

# Cross-Modal Transformer GAN: A Brain Structure-Function Deep Fusing Framework for Alzheimer’s Disease

Junren Pan, Shuqiang Wang

Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen  
518000, China  
sq.wang@siat.ac.cn

**Abstract.** Cross-modal fusion of different types of neuroimaging data has shown great promise for predicting the progression of Alzheimer’s Disease(AD). However, most existing methods applied in neuroimaging can not efficiently fuse the functional and structural information from multi-modal neuroimages. In this work, a novel cross-modal transformer generative adversarial network(CT-GAN) is proposed to fuse functional information contained in resting-state functional magnetic resonance imaging (rs-fMRI) and structural information contained in Diffusion Tensor Imaging (DTI). The developed bi-attention mechanism can match functional information to structural information efficiently and maximize the capability of extracting complementary information from rs-fMRI and DTI. By capturing the deep complementary information between structural features and functional features, the proposed CT-GAN can detect the AD-related brain connectivity, which could be used as a bio-marker of AD. Experimental results show that the proposed model can not only improve classification performance but also detect the AD-related brain connectivity effectively.

**Keywords:** Cross-modal fusion · Transformer · Bi-attention Mechanism  
· Brain Network · Generative Adversarial Strategy

## 1 Introduction

Alzheimer’s disease (AD), a chronic and irreversible neurodegenerative disease, is the main reason for dementia among aged people [1]. According to statistical data given by Alzheimer’s Association[2], at least 50 million people worldwide are suffering from AD. AD patients will gradually lose cognitive function such as remembering or thinking, and eventually become unable to take care of themselves [3]. The widespread incidence of AD makes a severe financial burden to both patients’ families and governments. With the development of artificial intelligence [4,5,6,7,8,9], researchers study AD from different angles using machine learning technology [10,11,12]. However, the cause of AD has not been fully revealed. One of the main reasons for the above difficulties is that brain is a highly

complex network, and completing cognitions requires specific coordination between regions-of-interest (ROIs).

A brain network can be characterized as a graph. The nodes of a brain network represent ROIs of the brain. The edges of a brain network represent the interaction relationship between ROIs of the brain. There are two basic connectivity categories of brain networks: functional connectivity (FC) and structural connectivity (SC). FC is defined as the interdependence between the blood-oxygen-level-dependent (BOLD) signals of two ROIs, where BOLD signals can be extracted from rs-fMRI. SC is defined as the neural fibers connection strength among ROIs, which can be extracted from DTI. Many studies [13,14,15] have used FC or SC to obtain some AD-related features that can not be discovered in traditional imaging methods. This shows that brain network methods have more advantages than the traditional imaging method in AD research. However, most existing brain network studies are based on a single modal, which can only focus on one of the brain structural information and brain functional information. Since single modality data may only contain complementary cross-modal information partially, it will lose an opportunity to take advantage of complementary cross-modal information to study AD more accurately. Therefore multimodal brain network methods [16,17,18,19] are gaining more and more attention in medical imaging computing. For the structure-function deep fusing task of multimodal brain networks, the key is how to efficiently use complementary cross-modal information that is heterogeneous and hidden in different types of neuroimaging data. Most existing structure-function fusion approaches just used linear relationships between different modalities. However, changes of brain structure and function can not be fully characterized by linear relationships. Previous studies[20,21] proved that strong SC inclines to be accompanied with strong FC, but not vice versa. On the other hand, clinical studies[22,23,24] show that when an SC between ROIs is reduced, some regions can increase functional activity to compensate for the reduced SC.

