

# Siamese Network-Based Vector Embeddings of MRI Scans for Monozygotic Twin Identification

**Matthew Kenely**

Supervisor: Dr Dylan Seychell

Co-Supervisor: Dr Claude Julien Bajada

September 2025

*Submitted in partial fulfilment of the requirements  
for the degree of Master of Science in Artificial Intelligence.*



**L-Università ta' Malta**  
Faculty of Information &  
Communication Technology



L-Universit`  
ta' Malta

## **University of Malta Library – Electronic Thesis & Dissertations (ETD) Repository**

The copyright of this thesis/dissertation belongs to the author. The author's rights in respect of this work are as defined by the Copyright Act (Chapter 415) of the Laws of Malta or as modified by any successive legislation.

Users may access this full-text thesis/dissertation and can make use of the information contained in accordance with the Copyright Act provided that the author must be properly acknowledged. Further distribution or reproduction in any format is prohibited without the prior permission of the copyright holder.

# Abstract

Monozygotic twins are identical twins that develop from a single fertilised egg that spontaneously splits, resulting in two individuals sharing 100% genetic material. Identifying monozygotic twins from brain MRI scans represents a frontier challenge in computational medical imaging with significant implications for understanding genetic influences on neuroanatomical structure through direct pattern recognition. While classical twin studies using ACE models decompose statistical variance to establish independent regional heritability estimates (60-80%), this study introduces a fundamentally different computational framework that learns directly from MRI data to rank neuroanatomical regions by their collective discriminative capacity for genetic similarity detection, complementing traditional statistical approaches through data-driven analysis.

A deep learning methodology employing Siamese networks with 3D CNN backbones is developed for automated twin identification using 138 genetically verified monozygotic twin pairs (276 subjects) from the Human Connectome Project S1200 dataset. Modified U-Net, ResNet, and DenseNet architectures generate 128-dimensional embeddings optimised via triplet loss with hard negative mining, forcing models to learn subtle genetic signatures by focusing on challenging discriminative examples that distinguish twins from their most similar morphological matches.

U-Net achieved superior computational performance with 92.0% F1-score ( $\sigma = 2.5\%$ ), 95.2% AUC-ROC, and 91.4% accuracy, while ResNet demonstrated competitive results (89.6% F1-score) and DenseNet showed greater variability (88.5% F1-score). Embedding analysis reveals clear bimodal separation between genetically related and unrelated individuals through learned morphological patterns.

Layer-Wise Relevance Propagation analysis provides the first data-driven ranking of neuroanatomical regions by discriminative importance for genetic relatedness detection. Statistical analysis reveals pronounced subcortical dominance with large effect size (Cohen's  $d = 2.80$ ,  $p = 3.89e-6$ ), with six subcortical structures occupying top positions, including the thalamus (0.955), brainstem (0.875), and hypothalamus (0.707). This computational hierarchy contrasts with traditional ACE studies reporting highest heritability in cortical areas (frontal 78-95%, temporal 77-89%), demonstrating that direct pattern recognition from MRI data identifies different neuroanatomical signatures than statistical variance decomposition. Notably, models utilise practically all brain regions (most importance scores  $> 0.2$ ), indicating distributed multivariate processing rather than selective regional dependence.

Ablation studies confirm data augmentation's critical role, with substantial

performance improvements across CNN architectures. Clinical integration through standard neuroimaging formats in Connectome Workbench demonstrates immediate practical utility, positioning this computational approach for adoption in research and clinical environments requiring direct analysis of genetic influences in brain structure.

The framework advances precision neuroimaging by providing automated, quantitative genetic similarity detection through direct pattern recognition, revealing spatial insights that complement traditional heritability studies while offering methodological advances applicable to diverse medical imaging classification tasks requiring regional discriminative analysis.

# Acknowledgements

I would like to thank Dr Dylan Seychell, my supervisor, for his continuous support and guidance in deep learning and computer vision. I would also like to thank my co-supervisor, Dr Claude Bajada, for his expertise in neuroimaging research and his guidance on the neuroanatomical aspects of this study.

I would like to thank my friends and family for their constant support throughout the course of this study and for their continuous interest in the findings.

Finally, I would like to express my gratitude to the Human Connectome Project and the WU-Minn HCP Consortium for providing access to the S1200 dataset. I thank the 1,206 participants who contributed their neuroimaging data, in particular the 138 monozygotic twin pairs whose MRI data made this research possible.

This study was funded by the Pathfinder Malta Digital Innovation Authority Digital Scholarship Award Agreement.

Data were provided [in part] by the Human Connectome Project, WU-Minn Consortium (Principal Investigators: David Van Essen and Kamil Ugurbil; 1U54MH091657) funded by the 16 NIH Institutes and Centers that support the NIH Blueprint for Neuroscience Research; and by the McDonnell Center for Systems Neuroscience at Washington University.

# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>Contents</b>	<b>vi</b>
<b>List of Figures</b>	<b>vii</b>
<b>List of Tables</b>	<b>xi</b>
<b>List of Abbreviations</b>	<b>xiii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Problem Definition . . . . .	1
1.2 Proposed Solution . . . . .	2
1.3 Aims and Objectives . . . . .	2
1.3.1 Objectives . . . . .	3
1.3.2 Research Questions . . . . .	3
1.4 Document Structure . . . . .	3
<b>2 Background and Literature Review</b>	<b>5</b>
2.1 Magnetic Resonance Imaging (MRI) . . . . .	5
2.1.1 T1 weighting . . . . .	6
2.2 Human Connectome Project 1200 Subjects Data Release . . . . .	6
2.2.1 Dataset Overview . . . . .	6
2.2.2 Glasser MMP Cortical Atlases . . . . .	7
2.2.3 22 parcellated cortical regions . . . . .	7
2.3 Deep Learning for Medical Imaging . . . . .	8
2.3.1 3D Convolutional Neural Networks . . . . .	8
2.3.2 U-Net . . . . .	8
2.3.3 ResNet . . . . .	9
2.3.4 DenseNet . . . . .	10
2.3.5 Attention Mechanisms in 3D CNNs . . . . .	10

2.4	Metric Learning and Similarity Networks . . . . .	11
2.4.1	Siamese Networks . . . . .	11
2.4.2	Triplet Loss and Hard Negative Mining . . . . .	12
2.5	Model Interpretability . . . . .	13
2.5.1	Layer-wise Relevance Propagation (LRP) . . . . .	13
2.5.2	T-distributed Stochastic Neighbor Embedding (t-SNE) . . . . .	14
2.6	Twin Studies in Medical Imaging . . . . .	14
2.6.1	Heritability Studies and the ACE Model . . . . .	14
2.6.2	From Statistical Decomposition to Computational Pattern Recognition . . . . .	15
<b>3</b>	<b>Methodology</b>	<b>16</b>
3.1	Data Definition . . . . .	16
3.1.1	Description . . . . .	16
3.1.2	Data Security . . . . .	16
3.1.3	Preprocessing . . . . .	18
3.1.4	Data Augmentation . . . . .	21
3.1.5	Data Preparation . . . . .	21
3.2	Siamese Architecture . . . . .	23
3.2.1	Siamese Network Design . . . . .	23
3.3	3D CNN Backbones . . . . .	25
3.3.1	U-Net . . . . .	25
3.3.2	ResNet . . . . .	26
3.3.3	DenseNet . . . . .	27
3.4	Training . . . . .	28
3.4.1	Loss Function . . . . .	28
3.4.2	Training Algorithm . . . . .	30
3.5	Quantitative Evaluation . . . . .	32
3.5.1	Model Selection and Evaluation Protocol . . . . .	32
3.5.2	Performance Metrics . . . . .	33
3.5.3	ROC and Precision-Recall Curve Analysis . . . . .	34
3.5.4	Threshold Optimisation . . . . .	35
3.6	Qualitative Evaluation . . . . .	36
3.6.1	Embeddings Interpretation . . . . .	36
3.6.2	Twin vs Non-twin Similarity Analysis . . . . .	37
3.6.3	Ensemble Modelling . . . . .	39
3.6.4	Medical Format Conversion . . . . .	40
<b>4</b>	<b>Evaluation</b>	<b>45</b>

4.1	Experimental Setup . . . . .	45
4.1.1	Hardware and Software . . . . .	45
4.2	Quantitative Results . . . . .	45
4.2.1	Loss Graphs . . . . .	45
4.2.2	Performance Metrics . . . . .	47
4.2.3	ROC and Precision-Recall Analysis . . . . .	49
4.2.4	Embedding Distance Distribution Analysis . . . . .	49
4.3	Qualitative Results . . . . .	50
4.3.1	Dimensionality Reduction . . . . .	50
4.3.2	Reference Attribution Patterns . . . . .	51
4.3.3	Embedding Space Uniqueness . . . . .	53
4.3.4	Discriminative Heatmaps . . . . .	54
4.3.5	Ensemble Integration . . . . .	55
4.3.6	MMP Glasser Brain Regions . . . . .	57
4.3.7	Medical Format Conversion . . . . .	58
4.4	Ablation Study . . . . .	63
4.4.1	Augmentation . . . . .	63
4.5	Discussion . . . . .	70
4.5.1	Deep Learning for Genetic Similarity Detection . . . . .	70
4.5.2	Computational vs. Statistical Approaches to Genetic Neuroimaging . . . . .	71
4.5.3	Data-Driven Regional Importance Hierarchy . . . . .	71
4.5.4	Methodological Paradigm and Clinical Integration . . . . .	72
<b>5</b>	<b>Conclusion</b>	<b>73</b>
5.1	Summary of Contributions . . . . .	73
5.2	Impact on Computational Medical Imaging and Neurogenetics . . . . .	73
5.3	Limitations and Methodological Considerations . . . . .	73
5.4	Future Directions in Computational Neurogenetics . . . . .	74
<b>A</b>	<b>Cortical and Subcortical Mappings</b>	<b>82</b>
<b>B</b>	<b>Full Brain Region Importance Listing</b>	<b>84</b>

# List of Figures

Figure 3.1	Left to right: axial, coronal, and sagittal views of a random HCP T1-weighted Magnetic Resonance Imaging (MRI) scan from the HCP 1200 Subjects Data Release preprocessed using the HCP minimal preprocessing pipelines, showing the $260 \times 311 \times 260$ voxel dimensions . . . . .	17
Figure 3.2	Histogram showing age distribution of the 138 monozygotic twin pairs (276 subjects) ranging from 22-36 years. . . . .	17
Figure 3.3	Intensity normalisation comparison showing histograms for two different subjects before (left) and after (right) grand-mean intensity normalisation. . . . .	19
Figure 3.4	Spatial downscaling comparison showing axial, coronal, and sagittal views of the same brain before (top row: $260 \times 311 \times 260$ ) and after (bottom row: $86 \times 103 \times 86$ ) trilinear interpolation downscaling. . . . .	20
Figure 3.5	Siamese network architecture overview. Two identical 3D CNN branches process brain MRI volumes to generate 128-dimensional embeddings. Cosine distance between embeddings determines twin versus non-twin classification through threshold comparison. . . . .	23
Figure 4.1	Training and validation loss curves showing triplet loss convergence and embedding distance separation across 2000 epochs. The twin/non-twin separation lines represent average embedding distances for twin pairs versus non-twin pairs across all pairs at each epoch. U-Net demonstrates most stable convergence with minimal overfitting and clear bimodal separation. . .	47
Figure 4.2	Confusion matrices for ResNet, U-Net, and DenseNet architectures on twin identification task across 10 evaluation runs. . . . .	48
Figure 4.3	Receiver Operating Characteristic (ROC) and Precision-Recall (PR) curve evaluation across model architectures. . . . .	49
Figure 4.4	Embedding distance distributions demonstrating clear bimodal separation between twin and non-twin pairs. U-Net achieves optimal separation with minimal overlap, while DenseNet shows substantial distribution overlap. . . . .	50

Figure 4.5	t-distributed Stochastic Neighbor Embedding (t-SNE) visualisation of embedding space across model architectures showing representative twin and non-twin pairs. Green dots represent twin pairs connected by dashed lines, red crosses indicate non-twin pairs, and blue circles show overlap points where subjects have both twin and non-twin relationships in the visualisation.	51
Figure 4.6	Reference LRP maps averaged across all subjects and embedding dimensions for each architecture as described in Equation (3.29) (left: axial, middle: coronal, right: sagittal views). All models converge on similar brain regions with consistent focus on subcortical structures and brainstem. Colour scale (0.0-1.0) represents normalised LRP attribution values, where higher values indicate regions contributing more to embedding representation.	52
Figure 4.7	Peak Signal-to-Noise Ratio analysis of embeddings with maximum deviation from reference Layer-wise Relevance Propagation (LRP) patterns. PSNR of each embedding denoted in white text. High similarity across architectures indicates models dedicate embedding dimensions to subtle structural variations within consistent regions.	53
Figure 4.8	Architecture-specific discriminative heatmaps showing differential embedding patterns between twin and non-twin pairs as described in Equation (3.32). Darker regions ( $< 0$ ) indicate larger embedding distances in non-twin pairs; lighter regions ( $> 0$ ) indicate larger embedding distances in twin pairs. U-Net and ResNet demonstrate similar patterns, while DenseNet shows weaker differentiation. Colour scale represents embedding distance differences, where negative values indicate regions with greater twin/non-twin discriminative capacity.	54
Figure 4.9	Ensemble heatmap combining weighted contributions from all three architectures as described in Equation (3.34), showing enhanced spatial coherence in subcortical regions. Top row: axial slices; middle row: coronal slices; bottom row: sagittal slices. Colour scale (0.0-1.0) represents discriminative importance for genetic relatedness detection, where higher values indicate regions that contribute more to distinguishing genetically related individuals.	55
Figure 4.10	Subject-specific atlas heatmap demonstrating clinical applicability for personalised twin similarity assessment, generated by applying population-level regional statistics to individual anatomical parcellation. Top row: axial slices; middle row: coronal slices; bottom row: sagittal slices. Colour scale (0.0-1.0) represents regional importance scores for genetic similarity, where higher values indicate anatomical regions with greater discriminative capacity for twin identification.	56

Figure 4.11 Gaussian-smoothed ( $\sigma = 2$ ) version of the subject-specific atlas heatmap in Figure 4.10, optimised for medical format conversion. Smoothing preserves spatial patterns while reducing noise for integration with standard neuroimaging pipelines. . . . .	56
Figure 4.12 Subject-specific T1-weighted anatomical volume displayed in Connectome Workbench, showing sagittal, coronal, and axial views used as the reference space for atlas heatmap alignment and volumetric format conversion.	59
Figure 4.13 Unsmoothed subject-specific atlas heatmap overlaid on T1-weighted anatomy in Connectome Workbench, displaying raw regional importance values with visible pixelation artefacts before Gaussian smoothing optimisation.	60
Figure 4.14 Gaussian-smoothed subject-specific atlas heatmap ( $\sigma = 2$ ) overlaid on anatomical volume in Connectome Workbench, demonstrating enhanced spatial coherence and improved visualisation quality suitable for clinical interpretation. . . . .	60
Figure 4.15 Subject-specific cortical surface reconstruction displayed in Connectome Workbench, showing bilateral hemisphere views of the native anatomical mesh used for surface-based heatmap projection. . . . .	61
Figure 4.16 Unsmoothed atlas heatmap projected onto cortical surfaces in Connectome Workbench using trilinear interpolation, showing raw importance values with high spatial frequency artefacts and discontinuous regional boundaries. . . . .	61
Figure 4.17 Gaussian-smoothed atlas heatmap projected onto cortical surfaces in Connectome Workbench, demonstrating improved spatial continuity with enhanced visualisation quality and clear discrimination of high-importance regions. . . . .	62
Figure 4.18 Validation F1-score curves across training configurations. Augmentation consistently accelerates convergence and achieves superior performance, with U-Net demonstrating gradual improvement over 5,000 epochs without augmentation, though remaining inefficient compared to augmented training. . . . .	63
Figure 4.19 Confusion matrices comparing augmented and non-augmented training across architectures. Augmentation consistently improves classification accuracy, with DenseNet showing the most substantial performance degradation without augmentation. . . . .	64
Figure 4.20 ROC and PR analysis revealing augmentation’s critical role in achieving robust discriminative performance, particularly for DenseNet architecture. .	66

Figure 4.21 Embedding distance distributions across augmentation strategies. Augmentation consistently improves twin vs non-twin separation, while extended training without augmentation leads to embedding collapse with reduced discriminative capacity. . . . . 67

Figure 4.22 Reference LRP maps across augmentation strategies. Non-augmented training produces diffuse activation patterns across the entire brain, while augmentation enables focused feature extraction in relevant neuroanatomical regions. . . . . 68

# List of Tables

Table 3.1	Siamese Network Architecture . . . . .	24
Table 3.2	Siamese Network Embedding Head . . . . .	24
Table 3.3	U-Net Architecture Comparison . . . . .	25
Table 3.4	U-Net Backbone Architecture . . . . .	26
Table 3.5	ResNet Architecture Comparison . . . . .	27
Table 3.6	ResNet Backbone Architecture . . . . .	27
Table 3.7	DenseNet Architecture Comparison . . . . .	28
Table 3.8	DenseNet Backbone Architecture . . . . .	28
Table 4.1	Siamese Network Training Parameters . . . . .	46
Table 4.2	Performance metrics comparison across architectures showing mean $\pm$ standard deviation across 10 evaluation runs with different randomly generated combinations of non-twin pairs to assess robustness to negative sample selection. . . . .	48
Table 4.3	Relative ranking of parcellated cortices and subcortical structures (bold) from HCP-MMP 1.0 atlas sorted by discriminative importance for genetic relatedness detection. Bootstrap confidence intervals (95% CI, n=10,000) demonstrate ranking stability, with subcortical structures showing significantly higher discriminative importance than cortical regions. Subcortical group mappings available in Table A.2. . . . .	57
Table 4.4	Ablation study on data augmentation strategies showing performance across architectures. A2k: Augmented 2k epochs, 2k: Non-augmented 2k epochs, 5k: Non-augmented 5k epochs. . . . .	65
Table 4.5	Regional importance rankings across augmentation strategies. Augmented training shows strong subcortical dominance with 6 subcortical structures in top 7 positions, while non-augmented conditions exhibit more distributed cortical-subcortical patterns. Full cortex names available in Table A.1. . . . .	69
Table A.1	Mapping of cortex abbreviations to full names . . . . .	82
Table A.2	Mapping of subcortical areas to functional groups. All mappings refer to both areas in the left and right hemispheres with the exception of the brain stem which constitutes a single area. . . . .	83

Table B.1 All individual brain regions from HCP-MMP 1.0 atlas sorted by importance, showing cortex, importance scores, activation, and volume for each region. Subcortical regions denoted in bold. Full cortex names available in Table A.1 . . . . . 84

# List of Abbreviations

AUC-PR Area Under the Precision-Recall Curve.

AUC-ROC Area Under the Receiver Operating Characteristic Curve.

BOLD Blood Oxygenation Level-Dependent.

CNN Convolutional Neural Network.

DTI Diffusion Tensor Imaging.

DWI Diffusion-Weighted Imaging.

fMRI Functional Magnetic Resonance Imaging.

FPR False Positive Rate.

Grad-CAM Gradient-weighted Class Activation Mapping.

HCP Human Connectome Project.

IoU Intersection over Union.

LRP Layer-wise Relevance Propagation.

MRI Magnetic Resonance Imaging.

PR Precision-Recall.

PSNR Peak Signal-to-Noise Ratio.

ROC Receiver Operating Characteristic.

t-SNE t-distributed Stochastic Neighbor Embedding.

TPR True Positive Rate.

# 1 Introduction

Convolutional Neural Networks (CNNs) have revolutionised medical image analysis, demonstrating exceptional capabilities in extracting complex patterns from MRI scans. Their architectural design, particularly in 3D implementations, makes them ideally suited for capturing spatial relationships in volumetric neuroimaging data. Substantial evidence supports their efficacy in brain tumour detection, classification, neurological disorder diagnosis, and structural analysis [1–5].

Understanding genetic influences on brain structure remains a fundamental challenge in neuroimaging. Classical twin studies employ ACE models (statistical frameworks that partition observed trait variance into additive genetic effects, shared environmental influences, and unique environmental factors). These approaches have established regional heritability estimates of 60-80% across different brain areas, quantifying the proportion of structural variance attributable to genetic factors. High heritability has been documented in both cortical areas (the brain's outer surface responsible for higher-order processing) and subcortical structures (deeper regions including the thalamus, basal ganglia, and hippocampus) [6, 7]. However, these traditional methods analyse regions independently rather than identifying which combinations of neuroanatomical features collectively provide the strongest signatures of hereditary relationships. Computational identification of monozygotic twins from MRI data represents largely uncharted territory, with significant implications for understanding how genetic similarity manifests in brain structure.

The Human Connectome Project (HCP)'s comprehensive dataset containing MRI volumes from 1,206 healthy young adults, including 149 genetically-verified monozygotic twin pairs, provides an unprecedented opportunity to develop computational approaches that rank neuroanatomical areas by their multivariate capacity for genetic similarity detection. Current computer vision research on twin identification has focused exclusively on facial features [8–10], overlooking the rich three-dimensional data available through neuroimaging.

## 1.1 Problem Definition

Current computational approaches to hereditary relationships in neuroimaging face significant limitations. Classical twin studies quantify regional heritability independently, providing percentage estimates for individual brain areas but failing to identify which features work jointly to enable genetic relatedness detection. No computational frameworks have been developed to rank brain regions by their joint discriminative importance for distinguishing genetically related individuals from their

most similar-appearing counterparts.

Existing computer vision research on twin identification remains confined to facial analysis, overlooking the rich volumetric neuroanatomical information available in brain imaging. Furthermore, no studies have systematically investigated which cortical and subcortical areas contribute most significantly to automated identification, nor have interpretable deep learning approaches been developed that provide relative rankings of discriminative importance across brain structures. This represents a critical gap in understanding how genetic similarity manifests as detectable patterns in neural architecture.

## 1.2 Proposed Solution

This research investigates Siamese network architectures with 3D CNN backbones for monozygotic twin identification using the HCP S1200 dataset. The computational framework employs modified U-Net, ResNet, and DenseNet architectures to generate 128-dimensional L2-normalised embeddings optimised via triplet loss with hard negative mining.

The approach incorporates hard negative selection, where difficult examples are chosen from subjects belonging to different twin pairs, forcing models to learn subtle genetic signatures that distinguish twins from their most phenotypically similar non-related individuals. Classification employs cosine distance between normalised embeddings, with optimal thresholds determined through F1-score maximisation.

LRP analysis provides voxel-level attributions through backward propagation of relevance scores, generating interpretability maps that highlight brain regions contributing most to similarity decisions. Performance-weighted ensemble modelling combines patterns across architectures, with resulting importance maps registered to the HCP-MMP 1.0 parcellation atlas to identify specific areas associated with genetic relatedness detection.

## 1.3 Aims and Objectives

The primary aim of this research is to develop and validate deep learning approaches for identifying monozygotic twins from brain MRI data, while providing the first computational ranking of neuroanatomical regions by their collective capacity for genetic similarity detection, complementing traditional ACE models' independent regional heritability estimates.

### 1.3.1 Objectives

1. Develop and validate Siamese network architectures for twin identification from 3D MRI volumes using triplet loss optimisation with hard negative mining
2. Achieve robust classification performance with F1-scores exceeding 80% across multiple neural backbone architectures
3. Apply Layer-Wise Relevance Propagation to rank brain structures by their distinguishing power for hereditary relationships through statistical testing and map findings to established neuroanatomical atlases
4. Demonstrate clinical applicability through integration with standard neuroimaging workflows and medical format compatibility

### 1.3.2 Research Questions

- Which 3D CNN architectures are most effective for learning discriminative features from neuroimaging data for automated twin identification?
- Which neuroanatomical structures demonstrate highest distinguishing capacity when collectively differentiating monozygotic twins from phenotypically similar non-related individuals?
- How do computational rankings of regional discriminative power relate to established heritability estimates from classical twin studies?
- Can this framework be successfully integrated into standard clinical neuroimaging workflows for practical applications?

## 1.4 Document Structure

The rest of this study is structured as follows:

- Chapter 2 provides an introduction to magnetic resonance imaging and the Human Connectome Project S1200 dataset, reviews key deep learning architectures and metric learning approaches for medical imaging, presents model interpretability methods, and concludes with a review of twin studies contrasting ACE models with computational pattern recognition.
- Chapter 3 describes the computational framework for detecting genetic similarity in brain MRI data, detailing dataset preparation, the core Siamese architecture with three 3D CNN backbones, and training and evaluation strategies.

- Chapter 4 presents and discusses the experimental validation of the framework, covering quantitative and qualitative evaluation, ablation studies, regional importance analysis using the HCP-MMP atlas, and medical format conversion demonstrating clinical applicability.
- Chapter 5 provides a summary of the work carried out in this study, its limitations, and suggested areas for future work.

## 2 Background and Literature Review

This chapter outlines the foundations of the research. It introduces magnetic resonance imaging, focusing on T1-weighted structural scans and the Human Connectome Project S1200 dataset. Key deep learning architectures for medical imaging are reviewed, including 3D CNNs, U-Net, ResNet, and DenseNet, which form the basis of the Siamese framework. Metric learning approaches are discussed, with emphasis on Siamese networks and triplet loss with hard negative mining. Model interpretability methods, such as Layer-wise Relevance Propagation and t-SNE visualisation, are presented. The chapter concludes with a review of twin studies in medical imaging, contrasting ACE models with computational pattern recognition, setting the foundation for the proposed methodology.

### 2.1 Magnetic Resonance Imaging (MRI)

MRI has emerged as a non-invasive diagnostic tool in neurological assessment since the production of the first clinical images in 1980 [11]. The technology has become widely accessible [12] and offers significant advantages over alternative imaging methods by avoiding exposure to ionising radiation [13]. MRI functions on principles of nuclear spin and the response of nuclei to external magnetic fields [11], producing high-resolution images of soft tissues that are particularly valuable for brain structure analysis.

Various MRI techniques are employed in neurological investigations, including T1-weighted, T2-weighted, proton density (PD), T2\*-weighted gradient echo, and Diffusion-Weighted Imaging (DWI) [14]. DWI, which quantifies water molecule movement, has proven especially useful in the examination of neurological conditions including multiple sclerosis and brain tumours [11]. Diffusion Tensor Imaging (DTI) extends these capabilities by mapping white matter tracts and neural connectivity patterns. Functional techniques such as perfusion MRI (PW-MRI) and Blood Oxygenation Level-Dependent (BOLD) Functional Magnetic Resonance Imaging (fMRI) [14] further expand the analytical scope by capturing dynamic brain activity.

The relevance of MRI to monozygotic twin identification lies in its ability to generate detailed structural representations of brain morphology. Previous research has demonstrated similarities in certain brain regions between monozygotic twins [6, 7], although the precise regions that remain consistent despite environmental influences remain incompletely characterised. MRI's capacity to produce diverse sets of high-resolution images configuring different aspects of brain structure provides an ideal foundation for computer vision approaches to twin identification. The detailed visualisation capabilities of MRI position it as an optimal imaging modality for

extracting the subtle morphological features necessary for distinguishing between twin and non-twin pairs.

### 2.1.1 T1 weighting

T1 weighted images represent a fundamental MRI acquisition technique that exploits tissue relaxation properties to generate specific contrast patterns [11]. T1 weighting, characterised by short repetition (TR) and echo times (TE), measures the longitudinal relaxation of protons returning to equilibrium after RF excitation [11]. This results in images where fat appears bright, white matter appears brighter than grey matter, and cerebrospinal fluid (CSF) appears dark. T1-weighted sequences excel at depicting anatomical detail, particularly the grey-white matter interface, making them invaluable for morphological assessment and volumetric analysis of brain structures [12]. In the context of monozygotic twin studies, T1 weighting offers superior anatomical detail for structural comparison, providing the high-resolution morphological information necessary for a Siamese network to learn to identify twin pairs.

## 2.2 Human Connectome Project 1200 Subjects Data Release

### 2.2.1 Dataset Overview

The HCP launched in 2009 as a five-year project funded by sixteen NIH components, marking the first Grand Challenge of the NIH's Blueprint for Neuroscience Research. The project sought to map comprehensive neural networks (connectomes) characterising anatomical and functional brain connectivity patterns in healthy populations, producing research datasets to advance understanding of neurological and psychiatric conditions such as autism, Alzheimer's disease, and schizophrenia. The S1200 dataset contains behavioural and 3 Tesla (3T) MR imaging data from 1,206 healthy young adults (aged 22-35) collected from August 2012 to October 2015. This includes 3T structural scans for 1,113 subjects, with 889 subjects having complete data across all four 3T MRI modalities: structural images (T1w and T2w), resting-state fMRI, task fMRI, and high angular resolution diffusion imaging<sup>1</sup>.

The recruitment strategy targeted families with twins, producing genetic data on 1,142 participants, including 149 pairs of genetically-confirmed monozygotic twins

---

<sup>1</sup><https://www.humanconnectome.org/study/hcp-young-adult/document/announcing-1200-subject-data-release>

(298 participants) and 94 pairs of dizygotic twins (188 participants)<sup>2</sup>. Genetic verification identified 36 twin pairs originally self-reported as dizygotic but confirmed as genetically monozygotic through blood and saliva genotyping<sup>3</sup>.

This twin sample enables quantification of genetic and environmental variation in brain structural patterns. Monozygotic twins exhibit maximal similarity due to nearly identical genetics, providing a resource for understanding heritable brain characteristics required for computational twin identification methods.

## 2.2.2 Glasser MMP Cortical Atlases

The HCP-Multimodal Parcellation (HCP-MMP) atlas addresses a long-standing neuroscience goal: creating an accurate areal map of the cerebral cortex [15]. Researchers applied an objective semi-automated neuroanatomical methodology to multimodal magnetic resonance images from the HCP, delineating 180 cortical areas per hemisphere based on distinct transitions in architecture, function, connectivity, and topography within a group-averaged template derived from 210 healthy young adults<sup>4</sup>.

The atlas utilised an observer-independent methodology tailored for non-invasively acquired multi-modal MRI data, where computational algorithms defined areal boundaries while neuroanatomists provided multi-modal data interpretation. This parcellation discovered 97 previously unidentified areas alongside 83 areas consistent with post-mortem microscopic studies [15], integrating diverse neurobiological properties (architecture, function, connectivity, topography) to provide complementary anatomical information.

## 2.2.3 22 parcellated cortical regions

The HCP-MMP 1.0 atlas organises the 180 cortical areas per hemisphere into 22 distinct regions based on anatomy, function, and connectivity. Each hemisphere's 180 areas are numbered 1-180 on the left and 201-380 on the right, with each area assigned to one of 22 larger cortical partitions numbered 1-22 on the left and 101-122 on the right<sup>5</sup>.

In addition to the 22 cortical regions, six subcortical structures are included in the current analysis: thalamus, brainstem, hypothalamus, basal ganglia, cerebellum, and limbic structures. These subcortical regions are incorporated to provide a comprehensive representation of brain architecture beyond cortical areas alone.

<sup>2</sup>[https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs001364.v1.p1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001364.v1.p1)

<sup>3</sup><https://www.humanconnectome.org/study/hcp-young-adult/article/announcing-1200-subject-data-release>

<sup>4</sup><https://balsa.wustl.edu/study/show/RVVG>

<sup>5</sup><https://neuroimaging-core-docs.readthedocs.io/en/latest/pages/atlasses.html>

This combined framework provides an optimal balance between anatomical specificity and computational tractability for neuroimaging analyses using machine learning approaches. For twin identification research, this organisation enables systematic analysis of genetic and environmental influences across functionally distinct brain regions.

## 2.3 Deep Learning for Medical Imaging

### 2.3.1 3D Convolutional Neural Networks

CNNs have revolutionised deep learning, achieving remarkable results across domains including computer vision and natural language processing [16]. Inspired by biological visual perception, CNNs offer advantages over fully connected architectures through their use of local connections, weight sharing, and hierarchical feature extraction [16]. These networks function as integrated systems combining feature extraction and classification components [17], making them well-suited for medical imaging analysis.

Three-dimensional CNNs represent a specialised adaptation for analysing volumetric medical data such as MRI scans [18]. Unlike their 2D counterparts, which are limited to planar information, 3D CNNs can extract and process spatial relationships across all dimensions simultaneously [19]. This capability proves particularly valuable for neuroimaging applications requiring comprehensive volumetric analysis. Given that brain MRIs are inherently three-dimensional, employing convolutions along all axes allows models to fully leverage this structural information when making predictions.

