# An Explainable AI Approach to ADHD and Gender Classification
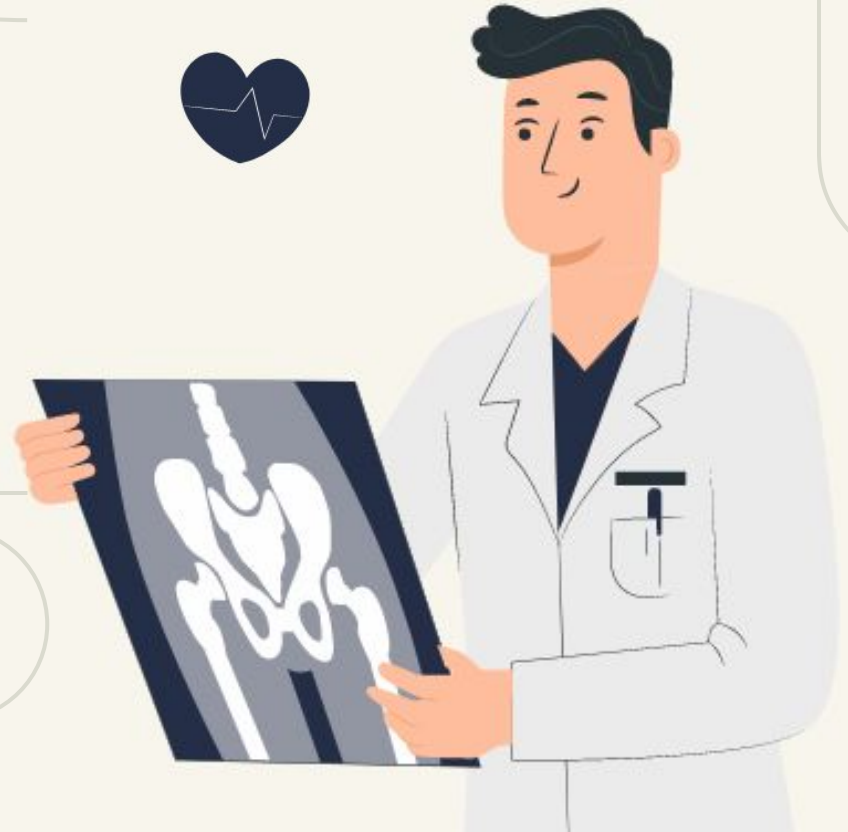
**Samia Rahman Misty**
Registration No: 2400570
Course No: CE- 888

# Introduction

## What is ADHD?

A neurodevelopmental disorder affecting attention, impulse control and activity levels.

## How Dangerous is ADHD?

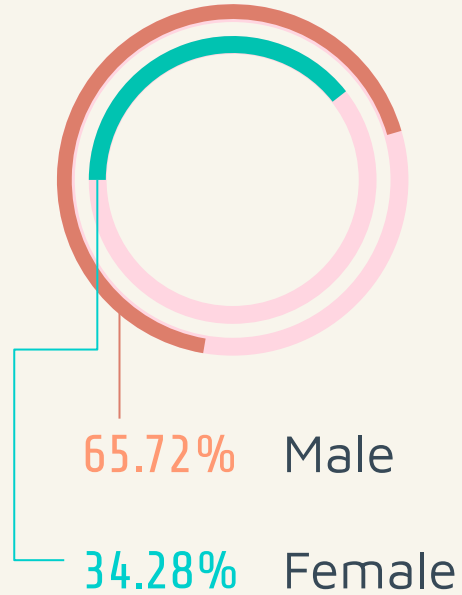If untreated, can lead to academic, relationship issues, low self-esteem, and mental health risks.

## How Deep Learning Helps

- Analyzes complex data (e.g., fMRI, behavior) for accurate ADHD diagnosis.
- Example: Deep learning achieved 70% accuracy in ADHD diagnosis using fMRI data[1].

1.Mao, Z., Su, Y., Xu, G., Wang, X., Huang, Y., Yue, W., ... & Xiong, N. (2019). Spatio-temporal deep learning method for adhd fmri classification. Information Sciences, 499, 1-11.

# Dataset

## GENDER



65.72% Male

34.28% Female

401 ADHD

157 Non-ADHD

174 ADHD

117 Non-ADHD

# Objective



Is it possible to build a model to predict both an individual's sex and their ADHD diagnosis?

