



CNN-based Classification of I-123 ioflupane dopamine transporter SPECT brain images to support the diagnosis of Parkinson's disease with Decision Confidence Estimation

Master Thesis

Master of Science in Applied Computer Science

Aleksej Kucerenko

November 9, 2023

Supervisor:

1st: Prof. Dr. Christian Ledig

2nd: Dr. Ralph Buchert, Universitätsklinikum Hamburg-Eppendorf

Chair of Explainable Machine Learning

Faculty of Information Systems and Applied Computer Sciences

Otto-Friedrich-University Bamberg

Abstract

Short summary of your thesis (max. 1 page) . . .

Abstract

Kurze Zusammenfassung Ihrer Abschlussarbeit (max. 1 Seite) ...

Acknowledgements

If you want to thank anyone (optional) . . .

Contents

List of Figures	vi
List of Tables	x
List of Acronyms	xi
1 Introduction	1
2 Background	4
2.1 DAT-SPECT for Detecting Parkinson’s Disease	4
2.2 Convolutional Neural Networks for Image Classification	5
3 Data Sources	5
3.1 Development dataset	5
3.2 Independent testing datasets	6
4 Methods	6
4.1 Software Tools and Libraries	7
4.2 Development Data Preparation	7
4.2.1 Data Preprocessing	7
4.2.2 Data Augmentation	7
4.2.3 Dataset Splitting	8
4.3 Univariate benchmark: Specific Binding Ratio	9
4.4 Multivariate benchmark: PCA-enhanced Random Forest	9
4.5 CNN-based classification	10
4.5.1 MVT-based and RLT-based methods	11
4.5.2 Regression-based method	12
4.6 Evaluation Metrics and Procedure	13
5 Evaluation	14
5.1 Baseline Performance	15
5.1.1 SBR Method	15
5.1.2 PCA-RFC Method	18
5.2 Experimental Methods Performance	20
5.2.1 CNN-MVT Method	20

5.2.2	CNN-RLT Method	22
5.2.3	CNN-Regression Method	24
5.3	Comparative Performance Analysis	26
5.3.1	Performance on test set of development dataset	26
5.3.2	Performance on PPMI dataset	27
5.3.3	Performance on MPH dataset	27
5.4	Conclusion	28
6	Discussion	28
6.1	Interpretation of Results	28
6.2	Practical Implications	29
6.3	Limitations of the Study	30
6.3.1	Metric	30
6.3.2	Site-Specific Development Data	30
7	Conclusion	30
A	Appendix	61
	Bibliography	62

List of Figures

1	DVR slabs for two sample cases from the development dataset, a healthy control case (above) and a PD case with reduced availability of DAT in the striatum (below). The two cases are presented in 12 different versions. In each version, attenuation and scatter corrections are either applied ('withASC') or not applied ('woASC'). Also, for each version, isotropic 3-dimensional Gaussian kernel smoothing with different FWHM values (10, 12, 14, 16, 18mm) was either performed or not performed ('original')	8
2	Principle components of the training set (development dataset) for the first random split.	10
3	Architecture of the CNN-based classification models.	12
4	Evaluation of the SBR method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.	16
5	Evaluation of the SBR method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).	17
6	Evaluation of the SBR method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). For better illustration the area under the mean of the balanced accuracy is highlighted.	31
7	Evaluation of the SBR method on PPMI dataset.	32
8	Evaluation of the SBR method on MPH dataset.	33
9	Evaluation of the PCA-RFC method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.	34
10	Evaluation of the PCA-RFC method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).	35

11	Evaluation of the PCA-RFC method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). The area under the mean of the balanced accuracy is highlighted for better illustration.	36
12	Evaluation of the PCA-RFC method on PPMI dataset.	37
13	Evaluation of the PCA-RFC method on MPH dataset.	38
14	Evaluation of the CNN-MVT method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.	39
15	Evaluation of the CNN-MVT method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).	40
16	Evaluation of the CNN-MVT method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). The area under the mean of the balanced accuracy is highlighted for better illustration.	41
17	Evaluation of the CNN-MVT method on PPMI dataset.	42
18	Evaluation of the CNN-MVT method on MPH dataset.	43
19	Evaluation of the CNN-RLT method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.	44
20	Evaluation of the CNN-RLT method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).	45

21	Evaluation of the CNN-RLT method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). The area under the mean of the balanced accuracy is highlighted for better illustration.	46
22	Evaluation of the CNN-RLT method on PPMI dataset.	47
23	Evaluation of the CNN-RLT method on MPH dataset.	48
24	Evaluation of the CNN-Regression method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.	49
25	Evaluation of the CNN-Regression method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).	50
26	Evaluation of the CNN-Regression method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). The area under the mean of the balanced accuracy is highlighted for better illustration.	51
27	Evaluation of the CNN-Regression method on PPMI dataset.	52
28	Evaluation of the CNN-Regression method on MPH dataset.	53
29	Comparison of different methods on test set of development data. Transferability of inconclusive intervals.	54
30	Comparison of different methods on test set of development data. Balanced accuracy over the percentage of observed inconclusive cases.	55
31	Comparison of different methods on PPMI dataset. Transferability of inconclusive intervals.	56
32	Comparison of different methods on PPMI dataset. Balanced accuracy over the percentage of observed inconclusive cases.	57
33	Comparison of different methods on MPH dataset. Transferability of inconclusive intervals.	58
34	Comparison of different methods on MPH dataset. Balanced accuracy over the percentage of observed inconclusive cases.	59

List of Tables

1	Evaluation of the SBR method on Development dataset (SBR cutoff mean \pm SD: 0.703 ± 0.009)	15
2	Evaluation of the PCA-RFC method on Development dataset. For evaluation, the natural sigmoid cutoff 0.5 was used.	18
3	Evaluation of the CNN-MVT method on Development dataset. For evaluation, the natural sigmoid cutoff 0.5 was used.	21
4	Evaluation of the CNN-RLT method on Development dataset. For evaluation, the natural sigmoid cutoff 0.5 was used.	23
5	Evaluation of the CNN-Regression method on Development dataset. For evaluation, the natural sigmoid cutoff 0.5 was used.	25

List of Acronyms

^{123}I	Iodine-123
[^{123}I]FP-CIT	Ioflupane Iodine-123; Datscan
AI	Artificial Intelligence
AUC	Area Under Curve
BCE	Binary Cross Entropy
CNN	Convolutional Neural Network
CUPS	Clinically Uncertain Parkinsonian Syndrome
DAT	Dopamine Transporter
DAT-SPECT	Dopamine Transporter Single-Photon Emission Computed Tomography
DVR	Distribution Volume Ratio
FWHM	Full Width at Half Maximum
ML	Machine Learning
MNI	Montreal Neurological Institute
MPH	Multiple-pinhole
MSE	Mean Squared Error
MVT	Majority vote training
NC	Normal Control case
PCA	Principal Component Analysis
PD	Parkinson's Disease
PIncObs	Percentage of Inconclusive cases observed in test set
PIncVal	Percentage of Inconclusive cases observed in validation set
PPMI	Parkinson's Progression Markers Initiative
PS	Parkinsonian Syndrome
ResNet	Residual Neural Network
RFC	Random Forest Classifier
RLT	Random label training
ROC	Receiver Operating Characteristic
SBR	Specific Binding Ratio
SNpc	Substantia Nigra pars compacta
SPECT	Single-Photon Emission Computed Tomography
UKE	University Medical Center Hamburg-Eppendorf

1 Introduction

Parkinson's disease (PD) is the second most common neurodegenerative disease after Alzheimer's disease (Twelves et al., 2003). It is expected to impose an increasing social and economic burden on societies as populations age (de Lau and Breteler, 2006). The prevalence of PD in industrialized countries is about 1% in people over 60 years of age (de Lau and Breteler, 2006). The standardized incidence rate of PD is estimated to range between about 10 and about 20 per 100,000 person-years (de Lau and Breteler, 2006). This results in the diagnosis of up to 100,000 new PD cases annually in the EU and up to 50,000 cases in the US.

PD is typically characterized by bradykinesia and variable expression of cardinal symptoms: resting tremor, rigidity, and postural instability (Tolosa et al., 2006; Gibb and Lees, 1988). However, this combination of symptoms, often referred to as 'parkinsonism' or 'parkinsonian syndrome' (PS), occurs not only in PD (and some rare 'atypical' neurodegenerative PS such as multiple system atrophy, progressive supranuclear palsy and corticobasal degeneration). It also occurs in so-called 'secondary' (non-neurodegenerative) PS that can be induced by drugs, head trauma, inflammatory or metabolic disorder, as well as other diseases such as essential tremor, dystonic tremor, or normal pressure hydrocephalus (Tolosa et al., 2006; Piccini and Whone, 2004). A particularly frequent cause of secondary PS is cerebrovascular disease (Funke et al., 2013). The differentiation between PD and secondary PS is highly relevant because secondary PS might be treated more effectively than PD and some secondary PS may be fully cured. Yet, the clinical, that is, symptom-based differentiation between PD and secondary PS is challenging in a significant fraction of patients, particularly at early disease stages with mild symptoms and in patients with atypical presentation (Hughes et al., 2002, 1992). These cases are often referred to as 'clinically uncertain parkinsonian syndromes' (CUPS) (Catafau et al., 2004).

Dopamine transporter single-photon emission computed tomography (DAT-SPECT) with ioflupane (^{123}I), also referred to as [^{123}I]FP-CIT, is an established nuclear medicine brain imaging procedure for Parkinson's disease diagnosis. The widespread use of the diagnostic procedure is due to its high accuracy, its relevant impact on patient management, and the strong guideline recommendations. In Europe, around 70,000 patients undergo DAT-SPECT scans annually, with 10,000 in Germany alone, and at the University Medical Center Hamburg-Eppendorf (UKE) currently around 400 per year (Marienhagen et al., 2017). The demographical change in industrial countries is expected to result in a further increase in the number of DAT-SPECT examinations, because age is the major risk factor for PD (Reeve et al., 2014). Furthermore, there are early signs of PD such as smell loss and idiopathic rapid eye movement sleep and behavioral disorder that can precede motor symptoms by several years but are not particularly specific for PD (Iranzo et al., 2017; Postuma and Berg, 2019; Postuma et al., 2019). It becomes increasingly important to detect PD at these early pre-motor stages because the earlier the treatment is initiated the better the chances of moderating the course of PD with disease-modifying drugs (Kim, 2017).

In clinical practice, the interpretation of DAT-SPECT is binary, that is, the nuclear medicine physician has to decide whether the SPECT images indicate degeneration of the dopaminergic neurotransmitter system (Parkinson’s disease) or not (secondary PS). This decision can be challenging by visual inspection of the tomographic SPECT images, particularly for less experienced readers (Schiebler et al., 2023). Thus, DAT-SPECT would benefit from methods for the automatic classification of the images that achieve similar (or better) performance as experienced readers. Convolutional neural networks (CNNs) appear particularly promising for this purpose (Wenzel et al., 2019; Chien et al., 2020; Magesh et al., 2020; Hathaliya et al., 2022; Nazari et al., 2022).

Yet, there are also ‘true’ borderline cases that cannot be classified with high certainty even by expert readers. In DAT-SPECT of CUPS, the proportion of visually inconclusive borderline cases ranges between 5 and 10% (Mäkinen et al., 2016; Albert et al., 2016). Automatic binary classification of these cases by a CNN might pretend a certainty of the diagnosis that is not actually given. It is important, therefore, to identify these cases to ensure that the user visually inspects these SPECT images in order to check the automatic categorization particularly carefully. The user will accept the CNN’s decision in some cases, overrule the CNN in other cases, and categorize the remaining cases as actually inconclusive (and might recommend follow-up DAT-SPECT after 6-12 months (Apostolova et al., 2017)).

The most obvious approach to identify borderline cases in CNN-based classification would be based on the distance of the CNN’s sigmoid output from a predefined decision threshold (e.g., 0.5). However, empirically, sigmoid outputs of CNN for classification of DAT-SPECT tend to cluster at the extreme values so that their utility for the identification of borderline cases seems limited. As a consequence, this approach is not recommended among practitioners, as it tends to overestimate the certainty of CNN-based classification (Ulmer and Cinà, 2021; Guo et al., 2017; Karimi and Gholipour, 2020).

Against this background, the current work aimed to propose and validate a CNN-based approach for the automatic classification of DAT-SPECT that allows reliable identification of inconclusive cases that might be misclassified by the CNN when the decision threshold is strictly applied. The ‘decision confidence’ of the classification model is evaluated on a metric, proposed in the following, that aims to maximize the classification performance of the model in conclusive cases while minimizing the potential effort of manual inspection originating from inconclusive cases.

Starting from the assumption that between-readers discrepancy in the binary visual interpretation of DAT-SPECT is much more likely in borderline cases than in conclusive cases, a standard CNN structure was trained for the automatic classification of DAT-SPECT using a large training dataset in which each SPECT image had been visually classified by three independent readers. During the model training phase, the standard-of-truth label was selected randomly from the three independent available reads. This way, the same borderline case image could be presented to the network with different standard-of-truth labels. The rationale was that this could allow the network to learn about the uncertainty of these cases, and that this would result

in sigmoid outputs close to the decision threshold. This “random label” training (RLT) approach was compared with the conventional majority vote training (MVT) approach. In the latter, the majority vote across the three readers was consistently used as standard-of-truth during the training phase. The MVT obviously “hides” the uncertainty associated with between-readers discrepancy from the network. To be able to better assess the performance of the CNN-based approaches, univariate and multivariate conventional methods were employed as benchmark methods. In addition, the performance of the approaches was also evaluated on independent external datasets.

The primary hypothesis put to test in this work was that the sigmoid output of the CNN is more effective for the identification of inconclusive cases (by an ‘inconclusive’ range around the decision threshold) when the network is trained using the RLT strategy compared to MVT.

To test this hypothesis, the proportion of inconclusive cases required to achieve a given balanced accuracy in the conclusive cases was proposed and used as a performance metric. More precisely, the area under the curve (AUC) of mean balanced accuracy in conclusive cases versus the mean proportion of inconclusive cases (observed in the test set) was used as a model-agnostic quality metric. The AUC does not depend on a specific operating point (target balanced accuracy). The rationale for this performance metric was that more inconclusive cases would require more attention and manual inspection by the attending physician which is considered ‘expensive’ (“90% inconclusive cases to achieve the required accuracy in the remaining 10% of cases is clearly useless”). Therefore the utility of the classifier for widespread use in clinical practice depends on the proportion of inconclusive cases to be accepted to achieve a predefined balanced accuracy in the remaining conclusive cases.

The following secondary hypotheses were put to test. First, CNN-based classification outperforms conventional baseline methods in terms of AUC of balanced accuracy, both univariate and multivariate baseline methods. The specific binding ratio (SBR) of the tracer uptake in the putamen was used for the univariate analysis as a benchmark method. Current procedure guidelines recommend the putaminal SBR to support the visual interpretation of DAT-SPECT in everyday clinical patient care (Morbelli et al., 2020). The putaminal SBR characterizes the contrast of the tracer uptake (= intensity) in the putamen relative to the mean tracer uptake in a reference region void of DAT (Buchert et al., 2019b). The putaminal SBR was assumed to be proportional to the density of DAT in the putamen (Buchert et al., 2019b). As a multivariate benchmark method, a random forest approach was implemented using the expression profile of a set of covariance patterns as input. The covariance patterns were identified by principal component analysis in the training dataset.

Second, CNN-based classification demonstrates enhanced generalizability, particularly in its robustness concerning varying image characteristics, such as spatial resolution. In particular, the need for robustness against variations in image characteristics arises from differences in acquisition hardware, such as various SPECT cameras and collimators, as well as from diverse reconstruction and correction meth-

ods, including those addressing photon attenuation, scatter recovery, and resolution recovery. To validate this hypothesis, the classification methods, trained on the training set of the development dataset, were assessed using two test datasets that were entirely separate from the development dataset.

The following research questions are addressed:

- When comparing the CNN-based classification approaches, how does the RLT approach perform compared to the MVT approach using the proposed performance metric?
- How do the CNN-based approaches perform on diverse testing data compared to conventional approaches? What conclusions can be made regarding the generalizability of the approaches under test?

Include thesis structure paragraph. TODO

2 Background

2.1 DAT-SPECT for Detecting Parkinson's Disease

Parkinson's disease (PD) and 'atypical' neurodegenerative Parkinsonian syndromes (PS) are both associated with the progressive loss of dopaminergic neurons in the substantia nigra pars compacta (SNpc) that project to the striatum (Piggott et al., 1999). The reduced availability of dopamine transporters (DAT) in the striatum is a well-validated biomarker for nigrostriatal degeneration in PD (Bernheimer et al., 1973; Fazio et al., 2018; Niznik et al., 1991). It can be detected by single photon emission computed tomography (SPECT) with dopamine transporter (DAT) ligands (Kuikka et al., 1995; Abi-Dargham et al., 1996). The reduction in striatal DAT availability is significantly advanced even in the earliest symptomatic (motor) stages of PD, as the degeneration of dopaminergic nerve endings in the striatum represents an early step in the pathological PD cascade (Bernheimer et al., 1973; Fazio et al., 2018; Niznik et al., 1991). The compensatory downregulation of the DAT expression in the remaining nerve endings leads to a more pronounced loss of striatal DAT (Lee et al., 2000; Saari et al., 2017; Honkanen et al., 2019). Secondary PS's are typically not associated with nigrostriatal degeneration or the loss of striatal DAT. To differentiate PD from secondary PS based on striatal DAT availability, the radiolabeled DAT ligand [¹²³I]FP-CIT (trade name: DaTscan[®]) has been approved as a SPECT tracer in both the US and Europe (Neumeyer et al., 1994).

A recent review, which involved a non-systematic meta-analysis of DAT-SPECT with [¹²³I]FP-CIT in patients with PS, confirmed that DAT-SPECT exhibits high sensitivity (median 93%) and high specificity (median 89%) in differentiating PD from secondary PS in patients with clinically uncertain parkinsonian syndrome (CUPS) (Buchert et al., 2019a). Moreover, the review demonstrated that DAT-SPECT results in a change in diagnosis for about 40% of patients with CUPS and

leads to a change in treatment for a similar proportion of these patients (Buchert et al., 2019a). Thus, DAT-SPECT with [^{123}I]FP-CIT is a highly accurate diagnostic method that significantly influences the diagnosis and treatment of patients with CUPS. Guidelines from professional neurological societies have therefore strongly emphasized the role of DAT-SPECT with [^{123}I]FP-CIT in recent years (Tatsch and Poepperl, 2013). For example, the current version of the S3 guideline “Idiopathic Parkinson syndrome” of the German Society of Neurology states that DAT-SPECT *should* be conducted at an early disease stage in CUPS patients.

2.2 Convolutional Neural Networks for Image Classification

3 Data Sources

The study retrospectively included 3 different datasets with a total of 3025 DAT-SPECT images. The primary dataset (“development dataset”) was used for both training and testing the models associated with the respective method, whereas the other two datasets, the *PPMI* dataset and the *MPH* dataset were used for testing only, not for training.

3.1 Development dataset

The development dataset consisted of 1740 consecutive DAT-SPECT scans obtained from clinical routine at the Department of Nuclear Medicine, University Medical Center Hamburg-Eppendorf (Schiebler et al., 2023). In brief, DAT-SPECT with [^{123}I]FP-CIT had been performed according to common procedures guidelines (Darcourt et al., 2010; Djang et al., 2012) with different double-head cameras equipped with low-energy-high-resolution or fan-beam collimators. The projection data were reconstructed using the iterative ordered-subsets-expectation-maximization (Hudson and Larkin, 1994) with attenuation and simulation-based scatter correction as well as collimator-detector response modeling as implemented in the Hybrid Recon-Neurology tool of the Hermes SMART workstation v1.6 (Hermes Medical Solutions, Stockholm, Sweden) (Diemling, 2021; Sohlberg and Kajaste, 2012; Solutions; Kangasmaa et al., 2016). All parameter settings were as recommended by Hermes (Diemling, 2021) for the EANM / EANM Research Ltd (EARL) ENC-DAT project (European Normal Control Database of DaTSCAN) (Tossici-Bolt et al., 2011; Dickson et al., 2010; Varrone et al., 2013; Tossici-Bolt et al., 2017; Dickson et al., 2012). More precisely, ordered-subsets-expectation-maximization was performed with 5 iterations and 15/16 subsets for 120/128 views. For noise suppression, reconstructed images were postfiltered by convolution with a 3-dimensional Gaussian kernel of 7 mm full-width-at-half-maximum. The ground-truth label, indicating either ‘normal’ or ‘Parkinson-typical’ reduction (‘reduced’) of the striatal signal, was obtained by visual interpretation of the DAT-SPECT images by three independent readers (Schiebler et al., 2023). The between-reader consensus on the label could not be

achieved for around 5% of dataset cases. The development dataset was utilized for both training and testing the classification models.

3.2 Independent testing datasets

The second dataset comprised 645 DAT-SPECT with [¹²³I]FP-CIT from the Parkinson’s Progression Markers Initiative (PPMI) (www.ppmi-info.org/data) (Parkinson Progression Marker Initiative, 2011). The external dataset included 438 patients with Parkinson’s disease and 207 healthy controls as described in Wenzel et al. (2019). Details of the PPMI DAT-SPECT protocol are given at <http://www.ppmi-info.org/study-design/research-documents-and-sops/> (Parkinson Progression Marker Initiative, 2011). Raw projection data has been transferred to the PPMI imaging core lab for central image reconstruction using an iterative (HOSEM) algorithm on a HERMES workstation. The clinical diagnosis was used as ground-truth label (Parkinson’s disease = “reduced”, healthy control = “normal”). The external dataset showed lower spatial resolution than the development dataset (lower striatum-to-background contrast).

The third dataset (“MPH dataset”) comprised 640 consecutive DAT-SPECT with [¹²³I]FP-CIT from clinical routine at UKE that had been acquired with a triple-head camera equipped with brain-specific multiple pinhole (MPH) collimators. Multiple pinhole SPECT concurrently improves count sensitivity and spatial resolution compared to SPECT with parallel-hole and fan-beam collimators (Mathies et al., 2022; Tecklenburg et al., 2020). The projection data were reconstructed with the Monte Carlo photon simulation engine and iterative one-step-late maximum-a-posteriori expectation-maximization implemented in the camera software (24 iterations, 2 subsets) (Tecklenburg et al., 2020; Magdics et al., 2010). Neither attenuation nor scatter correction was applied. The ground-truth label (“normal” or “reduced”) was obtained by the visual interpretation of an experienced reader (about 20 years of experience in clinical DAT-SPECT reading, $\geq 3,000$ cases). All SPECT images were interpreted twice (with different randomization) by the same reader. The delay between the reading sessions was 14 days. Cases with discrepant interpretations between the two reading sessions were read a third time by the same reader to obtain an intra-reader consensus as the ground-truth label. The MPH test dataset has not been described previously. Compared to the development dataset, the internal test dataset was characterized by better spatial resolution (resulting in higher striatum-to-background contrast) and less statistical noise.

4 Methods

TODO - δ What happens in chapter

It concludes with an examination of the significant performance metrics utilized for the evaluation of the research outcomes.

4.1 Software Tools and Libraries

The project was built on *Python 3.10*. A variety of widely adopted open-source libraries were used in the project. *NumPy* was utilized to perform efficient array operations and numerical calculations. The *NIBabel* library was used for reading and writing of medical image data stored in the Neuroimaging Informatics Technology Initiative (NIfTI) file format. *PyTorch*, a widely adopted deep learning framework, was employed for building and training the neural networks. The *Torchvision* package provided the machine learning models and image transformation capabilities utilized in this project. *Pandas* was used for efficient structured data manipulation and analysis. *Matplotlib* and *Seaborn* were employed for the creation of customized data visualizations. The *Scikit-Learn* library provided machine learning models and model evaluation tools utilized in this project, whereas *Scipy* was used for data interpolation.

The seeds of the random number generators in each package were initialized to ensure reproducibility.

4.2 Development Data Preparation

In the following, the data preparation techniques applied to the development dataset are explained in detail.

4.2.1 Data Preprocessing

Individual DAT-SPECT images were stereotactically normalized to the anatomical space of the Montreal Neurological Institute (MNI) using the Normalize tool of the Statistical Parametric Mapping software package (version SPM12) and a set of custom DAT-SPECT templates representative of normal and different levels of Parkinson-typical reduction of striatal uptake as target (Apostolova et al., 2023). The voxel size of the stereotactically normalized images was $2 \times 2 \times 2 \text{ mm}^3$. Intensity normalization was achieved by voxelwise scaling to the individual 75th percentile of the voxel intensity in a reference region comprising the whole brain without striata, thalamus, medial temporal lobe, brainstem, cerebellum, and ventricles (Kupitz et al., 2014). The resulting images are distribution volume (DVR) images. A 2-dimensional transversal DVR slab of 12mm thickness and 91x109 pixels with 2 mm edge length was obtained by averaging 6 transversal slices through the striatum (Buchert et al., 2006).

4.2.2 Data Augmentation

Data augmentation was applied to the development dataset to increase the heterogeneity of the data. To enhance robustness across various attenuation correction and

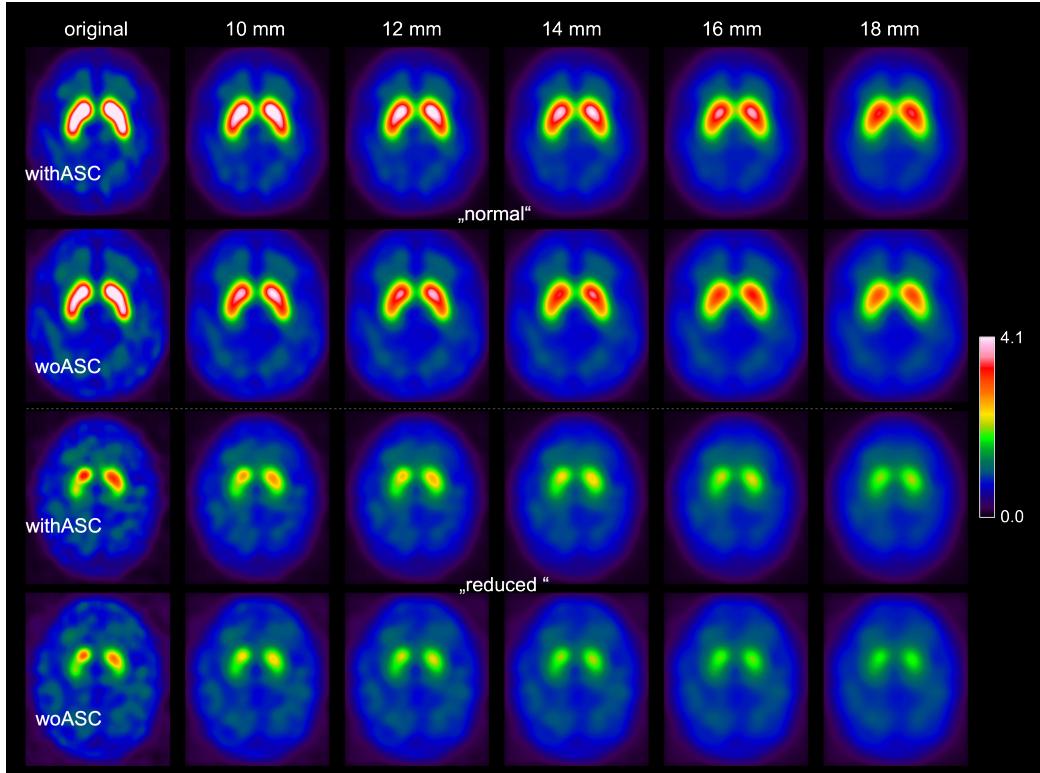


Figure 1: DVR slabs for two sample cases from the development dataset, a healthy control case (above) and a PD case with reduced availability of DAT in the striatum (below). The two cases are presented in 12 different versions. In each version, attenuation and scatter corrections are either applied ('withASC') or not applied ('woASC'). Also, for each version, isotropic 3-dimensional Gaussian kernel smoothing with different FWHM values (10, 12, 14, 16, 18mm) was either performed or not performed ('original').

scatter correction methods, each image was generated in a version with and without attenuation and scatter corrections applied (Schiebler et al., 2023). Also, 3D-smoothing of the 3-dimensional SPECT images in MNI space was performed before computing the 2-dimensional slabs as an augmentation technique. A 3-dimensional isotropic Gaussian kernel with various Full Width at Half Maximum (FWHM) values (FWHM = 10, 12, 14, 16, 18mm) was used for the smoothing. Thereby an augmented dataset of 20,880 images in total was constructed based on 1,740 cases. Two representative cases augmented using the described techniques are depicted in Figure 1.

