

**Differential diagnosis of parkinsonism based on deep metabolic imaging indices**

Ping Wu<sup>1,2†</sup>, Yu Zhao<sup>3,4,5†</sup>, Jianjun Wu<sup>6</sup>, Matthias Brendel<sup>7</sup>, Jiaying Lu<sup>1,2,3</sup>, Jingjie Ge<sup>1</sup>, Alexander Bernhardt<sup>8</sup>, Ling Li<sup>1</sup>, Ian Alberts<sup>3</sup>, Sabrina Katzdobler<sup>8</sup>, Igor Yakushev<sup>9</sup>, Jimin Hong<sup>3</sup>, Qian Xu<sup>1</sup>, Yimin Sun<sup>2,6</sup>, Fengtao Liu<sup>2,6</sup>, Johannes Levin<sup>8</sup>, Günter U Höglinger<sup>10</sup>, Claudio Bassetti<sup>11</sup>, Yihui Guan<sup>1,2</sup>, Wolfgang H Oertel<sup>12</sup>, Wolfgang Weber<sup>9</sup>, Axel Rominger<sup>3</sup>, Jian Wang<sup>2,6\*</sup>, Chuantao Zuo<sup>1,2\*</sup>, Kuangyu Shi<sup>3,5</sup>

**†Equally-contributed first authors**

**\*Equally-contributed senior authors**

1 PET Center, Huashan Hospital, Fudan University, Shanghai, China

2 National Research Center for Aging and Medicine & National Center for Neurological Disorders, Huashan Hospital, Fudan University, Shanghai, China

3 Department of Nuclear Medicine, University of Bern, Bern, Switzerland

4 AI Lab, Tencent, Shenzhen, China

5 Department of Informatics, Technische Universität München, Munich, Germany

6 Department of Neurology, Huashan Hospital, Fudan University, Shanghai, China

7 Department of Nuclear Medicine, University of Munich, Munich, Germany

8 Department of Neurology, University of Munich, Munich, Germany

9 Department of Nuclear Medicine, Technische Universität München, Munich, Germany

10 Department of Neurology, Hannover Medical School, Hannover, Germany

11 Department of Neurology, University of Bern, Bern, Switzerland

12 Department of Neurology, University of Marburg, Marburg, Germany

**Running title:** DMI indices for parkinsonism

**Correspondence to:**

Chuantao Zuo

PET Center, Huashan Hospital, Fudan University, 518 East Wuzhong Road, Shanghai, China

Telephone: +86-21-64280718

FAX: +86-21-64280718

E-mail: [zuochuantao@fudan.edu.cn](mailto:zuochuantao@fudan.edu.cn)

Or Jian Wang

Department of Neurology, Huashan Hospital, Fudan University, 12 Middle Wulumuqi Road, Shanghai,  
China

Telephone: +86-21-52888163

FAX: +86-21-62483421

E-mail: [wangjian\\_hs@fudan.edu.cn](mailto:wangjian_hs@fudan.edu.cn)

**First Author:**

Dr. Ping Wu and Dr. Yu Zhao

PET Center, Huashan Hospital, Fudan University, 518 East Wuzhong Road, Shanghai, China

Telephone: +86-21-64280718

Email: [wupingpet@fudan.edu.cn](mailto:wupingpet@fudan.edu.cn); [yu.zhao@tum.de](mailto:yu.zhao@tum.de)

## ABSTRACT

The clinical presentations of early idiopathic Parkinson's disease (PD) substantially overlap with those of atypical parkinsonian syndromes like multiple system atrophy (MSA) and progressive supranuclear palsy (PSP). This study aimed to develop metabolic imaging indices based on deep learning to support the differential diagnosis of these conditions. **Methods:** A benchmark Huashan parkinsonian PET imaging (HPPI, China) database including 1275 parkinsonian patients and 863 non-parkinsonian subjects with  $^{18}\text{F}$ -FDG PET images was established to support artificial intelligence development. A 3D deep convolutional neural network was developed to extract deep metabolic imaging (DMI) indices, which was blindly evaluated in an independent cohort with longitudinal follow-up from the HPPI, and an external German cohort of 90 parkinsonian patients with different imaging acquisition protocols. **Results:** The proposed DMI indices had less ambiguity space in the differential diagnosis. They achieved sensitivities of 98.1%, 88.5%, and 84.5%, and specificities of 90.0%, 99.2%, and 97.8% for the diagnosis of PD, MSA, and PSP in the blind test cohort. In the German cohort, They resulted in sensitivities of 94.1%, 82.4%, 82.1%, and specificities of 84.0%, 99.9%, 94.1% respectively. Employing the PET scans independently achieved comparable performance to the integration of demographic and clinical information into the DMI indices. **Conclusion:** The DMI indices developed on the HPPI database show potential to provide an early and accurate differential diagnosis for parkinsonism and is robust when dealing with discrepancies between populations and imaging acquisitions.

**Keywords:** Parkinson's disease; atypical parkinsonian syndrome; differential diagnosis; deep learning; deep metabolic imaging indices

## INTRODUCTION

Idiopathic Parkinson's disease (IPD) is one of the most common neurodegenerative disorders. Although extensively studied, its accurate diagnosis remains clinically challenging, particularly in early stage patients, since their symptoms overlap largely with atypical parkinsonian syndromes like multiple system atrophy (MSA) and progressive supranuclear palsy (PSP) (1). Approximately 20-30% patients with initial diagnoses of IPD were subsequently demonstrated to be either MSA or PSP at pathological examination (1). The development of accurate indices for parkinsonism's differential diagnosis is of importance and potential utility when determining therapeutic strategies.

$^{18}\text{F}$ -fluorodeoxyglucose positron emission tomography ( $^{18}\text{F}$ -FDG PET) detects a wide spectrum of neurobiological abnormalities and has been reported of advantage in the differential diagnosis of parkinsonism in advance of structural damage to brain (2). Metabolic patterns of IPD, MSA, and PSP identified by principal component analysis (PCA) (3,4), which were used as features for a machine learning method of logistic regression, have been found as effective surrogates for the early and accurate differential diagnosis (5). However, the PCA decomposition takes the 3D image volume of a subject as a squeezed 1D vector without considering the high-level spatial interrelation during the pattern extraction.

The differences among parkinsonism are reflected in the complex interaction of interrelated brain regions. The differential indices may be obscured by complexity within the metabolic imaging signal. We hypothesized that deep learning may reveal characteristic imaging indices from complex metabolic alterations and provide accurate classifications (6). Therefore, a 3D deep residual convolutional neural network termed PD Diagnosis Network (PDD-Net) was built for the automatic identification of imaging-related indices to support parkinsonism's differential diagnosis.

## MATERIALS AND METHODS

### Subjects and Study Protocol

*Huashan Parkinsonian PET Imaging (HPPI) Database.* A largest and unique HPPI database has been established to benchmark the imaging-based artificial intelligence development for parkinsonism. This database includes three cohorts with a total of 1275 parkinsonian patients (subset of PD Database and Samples Bank of Huashan Hospital) (Fig.1, Supplemental Table 1, Supplemental Table 2) (7-11). Among them, 85.7% performed dopaminergic imaging at the same time as  $^{18}\text{F}$ -FDG to assist the diagnosis and the remaining have been followed up for a long time to determine the diagnosis. A control cohort of 643 patients with various neurological disorders and 220 healthy subjects was also enrolled (Fig.1, Supplemental Table 3, Supplemental Table 4, Supplemental Fig.1 ).

The HPPI database includes *pre-training* (398 subjects with possible diagnoses), *training* (547 subjects with definite diagnoses), and *blind-test* (330 subjects with confirmative diagnoses with follow-up) *cohorts* (Fig.1, Table 1). These patients were routinely assessed by movement disorders specialists in Huashan Hospital before PET examination between June 2011 and April 2019. Routine MRI examinations were performed before PET scans and those with structural brain abnormalities were excluded. After PET examination, patients had at least one return visit and the movement disorders specialists made a clinical diagnosis according to the latest clinical criteria (9-11).

After a low-dose CT for attenuation correction, the emission data was acquired at 60-minute (lasting 10 min) post injection of approximately 185 MBq  $^{18}\text{F}$ -FDG with Biograph 64 HD PET/CT (Siemens, Germany). Following corrections for attenuation, scatter, dead time, and random coincidences, PET images were reconstructed using the ordered subset expectation maximization method.

*German Parkinsonian Cohort.* A German cohort with 34 IPD, 17 MSA and 39 PSP patients from the University Hospital of Munich was included for external validation. These patients were scanned on three different PET/CT systems (ECAT Exact HR+, GE Discovery 690, and Siemens Biograph 64) according to the EANM protocol (12) using a slow bolus injection of approximately 150 MBq  $^{18}\text{F}$ -FDG (Supplemental Table 5). The uptake difference between cohorts are presented in Supplemental Fig. 2.

The institutional review boards (IRB or equivalent from Huashan Hospital and University of Munich) approved this study and all subjects signed a written informed consent.

### **Image Preprocessing**

PET images were spatially normalized into Montreal Neurological Institute brain space and smoothed by a 3D Gaussian filter of 10 mm full width at half maximum by SPM5 software (Institute of Neurology, London, UK). Before inputting the PET image into the deep neural network, Z-score normalization was applied to convert PET image values into a certain range for facilitating the network training. Besides, the performance of utilizing the Z-score normalization and the global mean normalization were also compared (Supplemental Table 6).

### **Parkinsonism Differential Neural Network (PDD-Net) & Deep Metabolic Imaging (DMI) Indices**

The deep learning method contains two PDD-Nets (Supplemental Fig. 3). The PDD-Net-1 sought to exclude patients without parkinsonism. The PDD-Net-2 performed computation of deep metabolic imaging (DMI) indices and classification of IPD, MSA, or PSP. Both PDD-Nets were based on a 3D residual convolutional neural network. The PDD-Net 2 was trained preliminarily in the pre-training cohort, and then fine-tuned in the training cohort. The performance of the DMI indices was evaluated with cross-validation (six-fold) in the training cohort and then an independent test in the blind-test cohort and the external Germany cohort.

At the end of the PDD-Net, the extracted features were mapped to three classification probabilities of IPD, MSA, and PSP correspondingly, which were proposed as the DMI indices. The highest probability among the DMI indices was considered for the prediction of IPD, MSA, or PSP. An additional option of confidence inspection was provided to warn the predictions without sufficiently high probability. A confidence threshold can be customized. By default, a set of confidence thresholds were derived in the cross-validation stage based on the generalized Youden's index. Predictions lying below these thresholds were flagged as uncertain cases (Supplemental Table 7). We generated saliency maps using the full-gradient method (13) to assist the interpretation of the DMI indices. The saliency maps assign importance scores to both the input features and individual neurons in a network, which reflects the contribution of groups of pixels to the DMI probabilities.

### **Statistical Analysis**

The confidence intervals were calculated with DeLong's method (1988). The optimal cutoff points of the receiver operating characteristic curves were estimated using the generalized Youden's index. For continuous variables, the Wilcoxon test was used to compare two paired groups and the Kolmogorov-Smirnov test was used to compare two unpaired groups. While for categorical variables, the Chi-square test was used. Four standard metrics, i.e., sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV), were employed to illustrate the diagnostic performance of the DMI indices.

## RESULTS

### Performance of the DMI Indices in Cross-validation

The performance of the DMI indices in the cross-validation is illustrated in Fig. 2. The area under the curves were 0.986, 0.997, and 0.982 for IPD, MSA and PSP, respectively. The sensitivity, specificity, PPV, and NPV are summarized in Table 2, and all values were above 90% except for sensitivity and PPV for PSP with short symptom durations. Compared to those with short symptom durations, the specificity for those with long symptom durations slightly increased for IPD and MSA, while remained the same for PSP.

The probabilities of IPD, MSA, and PSP according to the DMI indices for individual subjects are plotted in 3D coordination space in Fig. 3. These probabilities tended to distribute aggregately to their expected centers: IPD for [1,0,0], MSA for [0,1,0], and PSP for [0,0,1]. If the probability for a category was high, the probabilities for the other two categories were much smaller. The aggregation distance, which is the mean distance of the probabilities to the corresponding expected centers, illustrates the determinability of the DMI indices. The probabilities of those with long symptom durations (aggregative distance=0.103) were more aggregated ( $P=0.020$ ) compared to the subjects with short symptom durations (aggregative distance=0.114). Overall, the probabilities among the DMI indices had less ambiguity space for differential diagnosis.

The saliency maps are showed in Supplemental Fig. 4-6 (13). Regions with relatively higher contribution to the DMI indices were putamen and midbrain for IPD, MSA, and PSP as well as cerebellum for MSA.



