

AI Alzheimer And Dementia Classification

Ταξινόμηση Αλτσχάιμερ και Άνοιας με Τεχνητή Νοημοσύνη

A Review of Neuroimaging and Deep Learning
Approaches

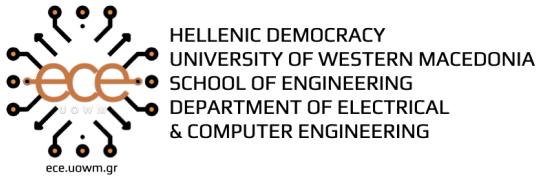
Gavriilidis Paraskevas Supervisor: Prof.

Fragulis Georgios

A thesis submitted in partial fulfillment for the
degree of MSc in Electrical Engineering

Department of Electrical and Computer Engineering
University of Western Macedonia

KOZANI / February / 2026



DECLARATION OF ORIGINALITY AND ASSUMPTION OF PERSONAL RESPONSIBILITY

I hereby declare that, in accordance with Article 8 of Law 1599/1986 and Articles 2, 4, 6 par. 3 of Law 1256/1982, this thesis entitled:

“AI Alzheimer And Dementia Classification”

as well as the electronic files and source code developed or modified within the framework of this work and explicitly mentioned in the accompanying text, which has been prepared at the Department of Electrical and Computer Engineering of the University of Western Macedonia, under the supervision of Prof. Fragulis Georgios, is exclusively the product of my personal work and does not infringe any intellectual property rights of third parties and is not a product of partial or complete copying.

The sources used are limited to the bibliographic references only. The points where I have used ideas, text, files and/or sources of other authors are clearly indicated in the text with appropriate citations and the relevant reference is included in the bibliographic references section with full description.

Copying, storage and distribution of this work, in whole or in part, for commercial purposes is prohibited. Reprinting, storage and distribution for non-profit, educational or research purposes is permitted, provided that the source is cited and this message is retained.

Inquiries regarding the use of this work for profit should be addressed to the author. The views and conclusions contained in this document express the author only.

Δήλωση Πνευματικών Δικαιωμάτων

Δηλώνω ρητά ότι, σύμφωνα με το άρθρο 8 του Ν. 1599/1986 και τα άρθρα 2,4,6 παρ. 3 του Ν. 1256/1982, η παρούσα Διπλωματική Εργασία με τίτλο

“AI Alzheimer and Dementia Classification”

καθώς και τα ηλεκτρονικά αρχεία και πηγαίοι κώδικες που αναπτύχθηκαν ή τροποποιήθηκαν στα πλαίσια αυτής της εργασίας και αναφέρονται ρητώς μέσα στο κείμενο που συνοδεύουν, και η οποία έχει εκπονηθεί στο Τμήμα Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών του Πανεπιστημίου Δυτικής Μακεδονίας, υπό την επίβλεψη του μέλους του Τμήματος κ. Φραγκούλη Γεώργιου.

αποτελεί αποκλειστικά προϊόν προσωπικής εργασίας και δεν προσβάλλει κάθε μορφής πνευματικά δικαιώματα τρίτων και δεν είναι προϊόν μερικής ή ολικής αντιγραφής, οι πηγές δε που χρησιμοποιήθηκαν περιορίζονται στις βιβλιογραφικές αναφορές και μόνον. Τα σημεία όπου έχω χρησιμοποιήσει ιδέες, κείμενο, αρχεία ή / και πηγές άλλων συγγραφέων, αναφέρονται ευδιάκριτα στο κείμενο με την κατάλληλη παραπομπή και η σχετική αναφορά περιλαμβάνεται στο τμήμα των βιβλιογραφικών αναφορών με πλήρη περιγραφή.

Απαγορεύεται η αντιγραφή, αποθήκευση και διανομή της παρούσας εργασίας, εξ ολοκλήρου ή τμήματος αυτής, για εμπορικό σκοπό. Επιτρέπεται η ανατύπωση, αποθήκευση και διανομή για σκοπό μη κερδοσκοπικό, εκπαιδευτικής ή ερευνητικής φύσης, υπό την προϋπόθεση να αναφέρεται η πηγή προέλευσης και να διατηρείται το παρόν μήνυμα. Ερωτήματα που αφορούν τη χρήση της εργασίας για κερδοσκοπικό σκοπό πρέπει να απευθύνονται προς τον συγγραφέα. Οι απόψεις και τα συμπεράσματα που περιέχονται σε αυτό το έγγραφο εκφράζουν τον συγγραφέα και μόνο.

Copyright (C) Ονοματεπώνυμο Φοιτητή & Επιβλέποντα/ες, Έτος, Πόλη

Copyright (C) Γαβριηλίδης Παρακευάς, Φραγκούλης Γεώργιος, 2026, Κοζάνη

Υπογραφή Φοιτητή:

Abstract

This thesis is a first-principles approach to the pipeline of Alzheimer’s disease and dementia classification using AI. Dementia is among the most prevalent conditions in older adults, with devastating economic implications. Projections indicate a continual rise in cases as the global population ages. This work aims to discuss each stage of the clinical process—from pathophysiology and biomarkers to the physics of imaging techniques and the machine learning approaches used in preprocessing and classification. The goal is to analyze current practices and provide a foundation grounded in the latest research. We also address key challenges and limitations, particularly class imbalance in medical datasets, and include code experiments comparing CNNs and Vision Transformers to validate and add robustness to our findings.

Keywords: Alzheimer’s Disease, Dementia, Deep Learning, MRI, PET, Neuroimaging, Classification, Convolutional Neural Networks, Medical Image Analysis

Περίληψη

Η παρούσα διπλωματική εργασία αποτελεί μια προσέγγιση από πρώτες αρχές στη διαδικασία ταξινόμησης της νόσου Αλτσχάιμερ και της άνοιας με τη χρήση τεχνητής νοημοσύνης. Η νόσος του Αλζημερ και η άνοια είναι από τις πιο διαδεδομένες παθήσεις στους ηλικιωμένους, με καταστροφικές οικονομικές συνέπειες, λόγω του κόστους περίθαλψης κατά τη διάρκεια της πάθησης. Οι προβλέψεις υποδεικνύουν συνεχή αύξηση των διαγνώσεων καθώς ο παγκόσμιος πληθυσμός συνεχίζει να γερνάει ραγδαία. Στόχος αυτής της εργασίας είναι η ανάλυση κάθε σταδίου της κλινικής διαδικασίας—από την παθοφυσιολογία και τους βιοδείκτες που μπορούν να χρησιμοποιηθούν έως τη φυσική των τεχνικών απεικόνισης και τις προσεγγίσεις μηχανικής μάθησης που χρησιμοποιούνται στην προεπεξεργασία και την ταξινόμηση. Ο στόχος είναι η ανάλυση των τρεχουσών πρακτικών και η παροχή μιας ανάλυσης βασισμένης στην πιο πρόσφατη έρευνα. Επίσης, εξετάζουμε βασικές προκλήσεις και περιορισμούς, ιδιαίτερα την ανισορροπία κλάσεων στα ιατρικά σύνολα δεδομένων, και περιλαμβάνουμε πειράματα κώδικα που συγχρίνουν συνελικτικά νευρωνικά δικτυα και άλλες τεχνικές για την επιβεβαίωση και την προσθήκη αξιοπιστίας στην έρευνα και τα ευρήματα μας.

Λέξεις Κλειδιά: Νόσος Αλτσχάιμερ, Άνοια, Βαθιά Μάθηση, Νευροαπεικόνιση, Ταξινόμηση

Acknowledgements

I would like to express my sincere gratitude to my supervisor, Prof. Fragulis Georgios, for his guidance and support throughout this research.

Contents

Abstract	iv
Acknowledgements	vi
1 Introduction	1
1.1 Contextual Background	1
1.2 Research Significance	2
1.3 Objective of the Review	3
2 Overview of Dementia and Alzheimer's disease imaging	5
2.1 Magnetic Resonance Imaging (MRI)	5
2.1.1 Spin	6
2.1.2 Properties of Spin	6
2.1.3 Hydrogen Nuclei in MRI	6
2.1.4 Relaxation	7
2.1.5 T1- and T2-Weighted Imaging	7
2.2 The Clinical Use of Structural MRI in Alzheimer Disease	8
2.3 Atrophy as a Neurodegeneration Marker	9
2.4 Alzheimer Disease Criteria and MRI	10
2.5 Computed Tomography (CT)	10
2.5.1 Mechanics of a CT Scan	11
2.5.2 Historical Development of CT	11
2.5.3 Basic Physical Principles of CT	11
2.5.4 Data Acquisition	12
2.5.5 Image Reconstruction	13
2.5.6 CT Numbers / Hounsfield Units	14
2.6 PET / PET-CT	14
2.6.1 Basic Physics	14
2.6.2 Detection of Annihilation Radiation	16
2.6.3 PET Spatial Resolution Limitations	16

2.6.4	PET-CT	16
2.6.5	The Use of PET in Alzheimer’s Disease	17
2.6.6	Processes Assessed by PET	17
3	Key Image Datasets for Dementia and Alzheimer’s Disease	19
3.1	ADNI	19
3.2	AIBL	20
3.3	OASIS	20
4	Pre-Processing and Feature Extraction Techniques	23
4.1	Intensity Normalization	23
4.1.1	Introduction	23
4.1.2	Necessity and Significance	24
4.1.3	Goals and Principles	24
4.1.4	Categories of Intensity Normalization Methods	25
4.1.5	Impact of Intensity Normalization	30
4.1.6	Challenges and Considerations	30
4.2	Denoising	31
4.2.1	Non-Local Means (NLM)	32
4.2.2	BM3D	32
4.2.3	Deep Learning Era	33
4.3	Skull Stripping	33
4.3.1	Skull Stripping Methods	35
4.3.2	Morphology-Based Methods	35
4.3.3	Intensity-Based Methods	36
4.3.4	Deformable Surface-Based Methods	36
4.3.5	Atlas / Template-Based Methods	37
4.3.6	Hybrid Methods	37
4.3.7	Deep Learning-Based Methods	38
4.3.8	Comparative Analysis	39
4.3.9	Multi-Site Dataset Considerations	39
4.3.10	Implications for Alzheimer’s Disease Research	40
4.3.11	Strengths and Weaknesses of Skull Stripping Approaches	40
4.3.12	Primary Challenges	41
4.4	Voxel-Based Morphometry (VBM)	41
4.4.1	Evolution of VBM	42
4.4.2	Alzheimer’s Disease and Atrophy Patterns	42
4.4.3	VBM as a Prognostic Tool	43
4.4.4	Differences Between Converters	43
4.4.5	Meta-Analytic Power of VBM	45

4.4.6	Subtypes of AD	46
4.4.7	Limitations of VBM	48
4.4.8	Machine Learning and VBM	49
4.4.9	Deep Learning	50
5	Classification Techniques for Dementia Image Analysis	51
5.1	Foundational Machine Learning Paradigms in Dementia Classification .	51
5.1.1	Support Vector Machines (SVM)	51
5.1.2	Ensemble Methods: Decision Trees and Random Forests	52
5.1.3	The Curse of Dimensionality: PCA and LDA	53
5.2	The Deep Learning Revolution	53
5.2.1	Convolutional Neural Networks (CNNs)	53
5.2.2	Generative Models: GANs	54
5.2.3	Vision Transformers (ViTs)	55
5.3	Hybrid Classification Approaches	55
6	Classification Performance and Accuracy Metrics	57
6.1	Evaluation Metrics	57
6.1.1	Accuracy, Sensitivity, Specificity	57
6.1.2	Precision and Recall	58
6.1.3	Model Validation Techniques	58
6.1.4	ROC Curve	58
6.1.5	Matthews Correlation Coefficient	59
7	Challenges in Dementia Image Analysis and Classification	61
7.1	Imbalanced Data	61
7.1.1	Threshold Metrics	61
7.1.2	Ranking Methods and Metrics	62
7.2	Limited Datasets	63
7.2.1	Context and Challenges	63
7.3	Small Sample Sizes	63
7.4	Image Quality and Variability	64
7.5	Interpretability in Medical Applications	64
7.6	Computational Costs	64
7.7	Data Distillation	65
8	Recent Advances And Innovations	67
8.1	Multi-Modal Integration	67
8.2	Explainable AI	68
8.3	Definitions and Taxonomy	68

8.3.1	LIME	69
8.3.2	SHAP	69
8.3.3	Other Approaches	70
8.4	Sanity Checks	70
8.5	Related Work	70
8.6	Data Augmentation & Transfer Learning	71
8.6.1	Transfer Learning	73
9	Gaps in Current Literature and Future Research Directions	75
9.1	Generalizability and Standardized Testing	75
9.2	Interpretable Models	75
9.2.1	Trust and Clinical Adoption	76
9.2.2	Reducing Bias and Increasing Fairness	76
9.2.3	Better Model Development and Debugging	76
9.3	Detection at Different Stages	76
10	Limitations	79
10.1	Data Leakage	79
10.2	The Accuracy Paradox and Class Imbalance	80
10.3	Domain Shift and Generalizability	81
10.3.1	Sources of Domain Shift in MRI Data	81
10.4	Shortcut Learning and Region of Interest Bias	82
10.4.1	Skull-Stripping Artifacts	82
10.4.2	Biological Plausibility of Learned Features	82
10.5	Ground Truth Uncertainty and Diagnostic Heterogeneity	83
10.6	Disease Heterogeneity as an Intrinsic Limitation	85
10.7	Transparency, Reproducibility, and the Black Box Problem	85
10.8	Limitations of Single-Modality Analysis	86
10.9	Summary of Limitations	86
11	Conclusion	87
Appendices		89
Pathophysiology of Alzheimer		91
		93
Source Code		94

List of Figures

2.1	(a) The first CT image produced at Atkinson Hospital. (b) A modern CT scan from an advanced scanner [22].	12
2.2	Illustration of X-ray attenuation through a sample composed of multiple materials. The final intensity corresponds to the cumulative effect of all attenuation coefficients along the beam path.	13
2.3	Hounsfield scale and representative values for different tissues [22].	14
4.1	Brain image depicting artifacts present in MR images and different tissue types. Adapted from [31].	25
4.2	Original image (a); Image obtained by Gaussian filtering with increasing σ values (b-d).	27
4.3	White Stripe Normalization concept.	28
4.4	Examples of automated skull stripping from [51].	34
4.5	Overview of skull-stripping categories [53].	35
4.6	Direct Comparisons using VBM show distinct patterns of atrophy. While Typical AD and FTLD patients showed loss in the hippocampal area , atypical subjects showcased severe damage in the left putamen and tempoparietal cortex. [71]	45
4.7	"The visualizations illustrate conjunction analyses , which show overlapping damage across different groups. The first one shows a comparison between AD typical and atypical damage to healthy controls. The second shows damaged that is shared between three specific variants: aphasic dementia, Corticobasal Syndrome (CBS), and behavioral variant Frontotemporal Dementia (bvFTD) " [71]	46
4.8	VBM shows distinct patterns of atrophy in between different gouprs. Typical AD and FTLD show more loss in the hippocampal area than atypical AD patients. Atypical patients show more severe damage in the left putamen (surpassing Typical AD) and in the tempoparietal cortex (surpassing atypical AD) . [71]	47
4.9	Shows the distinct atrophy for three clinical subgroups. Aphasic Dementia, Corticobasal Syndrome (CBS) and behavioral variant Frontotemporal Dementia (bvFTD) compared to healthy controls. [71]	47

4.10 "The images map statistically significant gray matter atrophy of the a-MCI patients compared to healthy controls. The top one shows damage on the brain outer surface , whiel the bottom show deep tissue loss in the cingulate cortex and medial temporal lobes" [72]	48
4.11 The figure shows specific brain areas in the aMCI-P group that showcase larger atrophy than those in the aMCI-S group. [72].	49
5.1 Example Grad-CAM heatmap highlighting salient regions [104].	54
8.1 Showcases all the available image augmentation techniques. Courtesy of [158].	73
8.2 Categorizations of transfer learning [161].	74
10.1 Class distribution showing the imbalance in the dataset.	80
10.2 Grad-CAM visualization showing original image, heatmap, and overlay.	84

List of Tables

4.1	Comparative strengths and weaknesses of skull-stripping methodologies.	40
4.2	Summary of regional atrophy patterns in AD and linked VBM-related figures.	44
6.1	Confusion Matrix showing Actual vs. Predicted values	57
6.2	Part 1: The Confusion Matrix and Row-based Rates (Sensitivity/Specificity)	59
6.3	Part 2: Derived Performance Metrics (Precision, Accuracy, etc.)	59
8.1	Image Augmentation Techniques	71

Chapter 1

Introduction

1.1 Contextual Background

The evolution of neuroimaging has transformed the diagnostic landscape for Alzheimer’s disease. Multiple imaging techniques, including MRI, PET, and CT, are used to support clinical evaluation.

Structural Magnetic Resonance Imaging (MRI) excels at quantifying brain atrophy, particularly in the hippocampus, where volumetric reductions serve as predictors of progression from Mild Cognitive Impairment (MCI) to AD.

Positron Emission Tomography (PET), with amyloid- and tau-specific tracers, enables visualization of amyloid- β ($A\beta$) plaques and tau pathology [1] (for more on AD pathology, see Appendix 1). Multimodal approaches integrating multiple imaging modalities provide superior sensitivity and specificity for early diagnosis and longitudinal disease tracking. AI- and deep-learning-enhanced multimodal frameworks now fuse neuroimaging data with genetic, clinical, and neuropsychological variables, achieving diagnostic accuracies exceeding 95% in diverse cohorts [2].

Neuroimaging allows for a non-invasive way to produce quantifiable and objective measurements. The combination of multiple modalities can help to detect structural changes, hypometabolism , connectivity loss and even tau or amyloid deposition. At a time where new treatments emerge at the preclinical stage the use of neuroimaging can provide objective information about the progression stages and the efficacy of the treatment. [3]

The latest innovations in imaging technologies and the impact on the diagnosis of AD have shifted diagnostic criterial. According to the National Institute of Aging (NIA) and the Alzheimer’s Association staging is based on amyloid , tau and neurodegeneration

profiles. Core-1 biomarkers include low CSF A β 42/40 ratios or positive amyloid PET, with Core-2 Biomarkers tracking progression - elevated tau PET or plasma p-tau217) . MRI protocols like diffusion tensor imaging and quantitative susceptibility mapping are used for entorhinal-hippocampal circuits where early tau deposition can begin even up to two decades early. [3]

For an expanded overview of AD pathology, see Appendix 1.

AD develops silently for 15-20 years, making early identification and diagnosis crucial. The biomarkers of tau PET and plasma p-tau217 trace early tau spread and predict neurodegeneration. Cognitive dysfunction for patients with amnesic MCI progresses to AD at a rate of 12-15% annually compared to 1-2% for cognitively normal adults [4] . Thus imaging-based tracking of amyloid and tau staging (Braak Stages) helps identify individuals at highest risk. [5]

1.2 Research Significance

Dementia poses a huge challenge for aged individuals , even more so with a very high percentage of aging population [6]. In the US , about 7.2 million individuals are living with AD as of 2025 , a number expected to double by 2060 (13.8 million) . The disease additionally creates a massive economic burden on healthcare systems and caregiving costs , reaching 384 billion dollars , far exceeding other chronic conditions. These trends highlight the need for improved diagnostic capabilities and treatment strategies. [7]

Understanding lifetime dementia risk is essential for prevention and resource planning. One study estimated a 42% baseline risk from age 55 to 95 for developing dementia, with higher prevalence observed in women (45–60%), Black adults, and carriers of the *APOE ε4* allele [8].

Another interesting aspect that can help prevention and enhance resource planning estimates is the understanding of lifetime dementia. In one study a 42% baseline risk from age 55-95 was found , with higher prevalence in women (45–60%), Black adults, and carriers of the *APOE ε4* allele [8].

Roughly 697 cases of dementia exist per 10.00 people over 50 as reported in [9] . From those 324 cases are of AD and 115 from vascular dementia. Also in this study it is mentioned that women present greater rates and that prevalence doubles every two years. As is understood population aging continues to be the most important cause for increased rates of dementia , despite trends not deviating significantly from historical statistics.

In order to enhance patient care, clinical decision-making, and research, highly precise

diagnostic techniques are necessary for dementia in general and Alzheimer's disease (AD) in particular. In order to diagnose AD-related brain disease earlier and develop more individualized treatment plans, neuroimaging is a key technique [10].

1.3 Objective of the Review

The goal of this review is to summarize existing work and provide an intuitive explanation of each step in the pipeline from the imaging tools to the algorithms used for classification, discuss existing image datasets available and highlight gaps for future research.

Chapter 2

Overview of Dementia and Alzheimer's disease imaging

2.1 Magnetic Resonance Imaging (MRI)

Magnetic Resonance Imaging (MRI) is a technique based on NMR principles. Even though the classical view is not fully correct, it is quite useful as an introduction to the technique. Atomic nuclei possess the property of spin, which can be imagined as a spinning top with that property being its angular momentum. When a large magnetic field is applied to the nuclei, they tend to align with the field while continuing to spin. Unlike other modalities such as CT scans, which can subject the patient to ionizing radiation, MRI does not cause any exposure-related hazards since it uses non-ionizing electromagnetic radiation [11]. Excitation in the case of MRI happens through radio-frequency (RF) pulses specifically tuned to the Larmor frequency. At this frequency, the nuclei resonantly absorb energy, causing the net magnetization vector to tip away from the applied magnetic field. The process analyzed classically is analogous to an oscillating body that performs a periodic movement (forced oscillator), receiving energy at a matching frequency to enhance the amplitude of its motion (e.g., a swing that gets pushed at the right intervals). It should be noted that pulse sequences can vary, but more detail will be discussed later. Once the pulse sequence stops, the nuclei relax back toward equilibrium and while doing so they emit a detectable signal captured by receiver coils [11–14]. Capturing, however, does not happen all at once. To achieve localized imaging of tissue, magnetic field gradients are employed. They alter the magnetic field so that it varies spatially, and in doing so only selected tissue precesses at the frequency matching the RF pulse. Overall, MRI offers a non-invasive way that possesses great accuracy in capturing the underlying tissue.

2.1.1 Spin

Atomic nuclei possess an intrinsic property known as *spin*. Spin occurs in multiples of $1/2$ and may be positive or negative. Not all particles possess spin; it appears in particles with an odd number of protons and neutrons ($P + N = \text{odd}$) [13].

2.1.2 Properties of Spin

When a particle with net spin is placed within a magnetic field of strength B , it can absorb a photon of frequency ν . This frequency depends on the gyromagnetic ratio γ of the particle [15] :

$$\nu = \gamma B.$$

For hydrogen atoms, the gyromagnetic ratio is [12]

$$\gamma = 42.58 \text{ MHz/T.}$$

MRI uses nuclei with either an unpaired proton or an unpaired neutron, allowing them to possess a net spinning charge—i.e., angular momentum. Because spin is associated with electric charge, these nuclei generate a magnetic field and behave as magnetic dipoles.

2.1.3 Hydrogen Nuclei in MRI

Hydrogen nuclei are the most abundant MR-visible nuclei in the human body (primarily in water, H_2O) and consist of a single proton with non-zero spin, allowing them to act as magnetic dipoles under a strong magnetic field [13, 16]. When a magnetic field is applied, all hydrogen spin angular momentum vectors tend to align either parallel (spin-up) or anti-parallel (spin-down) to the field. Since spin projection comes in quantized these energy states are distinct. Nuclei can transition between different energy states by absorbing or emitting energy , which in the case of MRI is provided by RF pulses at the (Larmor) resonance frequency.

The axes of spinning protons do not align perfectly parallel to the applied field; rather, they *precess*. The frequency of precession is called the *Larmor frequency*, which is proportional to the magnetic field strength:

$$\omega = \gamma B.$$

After applying the magnetic field, the sum of all nuclei yields a weak net magnetic moment or magnetization vector (MV) parallel to the field. Applying an RF pulse at the Larmor frequency in the presence of a magnetic gradient selects a slice of tissue and excites the nuclei, causing them to absorb energy and rotate into the plane of the RF pulse.

Longer RF pulses produce greater rotation angles. If the pulse has sufficient intensity, it flips the magnetization vector into the transverse (XY) plane, creating a *90° flip angle* in which all protons precess in phase.

At this moment, an RF signal is induced in a receiver coil. This signal depends on the presence and molecular environment of hydrogen atoms. Tightly bound hydrogen, such as in bone, produces weak or no usable signal, whereas loosely bound hydrogen in soft tissues produces a strong signal. The concentration of loosely bound hydrogen is known as the *proton density* or *spin density*.

2.1.4 Relaxation

When the RF pulse is turned off:

- The excited nuclei return to their lower-energy spin state, emitting energy. This decay of transverse magnetization is detected as the *free induction decay* (FID).
- Nuclei realign with the main magnetic field through energy transfer to surrounding molecules. This is the basis of *longitudinal* or *spin-lattice* relaxation.

The time constant describing the recovery of longitudinal magnetization is called the *T1 relaxation time*. Typical values are approximately 500 ms for short T1 tissues. The decay of transverse magnetization, governed by loss of phase coherence among precessing protons, is described by the *T2 relaxation time* (e.g., ~ 80 ms for many tissues).

2.1.5 T1- and T2-Weighted Imaging

Relaxation properties allow distinctions between tissue types. MRI commonly uses two contrast mechanisms:

T1-weighted imaging uses short repetition times(TR) and is detecting T1 relaxation. Tissues that appear bright have short T1 values (e.g. fat), while those with longer T1 values appear darker. This modality is specifically useful for visualizing anatomical detail due to its characteristic of high spatial resolution.

T2-Weighted Imaging uses long repetition times (TR) and long echo times (TE) . Tissues that are water rich or are abnormal usually have long T2 values and appear bright. [12]

By changing TR, TE values and the spatial localization of excitation using gradient fields, MRI can provide great anatomical detail while maintaining sensitivity to underlying pathological changes. These contrast mechanisms are the foundation of MRI and its use in clinical settings.

2.2 The Clinical Use of Structural MRI in Alzheimer Disease

In the case of Alzheimer's disease, we use the features of structural imaging that MRI can produce, and it has become an integral part of the clinical assessment process. New data from clinical studies showing changes in structural markers from preclinical to overt stages of the disease are reshaping the diagnostic landscape and influencing future diagnosis and treatment.