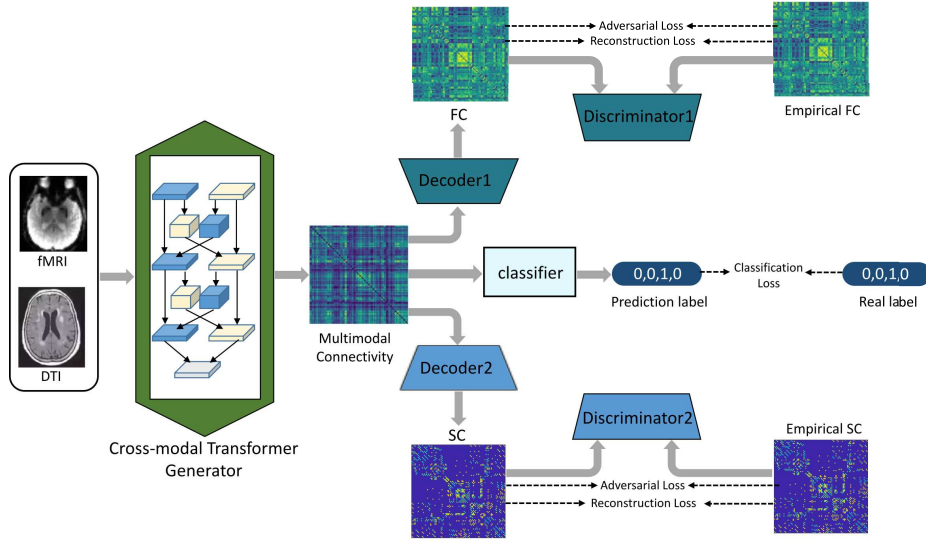
To overcome the above problem, a novel cross-modal transformer is proposed in this work to deal with the structure-function deep fusion task based on generative adversarial networks(GANs). GANs can be seen as variational-inference[25,26] based generative model. GANs are proved to be an efficient framework for learning complex distribution [27,28]. Currently, GANs are successfully used in many branches of medical image analysis [29,30,31,32,33,34]. Transformers [35] have shown their powerful capability for sequential analyzing in natural language processing (NLP). This is due to the self-attention mechanism which characterizes the nonlinear relationships between given inputs. Following their successful applications in the area of NLP, transformers have been adopted to image tasks very recently [36,37,38]. However, transformers have been few explored in the area of brain networks. In this study, a cross-modal transformer is proposed to fuse structure-function information of brain networks. The proposed cross-modal transformer consists of the following four modules: 1)convolutional neural networks(CNN) modules that are used to extract functional information from rs-fMRI; 2)graph convolutional networks(GCN) modules that

are used to extract structural information from DTI; 3) F2S-attention modules that transform from functional information to structural information; 4) S2F-attention modules that transform from structural information to functional information. The cross-modal transformer is based on a bi-attention mechanism where complementary information from rs-fMRI and DTI are fused layer by layer. On the other hand, a generative adversarial strategy is adopted to guide the proposed model’s training. The advantages of the proposed model are summarized as follows: 1) structural and functional information can be deeply fused; 2) complementary information from rs-fMRI and DTI can be effectively extracted; 3) due to the generative adversarial strategy, the proposed model does not need very deep nets, which make the model’s architecture and training more flexible.

## 2 Method

### 2.1 The Architecture

The proposed CT-GAN is illustrated in Fig. 1, which consists of four components: 1) a cross-modal transformer generator that outputs multimodal connectivity brain networks; 2) two decoders that decode multimodal connectivity to the corresponding SC and FC, respectively; 3) two discriminators, one of which determines whether an SC from learned by our proposed model or output by a software template, the other discriminator for FC is similar; 4) a classifier that predicts AD stages according to multimodal connectivity.



**Fig. 1.** The framework of the proposed CT-GAN.

## 2.2 Bi-Attention Mechanism

The key idea of this work is to exploit the bi-attention mechanism of transformers to fuse structure-function information for rs-fMRI and DTI given their complementary nature. A transformer mapping an input feature sequence as  $X = \mathbb{R}^{n \times d_x}$  to a target feature sequence as  $Y = \mathbb{R}^{n \times d_y}$ , where  $n$  is the total number of ROIs, can be described as follows. Firstly, a linear projections is used to compute a set of queries matrix  $Q$ , keys matrix  $K$ , and values matrix  $V$ ,

$$Q = XW^q, \quad K = XW^k, \quad V = XW^v \quad (1)$$

where  $W^q \in \mathbb{R}^{d_x \times d_q}$ ,  $W^k \in \mathbb{R}^{d_x \times d_k}$  with  $d_k = d_q$ , and  $W^v \in \mathbb{R}^{d_x \times d_v}$  are weight matrices. The attention of  $Q$ ,  $K$ , and  $V$  can be obtained as follows:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V. \quad (2)$$

