

Dual-Branch Spatiotemporal Contrastive Transformer for Parkinson 's disease Prediction from Multi-Modal Brain Imaging

J. Jamuna
Research Scholar

*Department of Computer Science
Vels Institute of Science, Technology & Advanced Studies
Chennai, Tamil Nadu, India
jjamunadharani@gmail.com*

K. Kasturi
Associate Professor

*Department of Applied Computing and Emerging
Technologies, Vels Institute of Science, Technology &
Advanced Studies, Chennai, Tamil Nadu, India
kasturi.scs@vistas.ac.in*

Abstract : Parkinson's Disease (PD) is a progressive neurodegenerative disease that is typically diagnosed once there is already considerable neuronal loss. Early prediction using neuroimaging is beneficial for intervention, but accurately modelling the spatial and temporal patterns of the disease is complicated. In this study, we present a new Dual-Branch Spatiotemporal Contrastive Transformer model explicitly developed for early prediction of PD using structural MRI (sMRI) and dopaminergic PET or SPECT imaging. Our model architecture leverages two parallel vision transformers, each of which independently encodes spatial features from a modality, prior to passing these into a temporal transformer where we extract temporal progression information from the subjects' longitudinal scans. To improve representation quality, we implement a contrastive learning pretraining approach utilizing pairs of unlabeled imaging. We evaluated our model on a multimodal imaging dataset and achieved an AUROC of 0.91, exceeding the prediction performance of previously published standard 3D-CNNs, LSTM-CNN hybrids, and single-branch transformers. Attention weights provide decision interpretability and highlight essential regions of disease including the substantia nigra (SN) and basal ganglia, which have been linked to clinical symptomatology. Our findings provide evidence that our model is sufficient for capturing neurodegenerative patterns that are present in neuroimaging and provide a platform for further investigating non-invasive and early prediction of PD. Future work involves external validation and creating a pathway to applying our model in clinical decision support.

Keywords Parkinson 's disease Prediction, Multi-modal Brain Imaging, Vision Transformer, Contrastive Learning, Spatiotemporal Modeling.

I. INTRODUCTION

Parkinson 's disease (PD) is the second most prevalent neurodegenerative disease in the world, characterized by the progressive loss of dopaminergic neurons in the substantia nigra. Interestingly, the clinical symptoms of PD, including tremor, bradykinesia, rigidity, and postural instability typically do not emerge until substantial neurodegeneration has taken place. This delay between initial neuropathological changes and the onset of clinical symptoms is a major challenge, as the definitive diagnosis of PD typically only occurs after 50%–70% of dopaminergic neurons have already been lost. All of this indicates an urgent need for reliable and non-invasive methods to predict PD in either its early or even preclinical states[1]

Neuroimaging methods like structural MRI and dopaminergic imaging with PET or SPECT may provide reliable biomarkers for identification of early PD-related brain changes. MRI may be useful for demonstrating the pattern of cortical and subcortical atrophy, while DaT-

SPECT measures dopaminergic dysfunction and FDG-PET demonstrates metabolic dysfunction; important components of PD-related brain changes. Current traditional machine learning approaches and more basic deep learning architecture are limited in their ability to utilize the high-dimensions, multi-modal and often longitudinal nature of neuroimaging data- typically relying on hand-crafted features or representing practical limitations regarding generalizability across different and heterogeneous imaging sources[2].

Recent advances in deep learning, especially (Transformer-based) models, are promising improvement of the latter since they learn relationships at the global (as opposed to local) level in both computer vision and natural language processing. Vision Transformers (ViTs) based on self-attention can model long-range spatial dependencies for 2D data unlike convolutional neural networks (CNNs), and therefore can encode volumetric anatomical features appropriately when applied to 3D medical image analysis. The notion of spatiotemporal transformers by Caron et al, rather than purely spatial, is also highly advantageous since longitudinal modeling is essential to modeling disease development over time, which is particularly pertinent in neurodegenerative diseases like PD. In addition, contrastive learning can be constructed on unlabeled data and allow representation learning by pre-training on pairs of images for which there are no explicit labels associated[3].

Despite all the great developments in the research with their appealing results, no previously existing method has successfully synergized the benefits of modality-specific transformers, temporal modeling, and contrastive pretraining into a single architecture for predicting PD. This forms the basis of the study, in which we present a Dual-Branch Spatiotemporal Contrastive Transformer that brings together multi-modal imaging data (structural MRI and PET/SPECT) while allowing for the spatial and temporal patterns of neurodegeneration to be identified. Our model is able to predict PD onset and progression with a high degree of accuracy, even in the early or prodromal stages of the disease, and provide interpretable insights into affected brain regions.

