controllable cardiac synthesis via disentangled anatomy arithmetic. In: International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI). pp. 160–170. Springer (2021)
25. Tokmakov, P., Wang, Y.X., Hebert, M.: Learning compositional representations for few-shot recognition. In: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6372–6381 (2019)
26. Yuan, X., Kortylewski, A., et al.: Robust instance segmentation through reasoning about multi-object occlusion. In: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 11141–11150 (2021)
27. Yuille, A.L., Liu, C.: Deep nets: What have they ever done for vision? *International Journal of Computer Vision* **129**, 781–802 (2021)
28. Zhang, Y., Kortylewski, A., Liu, Q., et al.: A light-weight interpretable compositional network for nuclei detection and weakly-supervised segmentation. arXiv preprint arXiv:2110.13846 (2021)