4.2.3 Dataset Splitting

Ten distinct random splits were created from the augmented development dataset, resulting in ten different combinations of training, validation, and test sets for the conducted experiments. In each random split, the data distribution was as follows: 60% for the training set, 20% for the validation set, and 20% for the test set.

While splitting the data it was ensured that all augmented images associated with a given patient were put into the same subset. Thereby, randomization into training, validation and test set was performed on the level of patients rather than on the level of single images. Thereby inter-subset leakage of images from the same patient was avoided.

4.3 Univariate benchmark: Specific Binding Ratio

The unilateral [¹²³I]FP-CIT specific binding ratio (SBR) was used as a benchmark classification method. Here, the SBR in left and right putamen was obtained by hottest voxels (HV) analysis of the stereotactically normalized DVR image using large unilateral putamen masks predefined in MNI space (Wenzel et al., 2019). The unilateral hottest voxel SBR was calculated as

$$\text{HV-SBR}_{\text{unilateral}} = \left(\frac{1}{K_{10\text{ml}}} \sum_k \hat{I}_{k,\text{ROI}} \right) - 1, \quad (1)$$

where $\hat{I}_{k,\text{ROI}}$ represent the *normalized* voxel intensities of the $K_{10\text{ml}}$ hottest voxels (i.e., voxels with the highest intensity) comprising a total volume of 10 ml within the unilateral putamen ROI in the DVR image. The voxel intensities of the hottest voxels are normalized to the 75th percentile of the voxel intensities in the reference region associated with non-specific binding (Wenzel et al., 2019). The minimum of the HV-SBR values from the left and right hemispheres of the brain was used for the analysis. An in-depth elaboration on SBR analysis can be found in Wenzel et al. (2019).

The SBR-based classifier was obtained for each of the random splits ($n = 10$) as follows. First, the optimal cutoff on the SBR was determined in the validation set using ROC analysis and the Youden criterion (Youden, 1950). The determined optimal cutoff was then used as the decision boundary between normal control cases and Parkinson's disease and evaluated on the test set of the development dataset. Also, the determined cutoff was evaluated on the PPMI and MPH test datasets described in Section 3.2. As a result, 10 optimal cutoffs on the SBR were determined and evaluated.

4.4 Multivariate benchmark: PCA-enhanced Random Forest

As a further benchmark, a random forest classifier was trained on PCA-transformed features of the training set of the development dataset.

To be comparable with CNN-based approaches, first, a 91x91 pixel square-shaped region of interest was defined within the 91x109 pixel DVR slab, and each development data image (of each subset) was cropped to this region. The square-shaped

region was determined by cropping an equal number of pixels from the top and bottom of the DVR slab along its second dimension (anterior-posterior direction).

Then a PCA model with 10 principle components was initialized and fit to the training set to obtain the principle components of the training set. The determined principle components were used to transform the training set into a lower-dimensional space, where each image was represented by a 10-vector characterizing the expression of each principal component. An example of the principle components of the training set for one of the random splits is depicted in Figure 2.

The training data transformed by the principle components was then used to train a random forest classifier with 100 decision trees. As hyperparameters, the Gini impurity was used to assess split quality, with a minimum of 2 samples required to split an internal node and 1 sample needed at a leaf node. The trained random forest classifier was evaluated on the test split of the development dataset for each of the 10 random splits. In addition, the trained model was tested on the PPMI and MPH datasets described in Section 3.2.

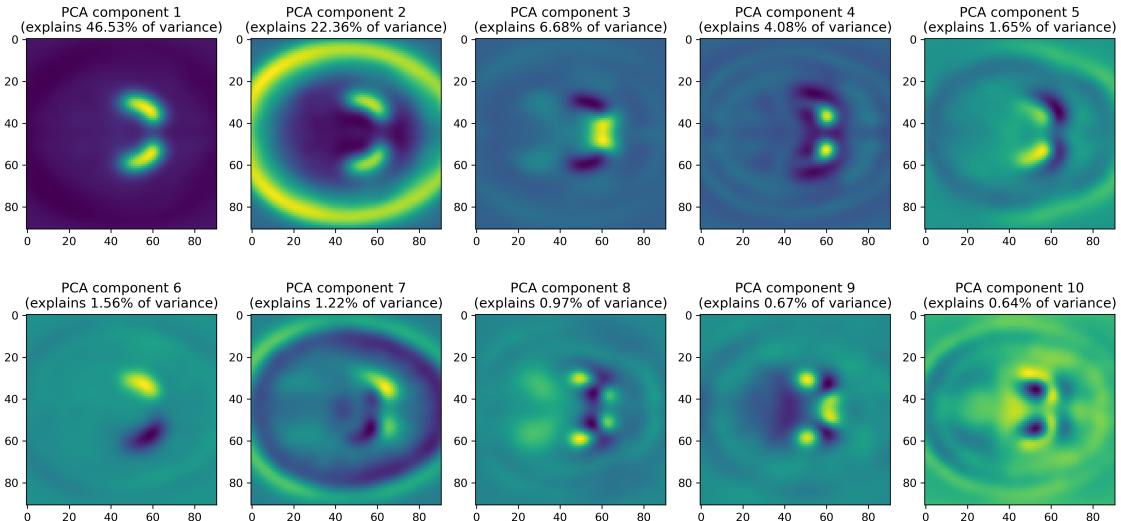


Figure 2: Principle components of the training set (development dataset) for the first random split.

4.5 CNN-based classification

The models of CNN-based classifiers were based on a Residual Network (ResNet) architecture. More precisely, the *ResNet-18* (He et al., 2015) model architecture consisting of 18 layers was used as basis. The non-pretrained weights of the ResNet-18 were used as initial weights. The ResNet-18 architecture expects input tensors of size (3, 224, 224), denoting images with 3 channels and spatial dimensions of 224 by 224 pixels. Since the development data has one color channel, the architecture was modified to expect one input channel at its first convolutional layer. Also the

dimensions of the last fully-connected layer of the architecture were modified to produce one output node in the output layer. The modified ResNet-18 model is depicted in Figure 3. To produce a probabilistic model output within the range of 0 to 1, the sigmoid function was applied to the output layer of the model.

Further development data preprocessing was performed to comply with the spatial input dimensions required by the model architecture. First, a 91x91 pixel square-shaped region of interest was defined within the 91x109 pixel DVR slab, and each development data image (of each subset) was cropped to this region. The cropping to a square shape was performed to preserve the aspect ratio while doing the subsequent upscaling. The square-shaped region was determined by cropping an equal number of pixels from the top and bottom of the DVR slab along its height dimension. Then the square-shaped images were resized to the target image size of 224x224 pixels using bicubic interpolation.

The CNN-based approaches were trained for 20 epochs using a batch size of 64. For the MVT and RLT approaches (described in Section 4.5.1) the Binary Cross Entropy (BCE) loss was employed for optimization, whereas for the Regression approach (described in Section 4.5.2) the Mean Squared Error (MSE) loss function was used. The Adam optimization algorithm was utilized with an initial learning rate of 0.0001. During the training of the model, the weights of the best epoch were saved for subsequent evaluations. Each CNN model was trained and evaluated separately for each of the 10 random splits of the development dataset. Additionally, the trained models were evaluated on the PPMI and MPH test datasets described in Section 3.2. No attempt was made to adapt the CNN models trained in the development dataset for these independent test datasets.

4.5.1 MVT-based and RLT-based methods

When training a CNN using the BCE loss function, one has to provide the ground truth label of each instance to the optimization algorithm. Given that each instance in the development data is labeled by three independent readers, a selection strategy must be determined. The following two label selection strategies are used for training the CNNs: Majority Vote training (MVT) and Random Label training (RLT). The labels chosen using one of the two strategies are then used, together with the model predictions, to compute the BCE loss.

Majority vote training involved selecting the label that received the majority of votes from the readers as the ground truth label. Since there are three available labels, a majority is reached when two out of the three readers agree on a particular label (e.g., the normal case (NC)). During the model training phase, the majority vote strategy was employed to select the labels for both the training and validation data instances.

In contrast to MVT, random label training involved choosing a random label from the three available options as the ground truth label. The seed of the random number generator (responsible for the random selection) is set only once at the start

Layer (type:depth-idx)	Output Shape	Param #
ResNet18	[64, 1]	--
└ResNet: 1-1	[64, 1]	--
└Conv2d: 2-1	[64, 64, 112, 112]	3,136
└BatchNorm2d: 2-2	[64, 64, 112, 112]	128
└ReLU: 2-3	[64, 64, 112, 112]	--
└MaxPool2d: 2-4	[64, 64, 56, 56]	--
└Sequential: 2-5	[64, 64, 56, 56]	--
└BasicBlock: 3-1	[64, 64, 56, 56]	73,984
└BasicBlock: 3-2	[64, 64, 56, 56]	73,984
└Sequential: 2-6	[64, 128, 28, 28]	--
└BasicBlock: 3-3	[64, 128, 28, 28]	230,144
└BasicBlock: 3-4	[64, 128, 28, 28]	295,424
└Sequential: 2-7	[64, 256, 14, 14]	--
└BasicBlock: 3-5	[64, 256, 14, 14]	919,040
└BasicBlock: 3-6	[64, 256, 14, 14]	1,180,672
└Sequential: 2-8	[64, 512, 7, 7]	--
└BasicBlock: 3-7	[64, 512, 7, 7]	3,673,088
└BasicBlock: 3-8	[64, 512, 7, 7]	4,720,640
└AdaptiveAvgPool2d: 2-9	[64, 512, 1, 1]	--
└Linear: 2-10	[64, 1]	513
Total params: 11,170,753		
Trainable params: 11,170,753		
Non-trainable params: 0		
Total mult-adds (G): 111.03		
Input size (MB): 12.85		
Forward/backward pass size (MB): 2543.32		
Params size (MB): 44.68		
Estimated Total Size (MB): 2600.85		

Figure 3: Architecture of the CNN-based classification models.

of the algorithm and is not reset between the model training epochs. Thereby a different label could be chosen as the ground truth label for each distinct training epoch. Here the random label selection strategy is applied both to the training and validation data.

4.5.2 Regression-based method

The regression-based approach aimed to incorporate the uncertainty regarding the ground truth label into the training algorithm. The ground-truth label was derived from the combination of the three available labels, resulting in a floating-point number. Each of the following states of certainty about the label was mapped to a distinct floating-point valued ground-truth label: *all readers agree on ‘normal’* (ground-truth label: 0.0), *majority of readers (two out of three) agree on ‘normal’* (ground-truth label: 1.0/3.0), *majority of readers (two out of three) agree on ‘reduced’* (ground-truth label: 2.0/3.0) and *all readers agree on ‘reduced’* (ground-truth

label: 1.0). This mapping of available labels to the ground-truth label was used for both the training and validation data during the model training phase.

During model training, the loss was computed using the Mean Square Error loss function which aims to minimize the mean of the squared differences between the model predictions and the ground-truth labels.

4.6 Evaluation Metrics and Procedure

In the following the performance metrics used for the evaluation of the different classification methods are explained in more detail.

First the mean \pm SD (standard deviation) of the following measures were calculated across the different random splits for each classification approach and subset (training, validation and testing) given a cutoff: Area under the Receiver Operating Characteristic curve (AUC-ROC), balanced accuracy, (overall) accuracy, sensitivity, specificity, Positive Predictive Value (PPV) and Negative Predictive Value (NPV). The natural cutoff of 0.5 was used for each classification approach except the SBR method. For the SBR method the optimal cutoff was determined using the Youden criterion (Youden, 1950) and was used for calculating the measures. In the test set of the development dataset (for each random split), the majority vote was used as ground truth in all cases.

Second, for each element within a set of considered percentages of inconclusive cases in the validation set (PIncVal) the corresponding inconclusive interval was determined. Inconclusive cases were defined as cases predicted within an inconclusive interval (bounded by lower and upper bound), while conclusive cases were those predicted outside this interval. The determination of the inconclusive interval was exclusively performed using the validation set for each random split and classification approach independently. The set of PIncVal values considered ranged from 0.2% to 20.0%, increasing in increments of 0.2%. For each target PIncVal value the lower and upper bounds of the inconclusive interval were independently determined in such a way that there was the same number of inconclusive cases (± 1 case) below and above the pre-defined cutoff. For the CNN-based classification methods and the multivariate benchmark the natural cutoff of 0.5 was used, whereas for the SBR-based univariate benchmark the optimal cutoff on the SBR was used.

To assess the stability of the determined inconclusive interval over the proportion of inconclusive cases, the determined upper and lower bounds (mean \pm SD across the 10 random splits) of the inconclusive interval were plotted against the corresponding PIncVal (%). The rate at which the lower (upper) bound decreases (increases) over the PIncVal reflects the density of inconclusive cases within a certain region of PIncVal. Specifically, higher function gradients indicate lower concentration of predictions, and vice versa. Also, a higher standard deviation indicates that the stable inconclusive interval determination is less robust within a certain region of PIncVal.

As the main performance metric (regarding the primary hypothesis of the project) we propose the area under the curve of mean balanced accuracy (AUC-bACC, %) on conclusive test cases as a function of the mean percentage of inconclusive test cases (mean PIncObs, %). More precisely the relative AUC-bACC (%) normalized to the maximum achievable area was used for the comparison. To obtain the relative AUC-bACC, first, the mean balanced accuracy function was interpolated using cubic spline interpolation. Then the area under the mean balanced accuracy curve was computed using the trapezoidal rule and then normalized to the maximum achievable area (100% balanced accuracy * (20% - 0.2% inconclusive cases)). The evaluation of each classification method with respect to this metric was conducted on the test set of the development dataset as well as on the independent datasets PPMI and MPH.

As a further metric, the mean \pm SD percentage of observed inconclusive cases in the test set (PIncObs, %) was plotted against the PIncVal (%). A mean of PIncObs(PIncVal) near the identity line is an indicator for a similar prediction distribution for validation set and test set on average. In case the mean of PIncObs(PIncVal) consistently lies over (under) the identity line the supposed prediction certainty on the test set, on average, is lower (higher) than on the validation set. Also a lower standard deviation of PIncObs over PIncVal indicates that PIncObs is less sensitive to the randomness of the inconclusive intervals across random splits. In particular, a lower standard deviation of PIncObs allows for a more reliable calculation of the main performance metric.

5 Evaluation

The preceding chapters have detailed the research methodology, data collection and sources, and the application of classification techniques to address the research questions posed in this study.

This chapter embarks on the evaluation of the research results, focusing on the performance and effectiveness of the methods employed, and the attainment of the research objectives.

The structure of this chapter has been designed to systematically lead readers through the assessment process. The chapter commences with the examination of the performance results obtained for the baseline methods. The core of this chapter subsequently unveils the results for the experimental methods evaluated using various test datasets and compared to the baseline performance. These findings are presented using performance summary tables for statistical measures and graphical representations. The chapter culminates with a comparative analysis, which seeks to assess and contrast the effectiveness and limitations of the research methods employed.

5.1 Baseline Performance

In this section, the performance of the SBR method is thoroughly evaluated. Furthermore the outcomes for the multivariate PCA-RFC method are also provided as additional baseline. The objective of this evaluation is twofold: to comprehend the inherent capabilities of the baseline methods, SBR and PCA-RFC, and to establish a clear point of reference for the CNN-based methodologies.

5.1.1 SBR Method

Test set of development dataset

Binary classification performance Table 1 presents the quantitative performance (balanced accuracy, accuracy, sensitivity, specificity, PPV, NPV and AUC-ROC) of the SBR-based classification on the particular subset of the Development dataset. In the evaluation process, the optimal SBR cutoff value of 0.703 (with a variation of ± 0.009 across random splits) was employed. The SBR method consistently achieved around 93% in balanced accuracy, accuracy, sensitivity, specificity, PPV and NPV on the validation set, with a variance between 0.5-1.5% across random splits. The performance on training and test sets is also similarly around 93% with respect to all the metrics. The comparable sensitivity and specificity imply a well-balanced SBR model that identifies both positive and negative cases similarly well. The SBR model achieved a stable AUC-ROC of 0.983 ± 0.002 .

Table 1: Evaluation of the SBR method on Development dataset (SBR cutoff mean \pm SD: 0.703 ± 0.009).

	train set	validation set	test set
Balanced Accuracy	0.936 ± 0.003	0.929 ± 0.008	0.935 ± 0.007
Accuracy	0.936 ± 0.003	0.930 ± 0.008	0.935 ± 0.007
Sensitivity	0.934 ± 0.006	0.924 ± 0.005	0.930 ± 0.014
Specificity	0.937 ± 0.003	0.935 ± 0.015	0.939 ± 0.012
PPV	0.933 ± 0.005	0.929 ± 0.014	0.930 ± 0.015
NPV	0.938 ± 0.005	0.930 ± 0.004	0.938 ± 0.018
AUC-ROC	0.983 ± 0.002		

Determined inconclusive intervals Figure 4 illustrates the determined lower and upper bounds on the SBR as a function of the percentages of inconclusive cases in the validation set (development dataset), along with the mean \pm SD of the optimal cutoff. Corroborating the intuitive expectation, the width of the inconclusive interval expands as the percentage of inconclusive cases increases. The close resemblance in slopes between the upper and lower bound functions indicates a nearly identical distribution of predictions both below and above the cutoff.

Transferability of inconclusive intervals In Figure 5 the correspondence between the percentage of inconclusive cases in the validation set (development dataset) and the mean \pm SD percentage of observed inconclusive cases (PIncObs) in the test set (development dataset) is demonstrated. The plot illustrates that the deviation of the mean PIIncObs in the test set from the identity line is negligibly small. This can be attributed to the nearly identical distribution of data in both the test and validation sets (due to random splitting) which results in a similar distribution of SBR model predictions.

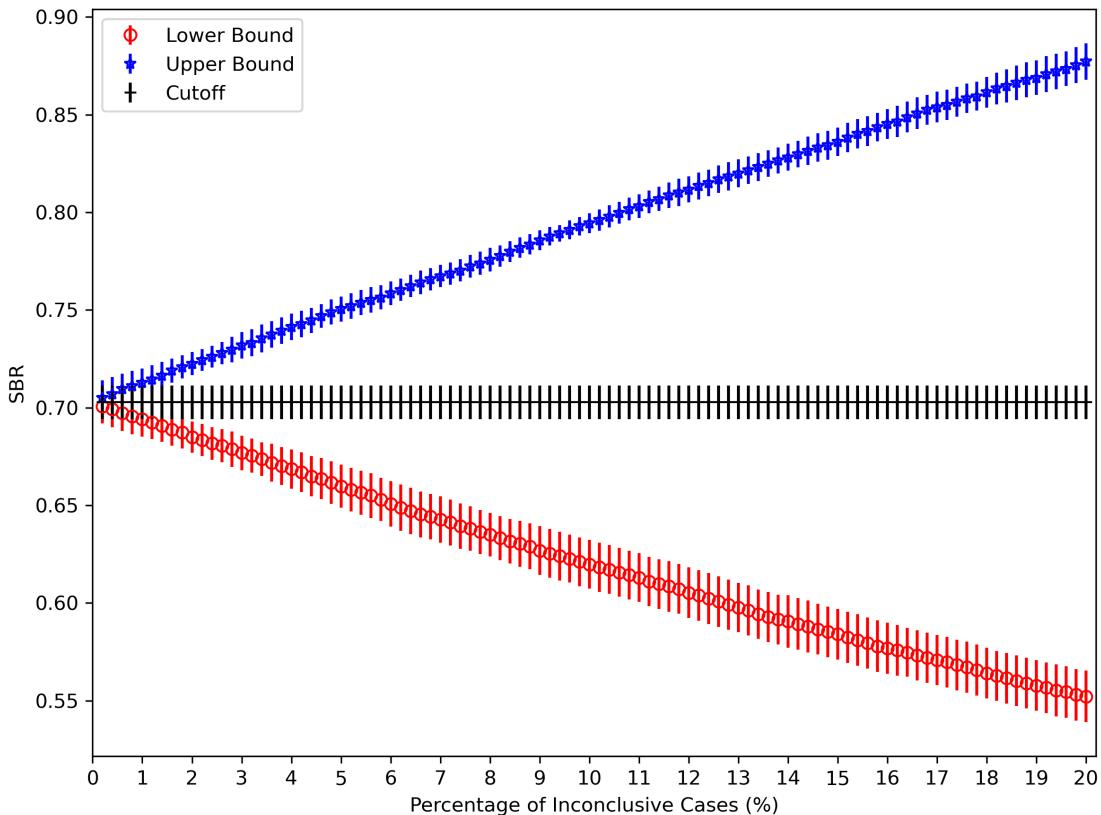


Figure 4: Evaluation of the SBR method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.

AUC-bACC for balanced accuracy over PIIncObs Figure 6a shows the balanced accuracy (mean \pm SD across random splits) on both conclusive and inconclusive cases as a function of the mean PIIncObs in the test set (development dataset). The balanced accuracy on inconclusive cases is not part of further performance analysis and comparison due to the emphasis on the balanced accuracy on conclusive cases as the basis for the main metric of this work. The balanced accuracy (mean \pm SD) on conclusive cases over the mean PIIncObs is depicted with enhanced clarity and precision in Figure 6b. The mean of the balanced accuracy rises from approximately

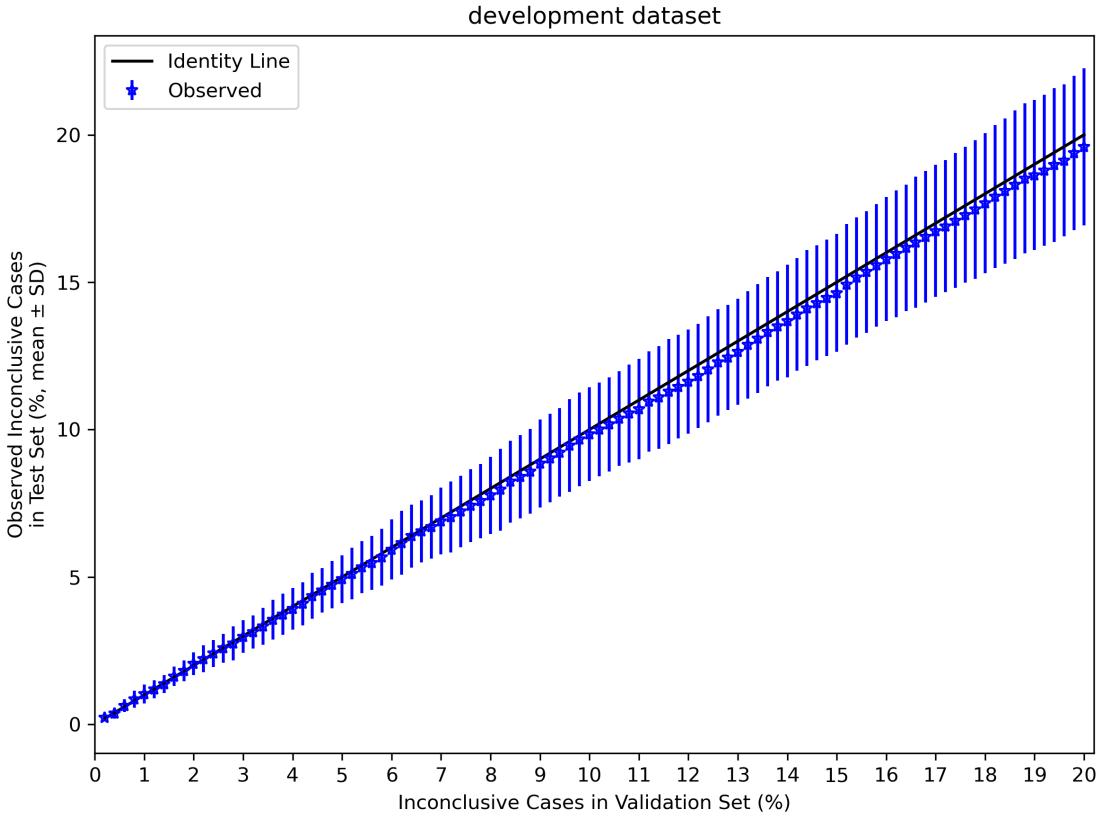


Figure 5: Evaluation of the SBR method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).

94% when there are around 1% of inconclusive cases in the test set to about 98% when there are around 20% of inconclusive cases in the test set. The SBR baseline method attains a relative AUC-bACC of 96.38% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the test set of the development dataset.

PPMI dataset The results obtained from evaluating the SBR method on the PPMI dataset are depicted in Figure 7. The mean \pm SD percentage of inconclusive cases observed (PIncObs) in the PPMI test dataset over the percentage of inconclusive cases in the validation set (development dataset) is consistently below the identity line, which can be seen in Figure 7a. That implies that, on average, the supposed prediction certainty on PPMI dataset is higher than on validation set (development dataset), regardless of the prediction accuracy. The balanced accuracy on conclusive cases over the mean PIncObs is shown in Figure 7b. The mean of the balanced accuracy rises from approximately 96% when there are around 1% of inconclusive cases in the PPMI test set to about 99% when there are around 20%

of inconclusive cases in the PPMI test set. The SBR baseline method achieved a relative AUC-bACC of 97.51% for the mean balanced accuracy on conclusive cases over the mean PIIncObs in the PPMI test dataset.