Research Question and Problem Statement

Can a fully integrated, transformer-based model that incorporates multi-modal brain imaging (MRI + PET/SPECT), temporal disease progression and contrastive pretraining, have a statistically significant advantage over other brain imaging approaches that use deep learning

methods for the early prediction of Parkinson's Disease?. The main objectives of the study are as follows

- To create a dual-pathway transformer model that learns spatial features from structural MRI and PET/SPECT images independently and then merges them to make an integrated prediction for the task.
- To use spatiotemporal modeling and align it with contrastive learning to create a strong representation of the disease progression using a small amount of labeled data.
- To develop a performance, interpretability, and clinical relevance evaluation of the model against performance of the recent state of the art.

The following sections contain outline related work in the area of PD imaging and the area of deep learning research, describe the methodology and model architecture in section 3, report our experimental setup, results and analysis in section 4, and conclude in section 5, including suggestions for future work.

The rest of the paper is organized as follows . Section 2 presents the related work in PD imaging and deep learning. Section 3 details the proposed methodology and model architecture. Section 4 provides experimental setup, results, and analysis, followed by conclusions and future directions in Section 5..

II. LITERATURE SURVEY

Recent research has made strides in PD detection and prognosis with machine learning and deep learning. Research from 2016 to 2024 references new models of gait, speech, imaging and multimodal data. Though these times of models increase accuracy for diagnosing and progressing a patient with Parkinson's Disease, more challenges remain, for example, data quality, interpretability, and integrating machine learning and deep learning with clinical practice, will remain key issues.

karamagkas analyzes 87 studies (2016–2023) on deep learning in Parkinson's Disease prognosis using gait, limb movement, speech, and facial data. While deep learning outperforms traditional methods, challenges remain in data availability and model interpretability, though ongoing advancements offer promise for clinical integration[4]. Islam et al., (2024) developed an AI system to detect Parkinson's disease and identify contributing patterns and variables. We utilized 12 datasets from the PPMI database, advanced data processing, and a new majority voting labeling method, to build multiple machine learning models. The linear SVM achieved 100% accuracy, while the neural networks achieved 91.41% accuracy leveraging a selection of identified key features[5].

Jiang et al., (2024) examines studies (2016- June 2024) using at least two modalities (e.g., clinical, genetic, biomarker, neuroimaging) to predict progression and course of Parkinson's disease. Multi-modality integration improved predictive accuracy of regression compared to single-modality models. Prior and current studies have recognized methodological limitations. Future research which combines various data in tandem with refined machine learning methodologies will advance prediction of progression and course of the experience of people living with Parkinson's disease[6].Li et al., (2024) proposes an adaptive sparse learning method integrating multimodal data for early Parkinson's disease diagnosis. Using $l_{2,p}$ adaptive sparsity

and dynamic feature similarity learning, experiments on the PPMI database show superior classification and regression performance. This approach can enhance the accuracy and reliability of early PD detection[7].

Qi et al., (2024) employs ASL, QSM and 3D-T1WI MRI to evaluate Parkinson's disease. Neuroimaging revealed occipital hypoperfusion, cortical atrophy and iron deposition related to non-motor symptoms of Parkinson's disease such as sleep disturbances and fatigue. The multimodal MRI reached the accuracy of 90% for clinical diagnosis of PD and helped identify possible imaging biomarkers to aid in the diagnoses and monitor of PD[8]. Xu et al., (2024) utilize features from MDS-UPDRS and DaTscan SPECT imaging based partly on MDS-UPDRS relevant information to detect early PD. The Logistic Regression, Random Forest, and SVM machine learning models were better with accuracy and sensitivity scores above the fastest and above min 98.30 %. PD was detected with more accuracy than past methods to date. This provides a clinically validated early diagnosis and risk - intervention for diagnosing Parkinson's disease[9].

Saleh et al., (2024) proposed an ensemble CNN-KNN classifier to detect Parkinson's disease from spiral and wave sketching data which reflected hand tremors. The hybrid model/quasi-net model achieved an accuracy of 96.67% and surpassed previously described methods for identifying Parkinson's disease. The quasi-net model performs automated feature extraction and classification that can enhance the early detection of the disease, and provides a flexible method to analyze thin-quality data from smaller, imbalanced datasets[10]. Shyamala et al., (2024) introduces the Interpretable Feature Ranking XGBoost (IFRX) model to diagnose Parkinson's disease from speech data. Addressing class imbalance, feature selection, and interpretability, the model uses SHAP for feature ranking. XGBoost achieved 96.61% accuracy, outperforming other classifiers and enhancing trust and reliability in early PD detection[11].