## **Performance of the DMI Indices in the Blind Test**

Table 3 illustrates the predictive accuracy of the DMI indices in the blind-test cohort. The image-based classification resulted in 98.1% sensitivity, 90.0% specificity, 94.5% PPV, and 96.4% NPV for PD and also accurate for MSA (88.5% sensitivity, 99.2% specificity, 96.4% PPV, and 97.4% NPV) and PSP (84.5% sensitivity, 97.8% specificity, 89.1% PPV, and 97.0% NPV). For the 108 patients in the blind-test cohort with follow-up PET scans, the DMI indices had slightly better performance comparing follow-up to baseline (P=0.017).

The probabilities among the DMI indices for subjects with follow-up imaging in the blind-test cohort are plotted in Fig. 4. The probabilities of MSA and PSP increased at follow-up imaging (MSA: P=0.028, PSP: P=0.002). The probabilities of IPD between at follow-up and baseline imaging were comparable (P=0.894), but the median and most of the IPD probabilities (38/66) increased. Nine cases presented relative significant lower probabilities of IPD at follow-up (over 0.1) compared with the baseline.

Besides, differential diagnosis performance of using the DMI indices only and using the combination of the DMI indices with demographic and clinical features were compared, and no difference was found (P=0.999) (Supplemental Table 8) (14). Besides, DMI indices made predictions inconsistent with the clinical diagnosis in six cases obvious probability decrease during follow-up (Supplemental Table 9).

## **Test on the External German Cohort**

The DMI indices achieved 94.1% sensitivity, 84.0% specificity, 78.0% PPV, and 95.9% NPV for the diagnosis of the IPD on the German cohort (Table 3). The diagnoses were also accurate for MSA (82.4% sensitivity, 99.9% specificity, 99.9% PPV, and 96.1% NPV) and PSP (82.1% sensitivity, 94.1% specificity, 91.4% PPV, and 87.3% NPV). Although the performance metrics were slightly lower than

those for the Chinese cohort, no significant difference has been observed in the performance of the diagnosis of IPD (P=0.14), MSA (P=0.25) and PSP (P=0.50).

## **DISCUSSION**

An effective imaging-based tool may contribute to earlier and more precise diagnosis in parkinsonian conditions and may help with the development and monitoring of individualized disease-modifying treatments (15,16). This study confirms that deep learning can identify accurate imaging-based indices from  $^{18}\text{F}$ -FDG PET.

Similar to pattern expression scores of PCA analysis (5), the DMI indices herein identified three probability scores from  $^{18}\text{F}$ -FDG PET for each individual and a prediction was generated by comparing these three probabilities. The conventional pattern related scores are derived from linear weightings of imaging intensities. In contrast, the DMI indices can reveal hyper-level inter-relations such as textures, which may better describe the complex heterogeneous pathogenesis of parkinsonian disorders. The extensive test in relatively large cohorts found that the DMI indices can achieve competitive or possibly better performance in the differential diagnosis of parkinsonism compared with previously reported studies (5).

The probabilities among the DMI indices have low ambiguity and a dominant maximal probability is definable for resulting in a robust diagnosis prediction. Nevertheless, we also support confidence inspection to differentiate predictions with different confidence levels. The confidence thresholds can be customized (Supplemental Table 7). For a default setting according to the optimization of generalized Youden's index, the confidence threshold for MSA was higher compared to IPD or PSP. In this study, the MSA patients were mixed with MSA-parkinsonian (MSA-P) and MSA-cerebellar (MSA-C) types and

have greater heterogeneity in metabolic pathological phenotype. Therefore, it could be posited that a higher confidence threshold is required to obtain a robust prediction.

The DMI indices can be combined with demographic and clinical information as well as other indices, such as impairment of olfactory function (for IPD versus MSA) or skin biopsy-positivity for phospho-alpha-synuclein aggregates (17) (for IPD and MSA versus PSP), to comprehensively generate diagnostic classifications. In our study, employing the PET scans independently achieved comparable performance to the integration of demographic and clinical information into the DMI indices, which indicated the most discriminative information for the parkinsonism diagnosis was included in the PET scan modality and could be extracted by the proposed method into the DMI indices. Besides, the two-stage design (Supplemental Fig. 3) (13,18,19) of our work allows the DMI indices to reduce the risk of erroneous predictions through excluding non-parkinsonian subjects in the control stage, which aims at further improving the robustness of diagnostic classifications.

In general, the DMI indices developed from the Chinese HPPI database achieved comparable performance in a German cohort. Indeed, there were substantial differences between the two cohorts: in contrast to the Chinese cohorts, the German cohorts were acquired on different scanners. The imaging protocols (i.e. acquisition time, reconstruction method, tracer dose) and patient preparation (i.e. eye patch and noise-cancelling differences) (Supplemental Table 5) varied. Significantly different metabolic uptakes were observed in the cerebellum, midbrain and caudate between these two cohorts (Supplemental Fig. 2), where population-based differences (3,20) may exist. The domain difference between data can present an obstacle to the wider clinical translation of conventional methods. A prerequisite for spatial covariance analysis in the established population-based patterns for IPD, MSA, and PSP is to bridge the difference between various populations (5). In contrast to pattern analysis, the hierarchical feature representation of deep learning is more flexible and affords migration of domain differences during the learning phase (21).

Similar to previous studies (22), our test confirmed that deep learning can be robust to the discrepancies inherent in molecular imaging acquisitions. This finding suggests the DMI marker extracted using deep learning in this study may be more generalizable and better suited for clinical translation.

Recently, concerns have been raised regarding the reproducibility or stability of deep learning methods: methods optimized in one cohort may have limited performance in other cohorts or in other applications (23). We subjected our DMI indices to a blind test as a means of independent in-depth validation (24). The performance of the DMI indices under conditions of a blind test was consistent with the cross-validation test. These results lead us to conclude that the DMI indices are reproducible. Another limitation of deep learning is the black-box nature of the derived model, which precludes the drawing of any links to the underlying pathophysiology. To address this concern, we employed saliency maps to understand the decision mechanism behind the neural networks. The saliency maps indicated that the DMI indices derived probabilities largely based on parkinsonism-related brain regions, which are consistent with the critical regions of IPD, MSA, and PSP-related covariance pattern (5,25).

Dopaminergic imaging is critical for diagnosing parkinsonian disorders, although it has not been confirmed to be suitable for the reliable differential diagnosis. Most patients with parkinsonism in our study underwent contemporary dopaminergic imaging as  $^{18}\text{F}$ -FDG. Therefore, this study can be regarded as performed based on dopaminergic imaging. Whether  $^{18}\text{F}$ -FDG imaging and deep learning can be used to diagnose parkinsonian disorders with blinded dopaminergic imaging results is an interesting future direction to explore.

One limitation of this study is that we did not employ MRI for partial volume correction and spatial normalization. Although MRI is generally included in the neurological work-up of these patients, many of them were scanned at external centers with a variety of protocols and the 3D images were not always retrievable. We conceded that the cortical thickness derived from MRI images might also assist the

differentiation of parkinsonism (26). The integration of these morphometries in any future study may further enhance the imaging-based indices. Besides, although performance on the training cohort, blind-test cohort and Germany cohort, which have different data distributions (IPD:MSA:PSP), has indicated the DMI indices have a certain level of ability to handle the distribution-different problems, different distributions may still be a factor influencing performance on another future cohort. It is worthy to conduct multi-center studies to further validate our method. Meanwhile, we only evaluated one possible multi-modality fusion method in this work. In the future, to further improve the diagnosis performance, other fusion methods such as gating-based attention-based late fusion will be evaluated.

## **CONCLUSION**

We developed a 3D deep residual convolutional neural network to extract DMI indices for the automated differential diagnosis of parkinsonism. The indices were evaluated with the cross-validation experiment and blind tests on both Chinese and German cohorts, demonstrating that the proposed method was both robust and accurate, which may complement diagnoses made by expert clinicians.

## **ACKNOWLEDGEMENTS**

This work was supported by the National Natural Science Foundation of China (81771483, 81671239, 81361120393, 81401135, 81971641, 81902282, 91949118, 81771372), the Ministry of Science and Technology of China (2016YFC1306504), Shanghai Municipal Science and Technology Major Project (No. 2017SHZDZX01, 2018SHZDZX03) and ZJ Lab, Youth Medical Talents - Medical Imaging Practitioner Program by Shanghai Municipal Health Commission and Shanghai Medical and Health Development Foundation (SHWRS(2020)\_087), Shanghai Sailing Program by Shanghai Science and Technology Committee (18YF1403100), the Swiss National Science Foundation (188350), Jacques & Gloria Gossweiler Foundation and Siemens Healthineers.

## **DISCLOSURE**

W.H.O is Hertie Senior Research Professor, supported by the Charitable Hertie Foundation, Frankfurt/Main, Germany. A.R. and K.S. received research support from Novartis and Siemens Healthineers. Other authors report no financial interests or potential conflicts of interest.

## **KEY POINTS**

**QUESTION:** Can deep learning effectively extract indices from brain glucose metabolic imaging ( $^{18}\text{F}$ -FDG PET) to improve the differential diagnosis of Parkinson's disease and atypical parkinsonian syndromes?

**PERTINENT FINDINGS:** The developed deep metabolic imaging (DMI) indices prediction using deep learning provides an early and accurate method for differential diagnosis which may complement diagnoses made by expert clinicians. The trustworthy artificial intelligence (AI) development was achieved by training on a largest benchmark data of  $^{18}\text{F}$ -FDG PET, extensive testing on longitudinal data and independent external data with different ethnicity or examination protocols.

**IMPLICATIONS FOR PATIENT CARE:** This developed DMI indices may assist early differential diagnosis of parkinsonism and the development of disease-modifying treatment strategies.

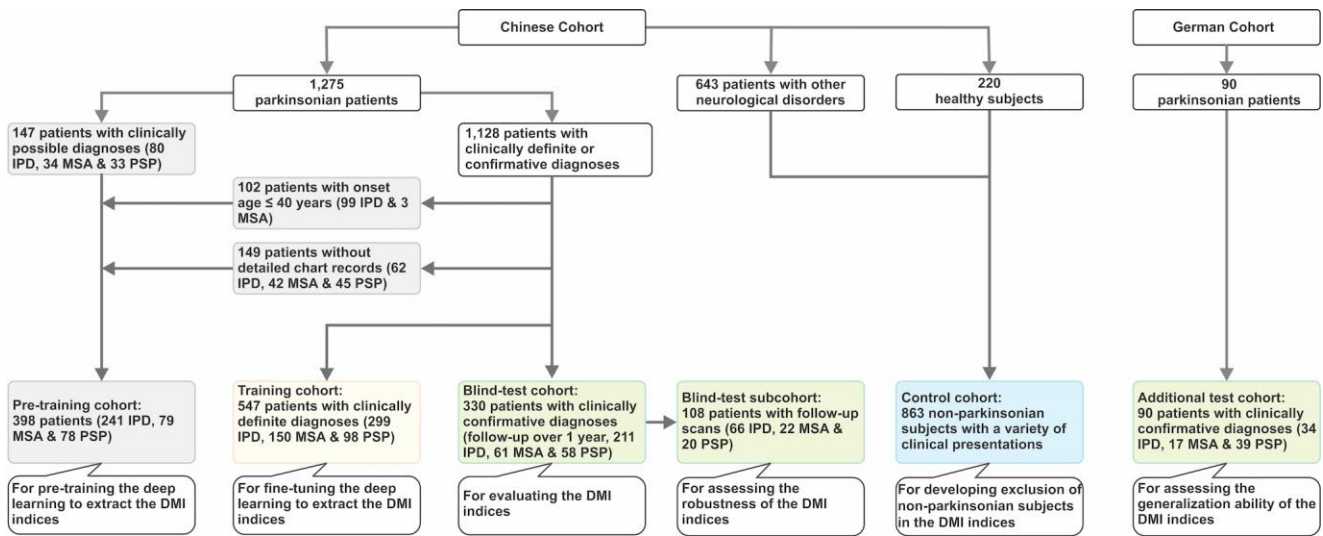
## REFERENCES

1. Hughes AJ, Daniel SE, Ben-Shlomo Y, Lees AJ. The accuracy of diagnosis of parkinsonian syndromes in a specialist movement disorder service. *Brain*. 2002;125:861-870.
2. Stoessel AJ, Lehericy S, Strafella AP. Imaging insights into basal ganglia function, Parkinson's disease, and dystonia. *Lancet*. 2014;384:532-544.
3. Ge J, Wu J, Peng S, et al. Reproducible network and regional topographies of abnormal glucose metabolism associated with progressive supranuclear palsy: Multivariate and univariate analyses in American and Chinese patient cohorts. *Hum Brain Mapp*. 2018;39:2842-2858.
4. Eckert T, Tang C, Ma Y, et al. Abnormal metabolic networks in atypical parkinsonism. *Mov Disord*. 2008;23:727-733.
5. Tang CC, Poston KL, Eckert T, et al. Differential diagnosis of parkinsonism: a metabolic imaging study using pattern analysis. *Lancet Neurol*. 2010;9:149-158.
6. Choi BW, Kang S, Kim HW, Kwon OD, Vu HD, Youn SW. Faster Region-Based Convolutional Neural Network in the Classification of Different Parkinsonism Patterns of the Striatum on Maximum Intensity Projection Images of [<sup>18</sup>F] FP-CIT Positron Emission Tomography. *Diagnostics*. 2021;11:1557.
7. Litvan I, Agid Y, Calne D, et al. Clinical research criteria for the diagnosis of progressive supranuclear palsy (Steele-Richardson-Olszewski syndrome) report of the NINDS-SPSP international workshop. *Neurology*. 1996;47:1-9.
8. Hughes AJ, Daniel SE, Kilford L, Lees AJ. Accuracy of clinical diagnosis of idiopathic Parkinson's disease: a clinico-pathological study of 100 cases. *J Neurol Neurosurg Psychiatry*. 1992;55:181-184.
9. Gilman S, Wenning GK, Low PA, et al. Second consensus statement on the diagnosis of multiple system atrophy. *Neurology*. 2008;71:670-676.
10. Höglinger GU, Respondek G, Stamelou M, et al. Clinical diagnosis of progressive supranuclear palsy: the movement disorder society criteria. *Mov Disord*. 2017;32:853-864.
11. Postuma RB, Berg D, Stern M, et al. MDS clinical diagnostic criteria for Parkinson's disease. *Mov Disord*. 2015;30:1591-1601.