3D CNNs have demonstrated effectiveness across various neurological applications, including brain tumour grading [17], Alzheimer's disease diagnosis [20], and haemorrhage detection [4]. Mzoughi *et al.* demonstrated that 3D CNN architectures can effectively merge both local and global contextual information with reduced weights when applied to whole volumetric T1-weighted MRI sequences for glioma classification [20]. Furthermore, studies have shown that 3D convolution approaches outperform 2D approaches (including slice-by-slice processing and multi-channel 2D methods) in brain MRI analysis tasks [21].

### 2.3.2 U-Net

The U-Net architecture, proposed by Ronneberger *et al.* in 2015 for biomedical image segmentation [22], has become a standard architecture for medical image analysis tasks. The original work demonstrated that the network can be trained end-to-end from limited training data while achieving superior performance compared to existing convolutional approaches on biomedical segmentation tasks. The architecture features

an encoder that systematically reduces spatial dimensions through convolutional and pooling operations to extract multi-scale features and contextual information, paired with a decoder that reconstructs the original spatial resolution using transposed convolutions in a symmetric configuration [22].

U-Net's encoder-decoder architecture with skip connections has proven particularly effective for brain MRI analysis, with recent implementations achieving remarkable performance metrics in brain tumour segmentation tasks. Lightweight variants have demonstrated mean Intersection over Union (IoU) scores of 89% while maintaining real-time processing capabilities [23]. The modular nature of U-Net allows for extensive customisation and improvement, leading to numerous successful variants including 3D U-Net for volumetric data [24], Attention U-Net [25] and U-Net++ [26].

For twin identification tasks, U-Net's dense feature propagation characteristics make it suitable for adaptation with embedding heads. Studies have shown that U-Net-based architectures can be effectively modified to extract discriminative features for similarity learning tasks [27]. The architecture's ability to preserve spatial information through skip connections is particularly valuable when learning embeddings that capture subtle neuroanatomical differences between individuals. Multi-scale context fusion approaches have achieved 90.56% Dice coefficient on medical image datasets [28], while transformer-enhanced variants show average Dice improvements of 1.06% over competitive baselines in multi-organ segmentation tasks [29].

### 2.3.3 ResNet

He *et al.* introduced ResNet in 2016 [30] to solve the degradation problem in deep networks using residual connections. Rather than learning complete transformations, layers learn additive residuals that modify their inputs. This residual learning is implemented through skip connections that bypass layers, allowing gradients to flow directly through the network during training and providing effective feature extraction for medical imaging applications requiring discriminative representations.

The architecture has demonstrated exceptional performance across diverse healthcare applications. ResNet-50 achieved  $95.23\% \pm 0.6\%$  classification accuracy for multi-class brain disease detection [31], while its ability to train very deep networks (up to 152 layers) while maintaining gradient stability proves advantageous for complex tasks requiring fine-grained feature discrimination [32]. The model serves effectively as a feature extractor when combined with embedding heads for similarity learning applications.

In neuroimaging specifically, ResNet-based models have achieved state-of-the-art performance in Alzheimer's disease classification using structural MRI

data [33]. Advanced implementations incorporating explainable AI techniques demonstrated 97.35% accuracy in brain tumour detection with Gradient-weighted Class Activation Mapping (Grad-CAM) interpretability [34]. Multi-modal brain analysis applications have shown superior performance in identifying glioblastoma versus solitary brain metastases [35] and effective evaluation of survival in high-grade glioma patients [36].

### 2.3.4 DenseNet

DenseNet, proposed by Huang *et al.* in 2017 [37], implements dense connectivity patterns where every layer concatenates feature maps from all previous layers as its input. This connectivity scheme creates  $L(L+1)/2$  direct pathways in an  $L$ -layer network compared to  $L$  sequential connections in traditional architectures. The dense connectivity pattern addresses vanishing gradient issues, enhances feature propagation and reuse across the network, while paradoxically requiring fewer parameters than conventional designs.

Medical imaging applications have benefited from this efficient parameter usage and enhanced feature extraction capabilities [38]. For brain MRI analysis, the model demonstrates high accuracy in capturing intricate structures valuable for clinical and research applications. Enhanced variants incorporating Squeeze-and-Excitation modules and dilated convolutions have outperformed established models including ResNet-101, VGG-19, and transformer-based architectures in brain tumour classification [39].

The feature reuse characteristics provide advantages for embedding-based tasks and discriminative representation learning. Applications span from anatomical brain segmentation achieving 0.673 Dice coefficient on IBSR datasets [38] to pathological brain detection with 97.1% accuracy in senile dementia diagnosis [40]. Multi-level feature aggregation from all preceding layers enables rich representations suitable for similarity learning, with liver lesion classification showing superior performance when combined with data augmentation techniques [41].

### 2.3.5 Attention Mechanisms in 3D CNNs

Attention mechanisms enhance the discriminative capabilities of 3D CNNs in medical imaging applications [42]. Spatial attention networks have been applied to 3D brain MRI analysis by incorporating attention blocks after convolutional layers to emphasise important regions in feature maps. These mechanisms allow models to focus on discriminative anatomical regions while suppressing irrelevant features, improving performance and interpretability [43].

Recent advances have demonstrated the effectiveness of attention mechanisms in brain tumour segmentation using multi-modal MRI data, where attention helps models focus on relevant parts of input images. Global attention mechanisms and efficient channel attention have been integrated with CNN architectures to improve accuracy in brain tumour classification tasks, addressing the challenge of high variability in tumour appearances and subtle early-stage manifestations [44].

These attention-enhanced architectures are particularly relevant for twin identification tasks, as they can learn to focus on brain regions that are most discriminative for genetic similarity while maintaining computational efficiency across large volumetric datasets.

## 2.4 Metric Learning and Similarity Networks

### 2.4.1 Siamese Networks

Siamese networks are specialised neural network architectures designed specifically for comparative analysis between input pairs, employing two or more identical network branches with shared weights [45]. This architectural approach is particularly advantageous when working with limited training data [46], as is often the case with specialised medical imaging datasets such as those containing monozygotic twin pairs.

Siamese networks have shown effectiveness across various medical imaging applications. Li *et al.* [47] applied convolutional Siamese networks for disease severity assessment and longitudinal progression tracking across patient visits, achieving correlations of  $\rho = 0.87$  and  $\rho = 0.89$  for retinopathy of prematurity and osteoarthritis respectively. Livieris *et al.* [48] showed that Siamese network variants outperform traditional models like Support Vector Machines or Discriminant Analysis when classifying high-dimensional radiomic features from T2 MRI images.

Xu *et al.* [49] developed Siamese networks incorporating node convolution operations (graph-based computations that aggregate information from neighbouring brain regions within connectivity networks) to process brain connectivity graphs derived from resting-state fMRI data, enabling personalised predictions based on individual-specific connectivity patterns rather than group-averaged templates. Complementary research has combined Siamese networks with gradient-based attention mechanisms such as Grad-CAM, generating interpretable similarity assessments that highlight which image regions drive the similarity decisions, thereby providing transparency as to the network's computational process for determining image similarity [48].

The fundamental operational principle of Siamese networks involves processing

two separate inputs through identical network branches to generate comparable feature representations. Similarity metrics like cosine similarity or Euclidean distance then measure correspondence between these representations. This approach is well-suited for twin identification, where the goal is to differentiate twin and non-twin pairs based on subtle structural brain differences.

### 2.4.2 Triplet Loss and Hard Negative Mining

Triplet loss has emerged as a powerful optimisation objective for learning discriminative embeddings in medical imaging applications. Recent work in brain tumour classification has applied triplet contrastive learning combined with unsupervised pre-training and data augmentation to address the lack of data problem in brain tumour imaging analysis [50]. This approach directly learns deep embeddings for brain tumour types that can be used for downstream classification tasks.

Hard negative mining is a training strategy that selectively chooses the most challenging negative examples during model optimisation, specifically those negative samples that are closest to the anchor sample in the embedding space while still belonging to different classes. This technique focuses computational resources on the most informative samples that contribute most significantly to gradient updates and model improvement. Hard negative mining strategies have proven effective in medical imaging applications, with recent work demonstrating their utility in lymphatic invasion detection for gastric cancer, achieving Area Under the Receiver Operating Characteristic (AUC-ROC) of 0.9738 and Area Under the Precision-Recall Curve (AUC-PR) of 0.9501 while reducing false-positive predictions [51]. The study showed that hard negative mining improves deep learning model performance by automatically identifying the most challenging negative examples during training, reducing the need for manual annotation of difficult cases and decreasing the computational expense of processing large numbers of easy negative samples that provide minimal learning value.

Advanced hard negative mining techniques like Bag of Negatives (BoN) have been developed to provide computationally efficient selection of relevant training samples, showing superior accuracy and training time when compared to state-of-the-art methods [52]. Recent research has also addressed gradient vanishing issues in triplet loss optimisation, proposing selective hard negative mining approaches that decide whether to mine the hardest negative samples according to gradient conditions [53].

Combining triplet loss with hard negative mining enables models to learn more discriminative feature representations by automatically identifying and prioritising the most challenging negative examples during training, forcing the network to focus on subtle differences that distinguish similar but unrelated samples rather than learning

from easily separable examples that provide limited discriminative value [54]. This approach is particularly valuable for twin identification tasks, where subtle morphological differences must be captured while maintaining discriminative power across the embedding space.

## 2.5 Model Interpretability

### 2.5.1 Layer-wise Relevance Propagation (LRP)

LRP explains deep neural network decisions by decomposing predictions through backward propagation of relevance scores. Introduced by Samek *et al.* in 2015, the method has become a standard technique in explainable AI for medical imaging applications [55]. LRP uses specialised propagation rules that distribute relevance backward through network layers [56].

The framework operates on a conservation law ensuring that relevance scores are preserved as they flow backward through the network, maintaining balance between the total relevance distributed from each layer and the cumulative relevance received by the preceding layer. Following a standard forward pass that generates predictions, the output relevance is propagated backward layer-by-layer to produce heatmaps showing voxel-level importance for the final prediction [55]. This approach can be formalised as a 'deep Taylor decomposition' using layer-specific propagation rules, with recent extensions developed for convolutional networks containing local renormalisation layers [56, 57].

Unlike other visualisation methods such as Grad-CAM, which provides broad class-discriminative localisation maps with relatively coarse visualisations, LRP delivers exact voxel-wise contributions and detailed relevance attributions [58, 59]. This superior granularity makes it particularly valuable for neurogenetic applications requiring fine-grained analysis of genetic influences on brain structure.

In twin identification studies, LRP enables identification of brain regions contributing most significantly to monozygotic twin classification, potentially revealing novel biomarkers of genetic influence. The generated relevance maps can be systematically mapped to established cortical parcellations such as the Glasser MMP atlas described in Subsection 2.2.3, enabling precise anatomical localisation of genetically-influenced regions. By analysing diversity in relevance patterns across twin pairs, the technique may help distinguish heritable neuroanatomical features from those shaped by environmental factors, advancing understanding of genetic contributions to brain structural characteristics.

## 2.5.2 T-distributed Stochastic Neighbor Embedding (t-SNE)

t-SNE is a dimensionality reduction technique that maps high-dimensional data to two or three dimensions while maintaining local neighbourhood relationships [60]. The algorithm transforms pairwise distances into conditional probabilities using Gaussian distributions in the original space and Student t-distributions in the reduced embedding space. Gradient descent optimisation minimises the Kullback-Leibler divergence between these probability distributions, preserving local clustering patterns while solving the crowding issues inherent in conventional dimensionality reduction approaches.

In medical imaging contexts, t-SNE enables visualisation of learned embeddings from neural networks, allowing researchers to assess whether models have captured meaningful similarities between input samples and identify clustering patterns corresponding to biological differences. However, practitioners should note that visual clusters can be strongly influenced by parameterisation, and clusters may appear even in data with no clear underlying structure.

## 2.6 Twin Studies in Medical Imaging

Twin studies are a fundamental method for separating genetic and environmental influences on human traits. MRI-based twin studies have revealed how genetics shape brain structure, establishing the foundation for computational approaches that directly identify genetic similarity from neuroimaging data [6, 7].

### 2.6.1 Heritability Studies and the ACE Model

The traditional twin study design leverages the genetic differences between monozygotic (identical) and dizygotic (fraternal) twins to partition phenotypic variance into distinct sources. The ACE model is a statistical framework used in behavioural genetics that decomposes observed variance in traits into three components: **A** represents additive genetic effects (the cumulative influence of multiple genes), **C** represents shared or common environmental influences (environmental factors experienced by both twins that make them more similar), and **E** represents unique or individual-specific environmental factors (experiences unique to each twin, including measurement error). Since monozygotic twins share identical DNA while dizygotic twins share approximately 50% of segregating genes [61], comparing their phenotypic similarities allows statistical decomposition of observed brain structural variance into these three components through the ACE modelling framework [62].

These traditional approaches have revealed heterogeneous patterns of genetic influence across brain regions through statistical variance decomposition. Subcortical structures demonstrate variable heritability, with the hippocampus showing 40-73% heritability, basal ganglia 64%, thalamus 42%, lateral ventricles 17%, and cerebellum 24% [63–65]. Cortical regions exhibit generally high heritability, particularly frontal (78-95%) and temporal lobes (77-89%), while parietal regions show 55-89% heritability [64, 66]. However, these studies analyse regions independently rather than identifying which combinations of brain areas collectively provide the strongest signatures for genetic relatedness detection.

### **2.6.2 From Statistical Decomposition to Computational Pattern Recognition**

Traditional twin neuroimaging studies have demonstrated high heritability for global brain volumes including total brain (90%), grey matter (91%) and white matter (84%) [61]. However, trained radiologists can visually distinguish between genetically identical twin pairs [6], revealing that subtle morphological differences exist despite shared genetic architecture.

This observation presents a fundamental opportunity for computational approaches that complement traditional statistical methods. While ACE models decompose variance to measure independent regional heritability, deep learning can directly learn from neuroimaging data to identify which brain regions collectively provide the strongest morphological signatures for distinguishing genetic relatedness. This represents a paradigm shift from statistical variance analysis to direct computational pattern recognition.

The computational framework ranks regions based on their multivariate discriminative capacity rather than statistical variance decomposition. Neural networks can systematically map which morphological patterns remain consistent between twins and which diverge, providing data-driven insights into neuroanatomical genetic influences that emerge directly from structural imaging analysis. The proposed Siamese network approach leverages this foundation to develop automated methods for genetic similarity detection, offering novel perspectives on how genetic relatedness manifests in brain structure through direct pattern analysis rather than traditional statistical modelling assumptions.

## 3 Methodology

This chapter describes the computational framework for detecting genetic similarity in brain MRI data. It details the Human Connectome Project S1200 dataset, including preprocessing, augmentation, and preparation for Siamese training. The core Siamese architecture is presented with three 3D CNN backbones (U-Net, ResNet, and DenseNet) adapted for triplet loss optimisation. Training strategies, including hard negative mining and optimisation algorithms, are outlined. Evaluation methods are also described, covering quantitative metrics, qualitative interpretation, Layer-wise Relevance Propagation, ensemble modelling, and medical format conversion for clinical use. All model weights and code are made publicly available<sup>1</sup>.

### 3.1 Data Definition

#### 3.1.1 Description

This study extracted 3T MRI scan data for 138 genetically-verified monozygotic twin pairs (276 subjects total) from the Human Connectome Project S1200 dataset [67]. Each MRI volume consists of high-resolution structural T1-weighted images with dimensions of  $260 \times 311 \times 260$  voxels, providing detailed anatomical information across axial, coronal, and sagittal planes.

The analysis utilises data processed through HCP's standardised acquisition protocols and preprocessing pipelines, including the updated diffusion preprocessing pipeline (v3.19.0) incorporating FSL's [68] enhanced EDDY tool for improved slice outlier detection and movement artifact removal.

The participant age distribution ranges from 22-35 years (Figure 3.2), representing a critical period when brain structural development has largely stabilised, minimising age-related developmental confounds in twin similarity analyses. This age range ensures that observed neuroanatomical similarities between monozygotic twins primarily reflect genetic influences rather than ongoing developmental changes that could vary between twin pairs.

#### 3.1.2 Data Security

Data security and participant privacy were maintained through strict adherence to HCP Restricted Access Data Use Terms. All subjects were assigned randomly generated anonymised identifiers, replacing original HCP subject IDs to prevent any

---

<sup>1</sup><https://mkenely.com/3d-siamese-twin-ranking>

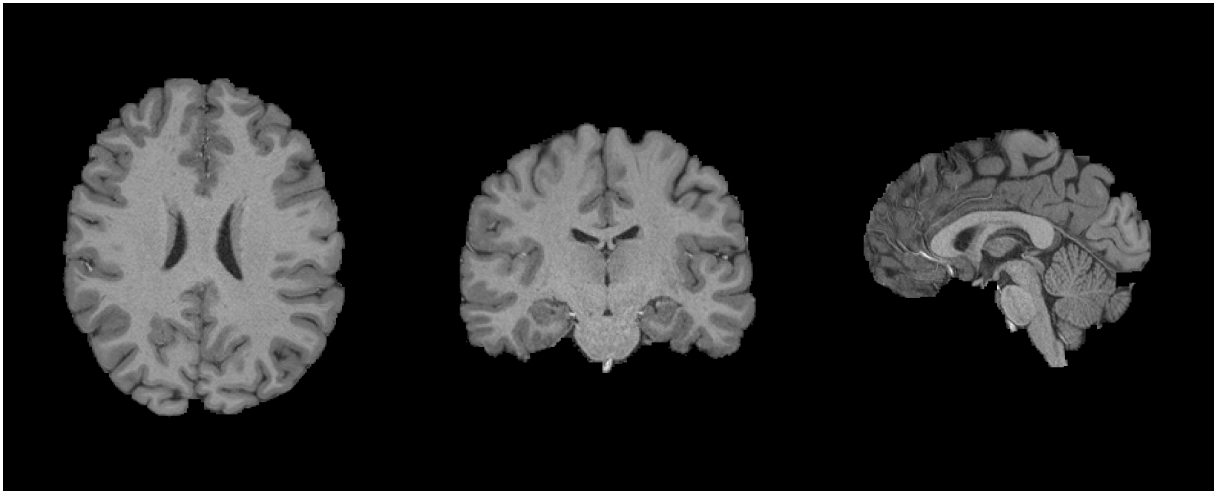


Figure 3.1 Left to right: axial, coronal, and sagittal views of a random HCP T1-weighted MRI scan from the HCP 1200 Subjects Data Release preprocessed using the HCP minimal preprocessing pipelines, showing the  $260 \times 311 \times 260$  voxel dimensions

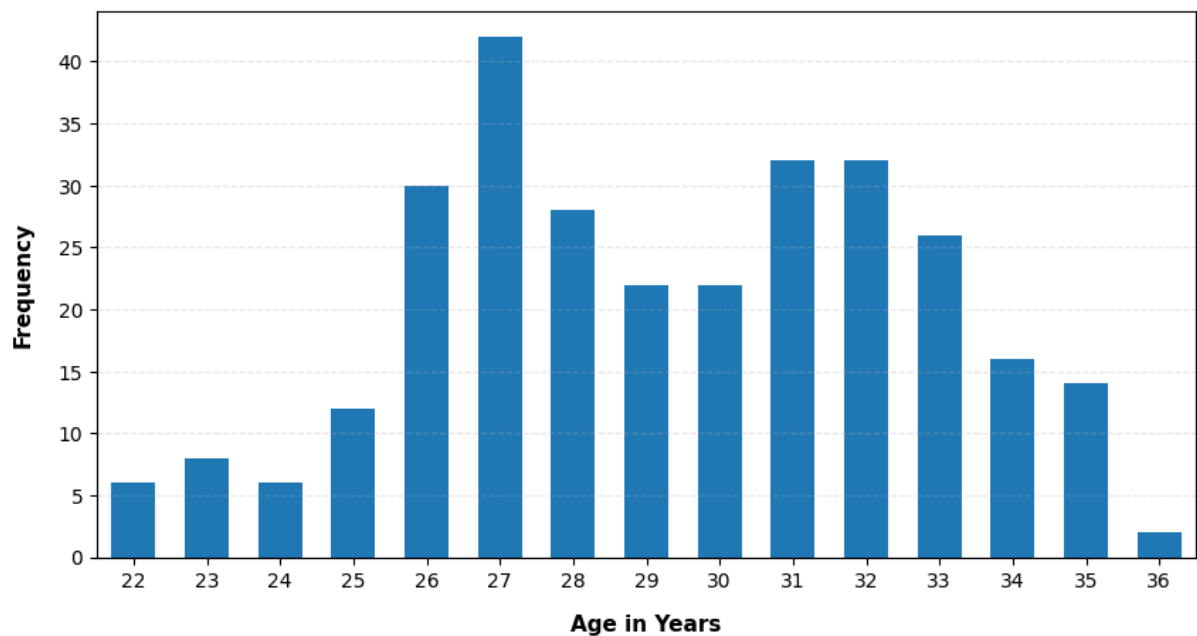


Figure 3.2 Histogram showing age distribution of the 138 monozygotic twin pairs (276 subjects) ranging from 22-36 years.

potential identification of individual participants. The randomisation process ensures no sequential patterns that could facilitate re-identification.

In accordance with HCP guidelines, all data was kept secure through password protection, with access limited to authorised personnel holding independent HCP Restricted Data approval. Family structure information (Family\_ID) and monozygotic twin status (ZygosityGT) was utilised for analysis while maintaining complete participant anonymity.

All procedures complied with institutional ethics requirements and HCP Data Use Terms, ensuring restricted data elements could not be combined to enable individual identification. The anonymisation was implemented prior to analysis, creating irreversible separation between participant identity and neuroimaging data.

### 3.1.3 Preprocessing

The HCP minimal preprocessing pipelines provide a comprehensive neuroimaging preparation framework specifically designed to preserve data integrity while accomplishing essential preprocessing tasks [69]. For the 138 monozygotic twin pairs analysed, these pipelines execute six core functions: correcting spatial artifacts and distortions, generating cortical surface reconstructions and tissue segmentations, enabling cross-modal registration within individual subjects, standardising data to common coordinate spaces, formatting outputs for visualisation software, and converting data to CIFTI grayordinate format.

Structural preprocessing utilises high-resolution 3D T1-weighted and T2-weighted images (0.7mm isotropic) to create accurate cortical surfaces and myelin maps [69]. The pipeline incorporates gradient nonlinearity correction to address spatial distortions from non-linear magnetic field gradients, maintaining anatomical precision necessary for cross-scanner compatibility.

FreeSurfer [70] performs cortical surface reconstruction, followed by registration to template spaces using MultiModal Surface Matching (MSM), specifically the MSMSulc approach that aligns surfaces based on cortical folding patterns. For twin identification applications, skull stripping and facial feature removal occur during FreeSurfer processing while maintaining data completeness, ensuring models extract features from genuine neuroanatomical patterns rather than confounding morphological characteristics. This conservative preprocessing strategy avoids destructive operations like extensive spatial smoothing or temporal filtering that eliminate substantial information content, preserving the subtle structural variations critical for detecting genetic similarities between twins.

A global bounding box was computed across all subjects to identify the minimal volume encompassing all brain tissue, removing excessive background regions. The

global bounding box coordinates were ((22, 236), (10, 304), (0, 223)) in the x, y, and z dimensions respectively. A global bounding box was employed rather than subject-specific cropping to ensure consistent volume dimensions across all subjects, avoiding alignment issues, and to enable efficient computation through a single precalculated operation when reverting to the original anatomical space.

Grand-mean intensity normalisation was applied with a target mean of 10,000 to standardise signal intensities across subjects:

$$I_{\text{normalised}} = I_{\text{original}} \times \frac{10000}{\text{mean}(I_{\text{masked}})} \quad (3.1)$$

where  $I_{\text{masked}}$  represents voxel intensities within the brain mask. Figure 3.3 demonstrates how this normalisation standardises intensity distributions between different subjects, reducing inter-subject variability that could confound twin identification.

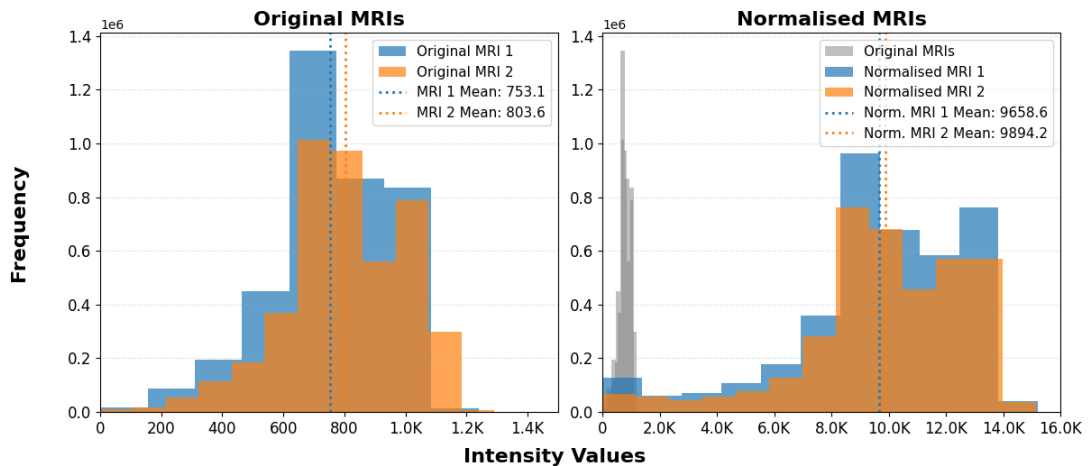


Figure 3.3 Intensity normalisation comparison showing histograms for two different subjects before (left) and after (right) grand-mean intensity normalisation.

Volumes were downsampled by a factor of 1/3 using trilinear interpolation to reduce computational requirements while maintaining spatial relationships critical for twin identification. Trilinear interpolation estimates values at new grid positions by performing linear interpolation along three orthogonal axes sequentially. For a target voxel at fractional coordinates  $(x, y, z)$ , the method first interpolates linearly between the eight nearest neighbour voxels in the original volume, computing weighted averages based on the distances to each neighbour.

Trilinear interpolation was employed for both downscaling original MRI volumes and upscaling generated relevance heatmaps. During downscaling, trilinear interpolation samples from the higher-resolution grid to create a coarser representation, averaging nearby voxel intensities to preserve the clear boundaries between anatomical regions that are more distinguishable at higher resolutions. For

upsampling, the method was specifically applied to relevance heatmaps rather than downscaling the HCP MMP atlases, which would risk losing smaller labelled regions. When upscaling heatmaps, trilinear interpolation estimates new voxel values by blending information from surrounding lower-resolution voxels, producing smooth transitions that avoid blocky artefacts which could otherwise leak into neighbouring anatomical regions. Since relevance heatmaps represent scores rather than anatomical structures, this smoothing is acceptable and preferable to the sharp discontinuities that nearest-neighbour interpolation would introduce. This approach balances computational efficiency with the preservation of anatomical boundaries in downsampled volumes and the generation of interpretable, artefact-free heatmaps in upsampled relevance maps.

This scaling factor was determined through experimentation to balance computational efficiency, essential for maximising batch size in triplet loss training, with neuroanatomical preservation. The final processed volumes had dimensions of  $86 \times 103 \times 86$  voxels, suitable for 3D CNN training while preserving essential neuroanatomical features. Figure 3.4 illustrates the preservation of anatomical detail despite the dimensional reduction.

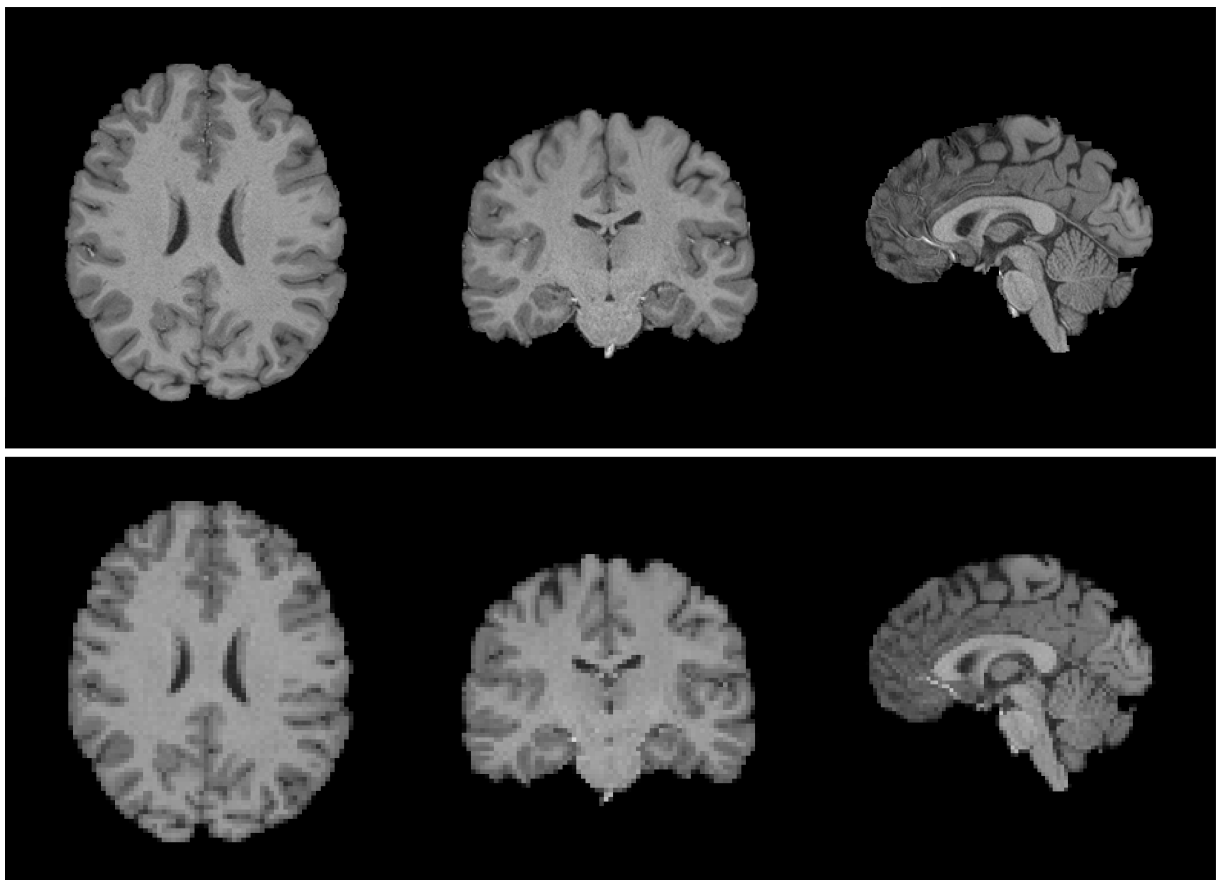


Figure 3.4 Spatial downscaling comparison showing axial, coronal, and sagittal views of the same brain before (top row:  $260 \times 311 \times 260$ ) and after (bottom row:  $86 \times 103 \times 86$ ) trilinear interpolation downscaling.

This comprehensive preprocessing approach ensures spatial accuracy and anatomical fidelity critical for detecting subtle structural similarities between monozygotic twins, while optimising data for efficient Siamese network training. Clinical integration was achieved through standard neuroimaging format conversion, with results exported to NIfTI and GIFTI formats compatible with Connectome Workbench [71] visualisation software, enabling immediate adoption in research and clinical neuroimaging workflows.

### 3.1.4 Data Augmentation

To enhance robustness and generalisation of the Siamese network, a targeted 3D MRI augmentation pipeline was implemented with probability  $p=0.5$  for each transformation. Four key augmentations were selected based on analysis of the HCP minimal preprocessing pipeline, representing good practice for model generalisability when applied to datasets with different preprocessing protocols.

Random rotations were applied along the sagittal-coronal, coronal-axial, or sagittal-axial planes with angles from  $-20^\circ$  to  $20^\circ$  using nearest-neighbour interpolation. The rotation range simulates natural head positioning variations during scanning while maintaining anatomical relationships, compensating for residual rotational differences between subjects that persist after cross-modal registration and standardisation [69].

Random flipping was applied along randomly selected spatial dimensions to enforce left-right symmetry invariance, ensuring that brain hemispheric differences do not influence twin identification while preserving anatomical fidelity.