MPH dataset The evaluation of the SBR method on the MPH dataset is shown in Figure 8. Figure 8a demonstrates the mean \pm SD percentage of inconclusive cases observed (PIIncObs) in the MPH test dataset over the percentage of inconclusive cases in the validation set (development dataset). Similar as in case of the PPMI dataset, here the PIIncObs in the MPH test dataset is also consistently below the identity line and thus the supposed prediction certainty on MPH dataset is higher than on validation set. The balanced accuracy on conclusive cases over the mean percentage of observed inconclusive cases (PIIncObs) is shown in Figure 8b. The mean of the balanced accuracy rises from approximately 91.5% when there are around 1% of inconclusive cases in the MPH test set to about 95% when there are around 20% of inconclusive cases in the MPH test set. The SBR baseline method achieved a relative AUC-bACC of 93.46% for the mean balanced accuracy on conclusive cases over the mean PIIncObs in the MPH test dataset.

5.1.2 PCA-RFC Method

Test set of development dataset

Binary classification performance Table 2 presents the quantitative performance (balanced accuracy, accuracy, sensitivity, specificity, PPV, NPV and AUC-ROC) of the PCA-RFC classification on the particular subset of the Development dataset. In the evaluation process, the natural sigmoid cutoff value of 0.5 was employed. The PCA-RFC method achieved around 96% in balanced accuracy, accuracy, sensitivity, specificity, PPV and NPV on the validation and test set, with a variance around 1% across random splits. The SBR model achieved a stable AUC-ROC of 0.994 ± 0.002 .

Table 2: Evaluation of the PCA-RFC method on Development dataset. For evaluation, the natural sigmoid cutoff 0.5 was used.

	train set	validation set	test set
Balanced Accuracy	1.000 ± 0.000	0.963 ± 0.010	0.966 ± 0.006
Accuracy	1.000 ± 0.000	0.963 ± 0.010	0.966 ± 0.006
Sensitivity	1.000 ± 0.000	0.957 ± 0.012	0.962 ± 0.010
Specificity	1.000 ± 0.000	0.969 ± 0.011	0.969 ± 0.009
PPV	1.000 ± 0.000	0.966 ± 0.012	0.965 ± 0.010
NPV	1.000 ± 0.000	0.961 ± 0.012	0.966 ± 0.011
AUC-ROC		0.994 ± 0.002	

Determined inconclusive intervals Figure 9 illustrates the determined lower and upper bounds on the probabilistic output as a function of the percentages of inconclusive cases in the validation set (development dataset), along with the natural cutoff of 0.5. The width of the inconclusive interval expands as the percentage of inconclusive cases increases and the visual resemblance in shape and slope between the curve a similar distribution of predictions both below and above the cutoff.

Transferability of inconclusive intervals The correspondence between the percentage of inconclusive cases in the validation set (development dataset) and the mean \pm SD percentage of observed inconclusive cases (PIncObs) in the test set (development dataset) is demonstrated in Figure 10. The deviation of the mean PIncObs in the test set from the identity line is small which can be attributed to the nearly identical distribution of data in both the test and validation sets (due to random splitting).

AUC-bACC for balanced accuracy over PIncObs Figure 11a shows the balanced accuracy (mean \pm SD across random splits) on both conclusive and inconclusive cases as a function of the mean PIncObs in the test set (development dataset). The balanced accuracy on inconclusive cases is not part of performance analysis and comparison due to the emphasis on the balanced accuracy on conclusive cases as the basis for the main metric of this work. The balanced accuracy (mean \pm SD) on conclusive cases over the mean PIncObs is depicted with enhanced clarity and precision in Figure 11b. The mean of the balanced accuracy rises from approximately 97% when there are around 1% of inconclusive cases in the test set to about 99.5% when there are around 20% of inconclusive cases in the test set. As a result, the PCA-RFC baseline method achieved a relative AUC-bACC of 98.71% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the test set of the development dataset.

PPMI dataset The following results were obtained when evaluating the PCA-RFC method on the PPMI dataset. Figure 12a shows the mean \pm SD percentage of inconclusive cases observed (PIncObs) in the PPMI test dataset over the percentage of inconclusive cases in the validation set (development dataset). The function is consistently above the identity line. Therefore, on average, the supposed prediction certainty of the PCA-RFC method on PPMI dataset is lower than on validation set, regardless of the prediction accuracy. The balanced accuracy on conclusive cases over the mean PIncObs is presented in Figure 12b. The mean of the balanced accuracy rises from approximately 98% when there are around 1% of inconclusive cases in the PPMI test set to about 99.5% when there are around 20% of inconclusive cases in the PPMI test set. The PCA-RFC baseline method achieved a relative AUC-bACC of 99.12% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the PPMI test dataset.

MPH dataset The evaluation of the PCA-RFC method on the MPH dataset shows the following results. In Figure 13a the mean \pm SD percentage of inconclusive cases observed (PIncObs) in the MPH test dataset over the percentage of inconclusive cases in the validation set (PIncVal) is illustrated. Here the mean of PIncObs in the MPH test dataset is also consistently above the identity line and its deviation from the identity line increases over PIncVal. Therefore the supposed prediction certainty on MPH dataset is lower than on validation set (development data). The balanced accuracy on conclusive cases over the mean PIncObs is shown in Figure 13b. The mean of the balanced accuracy rises from approximately 90.5% when there are around 1% of inconclusive cases in the MPH test set to about 94% when there are around 19% of inconclusive cases in the MPH test set. As a result, the PCA-RFC baseline method achieved a relative AUC-bACC of 92.42% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the MPH test dataset.

5.2 Experimental Methods Performance

This section presents the performance results for the CNN-based classification approaches separately and compares them to the results obtained by the baseline approaches. First the results for the CNN-MVT method are presented, whereafter the CNN-RLT method is evaluated. Finally the findings for the CNN-Regression method are showcased.

5.2.1 CNN-MVT Method

Test set of development dataset

Binary classification performance The quantitative performance results of the CNN-MVT classification on the particular subset of the Development dataset are presented in Table 3. In the evaluation process, the natural sigmoid cutoff value of 0.5 was employed. The CNN-MVT method achieved around 96.4% in sensitivity, 97.6% in specificity and a balanced accuracy of 97.0%, with a variance between 1-2% across random splits, on the test set. The performance results on the validation set are very similar. The method achieved a stable AUC-ROC of 0.996 ± 0.002 .

Determined inconclusive intervals In Figure 14 the determined lower and upper bounds on the probabilistic sigmoid output are plotted as a function of the percentages of inconclusive cases in the validation set (development dataset), along with the natural cutoff 0.5. The visual resemblance in shape and slope between the upper and lower bound curves indicates a similar distribution of predictions both below and above the cutoff. The width of the inconclusive interval increases more rapidly as the percentage of inconclusive cases increases when compared to the PCA-RFC baseline method. That implies that the CNN-MVT method produces relatively less inconclusive cases than the PCA-RFC baseline.

Table 3: Evaluation of the CNN-MVT method on Development dataset. For evaluation, the natural sigmoid cutoff 0.5 was used.

	train set	validation set	test set
Balanced Accuracy	0.999±0.003	0.970±0.014	0.970±0.008
Accuracy	0.999±0.003	0.970±0.014	0.970±0.008
Sensitivity	1.000±0.000	0.963±0.010	0.964±0.015
Specificity	0.997±0.006	0.976±0.023	0.976±0.013
PPV	0.997±0.006	0.975±0.024	0.972±0.018
NPV	1.000±0.000	0.966±0.010	0.968±0.014
AUC-ROC		0.996±0.002	

Transferability of inconclusive intervals The correspondence between the percentage of inconclusive cases in the validation set (development dataset) and the mean±SD percentage of observed inconclusive cases (PIncObs) in the test set (development dataset) is illustrated in Figure 15. As for the baseline cases, the deviation of the mean PIncObs in the test set from the identity line is small which can be attributed to the nearly identical distribution of data in both the test and validation sets (due to random splitting).

AUC-bACC for balanced accuracy over PIncObs The balanced accuracy (mean±SD) on conclusive cases over the mean PIncObs in the test set (development dataset) is depicted in Figure 16b. The mean of the balanced accuracy rises from about 97% when there are around 1% of inconclusive cases in the test set to about 99.5% when there are around 20% of inconclusive cases in the test set. As a result, the CNN-MVT method achieved a relative AUC-bACC of 98.95% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the test set of the development dataset. The achieved relative AUC-bACC is approximately 2.5% higher than that of the SBR baseline method and around 0.2% higher than the PCA-RFC baseline.

PPMI dataset The following results were obtained when evaluating the CNN-MVT method on the PPMI dataset. Figure 17a depicts the mean±SD percentage of inconclusive cases observed (PIncObs) in the PPMI test dataset over the percentage of inconclusive cases in the validation set (PIncVal) of development dataset. For lower PIncVal the corresponding PIncObs in the PPMI test dataset are similar. However as PIncVal increases (corresponding to increasing inconclusive intervals) the supposed prediction certainty on PPMI dataset decreases when compared to the certainty on validation set, on average. The balanced accuracy on conclusive cases over the mean PIncObs is illustrated in Figure 17b. The mean of the balanced accuracy rises from approximately 98% when there are around 1% of inconclusive cases in the PPMI test set to about 99.5% when there are around 20% of inconclusive cases in the PPMI test set. The CNN-MVT method achieved a relative AUC-bACC of 99.23% for the mean balanced accuracy on conclusive cases over the mean PIncObs

in the PPMI test dataset. The achieved relative AUC-bACC is approximately 1.7% higher than that of the SBR baseline method and around 0.1% higher than the PCA-RFC baseline.

MPH dataset The evaluation of the CNN-MVT method on the MPH dataset produced the following results. Figure 18a presents the mean \pm SD percentage of inconclusive cases observed (PIncObs) in the MPH test dataset over the percentage of inconclusive cases in the validation set (PIncVal) of development dataset. The mean of PIncObs in the MPH test dataset is consistently above the identity line and the deviation from the identity line increases over PIncVal. The standard deviation of PIncObs also increases over PIncVal. When compared to the mean PIncObs of the SBR baseline the mean PIncObs of the CNN-MVT method is higher which indicates that CNN-MVT is supposedly less certain about the MPH set predictions than the SBR method. Also the PIncObs of the CNN-MVT has a much higher standard deviation compared to the PIncObs of the SBR method. The balanced accuracy on conclusive cases over the mean percentage of observed inconclusive cases (PIncObs) is depicted in Figure 18b. The mean of the balanced accuracy increases from approximately 95% when there are around 1% of inconclusive cases in the MPH test set to about 96% when there are around 20% of inconclusive cases in the MPH test set. As a result, the CNN-MVT method achieved a relative AUC-bACC of 95.73% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the MPH test dataset. The achieved relative AUC-bACC is approximately 2.3% higher than that of the SBR baseline method and around 3.3% higher than the PCA-RFC baseline.

5.2.2 CNN-RLT Method

Test set of development dataset

Binary classification performance The quantitative performance results of the CNN-RLT classification on the particular subset of the Development dataset are presented in Table 4. In the evaluation process, the natural sigmoid cutoff value of 0.5 was employed. The CNN-RLT method achieved around 96.1% in sensitivity, 98.5% in specificity and a balanced accuracy of 97.3%, with a variance between 0.5-1.5% across random splits, on the test set. The performance results on the validation set are similar. The method achieved a stable AUC-ROC of 0.994 ± 0.002 .

Determined inconclusive intervals Figure 19 shows the determined lower and upper bounds on the probabilistic sigmoid output as a function of the percentages of inconclusive cases in the validation set (development dataset), along with the natural cutoff 0.5. The upper bound curve increases and saturates faster than the lower bound curve with a lower variance across the random splits. First this suggests a disparity in the distribution of predictions below and above the cutoff point. Also

Table 4: Evaluation of the CNN-RLT method on Development dataset. For evaluation, the natural sigmoid cutoff 0.5 was used.

	train set	validation set	test set
Balanced Accuracy	0.982±0.003	0.967±0.008	0.973±0.005
Accuracy	0.982±0.003	0.968±0.008	0.973±0.005
Sensitivity	0.980±0.008	0.951±0.013	0.961±0.014
Specificity	0.983±0.008	0.984±0.005	0.985±0.010
PPV	0.983±0.009	0.982±0.006	0.982±0.012
NPV	0.981±0.008	0.956±0.012	0.966±0.013
AUC-ROC	0.994±0.002		

the determination of stable lower bounds across the random splits is more difficult than the determination of stable upper bounds. When compared to the PCA-RFC baseline method the width of the inconclusive interval increases more rapidly as the percentage of inconclusive cases increases. That implies that the CNN-RLT method tends to produce relatively less inconclusive cases than the PCA-RFC baseline.

Transferability of inconclusive intervals The correspondence between the percentage of inconclusive cases in the validation set (development dataset) and the mean±SD percentage of observed inconclusive cases (PIncObs) in the test set (development dataset) is depicted in Figure 20. As for the baseline cases, the deviation of the mean PIncObs in the test set from the identity line is small which can be attributed to the nearly identical distribution of data in both the test and validation sets (due to random splitting).

AUC-bACC for balanced accuracy over PIncObs The balanced accuracy (mean±SD) on conclusive cases over the mean PIncObs in the test set (development dataset) is depicted in Figure 21b. The mean of the balanced accuracy rises from about 97.5% when there are around 1% of inconclusive cases in the test set to about 99.5% when there are around 20% of inconclusive cases in the test set. As a result, the CNN-RLT method achieved a relative AUC-bACC of 99.02% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the test set of the development dataset. The achieved relative AUC-bACC is approximately 2.6% higher than that of the SBR baseline method and around 0.3% higher than the PCA-RFC baseline.

PPMI dataset The following results were obtained when evaluating the CNN-RLT method on the PPMI dataset. Figure 22a shows the mean±SD percentage of inconclusive cases observed (PIncObs) in the PPMI test dataset over the percentage of inconclusive cases in the validation set (PIncVal) of the development dataset.

Here the mean PIncObs in the PPMI test dataset deviates only slightly from the identity line. For PIncVal less than 6% the mean PIncObs is slightly below the

identity line. Subsequently the mean PIncObs rises slightly above the identity line with an increasing standard deviation of PIncObs. The balanced accuracy on conclusive cases over the mean percentage of observed inconclusive cases (PIncObs) is presented in Figure 22b. The mean of the balanced accuracy rises from approximately 98.5% when there are around 1% of inconclusive cases in the PPMI test set to about 99.5% when there are around 20% of inconclusive cases in the PPMI test set. The CNN-RLT method achieved a relative AUC-bACC of 99.31% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the PPMI test dataset. The achieved relative AUC-bACC is approximately 1.8% higher than that of the SBR baseline method and around 0.2% higher than the PCA-RFC baseline.

MPH dataset The evaluation of the CNN-RLT method on the MPH dataset produced the following results. Figure 23a illustrates the mean \pm SD percentage of inconclusive cases observed (PIncObs) in the MPH test dataset over the percentage of inconclusive cases in the validation set (PIncVal) (development dataset). Here the mean of the PIncObs in the MPH test dataset is slightly above the identity line and the standard deviation increases over the PIncVal. The balanced accuracy on conclusive cases over the mean PIncObs is depicted in Figure 23b. The mean of the balanced accuracy slightly increases from approximately 95% when there are around 1% of inconclusive cases in the MPH test set to about 96% when there are around 20% of inconclusive cases in the MPH test set. As a result, the CNN-RLT method achieved a relative AUC-bACC of 96.12% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the MPH test dataset. The achieved relative AUC-bACC is approximately 2.7% higher than that of the SBR baseline method and around 3.7% higher than the PCA-RFC baseline.

5.2.3 CNN-Regression Method

Test set of development dataset

Binary classification performance The quantitative performance results of the CNN-Regression classification on the particular subset of the Development dataset are presented in Table 5. In the evaluation process, the natural sigmoid cutoff value of 0.5 was employed. The CNN-Regression method achieved around 96.1% in sensitivity, 98.5% in specificity and a balanced accuracy of 97.5%, with a standard deviation between 0.6-1.1% across random splits, on the test set. The performance results on the validation set are a balanced accuracy of 97.7%, a sensitivity of 98.3% and a specificity of 97.2%. The method achieved a stable AUC-ROC of 0.998 ± 0.001 .

Determined inconclusive intervals Figure 24 presents the determined lower and upper bounds on the probabilistic sigmoid output as a function of the percentages of inconclusive cases in the validation set (PIncVal) of the development dataset, along with the natural cutoff 0.5. Similar to the CNN-RLT method, here the upper

Table 5: Evaluation of the CNN-Regression method on Development dataset. For evaluation, the natural sigmoid cutoff 0.5 was used.

	train set	validation set	test set
Balanced Accuracy	0.982+/-0.003	0.977+/-0.006	0.975+/-0.006
Accuracy	0.980+/-0.003	0.977+/-0.007	0.976+/-0.006
Sensitivity	1.000+/-0.000	0.983+/-0.009	0.961+/-0.011
Specificity	0.963+/-0.005	0.972+/-0.009	0.988+/-0.008
PPV	0.960+/-0.005	0.967+/-0.011	0.986+/-0.009
NPV	1.000+/-0.000	0.985+/-0.008	0.967+/-0.010
AUC-ROC		0.998+/-0.001	

bound curve increases and saturates slightly faster than the lower bound curve with a lower variance across the random splits. This suggests a slight disparity in the distribution of predictions below and above the cutoff point. Since both the upper and lower bound functions exhibit a significant standard deviation across the random splits the determination of stable lower and upper bounds is difficult. When compared to the PCA-RFC baseline method the width of the inconclusive interval increases more rapidly over the PIncVal. Therefore the CNN-Regression method also tends to produce relatively less inconclusive cases than the PCA-RFC baseline.

Transferability of inconclusive intervals The correspondence between the PIncVal of the development dataset and the mean \pm SD percentage of observed inconclusive cases (PIncObs) in the test set of the development dataset is depicted in Figure 25. As for the baseline cases, the deviation of the mean PIncObs in the test set from the identity line is small which can be attributed to the nearly identical distribution of data in both the test and validation sets (due to random splitting).

AUC-bACC for balanced accuracy over PIncObs The balanced accuracy (mean \pm SD) on conclusive cases over the mean PIncObs in the test set (development dataset) is depicted in Figure 26b. The mean of the balanced accuracy rises from about 98% when there is a PIncObs of 1% in the test set to about 99.5% when there is a PIncObs around 20% in the test set. As a result, the CNN-Regression method achieved a relative AUC-bACC of 99.23% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the test set of the development dataset. The achieved relative AUC-bACC is approximately 2.8% higher than that of the SBR baseline method and around 0.5% higher than the PCA-RFC baseline.

PPMI dataset The following results were obtained when evaluating the CNN-Regression method on the PPMI dataset. Figure 27a illustrates the mean \pm SD percentage of inconclusive cases observed (PIncObs) in the PPMI test dataset over the percentage of inconclusive cases in the validation set (PIncVal) of the development dataset. Here for lower PIncVal values (less than 5%) the corresponding mean PIncObs in the PPMI test dataset is near the identity line. However for higher PIncVAL

values the mean of PIncObs increasingly rises above the identity line and the standard deviation of PIncObs increases strongly. The balanced accuracy on conclusive cases over the mean PIncObs is presented in Figure 27b. The mean of the balanced accuracy rises from approximately 98.5% when there are around 1% of inconclusive cases in the PPMI test set to about 99.5% when there are around 20% of inconclusive cases in the PPMI test set. The CNN-Regression method achieved a relative AUC-bACC of 99.38% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the PPMI test dataset. The achieved relative AUC-bACC is approximately 1.9% higher than that of the SBR baseline method and around 0.3% higher than the PCA-RFC baseline.

MPH dataset The evaluation of the CNN-Regression method on the MPH dataset produced the following results. Figure 28a demonstrates the mean \pm SD percentage of inconclusive cases observed (PIncObs) in the MPH test dataset over the percentage of inconclusive cases in the validation set (PIncVal) (development dataset). Here the mean of the PIncObs in the MPH test dataset is above the identity line and deviates stronger from the identity line as the PIncVal increases. Also the standard deviation of the PIncObs is high and increases over the increasing PIncVal. The balanced accuracy on conclusive cases over the mean PIncObs is depicted in Figure 28b. The mean of the balanced accuracy slightly increases from approximately 95% when there are around 1% of inconclusive cases in the MPH test set to about 96.5% when there are around 20% of inconclusive cases in the MPH test set. As a result, the CNN-Regression method achieved a relative AUC-bACC of 96.24% for the mean balanced accuracy on conclusive cases over the mean PIncObs in the MPH test dataset. The achieved relative AUC-bACC is approximately 2.8% higher than that of the SBR baseline method and around 3.8% higher than the PCA-RFC baseline.

5.3 Comparative Performance Analysis

In this section, a summary comparison of the performance between the baseline and experimental methods is presented. The comparison focuses on two aspects: transferability of inconclusive intervals (in both validation and test sets) and the AUC-bACC of balanced accuracy on conclusive cases across varying percentages of observed inconclusive cases (PIncObs). To support the analysis visually one comparison figure is used for each aspect tested on a specific dataset. The comparison is carried out for the test set of the development data, the PPMI dataset and the MPH dataset, respectively.

5.3.1 Performance on test set of development dataset

Figure 29 provides a comparison of the transferability of the inconclusive intervals from the validation set to the test set (development data) along the baseline and experimental methods. For each considered method the mean of the percentage of

observed inconclusive cases (PIncObs) hardly deviates from the identity line. The similarity in data distribution of the validation and test set due to random splitting is an explanation for that. The standard deviation of PIIncObs is also similarly low across the methods. Thus the PIIncObs is hardly affected by the randomness of the inconclusive intervals across random splits for each method. The mean of PIIncObs can be reliably used for the calculation of the main metric of this work compared in the following.

In Figure 30, the performance comparison of the methods on the test set (development dataset) concerning the main metric of this work, the relative AUC-bACC for the mean balanced accuracy on conclusive cases over the mean PIIncObs in the test set, is shown. In general, the CNN-based methods outperform both baseline methods, the SBR method and PCA-RFC. The highest performance is achieved by the CNN-Regression (relative AUC: 99.23%) method whereas the lowest AUC-bACC is that for the SBR-based method (relative AUC: 96.38%). The CNN-RLT method achieved slightly higher performance than the CNN-MVT.

5.3.2 Performance on PPMI dataset

Figure 31 shows a comparison of the transferability of the inconclusive intervals from the validation set to the PPMI test dataset along the baseline and experimental methods. The percentage of observed inconclusive cases (PIIncObs) of CNN-based methods shows a higher standard deviation compared to the baseline methods. A possible explanation for that is the higher sensitivity of the CNN-based methods to the variability in inconclusive intervals across random splits. Also the mean of PIIncObs deviates stronger from the identity line (above identity line) for the CNN-MVT and CNN-Regression methods compared to the baseline methods and CNN-RLT. That indicates that, on average, the CNN-MVT and CNN-Regression methods are supposedly less certain about the PPMI dataset predictions compared to the other methods.

The performance comparison of the methods concerning the relative AUC-bACC for the mean balanced accuracy on conclusive cases over the PIIncObs in the PPMI dataset is depicted in Figure 32. The CNN-based methods outperform both baseline methods, the SBR method and PCA-RFC. The highest performance is achieved by the CNN-Regression (relative AUC: 99.38%) method whereas the lowest AUC-bACC is that for the SBR-based method (relative AUC: 97.51%). The CNN-RLT method achieved slightly higher performance (relative AUC: 99.31%) than the CNN-MVT (relative AUC: 99.23%).

5.3.3 Performance on MPH dataset

Figure 33 illustrates a comparison of the transferability of the inconclusive intervals from the validation set to the MPH test dataset along the baseline and experimental methods. As for the PPMI dataset, on the MPH dataset the percentage of observed inconclusive cases (PIIncObs) of CNN-based methods shows a higher

standard deviation compared to the baseline methods. The higher sensitivity of the CNN-based methods to the variability in inconclusive intervals across random splits is a possible explanation. Here also the mean of PIncObs deviates stronger from the identity line (above identity line) for the CNN-MVT and CNN-Regression methods compared to the CNN-RLT. That indicates that, on average, the CNN-MVT and CNN-Regression methods are supposedly less certain about the MPH dataset predictions compared to the CNN-RLT method. However the highest deviation of the mean PIncObs from the identity line shows the baseline PCA-RFC method and thus shows the lowest supposed certainty about the MPH dataset predictions, on average.

In Figure 34 the performance comparison of the methods concerning the relative AUC-bACC for the mean balanced accuracy on conclusive cases over the PIncObs in the MPH dataset is presented. As for the other test datasets, CNN-based methods outperform both baseline methods, the SBR method and PCA-RFC. The highest performance is achieved by the CNN-Regression (relative AUC: 96.24%) method whereas the lowest AUC-bACC is that for the baseline PCA-RFC method (relative AUC: 92.42%). The CNN-RLT method achieved slightly higher performance (relative AUC: 96.12%) than the CNN-MVT (relative AUC: 95.73%).

5.4 Conclusion

A concluding overview of the performance of the methods on different test datasets concerning the relative AUC-bACC for the mean balanced accuracy on conclusive cases over the PIncObs is summarized in Figure 35. The AUC-bACC performance on the MPH test dataset is lower than that on the development data test set and PPMI test dataset across all considered classification methods. The CNN-based methods outperform the baseline methods, most remarkably on the MPH dataset. The best performance is achieved by the CNN-Regression method, followed by the CNN-RLT and CNN-MVT methods.

6 Discussion

6.1 Interpretation of Results

The primary research hypothesis proposed that a CNN trained using the RLT strategy is more effective in identifying inconclusive cases compared to the MVT strategy. The main metric proposed to evaluate and compare the performance was the AUC-bACC of balanced accuracy on conclusive cases over the percentage of observed inconclusive cases (PIncObs) in the test set. The performance results summarized in Figure 35 demonstrate that the RLT strategy leads to slightly better performance (higher by 0.05 – 0.1%) on the development data test set and the external PPMI test dataset. On the internal MPH test dataset the performance of the CNN using RLT strategy is higher by 0.4% compared to the MVT strategy. Since the MPH

dataset cases exhibit better spatial resolution than the development dataset and PPMI dataset cases and thus are more difficult to classify the clear superiority of the RLT strategy on this test dataset is particularly remarkable. Hence the findings support the primary hypothesis of this work.