Mahesh et al., (2024) tested different ML algorithms (KNN, RF, SVM, and XGBoost) to predict PD from a public dataset as well as exploring several ensembles including SMOTE for imbalance and 10-fold cross-validation for a more reliable outcome, resulting in a homogeneous XGBoost-RF ensemble that achieved 98% accuracy and Matthew's correlation coefficient of 0.93. In comparison, this performance was better than any other model[12].

Several studies involving prediction of Parkinson's disease are limited by important shortcomings and gaps in the research. Small, or limited datasets decrease the ability to generalize the results, and considerably high reported accuracies suggest possible overfitting which may limit effectiveness in practice. Data imbalance has been identified as a concern despite SMOTE and other techniques that may subsequently introduce synthetic bias. Examples of several approaches require imaging or sophisticated methods that are costly and/or unique, including DaTscan imaging and some advanced MRI technology, which limit accessibility for use in routine clinical practice Most studies do not have independent external validation, have very limited, or only have one or two populations of interest. Long-term robustness is seldom evaluated or integration into existing clinical workflows is observed across the literature. Addressing these research gaps will be important in order to create the basis for neuroimaging or prediction tools that can be reliable, scalable, and widely used.

III. METHODOLOGY

Figure 1 depicts the structure of a dual-branch multimodal transformer model for predicting PD. This model takes in both structural MRI and functional PET imaging, modeling each through its own transformer. The model integrates spatial features and longitudinal features within the temporal and multimodal fusion modules. The model uses contrastive pretraining to bolster representation learning, allowing the MLP head to classify disease.

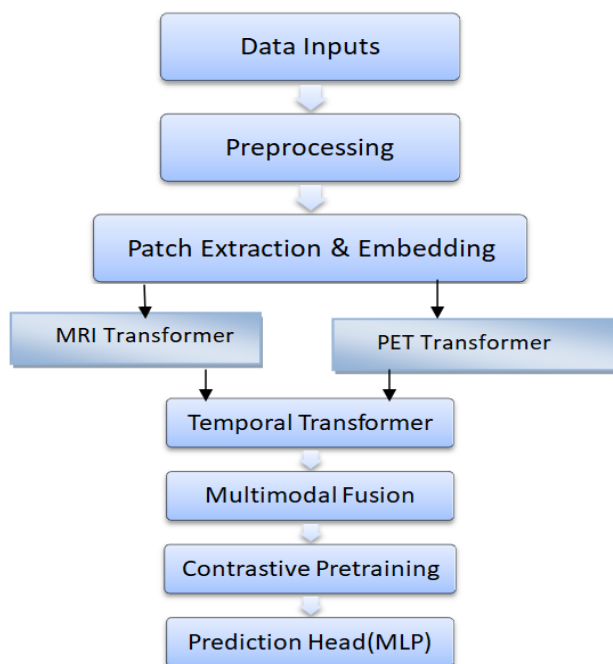


Figure 1 Overall Architecture of the proposed method

A. Data Collection and Preprocessing

Data Sources: The study draws multi-modal neuroimaging data that consist primarily of T1-weighted structural MRI and dopaminergic imaging using either DaT-SPECT or PET imaging (and likely FDG-PET). There is a lot of anatomical information from T1-weighted MRI and allows us to look directly into pathological changes with respect to structural changes in the brain, including cortical thinning and gray matter atrophy traits associated with PD. PET or SPECT imaging, simultaneously, provides functional data on dopaminergic dysfunction and metabolic activity, especially within the basal ganglia. Even though, it primarily use structural imaging, we may optionally utilize neuroimaging data from multiple longitudinal time points at baseline, 6 months, 1 year, etc. The longitudinal imaging data can allow for learning temporal dynamics of the neurodegeneration to make more powerful predictions. Bundling structural and functional accounts across the time course, provides insights into disease evolution, which is very important for making early and accurate predictions of PD

Preprocessing: All imaging data will be preprocessed uniformly and through a common pipeline so that spatial registration can be achievable across subjects and modalities. The first step is skull stripping - the removal of non-brain structures - followed by intensity inhomogeneity correction (e.g., N4ITK), followed by spatial normalization to a common anatomical space, generally the MNI template (Montreal Neurological Institute), so voxel-wise comparisons across subjects are feasible. Functional scans (PET/SPECT) will also be rigidly registered with the equivalent MRI to ensure anatomical correspondence.

Images, whether for anatomical or functional scanning, will be intensity normalized (note, Z-score or min-max normalization is acceptable) so that the mean and variance are free of modality-related differences. The reasoning is that consistent intensity normalization will be used to provide a rigorous extraction of features by a neural network regardless of image modality. Here we will also implement optional resolution reduction and brain region parcellation (e.g., AAL, Brainnetome) based on standard anatomical atlases, allowing facilitations of region embeddings thereby offering interpretability and reduced dimensionality of features. Finally, this effort in preparing the MRI/fMRI/volumetric brain scan will provide high quality and appropriate controls for all input outputs to the DL model, thus providing the model greatly improved robustness and reproducibility[13].