Gaussian noise was added to entire volumes with maximum standard deviation of 0.02, simulating scanner noise and acquisition variations. This accounts for subtle noise variations between acquisition sessions without compromising the signal quality of the high-resolution imaging [72].

Intensity scaling was employed to account for minor variations in signal intensity across different scanning sessions. Intensity shifting was excluded to preserve the careful intensity normalisation and bias field correction already applied in HCP preprocessing, which ensures consistent intensity profiles across subjects. However, intensity scaling remains beneficial for model generalisability when applied to datasets with different acquisition protocols.

### 3.1.5 Data Preparation

The HCP data required transformation into a format suitable for triplet loss training. Preprocessed brain volumes were extracted from twin pair directories and

consolidated into a centralised repository with individual subject files to eliminate data redundancy while maintaining efficient access patterns.

Twin pairs were randomly partitioned into train (60%), validation (20%), and test (20%) sets with strict family separation to prevent data leakage. The validation set monitors training progress and prevents overfitting, while the separate test set provides unbiased final evaluation. Split-specific metadata maintained references to twin relationships and negative sampling candidates without duplicating volume data.

The triplet sampling strategy employs hard negative mining to maximise training effectiveness by identifying the most challenging negative samples using cosine similarity in embedding space. This selects non-family subjects with the highest similarity to anchor subjects through vectorised embedding computations and precomputed similarity matrices. Each twin pair generates multiple triplets through bidirectional anchor-positive relationships, with caching mechanisms storing computed embeddings to optimise data loading efficiency.

This approach ensures balanced representation of twin relationships while providing maximally informative negative examples that become progressively more challenging as the model improves, maintaining strong learning signals essential for effective metric learning convergence.

### **Triplet Loss vs ACE Model**

This framework learns to distinguish monozygotic twin brains from unrelated individuals rather than quantifying heritability through ACE variance decomposition. The HCP dataset contains 149 monozygotic pairs (of which 138 had valid 3T MRI data) but only 94 dizygotic pairs, creating insufficient statistical power for reliable ACE modelling in deep learning contexts. The triplet loss approach circumvents this limitation by utilising all available twin pairs with negatives generated from different families.

Unlike classical twin studies that determine brain structure is 60-80% heritable through variance decomposition, this method cannot quantify heritability percentages or distinguish genetic effects from shared environmental factors. The framework discards information about intermediate genetic relatedness that dizygotic twins provide, focusing instead on the binary distinction between genetically identical and unrelated individuals.

However, successful implementation would demonstrate that genetic signatures in brain structure are sufficiently distinct for automated detection while providing complementary spatial insights to ACE models. High performance would validate that monozygotic twins share detectable neural patterns distinguishing them from even their most similar-appearing unrelated individuals. Layer-wise relevance propagation

would reveal which brain regions are most discriminative for twin identification, providing a relative ranking based on discriminative importance rather than independent heritability estimates. While ACE studies report regional heritability percentages, this approach would rank regions from most to least informative for genetic relatedness detection across the whole brain.

## 3.2 Siamese Architecture

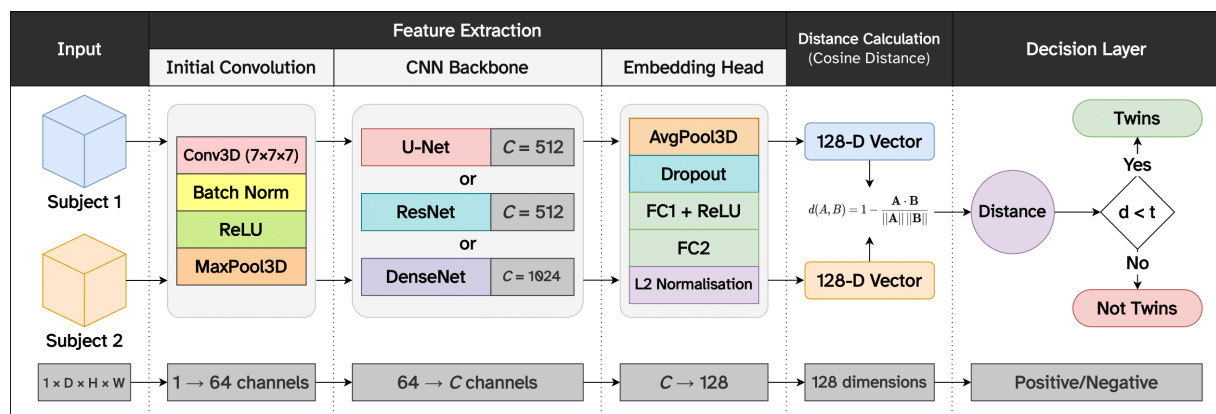


Figure 3.5 Siamese network architecture overview. Two identical 3D CNN branches process brain MRI volumes to generate 128-dimensional embeddings. Cosine distance between embeddings determines twin versus non-twin classification through threshold comparison.

### 3.2.1 Siamese Network Design

The Siamese network outlined in Figure 3.5 processes two 3D brain volumes through two identical CNN backbone architectures to generate 128-dimensional discriminative embeddings for twin identification. Weight sharing between branches ensures consistent feature extraction while enabling direct similarity comparison.

Each branch begins with a  $7 \times 7 \times 7$  convolutional layer capturing low-level volumetric features, followed by hierarchical feature extraction through the chosen backbone. Channel dimensions expand progressively while spatial dimensions reduce, transforming detailed volumetric data into abstract 128-dimensional anatomical embeddings.

Table 3.1 Siamese Network Architecture

Component	Input Shape	Operation	Output Shape
Input	$(B, 2, D, H, W)$	3D Brain Volumes	$(B, 2, D, H, W)$
Branch 1	$(B, 1, D, H, W)$	CNN Backbone	$(B, 128)$
Branch 2	$(B, 1, D, H, W)$	CNN Backbone (shared weights)	$(B, 128)$
Embeddings	$(B, 128), (B, 128)$	Cosine Distance: $d = 1 - \cos(e_1, e_2)$	$(B, 2)$

Channel and spatial attention mechanisms are applied to deeper layers (2, 3, 4), enhancing focus on discriminative anatomical regions while suppressing irrelevant features.

The Siamese network concludes with an embedding head outlined in Table 3.2:

Table 3.2 Siamese Network Embedding Head

Step	Operation	Input	Output	Parameters
1	AdaptiveAvgPool3d	$(B, C, D, H, W)$	$(B, C, 1^3)$	-
2	View	$(B, C, 1^3)$	$(B, C)$	-
3	Dropout	$(B, C)$	$(B, C)$	$p = 0.3$
4	FC1 + ReLU	$(B, C)$	$(B, 256)$	$C \rightarrow 256$
5	FC2	$(B, 256)$	$(B, 128)$	$256 \rightarrow 128$
6	L2 Normalise	$(B, 128)$	$(B, 128)$	-

The choice of 128-dimensional embeddings represents an optimal balance between representational capacity and computational efficiency. Experiments with 64-dimensional embeddings showed significantly slower convergence, indicating insufficient capacity for complex anatomical pattern discrimination. Peak Signal-to-Noise Ratio (PSNR) analyses revealed the start of embedding redundancy at 128 dimensions, demonstrating that higher-dimensional representations (256+) would be computationally excessive without meaningful performance gains. The supporting PSNR analysis is detailed in Subsection 4.3.3.

Twin classification uses cosine distance between normalised 128-dimensional embeddings, with distances computed as  $d = 1 - \cos(e_1, e_2)$ . Lower distances indicate higher similarity.

### 3.3 3D CNN Backbones

#### 3.3.1 U-Net

The modified U-Net backbone is derived from the `pytorch-3dunet` repository<sup>2</sup>, originally designed for 3D biomedical image segmentation. The architecture has been fundamentally restructured to function as a feature extractor for twin identification rather than segmentation. The encoder portion is retained with four progressive downsampling layers that extract multi-scale features using double convolution or residual blocks. Each encoder module combines optional max pooling for spatial downsampling with a basic convolution block that performs feature extraction, where the first encoder operates without pooling while subsequent encoders apply  $2 \times 2 \times 2$  max pooling before convolution. The decoder path is completely eliminated and replaced with global average pooling followed by fully connected layers that generate 128-dimensional normalised embeddings. Channel and spatial attention mechanisms are integrated into encoder layers 2-4 to enhance focus on discriminative neuroanatomical features critical for twin identification. The architecture is limited to four encoder layers to prevent over-pooling of the downscaled  $86 \times 103 \times 86$  volumes while maintaining spatial resolution necessary for detecting subtle anatomical differences between twins.

Table 3.3 U-Net Architecture Comparison

Component	Original	Modified
Purpose	3D segmentation	Embedding generation
Decoder	Full decoder path	Removed
Output	Segmentation masks	128-dim embeddings
Attention	None	Layers 2-4
Final layers	$1 \times 1 \times 1$ conv + activation	Global pool + FC
Depth	Variable	4 encoder layers

<sup>2</sup><https://github.com/wolny/pytorch-3dunet>

Table 3.4 U-Net Backbone Architecture

Layer	Input Ch.	Output Ch.	Kernel	Stride	Padding	Attention
Conv1	1	64	$7^3$	$1^3$	$3^3$	-
BN1 + ReLU	64	64	-	-	-	-
MaxPool	64	64	$3^3$	$2^3$	$1^3$	-
Encoder1	64	64	$3^3$	$1^3$	$1^3$	-
Encoder2	64	128	$3^3$	$2^3$	$1^3$	●
Encoder3	128	256	$3^3$	$2^3$	$1^3$	●
Encoder4	256	512	$3^3$	$1^3$	$1^3$	●
Embedding	512	128	-	-	-	-

### 3.3.2 ResNet

The 3D ResNet backbone originates from the 3D-ResNets-PyTorch repository<sup>3</sup>, initially developed for video action recognition with temporal modelling capabilities. The architecture has been adapted for neuroimaging analysis by replacing the classification head with a two-layer fully connected network that produces 128-dimensional normalised embeddings. Each layer consists of sequential BasicBlocks (BB) that implement residual learning through skip connections, where Layer1 contains 2 BasicBlocks operating at 64 channels without downsampling, Layer2 applies spatial downsampling to 128 channels, Layer3 maintains spatial resolution at 256 channels using dilated convolutions with dilation rate 2, and Layer4 operates at 512 channels with dilation rate 4. Each BasicBlock contains two  $3 \times 3 \times 3$  convolutions with batch normalisation and ReLU activation, connected by a residual skip connection that adds the input to the output before final activation. Dilated convolutions are implemented in layers 3-4 with dilation rates of 2 and 4 respectively to maintain spatial resolution while expanding receptive fields, which is crucial for the downscaled neuroimaging volumes. Channel and spatial attention modules are integrated into layers 2-4 to enhance focus on twin-discriminative features. The temporal stride configurations are modified to spatial-only strides since neuroimaging data lacks temporal dimension. ResNet-18 configuration is employed with 2 BasicBlocks per layer to minimise overfitting on the relatively small twin dataset while maintaining the residual learning benefits for complex anatomical pattern recognition.

<sup>3</sup><https://github.com/kenshohara/3D-ResNets-PyTorch>

Table 3.5 ResNet Architecture Comparison

Component	Original	Modified
Purpose	Video action recognition	Embedding generation
Convolutions	Standard	Dilated (layers 3-4)
Dilation rates	1	1, 1, 2, 4
Attention	None	Layers 2-4
Output	Classification logits	128-dim embeddings
Stride pattern	Temporal+spatial	Spatial only

Table 3.6 ResNet Backbone Architecture

Layer	Input Ch.	Output Ch.	Kernel	Stride	Padding	Attention
Conv1	1	64	$7^3$	(1, 2, 2)	$3^3$	-
BN1 + ReLU	64	64	-	-	-	-
MaxPool	64	64	$3^3$	$2^3$	$1^3$	-
Layer1 (2×BB)	64	64	$3^3$	$1^3$	$1^3$	-
Layer2 (2×BB)	64	128	$3^3$	$2^3$	$1^3$	•
Layer3 (2×BB)	128	256	$3^3$	$1^3$	$2^3$	•
Layer4 (2×BB)	256	512	$3^3$	$1^3$	$4^3$	•
Embedding	512	128	-	-	-	-

### 3.3.3 DenseNet

The 3D DenseNet backbone is derived from the 3D-ResNets-PyTorch repository<sup>4</sup>, adapted from 2D image classification to 3D video analysis. Dense block connections allow direct information flow between all layer pairs within blocks, reducing parameter requirements while mitigating vanishing gradient problems and enhancing representational diversity across network depths. Dense blocks implement this connectivity through concatenation of feature maps, where DenseBlock1 contains 6 layers that expand from 64 to 256 channels, DenseBlock2 contains 12 layers expanding from 128 to 512 channels, DenseBlock3 contains 24 layers expanding from 256 to 1024 channels, and DenseBlock4 contains 16 layers expanding from 512 to 1024 channels. Each dense layer within a block applies batch normalisation, ReLU activation, and  $1 \times 1 \times 1$  convolution for bottleneck compression, followed by batch normalisation, ReLU activation, and  $3 \times 3 \times 3$  convolution for feature extraction, with the growth rate of 32 determining the number of new features added per layer. Transition layers between dense blocks perform feature map reduction and spatial downsampling through  $1 \times 1 \times 1$

<sup>4</sup><https://github.com/kenshohara/3D-ResNets-PyTorch>

convolution that halves the channel dimension, followed by  $2 \times 2 \times 2$  average pooling for spatial reduction. The classification layer is replaced with the same two-layer fully connected architecture used in other models, outputting 128-dimensional normalised embeddings. Channel and spatial attention modules are integrated into dense blocks 2-4, applied after the dense connectivity concatenation to refine feature representations. DenseNet-121 configuration with layer arrangement (6,12,24,16) is employed as the smallest available configuration, utilising a growth rate of 32 to balance feature richness with computational efficiency for the neuroimaging task.

Table 3.7 DenseNet Architecture Comparison

Component	Original	Modified
Purpose	2D/3D classification	Embedding generation
Attention	None	Dense blocks 2-4
Output	Classification logits	128-dim embeddings
Depth mapping	DenseNet depths only	ResNet depth compatible
Configuration	Standard DenseNet	DenseNet-121 (minimal)
Final layers	Linear classifier	Global pool + FC

Table 3.8 DenseNet Backbone Architecture

Layer	Input Ch.	Output Ch.	Growth Rate	Layers	Attention
Conv1	1	64	-	-	-
BN1 + ReLU	64	64	-	-	-
MaxPool	64	64	-	-	-
DenseBlock1	64	256	32	6	-
Transition1	256	128	-	-	-
DenseBlock2	128	512	32	12	●
Transition2	512	256	-	-	-
DenseBlock3	256	1024	32	24	●
Transition3	1024	512	-	-	-
DenseBlock4	512	1024	32	16	●
Embedding	1024	128	-	-	-

## 3.4 Training

### 3.4.1 Loss Function

Siamese networks commonly employ contrastive loss or triplet loss for learning discriminative embeddings. While contrastive loss optimises pairs of similar/dissimilar

samples, triplet loss provides superior learning by simultaneously optimising three samples: an anchor ( $\mathbf{a}$ ), positive ( $\mathbf{p}$ , the anchor's corresponding twin), and negative ( $\mathbf{n}$ , a subject which is not the anchor's twin). This approach compels the model to develop more discriminative feature representations by explicitly contrasting against challenging negative examples, which is crucial for detecting the subtle anatomical variations that distinguish twins.

Cosine similarity calculates the angular relationship between two vectors, providing an orientation-based measure independent of magnitude differences. For L2-normalised embeddings  $\mathbf{a}$  and  $\mathbf{p}$ :

$$\cos(\mathbf{a}, \mathbf{p}) = \frac{\mathbf{a} \cdot \mathbf{p}}{|\mathbf{a}||\mathbf{p}|} = \mathbf{a} \cdot \mathbf{p} \quad (3.2)$$

Since embeddings are unit vectors, the dot product directly yields values ranging from -1 (completely dissimilar) to 1 (identical). Converting to distance metric form:

$$d(\mathbf{a}, \mathbf{e}) = 1 - \cos(\mathbf{a}, \mathbf{e}) \quad (3.3)$$

transforms the similarity range to distance range  $[0, 2]$ , where 0 indicates identical embeddings and 2 indicates completely opposite embeddings.

Triplet loss optimises these distances by reducing related pair (anchor-positive) separation while increasing non-related pair (anchor-negative) separation by margin  $\alpha$ :

$$\mathcal{L} = \max(0, d(\mathbf{a}, \mathbf{p}) - d(\mathbf{a}, \mathbf{n}) + \alpha) \quad (3.4)$$

Operating on  $[batch\_size, 128]$  dimensional feature embeddings, the loss computes cosine distances between anchor-positive and anchor-negative pairs. Active triplets (loss > 0) identify challenging samples requiring learning, with all triplets contributing to gradient updates.

### Hard Negative Mining Strategy

The framework employs hard negative mining by selecting the most challenging negative samples to improve model discriminative capability. Given an anchor embedding  $\mathbf{a}$  and a set of candidate negative embeddings  $\mathcal{N} = \{\mathbf{n}_1, \mathbf{n}_2, \dots, \mathbf{n}_k\}$ , the hard negative  $\mathbf{n}^*$  is selected as:

$$\mathbf{n}^* = \arg \min_{\mathbf{n}_i \in \mathcal{N}} d(\mathbf{a}, \mathbf{n}_i) \quad (3.5)$$

where the cosine distance is computed as:

$$d(\mathbf{a}, \mathbf{n}_i) = 1 - \cos(\mathbf{a}, \mathbf{n}_i) = 1 - \mathbf{a} \cdot \mathbf{n}_i \quad (3.6)$$

This selection identifies the negative sample with the highest cosine similarity (lowest distance) to the anchor, representing the most confusing negative case that challenges the model’s ability to distinguish between twin and non-twin pairs. The vectorised implementation efficiently computes these selections using batch matrix operations:

$$\mathbf{S} = \mathbf{a}\mathbf{N}^T \quad (3.7)$$

$$\mathbf{D} = \mathbf{1} - \mathbf{S} \quad (3.8)$$

$$j^* = \arg \min_j \mathbf{D}_j \quad (3.9)$$

where  $\mathbf{N} \in \mathbb{R}^{k \times d}$  is the matrix of candidate negative embeddings,  $\mathbf{S}$  contains cosine similarities,  $\mathbf{D}$  contains distances, and  $j^*$  identifies the hardest negative index.

By consistently training on these difficult cases, the network develops more robust representations that generalise better to subtle twin pair distinctions during inference. This strategy identifies the most challenging negative examples that lie closest to the decision boundary, forcing the model to learn fine-grained discriminative features necessary to distinguish between anatomically similar non-twin pairs.

### 3.4.2 Training Algorithm

The training procedure implements triplet loss optimisation with hard negative mining:

Training utilises AdamW, an optimiser that pairs Adam’s adaptive learning rates with separated weight decay for enhanced generalisation versus standard Adam. A learning rate of  $10^{-4}$  is used to balance convergence speed with stability for deep 3D networks.

OneCycleLR scheduling implements a cyclical learning rate that starts low, peaks at maximum, then decreases below the initial rate. The 30% warmup period allows gradual adaptation to the triplet loss landscape, while final low rates enable fine-tuning of embeddings.

Large effective batch sizes are crucial for triplet loss as they provide more negative examples for hard mining within each batch. Larger batches improve negative sample diversity and stabilise training dynamics. Gradient accumulation enables these large effective batch sizes despite GPU memory constraints - U-Net and ResNet use 2 accumulation steps with batch size 256, while DenseNet requires 8 steps with batch size 64 due to higher memory requirements.

Mixed precision training with FP16 reduces memory usage by  $\sim 50\%$  and accelerates training, enabling larger batch sizes necessary for effective triplet mining. However, FP16’s limited numerical range can cause gradient underflow, where very

---

**Algorithm 1** Siamese Network Training with Triplet Loss
 

---

**Input:** Model  $f_\theta$ , Training triplets, Validation triplets, Epochs  $N$   
 Initialise embedding layers with Xavier uniform (gain=0.5)  
 Initialise AdamW optimiser:  $lr = 10^{-4}$ ,  $\beta = (0.9, 0.999)$ , weight decay  $5 \times 10^{-4}$   
 Initialise OneCycleLR scheduler with  $max\_lr = 10^{-4}$ , warmup 30%  
 Initialise mixed precision scaler and gradient accumulation  
 Set margin  $\alpha = 0.1$ , regularisation  $\lambda = 0.01$   
**for** each epoch  $e = 1$  to  $N$  **do**  
   Set model to training mode  
   Update hard negative mining strategy  
   **for** each batch  $(a, p, n)$  in training data **do**  
     Forward pass:  $(e_a, e_p, e_n) = f_\theta(a, p, n)$   
     Compute triplet loss:  $\mathcal{L}_{triplet} = \max(0, d(e_a, e_p) - d(e_a, e_n) + \alpha)$   
     Add regularisation:  $\mathcal{L}_{reg} = \lambda \sum ||e_i||_2^2$   
     Total loss:  $\mathcal{L} = \mathcal{L}_{triplet} + \mathcal{L}_{reg}$   
     Scale and accumulate gradients  
     **if** accumulation step reached **then**  
       Clip gradients (max norm 0.5)  
       Update parameters and scheduler  
       Clear gradients  
     **end if**  
   **end for**  
   Evaluate on validation set  
   Track metrics: AUC-ROC, F1, distance gap, separability  
   Save best models based on multiple criteria  
   Clear embedding cache every 3 epochs  
**end for**

---

small gradient values are rounded to zero and stop parameter updates. Automatic loss scaling prevents this by multiplying the loss before backpropagation to keep gradients within FP16's representable range, then scaling them back down for parameter updates.

Xavier uniform initialisation with reduced gain (0.5) for final embedding layers promotes stable training and prevents initial embedding collapse. L2 regularisation on embeddings ( $\lambda = 0.01$ ) prevents feature collapse by penalising small pairwise distances between different samples. Gradient clipping with maximum norm 0.5 prevents exploding gradients common in deep 3D networks, ensuring stable convergence.

## 3.5 Quantitative Evaluation

A comprehensive evaluation protocol was implemented to assess the Siamese network's twin classification performance across three CNN architectures. The evaluation focuses on F1-score optimisation while monitoring multiple performance metrics throughout training.

The evaluation methodology accounts for the natural class imbalance in twin identification, where potential non-twin combinations significantly outnumber available twin pairs. While twin pairs are constrained by the available twin subjects in the dataset, non-twin pairs can be formed from any random selection of unrelated individuals, creating a significantly larger pool of negative examples. To ensure robust and unbiased evaluation, 10 independent datasets are generated, each containing balanced numbers of twin and non-twin pairs. For each dataset, non-twin pairs are randomly sampled from the pool of unrelated subjects, creating different negative example compositions across the 10 evaluation sets.

These 10 datasets remain consistent across all model evaluations, ensuring fair comparison between architectures. Each model (ResNet, U-Net, DenseNet) is evaluated on the same 10 datasets, and performance metrics are aggregated to compute means and standard deviations. This approach provides confidence intervals for model performance while accounting for variability in non-twin pair selection, ensuring that reported performance metrics are not biased toward specific non-twin pair combinations and providing reliable estimates of model generalisation capability.

### 3.5.1 Model Selection and Evaluation Protocol

During training, multiple metrics are continuously monitored on validation data to track model performance:

- **Validation Loss:** Triplet loss on validation data

- **Distance Gap:** Difference between average negative pair distance and average positive pair distance
- **Separability:** Proportion of negative distances exceeding positive distances
- **AUC-ROC:** Area under the ROC curve
- **F1-Score:** Harmonic mean of precision and recall

Model checkpoints are saved based on improvements in these validation metrics throughout training. After training completion, all saved checkpoints for each architecture are evaluated on the test set. The checkpoint achieving the highest F1-score on test data is selected as the final model for that architecture, yielding optimal ResNet, U-Net, and DenseNet candidates.

F1-score serves as the primary selection criterion because it balances precision and recall, ensuring models neither miss twin pairs nor generate excessive false positives. This balanced approach maximises overall embedding quality and is crucial for downstream brain region heatmap analysis, guaranteeing high-quality learned features that reliably capture twin relationships while maintaining discriminative power for interpretability studies.

### 3.5.2 Performance Metrics

Classification performance is assessed through multiple complementary metrics, each providing distinct insights into model behaviour.

Accuracy measures overall correctness across all pairs:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (3.10)$$

Although intuitive, accuracy can be deceptive in imbalanced datasets where models achieve high scores by correctly classifying the majority class while performing poorly on the minority class.

Precision quantifies exactness when predicting twin pairs:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (3.11)$$

High precision indicates that positive predictions are predominantly correct, crucial when false positive twin identifications have significant practical consequences.

Recall measures completeness in identifying actual twin pairs:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3.12)$$

High recall ensures most true twin relationships are captured, important when missing twin pairs would be problematic for downstream applications.

Specificity quantifies the ability to correctly identify non-twin pairs:

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (3.13)$$

High specificity indicates accurate negative predictions, preventing false twin identifications.

Negative Predictive Value (NPV) measures the reliability of negative predictions:

$$\text{NPV} = \frac{TN}{TN + FN} \quad (3.14)$$

High NPV ensures that when the model predicts non-twin pairs, the prediction is likely correct.

F1-score provides balanced assessment of precision and recall:

$$\text{F1-Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (3.15)$$

F1-score provides an assessment of the model's balanced performance across both classes. This indicates the overall quality of the predictions, which is crucial in this context as the models will be used for the downstream task of embedding generation, requiring reliable representation learning.

### 3.5.3 ROC and Precision-Recall Curve Analysis

ROC curves plot True Positive Rate (TPR) against False Positive Rate (FPR) across all classification thresholds. TPR is equivalent to recall, while FPR is equivalent to  $1 - \text{specificity}$ :

$$\text{TPR} = \frac{TP}{TP + FN} \quad (3.16)$$

$$\text{FPR} = \frac{FP}{FP + TN} \quad (3.17)$$

AUC-ROC evaluates performance across all threshold values without requiring threshold selection:

$$\text{AUC-ROC} = \int_0^1 \text{TPR}(\text{FPR}^{-1}(t)) dt \quad (3.18)$$

PR curves plot precision against recall across thresholds:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (3.19)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3.20)$$

ROC analysis measures the model's fundamental ability to discriminate between twin and non-twin pairs across all possible distance thresholds, providing comprehensive assessment of ranking quality independent of specific decision boundaries. In contrast, PR analysis emphasises precision-recall trade-offs when identifying twin pairs, offering focused evaluation at operationally relevant thresholds where minimising false twin identifications is critical.

Given the balanced dataset structure where non-twin pairs are randomly sampled from the subject pool, both metrics provide meaningful performance assessment. ROC curves offer broad discriminative evaluation suitable for comparing overall model capabilities, while PR curves provide targeted analysis for deployment scenarios requiring high-confidence twin identification.

### 3.5.4 Threshold Optimisation

Twin prediction follows a distance-based classification approach:

$$\hat{y} = \begin{cases} \text{twins,} & \text{if } d(\mathbf{e}_1, \mathbf{e}_2) < \tau \\ \text{non-twins,} & \text{otherwise} \end{cases} \quad (3.21)$$

where  $d(\mathbf{e}_1, \mathbf{e}_2) = 1 - \cos(\mathbf{e}_1, \mathbf{e}_2)$  represents cosine distance between L2-normalised 128-dimensional embeddings, and  $\tau$  is the optimal threshold determined through F1-score maximisation:

$$\tau^* = \arg \max_{\tau_i \in \mathcal{T}} \text{F1}(\tau_i) \quad (3.22)$$

where  $\tau^*$  is the optimal threshold,  $\tau_i$  represents candidate threshold values from the set  $\mathcal{T}$  of all evaluated thresholds, and  $\text{F1}(\tau_i)$  is the F1-score function evaluated at threshold  $\tau_i$ . The optimal classification threshold is determined by evaluating F1-scores across the distance range observed in test data and selecting the threshold that yields maximum harmonic mean of precision and recall.

## 3.6 Qualitative Evaluation

### 3.6.1 Embeddings Interpretation

#### t-Distributed Stochastic Neighbour Embedding (t-SNE)

t-SNE projects the 128-dimensional embeddings into 2D visualisation space while preserving local neighborhood relationships. The algorithm accomplishes this by minimising divergence between probability distributions in the original high-dimensional space and the reduced low-dimensional representation, ensuring that closely related data points remain clustered while unrelated points are separated [60].

In the original 128-dimensional embedding space, the probability that embedding  $x_j$  is a neighbour of  $x_i$  is:

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2/2\sigma_i^2)} \quad (3.23)$$

where  $x_i$  and  $x_j$  represent high-dimensional embeddings, and  $\sigma_i$  is the variance of the Gaussian centred at  $x_i$ .

In the 2D visualisation space, t-SNE uses a Student's t-distribution:

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|y_k - y_i\|^2)^{-1}} \quad (3.24)$$

where  $y_i$  and  $y_j$  are the corresponding 2D coordinates of points  $x_i$  and  $x_j$ .

The algorithm minimises the Kullback-Leibler divergence between these distributions:

$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (3.25)$$

where  $P$  and  $Q$  represent the high and low-dimensional probability distributions.

It should be noted that t-SNE visualisations are highly sensitive to hyperparameter selection, particularly perplexity, which controls the effective number of neighbours considered during embedding. Additionally, distances between clusters in the 2D projection do not necessarily reflect true distances in the original high-dimensional space, and visual clustering patterns can emerge even in data with no underlying structure. Therefore, t-SNE serves as a complementary qualitative tool to support quantitative distance metrics rather than as standalone evidence for embedding quality.

### Peak Signal-to-Noise Ratio (PSNR)

PSNR is used to assess the spatial uniqueness of individual embedding dimensions by measuring how each dimension's LRP heatmap deviates from the mean embedding LRP heatmap pattern across all dimensions. For each of the 128 embedding dimensions, PSNR quantifies whether dimensions capture distinct neuroanatomical regions or converge on similar spatial areas with only minor variations.

PSNR is calculated as:

$$\text{PSNR} = 10 \log_{10} \left( \frac{\text{MAX}^2}{\text{MSE}} \right) \quad (3.26)$$

where MAX denotes the maximum possible intensity value and MSE represents the mean squared error between individual embedding dimensions' LRP patterns and the overall mean pattern:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (I_{\text{embedding}}(i) - I_{\text{mean}}(i))^2 \quad (3.27)$$

where  $N$  is the total voxel count in the brain volume,  $I_{\text{embedding}}(i)$  represents the LRP intensity at voxel position  $i$  for a specific embedding dimension, and  $I_{\text{mean}}(i)$  is the corresponding intensity at position  $i$  in the mean LRP pattern averaged across all embedding dimensions.

Higher PSNR values indicate embedding dimensions that focus on similar brain regions as the population average, suggesting convergence on common spatial patterns. Lower PSNR values identify dimensions with spatially distinct activation patterns, revealing embeddings that capture unique neuroanatomical regions rather than subtle variations within overlapping areas. If even low-PSNR embeddings show spatial similarity to the population average, this indicates the model dedicates embedding space to capturing minor structural differences within the same general brain regions rather than encoding fundamentally different anatomical areas.

## 3.6.2 Twin vs Non-twin Similarity Analysis

### Layer-Wise Relevance Propagation Attribution

LRP generates embedding-specific attribution maps through gradient-based analysis of the trained Siamese network. For each subject, the process computes attribution scores across all 128 embedding dimensions using integrated gradients, then creates a reference LRP map by aggregating across all subjects to minimise individual variation.

The attribution computation traces relevance from each embedding dimension

back to input voxels using integrated gradients:

$$IG_{d,i,j,k} = (x_{i,j,k} - x'_{i,j,k}) \times \int_{\alpha=0}^1 \frac{\partial e_d}{\partial x'_{i,j,k} + \alpha(x_{i,j,k} - x'_{i,j,k})} d\alpha \quad (3.28)$$

where  $IG_{d,i,j,k}$  is the integrated gradient value for embedding dimension  $d$  at spatial coordinates  $(i, j, k)$ ,  $x_{i,j,k}$  represents the original input voxel intensity at coordinates  $(i, j, k)$ ,  $x'_{i,j,k}$  represents the baseline (zero) input at the same coordinates,  $e_d$  denotes the  $d$ -th embedding dimension output, and  $\alpha$  is the integration variable ranging from 0 to 1.