Is the model biased in its predictions or training data?

Can the model misdiagnose females due to bias?

# Data Integrity and Optimization

### Data Format Standardization

Variations in data formats and types were observed during collection. for example, the Edinburgh Handedness Questionnaire was aligned with the standard format[1].

### Missing Data

Some data entries were missing during collection, which were carefully handled through appropriate data imputation techniques.

### fMRI Data Handling

The fMRI brain imaging data was large and complex. It was processed using appropriate techniques to ensure efficiency and maintain data integrity.

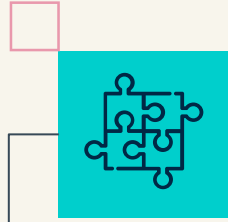### Feature Redundancy

In the Strength and Difficulties Questionnaire, the Total Difficulties Score, Externalizing Score, and Hyperactivity Scale were highly correlated. To avoid redundancy, only one of these was retained.

# Selected Models

## 01

### Logistic Regression

Its simplicity reduces the risk of errors when working with limited data.

## 02

### CNN + MLP Hybrid Architecture

1.CNN is effective at finding patterns in numbers, especially when those numbers are arranged in a meaningful order, like images or signals.

2.MLP part is good at understanding group-based information, such as age, gender, or other categories.
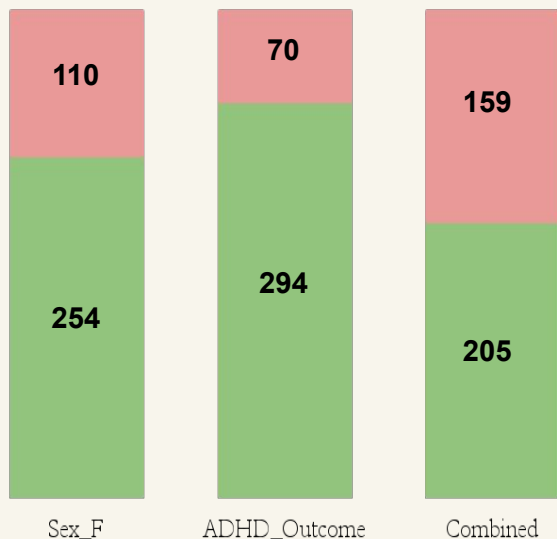
## 03

### CNN + ANN Hybrid Architecture

1.CNNs are effective at recognizing patterns in data that have a natural order or structure, helping to make sense of complex numerical information.
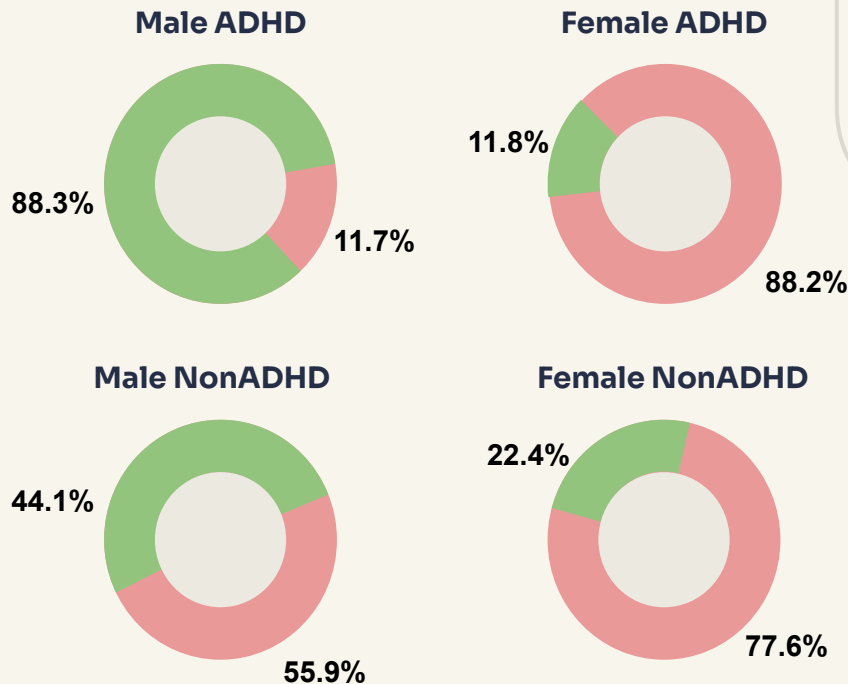
2.ANNs offer a more flexible way to understand grouped or labeled information, like age or gender, compared to simpler models.
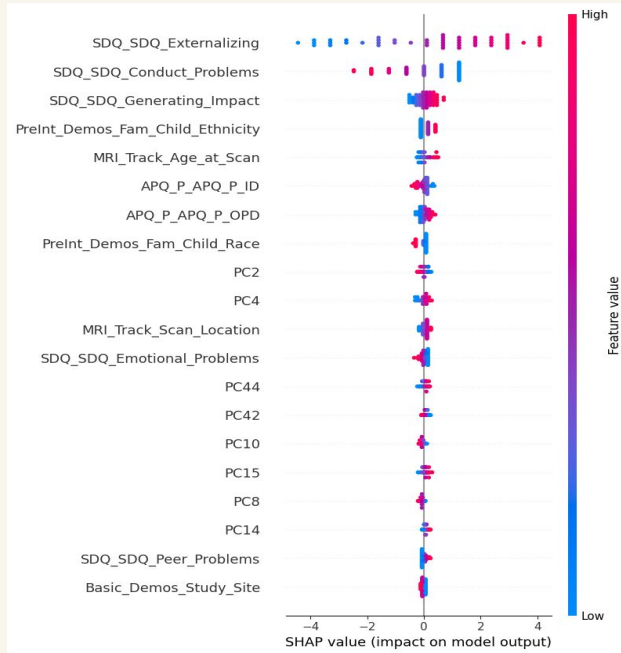
# Logistic Regression Prediction