And then, a fully connected layer(FC) is used to transform the attention of  $Q$ ,  $K$ , and  $V$  into a feature sequence with the same dimension as the target feature sequence  $Y$ . Finally, the output of a transformer is

$$X^{out} = \text{FC}(\text{Attention}(XW^q, XW^k, XW^v)) + \lambda Y, \quad (3)$$

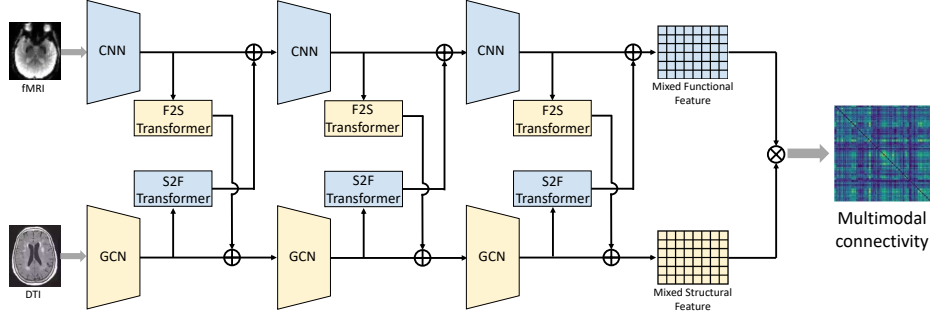
where  $\lambda$  is a hyper-parameter between 0 to 1.

There are two types of transformers in the proposed cross-modal generator. The one is used to transform functional information into structural information, abbreviated as F2S. The other is used to transform structural information into functional information, abbreviated as S2F. We can now introduce the bi-attention mechanism of transformers. BOLD signals extracted from rs-fMRI and empirical structural connectivity extracted from DTI are inputs to our cross-modal generator that uses several pairs of F2S and S2F modules for fusing intermediate features between the inputs. The intermediate features of BOLD signals and empirical structural connectivity are extracted by CNN and GCN, respectively. The detailed architecture of generator is shown in Fig. 2. Let us denote the mixed functional feature sequence and mixed structural feature sequence in Fig. 2 to be  $F$  and  $S$ , respectively. It is worth mentioning that  $F \in \mathbb{R}^{n \times d}$  represents the feature sequence of each ROI and  $S \in \mathbb{R}^{n \times n}$  represents the connective feature between ROIs where  $n$  is the total number of ROIs. The multimodal connectivity matrix MC as the final output of the proposed generator is obtained by

$$\text{MC} = SF^T S^T. \quad (4)$$

## 2.3 Loss Function

Success in structure-function deep fusion tasks requires semantic reasoning. Therefore a multimodal connectivity matrix learned by the proposed model should be able to decouple to be the corresponding empirical FC matrix and empirical



**Fig. 2.** The network architecture of proposed generator.

SC matrix. A generative adversarial strategy is used to achieve the above process. To guarantee the quality of generated multimodal connectivity, a hybrid loss function is used to optimize the network, including three types of terms: the adversarial loss, the classification loss, and the pair-wise connectivity reconstruction loss. The details are given as follows.