Individual subject LRP arrays with shape  $(128 \times W \times H \times D)$  are averaged to create the reference LRP:

$$R_{d,i,j,k} = \frac{1}{N} \sum_{s=1}^N L_{s,d,i,j,k} \quad (3.29)$$

where  $R_{d,i,j,k}$  is the reference LRP value at embedding dimension  $d$  and spatial coordinates  $(i, j, k)$ ,  $L_{s,d,i,j,k}$  is the LRP attribution for subject  $s$  at dimension  $d$  and coordinates  $(i, j, k)$ , and  $N$  is the total number of subjects.

This reference map provides generalised spatial patterns showing how each embedding dimension relates to brain anatomy, removing subject-specific variability. The reference LRP acts as a spatial mask that is weighted by subject-specific embedding distances, ensuring that embedding distance patterns drive region importance rather than individual LRP variations.

### Discriminative Brain Region Analysis

The discriminative region analysis uses embedding distance differences between subject pairs to weight the reference LRP map, identifying brain regions that contribute most strongly to twin versus non-twin discrimination. This approach uses the learned embedding space as a proxy for genetic similarity.

For each subject pair, the cosine distance between their embeddings quantifies their dissimilarity in the learned feature space:

$$d = 1 - \cos(\mathbf{e}_{s1}, \mathbf{e}_{s2}) \quad (3.30)$$

where  $\mathbf{e}_{s1}$  and  $\mathbf{e}_{s2}$  are the normalised embedding vectors for subjects  $s1$  and  $s2$  respectively. Twin pairs typically exhibit smaller distances than unrelated pairs.

The embedding distance weights the reference LRP map to create pair-specific attribution patterns:

$$W_{i,j,k} = d \times \frac{1}{D} \sum_{d=1}^D R_{d,i,j,k} \quad (3.31)$$

where  $W_{i,j,k}$  represents the weighted attribution value at spatial coordinates  $(i, j, k)$ ,  $d$  is the embedding distance, and  $D = 128$  is the embedding dimensionality. Higher embedding distances amplify the spatial attribution pattern.

Population-level discriminative patterns emerge through separate averaging of weighted attribution maps across twin and non-twin groups:

$$H_{disc} = \overline{W_{twin}} - \overline{W_{non-twin}} \quad (3.32)$$

where  $H_{disc}$  represents the discriminative heatmap,  $W_{twin}$  and  $W_{non-twin}$  are the weighted attribution maps for twin and non-twin pairs respectively, and the overbar denotes the mean operation across all pairs within each group.

Negative values indicate regions where non-twin pairs exhibited larger embedding distances than twin pairs, representing brain areas that contribute strongly to discrimination. Positive values reveal regions where twin pairs unexpectedly showed larger distances, potentially indicating areas where the model learned atypical patterns. This spatial attribution framework provides interpretable insights into which neuroanatomical structures the trained Siamese network prioritises when making identification decisions.

### 3.6.3 Ensemble Modelling

The ensemble approach combines discriminative heatmaps from all three architectures using performance-weighted averaging to leverage the complementary strengths of each model. Individual model heatmaps are weighted according to their F1-score performance, with exponential scaling to emphasise superior models while maintaining contributions from weaker performers.

Performance-based weights are computed using exponential scaling:

$$w_m = \frac{F1_m^\tau}{\sum_{i=1}^M F1_i^\tau} \quad (3.33)$$

where  $w_m$  is the weight for model  $m \in \{1, 2, 3\}$ ,  $F1_m$  is the F1-score for model  $m$ ,  $\tau = 10$  controls the scaling steepness to emphasise top-performing models, and  $M = 3$  is the total number of architectures.

Since negative regions in the discriminative heatmap indicate greater twin/non-twin separation, the ensemble applies negation to convert separation strength into positive importance values:

$$H_{ensemble} = - \sum_{m=1}^M w_m \cdot H_{disc,m} \quad (3.34)$$

where  $H_{disc,m}$  is the discriminative heatmap from model  $m$ . This negation ensures that higher importance values in the ensemble correspond to regions showing greater differentiating power between twin and non-twin pairs, providing an intuitive interpretation where larger values indicate more discriminative brain regions.

### 3.6.4 Medical Format Conversion

The conversion process transforms the ensemble heatmap into standard neuroimaging formats for clinical interpretation and visualisation using subject-specific atlases to maximise anatomical accuracy.

#### Subject-Specific Atlas Generation

The atlas generation begins by mapping surface-based cortical labels into native T1w volume space using ribbon-constrained sampling. This process maps each hemisphere separately, ensuring that labels are only assigned to voxels within the cortical ribbon between white matter and pial surfaces:

```

1 # Map left hemisphere cortical labels to volume
2 wb_command -label-to-volume-mapping cortex_L.label.gii \
3     subject.L.pial.32k_fs_LR.surf.gii \
4     T1w_acpc_dc_restore.nii.gz \
5     lh_cortical_labels.nii.gz \
6     -ribbon-constrained subject.L.white.32k_fs_LR.surf.gii \
7     subject.L.pial.32k_fs_LR.surf.gii
8
9 # Map right hemisphere cortical labels to volume
10 wb_command -label-to-volume-mapping cortex_R.label.gii \
11     subject.R.pial.32k_fs_LR.surf.gii \
12     T1w_acpc_dc_restore.nii.gz \
13     rh_cortical_labels.nii.gz \
14     -ribbon-constrained subject.R.white.32k_fs_LR.surf.gii \
15     subject.R.pial.32k_fs_LR.surf.gii

```

Once both hemispheres are mapped to volumetric space, they are combined into a single cortical parcellation volume. Since the hemispheres are anatomically separate, simple addition safely merges the label volumes without creating conflicts:

```

1 # Combine left and right hemisphere volumes
2 fslmaths lh_cortical_labels.nii.gz -add rh_cortical_labels.nii.gz \
3     cortical_combined.nii.gz

```

Subcortical structures require a different approach since they originate from MNI152 standard space. The pre-extracted subcortical label volume is warped into the subject's native T1w space using the subject-specific nonlinear transformation.

Nearest-neighbour interpolation preserves the discrete integer label values throughout this transformation:

```

1 # Warp subcortical labels from MNI space to native T1w space
2 applywarp --ref=T1w_acpc_dc_restore.nii.gz \
3           --in=mni_subcortical_labels.nii.gz \
4           --warp=standard2acpc_dc.nii.gz \
5           --out=subcortical_warped.nii.gz \
6           --interp=nn --datatype=int

```

Before combining cortical and subcortical parcellations, spatial conflicts must be resolved. Cortical voxels that overlap with subcortical regions are removed to prevent erroneous label summation. This is accomplished by creating an inverted mask of subcortical regions and applying it to the cortical volume:

```

1 # Remove cortical voxels that overlap with subcortical regions
2 fslmaths cortical_combined.nii.gz -mas subcortical_warped.nii.gz \
3       -binv -mul cortical_combined.nii.gz cortical_no_overlap.nii.gz

```

The final step merges the warped subcortical labels with the conflict-resolved cortical labels to create a unified volumetric atlas. Integer data type is explicitly specified to preserve the original label identities from the atlas lookup table:

```

1 # Merge subcortical and non-overlapping cortical labels
2 fslmaths subcortical_warped.nii.gz -add cortical_no_overlap.nii.gz \
3       subject_Glasser_MMP_Combined_Atlas_T1w.nii.gz -odt int

```

This pipeline transforms surface-based cortical parcellations and MNI-space subcortical regions into a unified volumetric atlas in the subject's native T1w space.

### Population-Level Statistical Analysis

The ensemble heatmap is applied to each individual subject atlas to extract region-specific importance values. For each subject  $s$  and brain region  $r$ , the regional importance is computed as:

$$I_{s,r} = \frac{1}{|\mathcal{V}_r^s|} \sum_{v \in \mathcal{V}_r^s} H_{ensemble}(v) \quad (3.35)$$

where  $\mathcal{V}_r^s$  represents the set of voxels belonging to region  $r$  in subject  $s$ 's atlas, and  $|\mathcal{V}_r^s|$  is the volume of that region. This process generates subject-specific regional importance values that account for individual anatomical variations.

Population-level statistics are computed by averaging across all subjects:

$$\bar{I}_r = \frac{1}{N} \sum_{s=1}^N I_{s,r} \quad (3.36)$$

$$\sigma_{I_r} = \sqrt{\frac{1}{N-1} \sum_{s=1}^N (I_{s,r} - \bar{I}_r)^2} \quad (3.37)$$

where  $N$  is the total number of subjects. These aggregated statistics provide the foundation for creating subject-specific visualisation heatmaps while maintaining consistency across the population.

### Subject-Specific Heatmap Generation and Medical Format Conversion

For medical format visualisation, subject-specific brain importance atlas heatmaps are generated by applying the population-level regional statistics to individual anatomical atlases. This process maps the aggregated regional importance values to their corresponding atlas regions, creating a subject-specific 3D volume where each voxel within a brain region receives the population-averaged importance value for that region:

$$H_{subject}(v) = \bar{I}_r \text{ where } v \in \text{Region}_r \quad (3.38)$$

where  $H_{subject}(v)$  represents the importance value at voxel  $v$ ,  $\text{Region}_r$  denotes all voxels belonging to atlas region  $r$ , and  $\bar{I}_r$  is the population-averaged importance value for region  $r$ .

The subject-specific atlas heatmap generation ensures that importance values are properly localised to individual anatomical variations while maintaining consistency with the population-level discriminative patterns. Gaussian smoothing is applied to each subject-specific importance atlas heatmap to reduce pixelation and enhance spatial coherence:

$$H_{smoothed} = \mathcal{G}_{\sigma=2} * H_{subject} \quad (3.39)$$

where  $H_{smoothed}$  is the smoothed subject-specific atlas heatmap,  $\mathcal{G}_{\sigma=2}$  represents the Gaussian kernel with standard deviation  $\sigma = 2$ , and  $*$  denotes the convolution operation.

This smoothing step removes high-frequency artefacts while preserving meaningful spatial patterns within the individual's anatomical framework. The smoothed subject-specific atlas heatmap is converted to standard neuroimaging formats using nibabel for volumetric representation. The volumetric conversion process ensures proper spatial alignment with the original brain anatomy by copying the affine transformation and header information from the reference T1-weighted scan:

```
1 import nibabel as nib
```

```

2
3 # Load reference T1-weighted brain scan
4 example_brain_nii = nib.load('T1w_acpc_dc_restore_brain.nii.gz')
5
6 # Copy affine transformation from reference image
7 def copy_affine(source_img, reference_img):
8     # Create new image with source data but reference affine
9     aligned_img = nib.Nifti1Image(
10         source_img,
11         reference_img.affine,
12         reference_img.header
13     )
14     return aligned_img
15
16 # Apply spatial alignment to smoothed subject heatmap
17 heatmap_aligned_nii = copy_affine(subject_heatmap_smooth, example_brain_nii
    )

```

This alignment process ensures that the subject-specific atlas heatmap maintains proper anatomical correspondence with the original brain scan coordinate system. The resulting volumetric NIfTI file preserves spatial precision while enabling compatibility with standard neuroimaging analysis pipelines.

The aligned subject-specific volumetric atlas is then mapped to cortical surfaces using `wb_command` for detailed visualisation:

```

1 # Map subject-specific heatmap to right hemisphere surface
2 wb_command -volume-to-surface-mapping \
3 subject_heatmap_smooth.nii \
4 subject.R.pial.native.surf.gii \
5 heatmap_surfaceR.func.gii -trilinear
6 # Map subject-specific heatmap to left hemisphere surface
7 wb_command -volume-to-surface-mapping \
8 subject_heatmap_smooth.nii \
9 subject.L.pial.native.surf.gii \
10 heatmap_surfaceL.func.gii -trilinear

```

The trilinear interpolation method ensures smooth projection of volumetric values onto the cortical surface while preserving spatial relationships. The resulting outputs include subject-specific volumetric NIfTI files with proper anatomical alignment and surface-based functional overlays suitable for neuroimaging software.

Quantitative tables provide comprehensive regional analysis with cortical system classifications from the HCP-MMP 1.0 atlas, systematically organising the 28 brain regions identified: 22 cortical systems and 6 subcortical structures, sorted by aggregate importance scores from the ensemble model analysis.

The regional importance metrics represent normalised attribution values quantifying each region's contribution to twin identification decisions, computed as

mean and standard deviation of importance scores across all subjects within each anatomical region. Lower standard deviations indicate consistent importance across the population, while higher values suggest variable contribution patterns. Region activation statistics provide the absolute magnitude of model attention within each area, measured as the sum of attribution values before normalisation.

Volume statistics express each region's spatial extent as mean and standard deviation percentages of total brain volume, enabling assessment of whether model importance correlates with anatomical size or reflects genuine discriminative capacity. This facilitates systematic comparison of genetic influence patterns across functionally distinct brain systems.

### **Statistical Significance Testing**

Regional importance rankings were validated using bootstrap resampling with replacement ( $n = 10,000$  iterations) to generate 95% confidence intervals for each region, providing robust uncertainty quantification.

Statistical significance of subcortical dominance was assessed using Welch's independent samples  $t$ -test comparing mean importance scores between six subcortical systems (thalamus, brainstem, hypothalamus, basal ganglia, cerebellum, limbic) and 22 cortical systems. This system-level approach avoids Type I error inflation from comparing a small number of subcortical regions against the larger pool of cortical regions. Effect sizes were quantified using Cohen's  $d$ , calculated as the standardised mean difference between cortex group means using pooled standard deviation.

## 4 Evaluation

This chapter reports the experimental validation of the framework. It begins with the experimental setup, including hardware and software details. Quantitative evaluation covers performance metrics, ROC analysis, and embedding distance distributions for all backbones. Qualitative evaluation includes dimensionality reduction, heatmap analysis, and regional importance using the HCP-MMP atlas. Ablation studies assess the effects of augmentation on performance and regional patterns. Medical format conversion demonstrates clinical applicability via Connectome Workbench. The chapter concludes by contrasting computational approaches with traditional statistical models, highlighting the shift from variance decomposition to direct morphological pattern analysis in neurogenetics.

### 4.1 Experimental Setup

#### 4.1.1 Hardware and Software

Experiments were conducted on a high-performance workstation with an AMD Ryzen 9 7900X3D 12-core processor (4.40 GHz) and 128 GB RAM. Training utilised an NVIDIA GeForce RTX 4090 with 24 GB VRAM, supporting the memory-intensive 3D CNN architectures and enabling large batch sizes and mixed precision training. The system ran 64-bit Windows with CUDA 12.1 and PyTorch for deep learning operations.

### 4.2 Quantitative Results

#### 4.2.1 Loss Graphs

The loss curves in Figure 4.1 display training and validation triplet losses alongside mean embedding distances for twin and non-twin pairs throughout training. The green line represents mean cosine distances between twin pair embeddings at each epoch, while the red line shows mean distances for non-twin pairs. Successful training is indicated by decreasing twin distances (green trending toward low values) and increasing non-twin distances (red trending toward higher values), creating clear separation between the two distributions.

Notably, embedding quality and separation can continue to improve even when validation loss appears to plateau, as the loss function only enforces margin-based separation rather than optimal clustering within each class, allowing for continued

Table 4.1 Siamese Network Training Parameters

Parameter	Value	Description
<b>Training Configuration</b>		
Epochs	2000	Total training iterations
Triplet Loss Margin	0.1	Minimum separation distance
L2 Regularisation	0.01	Embedding regularisation weight
Gradient Clipping	0.5	Maximum gradient norm
<b>Optimiser (AdamW)</b>		
Learning Rate	$1 \times 10^{-4}$	Base learning rate
Beta 1	0.9	First momentum coefficient
Beta 2	0.999	Second momentum coefficient
Weight Decay	$5 \times 10^{-4}$	Decoupled weight decay
<b>Learning Rate Scheduler (OneCycleLR)</b>		
Maximum LR	$1 \times 10^{-4}$	Peak learning rate
Warmup Period	30%	Percentage for gradual increase
<b>Batch Configuration</b>		
Effective Batch Size	512	Large batches for triplet mining
U-Net/ResNet Accumulation	2 steps	Gradient accumulation steps
DenseNet Accumulation	8 steps	Higher memory requirement
<b>Data Augmentation</b>		
Augmentation Probability	0.5	Per-transformation application rate
Random Rotation Range	-20° to 20°	Multi-plane rotation angles
Random Flipping	Random axes	Left-right symmetry invariance
Gaussian Noise Std	0.02	Maximum noise standard deviation
Intensity Scaling	0.9 to 1.1	Signal intensity normalisation
<b>Initialisation &amp; Precision</b>		
Embedding Init	Xavier Uniform	Weight initialisation method
Xavier Gain	0.5	Reduced gain for stability
Precision	FP16 (Mixed)	Memory optimisation
Cache Clear Frequency	Every 3 epochs	Memory management

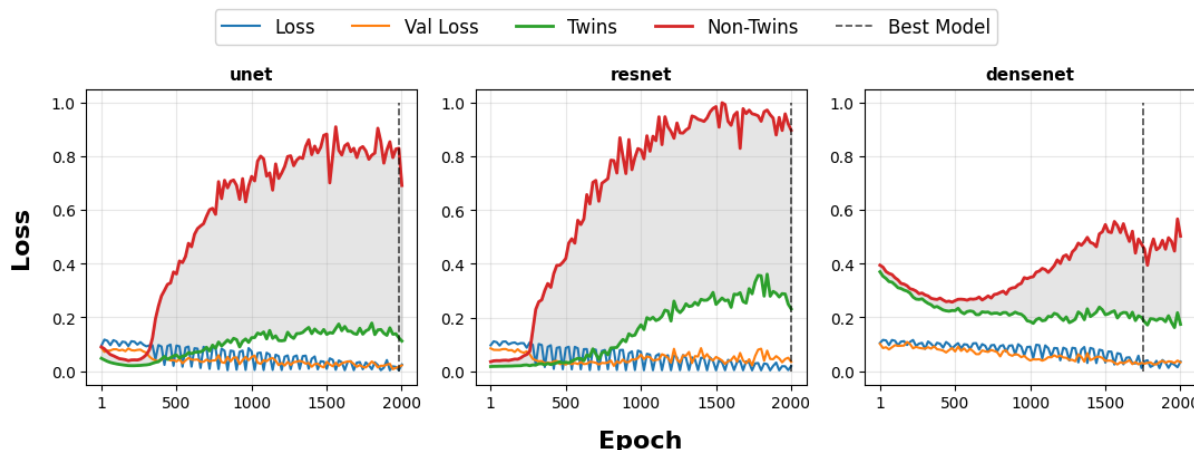


Figure 4.1 Training and validation loss curves showing triplet loss convergence and embedding distance separation across 2000 epochs. The twin/non-twin separation lines represent average embedding distances for twin pairs versus non-twin pairs across all pairs at each epoch. U-Net demonstrates most stable convergence with minimal overfitting and clear bimodal separation.

refinement of the embedding space geometry beyond the minimum margin requirement. This is particularly noticeable at the end of the training graph for ResNet.

U-Net demonstrates the most stable convergence with training and validation losses closely aligned, indicating minimal overfitting and robust learning. The twin (green) and non-twin (red) pair distances show clear separation by epoch 500, with the twin distances stabilising around 0.15 and non-twin distances reaching approximately 0.85. The best model selected is indicated by the vertical dashed line near epoch 2000.

ResNet exhibits similar convergence patterns but with slightly higher variability in validation loss during later epochs. The twin/non-twin separation emerges gradually, with stable separation achieved around epoch 1000. Twin distances stabilise around 0.25 while non-twin distances reach approximately 0.95, indicating effective but slightly less compact embedding space compared to U-Net.

DenseNet shows the most challenging training dynamics with higher overall loss values and continued oscillation throughout training. The twin and non-twin pair distances demonstrate less distinct separation, with twin distances fluctuating around 0.2 and non-twin distances reaching only approximately 0.5. This behaviour aligns with DenseNet's consistently lower performance metrics and suggests difficulty in learning discriminative embeddings for the twin identification task.

## 4.2.2 Performance Metrics

The confusion matrices in Figure 4.2 represent cumulative results from all 10 evaluation runs, showing the models' overall classification patterns. U-Net achieved the highest performance with 274 correctly identified twin pairs and 238 correctly

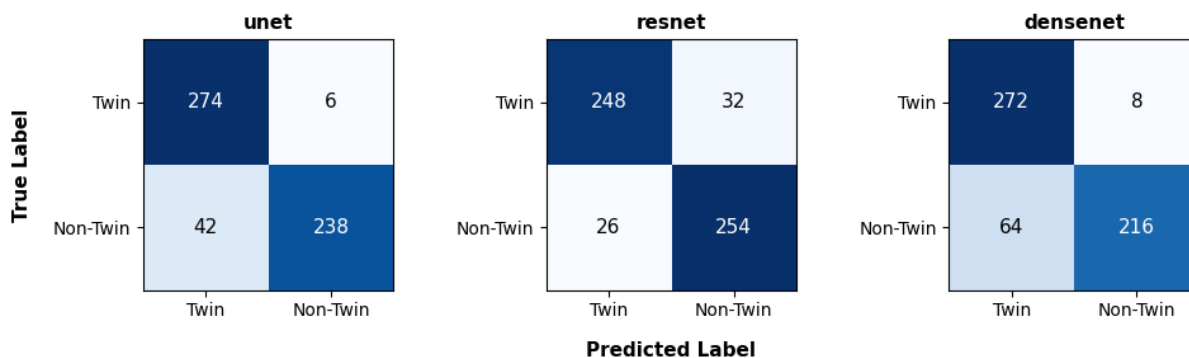


Figure 4.2 Confusion matrices for ResNet, U-Net, and DenseNet architectures on twin identification task across 10 evaluation runs.

identified non-twin pairs, resulting in only 48 total misclassifications. ResNet showed 248 correct twin identifications and 254 correct non-twin identifications with 58 total errors, while DenseNet achieved 272 correct twin identifications but only 216 correct non-twin identifications with 72 total errors.

Table 4.2 Performance metrics comparison across architectures showing mean  $\pm$  standard deviation across 10 evaluation runs with different randomly generated combinations of non-twin pairs to assess robustness to negative sample selection.

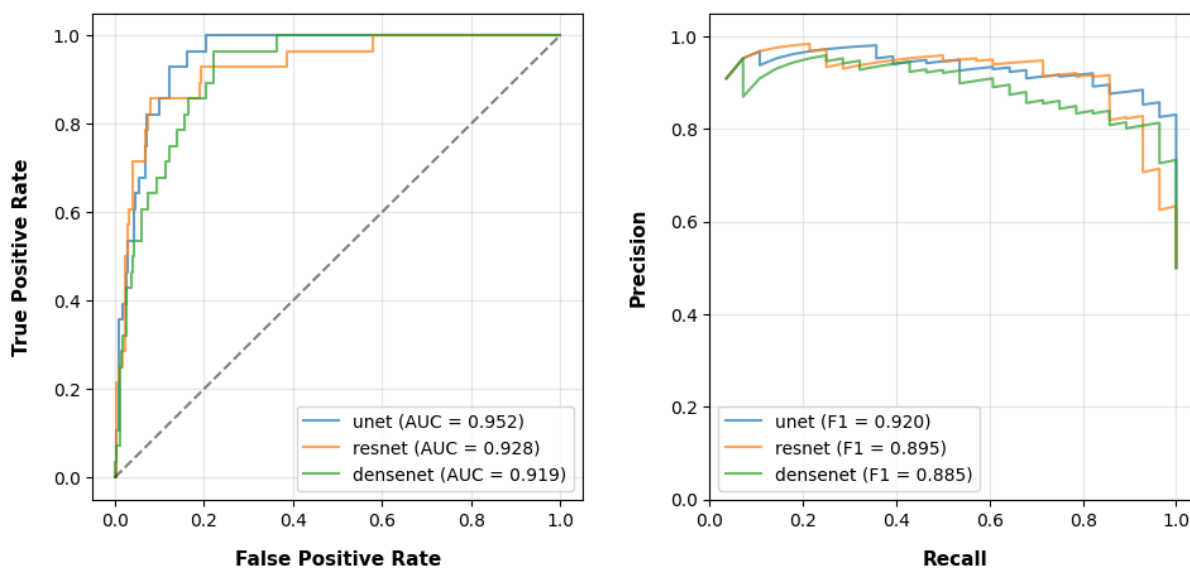
Metric	U-Net	ResNet	DenseNet
Accuracy	<b>91.4 <math>\pm</math> 2.9</b>	89.6 $\pm$ 2.2	87.1 $\pm$ 4.8
Precision	87.1 $\pm$ 5.1	<b>90.8 <math>\pm</math> 4.7</b>	81.6 $\pm$ 6.5
Recall	<b>97.9 <math>\pm</math> 2.9</b>	88.6 $\pm$ 3.5	97.1 $\pm$ 1.4
Specificity	85.0 $\pm$ 6.6	<b>90.7 <math>\pm</math> 5.1</b>	77.1 $\pm$ 10.4
F1	<b>92.0 <math>\pm</math> 2.5</b>	89.6 $\pm$ 2.1	88.5 $\pm$ 3.6
AUC	<b>95.2 <math>\pm</math> 2.3</b>	92.8 $\pm$ 2.4	91.9 $\pm$ 3.1
NPV	<b>97.7 <math>\pm</math> 2.9</b>	89.0 $\pm$ 2.9	96.6 $\pm$ 1.8

U-Net provides the most balanced and stable performance across all metrics, achieving the highest F1-score (92.0%,  $\sigma = 2.5\%$ ) and AUC (95.2%,  $\sigma = 2.3\%$ ). The consistently low standard deviations across metrics indicate robust generalisation capabilities and reduced sensitivity to non-twin pair variations.

ResNet demonstrates competitive performance with the highest precision (90.8%,  $\sigma = 4.7\%$ ) and specificity (90.7%,  $\sigma = 5.1\%$ ), indicating superior accuracy in identifying non-twin pairs. However, ResNet shows lower recall (88.6%) compared to U-Net, suggesting reduced sensitivity for identifying twin pairs.

DenseNet shows consistently lower performance across most metrics with higher variability, particularly in specificity (77.1%,  $\sigma = 10.4\%$ ), indicating challenges in distinguishing non-twin pairs. Despite achieving high recall (97.1%), the model's lower precision (81.6%) results in more false positive classifications.

### 4.2.3 ROC and Precision-Recall Analysis



(a) ROC curves comparing model performance across architectures. Architectures sorted by AUC-ROC.

(b) PR curves showing threshold sensitivity. Architectures sorted by F1-score.

Figure 4.3 ROC and PR curve evaluation across model architectures.

The ROC curves in Figure 4.3a demonstrate strong discriminative performance across all architectures. U-Net achieves the highest AUC-ROC (95.2%), followed by ResNet (92.8%) and DenseNet (91.9%). All models substantially outperform random classification (50%), with U-Net and ResNet showing particularly steep initial curves indicating effective separation of twin and non-twin pairs at conservative thresholds.

The PR curves in Figure 4.3b reveal architecture-specific characteristics. U-Net maintains the highest precision across most recall values, achieving optimal F1-score (92.0%) through superior precision-recall balance. ResNet shows competitive performance with F1-score (89.5%), while DenseNet exhibits earlier precision degradation as recall increases, resulting in lower F1-score (88.5%).

U-Net demonstrates the most robust performance across thresholds, maintaining high precision even at elevated recall levels. The curves indicate that U-Net provides optimal balance for twin identification tasks, while ResNet offers competitive performance with slightly reduced precision at higher recall values. DenseNet shows the steepest precision decline, indicating reduced reliability at higher recall thresholds.

### 4.2.4 Embedding Distance Distribution Analysis

Figure 4.4 reveals distinct separation characteristics across architectures. U-Net achieves the clearest bimodal distribution with twin pairs strongly concentrated at low

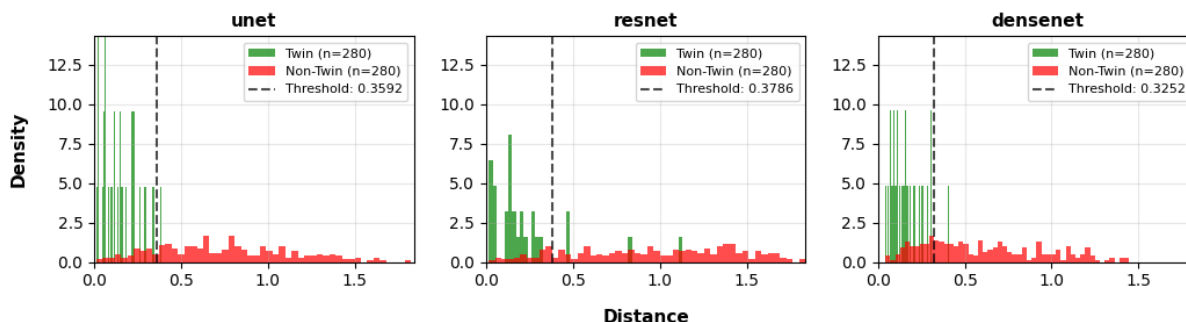


Figure 4.4 Embedding distance distributions demonstrating clear bimodal separation between twin and non-twin pairs. U-Net achieves optimal separation with minimal overlap, while DenseNet shows substantial distribution overlap.

distances (peak near 0.1) and minimal overlap with non-twin pairs. The optimal threshold of 0.3592 effectively separates the two distributions, supporting its superior classification performance.

ResNet demonstrates strong twin clustering with a sharp peak near zero distance, but exhibits broader distribution overlap in intermediate ranges. The threshold of 0.3786 reflects moderate separation between twin and non-twin clusters, indicating reasonably compact embedding space with some overlapping regions.

DenseNet shows the poorest separation quality with substantial overlap between twin and non-twin distributions across most distance ranges. The threshold of 0.3252 indicates difficulty in establishing clear decision boundaries, with twin and non-twin distributions showing significant overlap around the 0.2-0.6 distance range.

The distribution analysis confirms U-Net’s superior embedding quality through the most compact twin clustering and clearest class separation, directly supporting its highest F1-score and balanced precision-recall performance. The clear bimodal separation in U-Net’s distribution explains its robust performance across different threshold values.

## 4.3 Qualitative Results

### 4.3.1 Dimensionality Reduction

t-SNE visualisation of the 128-dimensional embeddings provides qualitative assessment of the learned embedding space structure. Figure 4.5 displays the embedding space across all three architectures, where eight representative twin and non-twin pairs per model are connected by dashed lines and colour-coded by relationship type. A perplexity of 10 was selected to balance local and global structure preservation given the small sample size, with a fixed random seed of 42 ensuring

reproducibility.

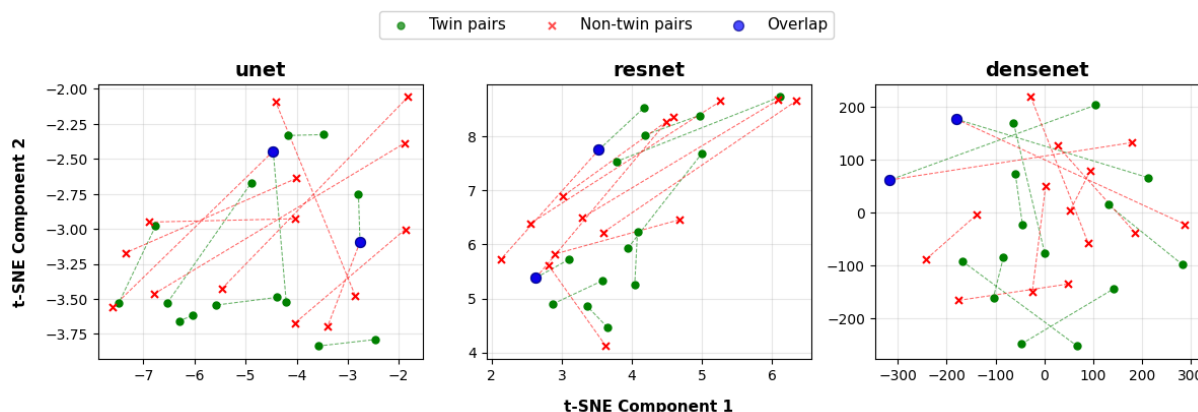


Figure 4.5 t-SNE visualisation of embedding space across model architectures showing representative twin and non-twin pairs. Green dots represent twin pairs connected by dashed lines, red crosses indicate non-twin pairs, and blue circles show overlap points where subjects have both twin and non-twin relationships in the visualisation.