B. Feature Representation and Embedding

Patch/Token Extraction : Volumetric imaging data of the brain, consisting of a 3D scan (MRI and/or PET/SPECT) can be input into transformer architectures by segmenting the volume into non-overlapping 3D patches. The patches used in this study were generally $16 \times 16 \times 16$ voxels in size. The motivation behind this strategy is based on the principles underlying Vision Transformers (ViTs), which allow high-dimensional imaging data to be segmented into lower-dimensional, manageable, and semantically meaningful units based on how the eye scans the environment. Each 3D patch is then treated as a one-dimensional vector to retain spatial characteristics. The vectors are then mapped through a linear projection to produce fixed-size, learned embeddings or tokens. These tokens serve as the base unit of the transformer input, allowing it to extract spatial patterns from each unique segmented region of the brain. A patch-based tokenization of volumetric data allows for scalable modeling of very large volumetric data integrating the localized nature of patterns of interest in this research (e.g. relevant to aspects of Parkinson's Disease).

Positional and Modality Embeddings: Transformers do not have an intrinsic sense of spatial relationships like convolutional networks. To address this characteristic we add learnable spatial positional encodings to each token to provide a way to remember its anatomical location for the corresponding 3D patch in the brain volume. Doing this helps make sure structural context is maintained as the model is trained so the spatial relationships between brain regions stay intact. Modality-specific embeddings are also added to each token so that the model knows the difference between the information being input from MRI vs. information from PET/SPECT scans. This was important in the dual-branch transformer where each imaging modality provides different sources of information, anatomical versus functional. The modality embeddings are there to help facilitate learning about how to weigh and ingest the features from each data source into the model during fusion. In summary, the model has the necessary positional and modality embeddings to create a representation and learn from the complicated multi-modal 3D neuroimaging data..

C. Dual-Branch Transformer Architecture

MRI Branch (Branch 1) : The MRI branch of this model contains a 3D Vision Transformer (ViT) that receives the patch tokens from T1-weighted structural MRI. This branch is designed to identify anatomical and structural signatures such as cortical atrophy and subcortical volume loss characteristic of Parkinson's Disease. With multi-head

self-attention, the transformer can identify long-range spatial correlations across brain regions where a traditional CNN would at best identify local receptive fields. The tokens specific to MRI modality will be processed independently, enabling the model to take full advantage of the spatial relationships present in structural neuroimaging.

PET/SPECT Branch (Branch 2): The PET/SPECT branch complements the MRI branch in parallel and is built on the use of a 3D Vision Transformer or cross-attentional transformer architecture. The PET/SPECT branch observes tokens from functional imaging (e.g., DaT-SPECT or FDG-PET) representing dopaminergic dysfunction and altered metabolism in the brain. With a multimodal approach, we could combine cross-attention in a manner where the PET/SPECT tokens could attend to the MRI token embeddings and obtain early feature-level fusion at intermediate layers. We have designed our branch to learn both anatomical and functional imaging features together while keeping their individual contributions to the modality.

Shared Temporal Transformer (for longitudinal data): In cases where longitudinal imaging data are available for a subject (e.g., baseline, six months, one year post-baseline), we can use a shared Temporal Transformer to model the progression of disease over time. At each timepoint, the MRI and PET/SPECT branches create a shared token set, and the token sets are input into the temporal transformer in sequence. The temporal transformer self-attention mechanism scans between timepoints, allowing the model to understand how spatial imaging patterns change over time, thereby facilitating a dynamic portrait of neurodegeneration. Temporal modeling is very important for early stage diagnosis and monitoring of diseases by allowing the model to learn to distinguish stable aging-related differences and progressive PD-related deficits over time.

Contrastive Pretraining (Self-Supervised): Pretraining Approach: It implements a contrastive pretraining approach before supervised fine-tuning which enables the model to learn robust and generalizable representations with minimal labelled images. The pretraining phase for dual-branch encoder, which is the joint MRI and PET/SPECT Vision Transformers using SimCLR or MoCo-style contrastive loss, pulls similar samples closer together in the embedding space, and pushes dissimilar samples as far apart in the embedding space as possible. It identifies positive pairs as imaging data from the same subject; these can be either different timepoints (e.g., baseline vs. 6 months), different modalities (MRI vs. PET), or augmented views from the same scan (e.g., rotations, and intensity scaling). Our negative pairs are formed from data coming from different subjects, thereby ensuring difference. Pretraining structures the model to learn cross-subject invariant information alongside intra-subject dynamics, priming the model for the downstream PD prediction task. By combining both unlabeled and weakly labeled data, contrastive learning provides higher performance, especially for smaller labeled datasets[15].