**Adversarial Loss.** To make the FC matrix and SC matrix decoded by multimodal connectivity matrix as close as possible to empirical FC matrices and SC matrices, the adversarial losses are defined as

$$\mathcal{L}_{adv}^{SC} = \mathbb{E}_x[\log D_1(SC_x)] + \mathbb{E}_x[\log(1 - D_1(Dec_1(G(x))))], \quad (5)$$

$$\mathcal{L}_{adv}^{FC} = \mathbb{E}_x[\log D_2(FC_x)] + \mathbb{E}_x[\log(1 - D_2(Dec_2(G(x))))], \quad (6)$$

where  $G$  is the generator,  $D_1$  and  $D_2$  are the discriminator1 and discriminator2,  $Dec_1$  and  $Dec_2$  are the decode1 and decode2 in Fig. 1.  $SC_x$  and  $FC_x$  represent the empirical SC matrix and empirical FC matrix of subject  $x$ . The generator  $G$  generates a multimodal connectivity matrix  $G(x)$  from subject  $x$ 's rs-fMRI and DTI. The decoders  $Dec_1$  and  $Dec_2$  decode  $G(x)$  into an FC matrix and an SC matrix, respectively. While the discriminator  $D_1$  attempts to distinguish between the FC matrix decoded by  $Dec_1$  and empirical FC matrices ( $D_2$  similarly works on SC matrices). The generative adversarial strategy is that  $G$  tries to minimize above adversarial losses while  $D$  tries to maximize it.

**Classification Loss.** For multimodal connectivity matrices, an important index to judge the effect of cross-modal fusing is whether they can achieve high accuracy in predicting AD stages. The classification loss is imposed when optimizing the classifier and the generator  $G$ . The formula of classification loss is given by

$$\mathcal{L}_{cls} = \mathbb{E}_{(x,y) \sim p_{real}(x,y)}[-\log p_c(y|G(x))], \quad (7)$$

where  $y$  represents AD stages, including normal controls (NC), early mild cognitive impairment (EMCI), late mild cognitive impairment (LMCI), and Alzheimer's Disease (AD). The  $p_c(y|G(x))$ , output by the classifier with the input  $G(x)$ , represents the probability that the subject  $x$  is now under stage  $y$ . By classification loss, the generator is trained to extract and fuse features, which contains more

disease information, from rs-fMRI and DTI. Meanwhile, the classifier can achieve the highest predicting accuracy through multimodal connectivity matrices.

**Pair-wise Connectivity Reconstruction Loss.** The  $L^1$  pair-wise connectivity reconstruction loss is used to impose an additional topological constraint on the generator  $G$ . It means that the generator  $G$  and the decoders  $Dec1, Dec2$  are not only needed to fool the discriminators  $D1$  and  $D2$ . In addition, they also need to minimize the sum of pair-wise connectivity difference between FC/SC matrices decoded by  $Dec1/Dec2$  and empirical FC/SC matrices. The  $L^1$  pair-wise reconstruction losses are formalized as follows:

$$\mathcal{L}_{rcs}^{FC} = \mathbb{E}_x \|FC_x - Dec_1(G(x))\|_1, \quad (8)$$

$$\mathcal{L}_{rcs}^{SC} = \mathbb{E}_x \|SC_x - Dec_2(G(x))\|_1. \quad (9)$$

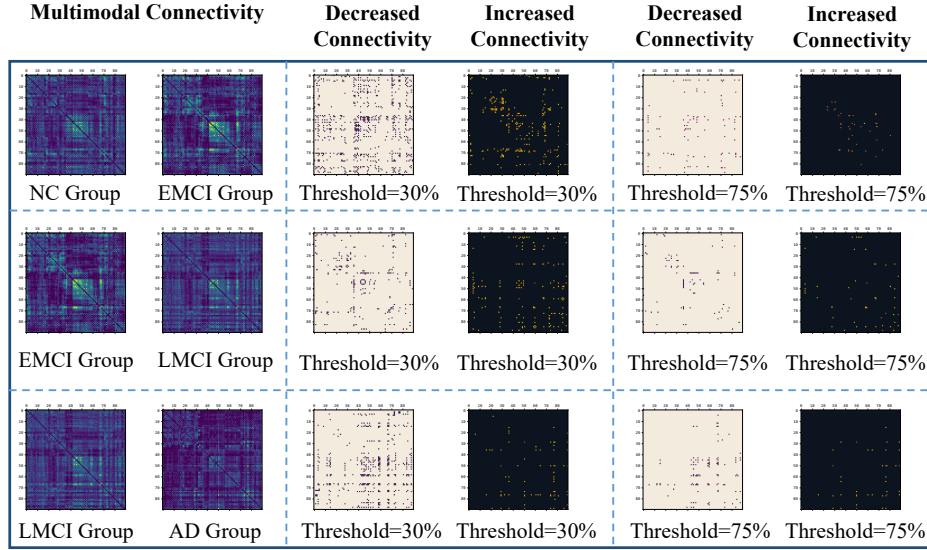
### 3 Experiments

DTI and rs-fMRI used to validate the proposed framework are from the public dataset ADNI(Alzheimer’s Disease Neuroimaging Initiative). There are 268 subjects’ data we used in this study whose detailed information can be seen in Table 1. AAL 90 atlas is used in the preprocessing. A string of preprocessing steps of rs-fMRI using the DPARSF toolbox is consisted of discarding the first 20 volumes, head motion correction, band-pass filtering, Gaussian smoothing, and extracting time series of all voxels. The structural connectivity is obtained by tracking fiber bundles between ROIs. The fiber tracking stopping conditions set in PANDA, the toolbox used to preprocess DTI, are following that: (1) the crossing angle between two consecutive moving directions is more than 45 degrees. (2) the fractional anisotropy value is not in the range of [0.2, 1.0].

**Table 1.** Subjects’ information in this study

Group	AD(63)	LMCI(41)	EMCI(80)	NC(84)
Male/Female	39M/24F	20M/21F	48M/32F	38M/46F
Age(mean $\pm$ SD)	75.3 $\pm$ 5.5	74.9 $\pm$ 5.3	75.8 $\pm$ 6.1	74.0 $\pm$ 5.9

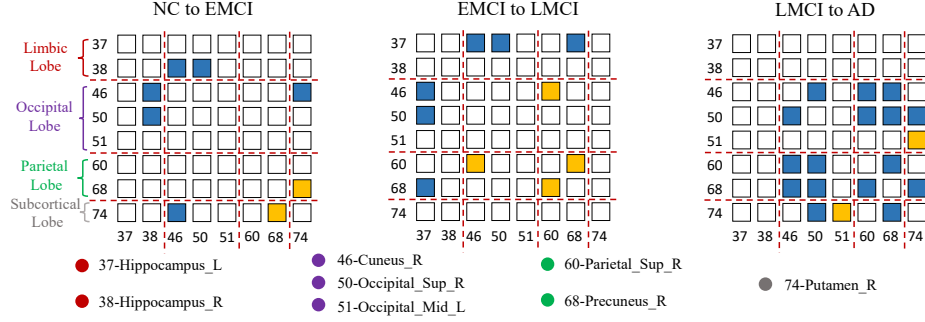
The decoders  $Dec_1$  and  $Dec_2$ , the discriminators  $D_1$  and  $D_2$ , the classifier are implemented by multilayer perceptron. The generator  $G$  consists of three blocks. Such blocks contain a CNN layer with kernel size (1,10), a GCN layer, an F2S transformer, and an S2F transformer. The hyper-parameter  $\lambda$  in transformers is set to be 0.1. The whole model is trained in an end-to-end manner. The Adam optimizer is chosen to update the model’s parameter during the training process. The hyper-parameter of the Adam optimizer, including learning rate, weight decay, and momentum rates, are set to be 0.001, 0.01, and (0.9,0.99), respectively.



**Fig. 3.** Averaged multimodal connectivity of different groups with the same disease stage and the changes of such connectivity under the proceeding of AD stages.