The visualisations demonstrate wide distribution of points across the embedding space, with each subject occupying distinct regions reflecting their neuroanatomical characteristics. Twin pairs exhibit clear clustering behaviour, with connected green points positioned closer together than non-twin pairs.

Blue overlap points, appearing where subjects have both twin and non-twin relationships within the visualised sample, provide additional validation. In these cases, the corresponding twin consistently lies closer than the non-twin pair, supporting the quantitative finding that learned distance metrics prioritise genetic similarity. However, these clustering patterns should be interpreted in conjunction with the quantitative metrics rather than as definitive spatial relationships, given t-SNE's sensitivity to hyperparameters and its non-preservation of global distances.

U-Net achieves compact clustering with clear twin-pair groupings, aligning with its superior quantitative performance (F1-score: 92.0%, AUC-ROC: 95.2%). ResNet displays similar clustering quality with slightly broader distribution, consistent with its competitive performance (F1-score: 89.5%, AUC-ROC: 94.1%). DenseNet exhibits more dispersed clustering with greater overlap between twin and non-twin pairs, supporting its lower classification metrics (F1-score: 88.5%, AUC-ROC: 91.8%) and the embedding distance distributions presented previously.

### 4.3.2 Reference Attribution Patterns

Figure 4.6 shows that despite distinct architectural approaches, all models consistently focus on the same brain regions. The reference LRP represents mean attribution patterns across all subjects for each embedding dimension, as described in

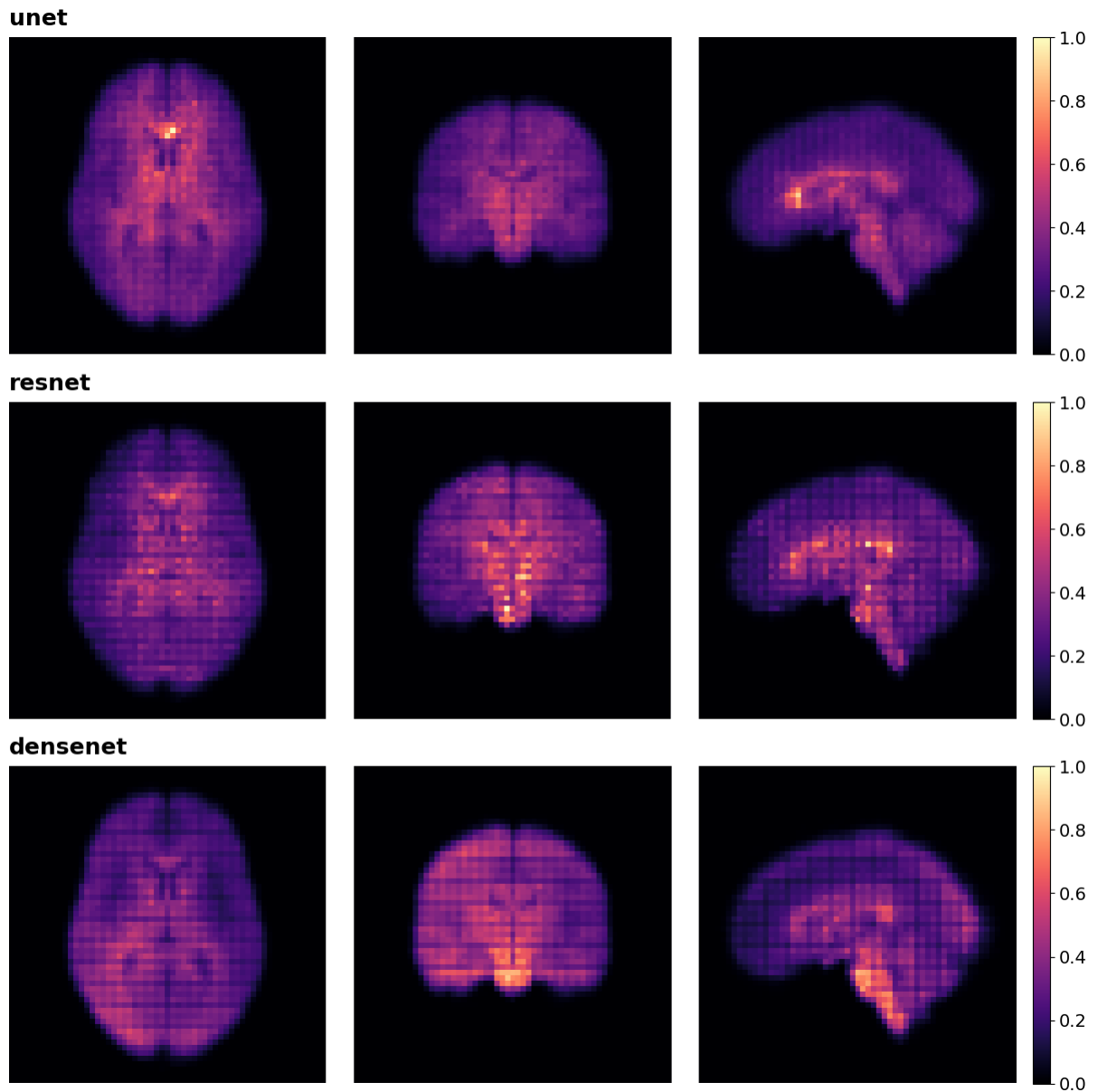


Figure 4.6 Reference LRP maps averaged across all subjects and embedding dimensions for each architecture as described in Equation (3.29) (left: axial, middle: coronal, right: sagittal views). All models converge on similar brain regions with consistent focus on subcortical structures and brainstem. Colour scale (0.0-1.0) represents normalised LRP attribution values, where higher values indicate regions contributing more to embedding representation.

## Section 3.6.2.

All three architectures converge on thalamic structures (importance as per Table 4.3: 0.955), brainstem (0.875), and hypothalamus (0.707), indicating these regions contain the most discriminative features for twin identification. The consistent attention to posterior cingulate cortex (PCC, 0.546) and early auditory cortex (EAC, 0.490) across models validates their importance for genetic similarity detection. While embedding quality varies between models, the relevant neuroanatomical regions remain constant.

### 4.3.3 Embedding Space Uniqueness

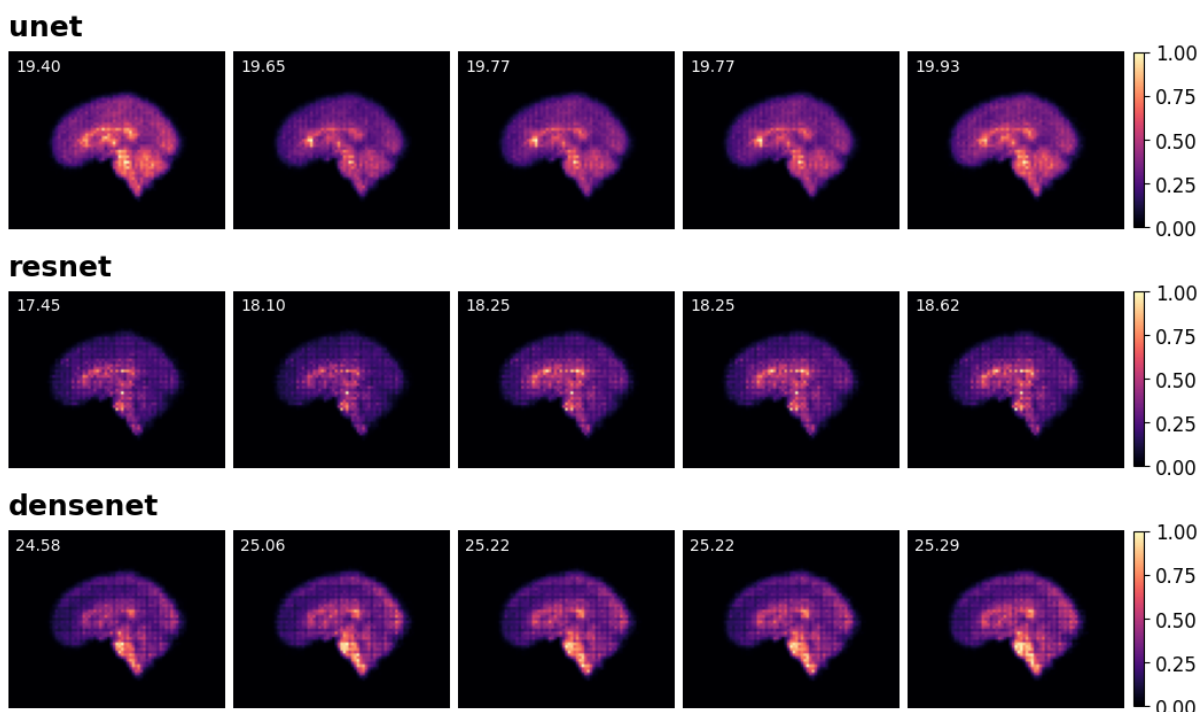


Figure 4.7 Peak Signal-to-Noise Ratio analysis of embeddings with maximum deviation from reference LRP patterns. PSNR of each embedding denoted in white text. High similarity across architectures indicates models dedicate embedding dimensions to subtle structural variations within consistent regions.

The Peak Signal-to-Noise Ratio (PSNR) analysis in Figure 4.7 examines embeddings that deviate most from reference patterns to assess embedding space utilisation. Even the most unique subject representations maintain focus on the same core brain regions identified in the reference maps.

The high similarity between architectures indicates that models allocate embedding dimensions to capture minor structural changes within established relevant areas rather than exploring different brain regions. This supports the hypothesis that

genetic similarity manifests through subtle morphological variations in consistent neuroanatomical structures.

#### 4.3.4 Discriminative Heatmaps

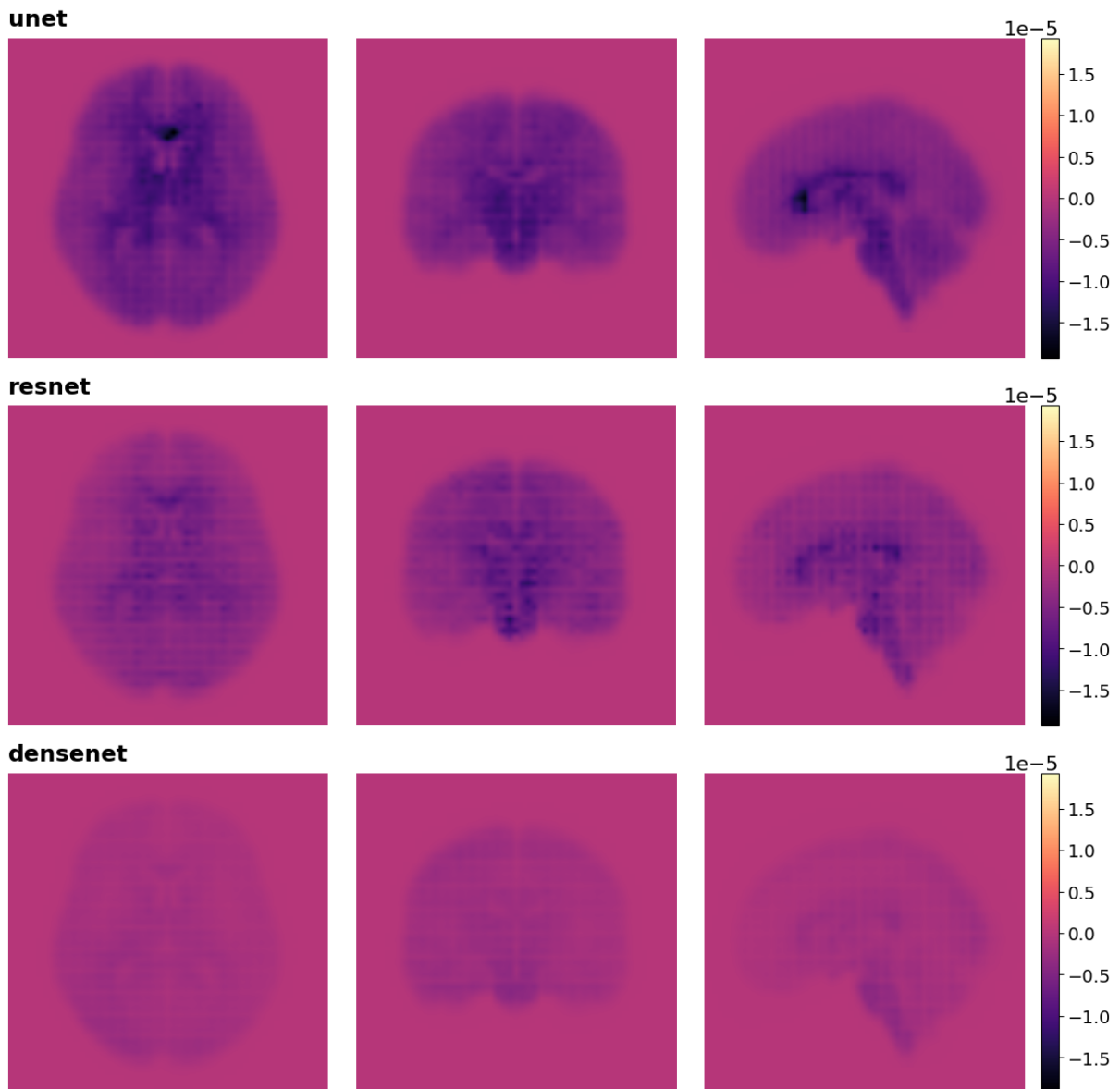


Figure 4.8 Architecture-specific discriminative heatmaps showing differential embedding patterns between twin and non-twin pairs as described in Equation (3.32). Darker regions ( $< 0$ ) indicate larger embedding distances in non-twin pairs; lighter regions ( $> 0$ ) indicate larger embedding distances in twin pairs. U-Net and ResNet demonstrate similar patterns, while DenseNet shows weaker differentiation. Colour scale represents embedding distance differences, where negative values indicate regions with greater twin/non-twin discriminative capacity.

The discriminative heatmap implementation follows the procedure detailed in Section 3.6.2. Figure 4.8 reveals striking similarities between U-Net and ResNet

discriminative patterns. Both architectures identify consistent regions where genetic similarity drives embedding proximity, with comparable discrimination strength across identified regions.

The regions showing the strongest discrimination between twin and non-twin pairs align closely with areas of high activation in the reference LRP heatmaps, demonstrating model coherence. This alignment indicates that the models converge on brain regions that exhibit minimal variability between twins while showing substantial differences between non-twins (precisely the neuroanatomical patterns most informative for genetic relatedness identification).

Strong negative values in subcortical areas suggest these regions contribute substantially to embedding similarity in genetically related individuals. DenseNet's discriminative heatmap shows substantially weaker discrimination patterns with reduced signal strength throughout the brain, directly reflecting its lower quantitative performance (F1-score: 88.5%). The diminished contrast and spatial coherence validate the embedding distance distribution analysis findings, demonstrating that DenseNet achieves less effective separation between twin and non-twin pairs.

### 4.3.5 Ensemble Integration

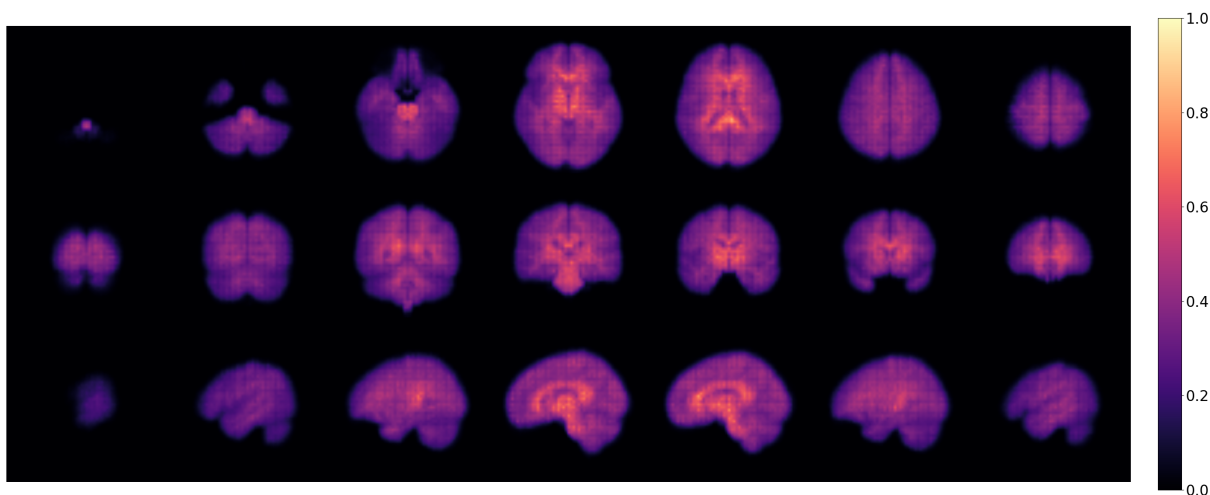


Figure 4.9 Ensemble heatmap combining weighted contributions from all three architectures as described in Equation (3.34), showing enhanced spatial coherence in subcortical regions. Top row: axial slices; middle row: coronal slices; bottom row: sagittal slices. Colour scale (0.0-1.0) represents discriminative importance for genetic relatedness detection, where higher values indicate regions that contribute more to distinguishing genetically related individuals.

The ensemble approach follows the methodology in Section 3.6.3. Performance weighting (U-Net: 0.41, ResNet: 0.31, DenseNet: 0.28) ensures superior models contribute proportionally more while maintaining information from all architectures.

The ensemble heatmap in Figure 4.9 shows enhanced spatial coherence and stronger discrimination signals compared to individual heatmaps.

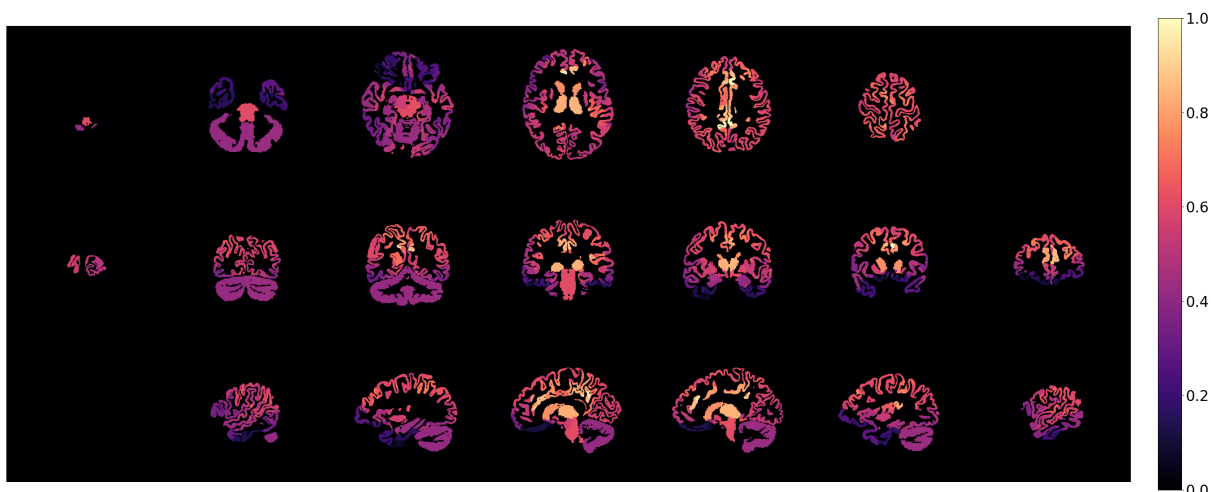


Figure 4.10 Subject-specific atlas heatmap demonstrating clinical applicability for personalised twin similarity assessment, generated by applying population-level regional statistics to individual anatomical parcellation. Top row: axial slices; middle row: coronal slices; bottom row: sagittal slices. Colour scale (0.0-1.0) represents regional importance scores for genetic similarity, where higher values indicate anatomical regions with greater discriminative capacity for twin identification.

Figure 4.10 demonstrates clinical translation by mapping ensemble results to individual subject anatomy using the atlas generation process described in Section 3.6.4. The subject-specific atlas heatmap maintains spatial precision while providing interpretable regional importance scores within individual anatomical contexts.

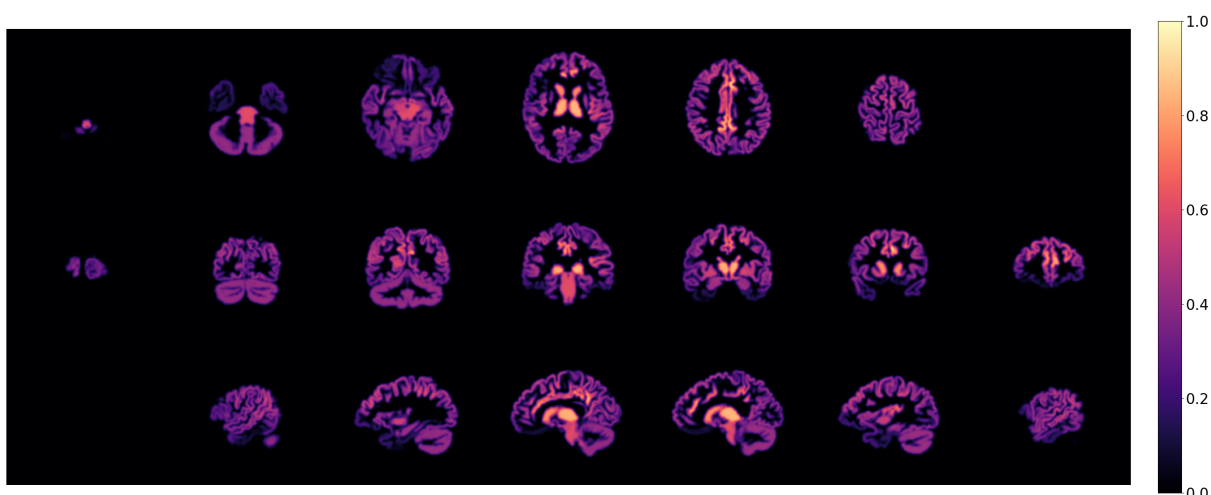


Figure 4.11 Gaussian-smoothed ( $\sigma = 2$ ) version of the subject-specific atlas heatmap in Figure 4.10, optimised for medical format conversion. Smoothing preserves spatial patterns while reducing noise for integration with standard neuroimaging pipelines.

The smoothed atlas heatmap in Figure 4.11 applies the Gaussian smoothing

procedure outlined in Section 3.6.4, providing the foundation for medical format conversion while preserving meaningful spatial discrimination patterns for clinical interpretation.

### 4.3.6 MMP Glasser Brain Regions

#### Brain Regions

Table B.1 provides a comprehensive listing of all 379 individual brain regions ranked by discriminative importance with their cortical/subcortical assignments, activation levels, and volumetric contributions, while Table 4.3 presents aggregate statistics across the 22 cortical areas and 6 subcortical structures to reveal the relative ranking of neuroanatomical discriminative importance for genetic relatedness detection.

Table 4.3 Relative ranking of parcellated cortices and subcortical structures (bold) from HCP-MMP 1.0 atlas sorted by discriminative importance for genetic relatedness detection. Bootstrap confidence intervals (95% CI, n=10,000) demonstrate ranking stability, with subcortical structures showing significantly higher discriminative importance than cortical regions. Subcortical group mappings available in Table A.2.

Cortex/Structure	Region Importance				Region Activation		Region Volume
	Count	Mean	Std	95% CI	Mean	Std	Sum (%)
<b>THALAMUS</b>	2	<b>0.955</b>	0.064	[0.910, 1.000]	11999	827	2.10
<b>BRAINSTEM</b>	1	<b>0.875</b>	-	[0.874, 0.876]	31595	-	3.02
<b>HYPOTHALAMUS</b>	2	<b>0.707</b>	0.067	[0.660, 0.755]	5316	390	1.26
<b>BASAL GANGLIA</b>	8	<b>0.585</b>	0.128	[0.498, 0.660]	3349	2425	3.59
Posterior Cingulate	26	<b>0.546</b>	0.141	[0.496, 0.601]	796	414	3.25
<b>CEREBELLUM</b>	2	<b>0.545</b>	0.036	[0.519, 0.570]	55808	2219	17.14
<b>LIMBIC</b>	4	<b>0.516</b>	0.077	[0.462, 0.591]	3206	1501	2.03
Temporo-Parieto Occipital Junction	1	<b>0.494</b>	-	[0.493, 0.495]	1203	-	0.20
Early Auditory	14	<b>0.490</b>	0.095	[0.444, 0.540]	613	355	1.42
Temporo-Parieto-Occipital Junction	9	<b>0.488</b>	0.084	[0.437, 0.539]	1043	375	1.60
Superior Parietal	20	<b>0.476</b>	0.051	[0.454, 0.497]	646	327	2.29
Inferior Parietal	20	<b>0.474</b>	0.059	[0.449, 0.498]	1511	821	5.29
Ventral Stream Visual	14	<b>0.464</b>	0.087	[0.422, 0.509]	833	585	2.06
MT+ Complex and Neighbouring Visual Areas	18	<b>0.463</b>	0.057	[0.437, 0.488]	633	370	2.13
Paracentral Lobular and Mid Cingulate	18	<b>0.461</b>	0.100	[0.417, 0.507]	1032	408	3.37
Insular and Frontal Opercular	24	<b>0.444</b>	0.105	[0.402, 0.484]	622	319	2.82
Dorsal Stream Visual	12	<b>0.440</b>	0.054	[0.414, 0.472]	537	219	1.26
Posterior Opercular	10	<b>0.431</b>	0.142	[0.349, 0.517]	654	323	1.36
Early Visual	6	<b>0.402</b>	0.009	[0.395, 0.408]	2656	599	3.31
Medial Temporal	16	<b>0.401</b>	0.190	[0.314, 0.495]	592	467	2.78
Dorsolateral Prefrontal	26	<b>0.377</b>	0.058	[0.356, 0.399]	1145	448	6.65
Somatosensory and Motor	10	<b>0.375</b>	0.047	[0.347, 0.402]	1867	1113	4.04
Auditory Association	16	<b>0.372</b>	0.099	[0.325, 0.419]	946	372	3.48
Primary Visual	2	<b>0.370</b>	0.007	[0.365, 0.375]	3804	136	1.72
Anterior Cingulate and Medial Prefrontal	32	<b>0.363</b>	0.122	[0.322, 0.405]	747	667	5.15
Premotor	14	<b>0.342</b>	0.066	[0.307, 0.374]	812	417	2.75
Inferior Frontal	18	<b>0.285</b>	0.082	[0.247, 0.322]	639	282	3.54
Lateral Temporal	16	<b>0.240</b>	0.101	[0.193, 0.288]	1102	733	6.35
Orbital and Polar Frontal	18	<b>0.212</b>	0.097	[0.171, 0.257]	528	281	4.01

The regional rankings in Table 4.3 reveal a clear discriminative hierarchy with

subcortical structures dominating the top positions. Six subcortical structures occupy the top 7 positions, led by thalamus (0.955), brainstem (0.875), and hypothalamus (0.707), with only posterior cingulate cortex (5th, 0.546) interrupting this dominance. Bootstrap confidence intervals ( $n = 10,000$  resamples) confirm this ranking stability, with the three highest-ranking regions maintaining consistent positions: thalamus ( $M = 0.955$ , 95% CI [0.910, 1.000]), brainstem ( $M = 0.875$ , 95% CI [0.874, 0.876]), and hypothalamus ( $M = 0.707$ , 95% CI [0.660, 0.755]).

Statistical analysis validates this observed subcortical dominance pattern at the cortex group level. Subcortical cortex groups ( $M_{subcortical} = 0.697$ ,  $n = 6$  cortex groups) demonstrated significantly higher mean discriminative importance compared to cortical cortex groups ( $M_{cortical} = 0.409$ ,  $n = 23$  cortex groups), yielding a large effect size and confidence (Cohen's  $d = 2.80$ ,  $t(27) = 5.77$ ,  $p = 3.89 \times 10^{-6}$ ). Notably, the deep learning models utilise practically all brain regions for twin identification, with importance scores showing a compressed distribution (Mean: 0.467, Median: 0.453, all cortices  $> 0.2$ ), indicating distributed multivariate processing of morphological features rather than excessive dependence on regional selection.

This computational approach from direct MRI analysis presents a fundamentally different perspective compared to traditional ACE heritability studies. While ACE models measure independent regional genetic influences through twin comparisons, our deep learning framework learns which brain regions collectively provide the strongest morphological signatures for distinguishing genetically identical individuals from structural imaging data.

The resulting hierarchy reveals subcortical dominance despite ACE studies reporting higher heritability in cortical areas (frontal 78-95%, temporal 77-89%) compared to subcortical structures (thalamus 42%, cerebellum 24%). This methodological contrast suggests that high individual regional heritability does not necessarily translate to discriminative utility for computational twin identification.

The deep learning models prioritise subcortical structures containing the most reliable patterns for genetic relatedness detection, complementing traditional heritability analysis by revealing how combinations of brain regions enable direct identification from neuroimaging data rather than statistical decomposition of genetic and environmental variance.

### 4.3.7 Medical Format Conversion

The medical format conversion successfully transforms the computational twin identification results into standard NIfTI and GIFTI formats compatible with clinical neuroimaging workflows. The conversion process demonstrates seamless integration with Connectome Workbench, the neuroimaging analysis platform developed by the

Human Connectome Project, validating the practical utility of the computational pipeline for clinical and research applications.

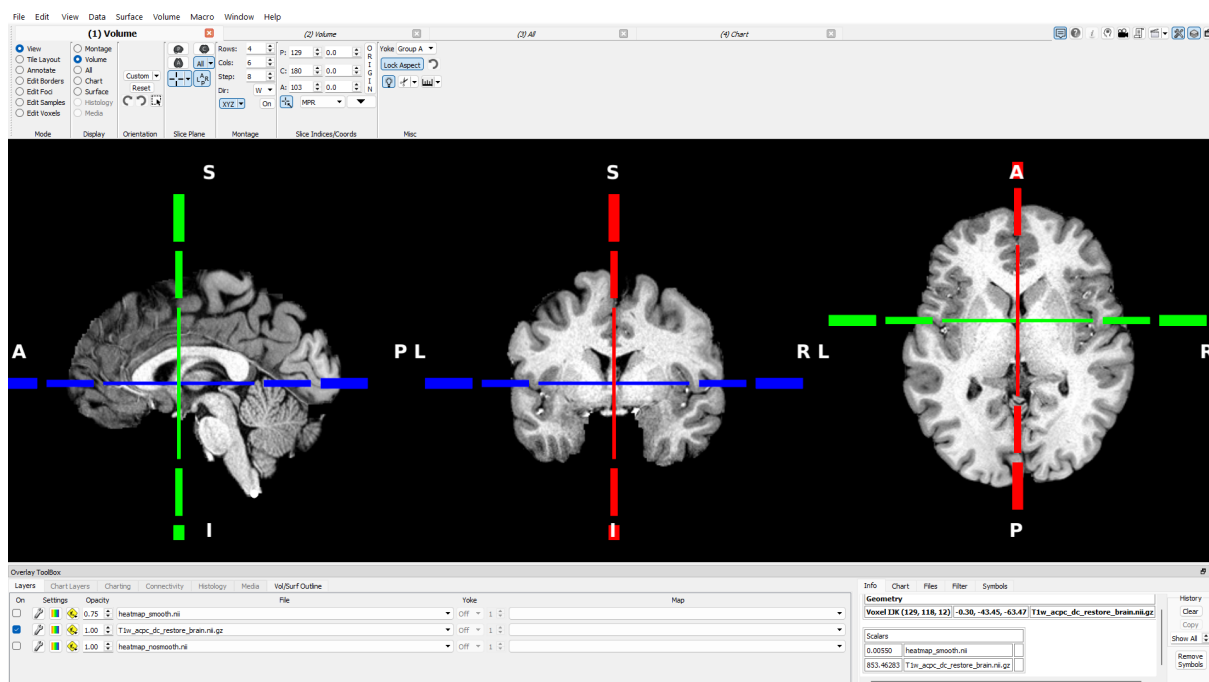


Figure 4.12 Subject-specific T1-weighted anatomical volume displayed in Connectome Workbench, showing sagittal, coronal, and axial views used as the reference space for atlas heatmap alignment and volumetric format conversion.

Figure 4.12 establishes the anatomical reference space using a subject's T1-weighted structural scan, providing the coordinate system foundation for accurate heatmap alignment. The three-plane volumetric display demonstrates Connectome Workbench's native capability to handle the converted data, with proper spatial registration and coordinate system preservation throughout the conversion process.