In AD Alzheimer can help find structural changes that take place in the brain of a patient , and it has become routine , with MRI being one of the most commonly used modalities among clinical settings . As has been mentioned before structural brain changes can take place even up to 20 years before without cognitive decline symptoms and become more obvious and accelerate decline as time goes on. One of the earliest and important signs is shrinkage in the medial temporal areas of the brain, which is used to now diagnose MCI (Mild Cognitive Impairment) . MRI can also allow to distinguish between AD and other types of dementia. Additionally , shrinkage in the hippocampus , or in the brain matter as a whole can be a sign that neurodegeneration is already underway , marking an important step in the progress of the disease.

MRI is also increasingly used in research to track the effectiveness of disease-modifying drugs based on the structural markers mentioned above. While MRI plays an important role in tracking disease progression and determining atrophy in specific brain regions, it is essential to point out that many other imaging and non-imaging techniques are used in clinical assessment.

The critical utility of MRI in the clinical domain has been discussed, particularly for AD and dementia. By adopting standardization protocols during acquisition or as part of further processing after the capture of a scan, we can ensure robustness and reliability. Any subsequent algorithmic task will thus be further improved [17]. The pattern of AD development specifically follows a progressive accumulation of abnormal proteins (e.g.,

$\text{A}\beta 42$), which in turn lead to synaptic, neuronal, and axonal damage. This pattern of accumulation occurs many years before any cognitive dysfunction becomes evident in patients. The progression in AD follows quite a typical pattern. It first presents itself in the medial temporal lobes (entorhinal cortex and hippocampus), which then spreads to inflict further neocortical damage [18, 19]. This delay, in fact, suggests that the toxic effects induced by abnormal protein deposition have to reach a certain threshold to initiate noticeable cognitive symptoms. For example, in amnesic Mild Cognitive Impairment (aMCI), deficits must emerge across multiple cognitive domains before the full AD diagnostic criteria are fulfilled. Additionally, the development of disease-modifying drugs that can decrease the rate of decline—or potentially halt it—requires early identification of at-risk individuals to evaluate treatment efficacy. Consequently, the value of early diagnosis has increased substantially. Studies have already shown that significant correlations exist between damage in regions such as the hippocampus and entorhinal cortex and the likelihood of MCI-to-AD progression. This constitutes the primary objective of this work: evaluating whether algorithmic approaches can reliably predict this conversion, which would have significant implications for treatment planning and disease management [17].

2.3 Atrophy as a Neurodegeneration Marker

Atrophy is a consequence of progressive neurodegeneration. Structural changes can be mapped onto the stages of tangle deposition based on the Braak staging system, as well as onto specific neuropsychological deficits. The earliest signs occur in the perforant pathway, producing memory impairments. Later changes in the parietal and frontal neocortices correspond to language deficits, visuospatial impairments, and behavioral changes.

Changes in structural measures such as whole-brain volume, entorhinal and hippocampal volume, temporal lobe volume, and ventricular enlargement can be identified using MRI, and these measures may serve as potential biomarkers. To be used clinically, such biomarkers must be identifiable across all patients and provide clear distinctions between disease stages.

In the asymptomatic stage, amyloid markers provide indirect evidence of disease pathology, whereas after the onset of Mild Cognitive Impairment (MCI), structural changes are more sensitive indicators [17].

2.4 Alzheimer Disease Criteria and MRI

The diagnostic approach to AD as proposed in the International Working Group (IWG-2) Criteria has shifted from a clinicopathological to a clinicobiological one. This shift reflects the incorporation of biomarkers into the clinical diagnosis process. The new framework requires two criteria to be satisfied: an appropriate clinical phenotype (typical or atypical) and the presence of a biomarker consistent with AD pathology.

The criteria are complete and can capture all disease stages, including typical AD, atypical variants, mixed AD, and preclinical states (asymptomatic at-risk and presymptomatic AD).

Pathophysiological markers are restricted to those that directly indicate amyloid or tau pathology and must be present for a diagnosis of AD.

Specific Criteria

- A CSF profile showing decreased $A\beta_{42}$ alongside increased T-tau or P-tau concentrations.
- Elevated tracer retention on amyloid PET imaging.
- Presence of an autosomal dominant AD mutation.

Topographical biomarkers such as volumetric MRI (e.g., hippocampal atrophy) and FDG-PET have been removed from the core diagnostic algorithm, as they lack sufficient pathological specificity for AD detection. Their new role is to monitor disease progression over time [20].

Despite this change, hippocampal atrophy remains one of the most established and validated MRI markers of AD. In vivo measurements correlate with Braak staging and neuronal counts, with volume reductions varying by disease stage (e.g., 15–30% in mild dementia, 10–15% in MCI with mild dementia) [17].

2.5 Computed Tomography (CT)

CT scans are a useful tool in the diagnosis of Alzheimer's disease (AD), providing images of anatomical structures in brain regions such as the medial temporal lobe, where atrophy is considered a marker for AD conversion. However, during the diagnostic process of cognitive complaints or deficiencies, MRI scans are generally preferred due to their superior soft-tissue contrast. CT scans are used primarily when MRI presents contraindications (e.g., pacemakers) [10].

The use of CT in combination with nuclear imaging techniques has increased, first with the introduction of PET/CT and later SPECT/CT in AD diagnosis. These combined modalities emerged because the anatomical localization of functional abnormalities in nuclear imaging alone was often imprecise. Adding CT provides accurate anatomical reference and resolves localization issues, provided the modalities are properly coregistered [21].

Given this context, the review of the basic principles of CT will remain brief, since this work does not focus on CT images directly, nor do CT scans provide structural detail comparable to MRI. Nevertheless, CT remains an important imaging method and should be acknowledged.

2.5.1 Mechanics of a CT Scan

2.5.2 Historical Development of CT

The discovery of X-ray radiation is attributed to Wilhelm Conrad Röntgen, who performed the first cathode-tube experiments on November 8, 1895. X-rays possess special properties: they can energize atoms (producing photons via fluorescence), and they can penetrate opaque materials.

Initially, X-ray projection imaging captured a two-dimensional image of a three-dimensional object. However, important structural information was lost due to overlapping tissues, while low-contrast regions were difficult to distinguish. Additionally, scattering produced noise and degraded image quality.

The term *computed tomography* reflects its two essential components: “computed,” referring to numerical reconstruction, and “tomography,” meaning cross-sectional slicing. Modern CT scanners use X-ray energies between 100–150 kV.

Image reconstruction in CT relies on the mathematical foundations of the Radon transform and its inverse. Radon demonstrated that a 2D image can be reconstructed from a set of projections collected at multiple rotation angles.

2.5.3 Basic Physical Principles of CT

X-ray Attenuation

As X-rays pass through tissue, some photons are absorbed or scattered while others reach the detector. Different tissue types attenuate X-rays differently—bone absorbs much more than soft tissue, for instance. This selective absorption is what creates contrast in the final image. The phenomenon is described by the exponential attenuation

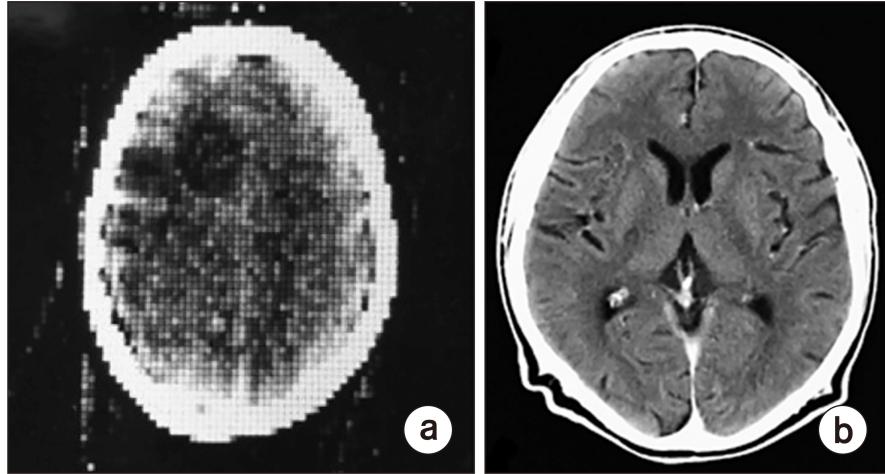


FIGURE 2.1: (a) The first CT image produced at Atkinson Hospital. (b) A modern CT scan from an advanced scanner [22].

law:

$$I_x = I_0 e^{-\mu x}, \quad (2.1)$$

- I_0 : initial intensity
- I_x : transmitted intensity at the detector
- μ : linear attenuation coefficient (tissue-dependent)
- x : material thickness

For multiple materials along the path:

$$I = I_0 \exp [-(\mu_1 x_1 + \mu_2 x_2 + \dots)], \quad (2.2)$$

known as the Lambert–Beer law.

Attenuation can also be expressed as a line integral:

$$\ln \left(\frac{I}{I_0} \right) = - \int \mu(s) ds. \quad (2.3)$$

2.5.4 Data Acquisition

The process of data acquisition in a CT scan begins by placing the patient or sample inside the scanner. A rotating X-ray source then emits radiation as it orbits the sample, while a detector on the opposite side captures the transmitted signal. These

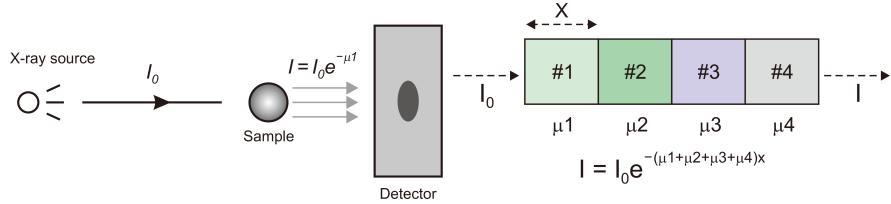


FIGURE 2.2: Illustration of X-ray attenuation through a sample composed of multiple materials. The final intensity corresponds to the cumulative effect of all attenuation coefficients along the beam path.

signal intensities are digitized by the Data Acquisition System (DAS) and prepared for reconstruction.

Acquisition criteria:

- Projections must be collected over many angles (typically 360° or 180° with symmetry).
- Each projection must fully include the object.
- The object must remain still.

2.5.5 Image Reconstruction

The goal of CT reconstruction is to compute the 2D attenuation map from the measured 1D projections (line integrals).

For an object described by a function $f(x, y)$, the Radon transform is:

$$p(s, \varphi) = \mathcal{R}f(x, y). \quad (2.4)$$

To convert (x, y) to rotated coordinates (s, u) :

$$\begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \cos \varphi & -\sin \varphi \\ \sin \varphi & \cos \varphi \end{pmatrix} \begin{pmatrix} s \\ u \end{pmatrix}. \quad (2.5)$$

Thus, the projection data are given by:

$$p(s, \varphi) = \int f(s \cos \varphi - u \sin \varphi, s \sin \varphi + u \cos \varphi) du. \quad (2.6)$$

Early CT systems applied a high-pass filter to sharpen images. Today, reconstruction methods include:

- matrix inversion,

- iterative algorithms,
- Fourier techniques,
- filtered backprojection (FBP),
- 3D Radon-based approaches.

The most common modern method is based on the *central slice theorem*, which links the Radon transform to the 2D Fourier transform.

2.5.6 CT Numbers / Hounsfield Units

A CT image consists of voxels (3D volume elements). Each pixel's intensity reflects the mean attenuation coefficient in its voxel. CT uses a 12-bit grayscale (4096 levels). The Hounsfield Unit (HU) scale normalizes attenuation values relative to water:

$$HU = 1000 \cdot \frac{\mu_{\text{pixel}} - \mu_{\text{water}}}{\mu_{\text{water}}} \quad (2.7)$$

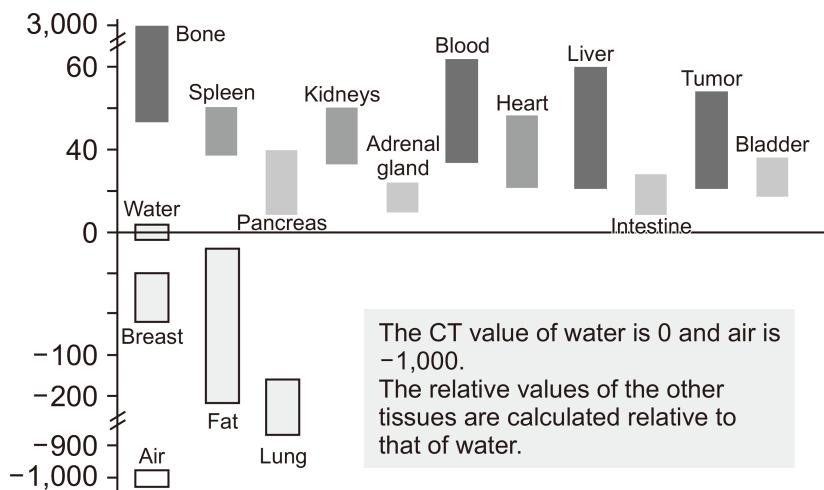


FIGURE 2.3: Hounsfield scale and representative values for different tissues [22].

2.6 PET / PET-CT

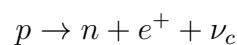
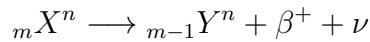
Positron Emission Tomography (PET) is an analytical process in which compounds labeled with radioisotopes are used as molecular probes to image and measure biochemical processes *in vivo*.

2.6.1 Basic Physics

Positron emission is based on proton-rich nuclei, called “emitters,” which are unstable. To stabilize themselves, they may either release excess protons and gain neutrons, or

capture electrons. The first process is known as *positron emission*, and the second as *electron capture (EC)*. Both processes are isobaric decays, meaning the mass number remains unchanged between the parent and daughter nuclei. In nuclei with low atomic weight, positron emission is more prevalent, while in heavier nuclei electron capture dominates.

It can be described by :



The nucleus by having too many protons , undergoes beta-plus decay . In this process a proton converts into a neutron, emitting a positron and a neutrino. This results in a new element with ;

- Same mass number (protons + neutrons unchanged)
- Atomic Number decreased by one (one fewer proton)

The neutrino leaves without interacting with any of the tissue, whereas the positron, being charged is slowed by scattering in the tissue.

The distance traveled (the *range*) depends on its energy:

$$E_{\text{positron}} + E_{\text{neutrino}} = \text{transition energy} - 1022 \text{ MeV}$$

As the positron slows, it eventually annihilates with an electron, producing either:

- a pair of 511 keV photons (direct annihilation), or
- a short-lived positronium, which also decays into two 511 keV photons.

The photons that are produced are emitted nearly 180° apart. However due to residual momentum , the emission angle varies slightly around 180° . Their energy is equal to that of a positron and electron at rest mass while they are more officially named "**annihilation photons** .

2.6.2 Detection of Annihilation Radiation

To detect the annihilation process the system identifies events through a coincidence detection window (also known as electronic collimation). This process works by creating a timing window (3-15ns) within which the photons need to strike two opposing detectors. If both of them arrive within this interval , the system records the interaction that occurred , somewhere along the line of response (LOR) , which is the direct path connecting the two detectors.

2.6.3 PET Spatial Resolution Limitations

1. **Positron Range:** Because the positron travels a finite distance before annihilation, the detected photons do not originate exactly where the positron was emitted. This introduces an intrinsic spatial uncertainty.
2. **Non-colinearity of Annihilation Photons:** The photons are not emitted at exactly 180° . Small deviations (typically $< 0.5^\circ$) cause blurring. The effect depends on detector ring diameter—about 1 mm for a 50 cm ring and 2 mm for a 90 cm whole-body system.
3. **Detector Size:** The intrinsic spatial resolution depends on the size of the individual detector crystals used in modern PET scanners, which typically consist of small scintillator arrays coupled to larger photodetectors [23].

2.6.4 PET-CT

PET and Anatomic Imaging

A major limitation of PET imaging is its low spatial resolution and lack of anatomical detail, making it difficult to accurately localize lesions with abnormal radiotracer uptake. Distinguishing physiological from pathological uptake can be challenging.

Initially, non-integrated scanners were used: PET and CT images were acquired separately and then aligned visually or using fusion algorithms. However, differences in patient position and organ motion significantly reduced accuracy. The dual-scan approach is also costly, time-consuming, and uncomfortable for patients.

Integrated PET/CT

PET scanner was introduced in 1998. It has the ability to show both body structure through the CT scan but also body function by utilizing the PET scan. This not only helps doctors make more accurate decisions but also helps patients to make the process faster , easier and cheaper for them. [21, 23].

2.6.5 The Use of PET in Alzheimer's Disease

The pathology of AD is defined by amyloid plaques in the first place (according to Braak staging) , tau tangles (neurofibrillary tangles) concentration , activated microglia, neurotransmitter changes, and neuronal loss.

The characteristic pathology of Alzheimer's Disease (AD) includes the accumulation of amyloid plaques in the brain and intracellular hyperphosphorylated tau protein in the form of neurofibrillary tangles; activation of microglial cells; alterations in neurotransmitter levels; and neuronal death. Changes in the brain can occur for a long time prior to the onset of clinical manifestations. Therefore, it is now important to consider CSF analysis and brain imaging as early biomarkers and useful tools for tracking AD progression. [20, 24].

New PET imaging technologies allow the identification of AD in prodromal stages and support the development and monitoring of disease-modifying treatments. PET enables measurement of multiple functional processes in the brain, including metabolism and neurotransmitter activity, providing region-specific insights that strengthen diagnostic accuracy and treatment evaluation.

2.6.6 Processes Assessed by PET

Brain Glucose Metabolism

To measure glucose metabolism, PET uses the tracer 2-fluoro-2-deoxy-glucose (F-FDG). A decline in glucose metabolism occurs long before clinical symptoms appear. This decline is region-specific, most prominently affecting the parietotemporal, frontal, and posterior cortices [25].

Patients with AD can be diagnosed with up to 90% sensitivity using FDG-PET [26]. Differentiation from other dementias is more difficult.

Longitudinal PET studies show:

- Conversion from healthy to MCI is best predicted by **medial temporal** glucose hypometabolism.
- Conversion from MCI to AD is best predicted by hypometabolism in the **posterior cingulate cortex**.

Functional MRI (fMRI) also helps distinguish patterns of brain activity. Its basis lies in the magnetic properties of deoxyhemoglobin and the fact that blood flow increases more than oxygen metabolism during neural activation. This produces a subtle MR

signal increase known as the *blood oxygenation level-dependent (BOLD)* effect. Modern fMRI achieves spatial resolution near 1 mm and temporal resolution around 1 s [27].

Chapter 3

Key Image Datasets for Dementia and Alzheimer's Disease

3.1 ADNI

The Alzheimer's Disease Neuroimaging Initiative (ADNI) is one of the largest and most significant research initiatives in the field of Alzheimer's disease. The ADNI is an example of a public – private collaboration, funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering (NIBIB), major pharmaceuticals such as Pfizer, Eli Lilly, Merck, GlaxoSmithKline and AstraZeneca, and non-profit organizations including the Alzheimer's Association.

The ADNI has been established to develop and validate biomarkers, and also to validate neuroimaging markers for detecting cognitive and functional decline in the elderly population, as well as those suffering from Mild Cognitive Impairment (MCI), and cognitively normal controls. The ADNI will collect data that are related to genetics, cerebrospinal fluid (CSF), clinical information, cognitive function, and a wide variety of neuroimaging techniques including, but not limited to, structural MRI, FDG-PET, amyloid PET, tau PET, and Diffusion Tensor Imaging (DTI).

This will allow researchers to have access to a diverse set of data regarding changes in brain structure and changes in cognitive function. The primary focus of the ADNI is to accelerate work in early detection of Alzheimer's disease, track the progression of the disease, and evaluate the effectiveness of treatments.

The ultimate goal of the ADNI is to create standardized MRI/PET imaging datasets from large cohorts at various centers across the U.S., and Canada, to identify sensitive biomarkers for the disease, and to determine the efficacy of potential therapeutic

interventions.

One other major attribute of ADNI is the establishment of a publicly available data repository containing all their collected longitudinal information on clinical outcomes , biomarkers , imaging and cognition. With its ability to provide a substantial body of data and a better understanding of different biomarkers , it has become a crucially important part of AD research. [28].

3.2 AIBL

The Australian Imaging, Biomarkers and Lifestyle (AIBL) study is an initiative to collect samples from 1000 individuals over 60 to collect data for AD research. The volunteers completed cognitive assessments , provided blood samples and went through detailed health and lifestyle questionnaires. One quarter of the participants received PET imaging while 10% received ActiGraph activity monitoring and body composition scanning.

In total the study collected data from 211 patients with AD, 133 participants with MCI, and 768 healthy controls. This dataset adds great value to ongoing AD research . The participants are reassessed at 18-month intervals to determine the predictive utility of available biomarkers. This longitudinal dataset is of immense value for tracking disease progression , evaluating predictive models and establishing clear labels for different disease stages. Additionally it helps to define the individual variability in both pathological and normal cognitive aging and to highlight the importance of such studies. One of the biggest points is that it also provides insights into the stage of MCI which is crucial as has been mentioned for early detection and targeted intervention. [29].

3.3 OASIS

The Open Access Series of Imaging studies (OASIS) is also a longitudinal imaging initiative similar to AIBL. It consists of MRI and PET images for 1098 participants collected over a period of 15 years. The age of participants ranges from 42 to 95 with 605 cognitively healthy individuals and 493 with cognitive decline .

The OASIS-3 Dataset consists of over 2000 MRI sessions featuring both structural and functional sequences, along with around 1500 raw PET scans. The dataset also includes additional products such as volumetric MRI segmentations and PET-derived metrics. Furthermore , OASiS provides data as APOE genotype information , and longitudinal clinical and cognitive assessments , making it a really valuable resource for research on

AD and Dementia. [30]

Chapter 4

Pre-Processing and Feature Extraction Techniques

4.1 Intensity Normalization

4.1.1 Introduction

Many processing steps rely on the foundational quality of the underlying image and the standardization of them across a dataset. Techniques such as registration, segmentation, and comparison or classification can be impacted significantly by variability in the intensity scale of the underlying image. Variability is inherent in MRI images , as scanner-related artifacts , different types of noise , acquisition protocols , and other parameters. [31]

To eliminate these issues and provide a standardized quality and scale across all images intensity normalization is employed as a critical pre-processing step. The main goal is to create a standardized scale that accurately represents the underlying tissue, by unifying the mean and variance within an image. The performance of this step directly influences the performance of processes responsible for tracking longitudinal disease progression. [31–33].

The goal is to standardize the scale so that the same tissue type will not be depicted in different values because of scanner-related differences and the same tissue should be depicted with the same intensity values across all images. [34, 35] . Overall this process is essential for downstream algorithms and reliability of further processing. [36, 37]

4.1.2 Necessity and Significance

The lack of tissue-specific intensity standardization for MR images creates variability, resulting in differences between two scans of the same subject, even when performed on the same scanner with the same pulse sequence [31]. As noted previously, this reduces the accuracy of subsequent processing algorithms. Key factors identified in the literature include:

Without a standardized intensity scale for specific tissue differences can occur even between two scans of the same subject , performed on the same scanner with the same pulse sequence [31] . As was also noted previously , this reduces the accuracy of any subsequent algorithm . Several contributing factors are identified in the literature :

- **Scanner Related Artifacts** Noise and scanner related artifacts interfere with subsequent processing tasks and must be addressed
- **Low Contrast and Partial Volume Effects** Different brain tissues should be ideally distinguishable - white matter (WM) , gray matter (GM) , cerebrospinal fluid (CSF) - . In practice , inside a voxel when two types of tissue exist boundaries are usually blurred [33] .
- **Intensity Variability** Differences in scanner hardware, acquisition protocols and data source can also create variations [38]. Those variations can occur either between different subjects or withing subsequent scans of the ame patient [39].
- **Lack of Tissue-Specific Meaning** Unlike CT , which uses standardized Hounsfield units , MRI intensities lack a physical meaning and so they depend on scanner specific settings. This limitation significantly affects longitudinal and quantitative analyses [37]

4.1.3 Goals and Principles

The normalization protocol follows the SPIN principles from Shinohara et. al [35]. Instead of using a reference histogram as [36] proposed, this method normalizes each scan based on the white matter intensities and characteristics of its own. The main goals are:

The primary goals are:

1. **Consistent measurements** across different scanners
2. **Reduced Variability** between subjects and scanning sessions
3. **Meaningful units** that relate to actual biology

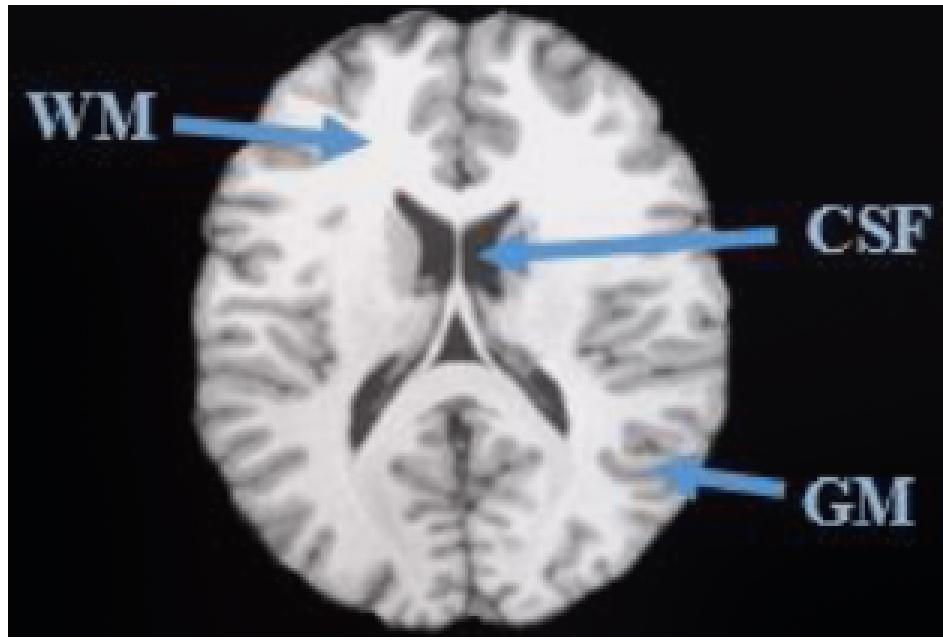


FIGURE 4.1: Brain image depicting artifacts present in MR images and different tissue types. Adapted from [31].