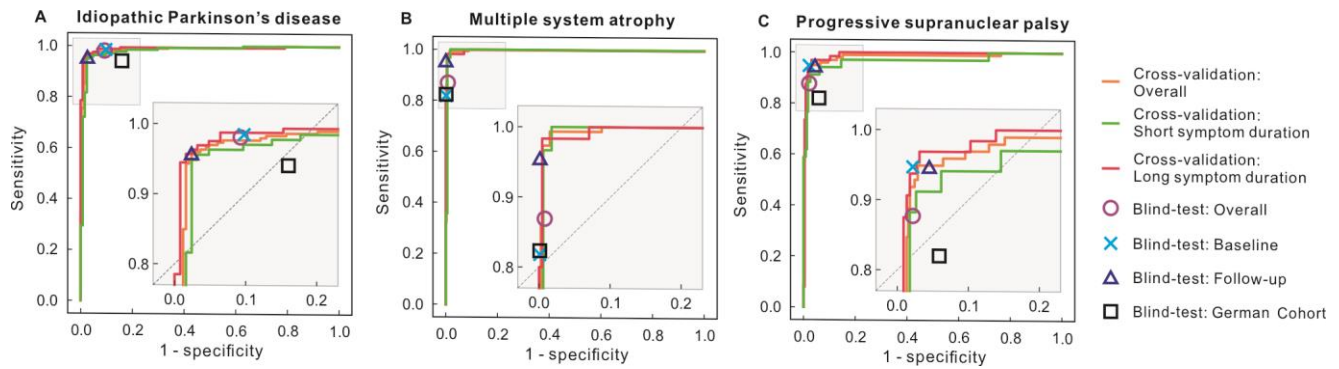


12. Varrone A, Asenbaum S, Vander Borgh T, et al. EANM procedure guidelines for PET brain imaging using [18F]FDG, version 2. *Eur J Nucl Med Mol Imaging*. 2009;36:2103-2110.
13. Srinivas S, Fleuret F. Full-gradient representation for neural network visualization. In: *The 33rd International Conference on Neural Information Processing Systems*. 2019:4126-4135.
14. Chen T, Guestrin C. Xgboost: A scalable tree boosting system. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 2016:785-794.
15. Strafella AP, Bohnen NI, Perlmutter JS, et al. Molecular imaging to track Parkinson's disease and atypical parkinsonisms: New imaging frontiers. *Mov Disord*. 2017;32:181-192.
16. Meles SK, Teune LK, de Jong BM, Dierckx RA, Leenders KL. Metabolic Imaging in Parkinson Disease. *J Nucl Med*. 2017;58:23-28.
17. Devine MJ, Gwinn K, Singleton A, Hardy J. Parkinson's disease and  $\alpha$ -synuclein expression. *Mov Disord*. 2011;26:2160-2168.
18. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition*. 2016:770-778.
19. Skrede O-J, De Raedt S, Kleppe A, et al. Deep learning for prediction of colorectal cancer outcome: a discovery and validation study. *Lancet*. 2020;395:350-360.
20. Shi L, Liang P, Luo Y, et al. Using large-scale statistical Chinese brain template (Chinese2020) in popular neuroimage analysis toolkits. *Front Hum Neurosci*. 2017;11:414.
21. Ghafoorian M, Mehrtash A, Kapur T, et al. Transfer Learning for Domain Adaptation in MRI: Application in Brain Lesion Segmentation. In: *2017 International Conference on Medical Image Computing and Computer-assisted Intervention*. 2017:516-524.
22. Wenzel M, Milletari F, Krüger J, et al. Automatic classification of dopamine transporter SPECT: deep convolutional neural networks can be trained to be robust with respect to variable image characteristics. *Eur J Nucl Med Mol Imaging*. 2019;46:2800-2811.
23. Maier-Hein L, Eisenmann M, Reinke A, et al. Why rankings of biomedical image analysis competitions should be interpreted with care. *Nat Commun*. 2018;9:1-13.

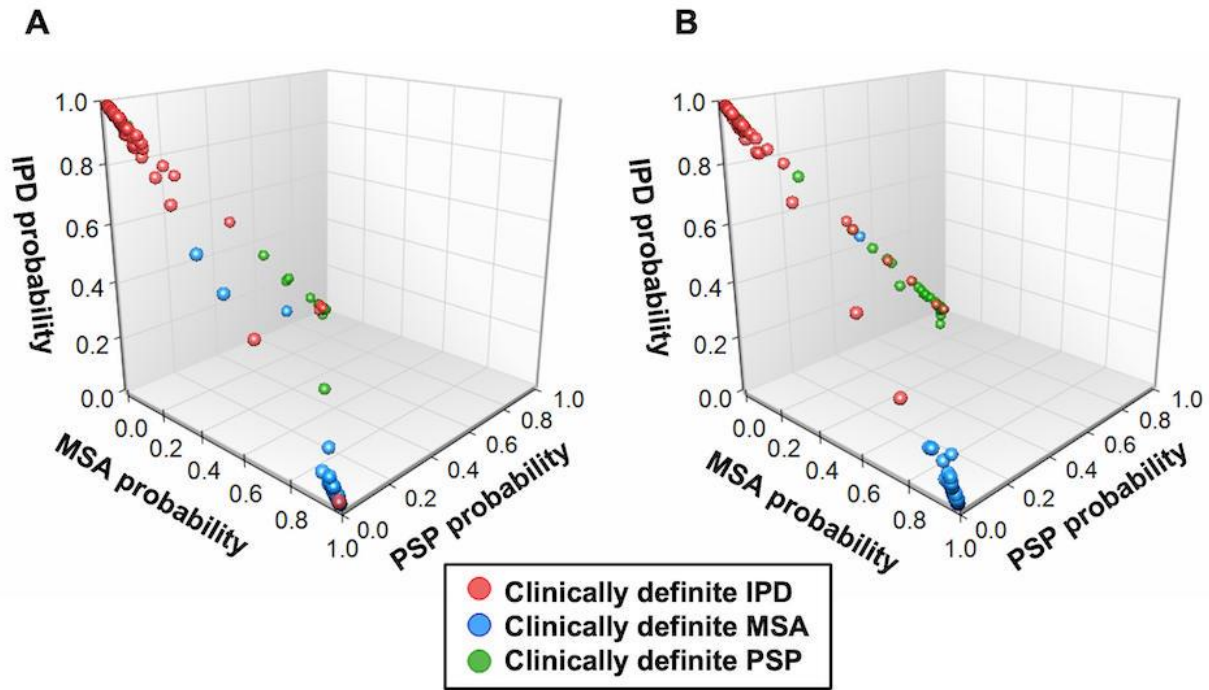
- 24.** Segler MH, Preuss M, Waller MP. Planning chemical syntheses with deep neural networks and symbolic AI. *Nature*. 2018;555:604-610.
- 25.** Matthews DC, Lerman H, Lukic A, et al. FDG PET Parkinson's disease-related pattern as a biomarker for clinical trials in early stage disease. *Neuroimage Clin*. 2018;20:572-579.
- 26.** Möller L, Kassubek J, Südmeyer M, et al. Manual MRI morphometry in Parkinsonian syndromes. *Mov Disord*. 2017;32:778-782.



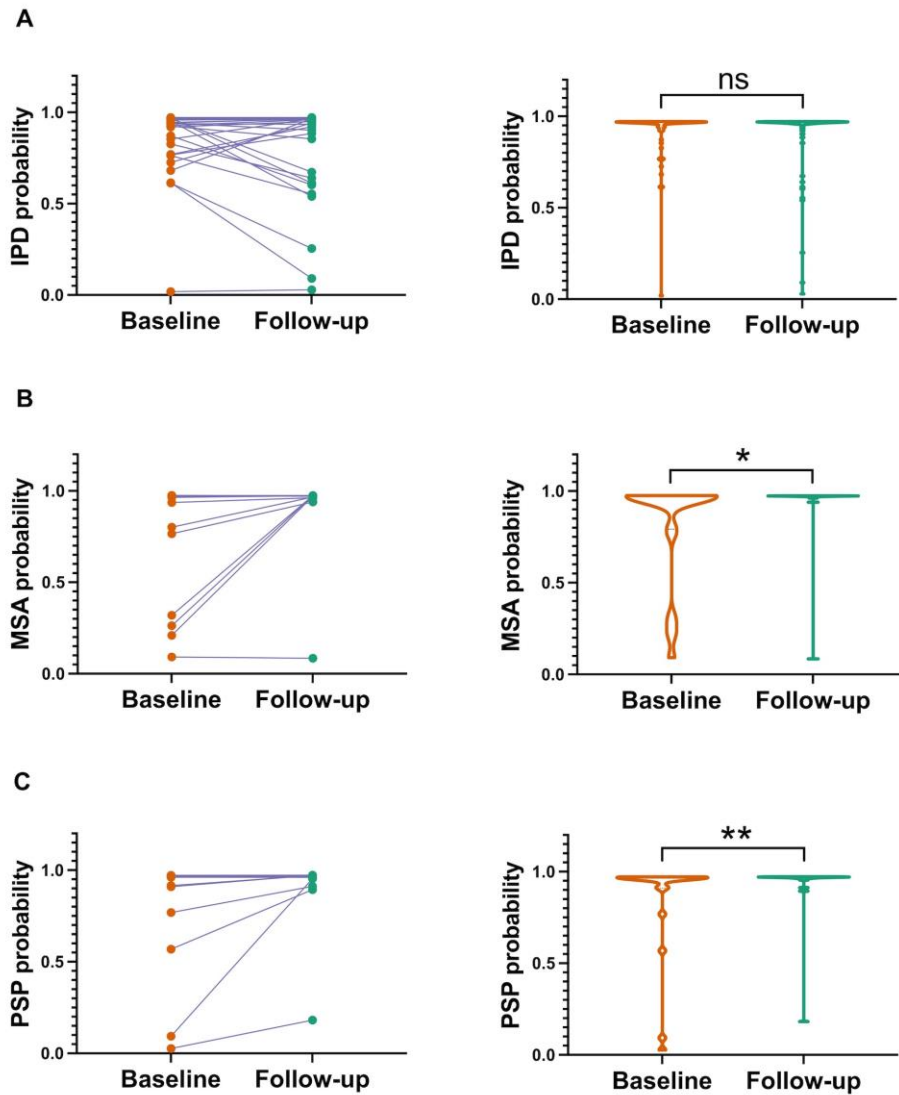
**Figure 1. Study profile.** A Chinese cohort (Huashan parkinsonian PET imaging dataset) and A Germany cohort were involved. IPD: idiopathic Parkinson's disease, MSA: multiple system atrophy, PSP: progressive supranuclear palsy, DMI: deep metabolic imaging, Clinically definite diagnoses: diagnoses by the clinical experts after return visit but without a formal clinical follow-up, Clinically confirmative diagnoses: diagnoses resulting from at least one formal clinical follow-up over one year after PET imaging.