After applying the trained model to groups of subjects with the same disease stage, the averaged multimodal connectivity matrices of different disease stages can be obtained. One of the most important biomarkers of AD is crucial brain connectomes which influence the proceeding of AD stages. We characterize such connectomes by analyzing the increase/decrease connectivity of the averaged multimodal connectivity matrices of different disease stages. The visualization of averaged multimodal connectivity matrices and the increase/decrease connectivity with different thresholds are illustrated in Fig. 3. Based on multimodal connectivity matrices, the top 8 ROIs that have the most significant connectivity changing under the AD development process are 37-Hippocampus\_L, 38-Hippocampus\_R, 46-Cuneus\_R, 50-Occipital\_Sup\_R, 51-Occipital\_Mid\_L, 60-Parietal\_Sup\_R, 68-Precuneus\_R, and 74-Putamen\_R, where the number before ROI is the corresponding AAL id. The connectivity changing between the 8 ROIs above is shown in Fig. 4.

To compare the ability to represent AD-related features with different fMRI-DTI fusion models, three binary classification experiments, including AD vs. NC, LMCI vs. NC, and EMCI vs. NC, are designed. Three index, detection accuracy(ACC), sensitivity(SEN), and specificity(SPEC), are used to evaluate the models performance. The experiments' result are shown in Table 2. The experiments' result shows that the proposed multimodal fusion model has the advantage of higher accuracy for predicting AD stages than other existing multimodal fusion models.



**Fig. 4.** The change of multimodal connectivity between the top 8 AD-related ROIs. The number before ROI is the corresponding AAL id for this ROI. The blue represents decreased connectivity; The yellow represents increased connectivity. The red dotted lines divide the 8 ROIs into their corresponding brain lobe.

**Table 2.** Prediction performance in AD vs. NC, LMCI vs. NC, and EMCI vs. NC. under different fMRI-DTI fusion models

Study	Method	AD vs. NC			LMCI vs. NC			EMCI vs. NC		
		ACC	SEN	SPE	ACC	SEN	SPE	ACC	SEN	SPE
Aderghal et al. [39]	CNN+ Transfer Learning	<b>94.44%</b>	<b>93.33%</b>	<b>95.24%</b>	87.1%	87.5%	86.96%	85.37%	88.89%	82.61%
Dyrba et al. [40]	SVM+ Multi-kernel	86.11%	85.71%	86.36%	83.87%	72.73%	90.0%	82.93%	84.21%	81.82%
Lu et al. [41]	GCN+ Deep learning	91.67%	87.5%	95.0%	90.32%	88.89%	90.91%	90.24%	90.0%	90.48%
Proposed	GAN+ Bi-attention	<b>94.44%</b>	<b>93.33%</b>	<b>95.24%</b>	<b>93.55%</b>	<b>90.0%</b>	<b>95.24%</b>	<b>92.68%</b>	<b>90.48%</b>	<b>95.0%</b>

## 4 Conclusion

In this paper, we proposed a novel CT-GAN to fuse rs-fMRI and DTI to multimodal connectivity of brain network. The key idea of this work is that we use a bi-attention mechanism to achieve the goal of mutual conversion between structural and functional information. Therefore, the bi-attention mechanism above can help the proposed model efficiently extracts the complementary information between rs-fMRI and DTI. The experiments' results proved our multimodal connectivity has higher accuracy of AD prediction than other multimodal fusion methods. Through analyzing our multimodal connectivity matrices, some AD-related connectomes are obtained. These connectomes are highly consistent with the results of clinical AD studies. Although this work focuses only on AD, it is worth mentioning that the proposed model can be easily extended to apply to other neurodegenerative diseases.