4. **Reliable Performance** even with lesions or imaging artifacts

4.1.4 Categories of Intensity Normalization Methods

Methods can be broadly categorized into:

1. Classical/Statistical Methods
2. Feature-Based Harmonization
3. Deep Learning-Based Approaches

Classical Methods

1. Min-Max Normalization This technique tries to squeeze all data values into a standard range, usually between 0 and 1. The smallest value becomes zero and the largest value becomes one, and everything else gets proportionately scaled. [40, 41]

$$v' = \frac{v - \min(A)}{\max(A) - \min(A)} \quad (4.1)$$

Where v' is the normalized value, v is the original data point, and $\min(A)$ and $\max(A)$ are the minimum and maximum feature values, respectively.

2. Z-Score Normalization This method adjusts each intensity value by subtracting from the voxel the average intensity value and dividing by the standard deviation (with-

out background). The result is a value that showcases how far it is from the average, measured in standard deviation units [38]. An important drawback of the technique can be that bright regions may be compressed towards the average , potentially losing information.

$$\mu = \frac{1}{|B|} \sum_{b \in B} I(b) \quad (4.2)$$

$$\sigma = \sqrt{\frac{\sum_{b \in B} (I(b) - \mu)^2}{|B| - 1}} \quad (4.3)$$

$$I_{\text{norm}}(x) = \frac{I(x) - \mu}{\sigma} \quad (4.4)$$

3. Mean and Standard Deviation Setting ($\mu \pm 3\sigma$ / M-std) This technique adjusts images in a certain way in order to fix their statistical value, typically set the mean to μ and the standard deviation σ to 1 [42, 43]. Some implementations confine intensity values within 3 standard deviations of the mean $\mu \pm \sigma$, This can work really well for pictures that follow a normal distribution $N(\mu, \sigma)$, but images that do not follow one may lose information due to value clipping [39].

4. Percentile Method This technique uses specific percentiles (usually the 5th and 95th percentiles) to define the range of the intensity values , which removes all the edge outliers [43].

5. Histogram Matching One of the most widely-used approaches , which works by trying to match the image to a reference histogram [36]

- **Nyul and Udupa’s Method:** At first the statistical landmarks of the training images are learned to create a standard distribution. Then , new images are transformed to match the standard using piecewise linear mapping.
- **Joint Histogram Matching:** Take pixel pairs and considers the relation of the input to the reference image. While it can potentially be more accurate it requires more computation citesunHistogrambasedNormalizationTechnique2015.

6. Homomorphic Filtering This technique tries to address uneven lighting that exists in images by separating the tissue from the illumination. The intensity of an image $I(x, y)$ can be represented as

$$I(x, y) = L(x, y) \cdot R(x, y) \quad (4.5)$$

where L represents illumination (lighting variations) and R represents the tissue reflectance. Taking the logarithm turns multiplication into addition:

$$\ln(I) = \ln(L) + \ln(R) \quad (4.6)$$

Turning the image into the frequency domain help remove slow-varying illumination (low frequencies) and keep tissue boundaries and details (high frequencies) [31]

7. Gaussian Filtering Used for smoothing, this requires careful selection of the σ parameter. A small σ may be inefficient for normalization, while a large σ can cause loss of edge information.

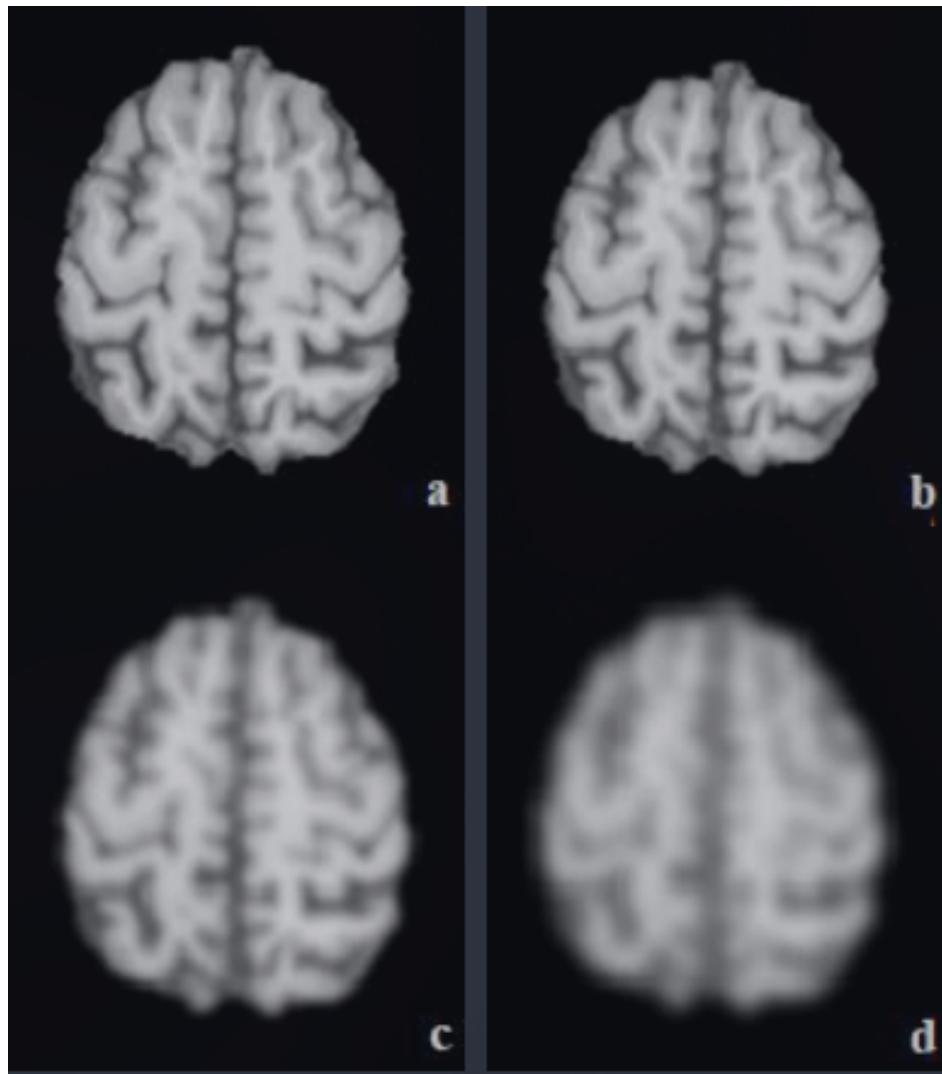


FIGURE 4.2: Original image (a); Image obtained by Gaussian filtering with increasing σ values (b-d).

8. White Stripe Normalization Proposed by Shinohara et al. [35], this method normalizes based on Normal-Appearing White Matter (NAWM). It identifies the NAWM

distribution ("white stripe") and matches its moments (mean and standard deviation) across subjects, providing biologically interpretable units robust to pathology.

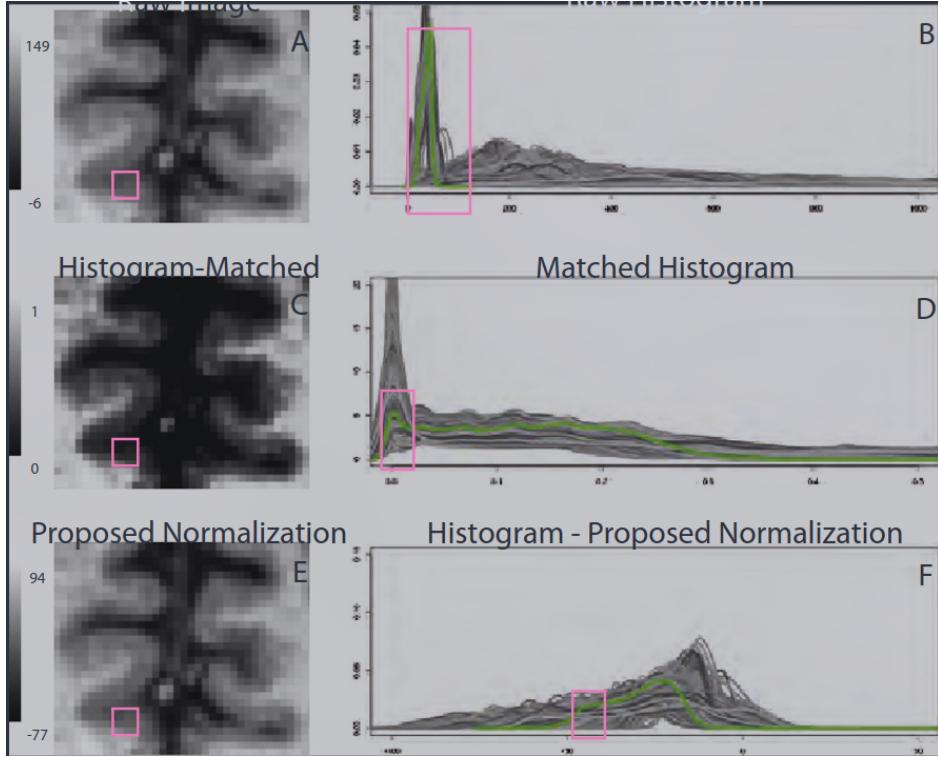


FIGURE 4.3: White Stripe Normalization concept.

9. Fuzzy C-Means (FCM) and Kernel Density Estimation (KDE) Fuzzy clustering assigns data points a membership function (0 to 1) indicating proximity to a cluster center. Normalization uses tissue segmentation to scale intensity relative to tissue means.

Single-Tissue Normalization: Normalize voxel x by the mean intensity of tissue T :

$$I_{\text{norm}}(x) = \frac{I(x)}{\mu_T} \quad (4.7)$$

Two-Tissue (Min-Max) Normalization: Using two tissue means μ_1 and μ_2 , define $a = \min(\mu_1, \mu_2)$ and $b = \max(\mu_1, \mu_2)$. The image is rescaled via:

$$I_{\text{norm}}(x) = \frac{I(x) - a}{b - a} \quad (4.8)$$

Kernel Density Estimation (KDE): KDE approximates the probability distribution function to smooth the histogram and locate peaks (e.g., white matter peak π) for

normalization:

$$p(x) = \frac{1}{NMLh} \sum_{i=1}^{NML} K\left(\frac{x - x_i}{h}\right) \quad (4.9)$$

The image is then normalized as $I_{\text{norm}}(x) = c \cdot \frac{I(x)}{\pi}$.

Feature-Based Harmonization

ComBat Originally for gene expression, ComBat is a feature-based method applied to imaging to eliminate batch effects [44]. It assumes feature similarity and allows only for shift or spread in the data.

Deep Learning Methods

Deep learning has emerged as a powerful tool for normalization, particularly in multi-site analysis.

1. Autoencoders (AEs) and Generative Adversarial Networks (GANs) An autoencoder is an algorithm designed to learn an informative representation of data by reconstructing input observations [45]. It encodes input x_i into a latent representation $h_i = g(x_i)$ and reconstructs it via $X_i = f(h_i)$. Training minimizes the reconstruction error:

$$\arg \min_{f,g} \Delta(X_i, f(g(x_i))) \quad (4.10)$$

GANs consist of two networks: a Generator that creates data and a Discriminator that evaluates it. In the context of normalization, authors have proposed using an autoencoder as the generator within a GAN framework. By utilizing adversarial training, the model learns to produce images that keep the structural integrity of brain anatomy while removing scanner related artifacts which the discriminator network identifies and flags , which improves its ability to work across different imaging sites [46, 47].

2. Adversarial and Task-Driven Normalization This approach uses two networks , one is a normalization network and the second a segmentation network. By combining the adversarial loss function with the loss function of the segmentation network it allows for a pipeline that can outperform generic domain adaptation as the normalization is made in order to maximize the segmentation of the tissue. [46]

3. Disentanglement Networks Disentanglement networks split the image content code (shapes , layout, geometry) with the style code (colors , textures) , where the first one is domain independent whereas the second one domain specific. This helps evade the assumption that unsupervised image-to-image translation has that there is a

one to one mapping. This allows frameworks like MUNIT to normalize content across different source domains. [47]

4.1.5 Impact of Intensity Normalization

Image Registration Intensity normalization can improve the performance of algorithmic registration because its goal is to ensure that similar tissue will have the same intensity values. This in turn increases the reliability of similarity metrics and leads to better image alignment [32, 34].

Image Segmentation One other subsequent task that benefits with intensity normalization is image segmentation. Since variations in tissue can negatively affect segmentation and volume estimation, normalization greatly enhances performance to reduce the variations within an image, even more so at multi-site datasets [33].

Texture Classification/Radiomics In texture analysis spatial gray-variations , are highly sensitive to acquisition protocols and the technique relies on them. Intensity normalization therefore is essential to ensure meaningful features. However they should be applied carefully , since their effects may depend on the imaging sequence (e.g. T1, T2 , FLAIR) [39].

Brain Template Construction Normalization can enable the calculation of different brain volume statistics for GM , WM and CSF. Post-Normalization brain templates show clearer tissue separation , making the boundaries between different types of tissue more distinct [32].

4.1.6 Challenges and Considerations

Even though its use can be critical for many subsequent tasks there are still challenges in the use of intensity normalization as a pre-processing step. FIrst of all the choice of method is mostly sequence-dependent [38]. Additionally deep learning frameworks outperform classical methods only in some specific cases , probably because of data scarcity , since deep learning even though very promisiisng is upper limited by the existence of diverse and large quantity datasets [43] . Overall the selection of an algorithm and an appropriate strategy remains one of the pivotal steps for an MR pipeline.

4.2 Denoising

Denoising is an important step in the pre-processing of images with the goal of using them as training data for a learning algorithm. Noise and artifacts can significantly impact the performance of such algorithms, which is why denoising is often included as a standard pre-processing step—especially when the quality of the raw images varies or when data come from multiple scanners (inherent variability).

In denoising, we assume there exists an “ideal” clean image represented as a vector $x \in \mathbb{R}^N$. However, what we observe is a noisy measurement y , which contains noise.

The most common mathematical noise model, due to the Central Limit Theorem and ease of use, is Additive White Gaussian Noise (AWGN). The relationship is:

$$y = x + v.$$

In this model, the noise is assumed to be drawn from a Gaussian distribution with mean zero and variance σ^2 :

$$v \sim \mathcal{N}(0, \sigma^2 I).$$

This means that each pixel is corrupted by a random value sampled from that distribution.

The goal is to find a function D such that $\hat{x} = D(y, \sigma)$, producing an estimate as close as possible to the original x .

Since we know the noisy image and the noise distribution, we can formulate the problem using Bayesian inference. We aim to maximize the posterior probability:

$$p(x | y) = \frac{p(y | x) p(x)}{p(y)}.$$

Taking the negative logarithm turns the product into a sum, resulting in the minimization objective:

$$-\log p(x | y) \propto \underbrace{-\log p(y | x)}_{\text{Data Fidelity}} + \underbrace{-\log p(x)}_{\text{Prior}}.$$

Because of the Gaussian noise model, the likelihood is:

$$p(y | x) \propto \exp\left(-\frac{\|y - x\|^2}{2\sigma^2}\right).$$

Taking the negative log removes the exponential:

$$-\log p(y | x) = \frac{1}{2\sigma^2} \|y - x\|^2.$$

The second term, $-\log p(x)$, serves as a prior and becomes a regularization term. The denoising objective becomes:

$$\hat{x} = \arg \min_x \underbrace{\|y - x\|^2}_{\text{Match the Data}} + \lambda \cdot \underbrace{\rho(x)}_{\text{Be a Real Image}}.$$

This formulation motivated the field's "evolution of priors," producing many classical denoising approaches.

Two of the most important classical algorithms are NLM and BM3D.

4.2.1 Non-Local Means (NLM)

NLM [48] grew from the idea of *spatial smoothness*. Earlier algorithms assumed images should be spatially homogeneous, so averaging pixels in local neighborhoods could remove noise while preserving signal. These methods often produced blurry results, especially near edges.

NLM instead uses neighborhoods (patches) of pixels and the concept of *self-similarity* across the entire image. For each pixel p , a patch P surrounding it is compared to patches Q around every other pixel q , using a similarity measure such as Euclidean distance $\|P - Q\|^2$.

The denoised estimate averages all pixels, but with weights proportional to patch similarity:

$$\hat{x}(p) = \frac{\sum_q w(p, q) y(q)}{\sum_q w(p, q)}.$$

4.2.2 BM3D

After NLM, the creation of algorithms such as BM3D improved denoising performance so dramatically that researchers questioned whether denoising had reached a theoretical limit.

BM3D stands for *Block-Matching and 3D Filtering*. Like NLM, it finds similar patches, but instead of averaging them directly, it stacks them into a 3D group. A transform (3D wavelet or DCT) is applied so that signal becomes concentrated in a few large co-

efficients while noise spreads into many small ones. Thresholding the small coefficients and inverting the transform yields high-quality results. [49]

4.2.3 Deep Learning Era

Before deep learning, humans designed the transforms (e.g., DCT) used to enforce sparsity. With deep learning, the network learns the appropriate transform directly from data.

We define a parametric function f_θ (a neural network) with parameters θ . The network takes a noisy image y and outputs a clean estimate x .

For supervised learning with an MSE loss:

$$\min_{\theta} \sum_{i=1}^N \|f_\theta(y_i) - x_i\|^2.$$

Since clean images are not always available, a common strategy is to corrupt clean images with Gaussian noise and train the model to recover them.

No single architecture is universally best; much progress arises from empirical experimentation to identify the best-performing denoiser.

A recently expanding research direction, highlighted in [50], uses denoising network priors as generative models capable of synthesizing new images. Using iterative methods such as Langevin dynamics, these priors can generate realistic images from noise. While not the primary focus of this work, such techniques can be beneficial in domains like medical imaging, where data scarcity is a major challenge.

4.3 Skull Stripping

In contrast with PET scans , MRI images capture non-brain anatomy - such as the eyes , skull and neck - that can create complications in automated analysis software. As a result, these extra tissues have to be digitally removed to ensure accurate brain measurements, making a critical preprocessing step in AD research. [51, 52].

The MRI system produces brain images as 3D volumetric data composed of 2D slices. Further computer-aided processing is essential in order to extract meaningful information, whether for research, diagnostic, or clinical purposes.

As noted earlier, quantitative morphometric studies require isolating brain tissue from non-brain tissue, a process referred to as *skull stripping*. Automated skull-stripping

enhances segmentation accuracy, making manual or automated segmentation methods more reliable [53]. Examples from [51] are shown in Figure 4.4.

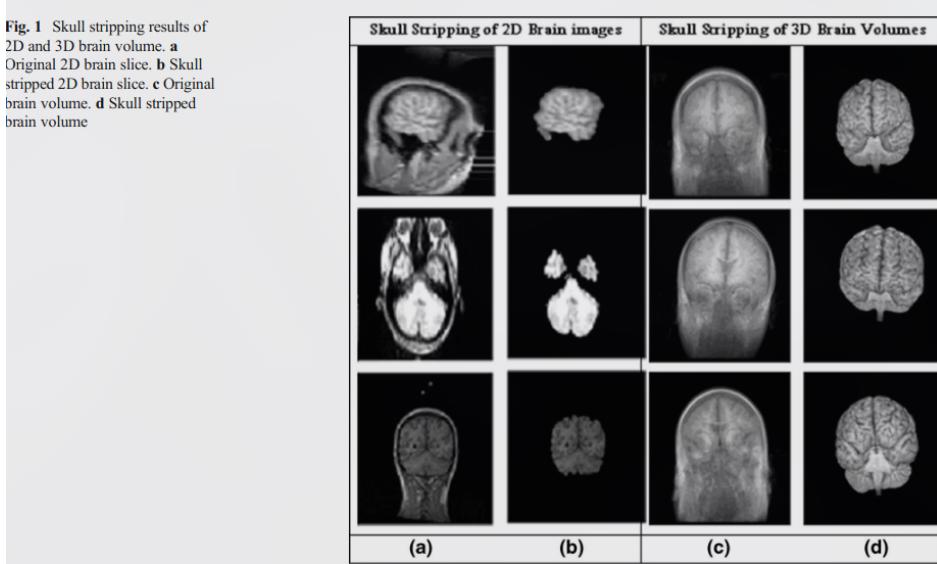


FIGURE 4.4: Examples of automated skull stripping from [51].

Voxel-based morphometry (VBM) results also showed significant improvements after skull stripping was applied [53]. Many brain imaging applications benefit from the precise segmentation of the brain from surrounding tissues, including cortical surface reconstruction, tissue classification, and registration to standard templates [54, 55].

Skull-stripping algorithms are typically evaluated by their speed and their impact on downstream automated tasks such as segmentation [56].

Multiple skull-stripping techniques have been developed due to their success and effectiveness in improving diagnostic and prognostic accuracy [51, 52].

Manual brain segmentation is considered the most robust and accurate approach, but it is extremely laborious, time-consuming, and subject to inter-clinician variability [57]. Manual masks are often treated as the ground truth against which automated methods are validated [53]. This has led to a strong need for automated skull stripping.

Despite significant progress, many skull-stripping approaches still perform well only on specific datasets and often require tuning of hyperparameters [56, 58]. According to [53], skull-stripping methods fall into three primary categories:

1. Manual skull stripping
2. Classical approaches
3. Deep learning–based approaches

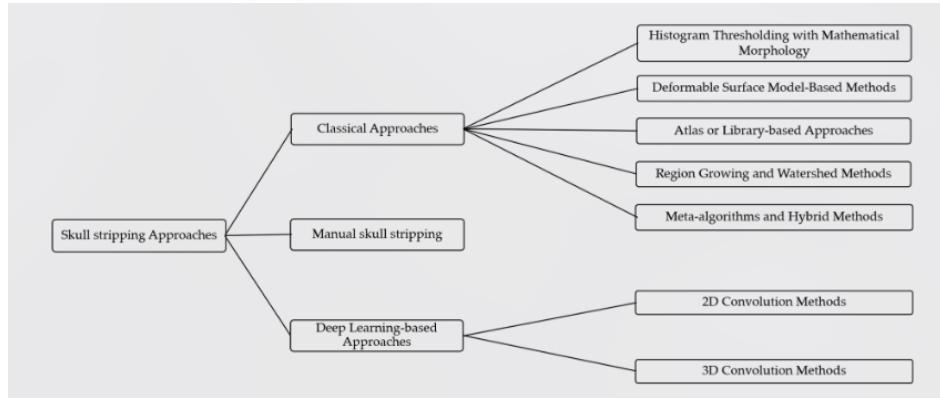


FIGURE 4.5: Overview of skull-stripping categories [53].

4.3.1 Skull Stripping Methods

Based on [51], skull-stripping techniques can be grouped as follows:

- Mathematical morphology-based methods
- Intensity-based methods
- Deformable surface-based methods
- Atlas-based methods
- Hybrid methods
- Deep learning-based methods

4.3.2 Morphology-Based Methods

These methods use morphological erosion and dilation operations—combined with thresholding and edge detection—to estimate the initial region of interest (ROI) and separate brain from non-brain tissues [51, 59].

A key drawback is that performance depends heavily on empirically tuned parameters related to the shape and size of the morphological structuring elements [52].

Examples

1. Brain Surface Extraction (BSE).

Uses anisotropic diffusion filtering, Marr–Hildreth edge detection, and morphological operations such as erosion and dilation. BSE is fast (≈ 3.5 seconds) [51, 59].

2. Brain Extraction Algorithm.

Uses diffusion, morphological operations, and connected component analysis for T1W and T2W MRI [53].

3. SMHASS (2012).

Combines deformable models with histogram analysis and morphological preprocessing [51].

4.3.3 Intensity-Based Methods

These methods classify brain vs. non-brain regions using pixel intensities. Examples include histogram-based, region-growing, and edge-based methods. Their primary limitation is sensitivity to MRI noise, low contrast, and bias field artifacts [51, 56].

Examples

1. Graph Cuts (GCUT, 2010).

Uses graph-theoretic segmentation to isolate brain tissue from dura, beginning with thresholding followed by connected submask selection [51].

2. Region Growing (RG).

Forms connected regions based on local intensity similarity. Variants include 2D approaches for coronal slices and MARGA for axial views and low-quality images [51].

4.3.4 Deformable Surface-Based Methods

These methods employ active contours (snakes) that evolve iteratively according to an energy functional. The contour adapts—shrinking or expanding—to match the brain boundary. Level-set formulations provide a robust mathematical framework [51, 54].

Their performance depends strongly on the initial contour and image quality, especially the clarity of edges [60].

Examples

1. BET (Brain Extraction Tool, 2002).

Widely used, fast, and freely available. It expands a tessellated sphere to match brain edges using adaptive forces. BET supports T1-weighted, T2-weighted, and proton density images, processing in 5–20 seconds [60, 61].

2. Model-Based Level Set (MLS, 2006).

Uses curvature and intensity-derived forces to evolve an active contour [51].

3. 3dSkullStrip (AFNI, 2005).

BET-like but includes improvements to avoid eyes and ventricles [51, 62].

4.3.5 Atlas / Template-Based Methods

These methods register the subject MRI to one or more anatomical atlases, enabling the transfer of brain masks to the subject [51, 55].

Examples

1. MAPS (2011).

Combines multiple atlas registrations to produce a consensus segmentation [51].

2. BEaST (2012).

Uses nonlocal segmentation with multi-resolution priors. BEaST achieves high accuracy through patch-based label fusion from a library of pre-labelled atlases [57, 61].

3. Pincram (2015).

Employs iterative refinement to propagate labels from multiple atlases [51].

4.3.6 Hybrid Methods

Hybrid approaches combine multiple algorithms or features to improve robustness and accuracy [51, 52].

Examples

1. SPECTRE (2011).

Integrates elastic registration, tissue segmentation, and morphological watershed operations [51].

2. HWA (Hybrid Watershed Algorithm, 2004).

Combines watershed segmentation with deformable surface models. HWA showed the highest sensitivity among compared methods and appeared more robust to parameter changes [54, 61].

3. BEMA (2004).

Runs multiple extractors (BET, BSE, Watershed, etc.) in parallel and merges the results [53].

4. ROBEX (2011).

Combines a Random Forest classifier with a point distribution model to ensure anatomically plausible results. The key advantage is parameter-free operation with consistent performance across datasets—achieving Dice scores of 95.6–97.0% on IBSR, LPBA40, and OASIS without tuning [58, 61].

4.3.7 Deep Learning-Based Methods

Deep learning approaches include 2D and 3D CNNs. While 3D CNNs capture volumetric context, they require higher computational resources. DL methods are generally categorized as [53, 63]:

- Patch-based CNNs
- Encoder–decoder CNNs (e.g., U-Net)

Encoder–decoder architectures typically perform better, are faster, and can capture global and local features [61].