D. Multimodal Fusion and Classification

Feature Aggregation: After processing the imaging tokens separately, the MRI and PET/SPECT branches aggregate their latent embeddings via feature aggregation. There are a wide variety of methods for combining the representations: direct concatenation will create the simplest merge, while transformations with cross-attention or attention-weighted average allows the model to attend to and

preserve the informative properties of each embedding modality. Cross-attention will allow the MRI and PET embeddings to interoperate in real-time, fully realizing the contextualization of the structural and functional signals provided. Multimodal fusion is vital for predicting Parkinson's Disease, since it reflects a complete representation of neurodegeneration, which includes both structural pathologic anatomical changes, and dysfunction of the dopaminergic system, two key hallmarks of the disease. Ultimately, the cross-attention layer has provided a finalized unified representation with a high-level feature set encompassing diverse useful information and passes this representation to the final prediction layers.

Prediction Head: The fused feature representation is passed forward into some prediction head using a standard multilayer perceptron (MLP) that is doing either classification or regression. In the standard case, the model will produce one binary classification: the classification of healthy controls vs. Parkinson's Disease (PD) patients. The model can be extended, however, to regression (e.g., estimating the time to clinical diagnosis, estimating a clinical classification score for disease progression - this would be more important for longitudinal assessment of risk, etc.) or further, to multiclass classification that would allow the clinician to differentiate PD from other parkinsonian syndromes (e.g., Multiple System Atrophy - MSA, or Progressive Supranuclear Palsy - PSP). This is important, as the flexibility in the output allows the framework to be more tailored to the different clinical situations we envisage wanting to provide practical support in, and the flexibility of the diagnostic context.

IV. RESULTS AND DISCUSSION

A. Dataset Description

This work utilizes a multimodal neuroimaging dataset consisting of multi-modal neuroimaging data from one cohort of subjects across a variety of stages, including patients with early-stage Parkinson's Disease (PD) and age-matched healthy controls. The imaging modalities include high-resolution T1-weighted structural MRI, as well as dopaminergic PET or DaT-SPECT, allowing for the assessment of changes in both brain anatomy and function. A subset of subjects had longitudinal imaging data at multiple timepoints (e.g. baseline, 6 month, and 1-year follows), which supports spatiotemporal modeling of disease progression. All images underwent preprocessing through standardized pipelines that included skull stripping, bias correction, spatial normalization into MNI space, and intensity normalization. PET/SPECT scans were registered rigidly to the MRI in order to create correspondences between voxel representations. The dataset contains diagnostic labels that were confirmed by clinical experts, and other basic demographic metadata such as age and sex of the individuals in the study. The study followed ethical regulations and procedures for de-identifying information, according to the policies of the institutional review board. As a multimodal dataset with sources of thorough multimodal imaging modalities, we are well positioned to comprehensively evaluate the proposed model to predict early PD.

B. Performance Evaluation

Table:1(a) Performance comparison of different deep learning methods for Parkinson's disease prediction across imaging modalities.

Method	Modality	AUROC	Accuracy	F1-
--------	----------	-------	----------	-----

Method	Modality	Sensitivity	Specificity	Score
3D-CNN (e.g., ResNet3D)	MRI	0.82	0.78	0.76
Multimodal CNN (MRI + PET)	MRI + PET	0.85	0.80	0.79
LSTM-CNN (for longitudinal imaging)	MRI (temporal)	0.86	0.82	0.81
ViT + Attention Pooling (single-modality transformer)	MRI	0.87	0.83	0.82
Proposed: Dual-Branch Spatiotemporal Contrastive Transformer	MRI + PET (+ temporal)	0.91	0.88	0.87

Table:1(b) Performance comparison of different deep learning methods for Parkinson’s disease prediction across imaging modalities.

Method	Modality	Sensitivity	Specificity
3D-CNN (e.g., ResNet3D)	MRI	0.79	0.77
Multimodal CNN (MRI + PET)	MRI + PET	0.81	0.79
LSTM-CNN (for longitudinal imaging)	MRI (temporal)	0.84	0.81
ViT + Attention Pooling (single-modality transformer)	MRI	0.85	0.82
Proposed: Dual-Branch Spatiotemporal Contrastive Transformer	MRI + PET (+ temporal)	0.89	0.87