The comparison between unsmoothed and smoothed atlas heatmaps reveals the critical importance of spatial processing for clinical visualisation in standard neuroimaging environments. The raw regional importance values, when converted to volumetric format, maintain anatomical precision while enabling integration with established clinical workflows.

Figure 4.13 demonstrates how the unprocessed volumetric conversion displays discrete regional boundaries with pixelation artefacts, while Figure 4.14 shows the effectiveness of Gaussian smoothing in creating clinically-appropriate volumes. Connectome Workbench's volume rendering capabilities properly interpret the converted data, with accurate intensity scaling and spatial navigation confirming successful format compliance.

Surface-based visualisation leverages Connectome Workbench's specialised functionality for cortical analysis, converting the heatmap data into functional surface

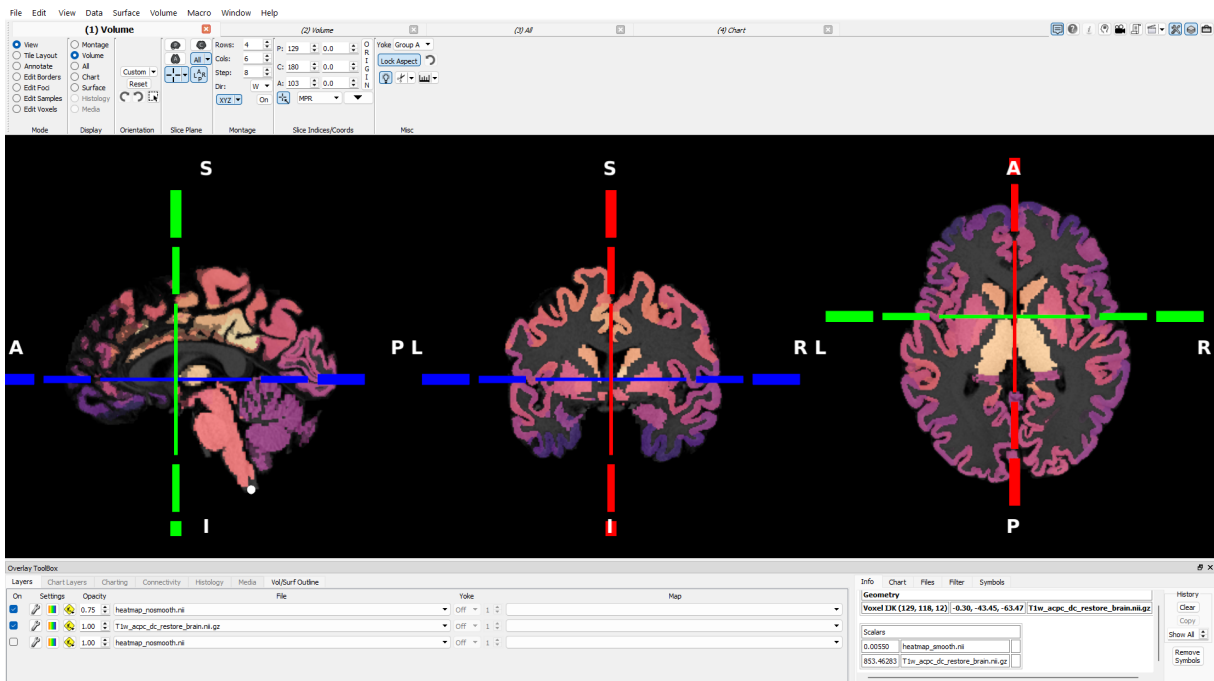


Figure 4.13 Unsmoothed subject-specific atlas heatmap overlaid on T1-weighted anatomy in Connectome Workbench, displaying raw regional importance values with visible pixelation artefacts before Gaussian smoothing optimisation.

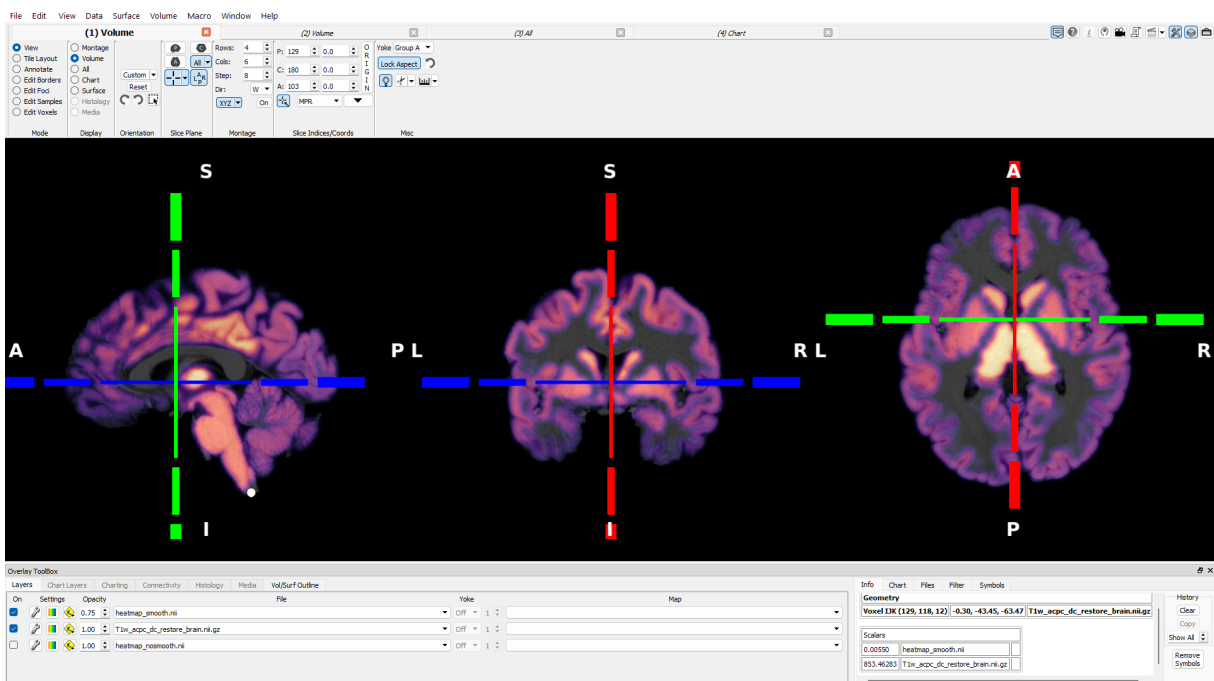


Figure 4.14 Gaussian-smoothed subject-specific atlas heatmap ( $\sigma = 2$ ) overlaid on anatomical volume in Connectome Workbench, demonstrating enhanced spatial coherence and improved visualisation quality suitable for clinical interpretation.

overlays that integrate seamlessly with the platform's surface visualisation pipeline.

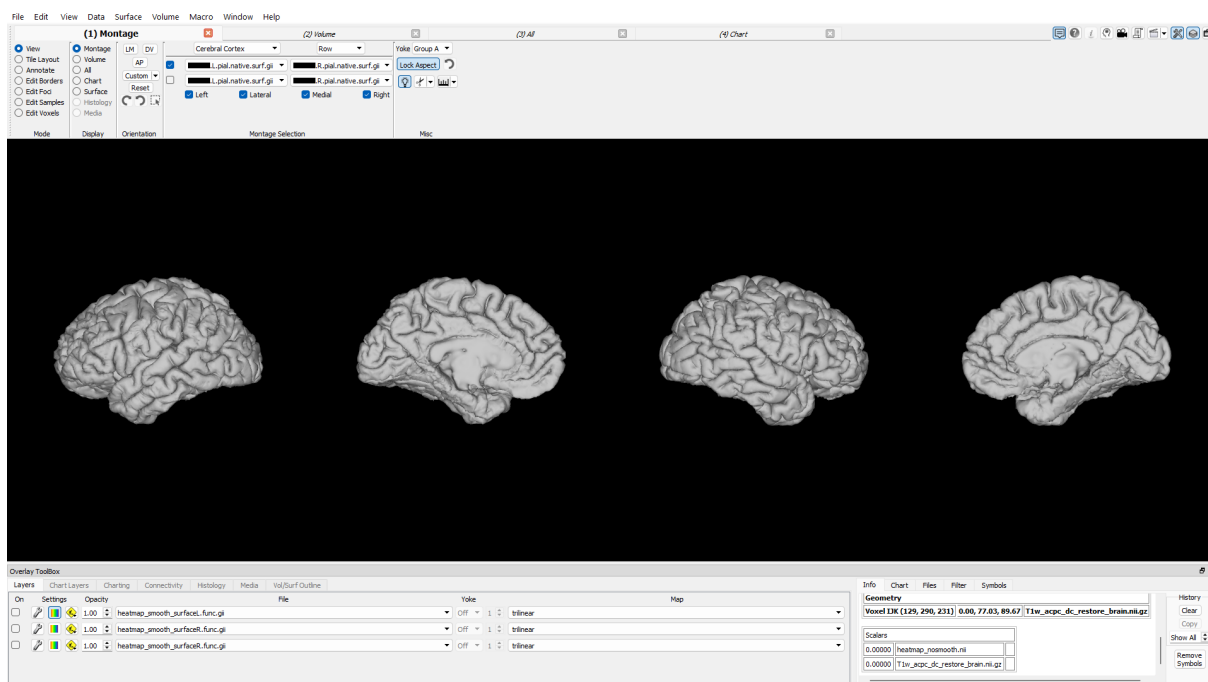


Figure 4.15 Subject-specific cortical surface reconstruction displayed in Connectome Workbench, showing bilateral hemisphere views of the native anatomical mesh used for surface-based heatmap projection.

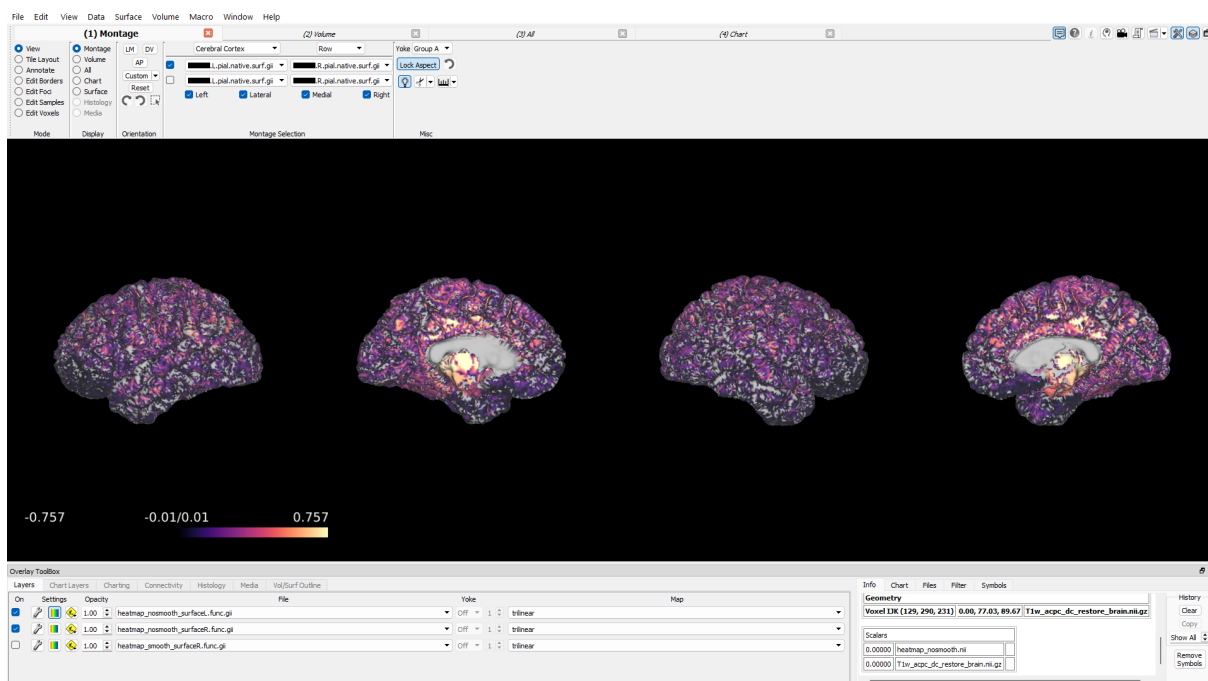


Figure 4.16 Unsmoothed atlas heatmap projected onto cortical surfaces in Connectome Workbench using trilinear interpolation, showing raw importance values with high spatial frequency artefacts and discontinuous regional boundaries.

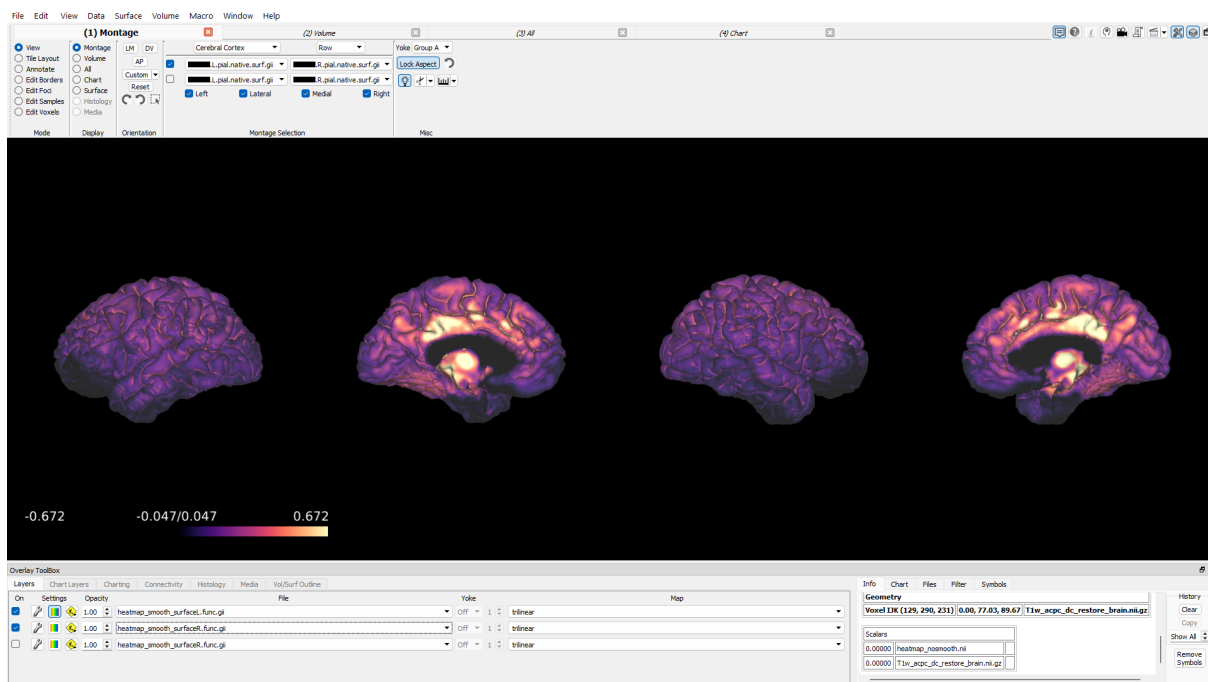


Figure 4.17 Gaussian-smoothed atlas heatmap projected onto cortical surfaces in Connectome Workbench, demonstrating improved spatial continuity with enhanced visualisation quality and clear discrimination of high-importance regions.

Figure 4.15 shows the high-quality cortical surface reconstruction that serves as the target for functional overlay projection. The surface mapping process demonstrates successful conversion to standard format, enabling Connectome Workbench's native surface visualisation capabilities. The comparison between unsmoothed (Figure 4.16) and smoothed (Figure 4.17) surface projections validates the conversion pipeline's preservation of spatial relationships while optimising for clinical interpretation.

The integration with Connectome Workbench confirms that the converted files meet standard neuroimaging format specifications. The platform's layer management system properly recognises the file formats, coordinate systems, and data ranges, while interactive navigation tools function correctly with the converted datasets. The colour mapping and intensity scaling demonstrate that the converted files maintain quantitative accuracy suitable for clinical analysis and research applications.

This successful format conversion bridges the gap between computational twin identification algorithms and established clinical neuroimaging workflows. The demonstrated compatibility with Connectome Workbench validates the practical utility of the approach for real-world applications, enabling clinicians and researchers to integrate these computational insights into existing analysis pipelines using familiar, standardised neuroimaging formats.

## 4.4 Ablation Study

### 4.4.1 Augmentation

The ablation study compares three training configurations across all architectures: augmented training for 2,000 epochs (AUG2K), non-augmented training for 2,000 epochs (NOAUG2K), and extended non-augmented training for 5,000 epochs (NOAUG5K). This analysis evaluates the impact of data augmentation on model convergence, performance stability, and learned feature representations.

#### Validation Graphs

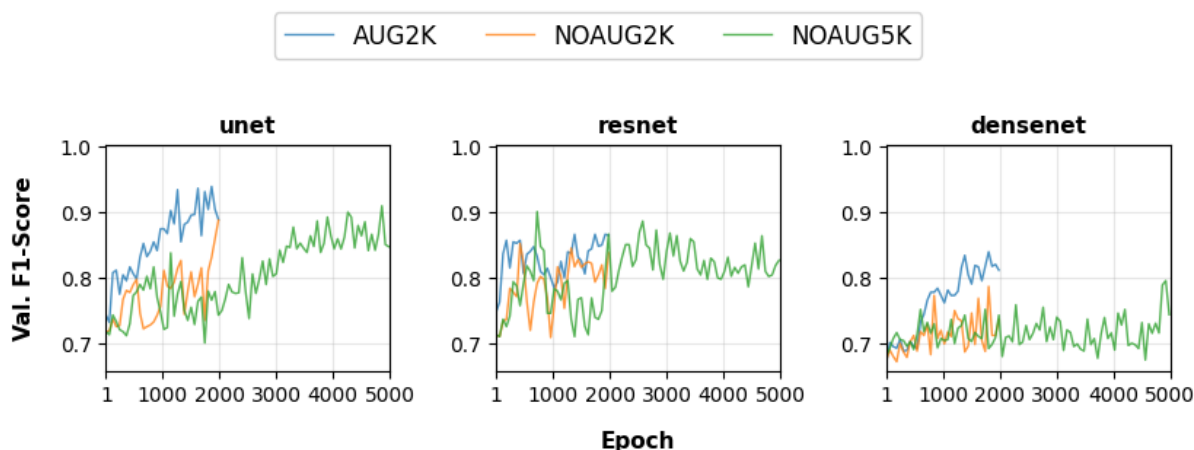


Figure 4.18 Validation F1-score curves across training configurations. Augmentation consistently accelerates convergence and achieves superior performance, with U-Net demonstrating gradual improvement over 5,000 epochs without augmentation, though remaining inefficient compared to augmented training.

Figure 4.18 demonstrates that augmentation consistently accelerates convergence across all architectures. U-Net with augmentation achieves stable F1-scores above 0.9 within 1,000 epochs, substantially outperforming non-augmented variants that plateau around 0.87. Extended training to 5,000 epochs shows gradual improvement for U-Net without augmentation, though remaining significantly less efficient than augmented training.

ResNet exhibits similar acceleration with augmentation, though the NOAUG5K configuration achieves comparable peak F1-scores through substantial fluctuations rather than stable convergence. This apparent performance represents random search behaviour rather than robust learning, as evidenced by the high variance throughout training. DenseNet shows the most dramatic augmentation dependence, with non-augmented variants struggling to exceed 0.74 F1-score regardless of training duration, while augmentation enables reliable performance.

## Performance Metrics

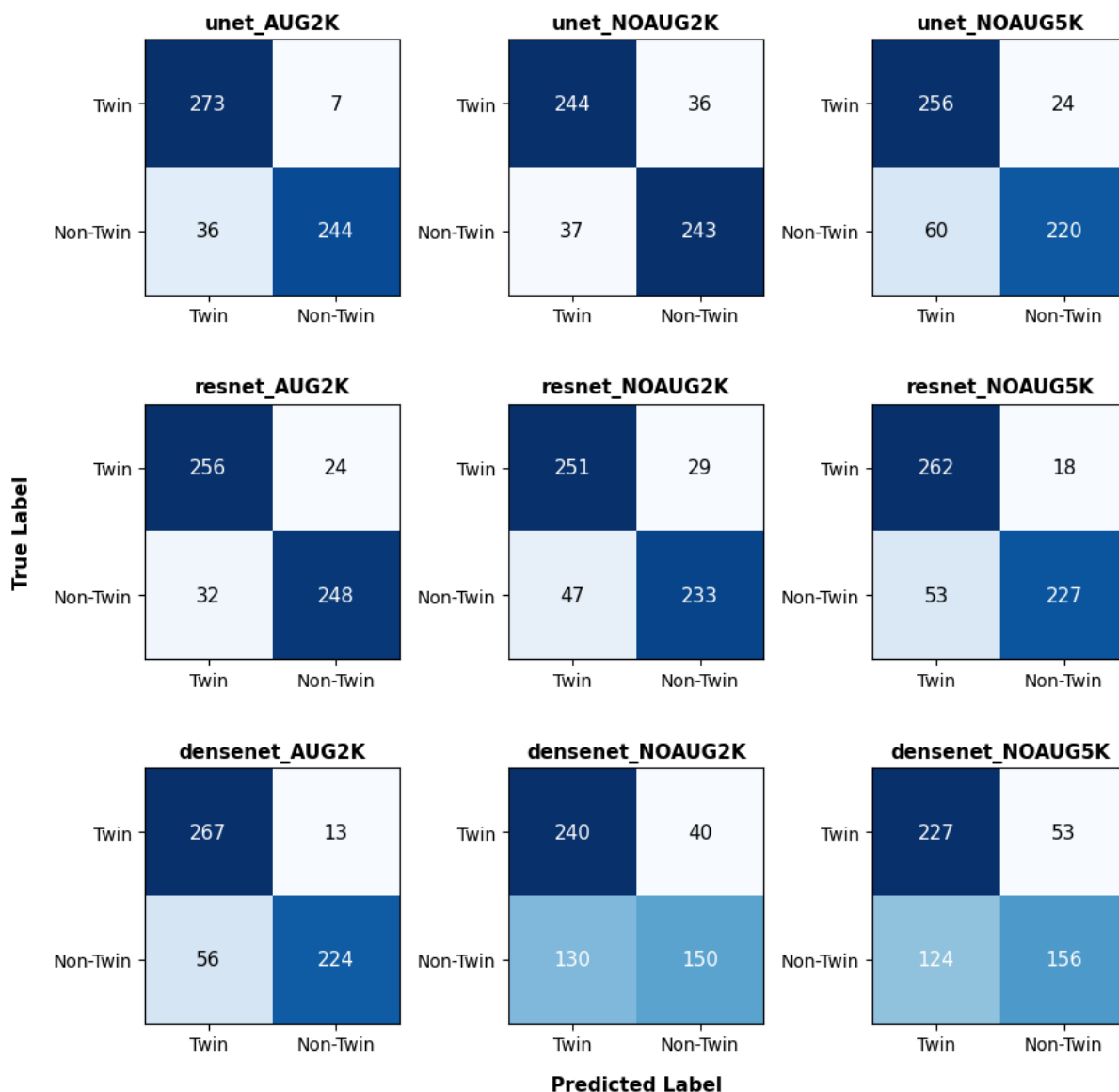


Figure 4.19 Confusion matrices comparing augmented and non-augmented training across architectures. Augmentation consistently improves classification accuracy, with DenseNet showing the most substantial performance degradation without augmentation.

The confusion matrices in Figure 4.19 reveal consistent performance improvements with augmentation. U-Net maintains strong twin identification (273 vs 244 correct classifications) and improved non-twin discrimination (244 vs 243) with augmentation. ResNet shows similar patterns with marginally better twin identification under augmentation. DenseNet exhibits the most dramatic differences, maintaining reasonable twin identification without augmentation but substantially degraded non-twin classification (150 vs 224 correct non-twin classifications), indicating poor generalisation capabilities.

Table 4.4 Ablation study on data augmentation strategies showing performance across architectures. A2k: Augmented 2k epochs, 2k: Non-augmented 2k epochs, 5k: Non-augmented 5k epochs.

Architecture	Accuracy ↑	F1 ↑	AUC ↑	Precision ↑	Recall ↑
U-Net (2k)	84.5	85.4	88.97	80.9	90.7
U-Net (5k)	83.9	85.4	90.05	79.2	92.9
<b>U-Net (A2k)</b>	<b>91.4</b>	<b>92.0</b>	<b>95.2</b>	<b>87.1</b>	<b>97.9</b>
ResNet (2k)	85.5	86.3	92.2	82.3	91.4
ResNet (5k)	81.8	83.9	87.5	76.0	<b>94.3</b>
<b>ResNet (A2k)</b>	<b>89.6</b>	<b>89.6</b>	<b>92.8</b>	<b>90.8</b>	88.6
DenseNet (2k)	68.9	74.2	73.1	64.9	88.2
DenseNet (5k)	69.8	71.9	71.6	68.0	76.8
<b>DenseNet (A2k)</b>	<b>87.1</b>	<b>88.5</b>	<b>91.9</b>	<b>81.6</b>	<b>97.1</b>

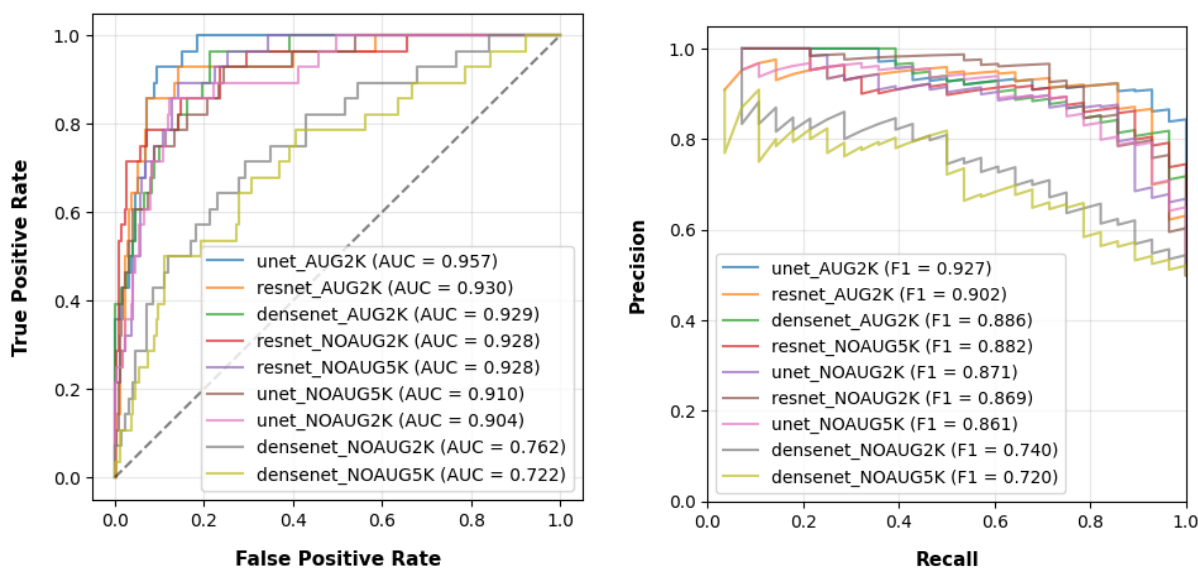
Table 4.4 demonstrates significant performance improvements with augmentation across all metrics. U-Net shows consistent 6-8% improvements in accuracy and F1-score with augmentation, while extended training without augmentation provides minimal gains. ResNet exhibits similar patterns with 4-5% improvements under augmentation. DenseNet shows the most dramatic dependence on augmentation, with 18.2% accuracy improvement and 14.3% F1-score enhancement, indicating that DenseNet performs reliably only with augmentation. Extended training without augmentation fails to recover performance losses, suggesting fundamental limitations in feature learning without data variability.

### ROC and Precision-Recall Analysis

The ROC and PR curves in Figure 4.20 emphasise augmentation’s critical role in achieving robust performance. Augmented models consistently achieve higher AUC values and maintain superior precision across recall ranges. DenseNet shows particularly stark differences, with non-augmented variants achieving only 73-76% AUC compared to 92% with augmentation. The precision-recall curves reveal that non-augmented training results in earlier precision degradation, indicating reduced reliability at higher sensitivity thresholds. Extended training provides marginal improvements but fails to match augmented performance.

### Embedding Distance Distribution Analysis

Figure 4.21 reveals that augmentation consistently produces superior embedding separation. Augmented models achieve clear bimodal distributions with minimal overlap between twin and non-twin pairs. Non-augmented training results in broader distributions with substantial overlap, while extended training paradoxically worsens separation by causing embedding collapse. DenseNet exhibits particularly poor



(a) ROC curves showing discriminative performance across augmentation strategies.

(b) PR curves demonstrating augmentation's impact on classification balance.

Figure 4.20 ROC and PR analysis revealing augmentation's critical role in achieving robust discriminative performance, particularly for DenseNet architecture.

separation without augmentation, with near-uniform distributions indicating failed feature learning. The embedding analysis directly correlates with quantitative performance metrics, validating augmentation's necessity for effective similarity learning.

### Reference Attribution Patterns

The LRP analysis in Figure 4.22 demonstrates that augmentation enables focused feature extraction. Augmented models show concentrated activation in specific brain regions with clear spatial coherence, while non-augmented variants exhibit diffuse activation patterns across the entire brain volume. Despite similar gross anatomical patterns emerging in non-augmented training, the quantitative metrics indicate substantially lower feature quality. This suggests that while models without augmentation identify relevant brain regions, they fail to extract discriminative features effectively. The widespread activation in non-augmented models indicates overfitting to training data specifics rather than learning generalise neuroanatomical patterns.

### MMP Glasser Brain Regions

Table 4.5 reveals systematic differences in regional focus across augmentation strategies. Augmented training demonstrates pronounced subcortical dominance, with 6 subcortical structures (thalamus, brainstem, hypothalamus, basal ganglia, cerebellum,

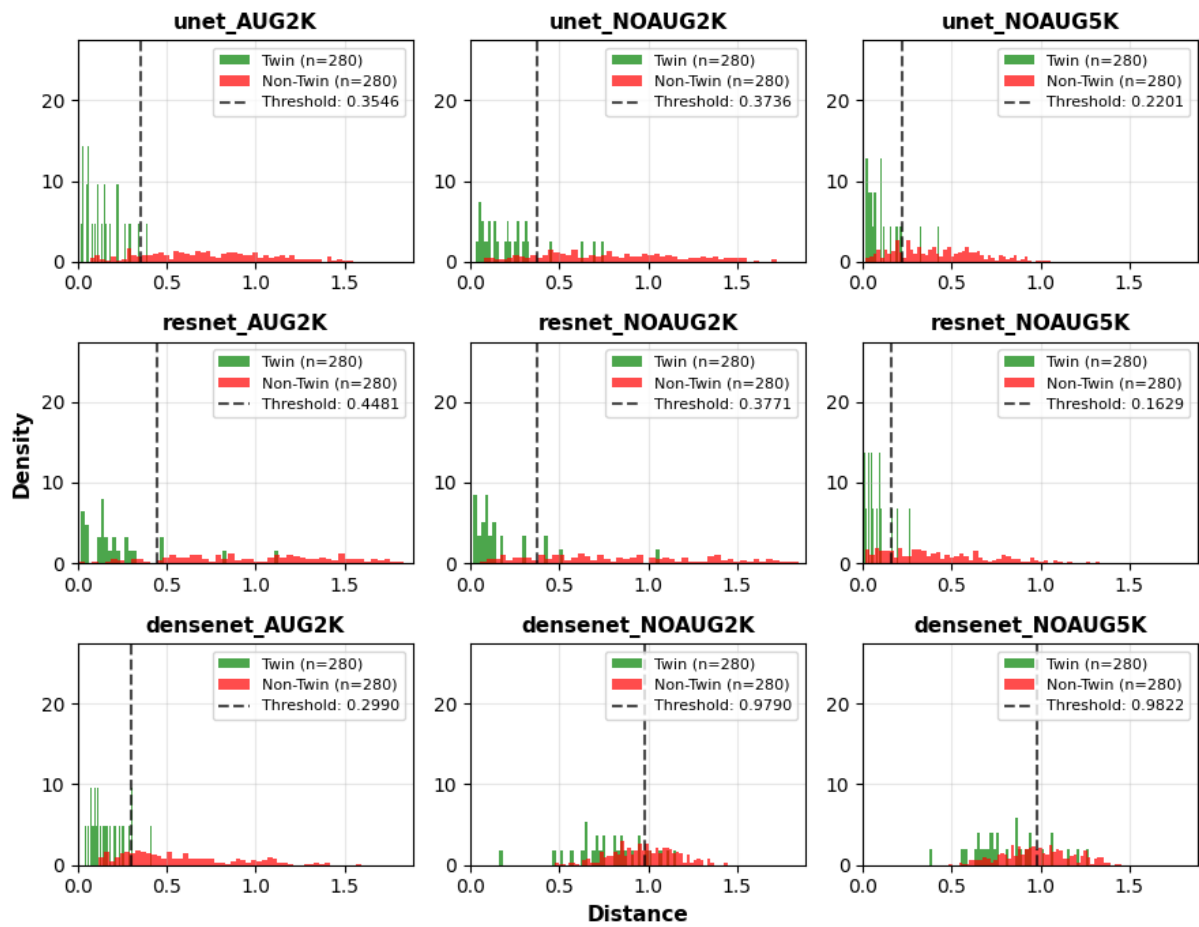


Figure 4.21 Embedding distance distributions across augmentation strategies. Augmentation consistently improves twin vs non-twin separation, while extended training without augmentation leads to embedding collapse with reduced discriminative capacity.