However despite their performance, deep learning methods showcase several limitations [52, 53]:

- The requirement of large annotated datasets
- The inability of hyperparameter tuning due to the black-box behavior of the networks
- The sensitivity that emerges when trained on healthy individuals

State-of-the-Art Deep Learning Methods

SynthStrip represents a paradigm shift through synthetic training data. The method trains a 3D U-Net entirely on synthetically generated images, randomising intensity distributions, artifacts, and deformations. This yields contrast-agnostic generalisation—a single model processes T1w, T2w, FLAIR, DWI, MRA, CT, and PET images with Dice scores of 96–98%. Processing takes under 2 seconds on GPU [64].

HD-BET (High Definition Brain Extraction Tool) was designed for clinical heterogeneity. Training used 6,586 MRI sequences from 372 patients across 37 European institutions, including brain tumours and resection cavities. The ensemble of five 3D U-Net models with test-time augmentation outperforms BET, 3dSkullStrip, BSE, ROBEX, BEaST, and MONSTR [61].

deepbet achieves the highest reported accuracy for T1-weighted images using a two-stage 3D LinkNet architecture. Trained on 7,837 images from 191 OpenNeuro datasets, deepbet achieves median Dice of 99.0% with processing times of only 0.5 seconds on GPU [65].

Deep MRI Brain Extraction by Kleesiek et al. pioneered CNN-based skull stripping, demonstrating that deep learning could handle pathological brains better than classical methods [63].

4.3.8 Comparative Analysis

Automated vs. Manual Approaches

- Automated methods are faster but may require parameter tuning [56].
- Semi-automated methods are accurate but slow and user-dependent [51].

Evaluation Metrics

Common metrics include [53, 61]:

- Dice coefficient
- Jaccard Index
- Sensitivity / Specificity
- False Positive Rate / False Negative Rate
- Hausdorff Distance
- Average Symmetric Surface Distance (ASSD)

Comparative findings from the literature include [51, 53, 61, 66]:

- McStrip (hybrid) outperforms BET and BSE.
- HWA shows high sensitivity and robustness.
- Deep learning achieves highest Dice and specificity; 3D U-Net has highest sensitivity.
- SynthStrip performs best on pediatric T2-weighted scans, with accuracy increasing with age.
- HD-BET outperformed all classical methods by +1.16 to +2.50 Dice points on the CC-359 multi-vendor dataset.

4.3.9 Multi-Site Dataset Considerations

Large-scale studies aggregate data from numerous imaging sites with different scanners, protocols, and field strengths. Souza et al. created the CC-359 dataset specifically to evaluate vendor and field-strength effects, comparing BET, 3dSkullStrip, FreeSurfer, BSE, ROBEX, BEaST, and OptiBET across 359 acquisitions from GE, Philips, and Siemens scanners at 1.5T and 3.0T. Results revealed statistically significant effects ($p < 0.001$) for both vendor and field strength [56].

4.3.10 Implications for Alzheimer’s Disease Research

The application of skull stripping in neurodegeneration research introduces domain-specific challenges. Brain atrophy—the hallmark of Alzheimer’s disease—alters the anatomical relationships that skull stripping algorithms exploit [55, 67].

Recent research by Tinauer et al. analysed 990 matched ADNI images and discovered that skull stripping introduces shortcut learning. CNNs trained on skull-stripped images learned brain contours introduced through preprocessing rather than clinically relevant atrophy markers—a “Clever Hans effect” inflating apparent classification accuracy [67].

Novosad and Collins evaluated skull stripping on ADNI subjects and found that brain masks include more subarachnoid CSF in atrophied brains, failure to remove non-brain tissue causes over-estimation of cortical thickness, and poor skull stripping propagates errors to regional volume and atrophy estimates [55].

It is also important to note that variability of anatomy, age, and extent of brain atrophy impacts skull stripping for volumetric Alzheimer’s analysis [58].

4.3.11 Strengths and Weaknesses of Skull Stripping Approaches

Methodology	Strengths	Weaknesses
Mathematical Morphology	Simple to implement; fast [59]	Parameter-dependent; noise-sensitive; risk of over/under-segmentation [51]
Intensity-Based	Uses fundamental image properties	Sensitive to noise, bias field, threshold choice; watershed over-segmentation [51]
Deformable Surface	Accurate and robust; boundary-aware [60]	Sensitive to initialization; noise-sensitive; may fail in non-standard cases [56]
Atlas-Based	Leverages anatomical priors; good when intensities are unreliable [57]	Dependent on registration quality; computationally intensive [55]
Hybrid	Combines strengths of multiple methods; often fully automatic [58]	Complex; inherits weaknesses of contributing methods [51]
Deep Learning	State-of-the-art performance; learns features automatically [61, 64]	Requires large datasets; expensive; black-box behavior [53]

TABLE 4.1: Comparative strengths and weaknesses of skull-stripping methodologies.

4.3.12 Primary Challenges

Major difficulties arise from [51, 52, 55]:

- MRI artifacts (noise, intensity bias, motion)
- Anatomical complexity and overlapping tissue intensities
- Presence of pathology (e.g., tumors) affecting segmentation accuracy
- Brain atrophy in neurodegeneration creating ambiguous tissue boundaries

4.4 Voxel-Based Morphometry (VBM)

VBM is a whole-brain statistical technique that replaced the manual procedure in which clinicians had to track a specific brain structure, such as the hippocampus, across sequential slices. This technique frees the user from being confined to a single region of interest and instead enables analysis of the entire brain, detecting even subtle structural changes.

The VBM pipeline consists of several key steps:

1. Image Acquisition

This typically involves acquiring a high-resolution T1-weighted MRI scan.

2. Spatial Alignment (Normalization)

All images must be aligned to a common template to allow voxel-wise statistical comparison.

3. Tissue Segmentation

Tissues are automatically classified into gray matter (GM), white matter (WM), and cerebrospinal fluid (CSF). This relies on spatial priors and voxel intensities.

Formally:

$$P(\text{Tissue Class} \mid \text{Intensity, Location})$$

4. Modulation

Modulation preserves the total amount of tissue after spatial transformations using the Jacobian determinant. The choice to use modulation depends on the research question:

- Without modulation: analysis reflects *concentration* (e.g., GM density per unit volume).
- With modulation: analysis reflects *absolute volume*.

Modulation may miss microscopic changes, and it is suggested that modulation should not be used to detect mesoscopic abnormalities such as cortical thinning [68].

5. Smoothing

An isotropic Gaussian kernel is applied to:

- (a) Increase signal-to-noise ratio by averaging neighboring voxels.
- (b) Compensate for minor spatial misalignments after normalization.
- (c) Render the data more normally distributed for parametric statistics.

6. Statistical Analysis

Voxel-wise statistical tests are typically performed using the General Linear Model (GLM). A t-value map (Statistical Parametric Map) is created by comparing group means and accounting for variance.

To control false positives:

- **Family-Wise Error (FWE)** correction
- **False Discovery Rate (FDR)**

4.4.1 Evolution of VBM

VBM initially faced criticism regarding registration and segmentation accuracy. Two major improvements emerged:

- **Optimized VBM**

Introduced study-specific templates and better skull-stripping to improve normalization accuracy.

- **DARTEL**

A diffeomorphic algorithm using exponentiated Lie algebra for high-dimensional registration.

The algorithm iteratively computes an increasingly accurate group average and stores subject-specific deformations as *individual flow fields*.

4.4.2 Alzheimer's Disease and Atrophy Patterns

1. Medial Temporal Lobe

The earliest and most critical changes occur here, involving structures essential for episodic memory.

2. Temporoparietal Association Cortex

Includes the temporal lobe, parietal lobe, and angular gyrus—responsible for memory, spatial cognition, and language processing.

3. Posterior Cingulate Cortex and Precuneus

Key nodes of the Default Mode Network (DMN), implicated in self-referential thought and autobiographical memory.

4. Frontal Lobe Involvement

VBM helped demonstrate that the frontal lobe is also affected, which has significant clinical implications for symptoms such as apathy and disinhibition.

5. Unaffected Regions

The primary visual cortex, primary sensorimotor cortex, and the cerebellum remain largely intact.

6. VBM and Braak Staging

Braak staging describes the sequential spread of neurofibrillary tangles (NFTs). VBM-detected atrophy patterns mirror those found in post-mortem studies, suggesting VBM can track macroscopic consequences of NFT accumulation.

7. Structural Integrity vs. Metabolism

Hypometabolism may precede structural loss, meaning that VBM cannot always capture early dysfunction [69, 70].

4.4.3 VBM as a Prognostic Tool

Mild Cognitive Impairment (MCI) is heterogeneous. VBM can distinguish:

- **MCI converters:** likely to progress to AD
- **MCI non-converters:** may remain stable or revert

[73]

4.4.4 Differences Between Converters

MCI converters show atrophy patterns similar to mild AD, involving:

- Entorhinal cortex
- Hippocampus
- Temporal lobe

Brain Region	Functional Role	Atrophy Pattern in AD	Supporting Imaging Reference	Figure Link
Medial Temporal Lobe (hippocampus, entorhinal cortex)	Episodic memory encoding & consolidation	Earliest and most severe atrophy on MRI/VBM; hallmark of typical AD progression.	VBM & MRI studies document significant grey matter loss here.	Fig. 4.6
Temporoparietal Cortex (posterior temporal & parietal lobes, angular gyrus)	Memory integration, spatial cognition, language	Substantial atrophy in typical AD; also key to posterior cortical involvement.	Typical AD shows loss in these regions relative to controls.	Fig. 4.7
Posterior Cingulate Cortex & Precuneus	Default Mode Network (DMN), self-referential processing	Significant atrophy , often seen with hypometabolism even before volume loss.	Posterior midline atrophy linked to DMN disruption.	Fig. 4.8
Frontal Lobes	Executive function, behavior regulation	Later-stage involvement with more diffuse atrophy; less prominent early.	Atrophy spreads later toward frontal regions in AD.	Fig. 4.9
Primary Visual, Sensorimotor Cortex, Cerebellum	Sensory/motor processing	Relatively spared compared to associative and memory networks in early/mild AD.	Primary regions less affected in typical AD.	Fig. 4.10
Structural vs Metabolic Changes	-	Hypometabolism can occur before detectable grey matter loss.	FDG PET shows metabolic changes preceding volume loss.	Fig. 4.11

TABLE 4.2: Summary of regional atrophy patterns in AD and linked VBM-related figures.

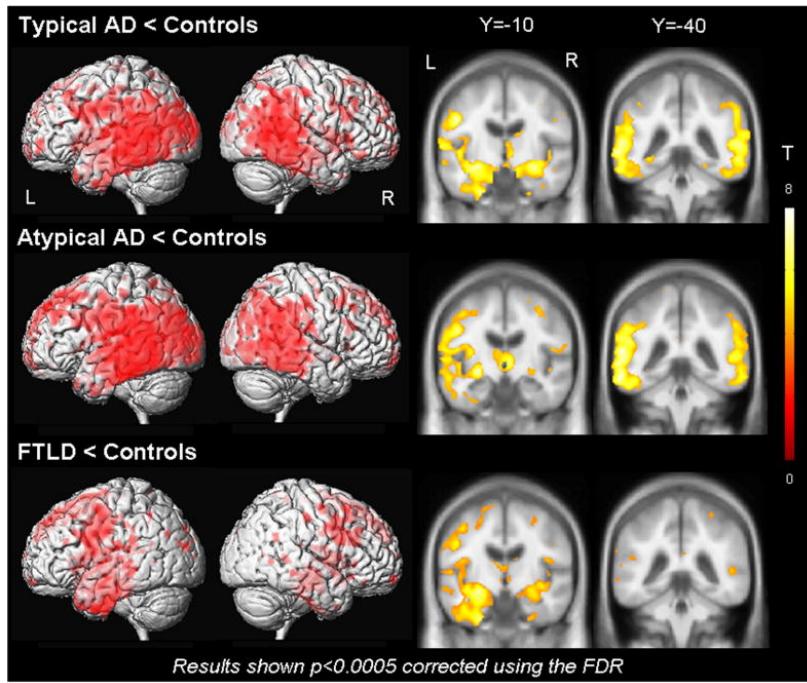


FIGURE 4.6: Direct Comparisons using VBM show distinct patterns of atrophy. While Typical AD and FTLD patients showed loss in the hippocampal area, atypical subjects showcased severe damage in the left putamen and tempoparietal cortex. [71]

- Parietal lobe
- Posterior cingulate cortex

Non-converters show more limited and localized GM loss [74].

4.4.5 Meta-Analytic Power of VBM

A major meta-analysis identified the **left hippocampus and parahippocampal gyrus** as the most consistent predictors of conversion to AD [75]. Another found robust atrophy in the **left amygdala** and **right hippocampus**.

Rate of Decline in AD

Longitudinal VBM shows that “fast decliners” have greater GM loss at baseline, demonstrating that initial atrophy predicts future trajectory [76].

Clinical Utility

Tools such as VSRAD provide clinicians with voxel-wise z-scores derived from large normative datasets, similar to how blood test references are used.

Accuracy

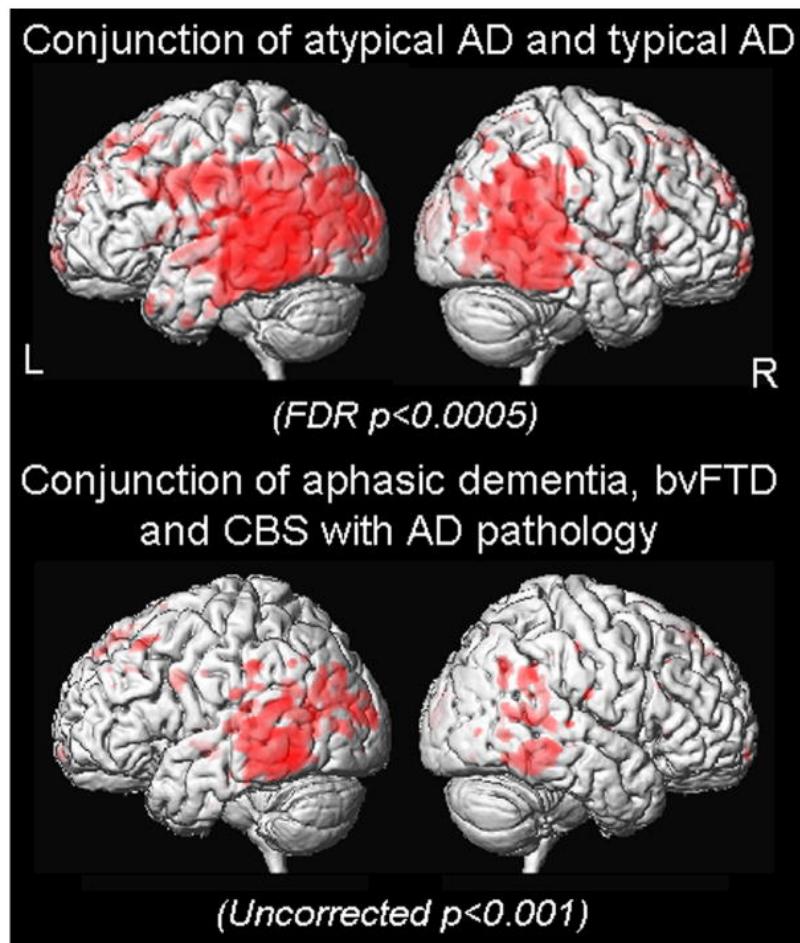


FIGURE 4.7: "The visualizations illustrate conjunction analyses , which show overlapping damage across different groups. The first one shows a comparison between AD typical and atypical damage to healthy controls. The second shows damaged that is shared between three specific variants: aphasic dementia, Corticobasal Syndrome (CBS), and behavioral variant Frontotemporal Dementia (bvFTD) " [71]

VBM often achieves an AUC exceeding 0.90 in distinguishing AD from controls [77], demonstrating high diagnostic accuracy.

4.4.6 Subtypes of AD

VBM helps distinguish:

- **Early-Onset AD (EOAD)** – more parietal and posterior cingulate atrophy
- **Late-Onset AD (LOAD)** – more medial temporal atrophy

[78]

Structure–Clinical Correlations

VBM correlates brain atrophy with clinical features:

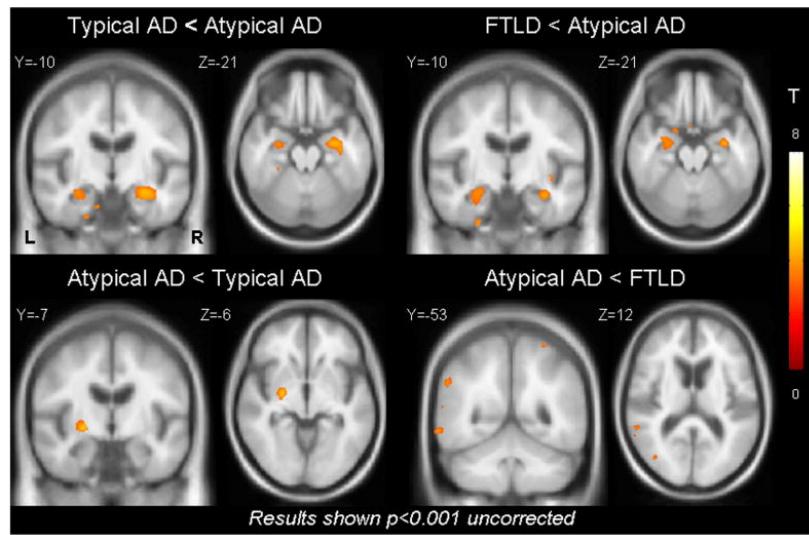


FIGURE 4.8: VBM shows distinct patterns of atrophy in between different groups. Typical AD and FTLD show more loss in the hippocampal area than atypical AD patients. Atypical patients show more severe damage in the left putamen (surpassing Typical AD) and in the tempoparietal cortex (surpassing atypical AD) . [71]

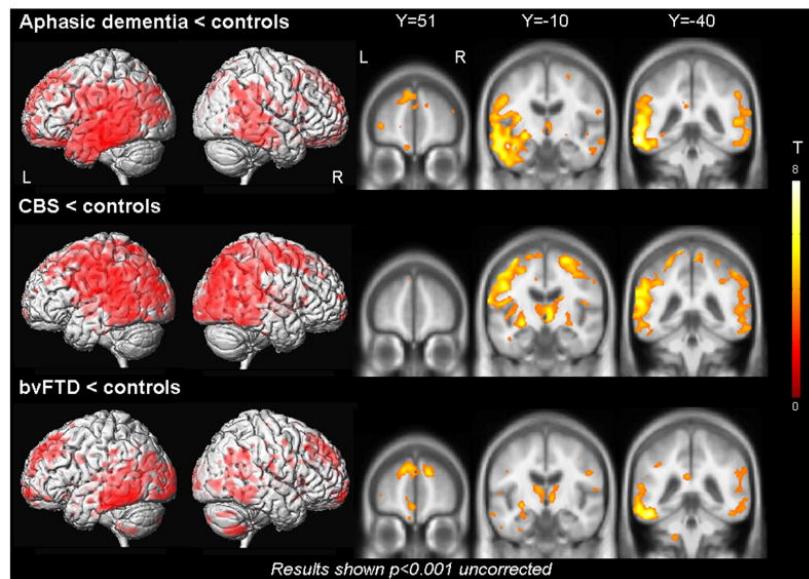


FIGURE 4.9: Shows the distinct atrophy for three clinical subgroups. Aphasic Dementia, Corticobasal Syndrome (CBS) and behavioral variant Frontotemporal Dementia (bvFTD) compared to healthy controls. [71]

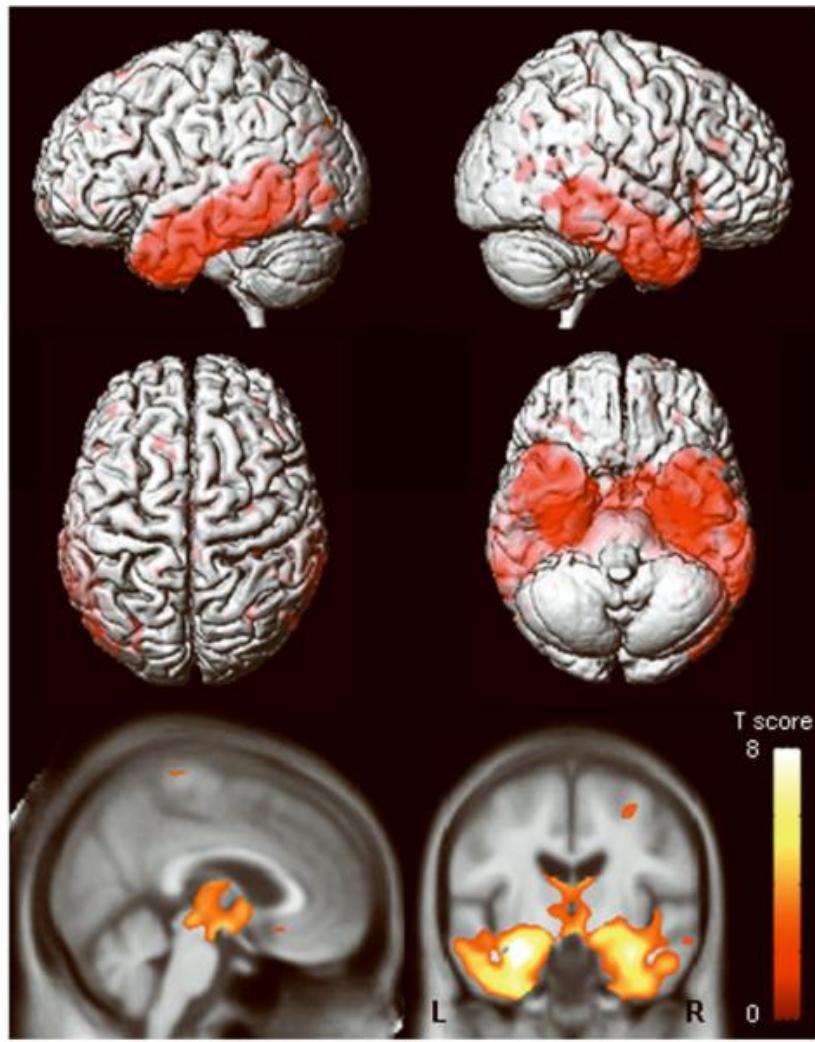


FIGURE 4.10: "The images map statistically significant gray matter atrophy of the a-MCI patients compared to healthy controls. The top one shows damage on the brain outer surface , whiel the bottom show deep tissue loss in the cingulate cortex and medial temporal lobes" [72]

- Frontal atrophy correlates with impaired daily living activities.
- Apathy correlates with atrophy in anterior cingulate and orbitofrontal cortex.
- Disinhibition correlates with orbitofrontal atrophy.

4.4.7 Limitations of VBM

1. Multiple Comparisons

Strict corrections (FWE) reduce false positives but may reduce sensitivity.

2. Preprocessing Sensitivity

Results depend heavily on choices in normalization, segmentation, and statistical thresholds.

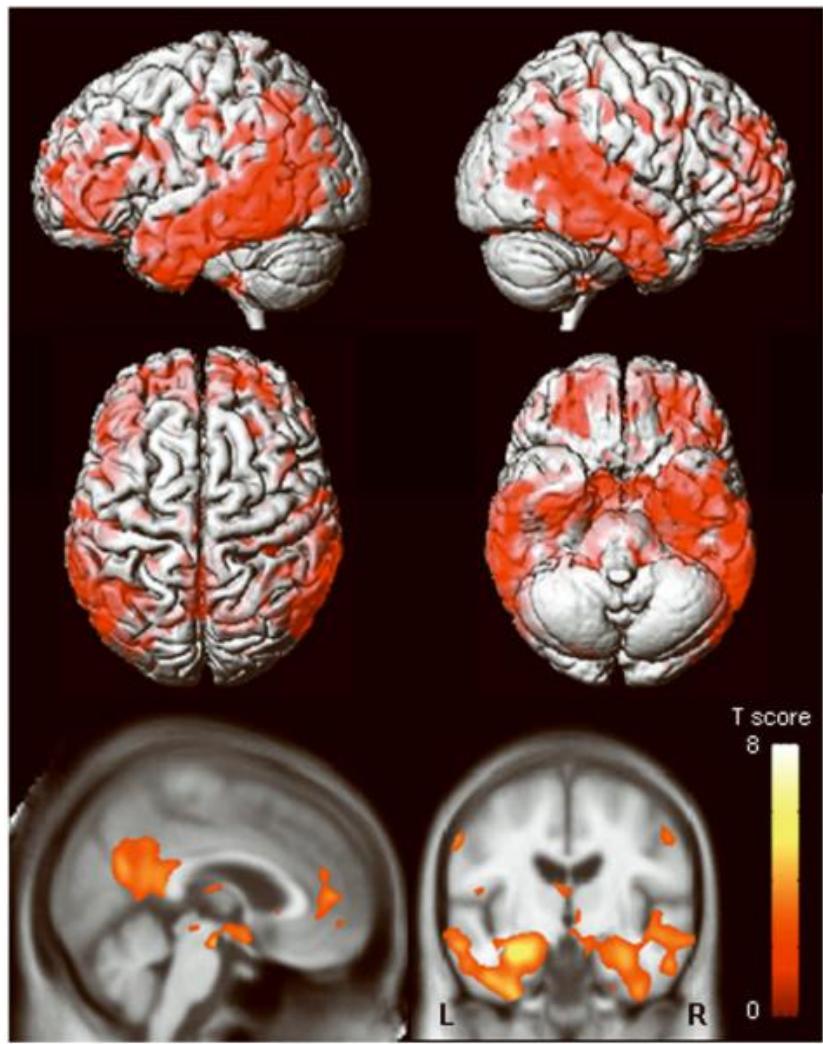


FIGURE 4.11: The figure shows specific brain areas in the aMCI-P group that showcase larger atrophy than those in the aMCI-S group. [72].

3. Registration Problems

Even DARTEL may misalign regions due to individual variability.

4. Segmentation Errors

Artifacts, low SNR, and partial volume effects can misclassify tissue.

5. Ambiguity

A decrease in voxel intensity could reflect thinning, reduced surface area, or folding changes.