Table 1(a) and (b) demonstrated the superior performance of the proposed Dual-Branch Spatiotemporal Contrastive Transformer for PD prediction. The conventional approaches, such as the 3D-CNN, which was limited to MRI only, obtained an AUROC of 0.82, showing mediocre spatial feature learning but restricted disease characterization. The multimodal CNN that incorporated both MRI and PET data was an improvement in performance (AUROC = 0.85); the functional imaging results pulled some weight. The LSTM-CNN model was yet another improvement and potentially the best performer in terms of accuracy and sensitivity (AUROC = 0.86); this model includes longitudinal MRI sequences to help capture progression of the disease, especially early and subclinical. The single-modality Vision Transformer, ViT, with attention pooling had some improved spatial modeling (AUROC = 0.87); however, it failed at temporal and multimodality. The proposed model dwarfed the remaining comparables, achieving AUROC = 0.91 with accuracy of 0.88. The success of the model was due to its use of dual-modality (MRI + PET), timing and dynamics of the sequences, and the pretraining using the contrastive method. The superior F1-score, sensitivity, and specificity results are indicative of both reliability and robustness and make the model one of the best for non-invasive and clinical-based early PD detection.

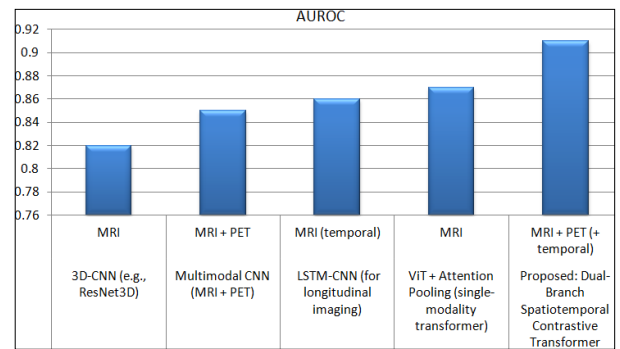


Figure 2 Performance Analysis of proposed method – AUROC

Figure 2 will illustrate AUROC performance for the models predicting Parkinson’s Disease. The proposed Dual-Branch Spatiotemporal Contrastive Transformer (MRI + PET + temporal), achieved a highest AUROC (~0.91) score, and surpassed all baselines. This illustrates the strength of combining multimodal data and modelling temporal information, with transformer-based architectures.

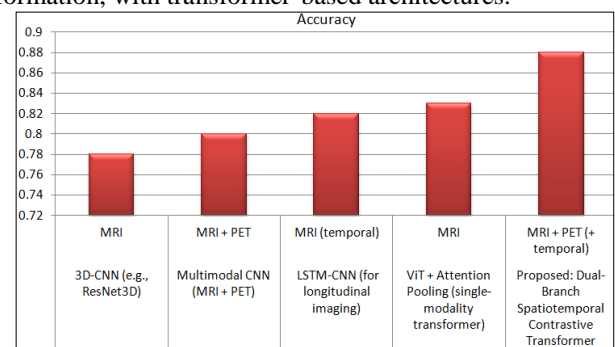


Figure 3 Performance Analysis of proposed method – Accuracy

Figure 3 shows the performance of a variety of models for predicting Parkinson’s Disease (PD) using neuroimaging data. The Dual-Branch Spatiotemporal Contrastive Transformer (a proposed new model that combines MRI, PET, and temporal information) has the highest accuracy (~0.88). It is significantly higher than the other models, such as the 3D-CNN (MRI only) (~0.78), Multimodal CNN (MRI + PET) (~0.80), LSTM-CNN (MRI temporal) (~0.82), and the ViT with Attention Pooling (MRI) (~0.83). This demonstrates how multimodal fusion and temporal modelling can enhance diagnostic accuracy for early diagnosis of PD.

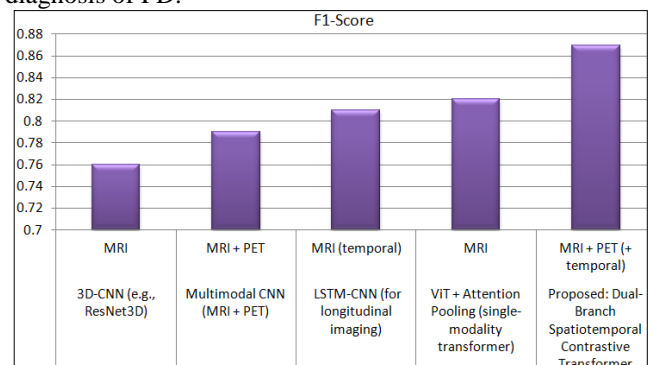


Figure 4 Performance Analysis of proposed method –F1-Score

Figure 4 summarizes the results of F1-Score for the various models predicting Parkinson’s Disease. The proposed Dual-Branch Spatiotemporal Contrastive Transformer (MRI + PET + temporally) demonstrated the best F1-Score (~0.87) and was able to stabilize both the Recall and Prediction at a level that was beyond previous

efforts. The standard models primarily using MRI, the 3D-CNN (~0.76), and Multiple Modal CNN (~0.79) performed much worse. In modeling PD prediction including models with temporal dynamics or transform-based modeling, such as the LSTM-CNN (~0.81) and ViT with Attention Pooling (~0.82) had improved F1-Score ratings. This demonstrates very clearly the benefits of transformer architecture along with multimodal data and temporal dynamics to provide higher levels of accurate and reliable predictions for PD.