**Figure 2. The accuracy of the deep metabolic imaging (DMI) indices in the development phase in the training cohort and blind-test phase on both Chinese and German test cohorts. The results in the cross-validation were plotted using receiver operating characteristic curves. The results in the Chinese blind-test cohort were illustrated as single points, where Overall represents the results of all the tested 330 patients. 108 patients in the blind test have follow-up scans and the performance of them at Baseline and Follow-up was plotted. The blind-test results in the German cohort (90 patients) are also included and denoted with the black rectangular for easy comparison.**

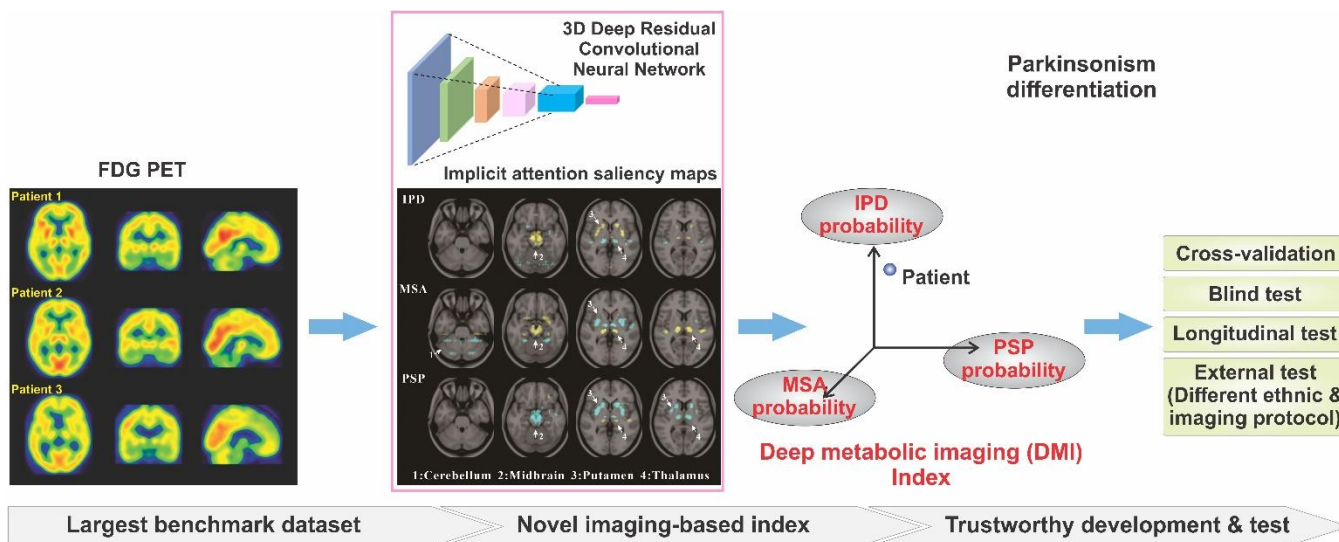


**Figure 3. 3D plot of the probabilities of idiopathic Parkinson's disease (IPD), multiple system atrophy (MSA) and progressive supranuclear palsy (PSP) of the deep metabolic imaging (DMI) indices in the training cohort. (A) patients with short symptom duration ( $\leq 2$  years), (B) patients with long symptom duration ( $> 2$  years).**



**Figure 4. The comparison of the probabilities of idiopathic Parkinson’s disease (IPD), multiple system atrophy (MSA), and progressive supranuclear palsy (PSP) of the deep metabolic imaging (DMI) indices on the 108 patients of the blind-test cohort with repeated PET scans. The left column compares the probability of the DMI indices extracted from baseline and follow-up PET for individuals. The violin plots in the right column demonstrate the statistical distribution of the probabilities of the DMI indices. \*:  $P \leq 0.05$ , \*\*:  $P \leq 0.01$ , ns: no significance.**

# Graphical Abstract



**Table 1. The demographic and clinical data.**

		Huashan parkinsonian PET imaging dataset (Chinese cohort)							German Cohort
		Pre-training Cohort	Training Cohort			Blind-test Cohort			
			Overall	Short Symptom Duration ( ≤ 2 years)	Long Symptom Duration ( > 2 years)	Overall	Baseline	Follow-up	
Idiopathic Parkinson Disease*	Patient number	241	299	136	163	211	66	66	34
	Sex (male/female)	154/87	166/133	73/63	93/70	130/81	43/23	43/23	21/13 (34/34)
	Age at PET (years)	50.0±15.5	60.2±8.5	59.1±9.0	61.0±8.0	60.0±7.6	60.0±7.9	62.1±7.9	72.9±9.5 (34/34)
	Symptom duration at PET (months)	/	45.3±46.0	13.0±5.9	72.3±47.4	39.0±41.3	26.0±24.1	53.4±24.2	44.5±32.9 (18/34)
	Hoehn and Yahr stage**	/	2.2±1.0	1.7±0.6	2.7±1.0	1.9±0.9	1.6±0.7	1.9±0.6	1.6±0.8 (22/34)
	UPDRS III	/	27.0±14.3	18.9±8.9	33.8±14.5	22.8±12.1	19.6±9.1	24.2±10.1	12.0±3.6 (3/34)
	Clinical follow-up (months)	/	/	/	/	46.8±30.4	/	64.5±25.3	19.1±21.8 (14/34)
Multiple System Atrophy	Patient number (MSA-C/MSA-P)	79	150 (57/93)	90 (39/51)	60 (18/42)	61 (21/40)	22 (8/14)	22 (8/14)	17 (8/8/1)
	Sex (male/female)	42/37	78/72	47/43	31/29	32/29	14/8	14/8	10/7 (17/17)
	Age at PET (years)	57.5±10.6	57.8±8.0	56.5±8.1	59.6±7.4	58.5±6.3	58.3±7.4	60.3±7.3	61.3±8.3 (17/17)
	Symptom duration at PET (months)	/	24.3±17.1	13.9±6.0	39.9±16.5	27.0±20.1	22.1±11.8	45.6±12.5	30.0±22.2 (17/17)
	Hoehn and Yahr stage**	/	3.1±0.8	3.0±0.8	3.5±0.7	2.9±0.8	2.6±0.6	3.4±0.8	2.4±1.1 (15/17)
	UPDRS III	/	30.6±14.5	25.9±12.4	37.6±14.7	29.3±14.4	23.5±8.2	36.4±11.1	34.6±12.8 (11/17)
	Clinical follow-up (months)	/	/	/	/	30.7±18.2	/	41.7±16.4	22.6±22.4 (17/17)
Progressive Supranuclear Palsy*	Patient number	78	98	34	64	58	20	20	39
	Sex (male/female)	45/33	60/38	23/11	37/27	39/19	17/3	17/3	21/18 (39/39)
	Age at PET (years)	64.6±8.6	67.2±8.0	65.0±9.3	68.5±6.9	65.1±6.6	64.8±7.5	67.0±7.2	70.0±7.1 (39/39)
	Symptom duration at PET (months)	/	35.0±20.7	15.3±5.4	45.5±18.0	34.1±22.7	32.4±22.0	58.8±22.8	22.4±15.7 (37/39)
	Hoehn and Yahr stage**	/	3.2±0.8	2.9±0.6	3.4±0.8	3.0±0.8	2.7±1.0	3.6±0.8	2.6±1.1 (37/39)
	UPDRS III	/	30.1±13.5	28.0±11.0	31.2±14.6	26.8±11.0	23.0±10.4	34.6±15.9	37.0±15.9 (20/39)
	Clinical follow-up (months)	/	/	/	/	25.1±15.7	/	37.5±12.9	22.2±13.8 (17/39)

Data are shown as mean ± standard deviation. In German cohort, the associated numbers of subjects with these items are provided together with the statistics information (subject number with certain item /total subject number).

\* Diagnosis information: Supplementary table 1

\*\*Detailed Hoehn and Yahr stage information: Supplementary table 2



**Table 2. Accuracy of the deep metabolic imaging (DMI) indices in the cross-validation on the training cohort.**

		Overall	Short Symptom Durations (≤ 2 years)	Long Symptom Durations (> 2 years)
Idiopathic Parkinson Disease	AUC	0.986 (0.977-0.996)	0.981 (0.965-0.997)	0.991 (0.981-1.000)
	Sensitivity	95.7% (92.7%-97.7%)	94.9% (89.7%-97.9%)	95.7% (91.4%-98.3%)
	Specificity	97.6% (94.8%-99.1%)	97.6% (93.1%-99.5%)	98.4% (94.3%-99.8%)
	PPV	97.9% (95.6%-98.9%)	97.7% (93.5%-99.1%)	98.7% (95.5%-99.5%)
	NPV	94.9% (91.5%-98.1%)	94.5% (89.1%-98.8%)	94.6% (89.2%-99.3%)
Multiple System Atrophy	AUC	0.997 (0.994-1.000)	0.996 (0.988-1.000)	0.998 (0.995-1.000)
	Sensitivity	97.3% (93.3%-99.3%)	100% (96.0%-100%)	98.3% (91.1%-100%)
	Specificity	99.5% (98.2%-99.9%)	98.2% (94.9%-99.6%)	99.6% (97.6%-100%)
	PPV	98.6% (95.3%-99.6%)	96.8% (91.0%-100%)	98.3% (91.3%-100%)
	NPV	99.0% (97.4%-99.9%)	100% (97.8%-100%)	99.6% (97.5%-100%)
Progressive Supranuclear Palsy	AUC	0.982 (0.965-0.998)	0.968 (0.925-1.000)	0.990 (0.980-1.000)
	Sensitivity	91.8% (84.5%-96.4%)	88.2% (72.5%-96.7%)	93.8% (84.8%-98.3%)
	Specificity	98.2% (96.5%-99.2%)	98.2% (95.5%-99.5%)	98.2% (95.5%-99.5%)
	PPV	91.8% (85.0%-96.4%)	88.2% (74.3%-96.7%)	93.7% (85.2%-98.3%)
	NPV	98.2% (96.4%-99.2%)	98.2% (95.1%-99.5%)	98.2% (95.3%-99.5%)

AUC: area under the curve, PPV: positive predictive value, NPV: negative predictive value.

**Table 3. Accuracy of the deep metabolic imaging (DMI) indices on the blind-test cohort from Huashan parkinsonian PET imaging dataset (Chinese cohort) and German cohort.**

		Huashan parkinsonian PET imaging dataset (Chinese cohort)			German cohort
		Overall	Baseline	Follow-up	
Idiopathic Parkinson Disease	Sensitivity	98.1%	98.5%	95.5%	94.1%
	Specificity	90.0%	88.1%	97.6%	84.0%
	PPV	94.5%	92.9%	98.4%	78.0%
	NPV	96.4%	97.4%	93.2%	95.9%
Multiple System Atrophy	Sensitivity	88.5%	81.8%	95.4%	82.4%
	Specificity	99.2%	99.9%	98.8%	99.9%
	PPV	96.4%	99.9%	95.5%	99.9%
	NPV	97.4%	95.6%	98.8%	96.1%
Progressive Supranuclear Palsy	Sensitivity	84.5%	90.0%	95.0%	82.1%
	Specificity	97.8%	97.7%	96.6%	94.1%
	PPV	89.1%	90.0%	86.4%	91.4%
	NPV	97.0%	97.7%	98.8%	87.3%

PPV: positive predictive value, NPV: negative predictive value

## Detailed information of the Chinese Cohort

In the Chinese cohort (Huashan parkinsonian PET imaging dataset), a total of 1275 parkinsonian patients were included. These patients were sorted into pre-training cohort, training cohort, and blind-test cohort according to whether their diagnosis was clinically definite and whether follow up clinical data (at least one year following PET imaging) were available. *Pre-training cohort* (241 IPD, 79 MSA, and 78 PSP): the patients with a clinically possible diagnosis of IPD, MSA, or PSP were used for preliminary training the PDD-Net. Considering that the purpose of this study was to obtain deep metabolic imaging (DMI) indices and make differential diagnoses of IPD, MSA, and PSP, and cognizant that the diagnostic standards of MSA and PSP have clear provisions on the age of onset, all the patients with an onset age younger than 40 years old were sorted into the pre-training cohort. In addition, the patients having definite clinical diagnosis but without detailed chart records were also grouped into the pre-training cohort. *Training cohort* (299 IPD, 150 MSA, and 98 PSP): the patients with a clinically definite diagnosis after return visit but without a formal clinical follow-up were used for fine-tuning and cross-validation of the PDD-Net to extract DMI indices. We distinguished between two subgroups of patients with short ( $\leq 2$  years) and long ( $> 2$  years) symptom duration for the test. *Blind-test cohort* (211 IPD, 61 MSA, and 58 PSP): the patients with a clinically confirmative diagnosis resulting from at least one formal clinical follow-up over one year after PET imaging were used for independently testing the DMI indices. The diagnosis of the individuals in the blind-test cohort was not disclosed to the algorithm developers who were blinded to clinical details. In the blind-test cohort, a subgroup of 108 patients had another PET scans at the time of follow-up in addition to the one at the time of first diagnosis (baseline). In this work, we denote FDG PET images at baseline of the all 330 patients on the blind-test cohort as “overall”, FDG PET images at baseline of the 108 patients with repeated PET scans as “baseline”, and FDG PET images at follow-up of the 108 patients with repeated PET scans as “follow-up” during analyzing the blind-test cohort.