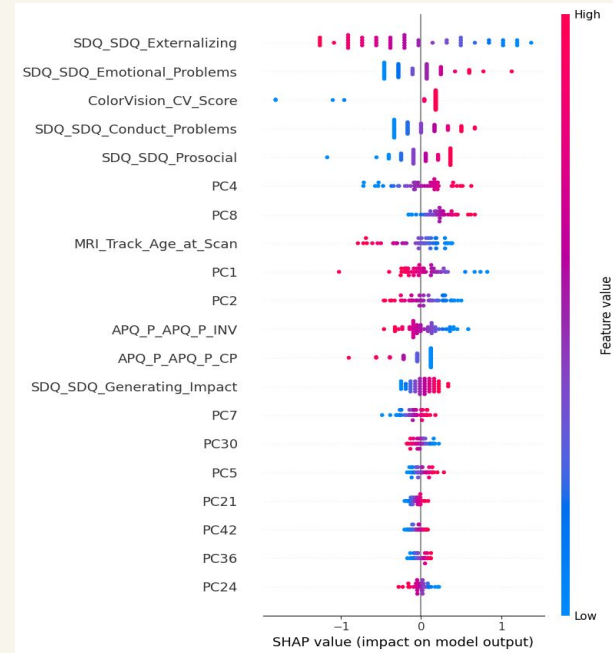Correct vs Incorrect Predictions Per Label

Sex_F: Incorrect 110, Correct 254
ADHD_Outcome: Incorrect 70, Correct 294
Combined: Incorrect 159, Correct 205