Two secondary hypotheses were also put to test in this work. The first secondary hypothesis was the superiority in performance of CNN-based classification methods compared to conventional baseline methods in terms of the main metric (AUC-bACC of balanced accuracy over PIncObs). The results shown in Figure 35 show that CNN-based classification methods have a consistent AUC-bACC performance advantage over the SBR univariate baseline method across all test datasets, achieving approximately a 2% higher AUC-bACC result. The AUC-bACC performance of the multivariate baseline method PCA-RFC closely approaches that of the CNN-based methods on the development data test set and PPMI test dataset. However, on the more difficult MPH test dataset, the AUC-bACC performance of the PCA-RFC method is over 3% lower compared to that of the CNN-based methods. The other secondary hypothesis stated that the CNN-based classification methods are more robust regarding varying image characteristics and thus exhibit better generalizability compared to baseline methods. The CNN-based methods outperform both baseline methods on the PPMI and MPH test datasets which differ in image characteristics. The performance advantage of the CNN-based methods is more prominent on the more difficult MPH dataset particularly compared to the PCA-RFC method. Consequently, both secondary hypotheses are confirmed by the findings.

6.2 Practical Implications

The findings of the study have practical implications for the classification of DAT-SPECT images. First, the study shows that random label selection as a ground-truth label selection strategy can lead to better performance results compared to the majority vote strategy when training a CNN classifier for Parkinson’s disease diagnosis based on DAT-SPECT. Second, the mean AUC-bACC of balanced accuracy on conclusive cases over the mean percentage of observed inconclusive cases (PIncObs) can be used as a metric for comparison of the actual underlying decision confidence of different binary classification models. The metric decouples the classification model performance from the arbitrarily chosen inconclusive interval bounds. Also, the balanced accuracy on conclusive cases over PIncObs can be used to decide for the model that produces the least relative amount of inconclusive cases for a specified target balanced accuracy. The applicability of the metric extends beyond DAT-SPECT classification to general binary classification problems. Third, the results once again confirm the superiority of CNN-based methods for DAT-SPECT classification compared to the widely adopted SBR method in clinical practice.

6.3 Limitations of the Study

6.3.1 Metric

There are several limitations to be considered that may impact the validity of the applied methods and results. The main metric used to compare the model performance, the AUC-bACC of mean balanced accuracy over mean PIncObs, depends on a set of inconclusive intervals determined within the validation set of the development dataset for each classification method and randomization individually. Since the balanced accuracy and PIncObs are averaged across the results for each random split the reliability of the metric is affected by the standard deviation of both variables across the random splits. To enhance the reliability of the metric a higher number of random splits can be used. Also, the resolution of the balanced accuracy over PIncObs decreases as the density of test set predictions around the cutoff increases in comparison to the validation set predictions. Finally the metric may be less intuitively understandable and requires more expertise when interpreting the results when compared to standard metrics such as balanced accuracy and AUC-ROC.

6.3.2 Site-Specific Development Data

To increase the heterogeneity of the training data and enhance the generalizability of the classification models the development dataset was augmented as described in Section 4.2.2. However the results, presented in Figure 35, show that AUC-bACC performance on the MPH test across all considered classification methods, is significantly lower than on the test set of the development dataset. The reason for that can be a site-specific bias introduced by training on data from one site (development dataset). Future research efforts should aim to include development data from multiple sites to enhance the overall robustness of the classification methods and results.

7 Conclusion

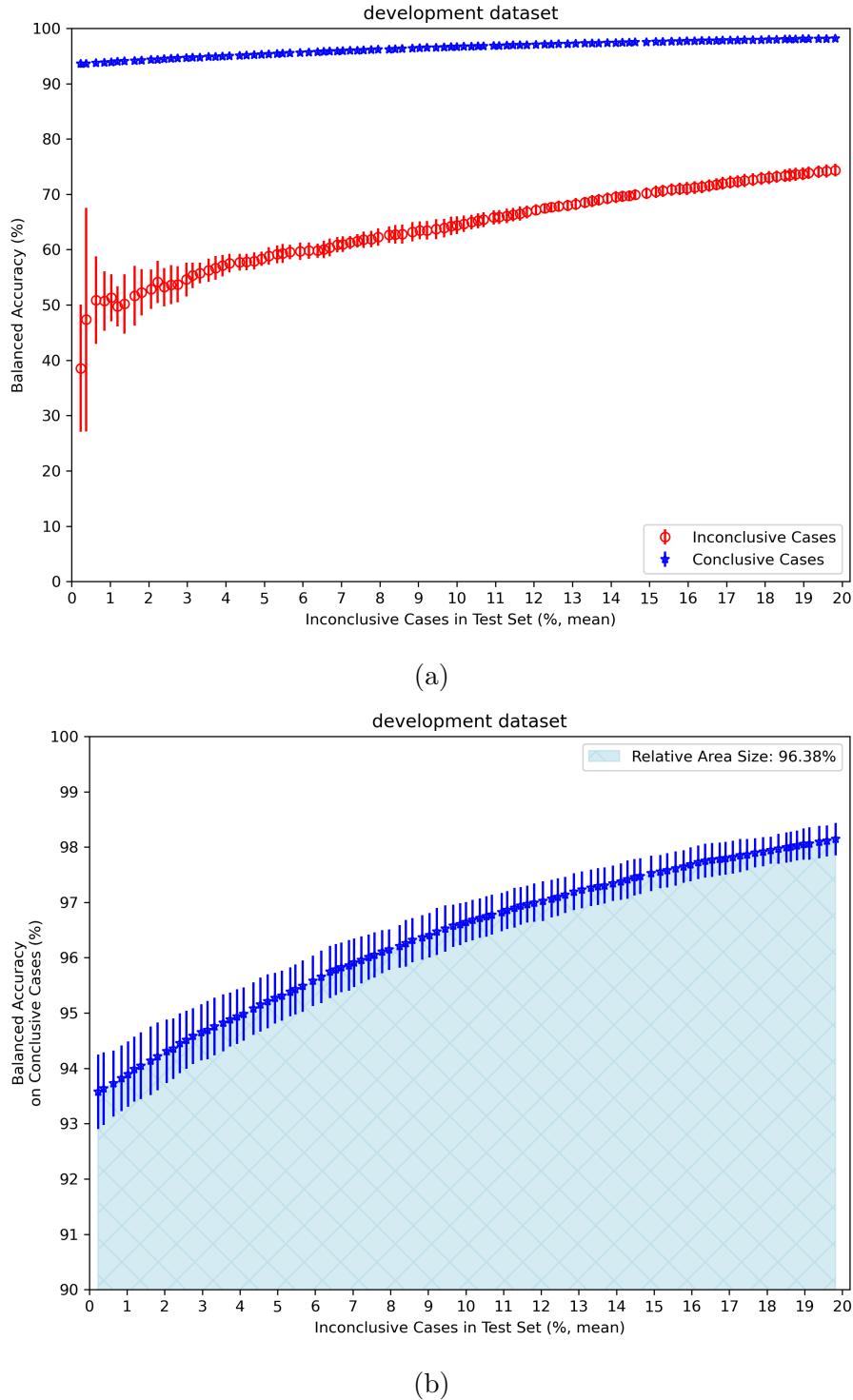
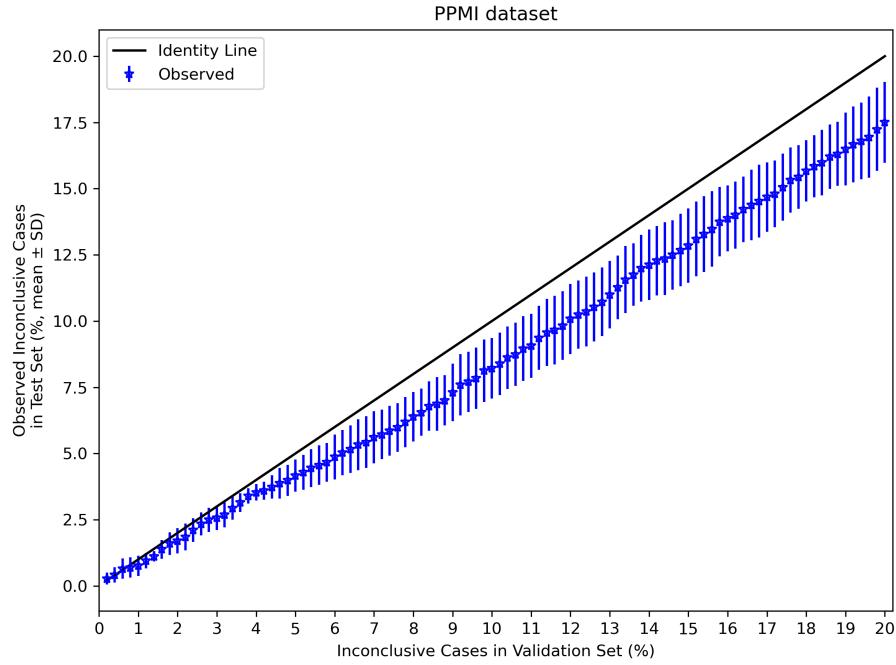
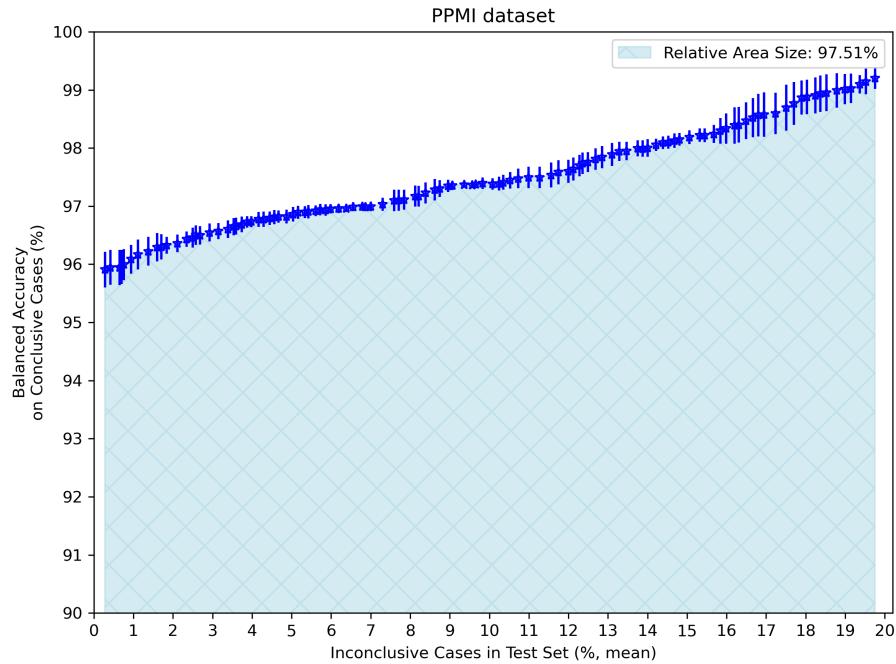


Figure 6: Evaluation of the SBR method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). For better illustration the area under the mean of the balanced accuracy is highlighted.

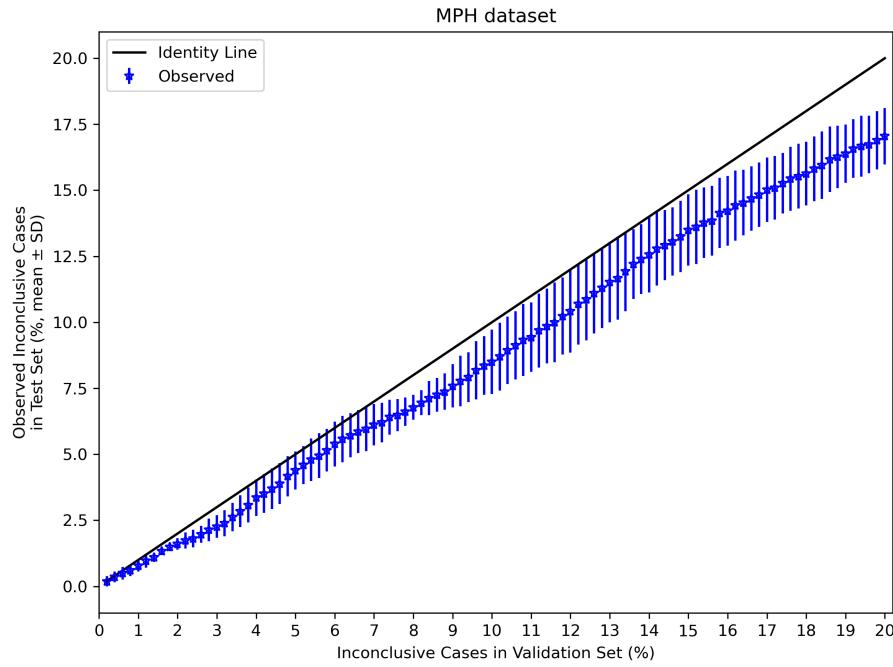


(a) Observed percentage of inconclusive cases in the PPMI dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).

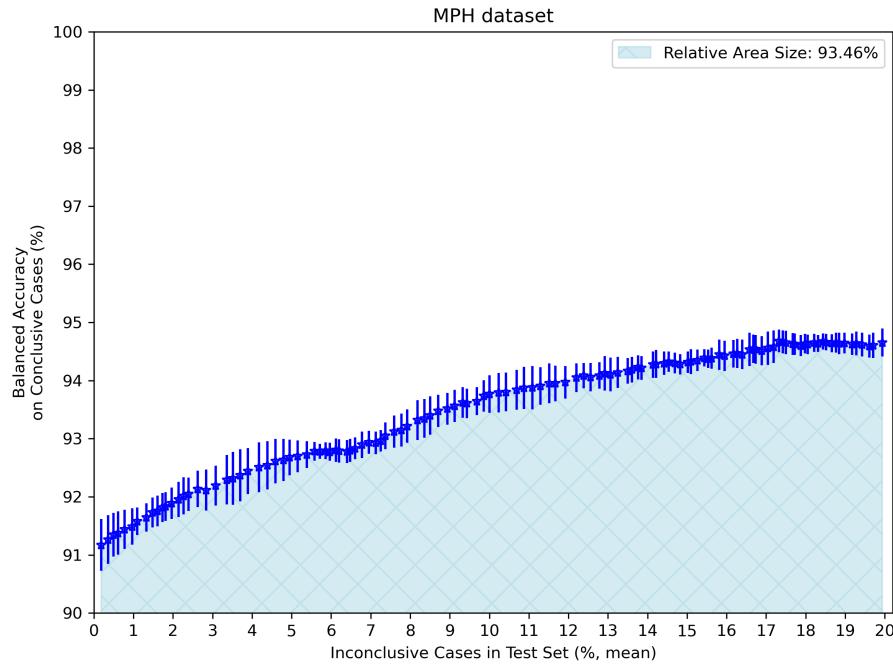


(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the PPMI dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 7: Evaluation of the SBR method on PPMI dataset.



(a) Observed percentage of inconclusive cases in the MPH dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).



(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the MPH dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 8: Evaluation of the SBR method on MPH dataset.

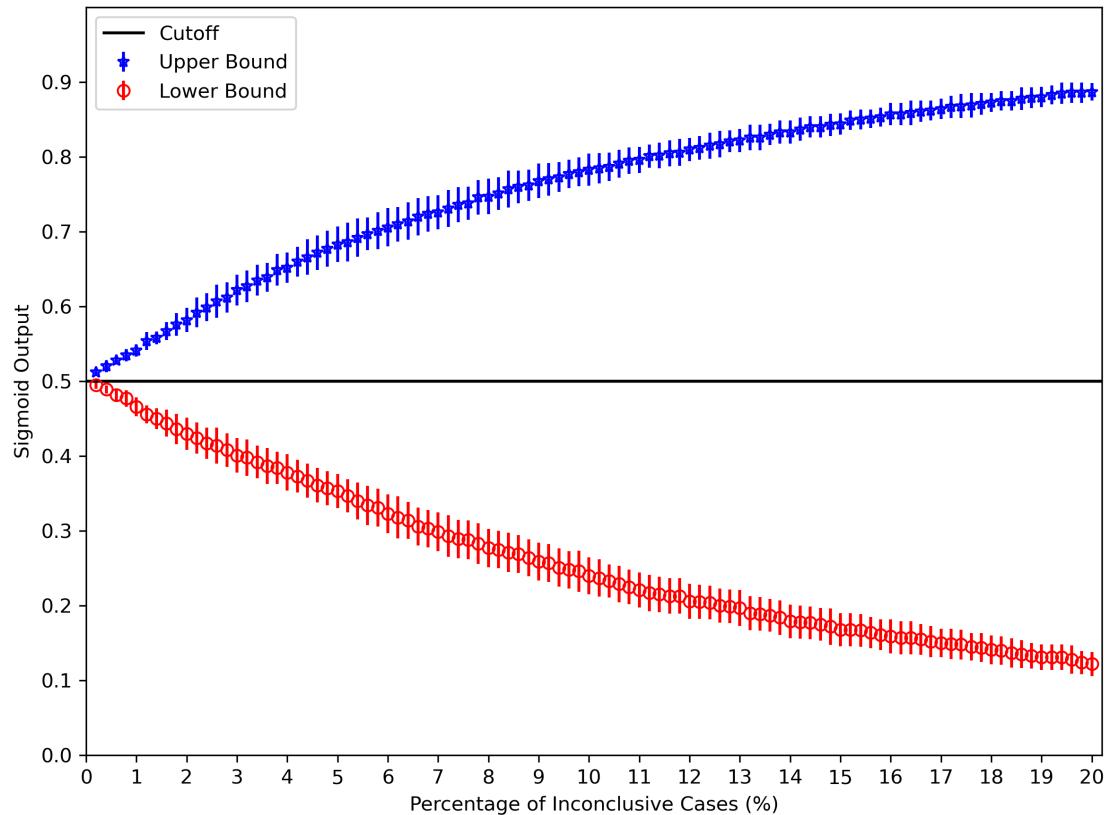


Figure 9: Evaluation of the PCA-RFC method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.

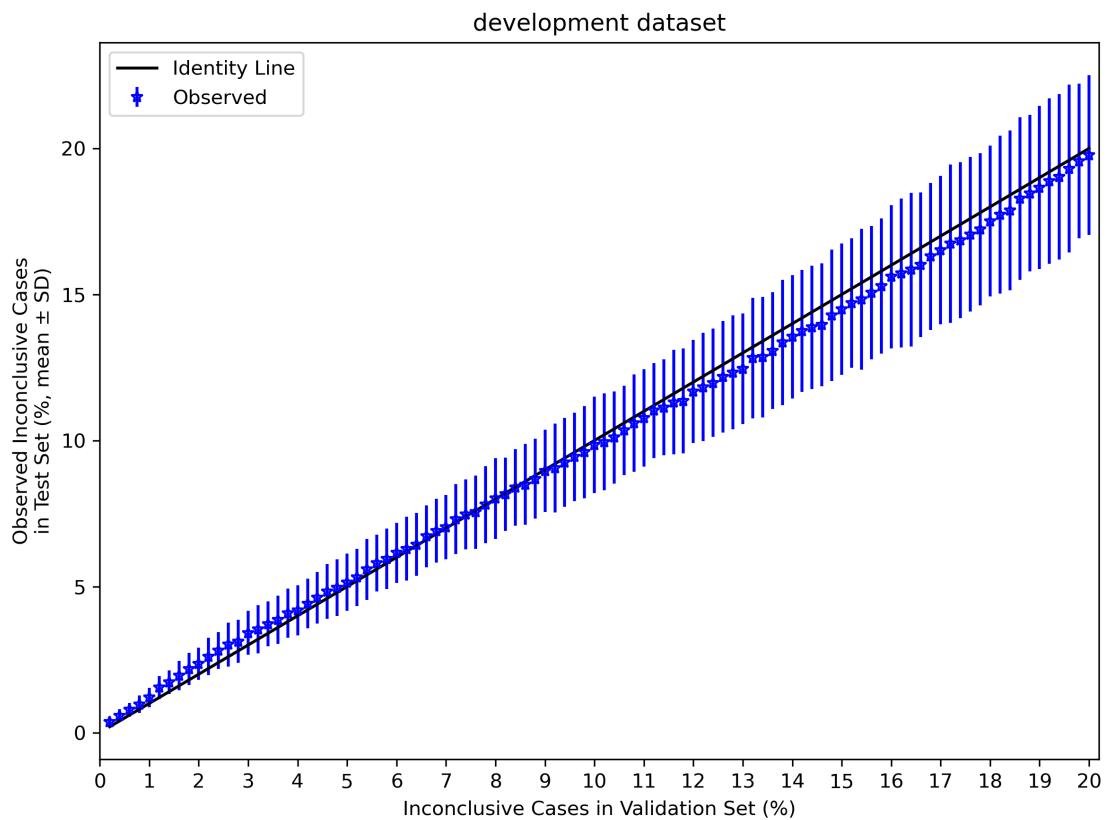


Figure 10: Evaluation of the PCA-RFC method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).

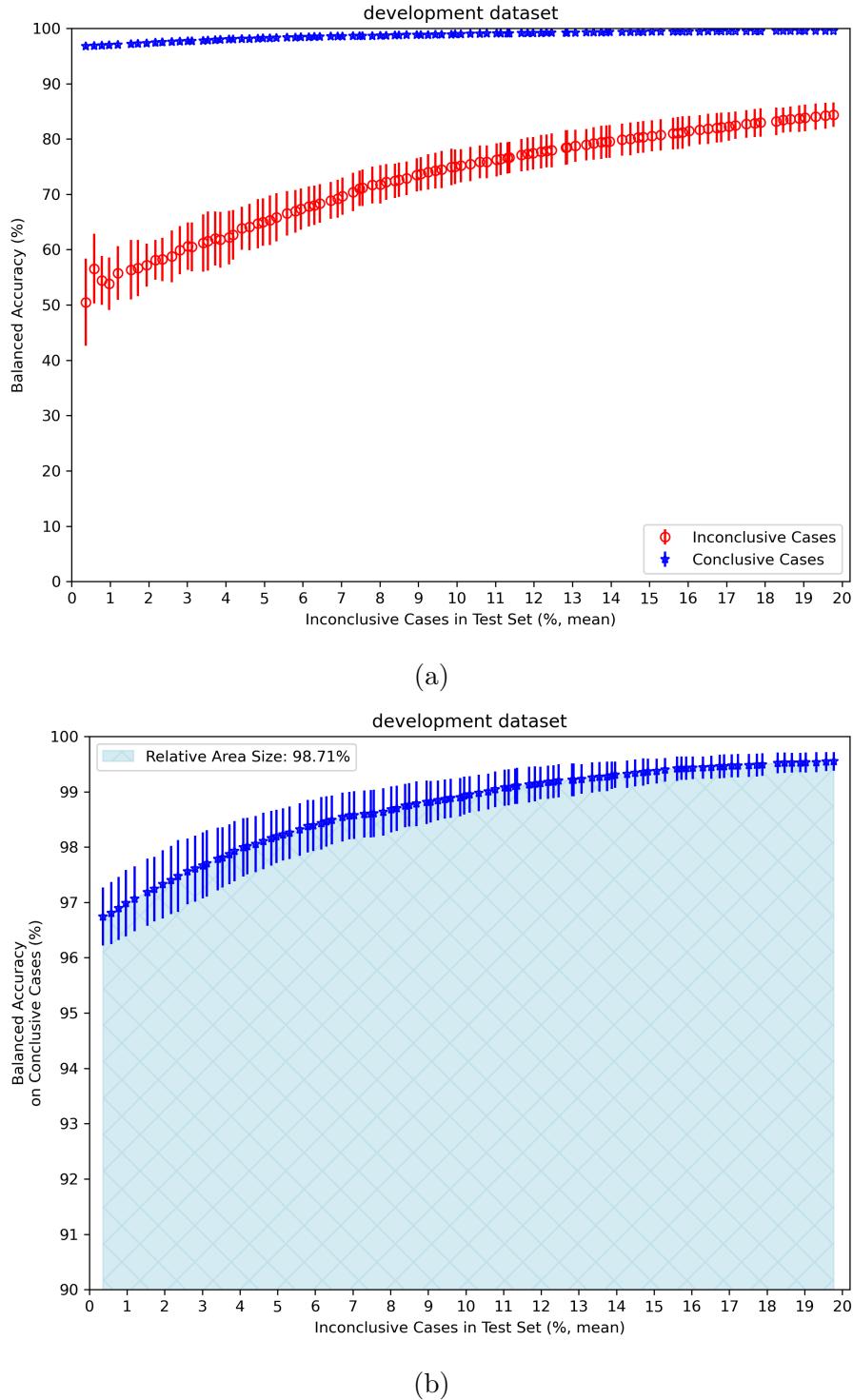
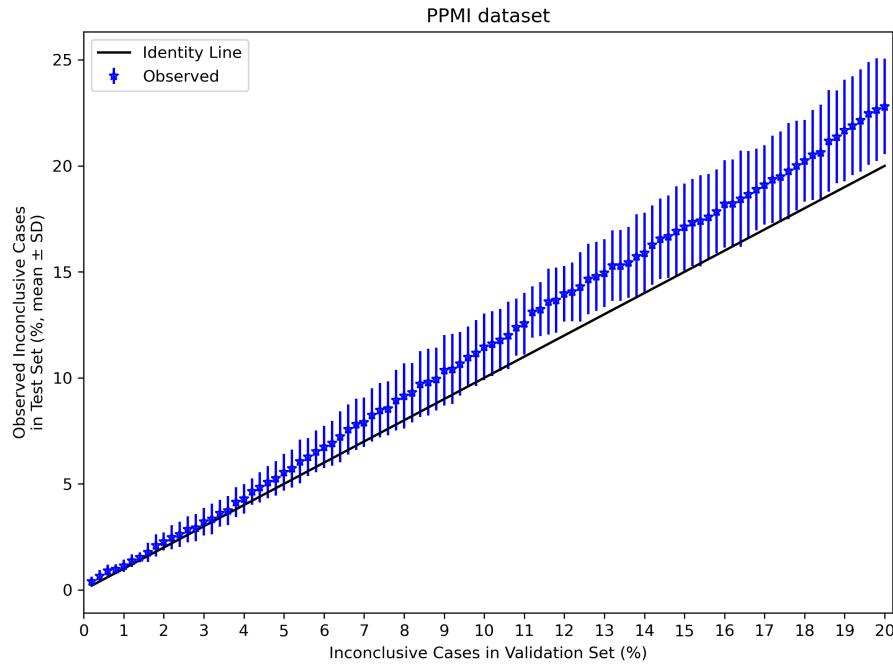
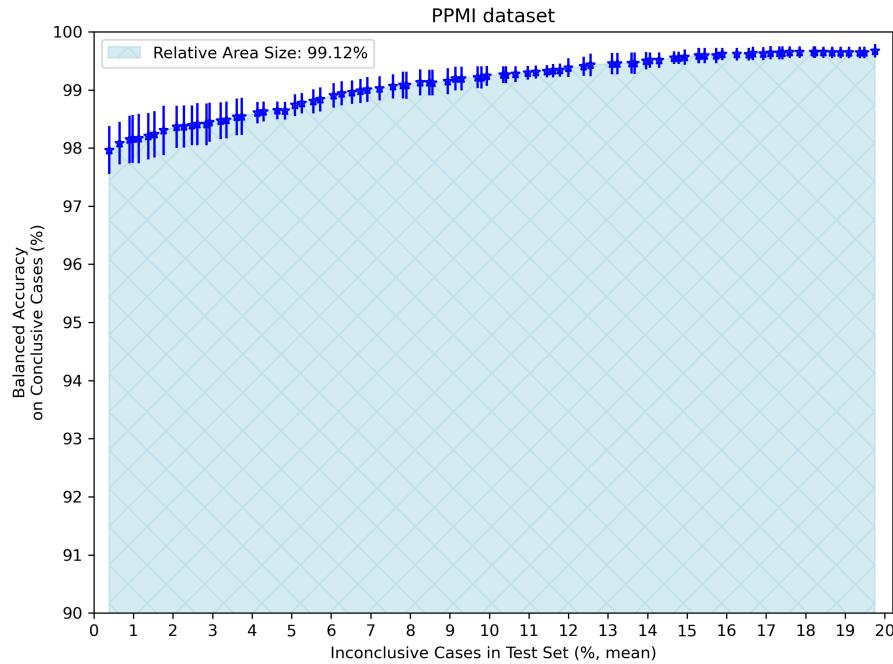


Figure 11: Evaluation of the PCA-RFC method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). The area under the mean of the balanced accuracy is highlighted for better illustration.