The clinical diagnosis of the patients in this study was according to the most recently published criteria (1-3). The diagnoses for idiopathic Parkinson’s disease (IPD) and progressive supranuclear palsy (PSP) made using the older criteria (4,5) in the training cohort and blind-test cohort were reconfirmed according to chart records or follow up using the latest criteria (1,2). The detailed information of the diagnosis according to different criteria are listed in Supplemental Table 1.

**Supplemental Table 1** The detailed information of the clinical diagnosis according to different versions of diagnostic criteria.

	Clinical Criteria	Pre-training Cohort	Training Cohort	Blind-test Cohort
<b>Idiopathic Parkinson Disease</b>	New <sup>(2,5)</sup>	112	185	77
	Old <sup>(5)</sup>	129	114	134
<b>Progressive Supranuclear Palsy*</b>	New <sup>(1)(4)</sup>	36	66	29
	Old <sup>(4)</sup>	42	32	29

\*PSP consists of 165 PSP-Richardson syndrome (PSP-RS) and 69 other subtypes

Note: All patients diagnosed with old criteria were reconfirmed with the new diagnosis criteria.

**Supplemental Table 2** The frequency distributions of Hoehn and Yahr stage

		Training Cohort <sup>1</sup>			Blind-test Cohort <sup>2</sup>		
		Overall	Short Symptom Duration	Long Symptom Duration	Overall	Baseline	Follow-up
Idiopathic Parkinson Disease	HY = 1	23.4%	36.0%	12.9%	37.3%	52.2%	23.9%
	HY = 2	46.2%	59.6%	35.0%	42.5%	34.3%	58.2%
	HY = 3	18.7%	4.4%	30.7%	14.2%	11.9%	17.9%
	HY = 4	9.4%	0.0%	17.2%	5.2%	1.5%	0.0%
	HY = 5	2.3%	0.0%	4.3%	0.9%	0.0%	0.0%
Multiple System Atrophy	HY = 1	2.0%	3.3%	0.0%	4.9%	4.6%	0.0%
	HY = 2	14.0%	22.2%	1.7%	23.0%	31.8%	4.6%
	HY = 3	58.0%	60.0%	55.0%	54.1%	63.6%	63.6%
	HY = 4	20.0%	12.2%	31.7%	14.8%	0.0%	18.2%
	HY = 5	6.0%	2.2%	11.7%	3.3%	0.0%	13.6%
Progressive Supranuclear Palsy	HY = 1	2.0%	2.9%	1.6%	3.5%	10.5%	0.0%
	HY = 2	7.1%	14.7%	3.1%	19.3%	31.6%	0.0%
	HY = 3	66.3%	73.5%	62.5%	57.9%	42.1%	63.2%
	HY = 4	17.4%	8.8%	21.9%	14.0%	10.5%	15.8%
	HY = 5	7.1%	0.0%	10.9%	5.3%	5.3%	21.1%

<sup>1</sup> The training cohort includes 547 patents with clinically definite diagnosis according to latest diagnostic criteria for the fine-tuning of the pre-trained deep neural network and the evaluation (cross-validation) during the development of the deep metabolic imaging (DMI) indices. Short symptom duration represents patients with symptom duration  $\leq 2$  years and long symptom duration means patients with symptom duration  $> 2$  years

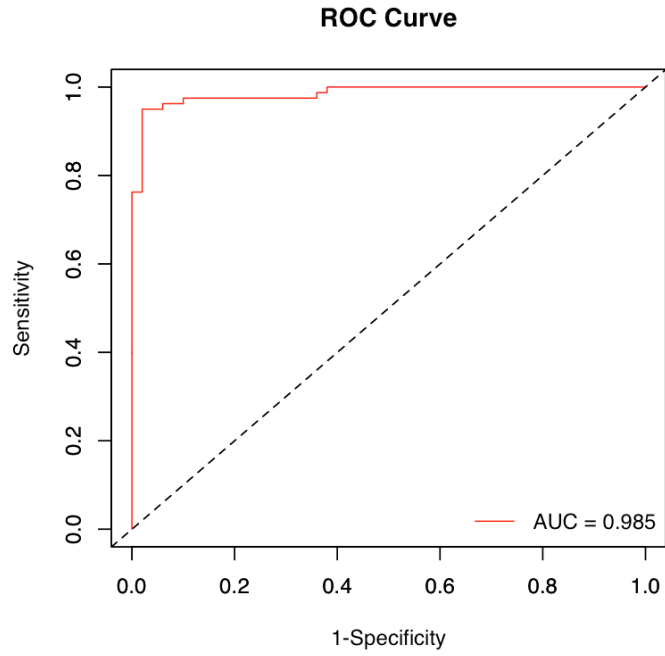
<sup>2</sup> The blind-test cohort includes 330 patients with clinically confirmative diagnosis after follow-up for independent and in-depth test of the developed deep metabolic imaging (DMI) indices. Among them, 108 patients had PET scans both at the time of first diagnosis (Baseline) and also at the time of follow up (Follow-up). During the development, the diagnosis of a patient is modified according to latest follow-up triggered out by the controversial diagnostic recommendation (PSP) compared to the diagnosis at the previous follow-up (IPD).

### **Excluding of non-parkinsonian patients**

The developed deep metabolic imaging (DMI) indices includes the option to pre-investigate the input PET images to avoid erroneous inclusion of non-parkinsonism subjects when calculating DMI indices. There was a control stage to exclude non-parkinsonism patients before the main classification stage. In this stage, a deep neural network was employed for the pre-investigation. Patients with either IPD, MSA, or PSP in pre-training cohort and training cohort were used as “positive” samples to train the network. A control cohort with 643 patients and 220 healthy subjects was collected (The detailed information of the control cohort is given in Supplemental Table 3), of whom 813 were randomly selected as “negative” samples to train the network and the remaining 50 patients were for testing. The performance of the control stage was then tested on 130 unseen patients (parkinsonian subjects: 80, non-parkinsonian subjects (including healthy people): 50). The network achieved ROCAUC of 0.985, sensitivity of 95.0%, specificity of 98.0%, PPV of 98.7%, and NPV of 92.5% for the exclusion of non-parkinsonian subjects (Supplemental Table 4).

**Supplemental Table 3** Control cohort to prevent the inappropriate computation of the DMI indices. In this stage, we trained a network to exclude non-parkinsonian subjects. Patients with IPD/MSA/PSP in pre-training cohort and training cohort were used as “positive samples” to train the network. Patients and healthy subjects in this table were used as “negative samples” to train the network.

<b>Disease Name</b>	<b>Number of Patients</b>
Alzheimer's disease (AD)	59
Posterior Cortical Atrophy	26
Semantic Dementia	25
Frontotemporal Dementia	19
Dementia of Unknown Origin	26
Mild Cognitive Impairment	7
Anorexia	44
Anxiety	12
Depression	30
Obsessive Compulsive Disorder	25
Drug Addiction	3
Cerebral Hemorrhage	7
Cerebral Infarction	8
Cerebral Small Vessel Disease	3
Encephalitis	175
Possible Creutzfeldt-Jakob Disease	22
Drug-Induced Parkinsonism	3
Dopa-Responsive Dystonia	3
Dystonia	2
Normal Pressure Hydrocephalus	2
Cerebral Palsy	32
Epilepsy	81
Motor Neuron Disease	3
Klein-Levin Syndrom	2
Narcolepsy	4
Healthy Persons	220



**Supplemental Figure 1** the ROC curve in exclusion of non-parkinsonian patients

**Supplemental Table 4** The performance of the proposed method in exclusion of non-parkinsonian patients based on FDG PET.

	<b>Other</b>	<b>MSA/IPD/PSP</b>
<b>Sensitivity</b>	98.0%	95.0%
<b>Specificity</b>	95.0%	98.0%
<b>PPV</b>	92.5%	98.7%
<b>NPV</b>	98.7%	92.5%

## Data difference between Chinese and German cohort

### PET/CT protocol difference between Chinese and German cohort

#### Chinese Cohort

After attenuation correction performed using low-dose CT, the emission scan was acquired at 60-minute post injection of approximately 185 MBq  $^{18}\text{F}$ -FDG and lasted 10 minutes (Siemens Biograph 64 HD PET/CT, Siemens, Germany). PET images were reconstructed by using the ordered subset expectation maximization method following corrections for scatter, dead time, and random coincidence.

#### German Cohort

##### (1) Siemens ECAT EXACT HR+ and GE Discovery 690

FDG-PET images were acquired on a GE Discovery 690 PET/CT scanner or a Siemens ECAT EXACT HR+ PET scanner. All patients had fasted for at least six hours and had a maximum plasma glucose level of 150 mg/dl at time of scanning. A single intravenous dose of  $140 \pm 7$  MBq FDG was administered while the patients rested in a room with dimmed light and low noise level, where they remained undisturbed for 20 minutes. After positioning in the scanner, a series of three static emission frames of five minutes each was acquired from 30 to 45 min p.i. on the GE Discovery 690 PET/CT, or from 30 to 60 min p.i. on the Siemens ECAT EXACT HR+ tomograph. A low-dose CT scan or a transmission scan with external  $^{68}\text{Ge}$ -source performed just prior to the static acquisition was used for attenuation correction. PET data were reconstructed iteratively (GE Discovery 690 PET/CT) or with filtered-back-projection (Siemens ECAT EXACT HR+ PET). After correction for movement between frames, the static scans were averaged.

##### (2) Siemens Biograph 64

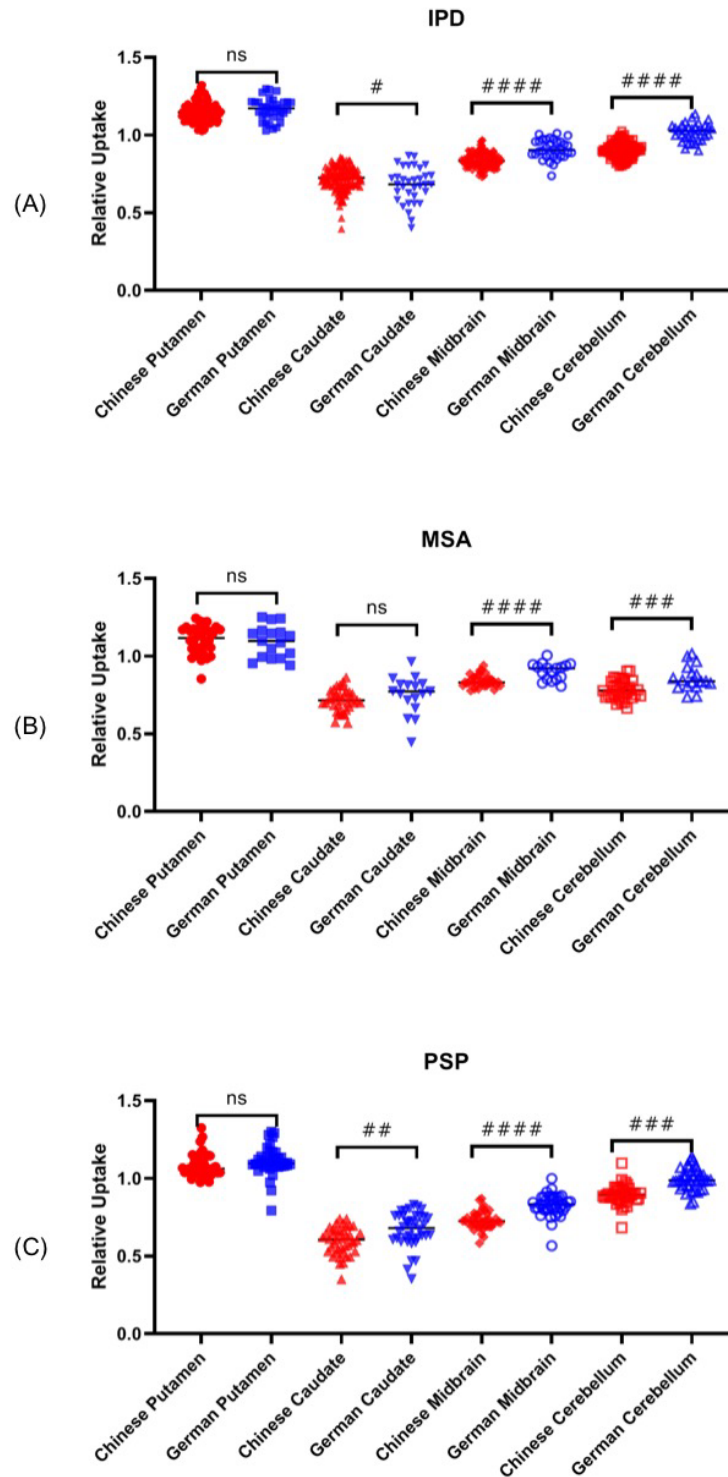
The PET data were acquired on a Siemens Biograph True point 64 PET/CT (Siemens, Erlangen, Germany). The dynamic brain PET data were acquired in 3-dimensional list-mode over 20min and reconstructed into a  $336 \times 336 \times 109$  matrix (voxel size:  $1.02 \times 1.02 \times 2.03$  mm<sup>3</sup>) using the built-in ordered subset expectation maximization (OSEM) algorithm with 4 iterations, 21 subsets and a 5mm Gaussian filter. A low dose CT served for attenuation correction.