4.4.8 Machine Learning and VBM

1. Supervised methods (SVM, Random Forests) classify AD vs. controls and predict cognitive decline.
2. Feature importance can reveal novel biomarkers.

4.4.9 Deep Learning

CNNs can classify AD and MCI without hand-crafted features, achieving high accuracy (e.g., 80.9% for MCI identification) [79].

Chapter 5

Classification Techniques for Dementia Image Analysis

5.1 Foundational Machine Learning Paradigms in Dementia Classification

Prior to the deep learning era, research in dementia and AD relied on traditional machine learning methods. These workflows typically followed a pipeline of preprocessing, manual feature extraction, selection, and classification

This pipeline formed the basis of early computational dementia research [80, 81].

5.1.1 Support Vector Machines (SVM)

Support Vector Machines are one of the most widely used algorithms and a foundational for classification. The algorithm works by solving an optimization problem to maximize the margin of a hyperplane between the projected data points, while separating the classes. For non linearly separable data it can employ the **kernel trick**, using polynomial or radial basis function to map data into a higher-dimensional space where they can be linearly separated.

In dementia research, SVMs were extensively applied for classifying Alzheimer's disease (AD) versus healthy controls (HC). Input features were typically derived from structural MRI, often by dividing the brain into regions of interest (ROIs) and computing measurements such as grey matter (GM) volume, cortical thickness, or voxel-based morphometry (VBM) values.

Several studies reported high accuracies. For example, one study using whole-brain

MRI-derived features achieved a mean classification accuracy of 94.5% (96.6% specificity, 91.5% sensitivity) for AD vs. controls [82]. Another reported an accuracy of 99.06% using 2D MRI slices [83].

However, performance degraded significantly when distinguishing more subtle clinical categories, such as healthy controls (HC), mild cognitive impairment (MCI), and in particular predicting MCI converters (MCI-C) versus non-converters (MCI-NC) [84].

The limitation lies not in the SVM algorithm itself, but in the nature of the disease and the constraints of manual feature engineering. Structural changes in late-stage AD—such as GM atrophy—are large and consistent across subjects, making the classes easily separable. MCI, however, is heterogeneous: not all MCI patients convert to AD, and subtle pathological differences are often not captured by handcrafted features such as ROI-based volumes.

This limitation highlighted the dependence of classical machine learning on high-quality engineered features, motivating the shift toward deep learning, where feature extraction is learned directly from the data.

5.1.2 Ensemble Methods: Decision Trees and Random Forests

Random Forests (RFs) emerged as a strong alternative to SVMs. An RF builds an ensemble of decision trees, each trained on a bootstrap sample of the data. At each node, splits consider only a subset of features, which reduces overfitting and is well-suited for high-dimensional, low-sample-size neuroimaging data [85].

The stability of RFs has been demonstrated in studies showing that even after substantial feature reduction, RF accuracy decreases less and remains more stable compared to multilayer perceptrons (MLPs) and convolutional neural networks (CNNs) [86].

A key advantage of RFs is their ability to compute feature importance via decreases in Gini impurity. This has repeatedly highlighted brain regions such as the hippocampus, amygdala, and inferior lateral ventricles as major contributors in NC/MCI/AD classification tasks [87].

RFs have achieved strong performance, including:

- 93.6% accuracy for predicting MCI-to-AD conversion from clinical data [88],
- 90.2% accuracy in three-class MRI-based classification (NC, MCI, AD) [87].

However, Gini-based importance scores may be biased toward continuous or high-cardinality features and can be misleading in the presence of correlated predictors

[89]. In neuroanatomy, many structures (e.g., left/right hippocampus) are strongly correlated, meaning RF importance rankings require careful interpretation.

5.1.3 The Curse of Dimensionality: PCA and LDA

A major challenge in neuroimaging machine learning is the *curse of dimensionality* [90]. When the number of features p greatly exceeds the number of subjects, models tend to overfit and generalize poorly. Dimensionality reduction methods help mitigate this issue.

Principal Component Analysis (PCA)

PCA is an unsupervised method that identifies orthonormal components capturing directions of maximal variance [90, 91]. Retaining only the top components reduces noise and redundancy [92]. PCA has been used to compress VBM maps and ROI volumes before feeding them into classifiers, often improving performance.

Linear Discriminant Analysis (LDA)

LDA is a supervised method that seeks projections maximizing the ratio of between-class to within-class variance [93]. Unlike PCA, which is label-agnostic, LDA explicitly optimizes separability among classes such as NC, MCI, and AD.

Studies consistently show LDA outperforming PCA for dementia classification [90]. However, as a supervised method, it may overfit if distributions shift between train and test sets.

5.2 The Deep Learning Revolution

Deep learning disrupted the traditional feature-engineering workflow by enabling end-to-end learning from raw or minimally processed images [94].

5.2.1 Convolutional Neural Networks (CNNs)

CNNs form the backbone of most imaging-based dementia classification systems [95]. Their strength lies in hierarchical feature extraction: early layers detect low-level patterns (edges, textures), while deeper layers capture structural patterns associated with neurodegeneration [95].

2D vs. 3D CNN Designs

Choosing between 2D and 3D architectures involves trade-offs:

2D CNNs. These operate on slice-based inputs, enabling computational efficiency and transfer learning from natural-image models such as VGG, ResNet, Inception, and DenseNet [96, 97]. However, slice-wise inputs discard 3D anatomical continuity [98].

3D CNNs. 3D CNNs preserve volumetric structure using 3D kernels [99]. They achieve strong performance on AD/NC and multi-class tasks but require large datasets and memory, making them susceptible to overfitting [100].

Hybrid approaches include:

- 2.5D slices (central slice plus neighbors) [101],
- CNN+RNN architectures for sequential slice modeling [95].

Multimodal Fusion

Structural MRI (T1) and FDG-PET offer complementary structural and metabolic information. CNNs can fuse these modalities using multi-stream architectures, where early layers learn modality-specific features before being merged for joint classification [102, 103].

Interpretability: Grad-CAM

Deep models face adoption challenges due to limited interpretability. Grad-CAM provides visual explanations by generating heatmaps indicating which regions contributed most to the model’s decision [104, 105].

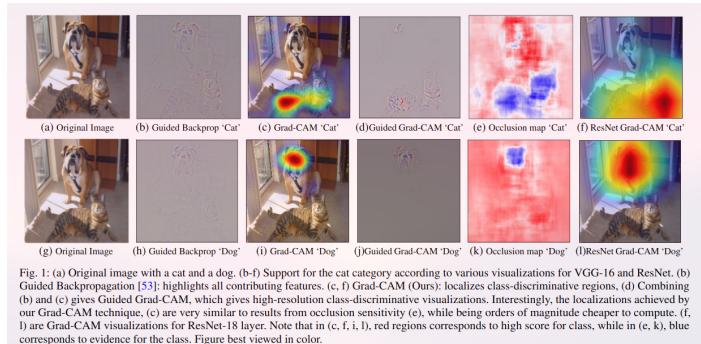


FIGURE 5.1: Example Grad-CAM heatmap highlighting salient regions [104].

Such maps enable researchers and clinicians to verify that models rely on plausible neuroanatomical biomarkers rather than noise or artifacts.

5.2.2 Generative Models: GANs

Deep learning is limited by the scarcity of large, annotated medical datasets [103, 106]. Generative Adversarial Networks (GANs) address this by producing synthetic

yet realistic MRI or PET images [107].

GAN-generated images can augment training datasets, alleviating class imbalance and improving classifier performance, with gains up to 11.68% reported [107]. GANs are also explored for:

- resolution enhancement (e.g., 1.5T → 3T MRI) [108],
- cross-modality synthesis (e.g., MRI-to-PET).

5.2.3 Vision Transformers (ViTs)

Transformers, originally developed for NLP, now challenge CNN dominance in medical imaging [109]. ViTs model long-range spatial dependencies using self-attention, enabling global context reasoning from the first layer [110]. This is appealing for dementia, where pathology is subtle and spatially distributed.

ViTs treat an image as a sequence of patches, which are embedded, combined with positional encodings, and processed through transformer layers [109]. Studies show ViTs achieving state-of-the-art AD classification performance [111].

5.3 Hybrid Classification Approaches

Hybrid methods combine CNNs for feature extraction with classical machine learning classifiers such as SVMs. These models leverage deep feature representations while benefiting from the robustness and theoretical properties of traditional classifiers. Such approaches have achieved up to 90% accuracy for AD vs. HC and up to 98% in four-class dementia classification tasks [112].

Chapter 6

Classification Performance and Accuracy Metrics

6.1 Evaluation Metrics

6.1.1 Accuracy, Sensitivity, Specificity

Evaluation metrics allow the visualization of the performance of a classifier and the comparison between different classifiers.

In Alzheimer's disease research, imbalanced data are extremely common. Imbalanced data occur when one class is significantly overrepresented or underrepresented.

To understand these metrics, we introduce the **confusion matrix**, which visualizes the relationship between predicted labels and actual labels:

TABLE 6.1: Confusion Matrix showing Actual vs. Predicted values

Total Population = P + N	Predicted Positive	Predicted Negative
Actual Positive	True Positive (TP)	False Negative (FN)
Actual Negative	False Positive (FP)	True Negative (TN)

Accuracy in machine learning is the proportion of correct classifications over the total number of classifications. However, accuracy can be misleading, especially when dealing with imbalanced datasets. Its weaknesses include lower informativeness, discriminability, distinctiveness, and a bias toward the majority class [113].

Accuracy is defined as:

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN}$$

Sensitivity is the ability of a classifier to detect actual positives. In medical applications, a sensitivity of 0.8 means that the method correctly identifies positive cases 80% of the time.

$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad (6.1)$$

Specificity measures the ability to correctly identify actual negatives:

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (6.2)$$

6.1.2 Precision and Recall

Precision (Positive Predictive Value) quantifies how many predicted positives are actually positive:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (6.3)$$

Recall is defined as the ratio of the actual positives (TP and FN).

$$\text{Recall} = \frac{TP}{TP + FN} \quad (6.4)$$

6.1.3 Model Validation Techniques

Model validation ensures that the model can generalize to unseen data. To improve robustness, methods such as **cross-validation** and **leave-one-out validation (LOOV)** are used.

- **Cross-validation:** Splits the dataset into k folds and trains the model k times, each time using one fold for testing and the rest for training.
- **Leave-One-Out Validation:** A special case where each data point is treated as a fold. Maximizes data usage but is computationally expensive.
- **Bootstrap Validation:** Repeatedly samples with replacement, trains on each sample, and tests on the remaining data. Effective for small or uncertain datasets [114].

6.1.4 ROC Curve

The Receiver Operating Characteristic (ROC) curve plots the True Positive Rate (TPR) against the False Positive Rate (FPR) at various threshold settings. The Area

Under the Curve (AUC) reflects classifier performance: the larger the AUC, the better the classifier.

AUC is robust to class imbalance [115, 116].

6.1.5 Matthews Correlation Coefficient

The Matthews Correlation Coefficient (MCC) is considered one of the most informative evaluation metrics [117]. Other related metrics include:

- Markedness
- Informedness (Bookmaker Informedness, BM)
- Balanced Accuracy (BA)

TABLE 6.2: Part 1: The Confusion Matrix and Row-based Rates (Sensitivity/Specificity)

		Predicted condition			
		Predicted Positive	Predicted Negative		
Actual condition	Total P	True Positive (TP)	False Negative (FN)	True Positive Rate (TPR)	False Negative Rate (FNR)
		<i>Hit</i> TP	<i>Miss, Type II Error</i> FN	Sensitivity, Recall $\frac{TP}{P} = 1 - FNR$	Miss Rate $\frac{FN}{P} = 1 - TPR$
	Total N	False Positive (FP)	True Negative (TN)	False Positive Rate (FPR)	True Negative Rate (TNR)
	<i>False Alarm, Type I Error</i> FP	<i>Correct Rejection</i> TN	Fall-out $\frac{FP}{N} = 1 - TNR$		Specificity $\frac{TN}{N} = 1 - FPR$

TABLE 6.3: Part 2: Derived Performance Metrics (Precision, Accuracy, etc.)

Prevalence	Precision (PPV)	False Omission Rate (FOR)	Positive Likelihood Ratio	Negative Likelihood Ratio
$\frac{P}{P+N}$	$\frac{TP}{TP+FP} = 1 - FDR$	$\frac{FN}{FN+TN} = 1 - NPV$	$LR+ = \frac{TPR}{FPR}$	$LR- = \frac{FNR}{TNR}$
Accuracy (ACC)	False Discovery Rate (FDR)	Neg. Predictive Value (NPV)	Markedness (MK)	Diagnostic Odds Ratio
$\frac{TP+TN}{P+N}$	$\frac{FP}{TP+FP} = 1 - PPV$	$\frac{TN}{TN+FN} = 1 - FOR$	$PPV + NPV - 1$	$DOR = \frac{LR+}{LR-}$

Chapter 7

Challenges in Dementia Image Analysis and Classification

7.1 Imbalanced Data

The performance of any machine learning model is fundamentally limited by the quality of its data [118]. A critical challenge that can degrade data quality is class imbalance, which occurs when one class contains many more sample than another - for example medical applications or fraud detection - [119]. This imbalance therefore represents a major challenge since standard ML algorithms are designed to maximize overall accuracy [120]. When trained on imbalanced data they tend to show bias towards the majority class while placing less significance to the minority class - often the class of most significance , since detecting the presence of disease is the primary interest- [121]. This performance degradation is documented , in decision tree classifiers [122].

When dealing with imbalanced datasets, standard evaluation metrics like accuracy can be misleading, as they are often biased toward the majority class. For a more robust and insightful assessment of classifier performance, it is crucial to employ metrics specifically designed to handle skewed class distributions. Based on [123], these can be divided into two primary categories: **Threshold Metrics** and **Ranking Methods**.

7.1.1 Threshold Metrics

Threshold metrics evaluate a classifier's performance at a fixed operational point. They provide a snapshot of performance but assume that the class distribution and error costs are known and constant.

Sensitivity and Specificity These metrics evaluate performance on each class independently.

- **Sensitivity (Recall, True Positive Rate):** The number of positive instances that are correctly identified as in Eq. (6.1)
- **Specificity:** Negative instances that are correctly identified as in Eq. (6.2)

Precision and Recall These two metrics are the quantification of precision and recall

- **Precision:** Precision is used to show the proportion of examples that are correctly identified Eq. (6.3)
- **Recall:** The ratio of true positives over the actual positives (Eq (6.4)).

Combined Metrics In order to simplify evaluation and provide a more comprehensive view some metrics combine all of the above.

- **G-Mean (Geometric Mean):** It balances the ability of the model both on positive and negative cases by taking the square root of the product between sensitivity and specificity.

$$\text{G-Mean} = \sqrt{\text{Sensitivity} \times \text{Specificity}} \quad (7.1)$$

- **F-Measure (F1-Score):** It computes the harmonic mean between recall and precision.

$$\text{F-Measure} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (7.2)$$

7.1.2 Ranking Methods and Metrics

In cases where class distribution varies such as imbalanced data or the error costs for making a wrong decision can vary , ranking methods are used to provide an overview of performance across all decision thresholds.

ROC Analysis Plots True Positive Rate vs. False Positive Rate across thresholds. A curve near the top-left corner indicates a strong classifier. Particularly effective for imbalanced data.

AUC Area Under the ROC Curve; represents the probability that a randomly chosen positive instance is ranked higher than a randomly chosen negative one.

Area Under the ROC Curve , is a metric for how well a classifier distinguished between classes. It actually quantifies the probability that a model assigns a higher probability to a randomly selected positive sample than a negative one .

Precision–Recall Curves Show precision–recall trade-offs across thresholds and are preferred for highly imbalanced datasets.

In conclusion, evaluating classifiers on imbalanced datasets requires combining threshold metrics with ranking-based metrics. While threshold metrics evaluate specific operational conditions, ranking methods such as ROC and PR curves provide a more holistic understanding of classifier performance.

7.2 Limited Datasets

The lack of large collections of labeled data is a major challenge in medical machine learning, as large annotated datasets require domain expertise and are both time consuming and expensive to produce [124, 125]. An additional challenge is domain shift - when the distribution of test data differs from the training data , leading to degraded model performance [126].

Data scarcity and dataset diversity remain the primary barriers for application in dementia research. Recent studies reinforce the hypothesis that limited data can impact model generalizability which in turn hurts clinical adoption and utility in dementia prediction. [127, 128].

7.2.1 Context and Challenges

All subsequent tasks such as early diagnosis , disease progression, and tracking are impacted heavily by the lack of data. The lack of large scale, heterogeneous datasets (neuroimaging , biomarkers , clinical records), has become the main bottleneck for accurate and general ML models.

7.3 Small Sample Sizes

Even though the application of DL on imaging challenges , especially combined with large image datasets can provide major opportunities , barriers such as data availability and logistical constraints persist [129] . Models often lack generalizability , lack of gold standards and they fail to address privacy concerns as well as safety [130]. Studies point to the lack of efficient methodologies to define an appropriate data sample [131]. Additionally smaller studies often report inflated accuracy metrics indicating that numbers

are upper limits or even superficial at times , meaning there exists systematic bias in the literature [132].

7.4 Image Quality and Variability

MRI lacks arbitrary units. As a consequence scans become non-comparable across sites. Even after normalization , images contain technical artifacts and specific scanner related effects [133]. For volumetric measurements the difference in variance can be up to 10x more depending on the scanner, protocol, and acquisition conditions [134]. More quantified metrics are reported in [135, 136] for scanner maker, model and protocol variability that limit the ability of CNN networks to generalize.

7.5 Interpretability in Medical Applications

Black-box models may be able to achieve high accuracy but they face regulatory scrutiny and difficulties in approval for clinical use. More interpretable methods like kernel-based or prototype-based models offer insights that clinicians can understand and rely upon for further decision making. On the same basis privacy-preserving, uncertainty estimation and fair techniques such as federated learning or diffusion based data generation can align more with the standards of the medical domain and issues that are most important [137] . Explainable AI is a term for explainable systems that improve transparency and support ethical standards , but integrating them into the clinical flow workflows and meeting regulatory standards still remains difficult [138] . ML has made progress and offered advancements in multiple medical domain areas but there are still concerns about bias , privacy and legal accountabilty [139].

7.6 Computational Costs

Deep Learning progress is strongly correlated with rising compute demands , which are quickly becoming economically and envrionmentally unsustainable [140]. Deep learning research explores how to compress, optimize and accelearate how the models run on the hardware. Newer techniques like pruning, quantization and knowledge distillation have emerged and can shrink models by more than 100 times with minimal loss in accuracy [141] . Compact models are necessary for running models on devices with limited computational capacity [142] . MobileNet and EfficientNet are two architectures designed for this purpose, still achieving comparable accuracy despite their small size [143]. Transfer learnign can also be used to adress computational costs in the training phase and it can improve results when data is scarce [144] . Additionally pruning and

knowledge distillation methods are key for reducing the size of larger models [145–147]. Finally there is work that applies these techniques to Alzheimer’s MRI classification [148].

7.7 Data Distillation

Data distillation provides a promising avenue for efficient and secure sharing of medical imaging datasets. Instead of sharing full datasets, distilled representations can preserve modeling performance while reducing data volume. Recent investigations show that a small, representative set of distilled images can achieve near-equivalent model performance, demonstrating potential for scalable clinical collaboration [149].

Chapter 8

Recent Advances And Innovations

8.1 Multi-Modal Integration

Most of the current research on Alzheimer's Disease (AD) classification relies on single modality data. A modality, as defined in [150], is an experience like sound, image, or touch.

Even though the scope of this work is not to review or examine works in data fusion, in order to provide a clearer understanding we should mention the challenges that are provided in one of the recent taxonomies of the field. The main challenges are:

1. Representation
2. Translation
3. Alignment
4. Fusion
5. Co-Learning

One study combined multipled modalities like MRI images, genetic (single Nucleotide Polymorphisms (SNP)) and clinical test scores into a unified deep learning framework. They used stacked denoising autoencoders to rrocess the clinical and genetic inputs, and a 3D CNN for the imaging data. Their deep learning models outperformed traditional classifiers (SVM , Decision Trees, Random Forests , KNN) across accuracy, precision , recall and F1 score. Additionaly multimodal fusion consistently beat single-modality approaches and the model was able to identify the hippocampus region - well established with AD literature- as key feature. [151].

A second study used a multimodal recurrent neural network (RNN) to predict progression from MCI to AD. Their inputs were neuroimaging markers from cross-sections, CSF biomarkers, and cognitive scores, all coming from ADNI. The combination of all modalities led to an improved accuracy by approximately 6 percentage points. The authors also suggested that multimodal approaches could help identify which mCI patients would benefit most from clinical trial enrollment [152].

Another study reported that its deep learning fusion network approach performed better overall in classifying NC, MCI, AD, and nADD across the range of clinical tasks [153].

From the above evidence, it is well understood that more and more research is being implemented in the field by integrating multiple modalities in order to achieve higher classification accuracy, since it seems likely that multi-modal data are more robust and provide a more holistic view of the disease.

8.2 Explainable AI

The role of explainability has become increasingly important in recent years, as deep learning achieves state-of-the-art results and exceeds previous approaches in domains where decision-making is high-stakes (e.g., healthcare, finance).

There are **three main pillars** that Explainable AI (XAI) aims to cover:

1. **Trust:** Allowing the end user to trust and understand the fundamental reasoning behind the decision-making.
2. **Reducing Bias & Increasing Fairness:** In models that produce wrong behaviors—either due to the overrepresentation of certain groups in human-curated data or algorithmic mistakes—the goal is to increase fairness and reduce bias.
3. **Model Improvement & Debugging:** Enabling developers to identify errors in the model's logic and improve performance.

8.3 Definitions and Taxonomy

In the field, there is no clear consensus regarding the definitions of certain terms. The two most debated definitions are **interpretability** and **explainability**. Even though there is a debate surrounding these terms, we will adopt a practical definition to move forward:

- **Interpretability** is considered the ability to understand the model through its parameters or a simple graph. The model can be easily understood by a human (intrinsic explainability).
- **Explainability** is considered the process of providing a valid explanation for the reasoning the model used to arrive at a specific decision (post-hoc explainability).

The **taxonomy** of techniques can be described as:

- **Local & Global:** (Scope of the explanation)
- **Perturbation & Gradient-Based Methods:** (Mathematical approach)
- **Model-Agnostic & Model-Specific:** (Applicability across different model architectures)

For a more concise approach, the most noteworthy techniques are **LIME** and **SHAP**. Both are primarily considered local approaches; however, SHAP can also be considered a global technique because it produces a sum of all feature attributions, allowing for a holistic view of the dataset.

8.3.1 LIME

In LIME, the goal is to produce a function g that approximates the output of the neural network f . It utilizes a weighting variable that is higher for small perturbations (samples that are closer to the original input are weighed more while those that are not are weighed less). The goal of the algorithm is to find a simpler and accurate surrogate model.

$$\xi(x) = \underset{g \in G}{\operatorname{argmin}} \mathcal{L}(f, g, \pi_x) + \Omega(g) \quad (8.1)$$

Where \mathcal{L} represents the **Loss** (fidelity of the surrogate model) and Ω represents the **Complexity** of the explanation model.

8.3.2 SHAP

This technique comes from Shapley Values from Game Theory. It works by assigning a score to each feature that reflects how much that feature contributes to the output when combined with all other features.

$$g(z') = \phi_0 + \sum_{j=1}^M \phi_j \cdot z'_j \quad (8.2)$$

Where z' is the simplified feature (coalition vector) and ϕ is the attribution score assigned to that feature. By aggregating the scores across different examples, SHAP can provide a holistic view of the importance of each feature in the network.

8.3.3 Other Approaches

Other memorable approaches include Saliency Maps (which are gradient-based) and Occlusion Sensitivity (which is perturbation-based).

For the final category involving concept-based explanations, a notable method is **TCAV** (Testing with Concept Activation Vectors).

8.4 Sanity Checks

Even though we can create explanations for a problem, we must be aware of the dangers and understand the capabilities of each approach. For this reason, it is suggested to perform sanity checks to ensure robustness. The sanity checks should include:

1. **Stability:** Models should provide the same (or very similar) explanation for small perturbations to the input, up to a constant ϵ .
2. **Robustness:** Models injected with adversarial examples should allow the explanation method to identify why those examples are widely different from the distribution.
3. **Determinism:** For the exact same input, models should strictly produce the same explanation.

8.5 Related Work

In this work, the authors tried to create a network by combining a predictor with an explainable tool in order to provide accurate diagnosis while using visualization maps to confirm prediction basis. They built a predictor based on an attention mechanism using multi-scale features to teach a network to predict the correct labels representing the input features. They state that the network shows state-of-the-art accuracy and explainability and is able to define critical areas more clearly and with less noise that matches the neuroscience background literature [154].

In another study, researchers used data-level fusion of clinical data, MRI segmentation data, and psychological data. They employed many algorithms with Random

Forest being the best and achieving the highest score value. They also used SHAP explanations for further explainability [155].

Additionally, in a systematic review on explainable AI in Alzheimer's research, it can be seen that most of the approaches contain post-hoc and model-agnostic approaches. Techniques such as SHAP, LIME, Grad-CAM, and LRP are shown to dominate the field. Also interesting are novel approaches using other modalities (speech & text), along with the current trade-off between accuracy and explainability in the domain of XAI [156].

8.6 Data Augmentation & Transfer Learning

Even though Deep Learning has performed tremendously in many computer vision tasks and has become a methodology of choice for analyzing medical images [157], these large networks usually require vast amounts of data in order to avoid overfitting.

Overfitting refers to the process where a model learns a function f with very high variance, perfectly modelling the training data and not being able to generalize to new instances that do not belong to the prior distribution of the training data. In many applications, such as medical imaging, data is limited due to a multitude of factors, but mainly due to the scarcity of datasets annotated by specialists, which is labor-intensive and has a high cost.

Data Augmentation is a collection of techniques used to increase both the quality and the size of an existing dataset to build better deep learning models. The image augmentation algorithms include the following [158] (See Table 8.1 and Figure 8.1).

TABLE 8.1: Image Augmentation Techniques

Technique	Description
1. Geometric transformations	This category includes basic operations like flipping , cropping , rotation, and translation. They are quite easy to apply and are especially beneficial in reducing positional biases. However transformation of this kind can change the labelr (e.g rotating a 6 can produce a 9) so care is needed on examples where the labels can change.
2. Color space augmentations	These transformations change pixel values in their respective color channels or histograms. Examples are changing the brightness or contrast to help the models handle different lighting conditions.