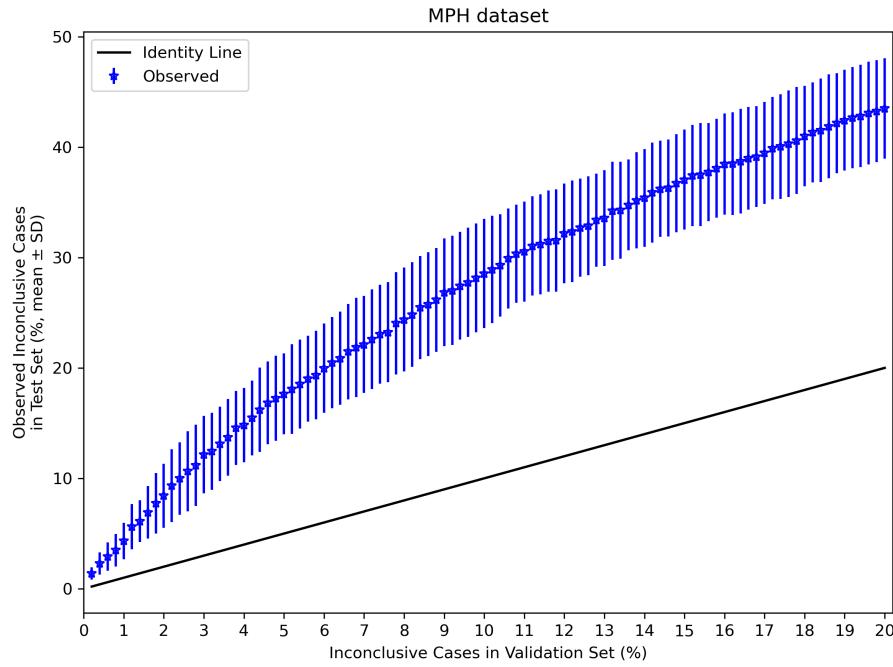


(a) Observed percentage of inconclusive cases in the PPMI dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).

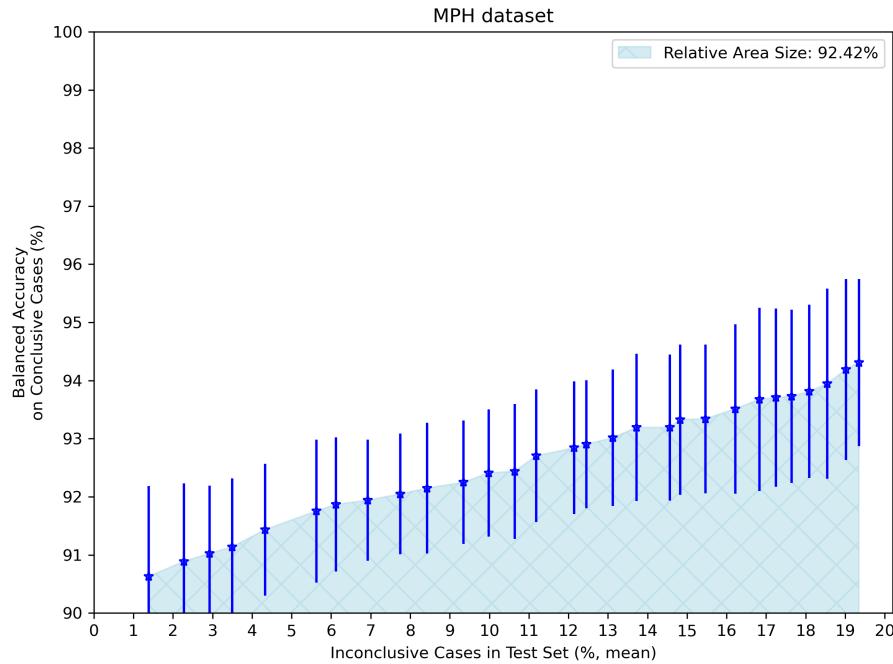


(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the PPMI dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 12: Evaluation of the PCA-RFC method on PPMI dataset.



(a) Observed percentage of inconclusive cases in the MPH dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).



(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the MPH dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 13: Evaluation of the PCA-RFC method on MPH dataset.

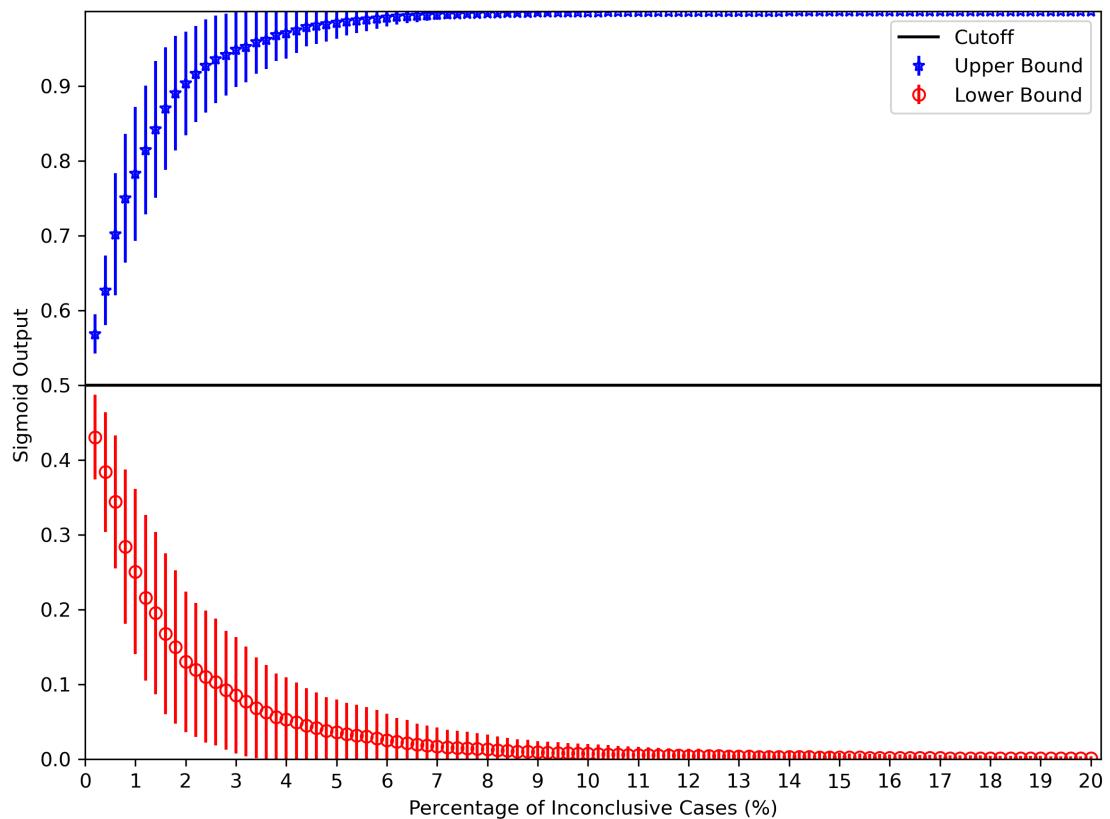


Figure 14: Evaluation of the CNN-MVT method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.

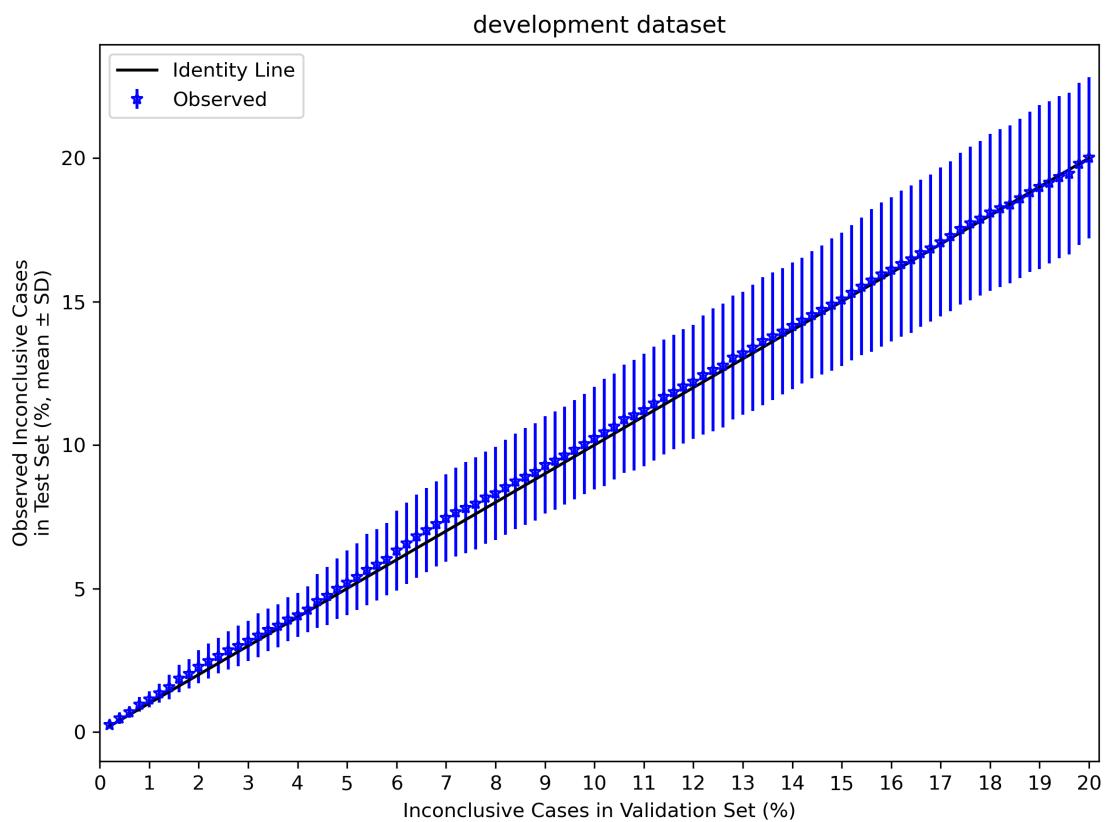


Figure 15: Evaluation of the CNN-MVT method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).

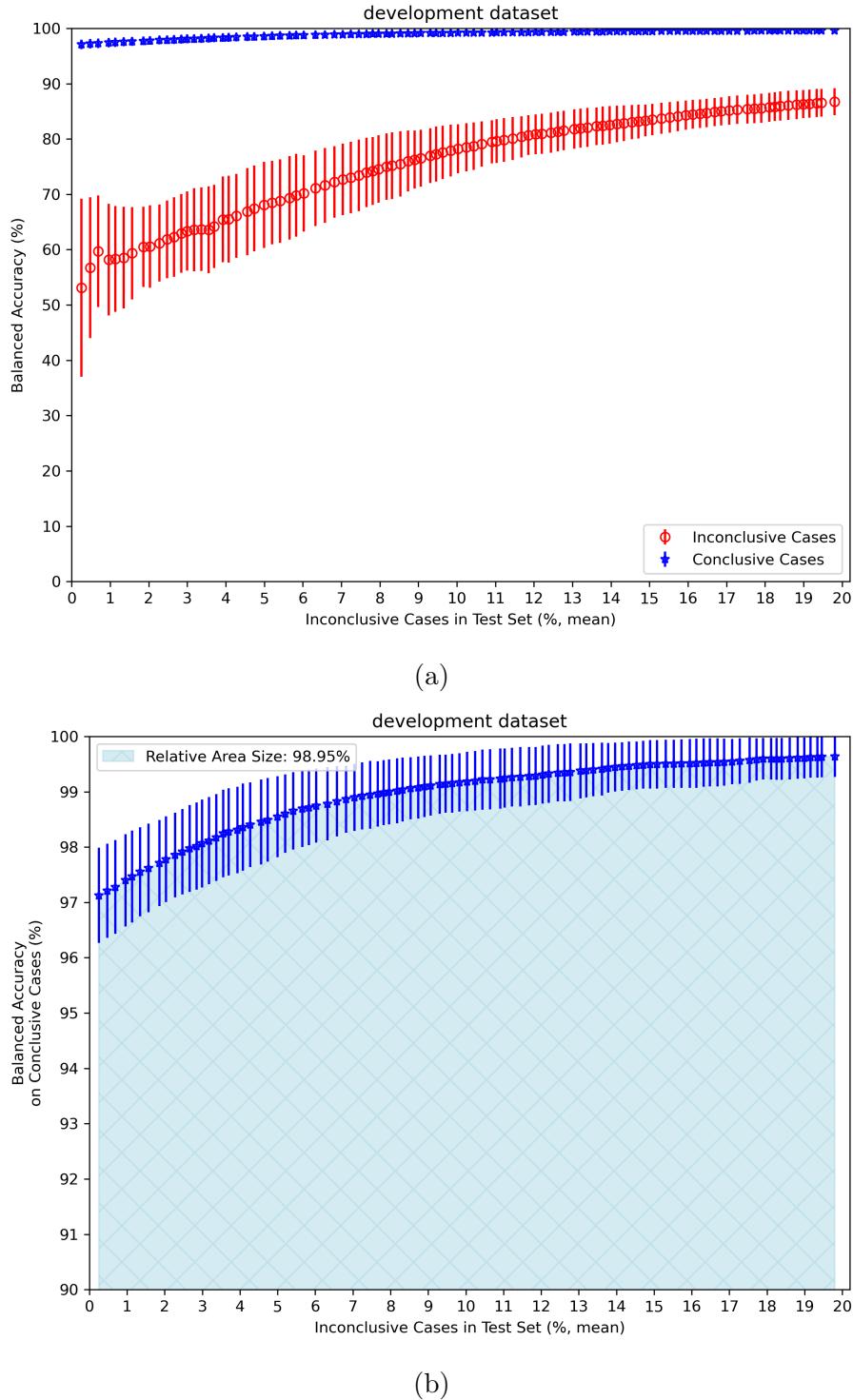
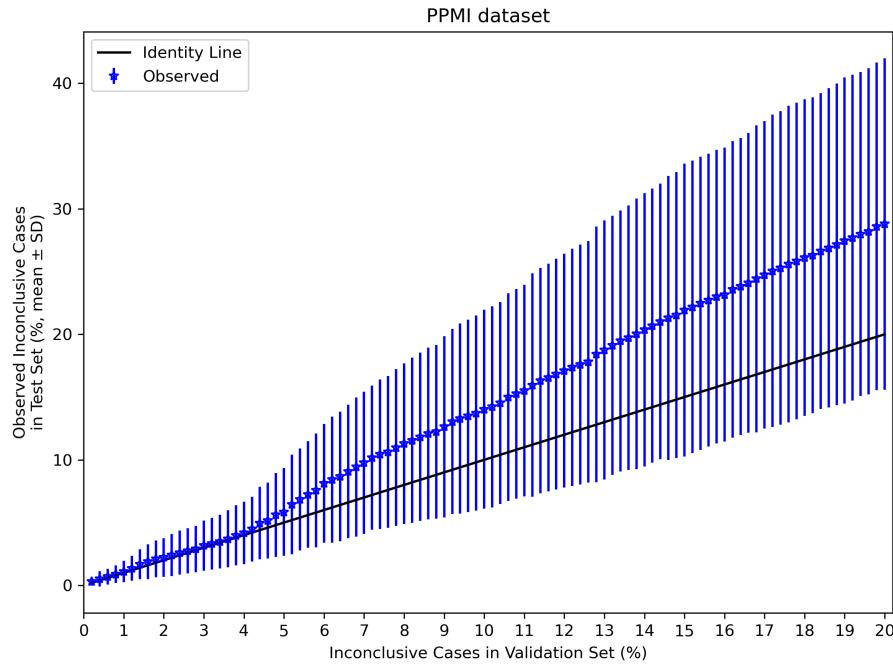
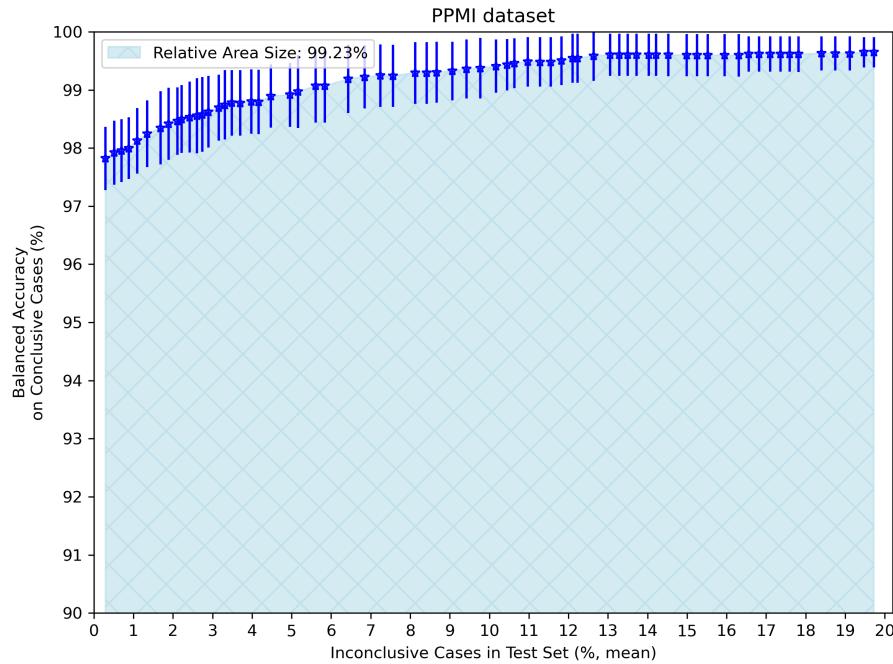


Figure 16: Evaluation of the CNN-MVT method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). The area under the mean of the balanced accuracy is highlighted for better illustration.

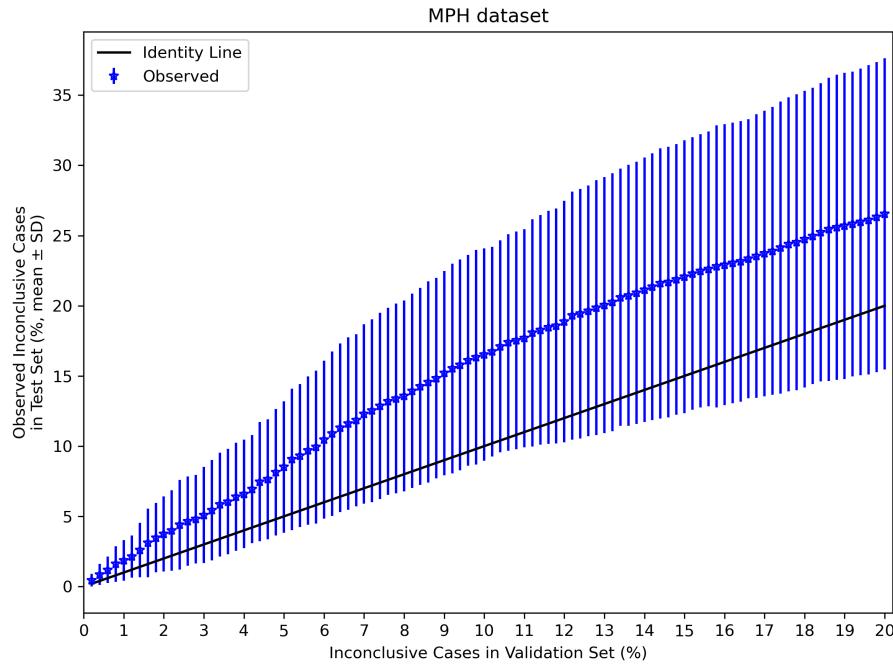


(a) Observed percentage of inconclusive cases in the PPMI dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).

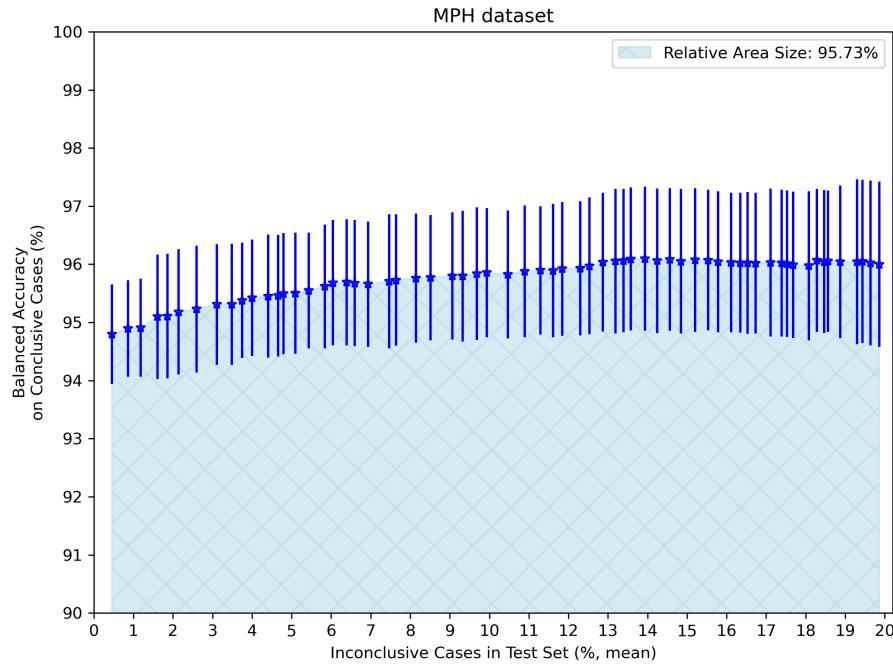


(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the PPMI dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 17: Evaluation of the CNN-MVT method on PPMI dataset.



(a) Observed percentage of inconclusive cases in the MPH dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).



(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the MPH dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 18: Evaluation of the CNN-MVT method on MPH dataset.

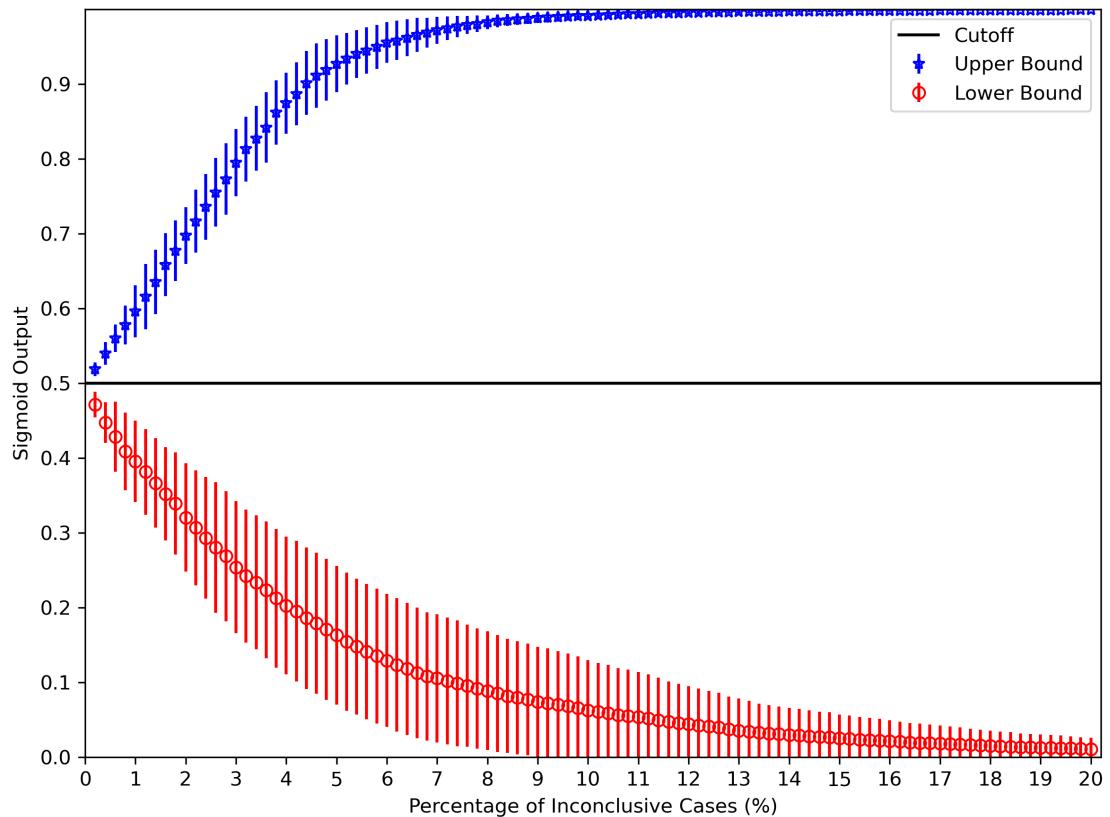


Figure 19: Evaluation of the CNN-RLT method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.