**Supplemental Table 5 The comparison of the PET/CT protocols between Chinese and German cohorts**

	Chinese cohort	German cohort		
	Siemens Biograph 64	Siemens ECAT Exact HR+	GE Discovery 690	Siemens Biograph 64
Sensitivity	4.5 kcps/MBq	6.65 kcps/MBq	7.5 cps/kBq	4.5 kcps/MBq
Transverse Resolution	4.2± 0.3 mm	4.39 mm	4.70	4.2± 0.3 mm
Axial Resolution	4.5± 0.3 mm	5.10 mm	5.06	4.5± 0.3 mm
Peak NEC	93 kcps	37 kcps	139.1 kcps	93 kcps
Scatter Fraction	32%	46.9%	37%	32%
Injection dose (MBq)	~185	140 ± 7	140 ± 7	/
Acquisition time p. i. (min)	60	95	30	30
Imaging duration (min)	10	20	15	20
Reconstruction method	OSEM	IFBP	Iterative	OSEM
Attenuation correction	CT	68 Ge transmission	CT	CT
Reconstructed voxel size	2.03×2.03×1.5 mm <sup>3</sup>	1.4×1.4×2.4 mm <sup>3</sup>	/	1.02×1.02×2.03 mm <sup>3</sup>
Smooth	Gaussian 10mm	/	/	Gaussian 5mm
Eye mask	yes	/	/	/
Fasting	>6 hour	>6 hour	>6 hour	>6 hour
Blood glucose level	<150 mg/dl	<150 mg/dl	<150 mg/dl	<150 mg/dl

**Uptake difference between Chinese cohort and German cohort**



**Supplemental Figure 2** The comparisons of relative uptake between Chinese cohort and German cohort of regions including Putamen, Caudate, Midbrain, Cerebellum (\* indicates  $P \leq 0.05$ , \*\* indicates  $P \leq 0.01$ , \*\*\* indicates  $P \leq 0.001$ , \*\*\*\*:  $P < 0.0001$ ).

## Test of Global Mean Normalization

To test the robustness of the deep metabolic imaging (DMI) indices, we tested the performance of the DMI indices extracted from FDG PET scans after the Global Mean Normalization. We removed all normalization layers in the Parkinson Differential Diagnosis Network (PDD-Net) to keep the intensity information of the original PET scans.

Different from the Z-score normalization which is defined as:

$$J_z(x) = (I(x) - \mu) / \sigma,$$

where  $J_z(x)$  represents the Z-score normalized PET image,  $x$  is a voxel,  $\mu$  is the average and  $\sigma$  is the standard deviation of the PET uptake computed within.

The global mean normalization is defined as:

$$J_G(x) = I(x) / u_G,$$

where  $J_G(x)$  represents the global-mean normalized PET image,  $x$  is a voxel, and  $u_G$  is the average of the PET uptake computed in the whole brain.

As shown in Supplemental Table 6, we found the DMI indices obtained similar performance between using two different normalization methods (ROCAUC P-value: 0.577 for IPD, 0.589 for MSA, and 0.617 for PSP). Generally, Z-score normalization resulted in slightly better ROCAUC than global mean normalization for MSA (0.001, 0.001, and 0.003 higher for overall, short symptom durations and long symptom durations respectively) but slightly lower ROCAUC for IPD (0.003 and 0.008 lower for overall and short symptom durations). For PSP, Z-score obtained slightly lower ROCAUC on overall (0.005 lower) and short symptom durations (0.020 lower) but slightly higher ROCAUC on long symptom durations (0.005 higher).

The individual Z-score normalization resulted in slightly high accuracy for the DMI indices. However, the choice of other intensity normalization methodologies may influence the data analysis and result in improved performance in the deep learning methods outlined herein, which future studies should consider.

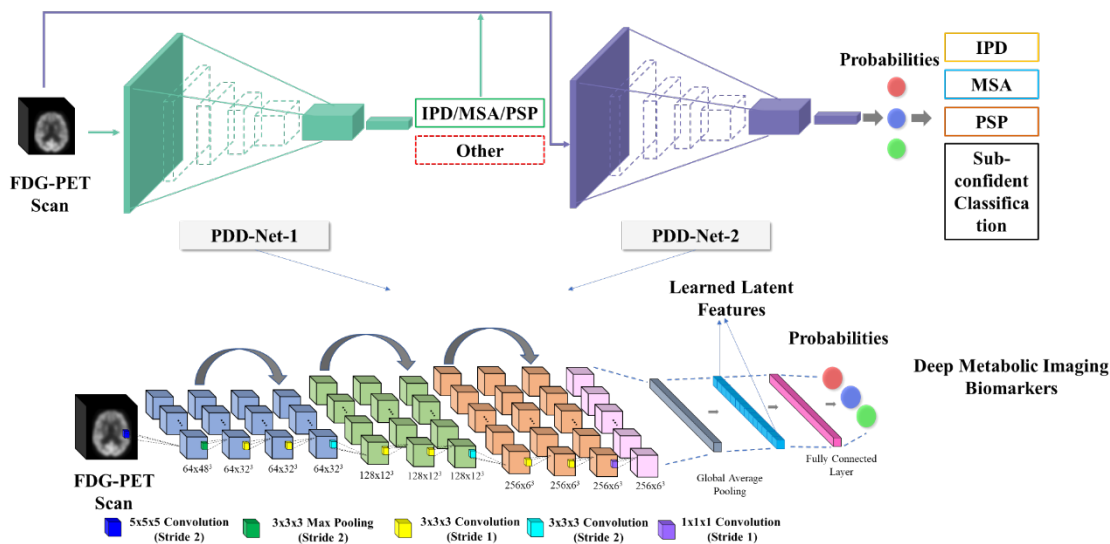
**Supplemental Table 6** Diagnostic accuracy of the DMI indices for parkinsonian disorders utilizing the Global Mean Normalization in data pre-processing step (Cross-validation, Training Cohort).

		<b>Overall</b>	<b>Short Symptom Durations (<math>\leq 2</math> years)</b>	<b>Long Symptom Durations (<math>&gt; 2</math> years)</b>
<b>Idiopathic Parkinson Disease</b>	ROCAUC	0.989 (0.983-0.996)	0.989 (0.979-0.999)	0.991 (0.983-0.998)
	Sensitivity	95.7% (92.7%-97.7%)	97.1% (92.6%-99.2%)	95.7% (91.4%-98.3%)
	Specificity	94.8% (91.2%-97.2%)	95.2% (89.8%-98.2%)	94.4% (88.7%-97.7%)
	PPV	95.7% (92.7%-97.7%)	95.7% (90.8%-98.8%)	95.7% (91.3%-98.3%)
	NPV	94.8% (91.2%-97.2%)	96.7% (91.8%-98.8%)	94.4% (88.8%-97.7%)
<b>Multiple System Atrophy</b>	ROCAUC	0.996 (0.991-1.000)	0.995 (0.988-1.000)	0.995 (0.988-1.000)
	Sensitivity	97.3% (93.3%-99.3%)	97.8% (92.2%-99.7%)	98.3% (91.1%-100%)
	Specificity	99.0% (97.4%-99.7%)	99.4% (96.8%-100%)	98.7% (96.2%-99.7%)
	PPV	97.3% (93.4%-99.3%)	98.9% (94.0%-100%)	95.2% (86.9%-99.9%)
	NPV	99.0% (97.4%-99.7%)	98.8% (95.8%-100%)	99.6% (97.5%-99.9%)
<b>Progressive Supranuclear Palsy</b>	ROCAUC	0.987 (0.978-0.995)	0.988 (0.975-1.000)	0.985 (0.974-0.997)
	Sensitivity	87.8% (79.6%-93.5%)	91.2% (76.3%-98.1%)	87.5% (79.2%-95.2%)
	Specificity	98.0% (96.2%-99.1%)	97.8% (94.9%-99.3%)	97.8% (97.0%-99.9%)
	PPV	90.5% (83.3%-95.1%)	86.1% (72.4%-96.9%)	91.8% (84.1%-98.3%)
	NPV	97.3% (95.2%-98.8%)	98.7% (95.8%-99.6%)	96.5% (93.9%-98.7%)

ROCAUC: the area under the receiver operating characteristic curve, PPV: positive predictive value, NPV: negative predictive value.

## The developed deep learning method.

The deep learning method contains two PDD-Nets. The PDD-Net-1 sought to exclude patients without parkinsonism. The PDD-Net-2 performed computation of deep metabolic imaging (DMI) indices and classification of IPD, MSA, or PSP. Both PDD-Net-1 and PDD-Net-2 are based on a 3D residual convolutional neural network (Supplemental Fig. 3).



**Supplemental Figure 3:** A sketch of the developed deep learning methods, which has two stages i.e., control stage and classification stage. In the control stage, The Parkinson Differential Neural Network-1(PDD-Net-1) works to exclude non-parkinsonian patients. In the classification stage, the Parkinson Differential Neural Network-2 (PDD-Net-2) extracts the deep metabolic imaging (DMI) indices to classify idiopathic Parkinson’s disease (IPD), multiple system atrophy (MSA), and progressive supranuclear palsy (PSP). Our network used the instance normalization in the architecture.

The employed deep neural network, i.e., Parkinson Differential Diagnosis Network (PDD-Net) comprised a down-sampling path including three repeated encoder stacks, a global average pooling and a fully connected layer with softmax activation. In each encoder stack, there were a residual module and a  $3 \times 3 \times 3$  convolutional layers with stride 2 for down-sampling the feature maps. Each residual module included two  $3 \times 3 \times 3$  convolutional layers and one dropout layer. The residual connections were employed for simplifying the optimization of the network and alleviating the vanishing gradient problem (6). We employed leaky rectified linear units (ReLU) as the activation function following the convolution layers and utilized categorical cross-entropy loss to train the network.

We implemented the network with the Keras library. Adam optimizer was used during training with an initial learning rate  $lr_{init} = 10^{-4}$ . The learning rate was reduced by a factor of 2 once learning stagnates. To regularize the network, we utilized the early stopping strategy with the patience of 10, which is a method employed to detect the convergence of training thereby avoiding overfitting. We implemented the full-gradient saliency map method by referring the library in (7) based on Pytorch.

The validation of the deep learning method was performed in two ways, using six-fold cross-validation in the training cohort and conducting an independent test in the blind-test cohort. As mentioned above, we first pre-trained the Parkinson Differential Diagnosis Network (PDD-Net) on 397 patients (the pre-training cohort). Then, we further trained the network and conducted

six-fold cross-validation in the training cohort. Finally, we utilized the blind-test cohort of the dataset to further evaluate the effectiveness of our method. In this blind-test stage, we employed a model ensemble procedure (8) to allow all six trained models in the cross-validation phase to jointly contribute to the differential diagnosis of parkinsonism. The obtained deep metabolic imaging (DMI) indices were the average DMI indices of six obtained models. The ground-truth labels of the samples in blind-test cohort were remained unseen for the algorithm developers. The obtained diagnosis classifications and related DMI indices of the obtained network was sent to our clinical co-authors (nuclear medicine physician) for independent evaluation. These clinical co-authors did not have access to or played a role in developing the algorithm.