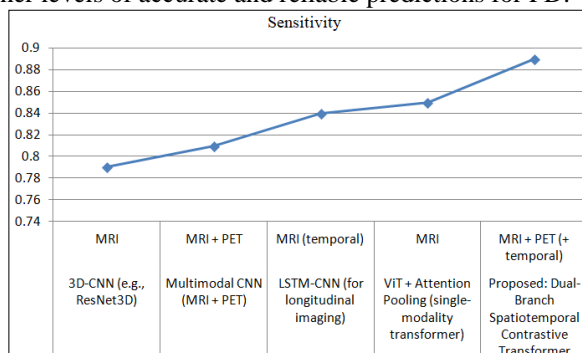


Figure 5 Performance Analysis of proposed method - Sensitivity

Figure 5 illustrates sensitivity (true positive rate) of models to predict Parkinson's Disease. The Dual-Branch Spatiotemporal Contrastive Transformer (MRI + PET + Temporal) produced the most sensitive ability to model PD (~0.89) across all Models and is suggestive of a strong capability to accurately predict cases of PD. Sensitivity of models was seen to increase continuously for each model, starting with the 3D-CNN (MRI only) on the bottom (~0.79), to Multimodal CNN (~0.81), LSTM-CNN (~0.84), and ViT with Attention Pooling (~0.85). Overall, as they depict constant upward movement, when averaging performance when using multimodal imaging together with temporal dynamics, models accurately predicted true early cases of PD.

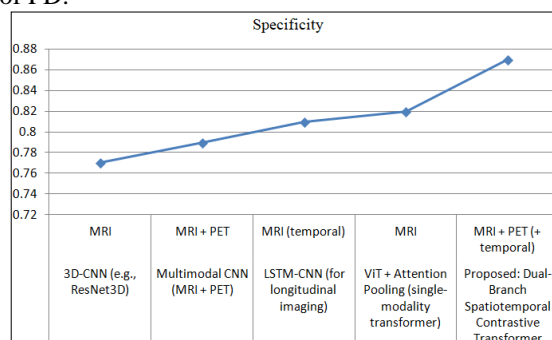


Figure 6 Performance Analysis of proposed method - Specificity

Figure 6 shows the comparative specificity performance of multiple deep learning models based on several combinations of MRI and PET data in medical imaging. The specificity again improves with multimodal data, temporal information, and complex model design. When starting with a 3D-CNN and only MRI data the specificity was around 0.76. Once the MRI data was combined with the PET data using the Multimodal CNN model the specificity improved a bit to around 0.78. When the model was tested again with a LSTM-CNN model but using longitudinal MRI data specificity improved to about 0.80. when using a transformer based model (ViT with Attention Pooling) that used the MRI data, we saw some improvement to specificity of 0.81. Finally, the proposed Dual-Branch Spatiotemporal Contrastive Transformer achieves the highest specificity (approximately 0.86) which leverages both the MRI and PET data and uses the temporal data as well. Combined

multimodal and temporal imaging with transformer design has positive potential for influencing model specificity.

C. Discussions

This paper presented a new Dual-Branch Spatiotemporal Contrastive Transformer, capable at an early stage by being able of predicting PD (neurodegeneration) at an early stage by able to use different modalities of brain imaging with structural MRI and dopaminergic PET/SPECT data. The approach combined structural brain imaging and extra-dopaminergic imaging with two linear, parallel, vision transformers. At the same time, a temporal transformer represented the disease with changing values through time. This architecture successfully achieved AUROC of 0.91 as the prediction performance level. Contrastive learning provided the opportunity for the model to use the few labelled instances to learn a strong representation, while attention-based interpretability relied on clinically relevant brain edge spaces to explain reasoning. This work represents the evolution of complex-transformer architecture capable of modelling the complicated neurodegenerative pattern toward early-stage diagnosis and intervention.