**Male ADHD** — 88.3% / 11.7%

**Female ADHD** — 11.8% / 88.2%

**Male NonADHD** — 44.1% / 55.9%

**Female NonADHD** — 22.4% / 77.6%

Correct vs Incorrect Predictions Per Category

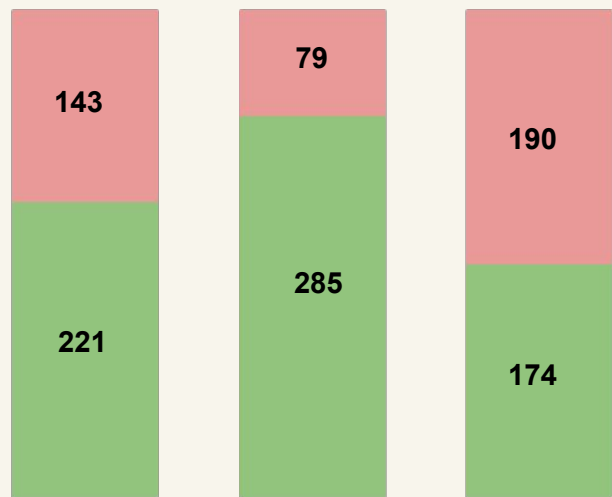■ Correct   ■ Incorrect

# Logistic Regression SHAP Analysis



SHAP Analysis for Sex F
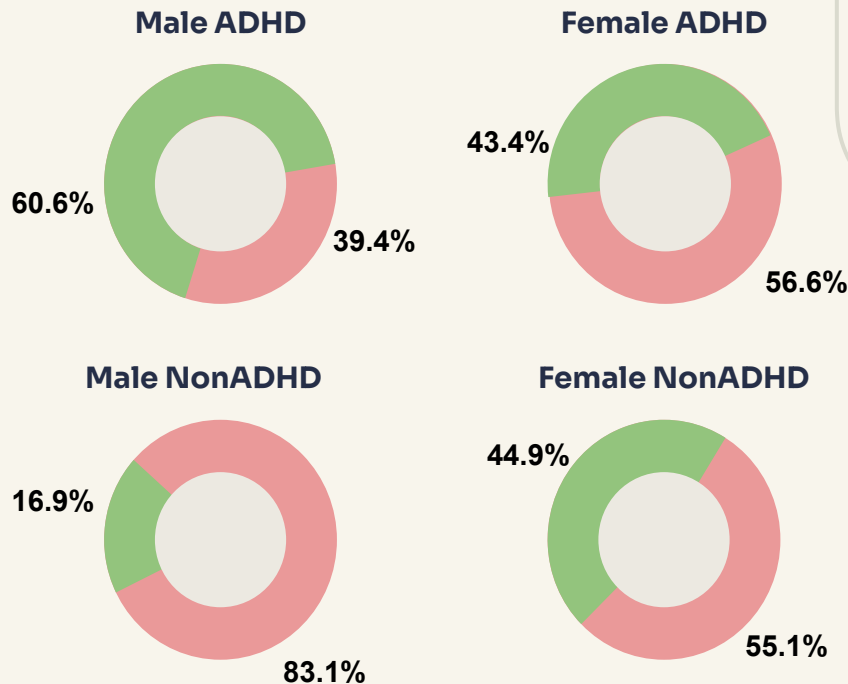


SHAP Analysis for ADHD Outcome

# CNN + MLP Hybrid Architecture Prediction



Correct vs Incorrect Predictions Per Label

**Male ADHD**
60.6% | 39.4%

**Female ADHD**
43.4% | 56.6%

**Male NonADHD**
16.9% | 83.1%

**Female NonADHD**
44.9% | 55.1%

Correct vs Incorrect Predictions Per Category

Bar chart values:
- Sex_F: 143 (Incorrect), 221 (Correct)
- ADHD_Outcome: 79 (Incorrect), 285 (Correct)
- Combined: 190 (Incorrect), 174 (Correct)

■ Correct ■ Incorrect

# CNN + MLP Hybrid Architecture SHAP Analysis



SHAP Analysis for Sex on Categorical Data

- Study site, parental education & occupation, ethnicity, and scan location impact predictions the most.