The ensemble strategy can be further summarized as follows:

- (1) Obtaining six trained model from cross-validation stages.
- (2) These six models are utilized to directly predict the possibilities of IPD/MSA/PSP for the subject on blind-test cohort.
- (3) We calculate the average prediction possibilities of the six models as follows:

$$P_E[IPD, MSA, PSP] = \frac{1}{6} \sum_{i=1}^6 P_i[IPD, MSA, PSP],$$

Where  $P_E[IPD, MSA, PSP]$  is the ensembled possibilities and  $P_i[IPD, MSA, PSP]$  is  $i^{\text{th}}$  prediction possibilities from the  $i^{\text{th}}$  model.

- (4) Based on  $P_E$  and we referred the cut-off points in the cross-validation to determine the prediction diagnoses.
- (5) All prediction diagnoses were submitted to our clinical parameters for independently evaluation.

## Confidence inspection

The prediction according to the deep metabolic imaging (DMI) indices is generally derived based on the maximal probability of the three probabilities of idiopathic Parkinson’s disease (IPD), multiple system atrophy (MSA), and progressive supranuclear palsy (PSP). An option to warn the uncertain predictions is also provided if the maximal probability is below certain customized threshold. A default set of confidence thresholds (IPD: 0.51, MSA: 0.80, PSP: 0.56) were derived based on the generalized Youden’s index in the cross-validation stage. This set of optimal cut-off points utilized in this study were determined in the cross-validation stage and resulted in warning of eight predictions (IPD: 0, MSA: 6, PSP: 2) below the thresholds in the blind test, which were flagged as being uncertain cases. The users can alternatively customize the confidence thresholds. A set of more strict confidence thresholds of 0.8 for IPD, MSA, and PSP were tested. With this set of thresholds, more patients were warned (29/330 vs 8/330) as uncertain. If we only consider the confident predictions in the summary of accuracies, the statistics are shown in the following Supplemental table 7.

**Supplemental Table 7** Diagnosis accuracy of the DMI indices in only confident predictions for parkinsonian disorders utilizing confidence threshold of 0.8 for IPD, MSA, and PSP (blind test)

		<b>Overall<sup>1</sup></b>	<b>Baseline<sup>2</sup></b>	<b>Follow-up<sup>3</sup></b>
<b>Idiopathic Parkinson Disease</b>	Sensitivity	91.4%	87.8%	86.4%
	Specificity	94.1%	95.2%	99.9%
	PPV	96.5%	96.6%	99.9%
	NPV	86.2%	83.3%	82.4%
<b>Multiple System Atrophy</b>	Sensitivity	78.7%	77.3%	95.5%
	Specificity	99.3%	99.9%	99.9%
	PPV	96.0%	99.9%	99.9%
	NPV	95.4%	94.5%	98.9%
<b>Progressive Supranuclear Palsy</b>	Sensitivity	81.0%	80.0%	95.0%
	Specificity	98.5%	98.9%	97.7%
	PPV	92.2%	94.1%	90.5%
	NPV	96.0%	95.6%	98.9%

<sup>1</sup> The statistics of Overall summarizes the accuracy of all the 330 patients of the blind-test cohort based on the DMI indices extracted from the FDG PET imaging at baseline diagnosis.

<sup>2</sup> The statistics of Baseline summarizes the accuracy of 108 patients with repeated PET scans based on the DMI indices extracted from the baseline FDG PET imaging.

<sup>3</sup> The statistics of Follow-up summarizes the accuracy of 108 patients with repeated PET scans based on the DMI indices extracted from the follow-up FDG PET imaging.

PPV and NPV represent positive predictive value and negative predictive value.

## **Performance of the combining demographic and clinical features with deep metabolic imaging**

To evaluate the performance of leveraging multi-modality data by combining the DMI indices with demographic and clinical features, a decision tree-based classifier, Extreme Gradient Boosting (XGBoost) (9), was trained to combine the DMI indices with demographic information and clinical data (age, gender, symptom duration, unified Parkinson's disease rating scale-III (UPDRS-III), Hoehn and Yahr stage) to obtain combined diagnostic classifications.

Compared to the prediction based on the DMI indices only, the combination of the DMI indices with demographic and clinical features had almost the same accuracy in the blind-test cohort including overall 330 subjects ( $P=0.999$ ). Similarly, for the 108 patients in the blind-test cohort who had follow-up imaging available, there was almost no performance difference between the prediction of DMI indices only and the combination at baseline ( $P=0.999$ ) or at the follow-up ( $P=0.735$ ) (Details are in supplement 8). At follow-up, the sensitivity, PPV, and NPV increased for IPD (95.5% to 96.9%, 98.4% to 98.5%, 93.2% to 95.3% respectively) with the specificity remaining the same (97.6%) after the combination. For MSA, the sensitivity, specificity, PPV, and NPV all slightly increased (95.4% to 95.5%, 98.8% to 99.9%, 95.5% to 99.9%, 98.8% to 98.9% respectively) after the combination, but the metrics for PSP had no change. Overall, the performance at the follow-up did not change significantly ( $P=0.735$ ) comparing the combination with using the DMI indices only.



**Supplemental Table 8:** combining demographic and clinical features with deep metabolic imaging: Accuracy of the differentiation of the parkinsonian disorders based on the deep metabolic imaging (DMI) indices and clinical information (age, gender, symptom duration, UPDRS III, Hoehn and Yahr stage) in the blind-test cohort. (“Multi” denotes multi-modality representing combining demographic and clinical features with deep metabolic imaging, and “single” represents single modality meaning using deep metabolic imaging only involved here for easy comparison.)

		Overall		Baseline		Follow-up	
		Multi	Single	Multi	Single	Multi	Single
Idiopathic Parkinson Disease	Sensitivity	98.1%	98.1%	98.5%	98.5%	96.9%	95.5%
	Specificity	90.0%	90.0%	88.1%	88.1%	97.6%	97.6%
	PPV	94.5%	94.5%	92.9%	92.9%	98.5%	98.4%
	NPV	96.3%	96.4%	97.4%	97.4%	95.3%	93.2%
Multiple System Atrophy	Sensitivity	86.9%	88.5%	81.8%	81.8%	95.5%	95.4%
	Specificity	99.2%	99.2%	99.9%	99.9%	99.9%	98.8%
	PPV	96.4%	96.4%	99.9%	99.9%	99.9%	95.5%
	NPV	97.1%	97.4%	95.6%	95.6%	98.9%	98.8%
Progressive Supranuclear Palsy	Sensitivity	86.2%	84.5%	89.9%	90.0%	95.0%	95.0%
	Specificity	97.8%	97.8%	97.7%	97.7%	96.6%	96.6%
	PPV	89.3%	89.1%	90.1%	90.0%	86.4%	86.4%
	NPV	97.1%	97.0%	97.7%	97.7%	98.8%	98.8%

<sup>1</sup> The statistics of Overall summarizes the accuracy of all the 330 patients of the blind-test cohort based on the DMI indices extracted from the FDG PET imaging at baseline diagnosis.

<sup>2</sup> The statistics of Baseline summarizes the accuracy of 108 patients with repeated PET scans based on the DMI indices extracted from the baseline FDG PET imaging.

<sup>3</sup> The statistics of Follow-up summarizes the accuracy of 108 patients with repeated PET scans based on the DMI indices extracted from the follow-up FDG PET imaging.

PPV and NPV represent positive predictive value and negative predictive value.

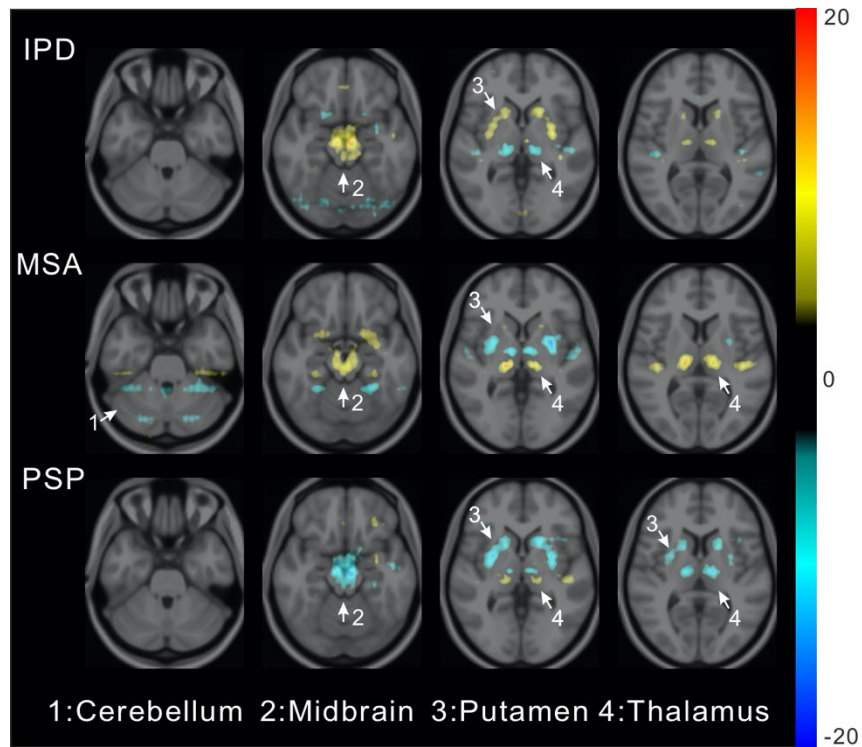
## Visualization of the deep metabolic imaging indices

We generated the saliency maps of input PET images using the full-gradient method (7) to assist the interpretation of the DMI indices. The saliency maps assign importance scores to both the input features and individual neurons in a network, which reflects the contribution of groups of pixels to the DMI probabilities. The full-gradient saliency map method (7) utilized in this work considers both the input importance indicating the contribution of individual input voxels and neuron importance reflecting the contribution of groups of voxels with specific structural information, which is sharper and more tightly confined to object regions compared to other existing methods. Thus, the full-gradient saliency map method mitigates against known issues with inaccuracy in location and provided a preliminary explanation of the learned model.

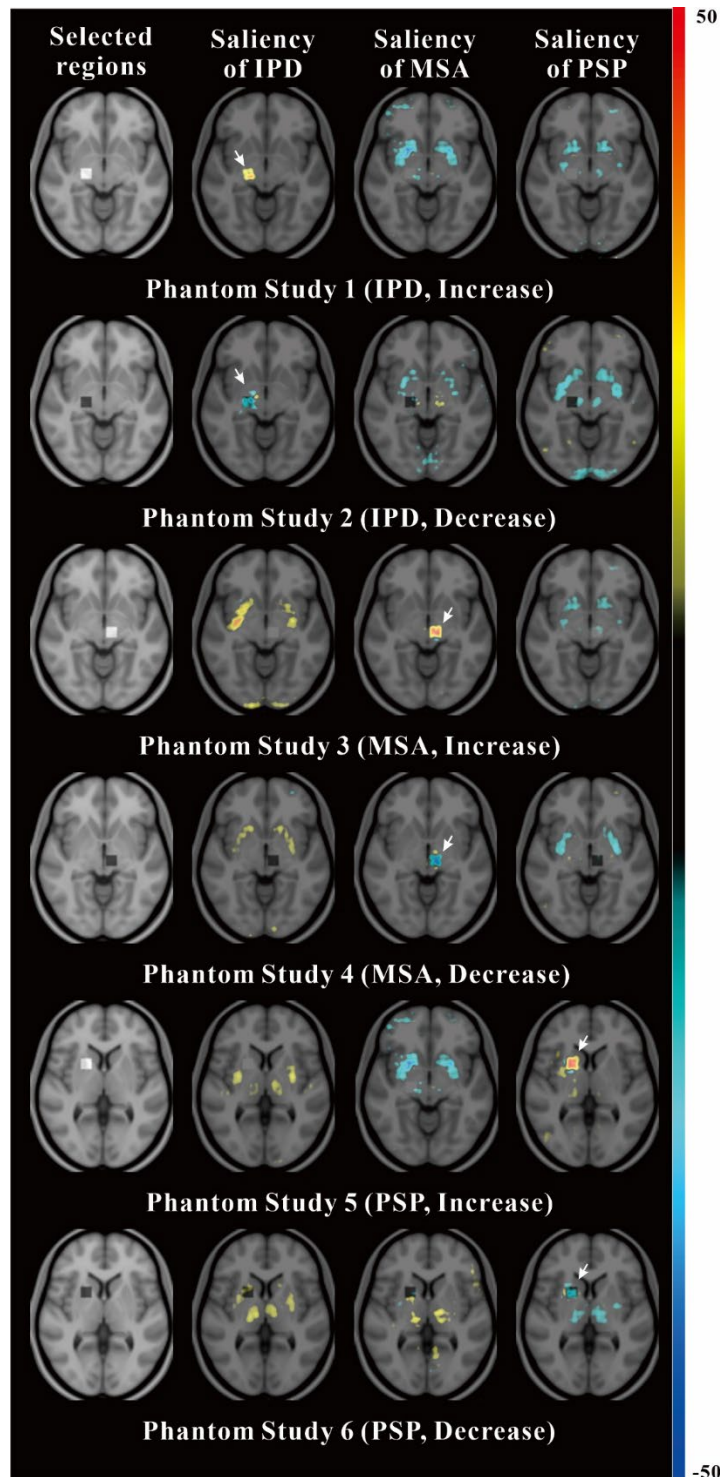
It should be noted that, due to the constraint that the sum of probability of IPD, MSA, and PSP should be equal to one, the saliency maps of IPD, MSA, and PSP are correlated, i.e., one factor leading to the increase of the IPD probability will result in the decrease of the probability of MSA and PSP simultaneously.