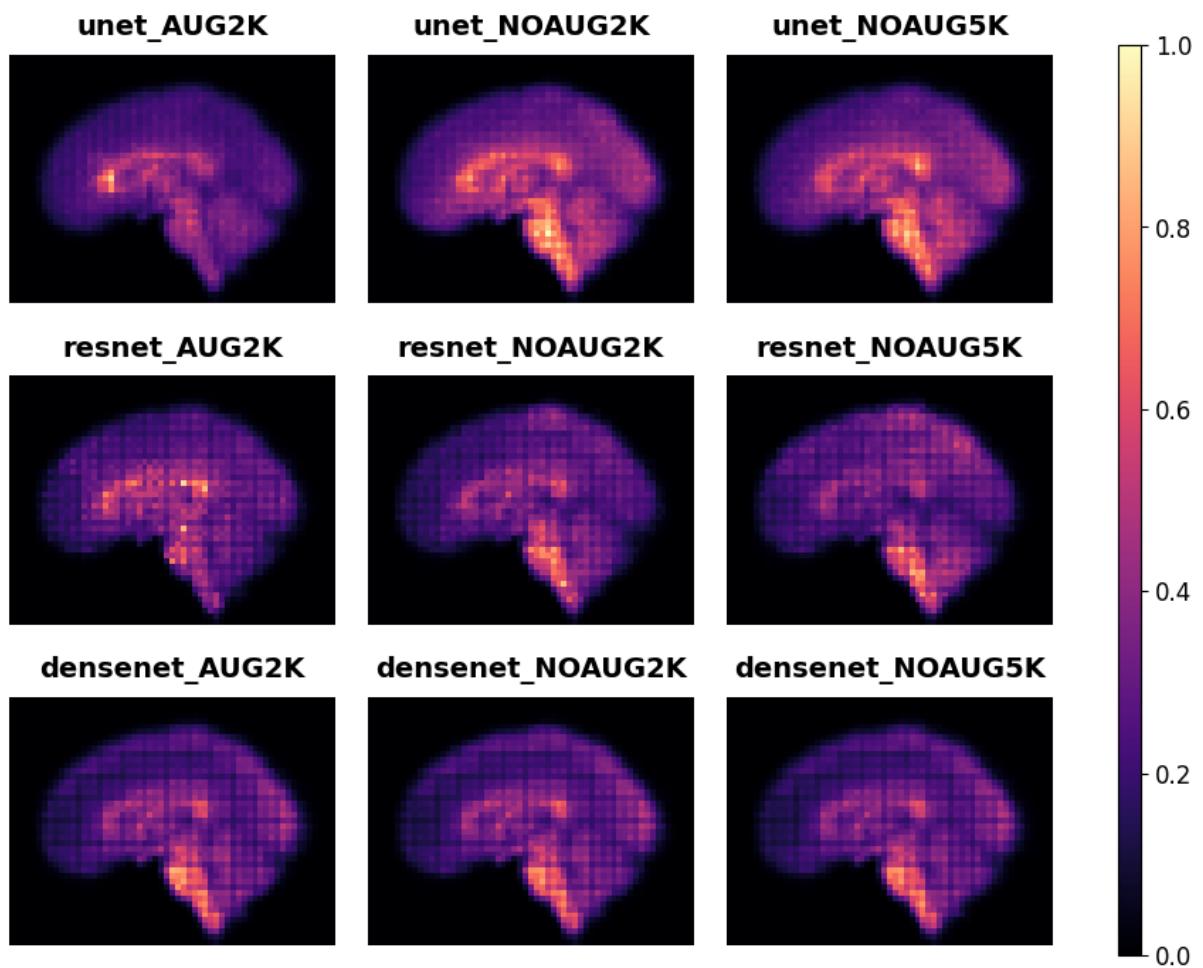


Figure 4.22 Reference LRP maps across augmentation strategies. Non-augmented training produces diffuse activation patterns across the entire brain, while augmentation enables focused feature extraction in relevant neuroanatomical regions.

Table 4.5 Regional importance rankings across augmentation strategies. Augmented training shows strong subcortical dominance with 6 subcortical structures in top 7 positions, while non-augmented conditions exhibit more distributed cortical-subcortical patterns. Full cortex names available in Table A.1.

AUG2K		NOAUG2K		NOAUG5K	
Structure	Imp	Structure	Imp	Structure	Imp
<b>THALAMUS</b>	0.955	<b>THALAMUS</b>	0.954	<b>THALAMUS</b>	0.879
<b>BRAINSTEM</b>	0.875	<b>BRAINSTEM</b>	0.848	<b>BRAINSTEM</b>	0.837
<b>HYPOTHALAMUS</b>	0.707	SPC	0.771	SPC	0.811
<b>BASAL GANGLIA</b>	0.585	<b>HYPOTHALAMUS</b>	0.745	PCC	0.752
PCC	0.546	PCC	0.745	PLMC	0.749
<b>CEREBELLUM</b>	0.545	PLMC	0.728	DSVC	0.733
<b>LIMBIC</b>	0.516	DSVC	0.716	IPC	0.700
EAC	0.490	EAC	0.705	EAC	0.698
TPOJ	0.488	<b>BASAL GANGLIA</b>	0.684	<b>HYPOTHALAMUS</b>	0.694
SPC	0.476	IPC	0.677	SMC	0.691
IPC	0.474	SMC	0.663	DLPC	0.669
VSVC	0.464	DLPC	0.659	MT+CNVA	0.661
MT+CNVA	0.463	TPOJ	0.650	TPOJ	0.656
PLMC	0.461	MT+CNVA	0.647	PMC	0.640
IFOC	0.444	PMC	0.633	<b>CEREBELLUM</b>	0.633
DSVC	0.440	IFOC	0.623	EVC	0.623
POC	0.431	<b>LIMBIC</b>	0.610	AAC	0.616
EVC	0.402	AAC	0.608	<b>BASAL GANGLIA</b>	0.611
MTC	0.401	EVC	0.606	IFOC	0.590
DLPC	0.377	<b>CEREBELLUM</b>	0.589	VSVC	0.585
SMC	0.375	ACMPC	0.585	PVC	0.581
AAC	0.372	PVC	0.564	<b>LIMBIC</b>	0.576
PVC	0.370	VSVC	0.563	POC	0.550
ACMPC	0.363	POC	0.546	ACMPC	0.544
PMC	0.342	MTC	0.491	MTC	0.488
IFC	0.285	IFC	0.475	IFC	0.470
LTC	0.240	LTC	0.388	LTC	0.405
OPFC	0.212	OPFC	0.367	OPFC	0.364

limbic) occupying the top 7 positions and achieving the highest importance scores (thalamus: 0.955 vs 0.954/0.879 in non-augmented conditions). This subcortical concentration suggests augmentation enables models to leverage deep brain structures more effectively for pattern recognition.

The enhanced subcortical focus likely reflects the biological robustness of these evolutionarily conserved structures, which maintain consistent morphological relationships due to their fundamental role in neural development and connectivity, allowing augmented models to extract invariant genetic signatures despite spatial transformations. Non-augmented training produces more distributed cortical-subcortical patterns, with cortical regions like SPC, PCC, and PLMC achieving higher rankings earlier in the hierarchy. Extended training without augmentation (NOAUG5K) partially redistributes toward cortical emphasis but maintains subcortical precedence, demonstrating augmentation's role in enhancing subcortical feature utilisation while filtering transformation-related variability that might obscure genetic signatures in more plastic cortical regions.

## 4.5 Discussion

### 4.5.1 Deep Learning for Genetic Similarity Detection

The experimental results demonstrate that Siamese networks with 3D CNNs achieve robust performance for computational genetic similarity detection directly from brain MRI data. U-Net's superior performance (92.0% F1-score, 95.2% AUC-ROC) validates that deep learning approaches can reliably extract subtle morphological signatures distinguishing genetically related individuals, complemented by strong performance from ResNet (89.6% F1-score) and DenseNet (88.5% F1-score).

Triplet loss optimisation demonstrated exceptional performance with limited training data, achieving 92% F1-score using only 83 monozygotic twin pairs (60% of 138 total pairs). This validates metric learning approaches for few-shot neuroimaging scenarios where large datasets are unavailable, with HCP's sophisticated preprocessing pipelines enabling robust discrimination despite data constraints. The success with limited data has significant implications for clinical applications where large-scale neuroimaging datasets may be unavailable.

The adaptation of these architectures to generate 128-dimensional embeddings through global pooling enables effective capture of multivariate genetic similarity patterns. Consistent performance across architectures validates the reliability of this computational approach, while triplet loss optimisation with hard negative mining ensures robust separation between genetically related and unrelated individuals in

embedding space.

## 4.5.2 Computational vs. Statistical Approaches to Genetic Neuroimaging

This work establishes the first computational framework for ranking brain regions by their multivariate discriminative importance for genetic relatedness detection. Unlike ACE studies that decompose statistical variance to determine independent regional heritability (typically 60-90%), our deep learning approach learns directly from MRI data to identify which combinations of brain regions collectively enable genetic similarity detection.

The computational hierarchy reveals pronounced subcortical dominance with large effect size (Cohen's  $d = 2.80$ ,  $t(27) = 5.77$ ,  $p = 3.89 \times 10^{-6}$ ), contrasting with traditional findings of higher cortical heritability (frontal 78-95%, temporal 77-89%) versus subcortical structures (thalamus 42%, cerebellum 24%). This demonstrates that high individual regional heritability does not translate to discriminative utility for computational identification.

Importantly, the models utilise practically all brain regions for twin identification, with compressed importance distributions (Mean: 0.467, Median: 0.453, most regions  $> 0.2$ ), indicating distributed multivariate processing rather than selective dependence. This addresses a fundamentally different question than statistical variance decomposition.

## 4.5.3 Data-Driven Regional Importance Hierarchy

The computational analysis provides the first relative ranking of neuroanatomical regions by their collective discriminative capacity. The systematic hierarchy, with subcortical structures dominating (thalamus 0.955, brainstem 0.875, hypothalamus 0.707), primary sensory-motor regions achieving intermediate rankings, and association cortices lowest, emerges directly from pattern recognition rather than statistical modelling assumptions.

This hierarchy demonstrates that evolutionarily ancient regulatory structures (those that developed earlier in brain evolution and are shared across many species) contain the most reliable morphological signatures for twin identification. These core structures, including the thalamus, brainstem, and hypothalamus, evolved to control fundamental biological processes and may retain more consistent genetic patterning compared to newer cortical regions that evolved later and show greater individual variation.

The volume-importance inversion exemplified by the cerebellum's massive size (17.14% of brain volume) but moderate ranking (6th, 0.545) illustrates that models identify precise structural patterns rather than anatomical prominence.

Frontal areas showed lower importance, potentially reflecting their later developmental trajectory and greater environmental susceptibility. Given that subjects are predominantly in their late 20s/early 30s (ages 22-35), these later-maturing regions may have been more shaped by individual experiences, reducing their genetic similarity detection utility despite high heritability.

#### 4.5.4 Methodological Paradigm and Clinical Integration

The integration of computational approaches with neuroimaging represents a paradigm shift from statistical modelling to direct pattern recognition in medical imaging. Successful clinical format conversion and Connectome Workbench compatibility demonstrates immediate practical utility for research and clinical environments.

The framework's generalisability extends to diverse neuroimaging classification tasks requiring regional importance analysis, including Alzheimer's detection, psychiatric disorder classification, and treatment response prediction. In neurodegenerative disease research, the approach could identify brain regions most discriminative for early-stage detection, potentially revealing novel biomarkers before clinical symptoms emerge. For psychiatric disorders, the methodology could map regional contributions to conditions like schizophrenia, depression, or bipolar disorder, providing insights into the neuroanatomical basis of mental health conditions.

Treatment response prediction represents another promising application, where the framework could analyse pre-treatment brain scans to identify regions most predictive of therapeutic outcomes. This could enable personalised treatment selection by revealing which patients are most likely to respond to specific interventions based on their individual brain morphology patterns. Drug development research could benefit from understanding which brain regions are most sensitive to pharmaceutical interventions.

The combination of Siamese networks with Layer-Wise Relevance Propagation provides interpretable deep learning suitable for medical applications requiring spatial understanding of classification decisions. This computational approach provides researchers with interpretable tools for understanding regional contributions through direct analysis of medical imaging data, advancing precision neuroimaging by enabling clinicians and researchers to complement statistical correlations with direct computational and regional analysis of morphological patterns underlying neurological and psychiatric conditions.

## 5 Conclusion

### 5.1 Summary of Contributions

This research established the first computational framework for ranking neuroanatomical regions by their collective discriminative importance for genetic relatedness detection. Siamese networks with 3D CNN backbones achieved robust performance for direct genetic similarity detection from brain MRI data using 138 monozygotic twin pairs from the Human Connectome Project, with U-Net demonstrating superior results (92.0% F1-score, 95.2% AUC-ROC).

Layer-Wise Relevance Propagation analysis revealed pronounced subcortical dominance with large effect size (Cohen's  $d = 2.80$ ), with six subcortical structures occupying the top seven positions: thalamus (0.955), brainstem (0.875), and hypothalamus (0.707). This computational hierarchy emerges from direct pattern recognition rather than statistical modelling assumptions, demonstrating successful clinical integration through standard neuroimaging formats.

### 5.2 Impact on Computational Medical Imaging and Neuro-genetics

The methodology provides a paradigmatic shift from statistical variance decomposition to direct computational pattern recognition for investigating genetic influences on brain structure. While ACE studies report highest heritability in cortical areas (frontal 78-95%, temporal 77-89%), the data-driven approach prioritises subcortical structures for actual genetic similarity detection, demonstrating that high regional heritability does not translate to discriminative utility for computational identification.

The automated processing capability coupled with interpretable spatial mapping provides unprecedented efficiency for investigating genetic influences in brain structure. The framework's broader applicability extends to neurodegenerative disorders, psychiatric conditions, and treatment response prediction, representing a significant advancement in precision neuroimaging.

### 5.3 Limitations and Methodological Considerations

The computational framework provides complementary but fundamentally different insights from classical twin studies, focusing on discriminative importance rather than statistical variance decomposition. Dataset constraints limit analysis to healthy young

adults (ages 22-35) from a single high-quality neuroimaging dataset, potentially affecting generalisability to other populations, age groups, or scanning protocols without appropriate domain adaptation strategies.

The focus on structural T1-weighted images overlooks potentially valuable multimodal information, while computational requirements demand significant processing resources. The cross-sectional design prevents investigation of temporal stability in discriminative importance patterns.

## 5.4 Future Directions in Computational Neurogenetics

Multimodal neuroimaging integration could provide comprehensive spatial maps of genetic influences across structural, functional, and connectivity domains, while longitudinal studies would enable investigation of temporal changes in computational discriminative importance patterns. Integration with genomic data represents a promising direction for connecting specific genetic variants to spatial patterns identified through computational analysis, potentially bridging molecular genetics with neuroanatomical phenotypes for targeted therapeutic interventions.

## References

- [1] W. Al-Dhabyani, M. Gomaa, H. Khaled, and A. Fahmy, "Deep Learning Approaches for Data Augmentation and Classification of Breast Masses using Ultrasound Images," *International Journal of Advanced Computer Science and Applications*, vol. 10, no. 5, 2019.
- [2] A. A. Akinyelu, F. Zaccagna, J. T. Grist, M. Castelli, and L. Rundo, "Brain Tumor Diagnosis Using Machine Learning, Convolutional Neural Networks, Capsule Neural Networks and Vision Transformers, Applied to MRI: A Survey," *Journal of Imaging*, vol. 8, no. 8, p. 205, Aug. 2022, Number: 8 Publisher: Multidisciplinary Digital Publishing Institute.
- [3] A. Chattopadhyay and M. Maitra, "MRI-based brain tumour image detection using CNN based deep learning method," *Neuroscience Informatics*, vol. 2, no. 4, p. 100 060, Dec. 2022.
- [4] H. Mzoughi *et al.*, "Deep Multi-Scale 3D Convolutional Neural Network (CNN) for MRI Gliomas Brain Tumor Classification," *Journal of Digital Imaging*, vol. 33, no. 4, pp. 903–915, Aug. 2020.
- [5] J. Miah, D. M. Cao, M. A. Sayed<sup>3</sup>, M. S. Taluckder, M. S. Haque, and F. Mahmud, *Advancing Brain Tumor Detection: A Thorough Investigation of CNNs, Clustering, and SoftMax Classification in the Analysis of MRI Images*, arXiv:2310.17720 [eess], Oct. 2023.
- [6] A. Biondi *et al.*, "Are the Brains of Monozygotic Twins Similar? A Three-Dimensional MR Study," 1998.
- [7] A. Mohr, M. Weisbrod, P. Schellinger, and M. Knauth, "The similarity of brain morphology in healthy monozygotic twins," *Cognitive Brain Research*, vol. 20, no. 1, pp. 106–110, Jun. 2004.
- [8] V. Sundaresan and S. Amala Shanthi, "Monozygotic twin face recognition: An in-depth analysis and plausible improvements," *Image and Vision Computing*, vol. 116, p. 104 331, Dec. 2021.
- [9] C. J. Parde *et al.*, "Twin Identification over Viewpoint Change: A Deep Convolutional Neural Network Surpasses Humans," *ACM Transactions on Applied Perception*, vol. 20, no. 3, pp. 1–15, Jul. 2023.
- [10] Khalid M.O. Nahar, Bilal Anas Abul-Huda, Abeer Fayez Al.bataineh, and Ra'ed M. Al-Khatib, *Twins and Similar Faces Recognition Using Geometric and Photometric Features with Transfer Learning*, 2021.

- [11] V. P. B. Grover, J. M. Tognarelli, M. M. E. Crossey, I. J. Cox, S. D. Taylor-Robinson, and M. J. W. McPhail, "Magnetic Resonance Imaging: Principles and Techniques: Lessons for Clinicians," *Journal of Clinical and Experimental Hepatology*, vol. 5, no. 3, pp. 246–255, Sep. 2015.
- [12] A. Panigrahy and S. Blüml, "Neuroimaging of Pediatric Brain Tumors: From Basic to Advanced Magnetic Resonance Imaging (MRI)," *Journal of Child Neurology*, vol. 24, no. 11, pp. 1343–1365, Nov. 2009, Publisher: SAGE Publications Inc.
- [13] A. Prasad, "Making Images/Making Bodies: Visibilizing and Disciplining through Magnetic Resonance Imaging (MRI)," *Science, Technology, & Human Values*, vol. 30, no. 2, pp. 291–316, Apr. 2005, Publisher: SAGE Publications Inc.
- [14] O. Sitburana and W. G. Ondo, "Brain magnetic resonance imaging (MRI) in parkinsonian disorders," *Parkinsonism & Related Disorders*, vol. 15, no. 3, pp. 165–174, Mar. 2009.
- [15] M. F. Glasser *et al.*, "A multi-modal parcellation of human cerebral cortex," *Nature*, vol. 536, no. 7615, pp. 171–178, Aug. 2016.
- [16] Z. Li, F. Liu, W. Yang, S. Peng, and J. Zhou, "A Survey of Convolutional Neural Networks: Analysis, Applications, and Prospects," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 12, pp. 6999–7019, Dec. 2022, Conference Name: IEEE Transactions on Neural Networks and Learning Systems.
- [17] S. Deepak and P. M. Ameer, "Automated Categorization of Brain Tumor from MRI Using CNN features and SVM," *Journal of Ambient Intelligence and Humanized Computing*, vol. 12, no. 8, pp. 8357–8369, Aug. 2021.
- [18] L. Zou, J. Zheng, C. Miao, M. J. Mckeown, and Z. J. Wang, "3D CNN Based Automatic Diagnosis of Attention Deficit Hyperactivity Disorder Using Functional and Structural MRI," *IEEE Access*, vol. 5, pp. 23 626–23 636, 2017, Conference Name: IEEE Access.
- [19] F. J. Dorfner, J. B. Patel, J. Kalpathy-Cramer, E. R. Gerstner, and C. P. Bridge, "A review of deep learning for brain tumor analysis in MRI," *npj Precision Oncology*, vol. 9, no. 1, p. 2, Jan. 2025, Publisher: Nature Publishing Group.
- [20] G. Folego, M. Weiler, R. F. Casseb, R. Pires, and A. Rocha, "Alzheimer's Disease Detection Through Whole-Brain 3D-CNN MRI," *Frontiers in Bioengineering and Biotechnology*, vol. 8, Oct. 2020, Publisher: Frontiers.
- [21] J. Wen *et al.*, "Convolutional neural networks for classification of Alzheimer's disease: Overview and reproducible evaluation," *Medical Image Analysis*, vol. 63, p. 101 694, Jul. 2020.

- [22] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," Springer, 2015, pp. 234–241, ISBN: 3-319-24573-2.
- [23] J. Walsh, A. Othmani, M. Jain, and S. Dev, "Using U-Net network for efficient brain tumor segmentation in MRI images," *Healthcare Analytics*, vol. 2, p. 100 098, 2022, Publisher: Elsevier.
- [24] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3D U-Net: Learning dense volumetric segmentation from sparse annotation," Springer, 2016, pp. 424–432, ISBN: 3-319-46722-0.
- [25] J. Schlemper *et al.*, "Attention gated networks: Learning to leverage salient regions in medical images," *Medical image analysis*, vol. 53, pp. 197–207, 2019, Publisher: Elsevier.
- [26] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "Unet++: Redesigning skip connections to exploit multiscale features in image segmentation," *IEEE transactions on medical imaging*, vol. 39, no. 6, pp. 1856–1867, 2019, Publisher: IEEE.
- [27] R. Yousef *et al.*, "U-Net-based models towards optimal MR brain image segmentation," *Diagnostics*, vol. 13, no. 9, p. 1624, 2023, Publisher: MDPI.
- [28] Y. Yuan and Y. Cheng, "Medical image segmentation with UNet-based multi-scale context fusion," *Scientific Reports*, vol. 14, no. 1, p. 15 687, 2024, Publisher: Nature Publishing Group UK London.
- [29] J. Chen *et al.*, "TransUNet: Rethinking the U-Net architecture design for medical image segmentation through the lens of transformers," *Medical Image Analysis*, vol. 97, p. 103 280, 2024, Publisher: Elsevier.
- [30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," 2016, pp. 770–778.
- [31] M. Talo, O. Yildirim, U. B. Baloglu, G. Aydin, and U. R. Acharya, "Convolutional neural networks for multi-class brain disease detection using MRI images," *Computerized Medical Imaging and Graphics*, vol. 78, p. 101 673, 2019, Publisher: Elsevier.
- [32] M. T. R, V. K. V, and S. Guluwadi, "Enhancing brain tumor detection in MRI images through explainable AI using Grad-CAM with Resnet 50," *BMC medical imaging*, vol. 24, no. 1, p. 107, 2024, Publisher: Springer.
- [33] S. Korolev, A. Safiullin, M. Belyaev, and Y. Dodonova, "Residual and plain convolutional neural networks for 3D brain MRI classification," IEEE, 2017, pp. 835–838, ISBN: 1-5090-1172-2.

- [34] K. Lakshmi, S. Amaran, G. Subbulakshmi, S. Padmini, G. P. Joshi, and W. Cho, "Explainable artificial intelligence with UNet based segmentation and Bayesian machine learning for classification of brain tumors using MRI images," *Scientific Reports*, vol. 15, no. 1, p. 690, 2025, Publisher: Nature Publishing Group UK London.
- [35] S. Lu, S.-H. Wang, and Y.-D. Zhang, "Detecting pathological brain via ResNet and randomized neural networks," *Heliyon*, vol. 6, no. 12, 2020, Publisher: Elsevier.
- [36] D. Nie *et al.*, "Multi-channel 3D deep feature learning for survival time prediction of brain tumor patients using multi-modal neuroimages," *Scientific reports*, vol. 9, no. 1, p. 1103, 2019, Publisher: Nature Publishing Group UK London.
- [37] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," 2017, pp. 4700–4708.
- [38] R. D. Gottapu and C. H. Dagli, "DenseNet for anatomical brain segmentation," *Procedia Computer Science*, vol. 140, pp. 179–185, 2018, Publisher: Elsevier.
- [39] Y. Mao, J. Kim, L. Podina, and M. Kohandel, "Dilated SE-DenseNet for brain tumor MRI classification," *Scientific Reports*, vol. 15, no. 1, p. 3596, 2025, Publisher: Nature Publishing Group UK London.
- [40] X. Zhang, Y. Yang, T. Li, H. Wang, and Z. He, "Discovering senile dementia from brain MRI using Ra-DenseNet," Springer, 2019, pp. 449–460.
- [41] A. Çinar and M. Yildirim, "Detection of tumors on brain MRI images using the hybrid convolutional neural network architecture," *Medical Hypotheses*, vol. 139, p. 109 684, Jun. 2020.
- [42] L. Munroe *et al.*, "Applications of interpretable deep learning in neuroimaging: A comprehensive review," *Imaging Neuroscience*, vol. 2, pp. 1–37, Jul. 2024.
- [43] R. Ranjbarzadeh, A. Bagherian Kasgari, S. Jafarzadeh Ghouschi, S. Anari, M. Naseri, and M. Bendeche, "Brain tumor segmentation based on deep learning and an attention mechanism using MRI multi-modalities brain images," *Scientific Reports*, vol. 11, no. 1, p. 10 930, May 2021, Publisher: Nature Publishing Group.
- [44] I. Pacal, O. Celik, B. Bayram, and A. Cunha, "Enhancing EfficientNetv2 with global and efficient channel attention mechanisms for accurate MRI-Based brain tumor classification," *Cluster Computing*, vol. 27, no. 8, pp. 11 187–11 212, Nov. 2024.
- [45] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah, "Signature Verification using a "Siamese" Time Delay Neural Network," in *Advances in Neural Information Processing Systems*, vol. 6, Morgan-Kaufmann, 1993.

- [46] S. Jindal, G. Gupta, M. Yadav, M. Sharma, and L. Vig, "Siamese Networks for Chromosome Classification," 2017, pp. 72–81.
- [47] M. D. Li *et al.*, "Siamese neural networks for continuous disease severity evaluation and change detection in medical imaging," *npj Digital Medicine*, vol. 3, no. 1, p. 48, Mar. 2020.
- [48] I. E. Livieris, E. Pintelas, N. Kiriakidou, and P. Pintelas, "Explainable Image Similarity: Integrating Siamese Networks and Grad-CAM," *Journal of Imaging*, vol. 9, no. 10, p. 224, Oct. 2023, Number: 10 Publisher: Multidisciplinary Digital Publishing Institute.
- [49] L. Xu *et al.*, "A Siamese Network With Node Convolution for Individualized Predictions Based on Connectivity Maps Extracted From Resting-State fMRI Data," *IEEE Journal of Biomedical and Health Informatics*, vol. 27, no. 11, pp. 5418–5429, Nov. 2023.
- [50] T. Y. Liu and J. Feng, "Triplet contrastive learning for brain tumor classification," *arXiv preprint arXiv:2108.03611*, 2021.
- [51] J. Lee, S. Ahn, H.-S. Kim, J. An, and J. Sim, "A robust model training strategy using hard negative mining in a weakly labeled dataset for lymphatic invasion in gastric cancer," *The Journal of Pathology: Clinical Research*, vol. 10, no. 1, e355, 2024, Publisher: Wiley Online Library.
- [52] B. Gajić, A. Amato, and C. Gatta, "Fast hard negative mining for deep metric learning," *Pattern Recognition*, vol. 112, p. 107 795, 2021, Publisher: Elsevier.
- [53] Z. Li, C. Guo, X. Wang, Z. Feng, and Z. Du, "Selectively hard negative mining for alleviating gradient vanishing in image-text matching," *IEEE Transactions on Circuits and Systems for Video Technology*, 2024, Publisher: IEEE.
- [54] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Boston, MA, USA: IEEE, Jun. 2015, pp. 815–823, ISBN: 978-1-4673-6964-0.
- [55] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, "On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation," *PLOS ONE*, vol. 10, no. 7, O. D. Suarez, Ed., e0130140, Jul. 2015.
- [56] G. Montavon, A. Binder, S. Lapuschkin, W. Samek, and K.-R. Müller, "Layer-Wise Relevance Propagation: An Overview," in *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K.-R. Müller, Eds., Cham: Springer International Publishing, 2019, pp. 193–209, ISBN: 978-3-030-28954-6.

- [57] A. Binder, G. Montavon, S. Lapuschkin, K.-R. Müller, and W. Samek, "Layer-Wise Relevance Propagation for Neural Networks with Local Renormalization Layers," in *Artificial Neural Networks and Machine Learning – ICANN 2016*, A. E. Villa, P. Masulli, and A. J. Pons Rivero, Eds., Cham: Springer International Publishing, 2016, pp. 63–71, ISBN: 978-3-319-44781-0.
- [58] R. R. Selvaraju, A. Das, R. Vedantam, M. Cogswell, D. Parikh, and D. Batra, *Grad-CAM: Why did you say that?* arXiv:1611.07450 [stat], Jan. 2017.
- [59] M. Böhle, F. Eitel, M. Weygandt, and K. Ritter, "Layer-Wise Relevance Propagation for Explaining Deep Neural Network Decisions in MRI-Based Alzheimer's Disease Classification," *Frontiers in Aging Neuroscience*, Jul. 2019, Place: Lausanne, Switzerland Publisher: Frontiers Research Foundation Section: Original Research ARTICLE.
- [60] L. v. d. Maaten and G. Hinton, "Visualizing Data using t-SNE," *Journal of Machine Learning Research*, vol. 9, no. 86, pp. 2579–2605, 2008.
- [61] M. W. Lukies *et al.*, "Heritability of brain volume on MRI in middle to advanced age: A twin study of Japanese adults," *PLOS ONE*, vol. 12, no. 4, e0175800, Apr. 2017, Publisher: Public Library of Science.
- [62] A. G. Jansen, S. E. Mous, T. White, D. Posthuma, and T. J. Polderman, "What twin studies tell us about the heritability of brain development, morphology, and function: A review," *Neuropsychology review*, vol. 25, no. 1, pp. 27–46, 2015, Publisher: Springer.
- [63] E. V. Sullivan, A. Pfefferbaum, G. E. Swan, and D. Carmelli, "Heritability of hippocampal size in elderly twin men: Equivalent influence from genes and environment," *Hippocampus*, vol. 11, no. 6, pp. 754–762, 2001.
- [64] W. Wen *et al.*, "Distinct Genetic Influences on Cortical and Subcortical Brain Structures," *Scientific Reports*, vol. 6, p. 32 760, Sep. 2016.
- [65] J. E. Schmitt *et al.*, "A multivariate analysis of neuroanatomic relationships in a genetically informative pediatric sample," *Neuroimage*, vol. 35, no. 1, pp. 70–82, 2007, Publisher: Elsevier.
- [66] G. L. Wallace *et al.*, "A pediatric twin study of brain morphometry," *Journal of Child Psychology and Psychiatry, and Allied Disciplines*, vol. 47, no. 10, pp. 987–993, Oct. 2006.
- [67] D. C. Van Essen *et al.*, "The Human Connectome Project: A data acquisition perspective," *NeuroImage*, vol. 62, no. 4, pp. 2222–2231, Oct. 2012.
- [68] M. Jenkinson, C. F. Beckmann, T. E. J. Behrens, M. W. Woolrich, and S. M. Smith, "FSL," *NeuroImage*, 20 YEARS OF fMRI, vol. 62, no. 2, pp. 782–790, Aug. 2012.