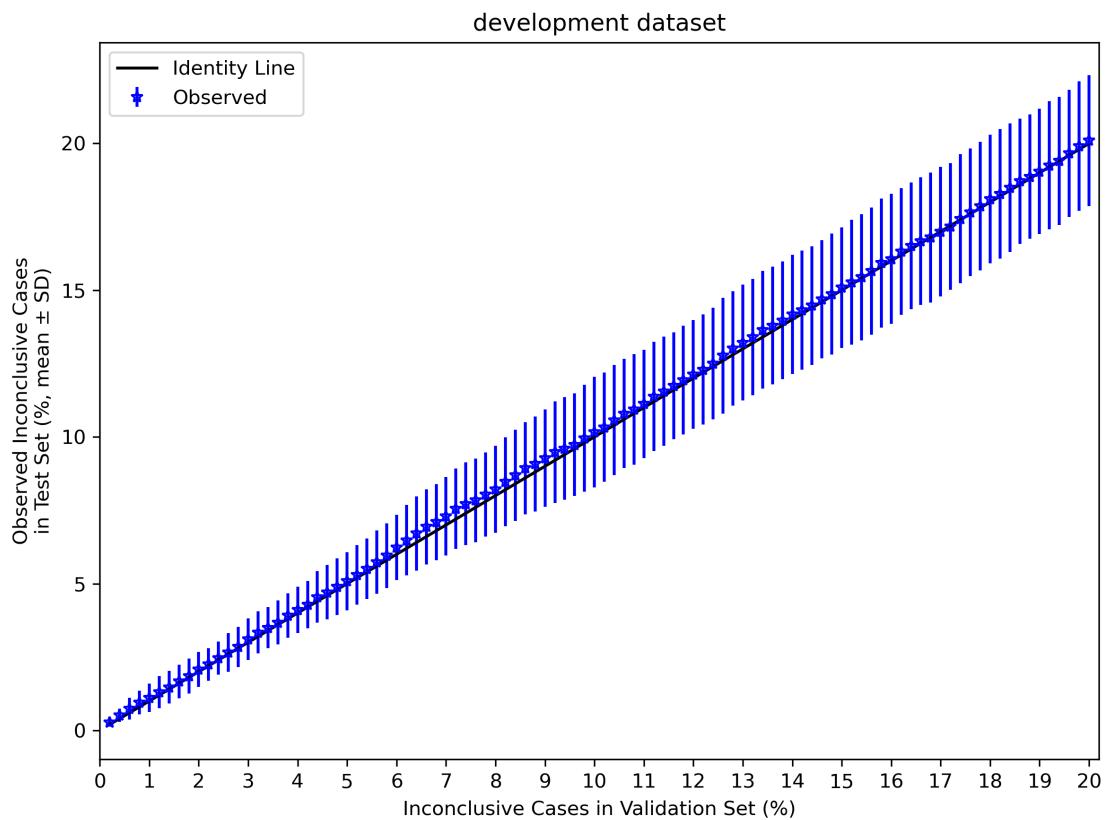


Figure 20: Evaluation of the CNN-RLT method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).

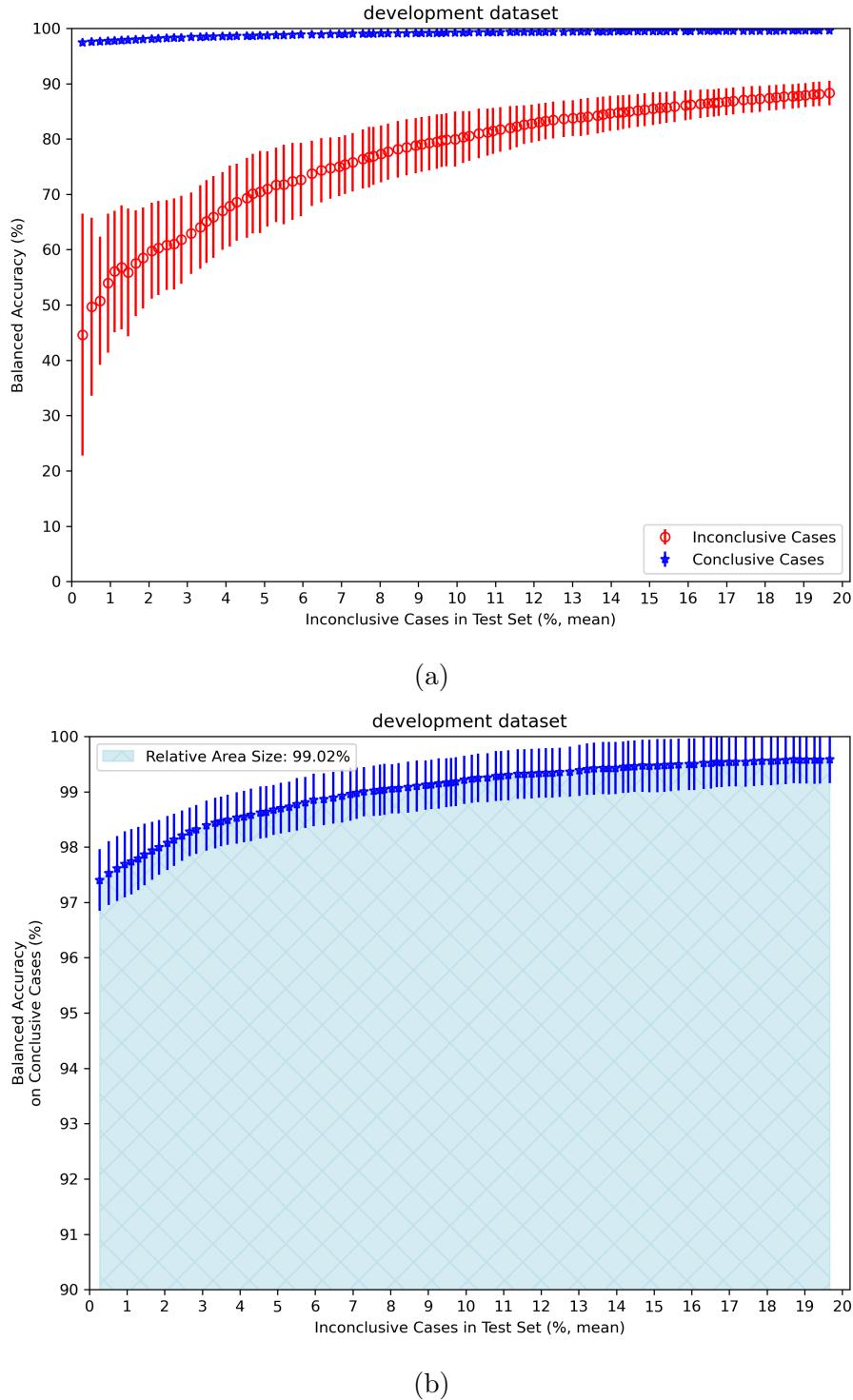
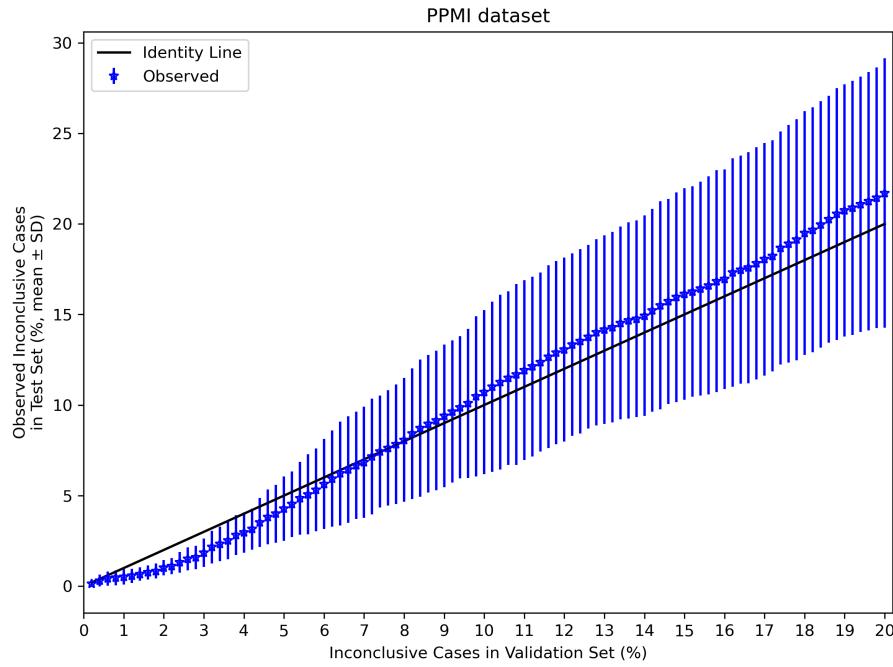
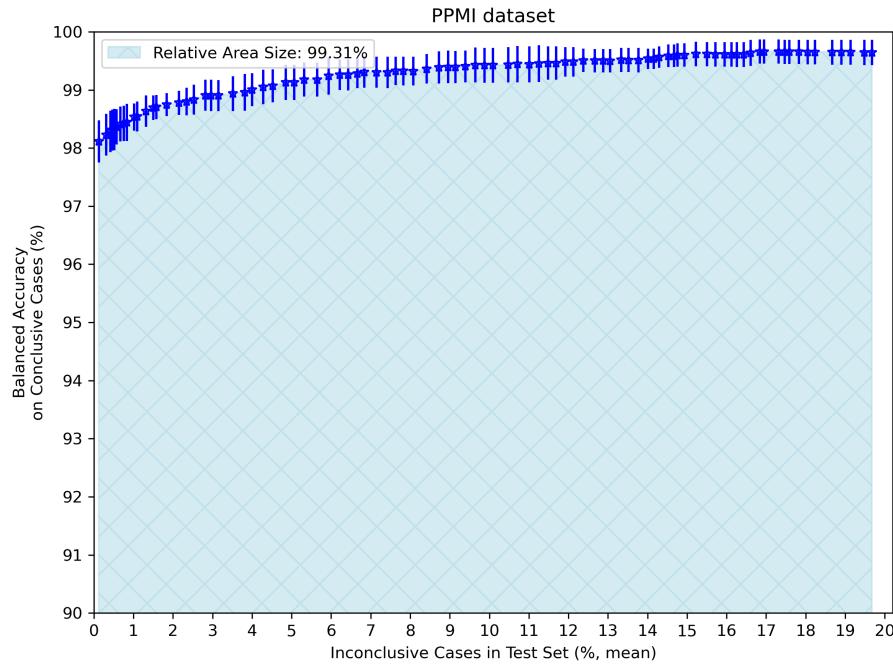


Figure 21: Evaluation of the CNN-RLT method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). The area under the mean of the balanced accuracy is highlighted for better illustration.

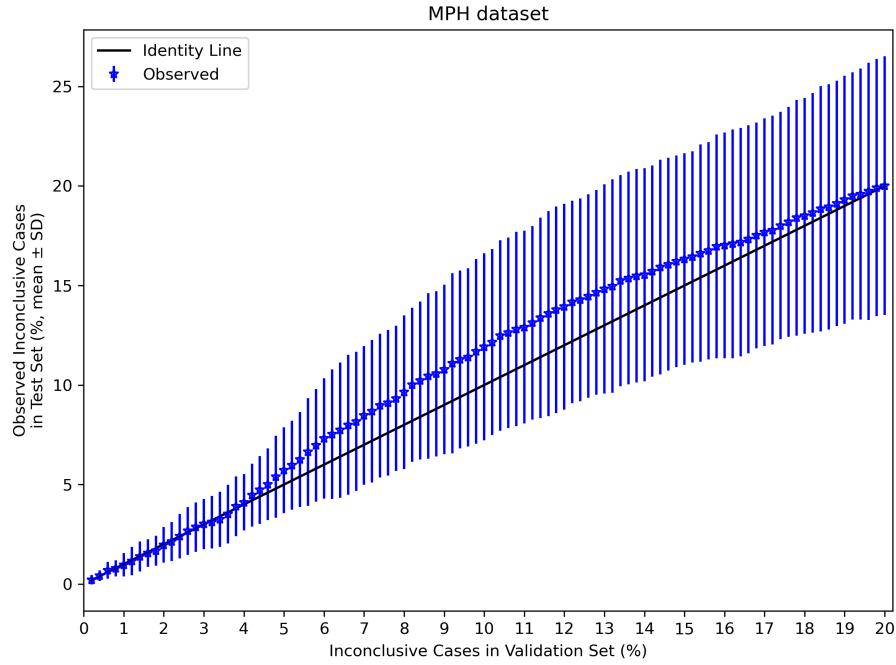


(a) Observed percentage of inconclusive cases in the PPMI dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).

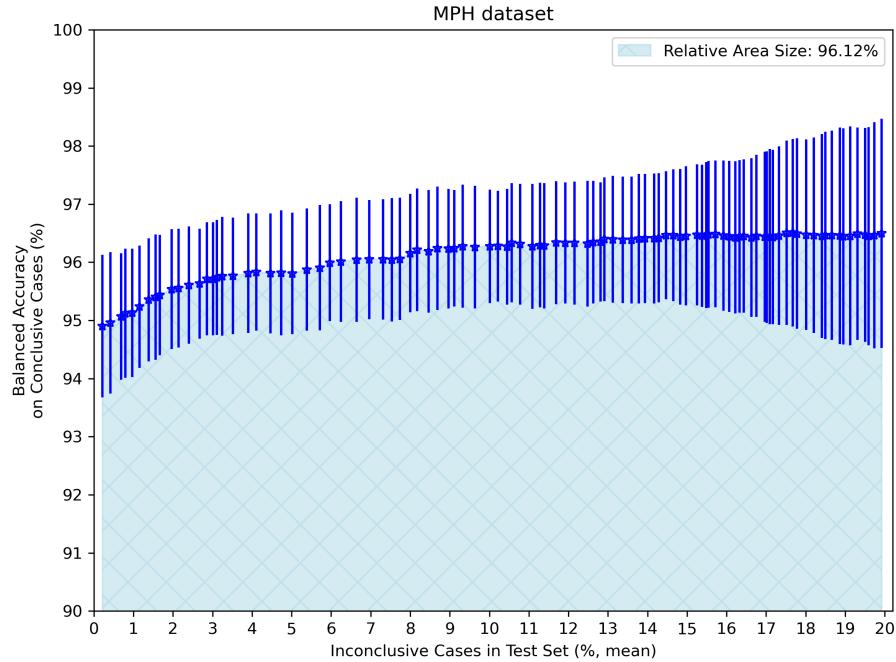


(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the PPMI dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 22: Evaluation of the CNN-RLT method on PPMI dataset.



(a) Observed percentage of inconclusive cases in the MPH dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).



(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the MPH dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 23: Evaluation of the CNN-RLT method on MPH dataset.

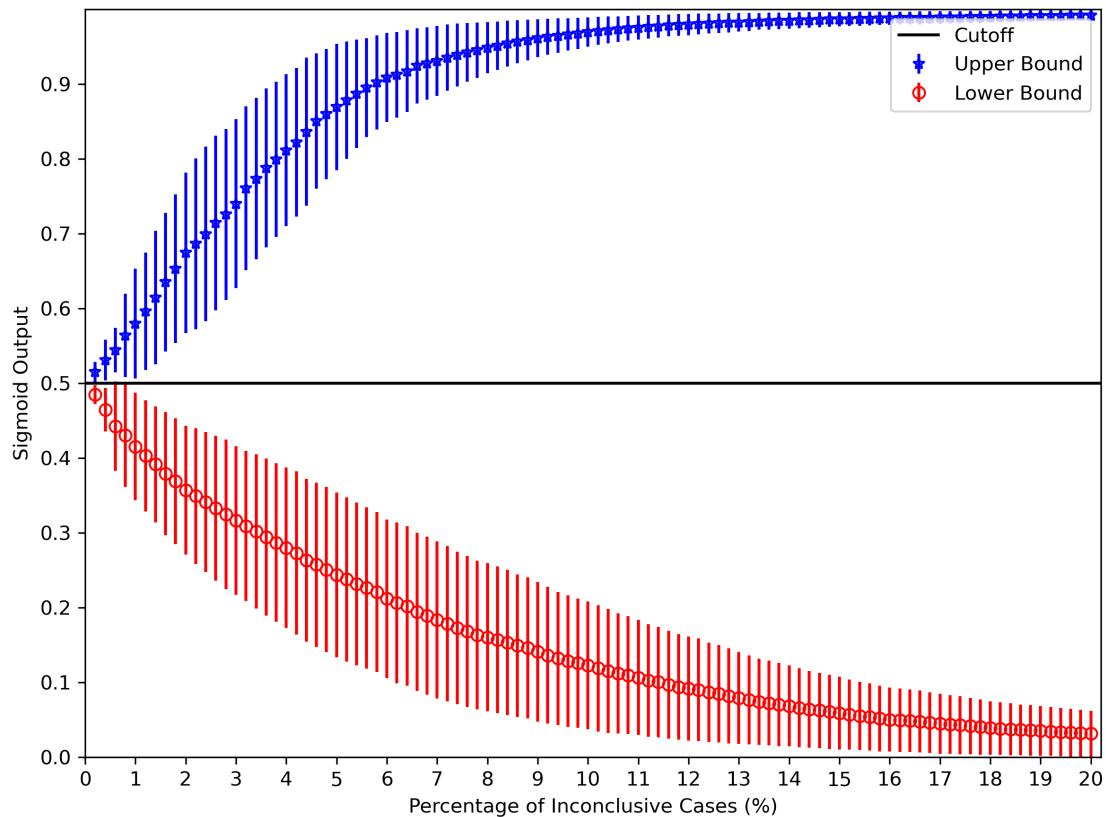


Figure 24: Evaluation of the CNN-Regression method on Test Set of Development dataset. Determined upper and lower bounds of the inconclusive interval as a function of the percentage of inconclusive cases.

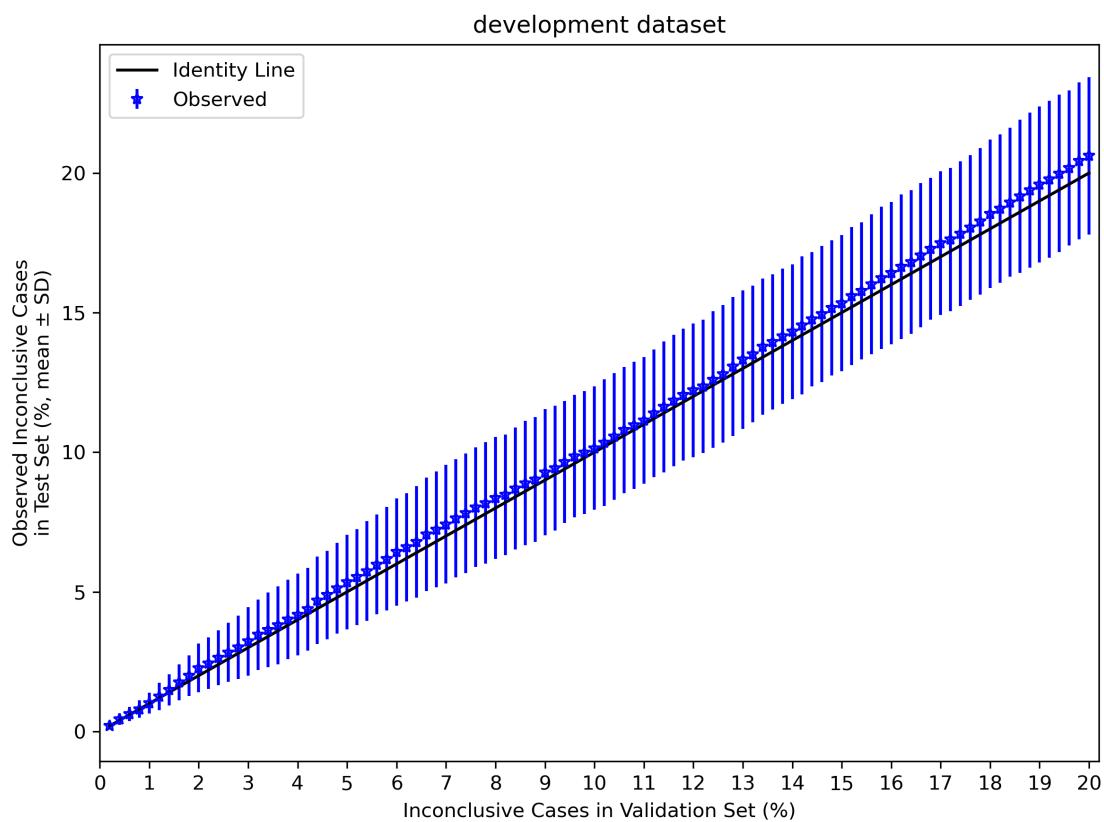


Figure 25: Evaluation of the CNN-Regression method on Test Set of Development dataset. Observed percentage of inconclusive cases in the test set for a given set of percentages of inconclusive cases in the validation set. Each of the percentages of inconclusive cases in the validation set is associated with an inconclusive range (determined in the validation set).

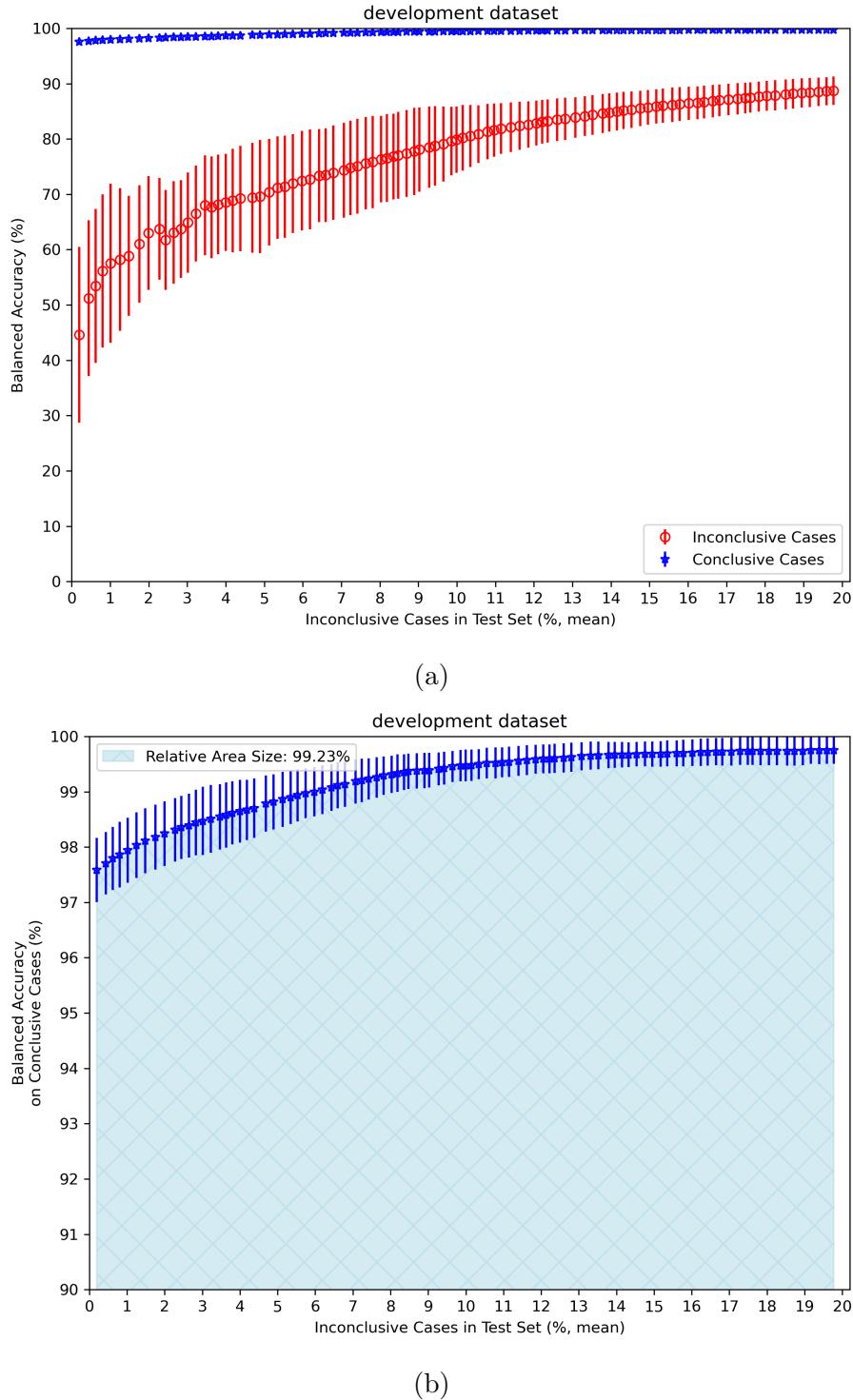
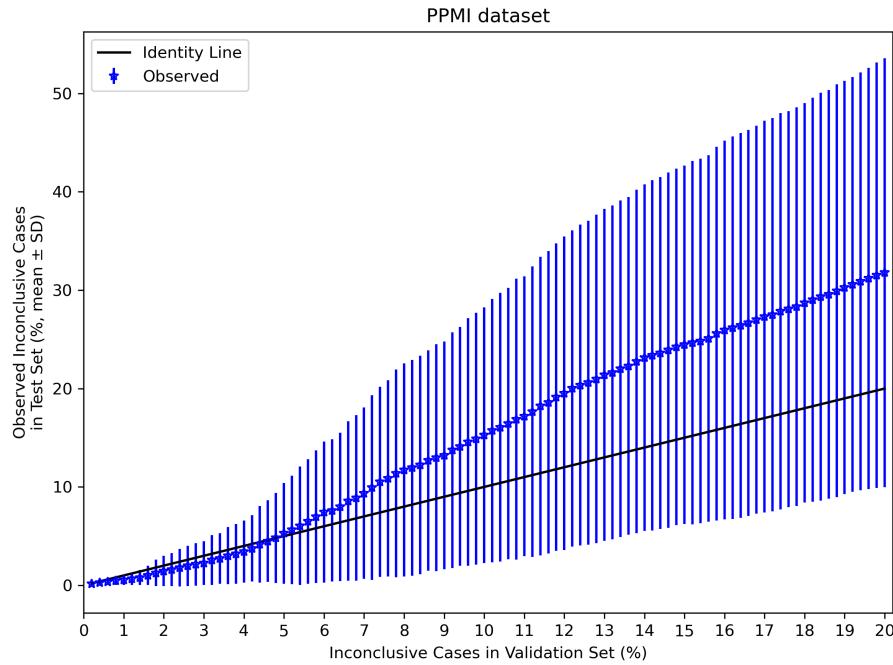
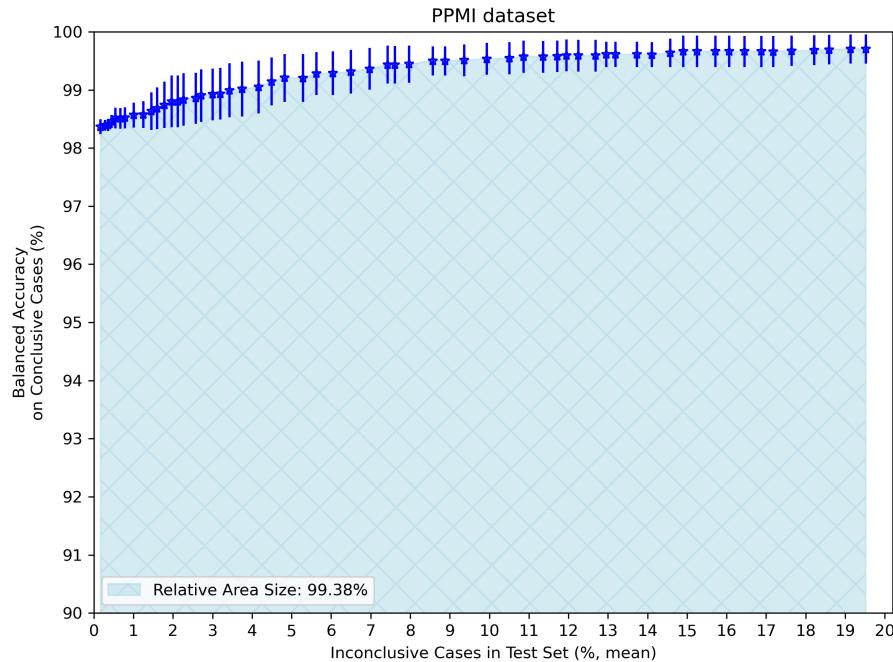


Figure 26: Evaluation of the CNN-Regression method on Test Set of Development dataset. Balanced accuracy for a given mean percentage of observed inconclusive cases in the test set on (a) both conclusive and inconclusive cases and (b) only conclusive cases. Each of the mean percentages of observed inconclusive cases is associated with an inconclusive range (determined in the validation set). The area under the mean of the balanced accuracy is highlighted for better illustration.