TABLE 8.1: Image Augmentation Techniques (continued)

Technique	Description
3. Kernel filters	This method works by sliding a kernel filter ($n \times n$ matrix) over an image to either blur or sharpen it. Blurring can help with motion blur resistance while sharpening can highlight details. A variant called Patchshuffle Regularization, shuffles pixel values within a window to encourage more robust learning.
4. Mixing images	This approach combines multiple images to create new samples. Techniques can either be average between two images (SamplePairing) or they can be concatenating random crops. Although the images can look unnatural, the method can reduce error rates.
5. Random erasing	Taking inspiration by dropout this method just masks an $n \times m$ patch, with zeros, maximum values, or noise. By hiding different parts it forces the model to look at the entire image, which can help with occlusion problems.
6. Feature space augmentation	Instead of manipulating input images, this technique operates on the lower-dimensional vector representations (feature space) found in high-level network layers. Techniques include adding noise to these vectors or performing interpolations between nearest neighbors (similar to SMOTE) to generate new instances.
7. Adversarial training	This involves using a rival framework to generate "adversarial attacks"—constrained noise injections that cause misclassifications. Using these adversarial examples during training acts as a search algorithm for augmentations, strengthening weak decision boundaries and improving model robustness.
8. GANs	Generative Adversarial Networks (GANs) use a "generator" network to create artificial images and a "discriminator" network to distinguish real from fake ones. This powerful oversampling technique can "unlock" additional information from a dataset and create synthetic training data to increase size and diversity.

TABLE 8.1: Image Augmentation Techniques (continued)

Technique	Description
9. Neural Style Transfer	This algorithm manipulates sequential representations in a CNN to transfer the artistic style or texture of one image to another while preserving content. It is particularly useful for randomizing environments (e.g., lighting and texture) when transferring models from simulations to the real world.
10. Meta-learning	This refers to using neural networks to optimize other neural networks, specifically for finding optimal augmentation strategies. Examples include "Neural Augmentation" (learning style transfer weights), "Smart Augmentation" (merging images via a network), and "AutoAugment" (using Reinforcement Learning to find optimal transformation policies).

FIGURE 8.1: Showcases all the available image augmentation techniques. Courtesy of [158].

Furthermore, distinct approaches have been developed to generate synthetic training data; one study utilized Deep Convolutional Generative Adversarial Networks (DCGANs) to synthesize Positron Emission Tomography (PET) images across three disease stages, effectively overcoming the lack of labeled data [159]. Similarly, other research employed data augmentation within a transfer learning framework to rectify severe class imbalance in 3D Magnetic Resonance Imaging (MRI) datasets from the OASIS dataset [160]. Both studies showed that the methods improved accuracy and diagnosis.

8.6.1 Transfer Learning

Transfer Learning aims as a technique to improve the performance of a learner in a targeted domain by leveraging knowledge of the learner in other domains that are related. By this method, the data of the target domain can be considerably decreased. In medical imaging, whereas it has been mentioned that data scarcity remains an issue, transfer learning seems promising to overcome this obstacle and produce state-of-the-art transfer learners and leverage domains that are similar with more data to construct better networks.

The different categorizations of transfer learning can be seen in Figure 8.2, provided by [161].

FIGURE 8.2: Categorizations of transfer learning [161].

Several approaches in Alzheimer's research leverage transfer learning to detect and classify Alzheimer's disease.

In one study, transfer learning was used to classify four stages of Alzheimer's Disease: Mild Demented, Moderate Demented, Non-Demented, and Very Mild Demented. The deep learning model made use of brain MRI scans and achieved an accuracy of 91.7%, outperforming previous approaches. More specifically, they used a modified AlexNet architecture which was trained on ImageNet datasets as a source domain to perform knowledge transfer [162].

Another study used a VGG architecture with already pre-trained weights. They used the foundational network to perform transfer learning and optimized the network using additional MRI images. They show through experiments that with a size almost 10 times smaller on the OASIS MRI dataset, they can perform comparable to or better than some current deep learning approaches [163].

Chapter 9

Gaps in Current Literature and Future Research Directions

9.1 Generalizability and Standardized Testing

Although AI is performing increasingly well in medical domains, it still lacks the generalizability testing required for wider adoption. This is due to researchers and clinicians having limited access to diverse data—different from the training data—which is needed to test model robustness and reliability. This article provides recommendations for creating “standardized tests” (benchmark datasets) to solve this problem [164].

Another article goes beyond validating models and proposes a framework that is standard across pre-processing and image acquisition protocols. In the case of Alzheimer’s disease, this is even more important, as different scanning protocols hinder the ability of models to generalize or limit the amount of training data. The work focuses on radiomics, but its broader implications for medical imaging and computer-aided diagnosis are relevant to Alzheimer’s and dementia [165].

Other studies in radiomics have demonstrated the impact that different acquisition protocols can have on diagnosis and model accuracy by extensively validating algorithms on diverse datasets. The findings highlight a significant decrease in performance when acquisition diversity is not properly addressed [166].

9.2 Interpretable Models

Taking influence from the pillars previously set for explainability, there is a need for better explanations of model decision-making, especially in domains that affect human life, such as medicine. The pillars discussed are three.

9.2.1 Trust and Clinical Adoption

The first pillar concerns the need to cultivate trust among patients and clinicians, who will be responsible for conducting computer-aided diagnosis. The ability to examine the decision-making process and access the model’s reasoning will vastly enhance clinicians’ trust in these systems. This, in turn, improves both the quality of diagnosis and the number of patients that clinicians can handle, while maintaining consistent performance across individuals.

9.2.2 Reducing Bias and Increasing Fairness

The second pillar focuses on reducing bias and improving fairness. As stated previously, models are constrained by the distribution of the data on which they are trained. The dataset defines the prior distribution and the upper limit of information captured by the model, even when accounting for generalization. If the data are skewed or do not represent a stratifiable distribution of the patient population, models may exhibit bias and lead to incorrect diagnoses. Explainable AI helps mitigate this by providing reasoning behind decisions for clinicians to assess.

9.2.3 Better Model Development and Debugging

The third pillar emphasizes improved model development and system debugging. Access to the model’s reasoning allows developers to identify limitations and provide rigorous debugging to improve performance more broadly.

These are the main pillars that make explainable AI crucial for the further development of systems, trust with end-users, and the avoidance of bias in decision-making. This is why explainability and interpretability remain major challenges in deploying deep learning for disease progression or classification outside research environments.

9.3 Detection at Different Stages

Based on the revised biomarker criteria [167], the integration of MRI and PET biomarkers—or their fusion—is of immense importance. Multimodal fusion can serve as the backbone of Alzheimer’s disease (AD) detection and enhance the ability to track disease progression, which is vital for developing drugs aimed at slowing or curing neurodegeneration.

AD is a complex pathology in which multiple factors contribute to neurodegeneration; early pathological signs, such as amyloid-beta ($A\beta_{42}$) deposition, can be observed up

to 20 years prior to clinical symptoms [168]. Nevertheless, **early prediction remains challenging due to subtle brain changes that are difficult to quantify** [169].

Deep learning and advanced methods—such as data augmentation, synthetic data generation, transfer learning, and multimodal fusion—can significantly enhance early disease classification, enable longitudinal patient monitoring, and help evaluate treatment efficacy.

Chapter 10

Limitations

While there are a lot of studies that showcase significant improvement in disease classification through MRI scans, it is necessary to accept that several limitations exist before any clinical utility can be recognized. These limitations arise from the inherent variability from multiple sources and the different acquisition protocols and systems, from the experiment architectures and the intricacies of different approaches, along with the heterogeneity of the disease. Understanding these limitations is essential for future research and building more robust, reliable and clinically translatable systems.

10.1 Data Leakage

The primary limitation that has been researched in the field of medical imaging and machine learning is that of data leakage, which presents a serious obstacle for clinical translation. As documented in [170], nearly half of the studies between 2017 and 2019 failed to separate same subject data from contaminating the testing set. This leakage creates artificially inflated accuracy metrics that do not reflect generalized diagnostic capability.

The most common cause of data leakage seems to happen through “slice-level” splitting, where MRI slices from the same patient are distributed across both sets (training, testing). This leads to model overfitting, a poor indicator of clinical performance and generalization or broader pattern recognition.

In [171] demonstrated this effect empirically within a single study, showing a 28-percentage-point accuracy drop when switching from slice-wise to subject-wise splitting. The magnitude of this drop—equivalent to the difference between a highly promising diagnostic tool and one barely exceeding chance—illustrates why methodological rigor is non-negotiable.

Furthermore, data contamination can persist in other forms. Longitudinal studies usually include data from the same subject at different disease stages, and if these visits are not carefully tracked, temporal correlations can produce data leakage. Similarly, datasets that aggregate data from multiple acquisition sites may include the same patient scanned at different facilities. In [171], it was found that only 4.5% of published studies implemented the complete “methodological triad” of subject-wise splitting, external validation, and confounder control, suggesting that the field’s reported performance metrics are systematically optimistic.

10.2 The Accuracy Paradox and Class Imbalance

The dataset employed in this study, drawn from the Alzheimer’s Disease Neuroimaging Initiative (ADNI), exhibits significant class imbalance typical of clinical populations. Healthy Controls (HC) and Mild Cognitive Impairment (MCI) cases substantially outnumber confirmed Alzheimer’s Disease cases. This imbalance renders overall accuracy an unreliable and potentially misleading performance metric, a phenomenon termed the “accuracy paradox” by [172].

The dataset employed in this study also exhibits significant class imbalance as can be seen in Figure 10.1.

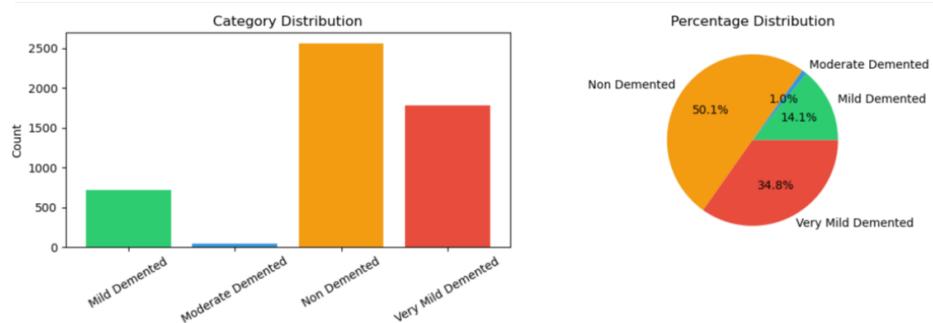


FIGURE 10.1: Class distribution showing the imbalance in the dataset.

This imbalance, as has already been mentioned, renders accuracy as an unreliable metric and more precise and robust metrics have been proposed (F1-score, AUC, MCC, etc.). The phenomenon of imbalance leading to misleading performance has been termed as “accuracy paradox” in [172].

To frame this in practical terms, imagine a medical test that correctly identifies 9 out of 10 people, but only because it tells everyone they are healthy. The 1 in 10 patients with Alzheimer’s Disease would receive false reassurance, potentially delaying critical interventions during the narrow window when treatments are most effective. This

scenario illustrates why clinical AI systems must be evaluated on metrics beyond raw accuracy.

For this reason, the results presented in this thesis prioritize sensitivity, specificity, F1-score, and Area Under the Receiver Operating Characteristic curve (AUC-ROC) alongside accuracy. We employed data augmentation using the MONAI (open source framework), weighted loss functions, and balanced sampling to address class imbalance during training. Moreover, additional frameworks like Grad-CAM can provide interpretability and were used in this study to assess model pattern recognition and be able to provide clinicians with an explanation of what the model evaluates as most important.

Nevertheless, residual bias toward the majority class likely persists. The model may still exhibit lower confidence and higher error rates when classifying true AD cases compared to healthy controls—a bias that would disproportionately affect the patients who most need accurate diagnosis.

10.3 Domain Shift and Generalizability

Most of the models where research is implemented are using curated datasets from research-grade MRI data (ADNI, OASIS) which employ standardized acquisition protocols across different sites fail when deployed on data acquired under different circumstances. While the standardized process reduces confounding variables during development, they introduce a critical limitation known as “domain shift.”

10.3.1 Sources of Domain Shift in MRI Data

Magnetic field strength significantly impacts image resolution and for higher-field strength images contain a better signal-to-noise ratio, but can also amplify certain artifacts. A model thus that has been trained on specific images will fail to accurately predict correctly because of unfamiliar noise patterns or reduced visibility of subtle atrophy.

Acquisition sequence parameters—including echo time, repetition time, flip angle, and slice thickness—vary across institutions and protocols. What constitutes a “standard” T1-weighted sequence differs between research consortia and clinical departments. The model has learned one version of normal and may misinterpret technically valid but differently acquired images.

In this study [173] they documented that models trained on research-grade images from one set frequently fail when applied to clinical-grade images from hospital settings. Dinsdale et al. (2024) [174] found that accuracy dropped to 71% when models were

tested on external datasets with different scanner configurations, compared to training data performance.

Critically, we did not perform external validation on a completely independent cohort (such as training on ADNI and testing on OASIS), meaning our reported metrics almost certainly overestimate real-world clinical performance.

10.4 Shortcut Learning and Region of Interest Bias

Another critical limitation is that deep learning models can learn shortcuts that are contained in medical images instead of actually learning the underlying conditions of the disease. This is called “shortcut learning.” This is also mentioned in interpretability of deep neural networks as a way to ensure model understanding at the programming phase and also enhance clinician trust.

10.4.1 Skull-Stripping Artifacts

Preprocessing pipelines usually include a step for non-brain tissue (skull stripping) to focus the model directly on the brain tissue and its structure. However, as noted in [170], imperfect skull stripping can leave informative artifacts. If there is a difference in the stripping algorithm between healthy brains and atrophic brains—for example, leaving more residual scalp on atrophic brains where the boundary of brain-skull is less distinct—then it is likely that the model may learn to classify images based on preprocessing residuals rather than brain tissue characteristics.

The gap between brain surface and inner skull table increases with cortical atrophy, a hallmark of Alzheimer’s Disease. If skull-stripping is incomplete, this enlarged cerebrospinal fluid space becomes visible as a classification cue. The model correctly identifies Alzheimer’s cases, but for the wrong reason—it has learned to detect atrophy artifacts rather than analyze hippocampal or cortical tissue directly.

10.4.2 Biological Plausibility of Learned Features

In our experiments, we used Gradient-weighted Class Activation Mapping (Grad-CAM) for model interpretability, generating heatmaps indicating image regions that influenced the model’s decision. However, it should be noted that interpretability methods are not validation methods, since a model can focus on ventricular enlargement which is a secondary issue but ignore the primary underlying pathology while still classifying a majority of samples correctly.

In [171] it was found that while 18.2% of studies used interpretability methods, only

12.5% of Grad-CAM implementations validated that highlighted regions corresponded to known neuropathological sites. This interpretability-validation chasm means that impressive heatmaps clustering around “reasonable” brain regions do not confirm that the model has learned disease-relevant features. Clinicians cannot trust explanations that lack rigorous validation against established pathological patterns. Figure 10.2 shows a Grad-CAM implementation along with the code.

```
from pytorch_grad_cam import GradCAM
from pytorch_grad_cam.utils.image import show_cam_on_image
from pytorch_grad_cam.utils.model_targets import ClassifierOutputTarget

def visualize_gradcam(model, img_tensor, true_label, target_layer):
    """Display GradCAM visualization for a single image."""
    img_tensor = img_tensor.to(device)

    # Create GradCAM object
    cam = GradCAM(model=model, target_layers=[target_layer])

    # Get prediction
    model.eval()
    with torch.no_grad():
        probs = F.softmax(model(img_tensor), dim=1)[0]
        pred = probs.argmax().item()

    # Generate GradCAM heatmap
    grayscale_cam = cam(input_tensor=img_tensor,
                         targets=[ClassifierOutputTarget(pred)])[0]

    # Prepare image for visualization
    img = img_tensor.cpu().squeeze().numpy()

    # Normalize image to [0, 1] range
    img_normalized = (img - img.min()) / (img.max() - img.min() + 1e-8)

    # Convert grayscale to RGB
    img_rgb = np.stack([img_normalized]*3, axis=-1).astype(np.float32)

    # Create overlay
    overlay = show_cam_on_image(img_rgb, grayscale_cam, use_rgb=True)
```

LISTING 10.1: Grad-CAM visualization implementation

10.5 Ground Truth Uncertainty and Diagnostic Heterogeneity

Another important limitation is the reliability of diagnostic labels that are used for training. Machine learning models are upper-limited by the labels and data they are given, and it seems that clinical diagnosis of Alzheimer’s Disease is imperfect. In [175] found that 71% of patients with dementia at autopsy had multiple coexisting pathologies, but only 67% had been clinically diagnosed with AD alone.

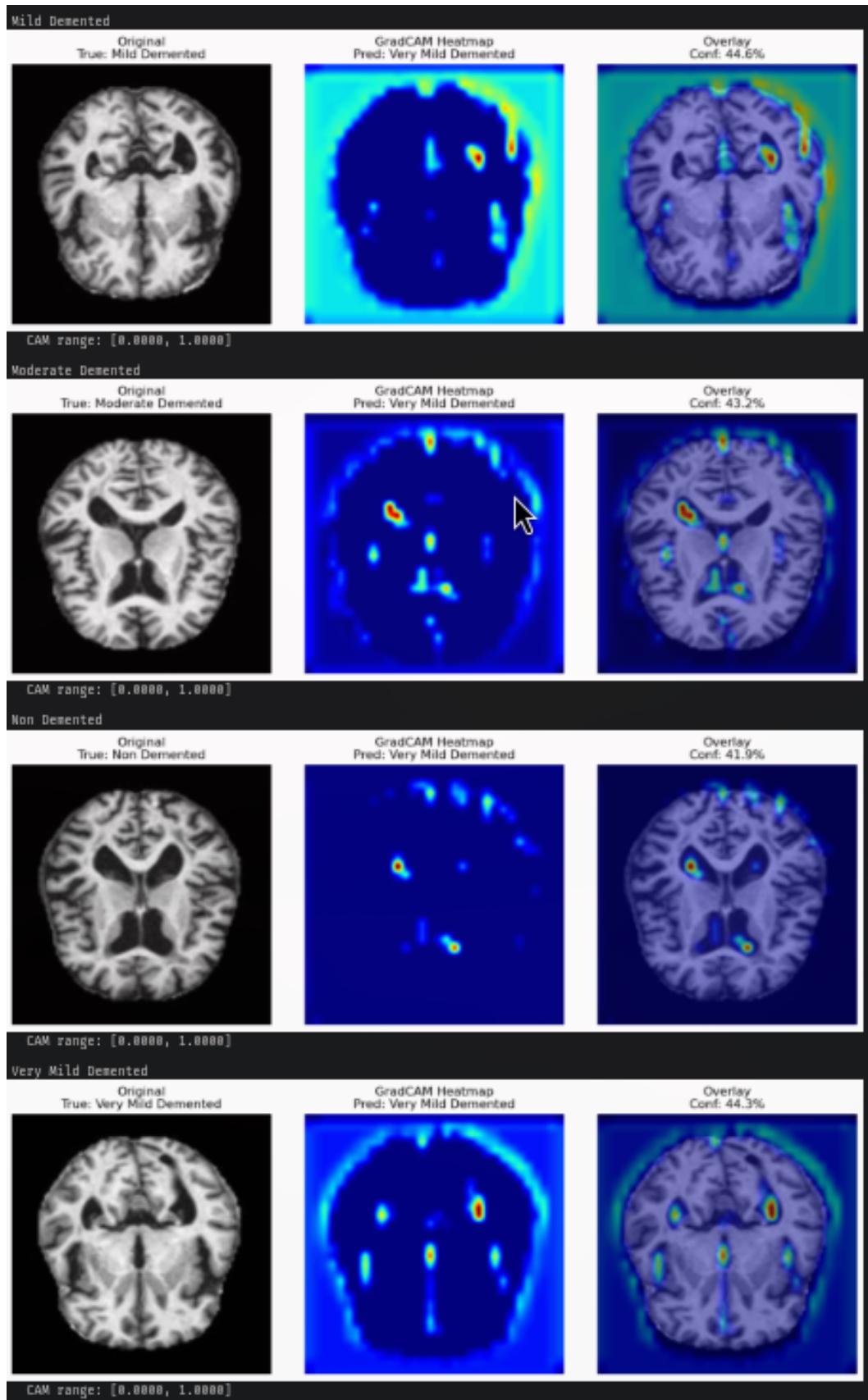


FIGURE 10.2: Grad-CAM visualization showing original image, heatmap, and overlay.

The same study demonstrated that when models were trained only on confirmed neuropathologically diagnoses instead of clinical labels, accuracy changed substantially—both for AD (84.4%) and vascular dementia (83.9%) but lower for Lewy body dementia (62.3%). This suggests that for different dementia subtypes there are different challenges for classification and that incorporating them under imprecise clinical labels degrades performance for all categories.

10.6 Disease Heterogeneity as an Intrinsic Limitation

Alzheimer’s disease itself shows substantial phenotypic heterogeneity. In [176], at least three distinct disease progression patterns were identified. Patients may present with predominantly hippocampal atrophy, predominantly cortical atrophy, or mixed patterns, with different sequences of regional involvement over time.

In another study, Kumar et al. (2024) [177] showed that even greater heterogeneity was found at more severe disease stages. Patients showcased distinct patterns of neurodegeneration, tau accumulation, and amyloid deposition. This presents a limitation where a binary classification between AD and control subjects might be inadequate since the disease manifests through multiple pathways with patient-specific trajectories.

10.7 Transparency, Reproducibility, and the Black Box Problem

Deep learning models operate as “black boxes,” extracting features through millions of parameters without providing human-interpretable explanations for their decisions. Jo et al. (2019) emphasized that this opacity is particularly problematic for medical applications, where clinicians must justify diagnostic decisions to patients, families, and institutional review processes.

Reproducibility presents a related challenge. [174] documented that none of the state-of-the-art studies they examined were fully reproducible. Authors rarely provide complete code, and critical implementation details—data augmentation strategies, exact model architectures, hyperparameter values, random seeds—are inconsistently reported. Sample sizes vary from 170 to 1,662 participants across studies, making cross-study benchmarking nearly impossible. Our own results, while documented as thoroughly as practical, may prove difficult to reproduce exactly due to inherent stochasticity in neural network training.

10.8 Limitations of Single-Modality Analysis

This study relies exclusively on structural MRI data, yet clinical diagnosis of AD integrates multiple information sources. In [178], it is mentioned that at least five modalities exist: structural imaging, functional imaging (PET, fMRI), cerebrospinal fluid biomarkers, genetic testing, and neuropsychological assessment. Relying solely on MRI can negatively impact information that can be leveraged for better model predictions.

However, even though the potential benefits of integrating multiple sources of information are documented [179, 180] and the combination of more modalities could improve accuracy, this can also introduce other complications. Data missingness, variable protocols, and dataset interoperability can increase model complexity and might compromise interpretability.

10.9 Summary of Limitations

The limitations presented in this chapter suggest that the metrics in our results represent upper bounds on clinical utility rather than deployment accuracy. Issues like data leakage may inflate accuracy. Class imbalance renders accuracy as misleading and masks reduced sensitivity for underrepresented classes. Domain shift implies that model performance will degrade if it encounters images from different scanners, protocols, or patient populations.

Shortcut learning means we cannot be confident that the model has learned disease-relevant features rather than confounding artifacts. Ground truth uncertainty implies that our training labels are themselves noisy, introducing irreducible error into the learning process. Disease heterogeneity suggests that the binary classification paradigm may be fundamentally mismatched to the biological reality of Alzheimer’s Disease.

These limitations do not invalidate the work presented in this thesis, but they do circumscribe its claims. The model demonstrates proof-of-concept that convolutional neural networks can learn discriminative features from brain MRI data, but substantial additional work—external validation, prospective clinical trials, integration with existing diagnostic workflows—would be required before any clinical deployment could be contemplated. Future research should prioritize the “methodological triad” identified by Young et al. (2025): rigorous subject-level data splitting, external validation on independent cohorts, and systematic confounder control. Only through such rigorous methodology can the field move from promising laboratory results toward genuinely useful clinical tools.

Chapter 11

Conclusion

The field of research in Dementia and Alzheimer Classification is clearly defined by the Deep Learning Revolution. Most of the research has translated into integrating Deep Learning into multiple steps of the machine learning pipeline to maximize classification accuracy and significant gains in diagnosis. Deep Learning networks are used throughout the pipeline, like in pre-processing (Intensity Normalization, Registration , Skull Stripping , Volumetric Difference Estimation, Denoising) but also as a tool to classify between stages of neurodegeneration.

Additionally we have seen networks trained to learn a prior, like in denoising , to be turned into generative networks for image generation.

Moreover deep neural networks are classified as universal approximators , meaning that they can resemble a function given enough parameters. So based on the ability to show gains connected to the power equations of scaling laws we can conclude that the main bottleneck of the field is data scarcity.

As has been mentioned efforts in the field have generated datasets for research purposes but there is a lack of benchmark datasets to check model generalizability.

Moreover even though the research front can generate competent enough models to detect and classify different stages of dementia the lack of trust and explainability , behind the model's reasoning blocks the usage in clinical settings.

Although these problems exist , methods that can overcome the limitations of data scarcity such as image fusion or data augmentation and even synthetic data have been explored showing promising results.

Finally even though limitations exist , by accounting for the progress that has happened in the field and the ability of models to become better serving as a prior , along

with the fact that models perform state-of-the-art in research, it is highly likely that progress and future research will produce innovation to overcome data scarcity and model generalization as well as generate trust between clinicians , patients and Artificial Intelligence.