- Higher values in certain ethnicity and education features positively shift predictions, while others lower it, indicating varied demographic influence.
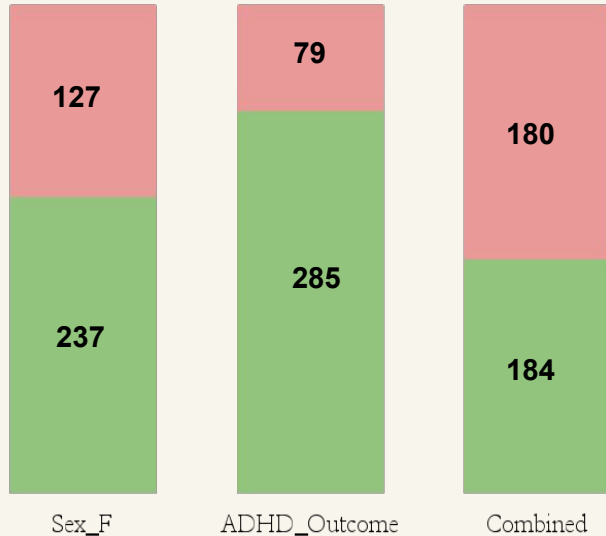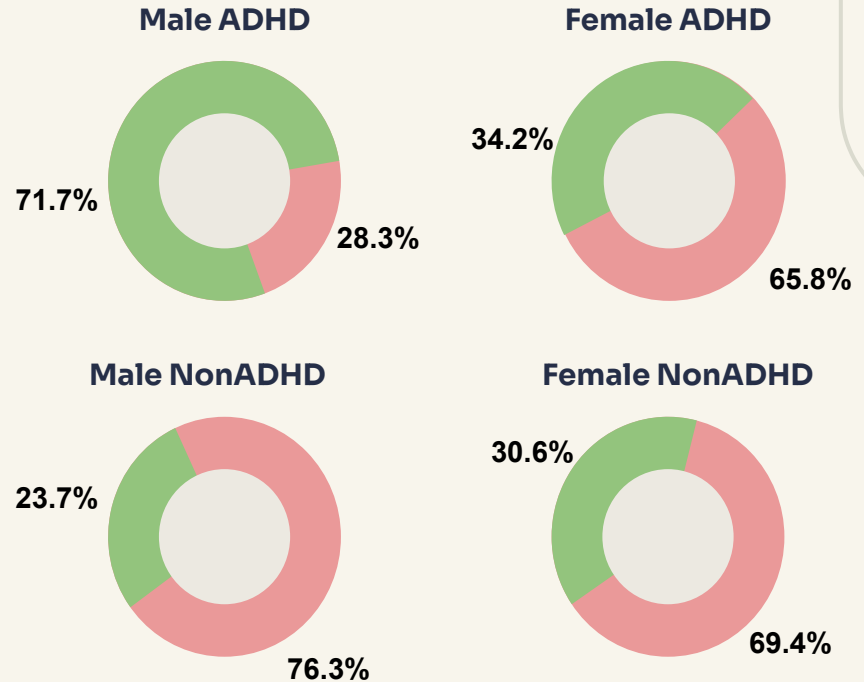
# CNN + MLP Hybrid Architecture SHAP Analysis



- Top influential features include ethnicity (PreInt_Demos_Fam_Child_Ethnicity_0), parental education (Barratt_Barratt_P1_Edu_18) and occupation (Barratt_Barratt_P2_Occ_45).

- MRI scan location contributes notably, suggesting site-specific effects in prediction.

- Demographic variables show strong influence, indicating the model is sensitive to socio-demographic context.

SHAP Analysis for ADHD Outcome on Categorical Data

# CNN + ANN Hybrid Architecture Prediction



Correct vs Incorrect Predictions Per Label

| | Correct | Incorrect |
|---|---|---|
| Sex_F | 237 | 127 |
| ADHD_Outcome | 285 | 79 |
| Combined | 184 | 180 |

**Male ADHD**

71.7% / 28.3%

**Female ADHD**

34.2% / 65.8%

**Male NonADHD**

23.7% / 76.3%

**Female NonADHD**

30.6% / 69.4%

Correct vs Incorrect Predictions Per Category

■ Correct   ■ Incorrect

12

# CNN + ANN Hybrid Architecture SHAP Analysis



SHAP Analysis for Sex on Categorical Data
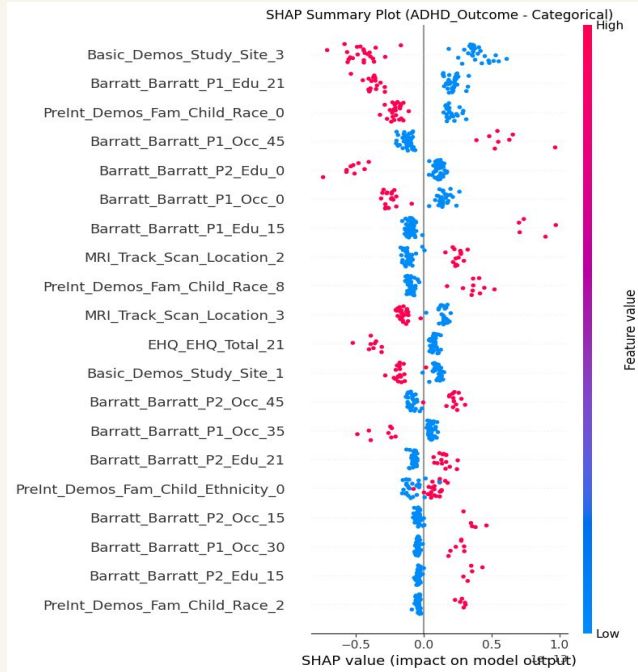
- Study site and parental education are the most impactful categorical features for predicting female sex.