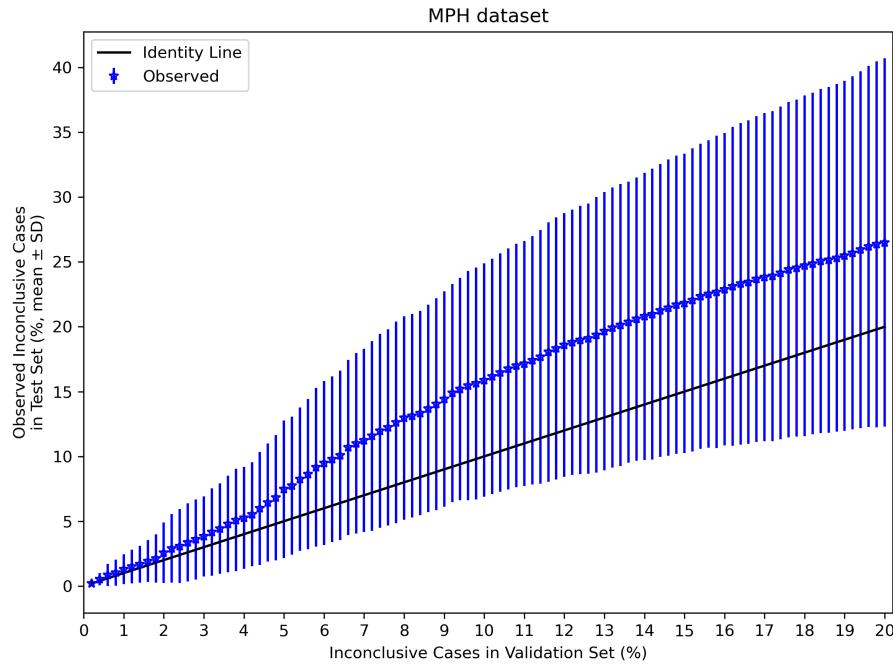


(a) Observed percentage of inconclusive cases in the PPMI dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).

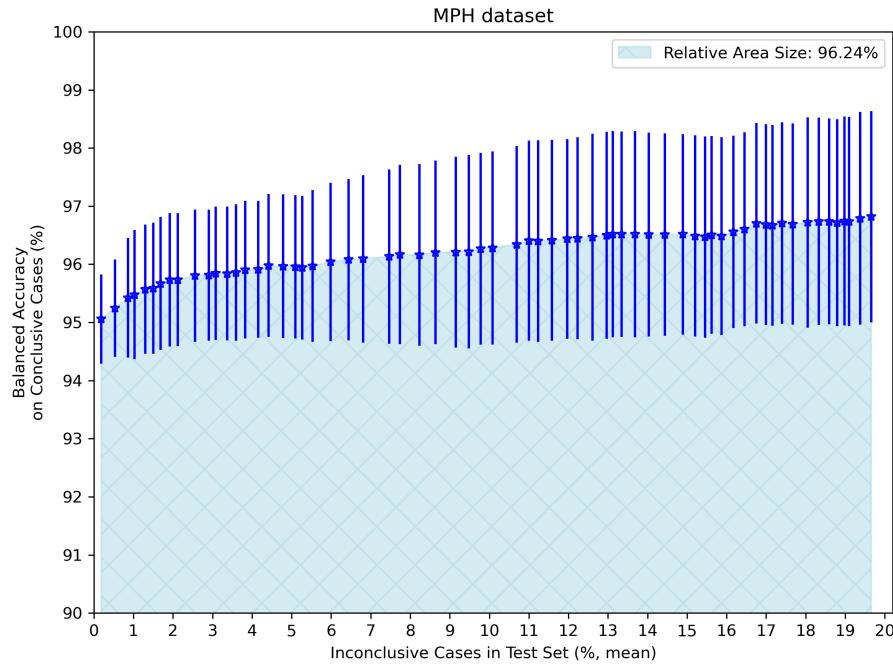


(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the PPMI dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 27: Evaluation of the CNN-Regression method on PPMI dataset.



(a) Observed percentage of inconclusive cases in the MPH dataset for a given set of percentages of inconclusive cases in the validation set (Development dataset).



(b) Balanced accuracy on conclusive cases for a given mean percentage of inconclusive cases observed in the MPH dataset. For better illustration the area under the mean of the balanced accuracy is highlighted.

Figure 28: Evaluation of the CNN-Regression method on MPH dataset.

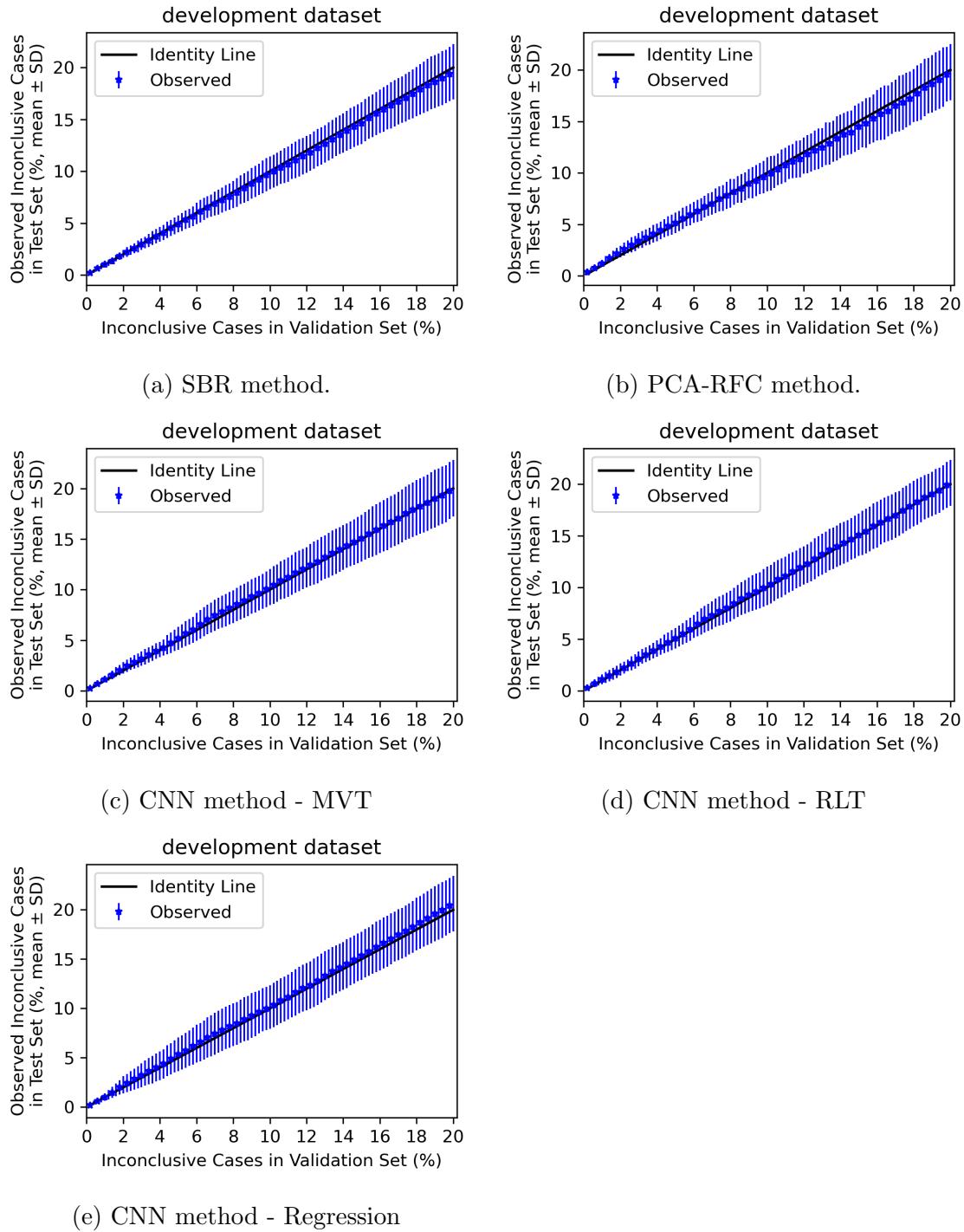


Figure 29: Comparison of different methods on test set of development data. Transferability of inconclusive intervals.

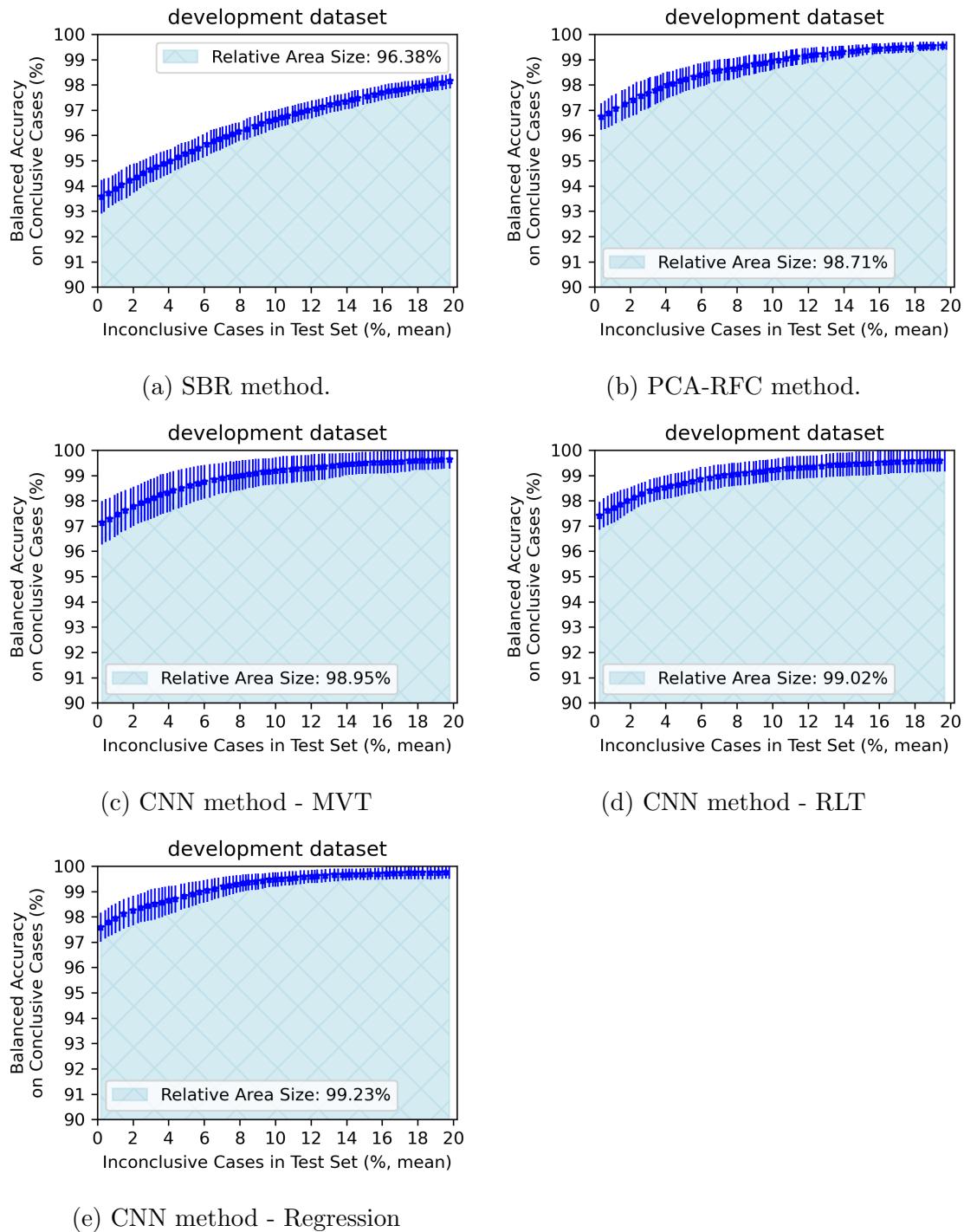


Figure 30: Comparison of different methods on test set of development data. Balanced accuracy over the percentage of observed inconclusive cases.

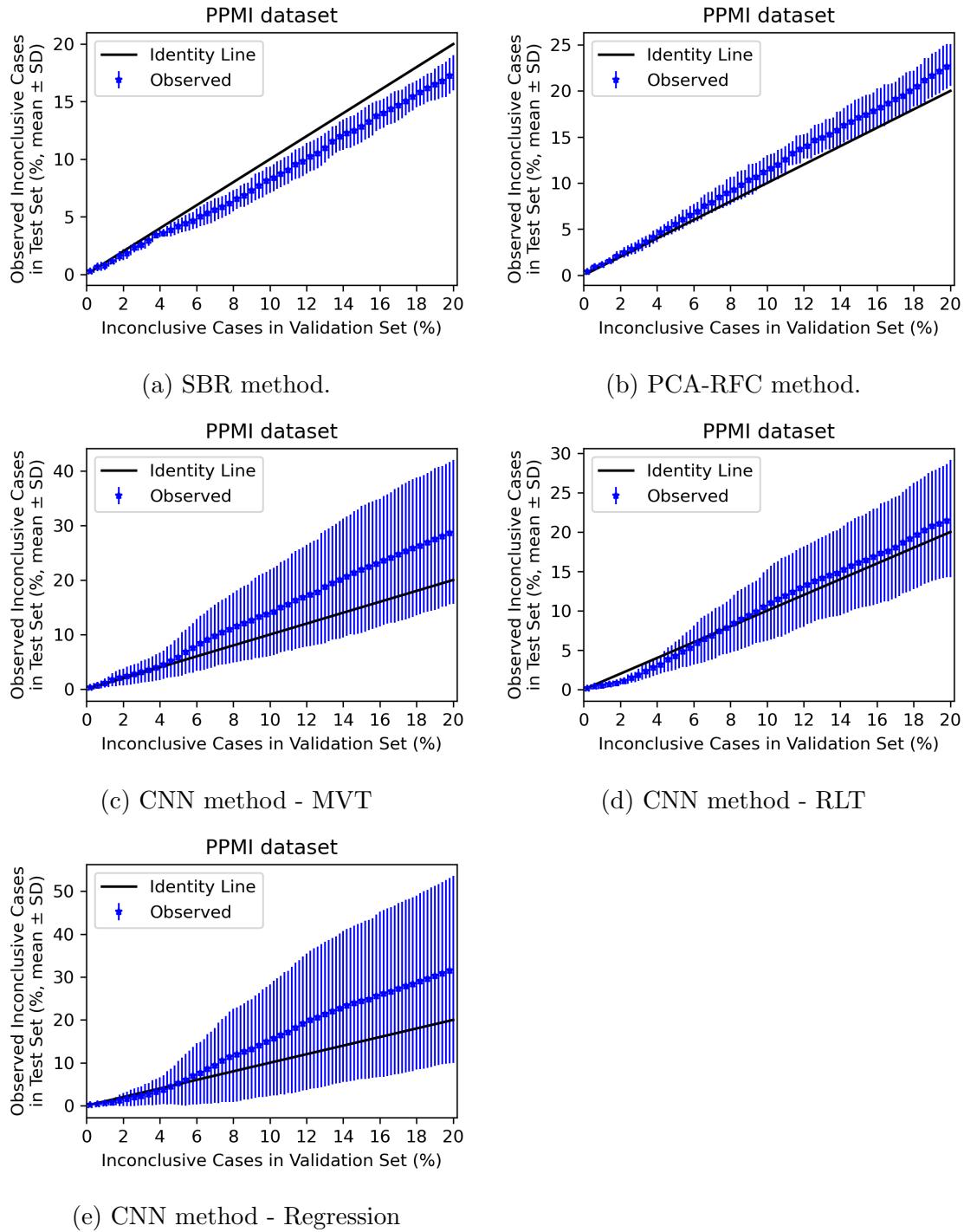


Figure 31: Comparison of different methods on PPMI dataset. Transferability of inconclusive intervals.

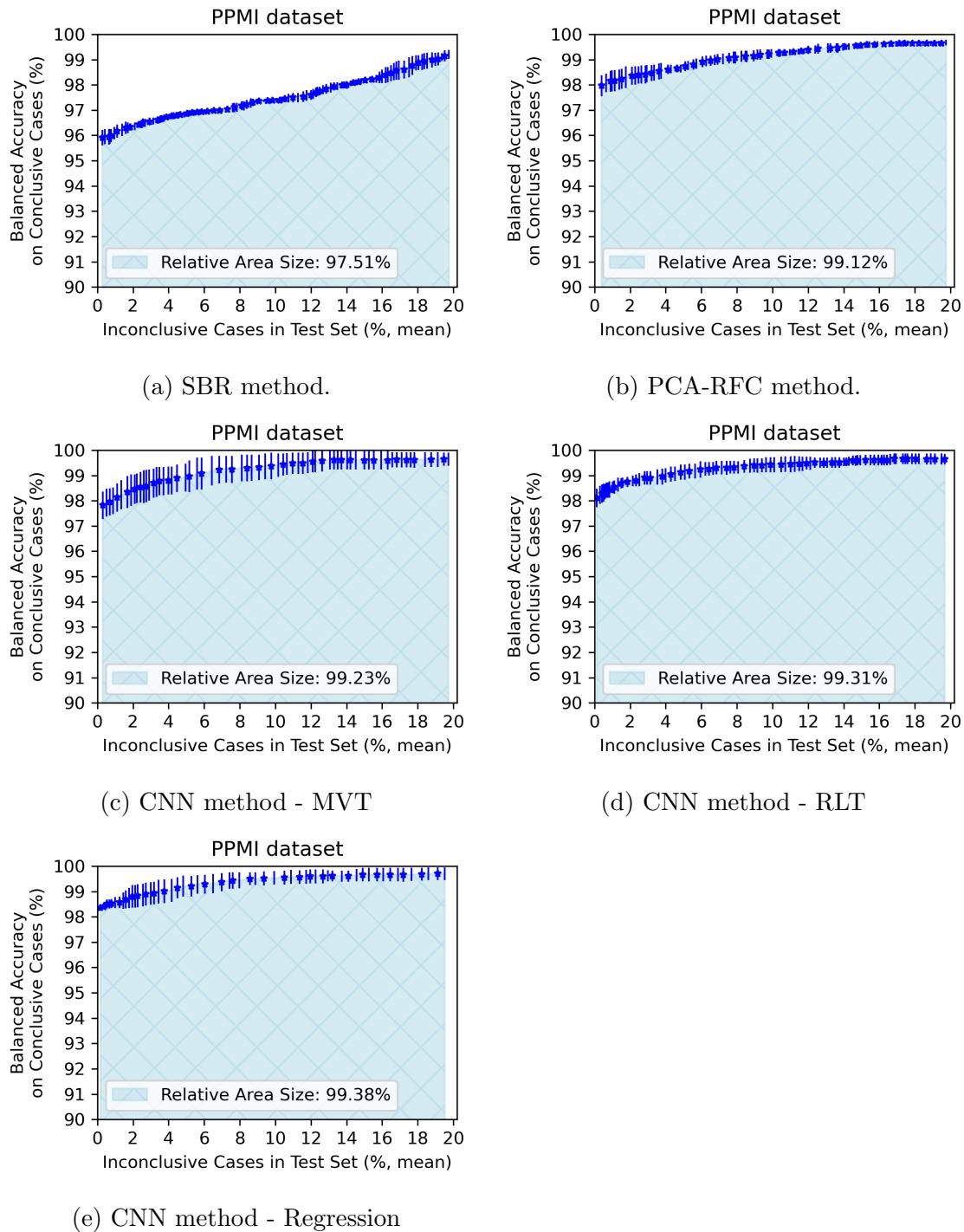


Figure 32: Comparison of different methods on PPMI dataset. Balanced accuracy over the percentage of observed inconclusive cases.

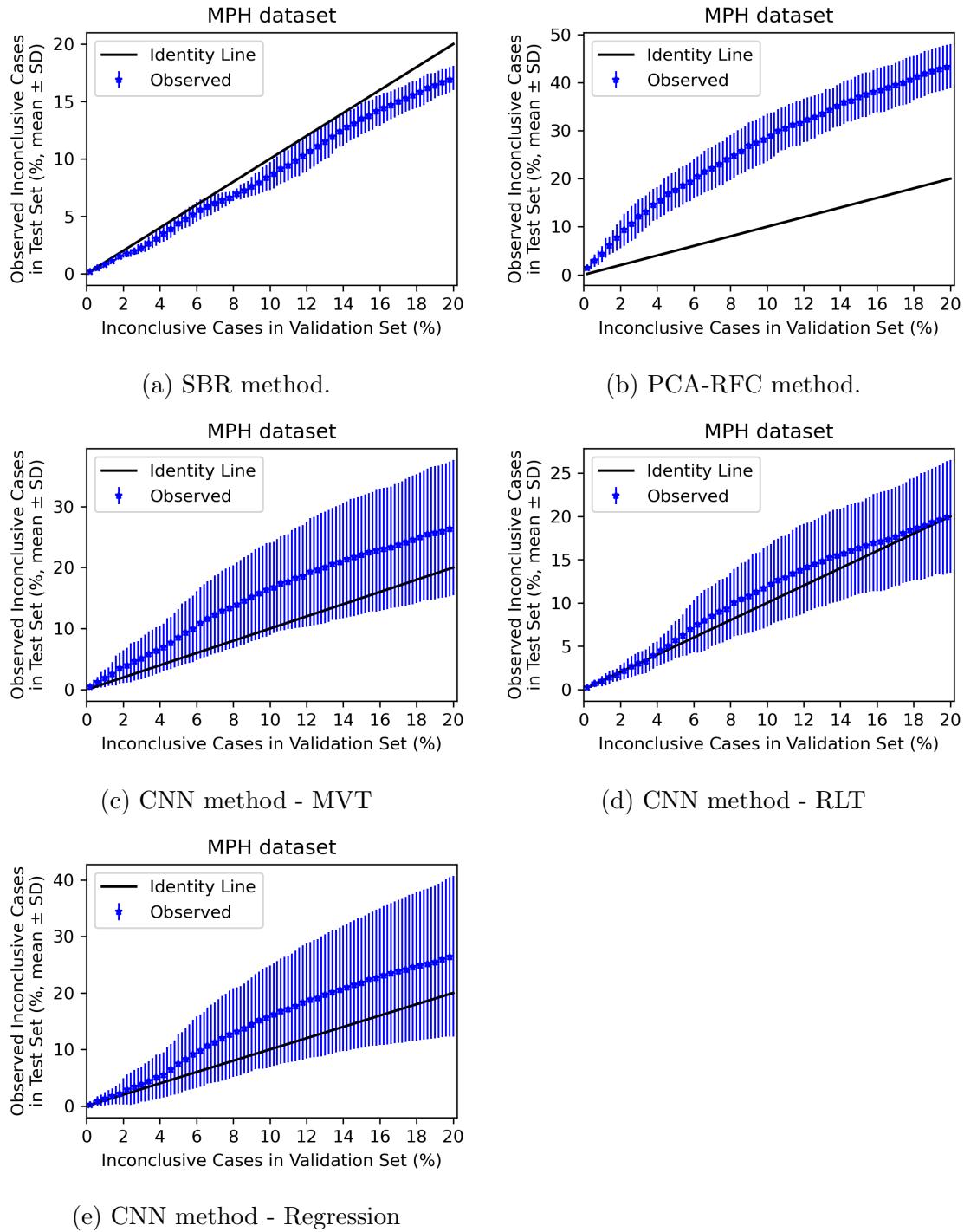


Figure 33: Comparison of different methods on MPH dataset. Transferability of inconclusive intervals.

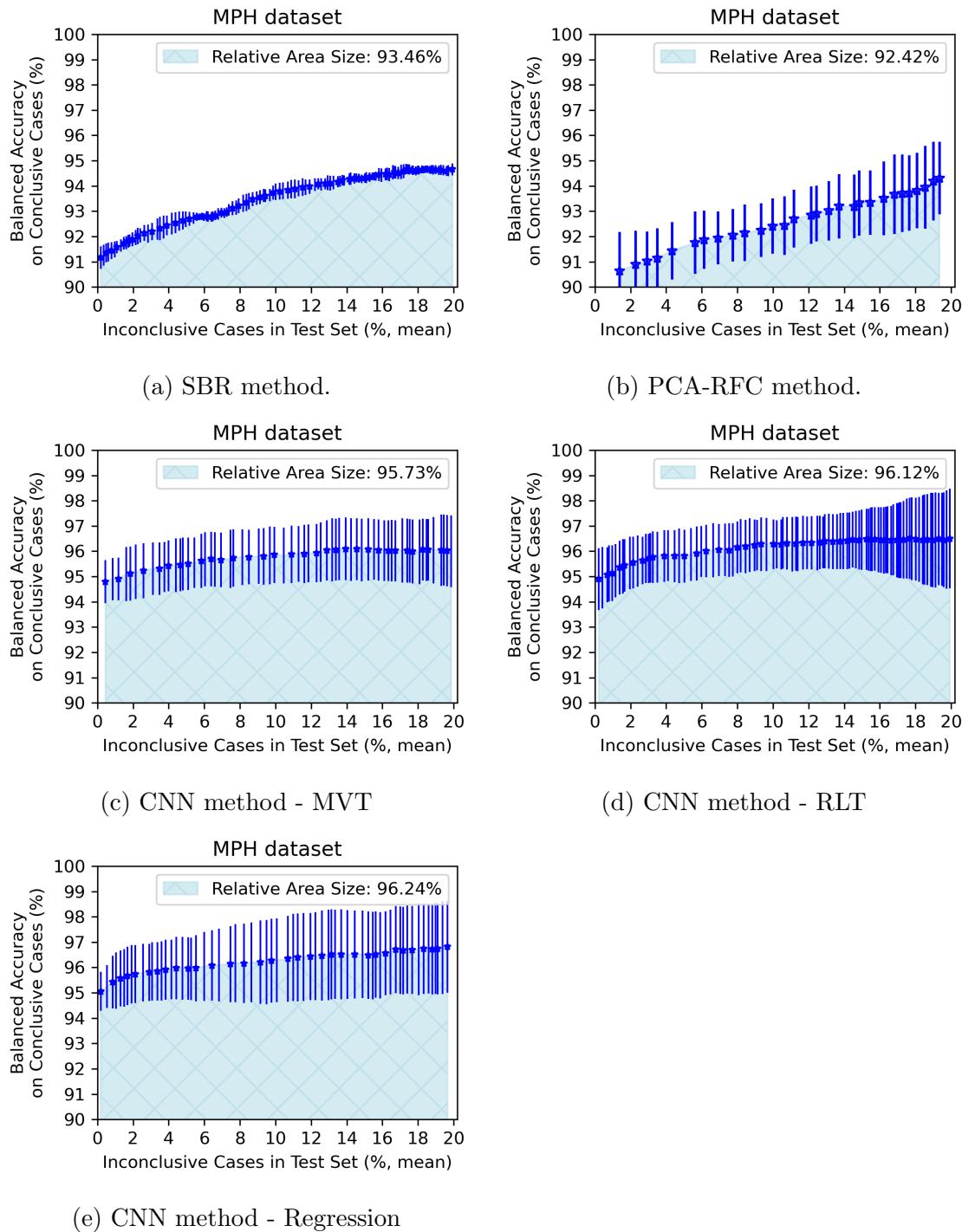


Figure 34: Comparison of different methods on MPH dataset. Balanced accuracy over the percentage of observed inconclusive cases.