Appendix

Pathophysiology of Alzheimer

1. Amyloid- β ($A\beta$) pathology is generally thought to begin with the early deposition of $A\beta_{42}$, a more aggregation-prone and fibrillogenic isoform of the peptide. As these deposits accumulate, they disrupt the surrounding neuronal environment and ultimately influence the stability and function of the tau protein network, which is essential for maintaining axonal structure and intracellular transport. The subsequent tau dysfunction and formation of neurofibrillary tangles contribute directly to progressive neuronal degeneration characteristic of Alzheimer's disease. Although the preferential production or accumulation of $A\beta_{42}$ over other forms such as $A\beta_{40}$ can be influenced by genetic factors, no single mechanism fully accounts for this shift. It is also important to note that many individuals produce $A\beta_{42}$ without developing Alzheimer's disease; pathology emerges when the peptide accumulates beyond the brain's capacity for clearance, leading to plaque formation. [181]

Source Code

The source code for this thesis is available at: <https://github.com/paris26/alzTheBatch>

Bibliography

- [1] P. A. Rowley, A. A. Samsonov, T. J. Betthauser, A. Pirasteh, S. C. Johnson, and L. B. Eisenmenger, “Amyloid and Tau PET Imaging of Alzheimer Disease and Other Neurodegenerative Conditions,” *Seminars in Ultrasound, CT and MRI*, vol. 41, pp. 572–583, Dec. 2020.
- [2] F. Márquez and M. A. Yassa, “Neuroimaging Biomarkers for Alzheimer’s Disease,” *Molecular Neurodegeneration*, vol. 14, p. 21, June 2019.
- [3] K. Kantarci and C. R. Jack, “Neuroimaging in Alzheimer disease: An evidence-based review,” *Neuroimaging Clinics*, vol. 13, no. 2, pp. 197–209, 2003.
- [4] M. S. Rafi and P. S. Aisen, “Detection and treatment of Alzheimer’s disease in its preclinical stage,” *Nature aging*, vol. 3, pp. 520–531, May 2023.
- [5] S. G. Mueller, M. W. Weiner, L. J. Thal, R. C. Petersen, C. Jack, W. Jagust, J. Q. Trojanowski, A. W. Toga, and L. Beckett, “The Alzheimer’s Disease Neuroimaging Initiative,” *Neuroimaging Clinics*, vol. 15, pp. 869–877, Nov. 2005.
- [6] U. N. D. o. E. a. S. Affairs, *World Population Ageing 2019*. United Nations, Oct. 2020.
- [7] “Alzheimer’s Association 2025 Alzheimer’s Disease Facts and Figures,”
- [8] M. Fang, J. Hu, J. Weiss, D. S. Knopman, M. Albert, B. G. Windham, K. A. Walker, A. R. Sharrett, R. F. Gottesman, P. L. Lutsey, T. Mosley, E. Selvin, and J. Coresh, “Lifetime risk and projected burden of dementia,” *Nature Medicine*, vol. 31, pp. 772–776, Mar. 2025.
- [9] M. Prince, G.-C. Ali, M. Guerchet, A. M. Prina, E. Albanese, and Y.-T. Wu, “Recent global trends in the prevalence and incidence of dementia, and survival with dementia,” *Alzheimer’s Research & Therapy*, vol. 8, p. 23, July 2016.
- [10] L. K. Ferreira and G. F. Busatto, “Neuroimaging in Alzheimer’s disease: Current role in clinical practice and potential future applications,” *Clinics*, vol. 66, pp. 19–24, Jan. 2011.

- [11] P. L. Davis, L. E. Crooks, A. R. Margulis, and L. Kaufman, “Nuclear Magnetic Resonance Imaging: Current Capabilities,” *Western Journal of Medicine*, vol. 137, pp. 290–293, Oct. 1982.
- [12] G. Katti, S. A. Ara, and A. Shireen, “Magnetic resonance imaging (MRI)—A review,” *International journal of dental clinics*, vol. 3, no. 1, pp. 65–70, 2011.
- [13] D. B. Plewes and W. Kucharczyk, “Physics of MRI: A primer,” *Journal of Magnetic Resonance Imaging*, vol. 35, pp. 1038–1054, May 2012.
- [14] A. Pai, R. Shetty, B. Hodis, and Y. S. Chowdhury, “Magnetic Resonance Imaging Physics,” in *StatPearls*, Treasure Island (FL): StatPearls Publishing, 2025.
- [15] R.-J. M. van Geuns, P. A. Wielopolski, H. G. de Bruin, B. J. Rensing, P. M. A. van Ooijen, M. Hulshoff, M. Oudkerk, and P. J. de Feyter, “Basic principles of magnetic resonance imaging,” *Progress in Cardiovascular Diseases*, vol. 42, pp. 149–156, Sept. 1999.
- [16] V. P. Grover, J. M. Tognarelli, M. M. Crossey, I. J. Cox, S. D. Taylor-Robinson, and M. J. McPhail, “Magnetic Resonance Imaging: Principles and Techniques: Lessons for Clinicians,” *Journal of Clinical and Experimental Hepatology*, vol. 5, pp. 246–255, Sept. 2015.
- [17] G. B. Frisoni, N. C. Fox, C. R. Jack Jr, P. Scheltens, and P. M. Thompson, “The clinical use of structural MRI in Alzheimer disease,” *Nature reviews neurology*, vol. 6, no. 2, pp. 67–77, 2010.
- [18] H. Braak and E. Braak, “Neuropathological staging of Alzheimer-related changes,” *Acta neuropathologica*, vol. 82, no. 4, pp. 239–259, 1991.
- [19] A. Delacourte, J.-P. David, N. Sergeant, L. Buee, A. Wattez, P. Vermersch, F. Ghozali, C. Fallet-Bianco, F. Pasquier, and F. Lebert, “The biochemical pathway of neurofibrillary degeneration in aging and Alzheimer’s disease,” *Neurology*, vol. 52, no. 6, pp. 1158–1158, 1999.
- [20] B. Dubois, H. H. Feldman, C. Jacova, H. Hampel, J. L. Molinuevo, K. Blennow, S. T. DeKosky, S. Gauthier, D. Selkoe, and R. Bateman, “Advancing research diagnostic criteria for Alzheimer’s disease: The IWG-2 criteria,” *The Lancet Neurology*, vol. 13, no. 6, pp. 614–629, 2014.
- [21] D. W. Townsend, J. P. Carney, J. T. Yap, and N. C. Hall, “PET/CT today and tomorrow,” *Journal of Nuclear Medicine*, vol. 45, no. 1 suppl, pp. 4S–14S, 2004.

- [22] H. Jung, “Basic Physical Principles and Clinical Applications of Computed Tomography,” *Progress in Medical Physics*, vol. 32, pp. 1–17, Mar. 2021.
- [23] S. Basu, T. C. Kwee, S. Surti, E. A. Akin, D. Yoo, and A. Alavi, “Fundamentals of PET and PET/CT imaging,” *Annals of the New York Academy of Sciences*, vol. 1228, pp. 1–18, June 2011.
- [24] B. Dubois, H. H. Feldman, C. Jacova, S. T. Dekosky, P. Barberger-Gateau, J. Cummings, A. Delacourte, D. Galasko, S. Gauthier, G. Jicha, K. Meguro, J. O’Brien, F. Pasquier, P. Robert, M. Rossor, S. Salloway, Y. Stern, P. J. Visser, and P. Scheltens, “Research criteria for the diagnosis of Alzheimer’s disease: Revising the NINCDS-ADRDA criteria,” *The Lancet. Neurology*, vol. 6, pp. 734–746, Aug. 2007.
- [25] L. Mosconi, “Brain glucose metabolism in the early and specific diagnosis of Alzheimer’s disease,” *European Journal of Nuclear Medicine and Molecular Imaging*, vol. 32, pp. 486–510, Apr. 2005.
- [26] G. W. Small, J. C. Mazziotta, M. T. Collins, L. R. Baxter, M. E. Phelps, M. A. Mandelkern, A. Kaplan, A. La Rue, C. F. Adamson, L. Chang, B. H. Guze, E. H. Corder, A. M. Saunders, J. L. Haines, M. A. Pericak-Vance, and A. D. Roses, “Apolipoprotein E Type 4 Allele and Cerebral Glucose Metabolism in Relatives at Risk for Familial Alzheimer Disease,” *JAMA*, vol. 273, pp. 942–947, Mar. 1995.
- [27] A. Nordberg, J. O. Rinne, A. Kadir, and B. Långström, “The use of PET in Alzheimer disease,” *Nature Reviews Neurology*, vol. 6, no. 2, pp. 78–87, 2010.
- [28] S. G. Mueller, M. W. Weiner, L. J. Thal, R. C. Petersen, C. R. Jack, W. Jagust, J. Q. Trojanowski, A. W. Toga, and L. Beckett, “Ways toward an early diagnosis in Alzheimer’s disease: The Alzheimer’s Disease Neuroimaging Initiative (ADNI),” *Alzheimer’s & Dementia*, vol. 1, no. 1, pp. 55–66, 2005.
- [29] K. A. Ellis, A. I. Bush, D. Darby, D. De Fazio, J. Foster, P. Hudson, N. T. Lautenschlager, N. Lenzo, R. N. Martins, P. Maruff, C. Masters, A. Milner, K. Pike, C. Rowe, G. Savage, C. Szoek, K. Taddei, V. Villemagne, M. Woodward, D. Ames, and AIBL Research Group, “The Australian Imaging, Biomarkers and Lifestyle (AIBL) study of aging: Methodology and baseline characteristics of 1112 individuals recruited for a longitudinal study of Alzheimer’s disease,” *International Psychogeriatrics*, vol. 21, pp. 672–687, Aug. 2009.
- [30] D. S. Marcus, A. F. Fotenos, J. G. Csernansky, J. C. Morris, and R. L. Buckner, “Open access series of imaging studies: Longitudinal MRI data in nondemented

and demented older adults,” *Journal of Cognitive Neuroscience*, vol. 22, pp. 2677–2684, Dec. 2010.

- [31] E. Goceri, “Fully Automated and Adaptive Intensity Normalization Using Statistical Features for Brain MR Images,” *Celal Bayar University Journal of Science*, vol. 14, pp. 125–134, Mar. 2018.
- [32] X. Sun, L. Shi, Y. Luo, W. Yang, H. Li, P. Liang, K. Li, V. C. T. Mok, W. C. W. Chu, and D. Wang, “Histogram-based normalization technique on human brain magnetic resonance images from different acquisitions,” *Biomedical Engineering Online*, vol. 14, p. 73, July 2015.
- [33] M. Shah, Y. Xiao, N. Subbanna, S. Francis, D. L. Arnold, D. L. Collins, and T. Arbel, “Evaluating intensity normalization on MRIs of human brain with multiple sclerosis,” *Medical Image Analysis*, vol. 15, pp. 267–282, Apr. 2011.
- [34] U. Bağcı, J. K. Udupa, and L. Bai, “The role of intensity standardization in medical image registration,” *Pattern Recognition Letters*, vol. 31, pp. 315–323, Mar. 2010.
- [35] R. T. Shinohara, E. M. Sweeney, J. Goldsmith, N. Shiee, F. J. Mateen, P. A. Calabresi, S. Jarso, D. L. Pham, D. S. Reich, and C. M. Crainiceanu, “Statistical normalization techniques for magnetic resonance imaging,” *NeuroImage: Clinical*, vol. 6, pp. 9–19, Jan. 2014.
- [36] L. G. Nyúl, J. K. Udupa, and X. Zhang, “New variants of a method of MRI scale standardization,” *IEEE transactions on medical imaging*, vol. 19, no. 2, pp. 143–150, 2000.
- [37] L. G. Nyúl and J. K. Udupa, “Method for standardizing the MR image intensity scale,” June 2003.
- [38] P. Salome, F. Sforazzini, G. Brugnara, A. Kudak, M. Dostal, C. Herold-Mende, S. Heiland, J. Debus, A. Abdollahi, and M. Knoll, “MR Intensity Normalization Methods Impact Sequence Specific Radiomics Prognostic Model Performance in Primary and Recurrent High-Grade Glioma,” *Cancers*, vol. 15, p. 965, Jan. 2023.
- [39] M. Kociołek, M. Strzelecki, and R. Obuchowicz, “Does image normalization and intensity resolution impact texture classification?,” *Computerized Medical Imaging and Graphics*, vol. 81, p. 101716, Apr. 2020.
- [40] A. Pandey and A. Jain, “Comparative analysis of KNN algorithm using various normalization techniques,” *International Journal of Computer Network and Information Security*, vol. 10, no. 11, p. 36, 2017.

- [41] S. G. K. Patro and K. K. Sahu, “Normalization: A Preprocessing Stage,” Mar. 2015.
- [42] G. Collewet, M. Strzelecki, and F. Mariette, “Influence of MRI acquisition protocols and image intensity normalization methods on texture classification,” *Magnetic Resonance Imaging*, vol. 22, pp. 81–91, Jan. 2004.
- [43] S. Albert, B. D. Wichtmann, W. Zhao, A. Maurer, J. Hesser, U. I. Attenberger, L. R. Schad, and F. G. Zöllner, “Comparison of Image Normalization Methods for Multi-Site Deep Learning,” *Applied Sciences*, vol. 13, p. 8923, Jan. 2023.
- [44] F. Orlhac, J. J. Eertink, A.-S. Cottreau, J. M. Zijlstra, C. Thieblemont, M. Meignan, R. Boellaard, and I. Buvat, “A Guide to ComBat Harmonization of Imaging Biomarkers in Multicenter Studies,” *Journal of Nuclear Medicine*, vol. 63, pp. 172–179, Feb. 2022.
- [45] U. Michelucci, “An Introduction to Autoencoders,” Jan. 2022.
- [46] P.-L. Delisle, B. Anctil-Robitaille, C. Desrosiers, and H. Lombaert, “Realistic image normalization for multi-Domain segmentation,” *Medical Image Analysis*, vol. 74, p. 102191, Dec. 2021.
- [47] C. Xu, Y. Sun, Y. Zhang, T. Liu, X. Wang, D. Hu, S. Huang, J. Li, F. Zhang, and G. Li, “Stain Normalization of Histopathological Images Based on Deep Learning: A Review,” *Diagnostics*, vol. 15, no. 8, p. 1032, 2025.
- [48] A. Buades, B. Coll, and J. M. Morel, “A Review of Image Denoising Algorithms, with a New One,” *Multiscale Modeling & Simulation*, vol. 4, pp. 490–530, Jan. 2005.
- [49] A. Danielyan, V. Katkovnik, and K. Egiazarian, “BM3D frames and variational image deblurring,” *IEEE Transactions on image processing*, vol. 21, no. 4, pp. 1715–1728, 2011.
- [50] M. Elad, B. Kawar, and G. Vaksman, “Image Denoising: The Deep Learning Revolution and Beyond—A Survey Paper,” *SIAM Journal on Imaging Sciences*, vol. 16, pp. 1594–1654, Sept. 2023.
- [51] P. Kalavathi and V. B. S. Prasath, “Methods on Skull Stripping of MRI Head Scan Images—a Review,” *Journal of Digital Imaging*, vol. 29, pp. 365–379, June 2016.
- [52] A. Fatima, A. R. Shahid, B. Raza, T. M. Madni, and U. I. Janjua, “State-of-the-Art Traditional to the Machine- and Deep-Learning-Based Skull Strip-

ping Techniques, Models, and Algorithms,” *Journal of Digital Imaging*, vol. 33, pp. 1443–1464, Dec. 2020.

- [53] H. Z. U. Rehman, H. Hwang, and S. Lee, “Conventional and Deep Learning Methods for Skull Stripping in Brain MRI,” *Applied Sciences*, vol. 10, p. 1773, Jan. 2020.
- [54] F. Ségonne, A. M. Dale, E. Busa, M. Glessner, D. Salat, H. K. Hahn, and B. Fischl, “A hybrid approach to the skull stripping problem in MRI,” *NeuroImage*, vol. 22, pp. 1060–1075, July 2004.
- [55] P. Novosad, V. Fonov, and D. L. Collins, “Accurate and robust segmentation of neuroanatomy in T1-weighted MRI by combining spatial priors with deep convolutional neural networks,” Feb. 2019.
- [56] R. Souza, O. Lucena, J. Garrafa, D. Gobbi, M. Saluzzi, S. Appenzeller, L. Rittner, R. Frayne, and R. Lotufo, “An open, multi-vendor, multi-field-strength brain MR dataset and analysis of publicly available skull stripping methods agreement,” *NeuroImage*, vol. 170, pp. 482–494, Apr. 2018.
- [57] S. F. Eskildsen, P. Coupé, V. Fonov, J. V. Manjón, K. K. Leung, N. Guizard, S. N. Wassef, L. R. Østergaard, D. L. Collins, and Alzheimer’s Disease Neuroimaging Initiative, “B_EaST: Brain extraction based on nonlocal segmentation technique,” *NeuroImage*, vol. 59, pp. 2362–2373, Feb. 2012.
- [58] J. E. Iglesias, C.-Y. Liu, P. M. Thompson, and Z. Tu, “Robust brain extraction across datasets and comparison with publicly available methods,” *IEEE transactions on medical imaging*, vol. 30, pp. 1617–1634, Sept. 2011.
- [59] D. W. Shattuck, S. R. Sandor-Leahy, K. A. Schaper, D. A. Rottenberg, and R. M. Leahy, “Magnetic resonance image tissue classification using a partial volume model,” *NeuroImage*, vol. 13, pp. 856–876, May 2001.
- [60] S. M. Smith, “Fast robust automated brain extraction,” *Human Brain Mapping*, vol. 17, pp. 143–155, Nov. 2002.
- [61] F. Isensee, M. Schell, I. Pflueger, G. Brugnara, D. Bonekamp, U. Neuberger, A. Wick, H.-P. Schlemmer, S. Heiland, W. Wick, M. Bendszus, K. H. Maier-Hein, and P. Kickingereder, “Automated brain extraction of multisequence MRI using artificial neural networks,” *Human Brain Mapping*, vol. 40, pp. 4952–4964, Dec. 2019.
- [62] R. W. Cox, “AFNI: Software for analysis and visualization of functional magnetic

- resonance neuroimages,” *Computers and Biomedical Research, an International Journal*, vol. 29, pp. 162–173, June 1996.
- [63] J. Kleesiek, G. Urban, A. Hubert, D. Schwarz, K. Maier-Hein, M. Bendszus, and A. Biller, “Deep MRI brain extraction: A 3D convolutional neural network for skull stripping,” *NeuroImage*, vol. 129, pp. 460–469, Apr. 2016.
- [64] A. Hoopes, J. S. Mora, A. V. Dalca, B. Fischl, and M. Hoffmann, “SynthStrip: Skull-stripping for any brain image,” *NeuroImage*, vol. 260, p. 119474, Oct. 2022.
- [65] L. Fisch, S. Zumdicke, C. Barkhau, D. Emden, J. Ernsting, R. Leenings, K. Sarink, N. R. Winter, B. Rissee, U. Dannowski, and T. Hahn, “Deepbet: Fast brain extraction of T1-weighted MRI using Convolutional Neural Networks,” Aug. 2023.
- [66] A. Schulz, E. Dragendorf, K. Wendt, A. Schomakers, E. Bültmann, and D. Wolff, “Skull stripping tools in pediatric T2-weighted MRI scans: A retrospective evaluation of segmentation performance,” *Frontiers in Neuroscience*, vol. 19, p. 1715514, 2025.
- [67] C. Tinauer, M. Sackl, R. Stollberger, R. Schmidt, S. Ropele, and C. Langkammer, “Skull-stripping induces shortcut learning in MRI-based Alzheimer’s disease classification,” *Insights into Imaging*, vol. 16, p. 283, Dec. 2025.
- [68] J. Radua, E. J. Canales-Rodríguez, E. Pomarol-Clotet, and R. Salvador, “Validity of modulation and optimal settings for advanced voxel-based morphometry,” *NeuroImage*, vol. 86, pp. 81–90, Feb. 2014.
- [69] M. T. Duong, S. R. Das, P. Khandelwal, X. Lyu, L. Xie, E. McGrew, N. Dehghani, C. T. McMillan, E. B. Lee, L. M. Shaw, P. A. Yushkevich, D. A. Wolk, I. M. Nasrallah, and Alzheimer’s Disease Neuroimaging Initiative, “Hypometabolic mismatch with atrophy and tau pathology in mixed Alzheimer’s and Lewy body disease,” *Brain*, vol. 148, pp. 1577–1587, May 2025.
- [70] G. Chételat, B. Desgranges, B. Landeau, F. Mézenge, J. B. Poline, V. de la Sayette, F. Viader, F. Eustache, and J.-C. Baron, “Direct voxel-based comparison between grey matter hypometabolism and atrophy in Alzheimer’s disease,” *Brain*, vol. 131, pp. 60–71, Jan. 2008.
- [71] J. L. Whitwell, R. J. Clifford, S. A. Przybelski, J. E. Parisi, M. L. Senjem, B. F. Boeve, D. S. Knopman, R. C. Petersen, D. W. Dickson, and K. A. Josephs, “Temporal-parietal atrophy: A marker of AD pathology independent of clinical diagnosis,” *Neurobiology of aging*, vol. 32, pp. 1531–1541, Sept. 2011.

- [72] J. L. Whitwell, M. M. Shiung, S. Przybelski, S. D. Weigand, D. S. Knopman, B. F. Boeve, R. C. Petersen, and C. R. Jack, “MRI patterns of atrophy associated with progression to AD in amnestic Mild Cognitive Impairment,” *Neurology*, vol. 70, pp. 512–520, Feb. 2008.
- [73] M. Bozzali, M. Filippi, G. Magnani, M. Cercignani, M. Franceschi, E. Schiatti, S. Castiglioni, R. Mossini, M. Falautano, G. Scotti, G. Comi, and A. Falini, “The contribution of voxel-based morphometry in staging patients with mild cognitive impairment,” *Neurology*, vol. 67, pp. 453–460, Aug. 2006.
- [74] L. G. Apostolova and P. M. Thompson, “Mapping Progressive Brain Structural Changes in Early Alzheimer’s Disease and Mild Cognitive Impairment,” *Neuropsychologia*, vol. 46, no. 6, pp. 1597–1612, 2008.
- [75] L. K. Ferreira, B. S. Diniz, O. V. Forlenza, G. F. Busatto, and M. V. Zanetti, “Neurostructural predictors of Alzheimer’s disease: A meta-analysis of VBM studies,” *Neurobiology of Aging*, vol. 32, pp. 1733–1741, Oct. 2011.
- [76] “VBM anticipates the rate of progression of Alzheimer disease | Neurology.” <https://www.neurology.org/doi/abs/10.1212/01.wnl.0000303960.01039.43>.
- [77] K. Ishii, T. Kawachi, H. Sasaki, A. K. Kono, T. Fukuda, Y. Kojima, and E. Mori, “Voxel-Based Morphometric Comparison Between Early- and Late-Onset Mild Alzheimer’s Disease and Assessment of Diagnostic Performance of Z Score Images,” *AJNR: American Journal of Neuroradiology*, vol. 26, pp. 333–340, Feb. 2005.
- [78] K. Ishii, T. Kawachi, H. Sasaki, A. K. Kono, T. Fukuda, Y. Kojima, and E. Mori, “Voxel-Based Morphometric Comparison Between Early- and Late-Onset Mild Alzheimer’s Disease and Assessment of Diagnostic Performance of Z Score Images,” *AJNR: American Journal of Neuroradiology*, vol. 26, pp. 333–340, Feb. 2005.
- [79] H. Huang, S. Zheng, Z. Yang, Y. Wu, Y. Li, J. Qiu, Y. Cheng, P. Lin, Y. Lin, J. Guan, D. J. Mikulis, T. Zhou, and R. Wu, “Voxel-based morphometry and a deep learning model for the diagnosis of early Alzheimer’s disease based on cerebral gray matter changes,” *Cerebral Cortex*, vol. 33, pp. 754–763, Feb. 2023.
- [80] T. Jo, K. Nho, and A. J. Saykin, “Deep Learning in Alzheimer’s Disease: Diagnostic Classification and Prognostic Prediction Using Neuroimaging Data,” *Frontiers in Aging Neuroscience*, vol. 11, p. 220, Aug. 2019.
- [81] J. Islam and Y. Zhang, “Brain MRI analysis for Alzheimer’s disease diagnosis

- using an ensemble system of deep convolutional neural networks,” *Brain Informatics*, vol. 5, p. 2, May 2018.
- [82] “(PDF) Support vector machine-based classification of Alzheimer’s disease from whole-brain anatomical MRI,” *ResearchGate*.
- [83] P. Niranjan Kumar, “SVM-Based Classifier For Early Detection Of Alzheimer’s Disease,” *Educational Administration: Theory and Practice*, pp. 1120–1131, May 2024.
- [84] E. Pellegrini, L. Ballerini, M. d. C. V. Hernandez, F. M. Chappell, V. González-Castro, D. Anblagan, S. Danso, S. Muñoz-Maniega, D. Job, and C. Pernet, “Machine learning of neuroimaging for assisted diagnosis of cognitive impairment and dementia: A systematic review,” *Alzheimer’s & Dementia: Diagnosis, Assessment & Disease Monitoring*, vol. 10, pp. 519–535, 2018.
- [85] A. Sarica, A. Cerasa, and A. Quattrone, “Random Forest Algorithm for the Classification of Neuroimaging Data in Alzheimer’s Disease: A Systematic Review,” *Frontiers in Aging Neuroscience*, vol. 9, Oct. 2017.
- [86] S. I. Dimitriadis and D. Liparas, “How random is the random forest? Random forest algorithm on the service of structural imaging biomarkers for Alzheimer’s disease: From Alzheimer’s disease neuroimaging initiative (ADNI) database,” *Neural Regeneration Research*, vol. 13, pp. 962–970, June 2018.
- [87] M. Song, H. Jung, S. Lee, D. Kim, and M. Ahn, “Diagnostic Classification and Biomarker Identification of Alzheimer’s Disease with Random Forest Algorithm,” *Brain Sciences*, vol. 11, p. 453, Apr. 2021.
- [88] M. Velazquez and Y. Lee, “Random forest model for feature-based Alzheimer’s disease conversion prediction from early mild cognitive impairment subjects,”
- [89] C. Strobl, A.-L. Boulesteix, A. Zeileis, and T. Hothorn, “Bias in random forest variable importance measures: Illustrations, sources and a solution,” *BMC Bioinformatics*, vol. 8, p. 25, Jan. 2007.
- [90] “(PDF) A Comparative Study of PCA and LDA for Dimensionality Reduction in a 4-Way Classification Framework.” https://www.researchgate.net/publication/378806266_A_Comparative_Study_of_PCA_and_Way_Classification_Framework.
- [91] L. Lazli, “Improved Alzheimer Disease Diagnosis With a Machine Learning Approach and Neuroimaging: Case Study Development,” *JMIRx Med*, vol. 6, p. e60866, Apr. 2025.