Supplemental Fig. 4 demonstrates average saliency maps (fused with template MRI) of patients with IPD, MSA, or PSP in the training cohort. Regions with relatively higher contribution to the DMI indices were putamen and midbrain for IPD, MSA, and PSP as well as cerebellum for MSA.

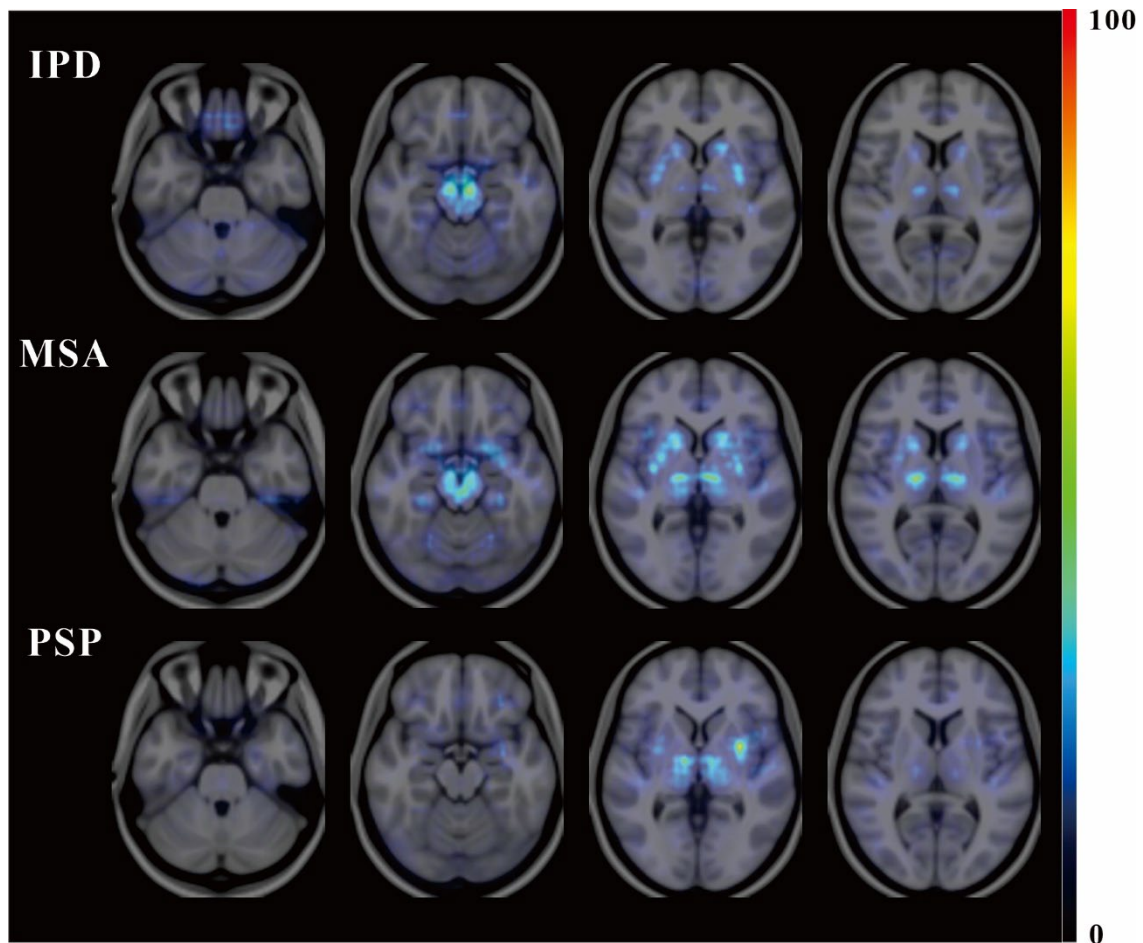
In order to interpret the results of the saliency map, we conducted 6 phantom studies. These phantom studies manipulate of the activities of the PET images of a set of 180 randomly selected patients of three categories, IPD (n=60), MSA (n=60), and PSP (n=60) from the training cohort. In each phantom study, we randomly selected a region on the PET scans (6\*6\*6 voxels), and then we artificially increased or decreased activities by 50% within this region of PET scans for patients in a category and kept PET scans of other two categories unchanged. For instance, in phantom study 1, we artificially increased activities of a selected region on PET scans in the IPD category and kept MSA and PSP categories the same as the original imaging data. Then we train the deep neural network on the artificially modified experiment datasets and calculated the saliency maps. In phantom study 2, we only decreased activities of the selected region in the IPD category and calculated the saliency maps. Similar procedures were employed in the computation of the saliency maps for MSA in phantom studies 3 and 4 and PSP in phantom studies 5 and 6. By manipulating the activities in each phantom study, the artificial regions with increased/decreased activities in one category were the most salient difference regions compared to the other two unchanged categories. The results are illustrated in the Supplemental Fig. 5, where we found that the saliency map recognized the selected regions with artificial characteristic activity-increase/-decrease as salient regions, which indicated the effectiveness and accuracy of the saliency map method.



**Supplemental Figure 4.** Visualization of average saliency maps of patients with idiopathic Parkinson’s disease (IPD), multiple system atrophy (MSA) and progressive supranuclear palsy (PSP) in the training cohort showing characteristic regions contributing to the deep metabolic imaging (DMI) indices. The colour corresponds to the importance score indicating the contribution of a region for the generated the deep metabolic imaging (DMI) indices. The colour directions (yellow and red vs cyan and blue) represent different influences on the DMI indices (Increase and Decrease the probability in the DMI indices). The arrows pointed to the most salient brain regions including 1: Cerebellum, 2: Midbrain, 3: Putamen, 4: Thalamus.



**Supplemental Figure 5** Interpretation of saliency map using on artificially designed experiment datasets. From left to right, the first column showed the artificially selected regions for activity manipulation. The region to increase the activity is marked as bright and the region to decrease the activity is marked as dark. The remaining columns showed the average saliency map of idiopathic Parkinson's disease (IPD), multiple system atrophy (MSA) and progressive supranuclear palsy (PSP) in phantom studies.



**Supplemental Figure 6** Visualization of the variance of saliency maps of the deep metabolic imaging indices for patients with idiopathic Parkinson's disease (IPD), multiple system atrophy (MSA), and progressive supranuclear palsy (PSP) in the training cohort.

The variance of saliency maps reflects the difference of the saliency maps at each voxel of patients within IPD, MSA, and PSP groups. The color corresponds to the variance scores. From this variance map, we can find that those regions with high variance locate at parkinsonism-related regions such as midbrain, putamen, and cerebellum which are in consist with salient regions in average saliency maps.

### **Cases of diagnostic classifications of the DMI indices inconsistent with the clinical diagnosis**

In contrast to the majority of cases in Fig. 4, there existed six cases where the DMI indices made predictions inconsistent with the clinical diagnosis and six cases with obvious probability decrease during follow-up (Supplement 9). Neurologists, who remained blind to the DMI indices predictions, were invited by nuclear medicine physicians to follow up the above-mentioned twelve patients along with the same number of randomly selected consistent samples. In one patient, at the post-AI follow-up, the diagnosis was updated from IPD to PSP. The DMI indices and neurologists both diagnosed the patient with IPD at baseline, but the DMI indices correctly diagnosed this patient as PSP at the first follow-up time. The DMI classification and the clinical diagnosis at different time point are listed in Supplemental Table 9.

**Supplemental Table 9** The diagnostic classifications of the deep metabolic imaging (DMI) indices and the clinical diagnosis at different time point of the cases where DMI classification were inconsistent with the clinical diagnoses and cases with significantly decreased IPD probability over 0.1.

Patient Order		Baseline Time (initial scan, blind-test cohort)		Follow-up Time (repeated scan, blind-test cohort)			Post AI Follow-up	
		DMI Diagnostic Classifications <sup>1</sup>	Clinically Definite Diagnosis	DMI Diagnostic Classifications <sup>1</sup>	Clinically confirmative Diagnosis	Follow-up Time (month)	Clinically confirmative Diagnosis	Follow-up Time (month)
Inconsistent Cases	1*	<b>IPD (0.83, 0.08, 0.09)</b>	<b>IPD</b>	<b>PSP (0.02, 0.02, 0.95)</b>	<b>IPD</b>	<b>24</b>	<b>PSP</b>	<b>66</b>
	2*	IPD (0.61, 0.13, 0.26)	IPD	PSP (0.09, 0.19, 0.72)	IPD	60	IPD	115
	3*	IPD (0.61, 0.35, 0.04)	IPD	MSA (0.25, 0.71, 0.04)	IPD	12	IPD	74
	4	PSP (0.02, 0.01, 0.97)	IPD	PSP (0.03, 0.02, 0.96)	IPD	36	IPD	84
	5	PSP (0.15, 0.09, 0.76)	MSA-P	PSP (0.04, 0.09, 0.88)	MSA-P	25	MSA-P	61
	6	IPD (0.96, 0.02, 0.03)	PSP	IPD (0.79, 0.03, 0.18)	PSP	24	PSP	51
IPD Probability Decreased Cased	1	IPD (0.95, 0.02, 0.03)	IPD	IPD (0.54, 0.03, 0.43)	IPD	25	IPD	70
	2	IPD (0.97, 0.02, 0.02)	IPD	IPD (0.61, 0.35, 0.04)	IPD	24	IPD	52
	3	IPD (0.97, 0.02, 0.02)	IPD	IPD (0.67, 0.29, 0.03)	IPD	26	IPD	45
	4	IPD (0.87, 0.02, 0.11)	IPD	IPD (0.60, 0.03, 0.37)	IPD	12	IPD	96
	5	IPD (0.76, 0.03, 0.21)	IPD	IPD (0.55, 0.04, 0.41)	IPD	36	IPD	95
	6	IPD (0.83, 0.10, 0.08)	IPD	IPD (0.64, 0.04, 0.32)	IPD	23	IPD	65

<sup>1</sup>DMI Diagnostic Classifications (Probability of IPD, MSA, PSP)

<sup>2</sup>HY: Hoehn and Yahr scale

<sup>3</sup>UPDRS III: Unified Parkinson's Disease Rating Scale-III.

\*Also belong to the cases with significantly decreased IPD probability over 0.1

## **Data availability**

The Huashan parkinsonian PET imaging database will be made available to the scientific community upon completion of the non-disclosure agreement (NDA) with the corresponding author according to international data protection regulations. Our code is available for download at: <https://github.com/Louis-YuZhao/deep-metabolic-imaging-indices.git>.



## REFERENCES

1. Höglinger GU, Respondek G, Stamelou M, et al. Clinical diagnosis of progressive supranuclear palsy: the movement disorder society criteria. *Mov Disord.* 2017;32:853-864.
2. Postuma RB, Berg D, Stern M, et al. MDS clinical diagnostic criteria for Parkinson's disease. *Mov Disord.* 2015;30:1591-1601.
3. Gilman S, others. Second consensus statement on the diagnosis of multiple system atrophy. *Neurology.* 2008;71:670-676.
4. Litvan I, Agid Y, Calne D, et al. Clinical research criteria for the diagnosis of progressive supranuclear palsy (Steele-Richardson-Olszewski syndrome) report of the NINDS-SPSP international workshop. *Neurology.* 1996;47:1-9.
5. Hughes AJ, Daniel SE, Kilford L, Lees AJ. Accuracy of clinical diagnosis of idiopathic Parkinson's disease: a clinico-pathological study of 100 cases. *J Neurol Neurosurg Psychiatry.* 1992;55:181-184.
6. He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition.* 2016:770-778.
7. Srinivas S, Fleuret F. Full-gradient representation for neural network visualization. In: *The 33rd International Conference on Neural Information Processing Systems.* 2019:4126-4135.
8. Skrede O-J, De Raedt S, Kleppe A, et al. Deep learning for prediction of colorectal cancer outcome: a discovery and validation study. *Lancet.* 2020;395:350-360.
9. Chen T, Guestrin C. Xgboost: A scalable tree boosting system. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.* 2016:785-794.