## References

1. Dadar, M., Pascoal, T.A., Manitsirikul, S., Fonov, V.S., et al.: Validation of a Regression Technique for Segmentation of White Matter Hyperintensities in Alzheimer's Disease. *IEEE Trans. Med. Imaging* **36**(8), 1758–1768 (2017)
2. Alzheimer's Association.: 2019 Alzheimer's disease facts and figures. *Alzheimer's Dementia*. vol. 15, no. 3, pp. 321–387, 2019.
3. Association, A.s.: Alzheimer's disease facts and figures. *Alzheimer's Dement* **14**, 367–429 (2018)
4. Wang, S., et al.: Skeletal maturity recognition using a fully automated system with convolutional neural networks. *IEEE Access*, 2018, 6: 29979–29993.
5. Wang, S., Shen, Y., Zeng, D., Hu, Y.: Bone age assessment using convolutional neural networks. *International Conference on Artificial Intelligence and Big Data (ICAIBD)*, 2018, 175–178.
6. Wang, S., et al.: Prediction of myelopathic level in cervical spondylotic myelopathy using diffusion tensor imaging. *Journal of Magnetic Resonance Imaging*, 2015, 41(6): 1682–1688.
7. Wang, S., et al.: An ensemble-based densely-connected deep learning system for assessment of skeletal maturity. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 2020, 52(1): 426–437.
8. Hu, S., Yuan, J., Wang, S.: Cross-modality synthesis from MRI to PET using adversarial U-net with different normalization. *2019 International Conference on Medical Imaging Physics and Engineering (ICMIPE)*, 2019: 1–5.
9. You, S., et al.: Fine perceptive gans for brain mr image super-resolution in wavelet domain. *IEEE transactions on neural networks and learning systems*, 2022.
10. Wang, S., Wang, H., Shen, Y., Wang, X.: Automatic recognition of mild cognitive impairment and alzheimers disease using ensemble based 3d densely connected convolutional networks. *2018 17th IEEE International conference on machine learning and applications (ICMLA)*, 2018, 517–523.
11. Wang, S., Wang, H., Cheung, A., Shen, Y., Gan, M.: Ensemble of 3D densely connected convolutional network for diagnosis of mild cognitive impairment and Alzheimer's disease. *Deep learning applications*, 2020: 53–73.
12. Yu, S., et al.: Multi-scale enhanced graph convolutional network for early mild cognitive impairment detection. *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2020: 228–237.
13. Wang, S., Hu, Y., Shen, Y., Li, H.: Classification of diffusion tensor metrics for the diagnosis of a myelopathic cord using machine learning. *International journal of neural systems*. **28** (02) , 1750036 (2018)
14. Jeon, E., Kang, E., Lee, J., Lee, J., Kam, T., Suk, H.: Enriched representation learning in resting-state fMRI for early MCI diagnosis. In: Martel A.L. et al. (eds) *MICCAI 2020, LNCS*, vol. 12267, pp. 397–406. Springer, Cham (2020). [https://doi.org/10.1007/978-3-030-59728-3\\_39](https://doi.org/10.1007/978-3-030-59728-3_39)
15. Wang, S., Shen, Y., Chen, W., Xiao, T.: Automatic recognition of mild cognitive impairment from mri images using expedited convolutional neural networks. *International Conference on Artificial Neural Networks*, 373–380 (2017)
16. Zhang, D., Shen, D.: Multi-modal multi-task learning for joint prediction of multiple regression and classification variables in Alzheimer's disease. *NeuroImage* 59, 895–907 (2012)
17. Li, L., Kang, J., Lockhart, S., et al.: Spatially adaptive varying correlation analysis for multimodal neuroimaging data. *IEEE transactions on medical imaging*, 2018, 38(1): 113–123.