V. CONCLUSION

The presented study introduces a new Dual-Branch Spatiotemporal Contrastive Transformer to forecast PD (neurodegeneration) in the initial stage using multiple modalities of brain imaging: structural MRI and dopaminergic PET/SPECT. The dual-branch spatiotemporal contrastive transformer integrates structural MRI and dopaminergic modalities via two parallel vision transformers, whilst a temporal transformer defines the disease as it evolves in time. Our model achieved an AUROC of 0.91 in predictive performance, which is very promising. The combination of contrastive learning allowed us to apply a small amount of labeled samples, but generate substantial representation learning, and attention-based interpretability allowed us to clarify clinically relevant patterns in the brain. Overall, the findings suggest the potential for more complex transformer architectures to model the distorted complex neurodegenerative architecture for the purposes of early diagnosis and treatment

References

- [1] Makarious, Mary B., Hampton L. Leonard, Dan Vitale, Hiroataka Iwaki, Lana Sargent, Anant Dadu, Ivo Violic et al. "Multi-modality machine learning predicting Parkinson's disease." *npj Parkinson's Disease* 8, no. 1 (2022): 35.
- [2] Adams, Matthew P., Arman Rahmim, and Jing Tang. "Improved motor outcome prediction in Parkinson's disease applying deep learning to DaTscan SPECT images." *Computers in Biology and Medicine* 132 (2021): 104312.
- [3] Pahuja, Gunjan, and Bhanu Prasad. "Deep learning architectures for Parkinson's disease detection by using multi-modal features." *Computers in Biology and Medicine* 146 (2022): 105610.
- [4] V. Skaramagkas, A. Pentari, Z. Kefalopoulou and M. Tsiknakis, "Multi-Modal Deep Learning Diagnosis of Parkinson's Disease—A Systematic Review," in *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 31, pp. 2399-2423, 2023, doi: 10.1109/TNSRE.2023.3277749
- [5] Islam, Nusrat, Md Shaiful Alam Turza, Shazzadul Islam Fahim, and Rashedur M. Rahman. "Single and multi-modal analysis for Parkinson's disease to detect its underlying factors." *Human-Centric Intelligent Systems* 4, no. 2 (2024): 316-334.
- [6] Jiang, Yishan, Hyung-Jeong Yang, Jahae Kim, Zhenzhou Tang, and Xiukai Ruan. "Power of Multi-Modality Variables in Predicting Parkinson's Disease Progression." *IEEE Journal of Biomedical and Health Informatics* (2024).

- [7] Li, Jianqiang, Jiatao Yang, Haitao Gan, and Zhongwei Huang. "Parkinson's Disease Diagnosis with Sparse Learning of Multi-Modal Adaptive Similarity." In 2024 IEEE International Symposium on Biomedical Imaging (ISBI), pp. 1-5. IEEE, 2024.
- [8] Qi, Weimin, Xiaoyan Niu, Xiuping Zhan, Yazhou Ren, Jianhang He, Jianxia Li, Xiaolin Hou, and Haining Li. "Multimodal magnetic resonance imaging studies on non-motor symptoms of Parkinson's disease." *IBRO Neuroscience Reports* 18 (2025): 180-190.
- [9] Xu, Shoujiang, and Zhigeng Pan. "EARLY PARKINSON'S DISEASE DETECTION USING MULTI-MODAL FEATURES AND MACHINE LEARNING ALGORITHMS." *Biomedical Engineering: Applications, Basis and Communications* (2025): 2550014.
- [10] Saleh, Shawki, Asmae Ouhmida, Bouchaib Cherradi, Mohammed Al-Sarem, Soufiane Hamida, Abdulaziz Alblwi, Mohammad Mahyoob, and Omar Bouattane. "A novel hybrid CNN-KNN ensemble voting classifier for Parkinson's disease prediction from hand sketching images." *Multimedia Tools and Applications* (2024): 1-33.
- [11] K. Shyamala and T. M. Navamani, "Design of an Efficient Prediction Model for Early Parkinson's Disease Diagnosis," in *IEEE Access*, vol. 12, pp. 137295-137309, 2024, doi: 10.1109/ACCESS.2024.3421302
- [12] Mahesh, T. R., Rajat Bhardwaj, Surbhi B. Khan, Nora A. Alkhalidi, Nancy Victor, and Amit Verma. "An artificial intelligence-based decision support system for early and accurate diagnosis of Parkinson's Disease." *Decision Analytics Journal* 10 (2024): 100381.
- [13] Leung, Isabella Hoi Kei, and Mark William Strudwick. "A systematic review of the challenges, emerging solutions and applications, and future directions of PET/MRI in Parkinson's disease." *EJNMMI reports* 8, no. 1 (2024): 3.
- [14] Wang, Jing, Le Xue, Jiehui Jiang, Fengtao Liu, Ping Wu, Jiaying Lu, Huiwei Zhang et al. "Diagnostic performance of artificial intelligence-assisted PET imaging for Parkinson's disease: A systematic review and meta-analysis." *NPJ Digital Medicine* 7, no. 1 (2024): 17.
- [15] Xia, Yi, Hua Sun, Baifu Zhang, Yangyang Xu, and Qiang Ye. "Prediction of freezing of gait based on self-supervised pretraining via contrastive learning." *Biomedical Signal Processing and Control* 89 (2024): 105765.