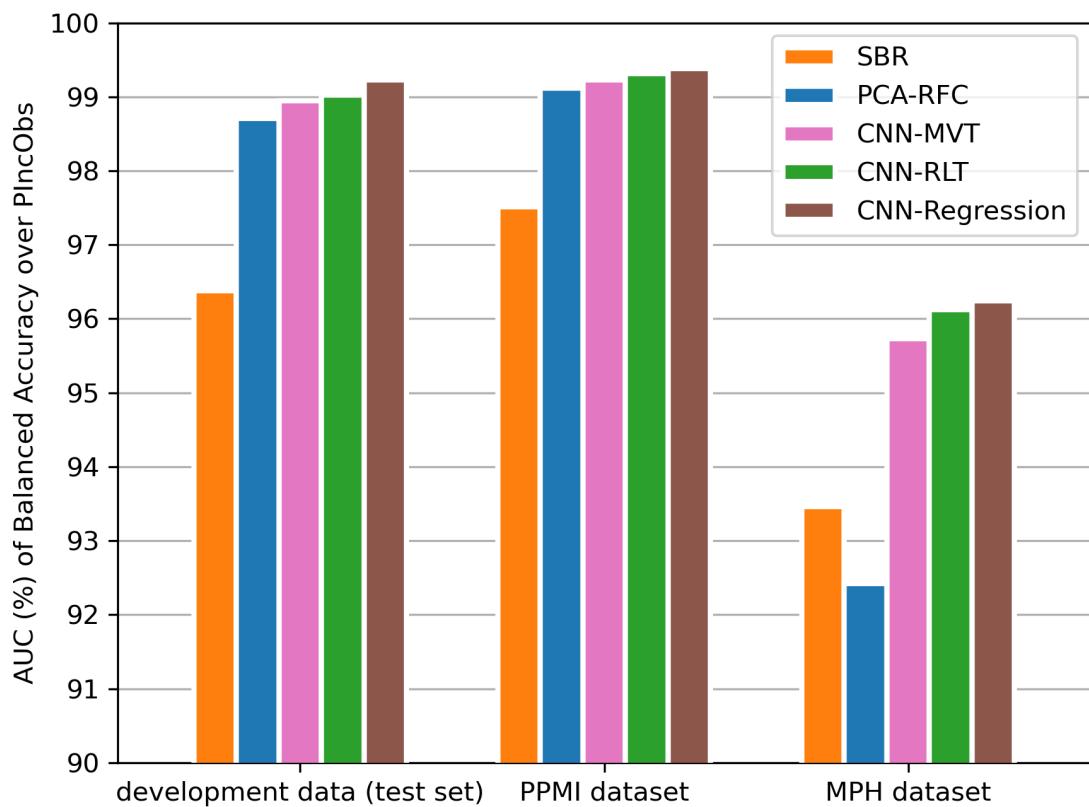


Figure 35: AUC-bACC achieved by baseline and experimental methods on different test data. The AUC-bACC was calculated for the mean balanced accuracy over the percentage of inconclusive cases in the considered test set.

A Appendix

If needed for supplementary material, such as detailed description of data collection, tables, or figures.

Bibliography

- A. Abi-Dargham, M. S. Gadelman, G. A. DeErausquin, Y. Zea-Ponce, S. S. Zoghbi, R. M. Baldwin, M. Laruelle, D. S. Charney, P. B. Hoffer, J. L. Neumeyer, and R. B. Innis. SPECT imaging of dopamine transporters in human brain with iodine-123-fluoroalkyl analogs of beta-CIT. *J. Nucl. Med.*, 37(7):1129–1133, July 1996.
- Nathalie L. Albert, Marcus Unterrainer, Markus Diemling, Guoming Xiong, Peter Bartenstein, Walter Koch, Andrea Varrone, John C. Dickson, Livia Tossici-Bolt, Terez Sera, Susanne Asenbaum, Jan Booij, L Özlem Atay Kapucu, Andreas Kluge, Morten Ziebell, Jacques Darcourt, Flavio Nobili, Marco Pagani, Osama Sabri, Swen Hesse, Thierry Vander Borght, Koen Van Laere, Klaus Tatsch, and Christian la Fougère. Implementation of the european multicentre database of healthy controls for [(123)I]FP-CIT SPECT increases diagnostic accuracy in patients with clinically uncertain parkinsonian syndromes. *Eur. J. Nucl. Med. Mol. Imaging*, 43(7):1315–1322, July 2016.
- Ivayla Apostolova, Daulat S. Taleb, Axel Lipp, Imke Galazky, Dennis Kupitz, Catharina Lange, Marcus R. Makowski, Winfried Brenner, Holger Amthauer, Michail Plotkin, and Ralph Buchert. Utility of follow-up dopamine transporter SPECT with 123I-FP-CIT in the diagnostic workup of patients with clinically uncertain parkinsonian syndrome. *Clin. Nucl. Med.*, 42(8):589–594, August 2017.
- Ivayla Apostolova, Tassilo Schiebler, Catharina Lange, Franziska Lara Mathies, Wencke Lehnert, Susanne Klutmann, and Ralph Buchert. Stereotactical normalization with multiple templates representative of normal and parkinson-typical reduction of striatal uptake improves the discriminative power of automatic semi-quantitative analysis in dopamine transporter SPECT. *EJNMMI Phys.*, 10(1):25, March 2023.
- H. Bernheimer, W. Birkmayer, O. Hornykiewicz, K. Jellinger, and F. Seitelberger. Brain dopamine and the syndromes of parkinson and huntington. clinical, morphological and neurochemical correlations. *J. Neurol. Sci.*, 20(4):415–455, December 1973.
- Ralph Buchert, Georg Berding, Florian Wilke, Brigitte Martin, Daniel von Borczyskowski, Janos Mester, Winfried Brenner, and Malte Clausen. IBZM tool: a fully automated expert system for the evaluation of IBZM SPECT studies. *Eur. J. Nucl. Med. Mol. Imaging*, 33(9):1073–1083, September 2006.
- Ralph Buchert, Carsten Buhmann, Ivayla Apostolova, Philipp T. Meyer, and Jürgen Gallinat. Nuclear imaging in the diagnosis of clinically uncertain parkinsonian syndromes. *Dtsch. Arztebl. Int.*, 116(44):747–754, November 2019a.
- Ralph Buchert, Catharina Lange, Timo S. Spehl, Ivayla Apostolova, Lars Frings, Cathrine Jonsson, Philipp T. Meyer, and Sabine Hellwig. Diagnostic performance

- of the specific uptake size index for semi-quantitative analysis of I-123-FP-CIT SPECT: harmonized multi-center research setting versus typical clinical single-camera setting. *EJNMMI Res.*, 9(1):37, May 2019b.
- Ana M. Catafau, Eduardo Tolosa, and DaTSCAN Clinically Uncertain Parkinsonian Syndromes Study Group. Impact of dopamine transporter SPECT using 123I-Ioflupane on diagnosis and management of patients with clinically uncertain parkinsonian syndromes. *Mov. Disord.*, 19(10):1175–1182, October 2004.
- Chung-Yao Chien, Szu-Wei Hsu, Tsung-Lin Lee, Pi-Shan Sung, and Chou-Ching Lin. Using artificial neural network to discriminate parkinson's disease from other parkinsonisms by focusing on putamen of dopamine transporter SPECT images. *Biomedicines*, 9(1):12, December 2020.
- Jacques Darcourt, Jan Booij, Klaus Tatsch, Andrea Varrone, Thierry Vander Borgh, Ozlem L. Kapucu, Kjell Någren, Flavio Nobili, Zuzana Walker, and Koen Van Laere. EANM procedure guidelines for brain neurotransmission SPECT using (123)i-labelled dopamine transporter ligands, version 2. *Eur. J. Nucl. Med. Mol. Imaging*, 37(2):443–450, February 2010.
- Lonneke M. L. de Lau and Monique M. B. Breteler. Epidemiology of parkinson's disease. *Lancet Neurol.*, 5(6):525–535, June 2006.
- John C. Dickson, Livia Tossici-Bolt, Terez Sera, Kjell Erlandsson, Andrea Varrone, Klaus Tatsch, and Brian F Hutton. The impact of reconstruction method on the quantification of DaTSCAN images. *Eur. J. Nucl. Med. Mol. Imaging*, 37(1):23–35, January 2010.
- John Caddell Dickson, Livia Tossici-Bolt, Terez Sera, Robin de Nijs, Jan Booij, Maria Claudia Bagnara, Anita Seese, Pierre Malick Koulibaly, Umit Ozgur Akdemir, Cathrine Jonsson, Michel Koole, Maria Raith, Markus Nowak Lonsdale, Jean George, Felicia Zito, and Klaus Tatsch. Proposal for the standardisation of multi-centre trials in nuclear medicine imaging: prerequisites for a european 123I-FP-CIT SPECT database. *Eur. J. Nucl. Med. Mol. Imaging*, 39(1):188–197, January 2012.
- M. Diemling. HERMES camera correction for the ENCDAT database using DaTscan. Technical report, Hermes Medical Solutions, 2021.
- David S. W. Djang, Marcel J. R. Janssen, Nicolaas Bohnen, Jan Booij, Theodore A. Henderson, Karl Herholz, Satoshi Minoshima, Christopher C. Rowe, Osama Sabri, John Seibyl, Bart N. M. Van Berckel, and Michele Wanner. SNM practice guideline for dopamine transporter imaging with 123i-ioflupane SPECT 1.0. *J. Nucl. Med.*, 53(1):154–163, January 2012.
- Patrik Fazio, Per Svenningsson, Zsolt Cselényi, Christer Halldin, Lars Farde, and Andrea Varrone. Nigrostriatal dopamine transporter availability in early parkinson's disease. *Mov. Disord.*, 33(4):592–599, April 2018.

- Elisabeth Funke, Andreas Kupsch, Ralph Buchert, Winfried Brenner, and Michail Plotkin. Impact of subcortical white matter lesions on dopamine transporter SPECT. *J. Neural Transm. (Vienna)*, 120(7):1053–1060, July 2013.
- W. R. Gibb and A. J. Lees. The relevance of the lewy body to the pathogenesis of idiopathic parkinson’s disease. *J. Neurol. Neurosurg. Psychiatry*, 51(6):745–752, June 1988.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1321–1330. PMLR, 06–11 Aug 2017.
- Jigna Hathaliya, Raj Parekh, Nisarg Patel, Rajesh Gupta, Sudeep Tanwar, Fayez Alqahtani, Magdy Elghatwary, Ovidiu Ivanov, Maria Simona Raboaca, and Bogdan-Constantin Neagu. Convolutional neural network-based parkinson disease classification using spect imaging data. *Mathematics*, 10(15), 2022. ISSN 2227-7390. doi: 10.3390/math10152566.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- Emma A. Honkanen, Laura Saari, Katri Orte, Maria Gardberg, Tommi Noponen, Juho Joutsa, and Valtteri Kaasinen. No link between striatal dopaminergic axons and dopamine transporter imaging in parkinson’s disease. *Mov. Disord.*, 34(10):1562–1566, October 2019.
- H.M. Hudson and R.S. Larkin. Accelerated image reconstruction using ordered subsets of projection data. *IEEE Transactions on Medical Imaging*, 13(4):601–609, 1994. doi: 10.1109/42.363108.
- A. J. Hughes, S. E. Daniel, L. Kilford, and A. J. Lees. Accuracy of clinical diagnosis of idiopathic parkinson’s disease: a clinico-pathological study of 100 cases. *J. Neurol. Neurosurg. Psychiatry*, 55(3):181–184, March 1992.
- Andrew J. Hughes, Susan E. Daniel, Yoav Ben-Shlomo, and Andrew J. Lees. The accuracy of diagnosis of parkinsonian syndromes in a specialist movement disorder service. *Brain*, 125(4):861–870, April 2002.
- Alex Iranzo, Joan Santamaría, Francesc Valldeoriola, Monica Serradell, Manel Salamero, Carles Gaig, Aida Niñerola-Baizán, Raquel Sánchez-Valle, Albert Lladó, Roberto De Marzi, Ambra Stefani, Klaus Seppi, Javier Pavia, Birgit Högl, Werner Poewe, Eduard Tolosa, and Francisco Lomeña. Dopamine transporter imaging deficit predicts early transition to synucleinopathy in idiopathic rapid eye movement sleep behavior disorder. *Annals of neurology*, 82(3):419—428, September 2017. ISSN 0364-5134. doi: 10.1002/ana.25026.

- Tuija S. Kangasmaa, Chris Constable, Eero Hippeläinen, and Antti O. Sohlberg. Multicenter evaluation of single-photon emission computed tomography quantification with third-party reconstruction software. *Nucl. Med. Commun.*, 37(9):983–987, September 2016.
- Davood Karimi and Ali Gholipour. Improving calibration and out-of-distribution detection in medical image segmentation with convolutional neural networks, 2020.
- Kwang-Soo Kim. Toward neuroprotective treatments of parkinson’s disease. *Proc. Natl. Acad. Sci. U. S. A.*, 114(15):3795–3797, April 2017.
- J. T. Kuikka, K. A. Bergstrom, A. Ahonen, J. Hiltunen, J. Haukka, E. Lansimies, S. Y. Wang, and J. L. Neumeyer. Comparison of i-123 labeled 2-beta-carbomethoxy-3-beta-(4-iodophenyl)tropane and 2-beta-carbomethoxy-3-beta-(4-iodophenyl)-n-(3-fluoropropyl)nortropane for imaging of the dopamine transporter in the living human brain. *European Journal of Nuclear Medicine and Molecular Imaging*, 22:356–360, 1995. ISSN 1619-7070.
- D. Kupitz, I. Apostolova, C. Lange, G. Ulrich, H. Amthauer, W. Brenner, and R. Buchert. Global scaling for semi-quantitative analysis in FP-CIT SPECT. *Nuklearmedizin*, 53(6):234–241, September 2014.
- C. S. Lee, A. Samii, V. Sossi, T. J. Ruth, M. Schulzer, J. E. Holden, J. Wudel, P. K. Pal, R. de la Fuente-Fernandez, D. B. Calne, and A. J. Stoessl. In vivo positron emission tomographic evidence for compensatory changes in presynaptic dopaminergic nerve terminals in parkinson’s disease. *Ann. Neurol.*, 47(4):493–503, April 2000.
- Milán Magdics, László Szirmay-Kalos, Ákos Szlavecz, Gábor Hesz, Balázs Benyó, Áron Cserkaszky, Judit Lantos, D. Légrády, S. Czifrus, András Wirth, et al. TeraTomo project: a fully 3D GPU based reconstruction code for exploiting the imaging capability of the NanoPET/CT system. *Mol Imaging Biol*, 12, 2010.
- Pavan Rajkumar Magesh, Richard Delwin Myloth, and Rijo Jackson Tom. An explainable machine learning model for early detection of parkinson’s disease using LIME on DaTSCAN imagery. *Comput. Biol. Med.*, 126(104041):104041, November 2020.
- Elina Mäkinen, Juho Joutsa, Jarkko Johansson, Maija Mäki, Marko Seppänen, and Valtteri Kaasinen. Visual versus automated analysis of [I-123]FP-CIT SPECT scans in parkinsonism. *J. Neural Transm. (Vienna)*, 123(11):1309–1318, November 2016.
- Jörg Marienhagen, Karin Menhart, Jirka Grosse, and Dirk Hellwig. Nuklearmedizin in deutschland. *Nuklearmedizin*, 56(02):55–68, 2017.
- Franziska Mathies, Ivayla Apostolova, Lena Dierck, Janin Jacobi, Katja Kuen, Markus Sauer, Michael Schenk, Susanne Klutmann, Attila Forgács, and Ralph

- Buchert. Multiple-pinhole collimators improve intra- and between-rater agreement and the certainty of the visual interpretation in dopamine transporter SPECT. *EJNMMI Res.*, 12(1):51, August 2022.
- Silvia Morbelli, Giuseppe Esposito, Javier Arbizu, Henryk Barthel, Ronald Boellaard, Nico I. Bohnen, David J. Brooks, Jacques Darcourt, John C. Dickson, David Douglas, Alexander Drzezga, Jacob Dubroff, Ozgul Ekmekcioglu, Valentina Garibotto, Peter Herscovitch, Phillip Kuo, Adriaan Lammertsma, Sabina Pappata, Iván Peñuelas, John Seibyl, Franck Semah, Livia Tossici-Bolt, Elsmarieke Van de Giessen, Koen Van Laere, Andrea Varrone, Michele Wanner, George Zubal, and Ian Law. EANM practice guideline/SNMMI procedure standard for dopaminergic imaging in parkinsonian syndromes 1.0. *Eur. J. Nucl. Med. Mol. Imaging*, 47(8):1885–1912, July 2020.
- Mahmood Nazari, Andreas Kluge, Ivayla Apostolova, Susanne Klutmann, Sharok Kimiae, Michael Schroeder, and Ralph Buchert. Explainable AI to improve acceptance of convolutional neural networks for automatic classification of dopamine transporter SPECT in the diagnosis of clinically uncertain parkinsonian syndromes. *Eur. J. Nucl. Med. Mol. Imaging*, 49(4):1176–1186, March 2022.
- J. L. Neumeyer, S. Wang, Y. Gao, R. A. Milius, N S. Kula, A. Campbell, R. J. Baldessarini, Y. Zea-Ponce, R. M. Baldwin, and R. B. Innis. N-omega-fluoroalkyl analogs of (1r)-2 beta-carbomethoxy-3 beta-(4-iodophenyl)-tropane (beta-CIT): radiotracers for positron emission tomography and single photon emission computed tomography imaging of dopamine transporters. *J. Med. Chem.*, 37(11):1558–1561, May 1994.
- H. B. Niznik, E. F. Fogel, F. F. Fassos, and P. Seeman. The dopamine transporter is absent in parkinsonian putamen and reduced in the caudate nucleus. *J. Neurochem.*, 56(1):192–198, January 1991.
- Parkinson Progression Marker Initiative. The parkinson progression marker initiative (PPMI). *Prog. Neurobiol.*, 95(4):629–635, December 2011.
- Paola Piccini and Alan Whone. Functional brain imaging in the differential diagnosis of parkinson's disease. *Lancet Neurol.*, 3(5):284–290, May 2004.
- M. A. Piggott, E. F. Marshall, N. Thomas, S. Lloyd, J. A. Court, E. Jaros, D. Burn, M. Johnson, R. H. Perry, I. G. McKeith, C. Ballard, and E. K. Perry. Striatal dopaminergic markers in dementia with lewy bodies, alzheimer's and parkinson's diseases: rostrocaudal distribution. *Brain*, 122(8):1449–1468, August 1999.
- Ronald B. Postuma and Daniela Berg. Prodromal parkinson's disease: The decade past, the decade to come. *Mov. Disord.*, 34(5):665–675, May 2019.
- Ronald B. Postuma, Alex Iranzo, Michele Hu, Birgit Högl, Bradley F. Boeve, Rafaële Manni, Wolfgang H. Oertel, Isabelle Arnulf, Luigi Ferini-Strambi, Monica

Puligheddu, Elena Antelmi, Valerie Cochen De Cock, Dario Arnaldi, Brit Mol-lenhauer, Aleksandar Videnovic, Karel Sonka, Ki-Young Jung, Dieter Kunz, Yves Dauvilliers, Federica Provini, Simon J Lewis, Jitka Buskova, Milena Pavlova, Anna Heidbreder, Jacques Y. Montplaisir, Joan Santamaria, Thomas R Barber, Ambra Stefani, Erik K. St Louis, Michele Terzaghi, Annette Janzen, Smandra Leu-Semenescu, Giuseppe Pazzoli, Flavio Nobili, Friederike Sixel-Doering, Petr Dusek, Frederik Bes, Pietro Cortelli, Kaylena Ehgoetz Martens, Jean-Francois Gagnon, Carles Gaig, Marco Zucconi, Claudia Trenkwalder, Ziv Gan-Or, Chris-tine Lo, Michal Rolinski, Philip Mahlknecht, Evi Holzknecht, Angel R. Boeve, Luke N. Teigen, Gianpaolo Toscano, Geert Mayer, Silvia Morbelli, Benjamin Dawson, and Amelie Pelletier. Risk and predictors of dementia and parkinsonism in idiopathic REM sleep behaviour disorder: a multicentre study. *Brain*, 142(3): 744–759, March 2019.

Amy Reeve, Eve Simcox, and Doug Turnbull. Ageing and parkinson's disease: Why is advancing age the biggest risk factor? *Ageing Res. Rev.*, 14:19–30, March 2014.

Laura Saari, Katri Kivinen, Maria Gardberg, Juho Joutsa, Tommi Noponen, and Valtteri Kaasinen. Dopamine transporter imaging does not predict the number of nigral neurons in parkinson disease. *Neurology*, 88(15):1461–1467, April 2017.

Tassilo Schiebler, Ivayla Apostolova, Franziska Lara Mathies, Catharina Lange, Su-sanne Klutmann, and Ralph Buchert. No impact of attenuation and scatter cor-rection on the interpretation of dopamine transporter spect in patients with clin-ically uncertain parkinsonian syndrome. *European Journal of Nuclear Medicine and Molecular Imaging*, 50(11):3302–3312, September 2023. ISSN 1619-7089. doi: 10.1007/s00259-023-06293-2.

Antti O. Sohlberg and Markus T. Kajaste. Fast monte carlo-simulator with full collimator and detector response modelling for SPECT. *Ann. Nucl. Med.*, 26(1): 92–98, January 2012.

Hermes Medical Solutions. Hybrid Recon. White paper, Hermes Medical Solutions.

Klaus Tatsch and Gabriele Poepperl. Nigrostriatal dopamine terminal imaging with dopamine transporter SPECT: an update. *J. Nucl. Med.*, 54(8):1331–1338, August 2013.

K. Tecklenburg, A. Forgács, I. Apostolova, W. Lehnert, S. Klutmann, J. Csirik, E. Garutti, and R Buchert. Performance evaluation of a novel multi-pinhole collimator for dopamine transporter SPECT. *Phys. Med. Biol.*, 65(16):165015, August 2020.

Eduardo Tolosa, Gregor Wenning, and Werner Poewe. The diagnosis of parkinson's disease. *Lancet Neurol.*, 5(1):75–86, January 2006.

Livia Tossici-Bolt, John C. Dickson, Terez Sera, Robin de Nijs, Maria Claudia Bag-nara, Catherine Jonsson, Egon Scheepers, Felicia Zito, Anita Seese, Pierre Mal-ick Koulibaly, Ozlem L. Kapucu, Michel Koole, Maria Raith, Jean George,

- Markus Nowak Lonsdale, Wolfgang Münzing, Klaus Tatsch, and Andrea Varrone. Calibration of gamma camera systems for a multicentre european ^{123}I -FP-CIT SPECT normal database. *Eur. J. Nucl. Med. Mol. Imaging*, 38(8):1529–1540, August 2011.
- Livia Tossici-Bolt, John C. Dickson, Terez Sera, Jan Booij, Susanne Asenbaum-Nan, Maria C. Bagnara, Thierry Vander Borght, Cathrine Jonsson, Robin de Nijs, Swen Hesse, Pierre M. Koulibaly, Umit O. Akdemir, Michel Koole, Klaus Tatsch, and Andrea Varrone. $[^{123}\text{I}]$ FP-CIT ENC-DAT normal database: the impact of the reconstruction and quantification methods. *EJNMMI Phys.*, 4(1):8, December 2017.
- Dominique Twelves, Kate S. M. Perkins, and Carl Counsell. Systematic review of incidence studies of parkinson’s disease. *Mov Disord*, 18(1):19–31, January 2003.
- Dennis Ulmer and Giovanni Cinà. Know your limits: Uncertainty estimation with relu classifiers fails at reliable ood detection. In Cassio de Campos and Marloes H. Maathuis, editors, *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, volume 161 of *Proceedings of Machine Learning Research*, pages 1766–1776. PMLR, 27–30 Jul 2021.
- Andrea Varrone, John C. Dickson, Livia Tossici-Bolt, Terez Sera, Susanne Asenbaum, Jan Booij, Ozlem L. Kapucu, Andreas Kluge, Gitte M. Knudsen, Pierre Malick Koulibaly, Flavio Nobili, Marco Pagani, Osama Sabri, Thierry Vander Borght, Koen Van Laere, and Klaus Tatsch. European multicentre database of healthy controls for $[^{123}\text{I}]$ FP-CIT SPECT (ENC-DAT): age-related effects, gender differences and evaluation of different methods of analysis. *Eur. J. Nucl. Med. Mol. Imaging*, 40(2):213–227, January 2013.
- Markus Wenzel, Fausto Milletari, Julia Krüger, Catharina Lange, Michael Schenk, Ivayla Apostolova, Susanne Klutmann, Marcus Ehrenburg, and Ralph Buchert. 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*, 46(13):2800–2811, December 2019.
- William J. Youden. Index for rating diagnostic tests. *Cancer*, 3(1):32–35, 1950.

Declaration of Authorship

Ich erkläre hiermit gemäß §9 Abs. 12 APO, dass ich die vorstehende Abschlussarbeit selbstständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe. Des Weiteren erkläre ich, dass die digitale Fassung der gedruckten Ausfertigung der Abschlussarbeit ausnahmslos in Inhalt und Wortlaut entspricht und zur Kenntnis genommen wurde, dass diese digitale Fassung einer durch Software unterstützten, anonymisierten Prüfung auf Plagiate unterzogen werden kann.

Place, Date

Signature