18. Hao, X., Bao, Y., Guo, Y., et al.: Multi-modal neuroimaging feature selection with consistent metric constraint for diagnosis of Alzheimer’s disease. *Medical image analysis*, 2020, 60: 101625.
19. Pan, J., Lei, B., Shen, Y., et al.: Characterization Multimodal Connectivity of Brain Network by Hypergraph GAN for Alzheimer’s Disease Analysis. *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*. Springer, Cham, 2021: 467-478.
20. Honey, C., Sporns, O., Cammoun, L., et al.: Predicting human resting-state functional connectivity from structural connectivity. *Proceedings of the National Academy of Sciences*, 2009, 106(6): 2035–2040.
21. Li, K., Guo, L., Zhu, D., et al.: Individual functional ROI optimization via maximization of group-wise consistency of structural and functional profiles. *Neuroinformatics*, 2012, 10(3): 225–242.
22. Daselaar, S., Iyengar, V., Davis., et al.: Less wiring, more firing: low-performing older adults compensate for impaired white matter with greater neural activity. *Cerebral cortex*, 2015, 25(4): 983–990.
23. Lei, B., Cheng, N., Frangi, A.F., Tan, E.-L., Cao, J., Yang, P., et al.: Self-calibrated brain network estimation and joint non-convex multi-task learning for identification of early Alzheimer’s disease. *Med. Image Anal.* 61, 101652 (2020)
24. Cao, P., et al.: Generalized fused group lasso regularized multi-task feature learning for predicting cognitive outcomes in Alzheimers disease. *Comput. Meth. Programs Biomed.* **162**, 19–45 (2018)
25. Wang, S.: A variational approach to nonlinear two-point boundary value problems. *Computers & Mathematics with Applications*, 2009, 58(11-12): 2452–2455.
26. Wang, S., He, J.: Variational iteration method for a nonlinear reaction-diffusion process. *International Journal of Chemical Reactor Engineering*, 2008, 6(1).
27. Goodfellow, I., et al.: Generative adversarial nets. In: *Advances in Neural Information Processing Systems*, pp. 2672–2680 (2014).
28. Wang, S., et al.: Diabetic retinopathy diagnosis using multichannel generative adversarial network with semisupervision. *IEEE Transactions on Automation Science and Engineering*, 2020.
29. Yu, W., et al.: Tensorizing GAN with high-order pooling for Alzheimer’s disease assessment. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
30. Hu, S., Shen, Y., Wang, S., Lei, B.: Brain mr to pet synthesis via bidirectional generative adversarial network. *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2020: 698–707.
31. Hu, S., Yu, W., Chen, Z., Wang, S.: Medical image reconstruction using generative adversarial network for Alzheimer disease assessment with class-imbalance problem. *2020 IEEE 6th International Conference on Computer and Communications (ICCC)*, 2020: 1323–1327.
32. Hu, S., et al.: Bidirectional mapping generative adversarial networks for brain MR to PET synthesis. *IEEE Transactions on Medical Imaging*, 2021, 41(1): 145–157.
33. Yu, W., et al.: Morphological feature visualization of Alzheimer’s disease via Multi-directional Perception GAN. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
34. Pan, J., et al.: DecGAN: Decoupling Generative Adversarial Network detecting abnormal neural circuits for Alzheimer’s disease. *arXiv preprint arXiv:2110.05712*.
35. Vaswani, A., Shazeer, N., Parmar, N., et al.: Attention is all you need. *Advances in neural information processing systems*, 2017, 30.

36. Lanchantin, J., Wang, T., Ordonez, V., et al.: General multi-label image classification with transformers. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021: 16478-16488.
37. Esser, P., Rombach, R., Ommer, B.: Taming transformers for high-resolution image synthesis. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2021: 12873-12883.
38. Hudson, D., Zitnick, L.: Generative adversarial transformers. *International Conference on Machine Learning*. PMLR, 2021: 4487-4499.
39. Aderghal, K., Khvostikov, A., Krylov, A., Benois-Pineau, J., Afdel, K., Catheline, G.: Classification of alzheimer disease on imaging modalities with deep cnns using cross-modal transfer learning. In: *2018 IEEE 31st International Symposium on Computer-Based Medical Systems (CBMS)*. IEEE, pp. 345–350.
40. Dyrba, M., Grothe, M., Kirste, T., Teipel, S.J.: Multimodal analysis of functional and structural disconnection in a lzheimer’s disease using multiple kernel svm. *Human brain mapping* 36, 2118–2131.
41. Zhang, L., Wang, L., Gao, J., et al.: Deep fusion of brain structure-function in mild cognitive impairment. *Medical image analysis*, 2021, 72: 102082.