- Ethnicity and race variables also play a notable role.

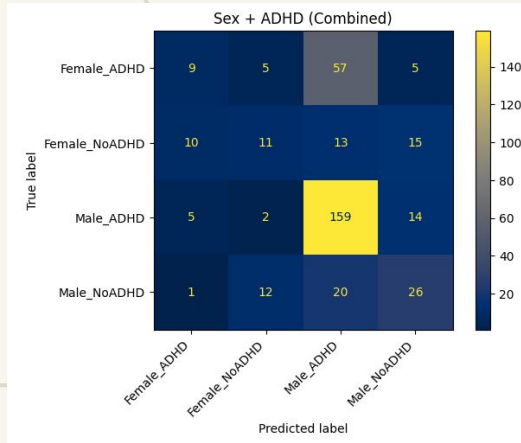- SHAP values closer to ±2 indicate stronger influence (either positive or negative) on the output.

# CNN + ANN Hybrid Architecture SHAP Analysis



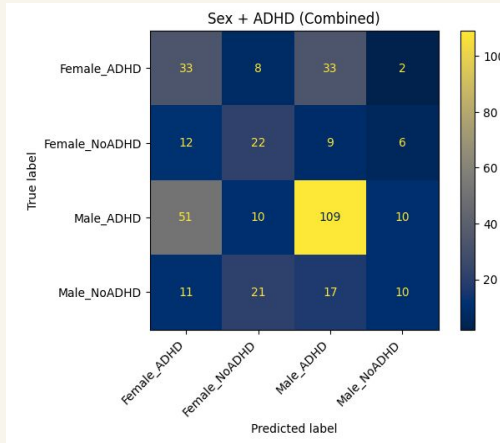SHAP Analysis for ADHD Outcome on Categorical Data

- Study Site 3 is the most influential factor in predicting ADHD.

- Parental education, especially having a graduate degree, strongly influences the model's predictions.

- The model heavily leans on socioeconomic, demographic, and site-related features to make ADHD predictions.
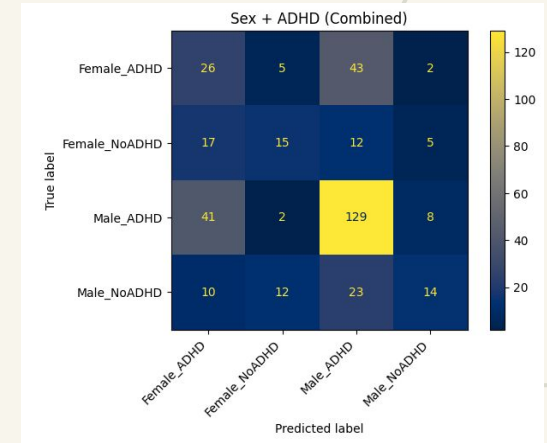
# True vs. Predicted Labels: Model–Wise Breakdown



Confusion Matrix of Logistic Regression

Confusion Matrix of CNN + MLP Hybrid Architecture

Confusion Matrix of CNN + ANN  Hybrid Architecture

# Conclusions

**01**

Is the model suitable for the company's needs and ready for deployment?

No, the model is not ready for deployment, as it fails to perform reliably across all classes.

**02**

Is the model showing any signs of bias, particularly gender-related?

Yes, the model shows gender related bias, frequently misclassifying female subjects.

**03**

Can insights from the data support future NHS projects beyond ADHD diagnosis?

Yes, the data insights can support future NHS projects, especially those involving behavioral patterns and gender disparities in diagnosis.

# Future Work

Bias Mitigation through Demographic-Specific Modeling:

- Large variation observed in parental education, occupation, and child's race.
- These factors may contribute to model bias or misclassification.
- Train separate models for specific demographic groups (e.g., based on education and occupation levels).
- Compare performance and bias metrics across groups.
- Investigate if demographic-specific training improves fairness and accuracy.

# Thank You!