- [69] M. F. Glasser *et al.*, "The minimal preprocessing pipelines for the Human Connectome Project," *NeuroImage*, Mapping the Connectome, vol. 80, pp. 105–124, Oct. 2013.
- [70] B. Fischl, "FreeSurfer," *NeuroImage*, 20 YEARS OF fMRI, vol. 62, no. 2, pp. 774–781, Aug. 2012.
- [71] D. Marcus *et al.*, "Informatics and Data Mining Tools and Strategies for the Human Connectome Project," *Frontiers in Neuroinformatics*, vol. 5, Jun. 2011, Publisher: Frontiers.
- [72] J. S. Elam *et al.*, "The human connectome project: A retrospective," *NeuroImage*, vol. 244, p. 118 543, 2021, Publisher: Elsevier.

# Appendix A Cortical and Subcortical Mappings

Table A.1 Mapping of cortex abbreviations to full names

Short Name	Full Name
AAC	Auditory Association
ACMPC	Anterior Cingulate and Medial Prefrontal
DLPC	Dorsolateral Prefrontal
DSVC	Dorsal Stream Visual
EAC	Early Auditory
EVC	Early Visual
IFC	Inferior Frontal
IFOC	Insular and Frontal Opercular
IPC	Inferior Parietal
LTC	Lateral Temporal
MT+CNVA	MT+ Complex and Neighboring Visual Areas
MTC	Medial Temporal
OPFC	Orbital and Polar Frontal
PCC	Posterior Cingulate
PMC	Premotor
PLMC	Paracentral Lobular and Mid Cingulate
POC	Posterior Opercular
PVC	Primary Visual
SMC	Somatosensory and Motor
SPC	Superior Parietal
TPOJ	Temporo-Parieto-Occipital Junction
VSVC	Ventral Stream Visual

Table A.2 Mapping of subcortical areas to functional groups. All mappings refer to both areas in the left and right hemispheres with the exception of the brain stem which constitutes a single area.

<b>Subcortical Area</b>	<b>Group</b>
THALAMUS (L/R)	THALAMUS
BRAIN_STEM	BRAINSTEM
DIENCEPHALON_VENTRAL (L/R)	HYPOTHALAMUS
CAUDATE (L/R)	BASAL GANGLIA
PUTAMEN (L/R)	BASAL GANGLIA
PALLIDUM (L/R)	BASAL GANGLIA
ACCUMBENS (L/R)	BASAL GANGLIA
CEREBELLUM (L/R)	CEREBELLUM
HIPPOCAMPUS (L/R)	LIMBIC
AMYGDALA (L/R)	LIMBIC

# Appendix B Full Brain Region Importance Listing

Table B.1 All individual brain regions from HCP-MMP 1.0 atlas sorted by importance, showing cortex, importance scores, activation, and volume for each region. Subcortical regions denoted in bold. Full cortex names available in Table A.1

Long Name	Cortex/Subcortex	Importance	Total Activation	Volume (%)
<b>THALAMUS_LEFT</b>	THALAMUS	1.000	12584.03	1.1
<b>THALAMUS_RIGHT</b>	THALAMUS	0.910	11413.84	1.0
<b>BRAIN_STEM</b>	BRAINSTEM	0.875	31595.38	3.0
Area_31a_L	PCC	0.840	871.89	0.1
Area_31a_R	PCC	0.806	644.11	0.1
ParaHippocampal_Area_1_L	MTC	0.781	878.85	0.1
<b>DIENCEPHALON_VENTRAL_LEFT</b>	HYPOTHALAMUS	0.755	5591.59	0.6
<b>PALLIDUM_LEFT</b>	BASAL GANGLIA	0.741	2336.38	0.3
Area_23d_R	PCC	0.723	1110.66	0.1
Area_31pd_L	PCC	0.713	679.36	0.1
Area_OP2-3-VS_L	POC	0.703	835.06	0.1
<b>PUTAMEN_LEFT</b>	BASAL GANGLIA	0.682	6989.74	0.9
Insular_Granular_Complex_L	IFOC	0.676	406.05	0.1
PreCuneus_Visual_Area_R	PCC	0.671	1317.74	0.2
RetrolInsular_Cortex_L	EAC	0.666	1004.60	0.1
<b>CAUDATE_RIGHT</b>	BASAL GANGLIA	0.666	5100.32	0.6
Area_dorsal_23_a+b_L	PCC	0.661	763.93	0.1
<b>DIENCEPHALON_VENTRAL_RIGHT</b>	HYPOTHALAMUS	0.660	5040.54	0.6
Area_23d_L	PCC	0.653	684.45	0.1
Area_23c_L	PLMC	0.651	1516.92	0.2
Area_anterior_32_prime_L	ACMPC	0.650	1717.33	0.2
Primary_Auditory_Cortex_L	EAC	0.637	662.70	0.1
ParaHippocampal_Area_2_L	MTC	0.634	463.01	0.1
<b>CAUDATE_LEFT</b>	BASAL GANGLIA	0.633	4378.87	0.6
<b>HIPPOCAMPUS_LEFT</b>	LIMBIC	0.626	4832.53	0.6
Area_TemporoParietoOccipital_Junction_1_L	TPOJ	0.622	1257.31	0.2
VentroMedial_Visual_Area_3_L	VSVC	0.609	725.20	0.1
Area_7m_R	PCC	0.609	1132.95	0.2
ParaBelt_Complex_L	EAC	0.608	1431.64	0.2
Frontal_Opercular_Area_3_L	IFOC	0.602	359.17	0.0
PreCuneus_Visual_Area_L	PCC	0.593	1233.22	0.2
Ventral_Area_24d_L	PLMC	0.589	892.30	0.1
<b>PALLIDUM_RIGHT</b>	BASAL GANGLIA	0.588	1730.57	0.2
Dorsal_Area_24d_L	PLMC	0.587	1754.76	0.3
Area_5m_ventral_L	PLMC	0.586	674.53	0.1
Area_OP1-SII_L	POC	0.586	889.28	0.1
Area_TemporoParietoOccipital_Junction_2_L	TPOJ	0.583	1413.57	0.2
ProStriate_Area_L	PCC	0.580	537.71	0.1
Area_31p_ventral_L	PCC	0.578	616.25	0.1
VentroMedial_Visual_Area_1_L	VSVC	0.578	647.61	0.1
Ventral_Visual_Complex_L	VSVC	0.572	1814.29	0.3
<b>CEREBELLUM_LEFT</b>	CEREBELLUM	0.571	57377.57	8.4
Seventh_Visual_Area_R	DSVC	0.570	346.42	0.1
Area_31p_ventral_R	PCC	0.567	610.76	0.1
Medial_Belt_Complex_L	EAC	0.566	640.33	0.1
Medial_Area_7P_R	SPC	0.562	713.86	0.1
ParaHippocampal_Area_1_R	MTC	0.561	629.13	0.1
Area_PGp_R	IPC	0.561	1545.61	0.2

Continued on next page

## B Full Brain Region Importance Listing

Table B.1 – continued from previous page

Long Name	Cortex	Importance	Total Activation	Volume (%)
PreSubiculum_L	MTC	0.558	398.77	0.1
VentroMedial_Visual_Area_1_R	VSVC	0.553	889.37	0.1
Anterior_Ventral_Insular_Area_L	IFOC	0.551	886.07	0.1
Posterior_Insular_Area_2_L	IFOC	0.549	963.63	0.1
Area_Posterior_24_prime_L	ACMPC	0.548	479.40	0.1
Area_Lateral_Occipital_3_L	MT+CNVA	0.545	693.62	0.1
Area_Lateral_IntraParietal_dorsal_R	SPC	0.545	403.93	0.1
<b>PUTAMEN_RIGHT</b>	<b>BASAL GANGLIA</b>	<b>0.543</b>	<b>5283.64</b>	<b>0.8</b>
Area_IntraParietal_0_L	IPC	0.543	811.80	0.1
Area_23c_R	PLMC	0.542	1498.36	0.2
Area_PGi_L	IPC	0.540	3138.13	0.5
Dorsal_Area_24d_R	PLMC	0.540	1579.77	0.2
Area_IntraParietal_2_L	IPC	0.538	1125.40	0.2
Area_Posterior_Insular_1_L	IFOC	0.535	1081.75	0.2
Middle_Temporal_Area_L	MT+CNVA	0.533	382.56	0.1
Area_ventral_23_a+b_L	PCC	0.532	548.75	0.1
Area_Lateral_Occipital_3_R	MT+CNVA	0.532	584.78	0.1
Area_7m_L	PCC	0.531	1056.95	0.2
Area_TemporoParietoOccipital_Junction_3_R	TPOJ	0.530	821.73	0.1
Area_ventral_23_a+b_R	PCC	0.528	532.20	0.1
Area_V3CD_L	MT+CNVA	0.528	844.09	0.1
Area_PGs_R	IPC	0.528	1679.62	0.3
Ventral_IntraParietal_Complex_R	SPC	0.526	372.62	0.1
ParaHippocampal_Area_3_L	MTC	0.526	839.56	0.1
Anterior_IntraParietal_Area_R	SPC	0.525	1831.41	0.3
Area_IntraParietal_1_R	IPC	0.525	1164.13	0.2
Area_Lateral_Occipital_1_L	MT+CNVA	0.524	266.92	0.0
Area_Frontal_Opercular_5_L	IFOC	0.524	940.92	0.2
Ventral_Visual_Complex_R	VSVC	0.524	1385.58	0.2
Area_31pd_R	PCC	0.523	290.22	0.0
Area_Lateral_IntraParietal_dorsal_L	SPC	0.523	489.06	0.1
Area_dorsal_32_R	ACMPC	0.523	1562.71	0.2
ParaHippocampal_Area_2_R	MTC	0.522	216.27	0.0
Area_IntraParietal_1_L	IPC	0.519	1078.99	0.2
<b>CEREBELLUM_RIGHT</b>	<b>CEREBELLUM</b>	<b>0.519</b>	<b>54238.99</b>	<b>8.7</b>
Area_STSd_posterior_L	AAC	0.518	1229.48	0.2
Retrolsular_Cortex_R	EAC	0.515	933.34	0.2
Area_TA2_L	AAC	0.514	716.40	0.1
Frontal_Opercular_Area_2_L	IFOC	0.510	444.10	0.1
Ventral_IntraParietal_Complex_L	SPC	0.509	568.00	0.1
Medial_Area_7P_L	SPC	0.507	475.71	0.1
Supplementary_and_Cingulate_Eye_Field_R	PLMC	0.507	1383.87	0.2
Middle_Insular_Area_L	IFOC	0.506	936.90	0.2
<b>AMYGDALA_LEFT</b>	<b>LIMBIC</b>	<b>0.505</b>	<b>2014.40</b>	<b>0.3</b>
Lateral_Area_7P_R	SPC	0.505	375.69	0.1
Area_TemporoParietoOccipital_Junction_1_R	TPOJ	0.502	1648.91	0.3
Area_OP2-3-VS_R	POC	0.501	324.45	0.1
Area_STSv_posterior_L	AAC	0.500	1647.95	0.3
Medial_Area_7A_L	SPC	0.500	1008.04	0.2
Area_V3CD_R	MT+CNVA	0.496	805.23	0.1
Parieto-Occipital_Sulcus_Area_2_R	PCC	0.495	1781.80	0.3
PeriSylvian_Language_Area_L	TPOJ	0.494	1202.68	0.2
Area_PFm_Complex_L	IPC	0.493	2946.54	0.5
Parieto-Occipital_Sulcus_Area_2_L	PCC	0.493	1692.83	0.3
Area_PFcm_L	EAC	0.491	806.88	0.1
Area_Lateral_IntraParietal_ventral_R	SPC	0.491	480.41	0.1

Continued on next page

## B Full Brain Region Importance Listing

Table B.1 – continued from previous page

Long Name	Cortex	Importance	Total Activation	Volume (%)
Superior_Temporal_Visual_Area_R	TPOJ	0.490	780.41	0.1
Insular_Granular_Complex_R	IFOC	0.490	618.32	0.1
Area_52_L	EAC	0.488	244.29	0.0
IntraParietal_Sulcus_Area_1_R	DSVC	0.487	752.79	0.1
Medial_Superior_Temporal_Area_L	MT+CNVA	0.487	420.44	0.1
<b>HIPPOCAMPUS_RIGHT</b>	LIMBIC	0.486	4128.72	0.7
<b>ACCUMBENS_LEFT</b>	BASAL GANGLIA	0.485	597.22	0.1
Area_STSd_anterior_L	AAC	0.485	1120.77	0.2
Area_PFm_Complex_R	IPC	0.483	3045.24	0.5
Area_TemporoParietoOccipital_Junction_2_R	TPOJ	0.483	1128.28	0.2
Area_FST_L	MT+CNVA	0.478	1002.28	0.2
Area_PFt_L	IPC	0.477	750.69	0.1
Area_PGs_L	IPC	0.476	2161.47	0.4
Frontal_Opercular_Area_3_R	IFOC	0.474	307.60	0.1
Frontal_Opercular_Area_4_L	IFOC	0.474	1445.37	0.3
Middle_Temporal_Area_R	MT+CNVA	0.474	316.93	0.1
Area_V6A_R	DSVC	0.473	385.83	0.1
Area_PGp_L	IPC	0.470	1172.56	0.2
ParaHippocampal_Area_3_R	MTC	0.470	372.71	0.1
Fusiform_Face_Complex_L	VSVC	0.470	1962.36	0.3
Posterior_Insular_Area_2_R	IFOC	0.470	660.33	0.1
Area_46_L	DLPC	0.468	1869.00	0.3
Area_8BM_R	ACMPC	0.467	1616.02	0.3
Lateral_Area_7A_R	SPC	0.467	532.86	0.1
Parieto-Occipital_Sulcus_Area_1_L	PCC	0.465	1167.48	0.2
Anterior_24_prime_R	ACMPC	0.461	461.34	0.1
Seventh_Visual_Area_L	DSVC	0.460	285.82	0.1
Area_anterior_9-46v_L	DLPC	0.460	1116.23	0.2
Area_V4t_R	VSVC	0.460	458.04	0.1
Area_V6A_L	DSVC	0.459	340.64	0.1
Area_p32_R	ACMPC	0.459	940.91	0.2
Supplementary_and_Cingulate_Eye_Field_L	PLMC	0.458	1215.05	0.2
Medial_IntraParietal_Area_R	SPC	0.457	680.10	0.1
Inferior_6-8_Transitional_Area_L	DLPC	0.457	613.97	0.1
Superior_6-8_Transitional_Area_L	DLPC	0.455	536.10	0.1
Area_IntraParietal_0_R	IPC	0.455	837.77	0.2
Lateral_Area_7A_L	SPC	0.453	752.77	0.1
Area_8BM_L	ACMPC	0.453	1559.48	0.3
Area_V4t_L	VSVC	0.452	402.30	0.1
Area_anterior_32_prime_R	ACMPC	0.452	790.36	0.1
Area_PGi_R	IPC	0.452	2287.87	0.4
Area_7PC_L	SPC	0.451	604.46	0.1
<b>AMYGDALA_RIGHT</b>	LIMBIC	0.448	1849.55	0.3
Area_p32_prime_L	ACMPC	0.448	527.27	0.1
Area_dorsal_23_a+b_R	PCC	0.448	228.33	0.0
Area_8C_L	DLPC	0.448	2220.13	0.4
Area_OP4-PV_L	POC	0.447	1369.11	0.3
Area_posterior_10p_L	OPFC	0.443	981.51	0.2
Area_Lateral_Occipital_2_L	MT+CNVA	0.443	472.21	0.1
Area_9_Posterior_L	DLPC	0.441	1323.54	0.3
Area_OP1-SII_R	POC	0.440	526.43	0.1
Hippocampus_L	MTC	0.440	107.87	0.0
Eighth_Visual_Area_L	VSVC	0.438	470.13	0.1
Area_FST_R	MT+CNVA	0.437	610.14	0.1
Area_2_L	SMC	0.437	2140.95	0.4
Area_Lateral_IntraParietal_ventral_L	SPC	0.436	401.25	0.1

Continued on next page

## B Full Brain Region Importance Listing

Table B.1 – continued from previous page

Long Name	Cortex	Importance	Total Activation	Volume (%)
Primary_Sensory_Cortex_L	SMC	0.436	1930.49	0.4
Area_V3B_R	DSVC	0.435	464.82	0.1
Lateral_Area_7P_L	SPC	0.435	449.61	0.1
Frontal_Opercular_Area_1_L	POC	0.433	313.58	0.1
Medial_IntraParietal_Area_L	SPC	0.433	503.53	0.1
Area_6_anterior_L	PMC	0.432	1742.16	0.3
Lateral_Belt_Complex_L	EAC	0.432	241.00	0.0
Area_PF_Complex_L	IPC	0.431	1900.60	0.4
Auditory_4_Complex_L	AAC	0.430	1487.94	0.3
Area_p32_L	ACMPC	0.430	898.23	0.2
Area_STSd_posterior_R	AAC	0.430	1084.44	0.2
Primary_Auditory_Cortex_R	EAC	0.428	365.87	0.1
Area_5m_R	PLMC	0.427	677.90	0.1
Area_dorsal_32_L	ACMPC	0.426	798.10	0.2
VentroMedial_Visual_Area_2_R	VSVC	0.424	301.13	0.1
Area_5m_ventral_R	PLMC	0.424	750.13	0.1
ParaBelt_Complex_R	EAC	0.423	612.38	0.1
Anterior_Agranular_Insula_Complex_L	IFOC	0.423	420.59	0.1
Area_7PC_R	SPC	0.422	848.03	0.2
Area_10r_R	ACMPC	0.421	412.65	0.1
ProStriate_Area_R	PCC	0.420	294.67	0.1
Frontal_Eye_Fields_L	PMC	0.420	888.25	0.2
RetroSplenial_Complex_L	PCC	0.419	576.73	0.1
Frontal_Opercular_Area_2_R	IFOC	0.418	364.17	0.1
Area_5L_L	PLMC	0.418	600.58	0.1
Area_52_R	EAC	0.417	167.85	0.0
Posterior_InferoTemporal_complex_L	MT+CNVA	0.416	532.59	0.1
Area_8Av_L	DLPC	0.415	1570.55	0.3
Area_IntraParietal_2_R	IPC	0.415	670.28	0.1
Area_Posterior_24_prime_R	ACMPC	0.415	467.85	0.1
Area_PFt_R	IPC	0.414	810.37	0.2
Frontal_Opercular_Area_4_R	IFOC	0.414	854.66	0.2
Area_IFJp_L	IFC	0.413	436.43	0.1
Superior_Temporal_Visual_Area_L	TPOJ	0.412	956.71	0.2
Medial_Area_7A_R	SPC	0.412	660.54	0.1
Fourth_Visual_Area_L	EVC	0.412	2486.73	0.5
VentroMedial_Visual_Area_2_L	VSVC	0.411	179.26	0.0
Area_PF_Opercular_L	IPC	0.411	772.88	0.2
Superior_6-8_Transitional_Area_R	DLPC	0.410	605.95	0.1
Area_2_R	SMC	0.410	1751.71	0.4
Sixth_Visual_Area_R	DSVC	0.409	681.52	0.1
PreSubiculum_R	MTC	0.409	195.59	0.0
Dorsal_area_6_R	PMC	0.409	1089.00	0.2
Area_Posterior_Insular_1_R	IFOC	0.408	635.17	0.1
Third_Visual_Area_L	EVC	0.408	2220.11	0.5
Dorsal_area_6_L	PMC	0.408	899.49	0.2
Area_PF_Complex_R	IPC	0.407	1625.13	0.3
Area_PFcm_R	EAC	0.407	736.75	0.2
Area_IFSa_L	IFC	0.406	981.04	0.2
Second_Visual_Area_L	EVC	0.406	3232.07	0.7
Area_9_Middle_L	DLPC	0.406	2580.39	0.5
Area_V3B_L	DSVC	0.406	332.56	0.1
Area_9_anterior_L	DLPC	0.406	1153.06	0.2
Area_10r_L	ACMPC	0.405	625.13	0.1
IntraParietal_Sulcus_Area_1_L	DSVC	0.404	584.40	0.1
Area_5m_L	PLMC	0.404	473.91	0.1

Continued on next page

## B Full Brain Region Importance Listing

Table B.1 – continued from previous page

Long Name	Cortex	Importance	Total Activation	Volume (%)
Parieto-Occipital_Sulcus_Area_1_R	PCC	0.403	995.77	0.2
Second_Visual_Area_R	EVC	0.402	3570.61	0.7
Area_V3A_R	DSVC	0.402	1009.95	0.2
Eighth_Visual_Area_R	VSVC	0.402	570.10	0.1
Area_PHT_L	LTC	0.402	1711.54	0.4
Fourth_Visual_Area_R	EVC	0.400	2131.33	0.4
Area_V3A_L	DSVC	0.400	685.97	0.1
Lateral_Belt_Complex_R	EAC	0.398	328.65	0.1
Inferior_6-8_Transitional_Area_R	DLPC	0.396	656.23	0.1
Area_Lateral_Occipital_2_R	MT+CNVA	0.396	381.99	0.1
Medial_Superior_Temporal_Area_R	MT+CNVA	0.395	425.55	0.1
Area_TemporoParietoOccipital_Junction_3_L	TPOJ	0.393	381.83	0.1
Area_posterior_9-46v_L	DLPC	0.393	954.28	0.2
Primary_Motor_Cortex_L	SMC	0.392	3941.55	0.8
Area_Lateral_Occipital_1_R	MT+CNVA	0.392	357.02	0.1
Area_PHT_R	LTC	0.391	1189.32	0.3
Area_9_Middle_R	DLPC	0.389	2629.33	0.6
Area_6m_anterior_L	PLMC	0.387	1055.72	0.2
Area_PH_L	MT+CNVA	0.386	1431.75	0.3
Middle_Insular_Area_R	IFOC	0.386	866.61	0.2
Third_Visual_Area_R	EVC	0.386	2294.93	0.5
Area_8Ad_L	DLPC	0.385	1269.44	0.3
Area_s32_R	ACMPC	0.384	306.23	0.1
Fusiform_Face_Complex_R	VSVC	0.384	1462.86	0.3
Area_6mp_R	PLMC	0.381	1187.78	0.3
Area_PH_R	MT+CNVA	0.380	1545.99	0.3
Medial_Belt_Complex_R	EAC	0.380	409.86	0.1
Area_8B_Lateral_R	DLPC	0.380	1816.28	0.4
Area_5L_R	PLMC	0.380	694.77	0.2
Area_IFSp_L	IFC	0.379	576.48	0.1
Area_posterior_47r_L	IFC	0.379	815.30	0.2
Anterior_24_prime_L	ACMPC	0.378	490.96	0.1
Area_1_L	SMC	0.378	1791.36	0.4
Primary_Motor_Cortex_R	SMC	0.376	3468.36	0.8
Sixth_Visual_Area_L	DSVC	0.376	576.33	0.1
Superior_Frontal_Language_Area_L	DLPC	0.375	1132.57	0.3
Primary_Visual_Cortex_R	PVC	0.375	3900.82	0.9
PeriSylvian_Language_Area_R	TPOJ	0.373	1002.16	0.2
Primary_Sensory_Cortex_R	SMC	0.373	1260.84	0.3
Superior_Frontal_Language_Area_R	DLPC	0.372	831.30	0.2
Premotor_Eye_Field_L	PMC	0.371	248.89	0.1
Area_TE1_posterior_L	LTC	0.369	2057.21	0.5
Posterior_InferoTemporal_complex_R	MT+CNVA	0.368	443.72	0.1
Pirform_Cortex_L	IFOC	0.367	341.54	0.1
Primary_Visual_Cortex_L	PVC	0.365	3708.14	0.9
Anterior_Ventral_Insular_Area_R	IFOC	0.364	739.74	0.2
Dorsal_Transitional_Visual_Area_L	PCC	0.364	541.15	0.1
Anterior_IntraParietal_Area_L	SPC	0.363	769.71	0.2
Area_8Ad_R	DLPC	0.363	1305.33	0.3
Area_6mp_L	PLMC	0.358	1121.17	0.3
Area_TA2_R	AAC	0.357	715.17	0.2
Area_STSv_posterior_R	AAC	0.355	1116.28	0.3
Area_STSv_anterior_L	AAC	0.350	675.32	0.2
Area_a24_L	ACMPC	0.349	661.07	0.2
Area_8B_Lateral_L	DLPC	0.347	1148.47	0.3
Area_IFJa_L	IFC	0.347	355.94	0.1

Continued on next page

## B Full Brain Region Importance Listing

Table B.1 – continued from previous page

Long Name	Cortex	Importance	Total Activation	Volume (%)
Ventromedial_Visual_Area_3_R	V SVC	0.347	276.86	0.1
Area_43_L	POC	0.346	566.23	0.1
Frontal_Opercular_Area_1_R	POC	0.345	423.89	0.1
Rostral_Area_6_L	PMC	0.345	1176.06	0.3
Area_45_L	IFC	0.345	1111.73	0.3
Area_posterior_9-46v_R	DLPC	0.344	1272.77	0.3
Area_8Av_R	DLPC	0.344	1137.09	0.3
Area_anterior_9-46v_R	DLPC	0.342	555.53	0.1
Area_TE2_posterior_L	LTC	0.341	713.46	0.2
Area_6_anterior_R	PMC	0.340	1272.43	0.3
Area_posterior_24_L	ACMPC	0.339	449.48	0.1
<b>ACCUMBENS_RIGHT</b>	BASAL GANGLIA	0.339	377.85	0.1
Para-Insular_Area_L	IFOC	0.339	415.65	0.1
Area_10d_L	OPFC	0.338	1097.69	0.3
Area_PF_Opercular_R	IPC	0.337	685.94	0.2
Area_1_R	SMC	0.337	1284.20	0.3
Area_STSv_anterior_R	AAC	0.336	433.77	0.1
Area_9-46d_R	DLPC	0.335	1603.98	0.4
Area_STSd_anterior_R	AAC	0.332	1038.65	0.3
Frontal_Eye_Fields_R	PMC	0.332	711.42	0.2
Ventral_Area_24d_R	PLMC	0.329	396.53	0.1
Area_55b_L	PMC	0.328	580.40	0.1
Auditory_5_Complex_L	AAC	0.328	848.57	0.2
Area_6m_anterior_R	PLMC	0.328	1092.36	0.3
Area_9-46d_L	DLPC	0.327	1513.34	0.4
Ventral_Area_6_L	PMC	0.327	464.95	0.1
Area_8C_R	DLPC	0.317	1167.68	0.3
Premotor_Eye_Field_R	PMC	0.316	420.78	0.1
Anterior_Agranular_Insula_Complex_R	IFOC	0.315	150.05	0.0
Area_10d_R	OPFC	0.315	931.22	0.2
Dorsal_Transitional_Visual_Area_R	PCC	0.315	442.93	0.1
Pirform_Cortex_R	IFOC	0.314	240.55	0.1
Area_anterior_10p_L	OPFC	0.314	584.70	0.2
Area_55b_R	PMC	0.314	536.50	0.1
Area_9_anterior_R	DLPC	0.313	960.30	0.3
Hippocampus_R	MTC	0.311	44.19	0.0
Area_3a_R	SMC	0.307	558.82	0.2
Area_3a_L	SMC	0.307	537.64	0.1
Area_a24_R	ACMPC	0.305	323.74	0.1
Area_44_L	IFC	0.302	808.16	0.2
Area_TE2_posterior_R	LTC	0.301	819.01	0.2
Area_33_prime_L	ACMPC	0.295	135.28	0.0
Area_10v_L	ACMPC	0.293	1040.49	0.3
Area_posterior_10p_R	OPFC	0.290	794.90	0.2
Area_TE1_posterior_R	LTC	0.290	1895.24	0.5
Area_47l_(47_lateral)_L	IFC	0.289	717.10	0.2
Auditory_5_Complex_R	AAC	0.287	1328.39	0.4
Area_posterior_24_R	ACMPC	0.283	327.86	0.1
Area_IFSa_R	IFC	0.283	846.42	0.3
RetroSplenial_Complex_R	PCC	0.276	350.54	0.1
Area_STGa_L	AAC	0.274	535.46	0.2
Area_IFSp_R	IFC	0.273	511.83	0.2
Area_OP4-PV_R	POC	0.270	783.31	0.2
Para-Insular_Area_R	IFOC	0.269	351.00	0.1
Area_p32_prime_R	ACMPC	0.267	368.89	0.1
Area_Frontal_Opercular_5_R	IFOC	0.267	493.51	0.2

Continued on next page

## B Full Brain Region Importance Listing

Table B.1 – continued from previous page

Long Name	Cortex	Importance	Total Activation	Volume (%)
Area_10v_R	ACMPC	0.267	779.44	0.2
Area_STGa_R	AAC	0.265	379.02	0.1
Area_46_R	DLPC	0.263	1014.43	0.3
Area_IFJa_R	IFC	0.263	300.53	0.1
Rostral_Area_6_R	PMC	0.263	968.38	0.3
Area_33_prime_R	ACMPC	0.262	162.87	0.1
Area_47s_L	OPFC	0.257	553.56	0.2
Area_s32_L	ACMPC	0.255	177.84	0.1
Area_9_Posterior_R	DLPC	0.253	408.91	0.1
Area_anterior_47r_L	IFC	0.252	1020.06	0.3
Area_TE1_Middle_L	LTC	0.252	467.08	0.2
Area_IFJp_R	IFC	0.248	133.27	0.0
Area_TG_dorsal_L	LTC	0.243	2740.91	0.9
Perirhinal_Ectorhinal_Cortex_L	MTC	0.239	1162.94	0.4
Area_43_R	POC	0.235	506.25	0.2
Area_posterior_47r_R	IFC	0.234	490.61	0.2
Entorhinal_Cortex_L	MTC	0.229	246.05	0.1
Area_TE1_anterior_L	LTC	0.224	745.89	0.3
Area_TF_L	MTC	0.220	1564.45	0.6
Area_anterior_10p_R	OPFC	0.216	318.70	0.1
Polar_10p_L	OPFC	0.208	499.38	0.2
Polar_10p_R	OPFC	0.207	572.75	0.2
Area_anterior_47r_R	IFC	0.206	934.31	0.4
Area_44_R	IFC	0.201	607.51	0.3
Area_TF_R	MTC	0.199	1196.47	0.5
Area_45_R	IFC	0.191	582.58	0.3
Area_TE2_anterior_L	LTC	0.189	1036.99	0.5
Auditory_4_Complex_R	AAC	0.186	775.73	0.3
Area_25_L	ACMPC	0.184	169.01	0.1
Ventral_Area_6_R	PMC	0.183	363.18	0.2
Area_47s_R	OPFC	0.174	417.87	0.2
Perirhinal_Ectorhinal_Cortex_R	MTC	0.170	1038.56	0.5
Area_TG_dorsal_R	LTC	0.168	1824.15	0.9
Area_47m_R	OPFC	0.166	122.29	0.1
Area_13l_L	OPFC	0.164	318.89	0.2
Area_TE1_Middle_R	LTC	0.160	469.29	0.2
Area_47m_L	OPFC	0.157	161.31	0.1
Area_25_R	ACMPC	0.155	145.69	0.1
Area_TG_Ventral_L	LTC	0.152	262.92	0.1
Area_TE1_anterior_R	LTC	0.150	501.01	0.3
Entorhinal_Cortex_R	MTC	0.147	114.69	0.1
Area_11l_L	OPFC	0.147	552.08	0.3
Posterior_OFC_Complex_L	ACMPC	0.143	177.61	0.1
Area_TE2_anterior_R	LTC	0.124	898.64	0.6
Area_13l_R	OPFC	0.119	159.43	0.1
Area_47l_(47_lateral)_R	IFC	0.111	277.63	0.2
Area_11l_R	OPFC	0.107	344.99	0.3
posterior_OFC_Complex_R	ACMPC	0.106	130.42	0.1
Orbital_Frontal_Complex_L	OPFC	0.093	517.33	0.5
Orbital_Frontal_Complex_R	OPFC	0.092	571.67	0.5
Area_TG_Ventral_R	LTC	0.077	291.06	0.3