- [92] V. A. Miller, S. Erlien, and J. Piersol, “Support vector machine classification of dimensionally reduced structural MRI images for dementia,” June 2014.
- [93] M. M. Dessouky, M. A. Elrashidy, and H. M. Abdelkader, “Selecting and Extracting Effective Features for Automated Diagnosis of Alzheimer’s Disease,” *International Journal of Computer Applications*, vol. 81, pp. 17–28, Nov. 2013.
- [94] D. Ozkan, O. Katar, M. Ak, M. A. Al-Antari, N. Yasan Ak, O. Yildirim, H. S. Mir, R.-S. Tan, and U. Rajendra Acharya, “Deep Learning Techniques for Automated Dementia Diagnosis Using Neuroimaging Modalities: A Systematic Review,” *IEEE Access*, vol. 12, pp. 127879–127902, 2024.
- [95] A. Ebrahimi and S. Luo, “Convolutional neural networks for Alzheimer’s disease detection on MRI images,” *Journal of Medical Imaging*, vol. 8, p. 024503, Mar. 2021.
- [96] V. S M., K. D., and N. V C., “Deep Learning-Driven Alzheimer’s Disease Classification: Custom CNN and Pretrained Architectures for Accurate MRI Analysis,” *Journal of Soft Computing Paradigm*, vol. 7, pp. 31–43, Mar. 2025.
- [97] “(PDF) Improved Classification of Alzheimer’s Disease With Convolutional Neural Networks,” in *ResearchGate*.
- [98] M. U. Ali, K. S. Kim, M. Khalid, M. Farrash, A. Zafar, and S. W. Lee, “Enhancing Alzheimer’s disease diagnosis and staging: A multistage CNN framework using MRI,” *Frontiers in Psychiatry*, vol. 15, June 2024.
- [99] B. Khagi and G.-R. Kwon, “3D CNN Design for the Classification of Alzheimer’s Disease Using Brain MRI and PET,” *IEEE Access*, vol. 8, pp. 217830–217847, 2020.
- [100] X. Xu, L. Lin, S. Sun, and S. Wu, “A review of the application of three-dimensional convolutional neural networks for the diagnosis of Alzheimer’s disease using neuroimaging,” *Reviews in the Neurosciences*, vol. 34, pp. 649–670, Aug. 2023.
- [101] W. Lin, T. Tong, Q. Gao, D. Guo, X. Du, Y. Yang, G. Guo, M. Xiao, M. Du, X. Qu, and T. A. D. N. Initiative, “Convolutional Neural Networks-Based MRI Image Analysis for the Alzheimer’s Disease Prediction From Mild Cognitive Impairment,” *Frontiers in Neuroscience*, vol. 12, Nov. 2018.
- [102] Y. Huang, J. Xu, Y. Zhou, T. Tong, and X. Zhuang, “Diagnosis of Alzheimer’s Disease via Multi-Modality 3D Convolutional Neural Network,” *Frontiers in Neuroscience*, vol. 13, p. 509, May 2019.

- [103] Z. Zhao, J. H. Chuah, K. W. Lai, C.-O. Chow, M. Gochoo, S. Dhanalakshmi, N. Wang, W. Bao, and X. Wu, “Conventional machine learning and deep learning in Alzheimer’s disease diagnosis using neuroimaging: A review,” *Frontiers in Computational Neuroscience*, vol. 17, p. 1038636, Feb. 2023.
- [104] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization,” *International Journal of Computer Vision*, vol. 128, pp. 336–359, Feb. 2020.
- [105] M. Wang, “Interpretable 2D and 3D Convolutional Neural Networks for Alzheimer’s Disease in Brain Scans,”
- [106] G. Marcus, “Deep Learning: A Critical Appraisal,”
- [107] F. Konidaris, T. Tagaris, M. Sdraka, and A. Stafylopatis, “Generative Adversarial Networks as an Advanced Data Augmentation Technique for MRI Data;” in *Proceedings of the 14th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, (Prague, Czech Republic), pp. 48–59, SCITEPRESS - Science and Technology Publications, 2019.
- [108] X. Zhou, S. Qiu, P. S. Joshi, C. Xue, R. J. Killiany, A. Z. Mian, S. P. Chin, R. Au, and V. B. Kolachalama, “Enhancing magnetic resonance imaging-driven Alzheimer’s disease classification performance using generative adversarial learning,” *Alzheimer’s Research & Therapy*, vol. 13, p. 60, Mar. 2021.
- [109] “(PDF) Vision Transformers in Medical Imaging: A Comprehensive Review of Advancements and Applications Across Multiple Diseases.” https://www.researchgate.net/publication/390372350_Vision_Transformers_in_Medical_Im
- [110] “Vision transformer architecture and applications in digital health: A tutorial and survey - PMC.” <https://pmc.ncbi.nlm.nih.gov/articles/PMC10333157/>.
- [111] P. T. Krishnan, P. Krishnadoss, M. Khandelwal, D. Gupta, A. Nihaal, and T. S. Kumar, “Enhancing brain tumor detection in MRI with a rotation invariant Vision Transformer,” *Frontiers in Neuroinformatics*, vol. 18, June 2024.
- [112] “Hybrid CNN-SVM for Alzheimer’s Disease Classification from Structural MRI and the Alzheimer’s Disease Neuroimaging Initiative (ADNI),” in *2018 International Conference on Biomedical Engineering, Machinery and Earth Science (BEMES 2018)*, Francis Academic Press, 2018.
- [113] H. M and S. M.N, “A Review on Evaluation Metrics for Data Classification

Evaluations,” *International Journal of Data Mining & Knowledge Management Process*, vol. 5, pp. 01–11, Mar. 2015.

- [114] S. Raschka, “Model Evaluation, Model Selection, and Algorithm Selection in Machine Learning,” Nov. 2020.
- [115] J. X. C. Ke, A. DhakshinaMurthy, R. B. George, and P. Branco, “The effect of resampling techniques on the performances of machine learning clinical risk prediction models in the setting of severe class imbalance: Development and internal validation in a retrospective cohort,” *Discover Artificial Intelligence*, vol. 4, p. 91, Nov. 2024.
- [116] E. Richardson, R. Trevizani, J. A. Greenbaum, H. Carter, M. Nielsen, and B. Peters, “The receiver operating characteristic curve accurately assesses imbalanced datasets,” *Patterns*, vol. 5, p. 100994, May 2024.
- [117] D. Chicco and G. Jurman, “The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation,” *BMC Genomics*, vol. 21, p. 6, Jan. 2020.
- [118] A. Jain, H. Patel, L. Nagalapatti, N. Gupta, S. Mehta, S. Guttula, S. Mujumdar, S. Afzal, R. Sharma Mittal, and V. Munigala, “Overview and Importance of Data Quality for Machine Learning Tasks,” in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, (Virtual Event CA USA), pp. 3561–3562, ACM, Aug. 2020.
- [119] J. M. Johnson and T. M. Khoshgoftaar, “Survey on deep learning with class imbalance,” *Journal of Big Data*, vol. 6, p. 27, Dec. 2019.
- [120] N. U. Niaz, K. N. Shahriar, and M. J. A. Patwary, “Class Imbalance Problems in Machine Learning: A Review of Methods And Future Challenges,” in *Proceedings of the 2nd International Conference on Computing Advancements*, (Dhaka Bangladesh), pp. 485–490, ACM, Mar. 2022.
- [121] X. Guo, Y. Yin, C. Dong, G. Yang, and G. Zhou, “On the class imbalance problem,” in *2008 Fourth International Conference on Natural Computation*, vol. 4, pp. 192–201, IEEE, 2008.
- [122] N. Japkowicz and S. Stephen, “The class imbalance problem: A systematic study1,” *Intelligent Data Analysis*, vol. 6, pp. 429–449, Nov. 2002.
- [123] N. Japkowicz, “Assessment Metrics for Imbalanced Learning,” in *Imbalanced Learning* (H. He and Y. Ma, eds.), pp. 187–206, Wiley, 1 ed., June 2013.

- [124] F. Yao, “Machine learning with limited data,” Jan. 2021.
- [125] A. Merkin, R. Krishnamurthi, and O. N. Medvedev, “Machine learning, artificial intelligence and the prediction of dementia,” *Current Opinion in Psychiatry*, vol. 35, p. 123, Mar. 2022.
- [126] S. A. Martin, F. J. Townend, F. Barkhof, and J. H. Cole, “Interpretable machine learning for dementia: A systematic review,” *Alzheimer’s & Dementia*, vol. 19, pp. 2135–2149, May 2023.
- [127] Y.-C. Huang, T.-C. Liu, and C.-J. Lu, “Establishing a machine learning dementia progression prediction model with multiple integrated data,” *BMC Medical Research Methodology*, vol. 24, p. 288, Nov. 2024.
- [128] Y. Wang, S. Liu, A. G. Spiteri, A. L. H. Huynh, C. Chu, C. L. Masters, B. Goudey, Y. Pan, and L. Jin, “Understanding machine learning applications in dementia research and clinical practice: A review for biomedical scientists and clinicians,” *Alzheimer’s Research & Therapy*, vol. 16, p. 175, Aug. 2024.
- [129] “Challenges for machine learning in clinical translation of big data imaging studies,” *Neuron*, vol. 110, pp. 3866–3881, Dec. 2022.
- [130] M. Aljuhani, A. Ashraf, and P. Edison, “Use of Artificial Intelligence in Imaging Dementia,” *Cells*, vol. 13, p. 1965, Nov. 2024.
- [131] “Sample-Size Determination Methodologies for Machine Learning in Medical Imaging Research: A Systematic Review,” *Canadian Association of Radiologists Journal*, vol. 70, pp. 344–353, Nov. 2019.
- [132] L. Lu, Q. S. Phua, S. Bacchi, R. Goh, A. K. Gupta, J. G. Kovoor, C. D. Ovenden, and M.-S. To, “Small Study Effects in Diagnostic Imaging Accuracy,” *JAMA Network Open*, vol. 5, p. e2228776, Aug. 2022.
- [133] T. L. S. Benzinger, T. Blazey, C. R. Jack, R. A. Koeppe, Y. Su, C. Xiong, M. E. Raichle, A. Z. Snyder, B. M. Ances, R. J. Bateman, N. J. Cairns, A. M. Fagan, A. Goate, D. S. Marcus, P. S. Aisen, J. J. Christensen, L. Ercole, R. C. Hornbeck, A. M. Farrar, P. Aldea, M. S. Jasielec, C. J. Owen, X. Xie, R. Mayeux, A. Brickman, E. McDade, W. Klunk, C. A. Mathis, J. Ringman, P. M. Thompson, B. Ghetti, A. J. Saykin, R. A. Sperling, K. A. Johnson, S. Salloway, S. Correia, P. R. Schofield, C. L. Masters, C. Rowe, V. L. Villemagne, R. Martins, S. Ourselin, M. N. Rossor, N. C. Fox, D. M. Cash, M. W. Weiner, D. M. Holtzman, V. D. Buckles, K. Moulder, and J. C. Morris, “Regional variability of imag-

ing biomarkers in autosomal dominant Alzheimer's disease," *Proceedings of the National Academy of Sciences*, vol. 110, Nov. 2013.

- [134] F. Kruggel, J. Turner, L. T. Muftuler, and A. D. N. Initiative, "Impact of scanner hardware and imaging protocol on image quality and compartment volume precision in the ADNI cohort," *Neuroimage*, vol. 49, no. 3, pp. 2123–2133, 2010.
- [135] I. Montero, S. Sotoudeh-Paima, E. Abadi, and E. Samei, "Intra- and inter-scanner CT variability and their impact on diagnostic tasks," *Proceedings of SPIE—the International Society for Optical Engineering*, vol. 13405, p. 134054C, Feb. 2025.
- [136] S. Bhosekar, P. Singh, D. Garg, V. Ravi, and M. Diwakar, "A Review of Deep Learning-based Multi-modal Medical Image Fusion,"
- [137] A. Begüm Bektaş and M. Gönen, "Machine Learning for Medicine Must Be Interpretable, Shareable, Reproducible and Accountable by Design," *arXiv e-prints*, pp. arXiv–2508, 2025.
- [138] "A survey of explainable artificial intelligence in healthcare: Concepts, applications, and challenges," *Informatics in Medicine Unlocked*, vol. 51, p. 101587, Jan. 2024.
- [139] "Machine Learning in Healthcare: A Review of Current Applications and Future Trends." <https://ieeexplore.ieee.org/document/10968281>.
- [140] N. C. Thompson, K. Greenewald, K. Lee, and G. F. Manso, "The Computational Limits of Deep Learning," July 2022.
- [141] "(PDF) Efficient Deep Learning: A Survey of Model Compression and Optimization Techniques for Resource-Constrained Environments," Aug. 2025.
- [142] H.-I. Liu, M. Galindo, H. Xie, L.-K. Wong, H.-H. Shuai, Y.-H. Li, and W.-H. Cheng, "Lightweight Deep Learning for Resource-Constrained Environments: A Survey," Apr. 2024.
- [143] K. B. Nampalle, P. Singh, U. V. Narayan, and B. Raman, "DeepMediX: A Deep Learning-Driven Resource-Efficient Medical Diagnosis Across the Spectrum," July 2023.
- [144] P. K. Dash and D. S. Sisodia, "Transfer learning based lightweight model for classification of Alzheimer's disease using brain MR images," in *2024 OPJU International Technology Conference (OTCON) on Smart Computing for Innovation and Advancement in Industry 4.0*, pp. 1–6, June 2024.

- [145] J. Gou, B. Yu, S. J. Maybank, and D. Tao, “Knowledge Distillation: A Survey,” *International Journal of Computer Vision*, vol. 129, pp. 1789–1819, June 2021.
- [146] S. Vadera and S. Ameen, “Methods for Pruning Deep Neural Networks,” *IEEE Access*, vol. 10, pp. 63280–63300, 2022.
- [147] T. Hoefer, D. Alistarh, T. Ben-Nun, and N. Dryden, “Sparsity in Deep Learning: Pruning and growth for efficient inference and training in neural networks,”
- [148] Y. Li, J. Luo, and J. Zhang, “Classification of Alzheimer’s disease in MRI images using knowledge distillation framework: An investigation,” *International Journal of Computer Assisted Radiology and Surgery*, vol. 17, pp. 1235–1243, July 2022.
- [149] M. Li, C. Cui, Q. Liu, R. Deng, T. Yao, M. Lionts, and Y. Huo, “Dataset Distillation in Medical Imaging: A Feasibility Study,” Feb. 2025.
- [150] C. Cui, H. Yang, Y. Wang, S. Zhao, Z. Asad, L. A. Coburn, K. T. Wilson, B. A. Landman, and Y. Huo, “Deep multimodal fusion of image and non-image data in disease diagnosis and prognosis: A review,” *Progress in biomedical engineering (Bristol, England)*, vol. 5, pp. 10.1088/2516–1091/acc2fe, Apr. 2023.
- [151] J. Venugopalan, L. Tong, H. R. Hassanzadeh, and M. D. Wang, “Multimodal deep learning models for early detection of Alzheimer’s disease stage,” *Scientific reports*, vol. 11, no. 1, p. 3254, 2021.
- [152] G. Lee, K. Nho, B. Kang, K.-A. Sohn, and D. Kim, “Predicting Alzheimer’s disease progression using multi-modal deep learning approach,” *Scientific reports*, vol. 9, no. 1, p. 1952, 2019.
- [153] S. Qiu, M. I. Miller, P. S. Joshi, J. C. Lee, C. Xue, Y. Ni, Y. Wang, I. De Anda-Duran, P. H. Hwang, J. A. Cramer, B. C. Dwyer, H. Hao, M. C. Kaku, S. Kedar, P. H. Lee, A. Z. Mian, D. L. Murman, S. O’Shea, A. B. Paul, M.-H. Saint-Hilaire, E. Alton Sartor, A. R. Saxena, L. C. Shih, J. E. Small, M. J. Smith, A. Swaminathan, C. E. Takahashi, O. Taraschenko, H. You, J. Yuan, Y. Zhou, S. Zhu, M. L. Alosco, J. Mez, T. D. Stein, K. L. Poston, R. Au, and V. B. Kolachalama, “Multimodal deep learning for Alzheimer’s disease dementia assessment,” *Nature Communications*, vol. 13, p. 3404, June 2022.
- [154] L. Yu, W. Xiang, J. Fang, Y.-P. P. Chen, and R. Zhu, “A novel explainable neural network for Alzheimer’s disease diagnosis,” *Pattern Recognition*, vol. 131, p. 108876, 2022.
- [155] S. Jahan, K. Abu Taher, M. S. Kaiser, M. Mahmud, M. S. Rahman, A. S. Hosen,

- and I.-H. Ra, "Explainable AI-based Alzheimer's prediction and management using multimodal data," *Plos one*, vol. 18, no. 11, p. e0294253, 2023.
- [156] V. Viswan, N. Shaffi, M. Mahmud, K. Subramanian, and F. Hajamohideen, "Explainable Artificial Intelligence in Alzheimer's Disease Classification: A Systematic Review," *Cognitive Computation*, vol. 16, pp. 1–44, Jan. 2024.
- [157] Geert Litjens, G. Litjens, Thijs Kooi, T. Kooi, Babak Ehteshami Bejnordi, B. E. Bejnordi, Arnaud Arindra Adiyoso Setio, A. A. A. Setio, Francesco Ciompi, F. Ciompi, Mohsen Ghafoorian, M. Ghafoorian, Jeroen van der Laak, J. van der Laak, Bram van Ginneken, B. van Ginneken, Clara I. Sánchez, and C. I. Sánchez, "A survey on deep learning in medical image analysis," *Medical Image Analysis*, vol. 42, pp. 60–88, Dec. 2017.
- [158] C. Shorten and T. M. Khoshgoftaar, "A survey on Image Data Augmentation for Deep Learning," *Journal of Big Data*, vol. 6, p. 60, Dec. 2019.
- [159] M. Sajjad, F. Ramzan, M. U. G. Khan, A. Rehman, M. Kolivand, S. M. Fati, and S. A. Bahaj, "Deep convolutional generative adversarial network for Alzheimer's disease classification using positron emission tomography (PET) and synthetic data augmentation," *Microscopy Research and Technique*, vol. 84, pp. 3023–3034, Dec. 2021.
- [160] S. Afzal, M. Maqsood, F. Nazir, U. Khan, F. Aadil, K. M. Awan, I. Mehmood, and O.-Y. Song, "A data augmentation-based framework to handle class imbalance problem for Alzheimer's stage detection," *IEEE access*, vol. 7, pp. 115528–115539, 2019.
- [161] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He, "A comprehensive survey on transfer learning," *Proceedings of the IEEE*, vol. 109, no. 1, pp. 43–76, 2020.
- [162] T. M. Ghazal, S. Abbas, S. Munir, M. A. Khan, M. Ahmad, G. F. Issa, S. Binish Zahra, M. Adnan Khan, and M. Kamrul Hasan, "Alzheimer Disease Detection Empowered with Transfer Learning," *Computers, Materials & Continua*, vol. 70, no. 3, pp. 5005–5019, 2022.
- [163] M. Hon and N. Khan, "Towards Alzheimer's Disease Classification through Transfer Learning," Nov. 2017.
- [164] N. Sourlos, R. Vliegenthart, J. Santinha, M. E. Klontzas, R. Cuocolo, M. Huisman, and P. van Ooijen, "Recommendations for the creation of benchmark

- datasets for reproducible artificial intelligence in radiology,” *Insights into Imaging*, vol. 15, p. 248, Oct. 2024.
- [165] M. Cobo, P. Menéndez Fernández-Miranda, G. Bastarrika, and L. Lloret Iglesias, “Enhancing radiomics and Deep Learning systems through the standardization of medical imaging workflows,” *Scientific Data*, vol. 10, p. 732, Oct. 2023.
- [166] K. Namdar, M. W. Wagner, B. B. Ertl-Wagner, and F. Khalvati, “Open-radiomics: A collection of standardized datasets and a technical protocol for reproducible radiomics machine learning pipelines,” *BMC Medical Imaging*, vol. 25, p. 312, Aug. 2025.
- [167] C. R. Jack, D. A. Bennett, K. Blennow, M. C. Carrillo, B. Dunn, S. B. Haeberlein, D. M. Holtzman, W. Jagust, F. Jessen, J. Karlawish, E. Liu, J. L. Molinuevo, T. Montine, C. Phelps, K. P. Rankin, C. C. Rowe, P. Scheltens, E. Siemers, H. M. Snyder, R. Sperling, Contributors, C. Elliott, E. Masliah, L. Ryan, and N. Silverberg, “NIA-AA Research Framework: Toward a biological definition of Alzheimer’s disease,” *Alzheimer’s & Dementia*, vol. 14, pp. 535–562, Apr. 2018.
- [168] R. J. Bateman, C. Xiong, T. L. Benzinger, A. M. Fagan, A. Goate, N. C. Fox, D. S. Marcus, N. J. Cairns, X. Xie, T. M. Blazey, D. M. Holtzman, A. Santacruz, V. Buckles, A. Oliver, K. Moulder, P. S. Aisen, B. Ghetti, W. E. Klunk, E. McDade, R. N. Martins, C. L. Masters, R. Mayeux, J. M. Ringman, M. N. Rossor, P. R. Schofield, R. A. Sperling, S. Salloway, and J. C. Morris, “Clinical and Biomarker Changes in Dominantly Inherited Alzheimer’s Disease,” *New England Journal of Medicine*, vol. 367, pp. 795–804, Aug. 2012.
- [169] S. Rathore, M. Habes, M. A. Iftikhar, A. Shacklett, and C. Davatzikos, “A review on neuroimaging-based classification studies and associated feature extraction methods for Alzheimer’s disease and its prodromal stages,” *NeuroImage*, vol. 155, pp. 530–548, 2017.
- [170] J. Wen, E. Thibreau-Sutre, M. Diaz-Melo, J. Samper-González, A. Routier, S. Bottani, D. Dormont, S. Durrleman, N. Burgos, and O. Colliot, “Convolutional neural networks for classification of Alzheimer’s disease: Overview and reproducible evaluation,” *Medical Image Analysis*, vol. 63, p. 101694, July 2020.
- [171] V. M. Young, S. Gates, L. Y. Garcia, and A. Salardini, “Data Leakage in Deep Learning for Alzheimer’s Disease Diagnosis: A Scoping Review of Methodological Rigor and Performance Inflation,” *Diagnostics*, vol. 15, p. 2348, Sept. 2025.
- [172] B. Tong, Z. Zhou, D. A. Tarzanagh, B. Hou, A. J. Saykin, J. Moore, M. Ritchie, and L. Shen, “Class-Balanced Deep Learning with Adaptive Vector Scaling Loss

- for Dementia Stage Detection,” *Machine learning in medical imaging. MLMI (Workshop)*, vol. 14349, pp. 144–154, 2024.
- [173] V. S. Diogo, H. A. Ferreira, D. Prata, and Alzheimer’s Disease Neuroimaging Initiative, “Early diagnosis of Alzheimer’s disease using machine learning: A multi-diagnostic, generalizable approach,” *Alzheimer’s Research & Therapy*, vol. 14, p. 107, Aug. 2022.
- [174] R. Turrisi, A. Verri, and A. Barla, “Deep learning-based Alzheimer’s disease detection: Reproducibility and the effect of modeling choices,” *Frontiers in Computational Neuroscience*, vol. 18, p. 1360095, Sept. 2024.
- [175] D. Wang, N. Honnorat, J. B. Toledo, K. Li, S. Charisis, T. Rashid, A. Benet Nir-mala, S. R. Brandigampala, M. Mojtabai, S. Seshadri, M. Habes, and the Alzheimer’s Disease Neuroimaging Initiative, “Deep learning reveals pathology-confirmed neuroimaging signatures in Alzheimer’s, vascular and Lewy body dementias,” *Brain*, vol. 148, pp. 1963–1977, June 2025.
- [176] A. L. Young, R. V. Marinescu, N. P. Oxtoby, M. Bocchetta, K. Yong, N. C. Firth, D. M. Cash, D. L. Thomas, K. M. Dick, J. Cardoso, J. van Swieten, B. Borroni, D. Galimberti, M. Masellis, M. C. Tartaglia, J. B. Rowe, C. Graff, F. Tagliavini, G. B. Frisoni, R. Laforce, E. Finger, A. de Mendonça, S. Sorbi, J. D. Warren, S. Crutch, N. C. Fox, S. Ourselin, J. M. Schott, J. D. Rohrer, and D. C. Alexander, “Uncovering the heterogeneity and temporal complexity of neurodegenerative diseases with Subtype and Stage Inference,” *Nature Communications*, vol. 9, p. 4273, Oct. 2018.
- [177] A. Kumar, J. Sidhu, F. Lui, and J. W. Tsao, “Alzheimer Disease,” in *StatPearls*, Treasure Island (FL): StatPearls Publishing, 2025.
- [178] M. J. Leming, E. E. Bron, R. Bruffaerts, Y. Ou, J. E. Iglesias, R. L. Gollub, and H. Im, “Challenges of implementing computer-aided diagnostic models for neuroimages in a clinical setting,” *NPJ Digital Medicine*, vol. 6, p. 129, July 2023.
- [179] M. L. Raza, S. T. Hassan, S. Jamil, N. Hyder, K. Batool, S. Walji, and M. K. Abbas, “Advancements in deep learning for early diagnosis of Alzheimer’s disease using multimodal neuroimaging: Challenges and future directions,” *Frontiers in Neuroinformatics*, vol. 19, May 2025.
- [180] S. Golriz Khatami, C. Robinson, C. Birkenbihl, D. Domingo-Fernández, C. T. Hoyt, and M. Hofmann-Apitius, “Challenges of Integrative Disease Modeling in Alzheimer’s Disease,” *Frontiers in Molecular Biosciences*, vol. 6, p. 158, Jan. 2020.

[181] V. W. Henderson and C. E. Finch, “The neurobiology of Alzheimer’s disease,” *Journal of neurosurgery*, vol. 70, no. 3, pp. 335–353, 1989.

Abbreviations and Acronyms

AD	Alzheimer’s Disease
AI	Artificial Intelligence
CAD	Computer-Aided Diagnosis
CNN	Convolutional Neural Network
CT	Computed Tomography
DL	Deep Learning
fMRI	Functional Magnetic Resonance Imaging
Grad-CAM	Gradient-weighted Class Activation Mapping
MCI	Mild Cognitive Impairment
ML	Machine Learning
MRI	Magnetic Resonance Imaging
NC	Normal Control
PET	Positron Emission Tomography
ROI	Region of Interest
VBM	Voxel-Based Morphometry