



University of
Strathclyde
Engineering



Royal Charter
since 1964
Useful Learning
since 1796

Deep Learning for Defect Analysis in Ultrasonic Non-Destructive Evaluation of Aerospace Composites

Shaun McKnight

Department of Electronic and Electrical Engineering
University of Strathclyde

A thesis submitted for the degree of
Doctor of Philosophy
August 2024

**THE QUEEN'S
ANNIVERSARY PRIZES**
1996, 2019, 2021 & 2023
For Higher and Further Education

**UNIVERSITY
OF THE YEAR**
2012 & 2019
Times Higher Education

**UNIVERSITY
OF THE YEAR**
2024 RUNNER-UP
Daily Mail University of the Year Awards

**SCOTTISH UNIVERSITY
OF THE YEAR**
2024
Daily Mail University of the Year Awards

**EUROPEAN ENTREPRENEURIAL
UNIVERSITY OF THE YEAR**
2023
Triple E Awards

Copyright

This thesis is the result of the author's original research. It has been composed by the author and has not been previously submitted for the examination which has led to the award of a degree.

The copyright of this thesis belongs to the author under the terms of the United Kingdom Copyright Acts as qualified by the University of Strathclyde Regulation 3.50. The due acknowledgment must always be made of the use of any material contained in, or derived from, this thesis.

Signed:

A handwritten signature in black ink, appearing to read 'Shan McKinnon', with a stylized flourish at the end.

Date: 8th August 2024

Acknowledgements

In memory of “Pat” Bingham and “Bobby” McKnight.

This work would not have been possible without the support and contributions of many people. I firstly extend my heartfelt thanks to my supervision team: Dr Ehsan Mohseni, Prof. Gareth Pierce, and Prof. Charles MacLeod. Your guidance and mentorship have been invaluable. I was incredibly fortunate to have an experienced supervisory team that not only guided my work but also granted me the creative freedom to explore different ideas.

I am grateful to all members of the SEARCH lab, past and present, who have given their time and helped facilitate this work. I would like to thank my friends and colleagues, particularly Dr. Euan Duernberger and Dr. Euan Foster, for their exceptional technical support and friendship, which helped me settle smoothly into the start of my PhD project. It has also been a privilege to work closely with Vedran Tunukovic, sharing many of the highs and lows of research.

I would like to thank my industrial sponsor Spirit AeroSystems and the Royal Academy of Engineering Research Chair for funding this work. Along with the wider Spirit team, I would specifically like to thank Tom O’Hare for being a key contact for any industrial queries.

Finally, I wish to express my appreciation to my parents, Joan and John, whose love and support have been immeasurable. I am forever grateful for the opportunities you have afforded me.

Abstract

Carbon fibre reinforced polymer (CFRP) composites are critical materials in aerospace components and are seeing ever increasing demand. Ensuring the integrity of these composites is critical and Non-Destructive Evaluation (NDE) is widely used to provide crucial insights about the component. While robotic sensor delivery has automated the physical aspect of sensor placement, the analysis of the resulting data remains a manual and labour-intensive task. This manual analysis is not only time-consuming but also susceptible to human error, thereby limiting manufacturing capabilities. Recent advancements in machine learning (ML) and deep learning (DL) offer potential solutions to automate this analysis process.

This thesis investigates the application of DL methods to analyse Ultrasonic Testing (UT) data from CFRP composites, one of the most common NDE techniques used for aerospace components. The initial phase of research addressed the challenge of data scarcity for training DL models by utilising synthetic datasets. Various methods for generating synthetic UT data were explored, with a deep generate method resulting in a trained classifier with a 24.2% improvement in defect detection accuracy when compared to the same classifier trained on simulated data. This helped to bridge the gap between simulated and experimental training data.

Building on this foundation, the research then focused on the automatic analysis of volumetric UT data. Initially targeting defect detection, the study progressed to volumetric segmentation. For defect detection, a novel architecture was developed, achieving a 22.6% increase in classification accuracy compared to established architectures. A fully supervised model was then employed to train a 3D U-Net for segmentation, which performed well in sizing and localising

defects within the training distribution. However, the model's performance declined when tested on out-of-distribution samples.

The final phase of research explored self-supervised learning to train a model for defect segmentation, reframing the problem as one of outlier prediction. This approach eliminates the need for defective examples during training and significantly enhances generalisability. The method was also tested on an industrial component and scan setup, demonstrating promising performance and applicability in real-world scenarios.

This thesis presents promising methods to automating the analysis of UT data from CFRP composites, highlighting the potential of DL methods to improve accuracy, reduce human error, and enhance manufacturing processes in the aerospace industry.

Contents

Copyright	2
Acknowledgements	3
Abstract	4
i. List of Figures	9
ii. List of Tables	14
iii. Abbreviations	16
1 Introduction:.....	18
1.1 Industrial Motivation and Research Context.....	18
1.2 Aims and Objectives	21
1.3 Outline of Thesis Structure	21
1.4 Contributions to knowledge	22
1.4.1 Lead Author Journal Publications	23
2 Literature and Background:	25
2.1.1 Ultrasonic Testing	26
2.1.2 Composites in Aerospace.....	36
2.1.3 Ultrasonic Testing of Composites	42
2.1.4 Artificial Intelligence	43
2.1.5 Machine Learning in Non-Destructive Evaluation	57
2.1.6 Summary of Key Challenges	64
2.2 Experimental Data Collection	69
2.3 Signal Processing	75
2.4 Dataset Introduction	79
3 Synthetic Data: A Solution to Training Data Scarcity.....	81
3.1 Introduction	81
3.2 Simulated Data	84
3.3 Image Generation	86
3.4 CNN Classification and Evaluation	88
3.4.1 Classification Evaluation with CNNs	88
3.4.2 Hyperparameter Optimisation on Experimental Data.....	90
3.5 Experimental and Simulated Data Classification Performance	93
3.5.1 Experimental Results	93

3.5.2	Simulated Results.....	94
3.6	Methods of Synthetic Noise Generation	94
3.6.1	Approach 1: Simulated to Experimental Domain Mapping with CycleGAN	95
3.6.2	Approach 2: Experimental C-Scan Noise Superposition.....	102
3.6.3	Approach 3: Simulated C-Scan Noise.....	103
3.6.4	Approach 4: Simulated Ultrasonic A-Scan Noise.....	106
3.7	Discussion	111
3.7.1	Comparison of Classification Results	111
3.7.2	Model Interpretability with Guided Grad-CAM.....	116
3.8	Conclusion	118
4	Volumetric Detection.....	121
4.1	Introduction	121
4.1.1	Procedure.....	125
4.2	Data	126
4.2.1	CIVA Simulations.....	126
4.2.2	Synthetic Data Generation	128
4.2.3	Augmentation.....	130
4.3	Network Architectures	133
4.3.1	Evaluation metrics.....	135
4.3.2	VoxNet: Baseline Architecture	135
4.3.3	Hand Designed Architecture	136
4.3.4	NAS Discovered: 3D ResNet based Neural Architecture Search.....	138
Results	144
4.4	Discussion	146
4.5	Conclusion	149
5	Supervised Volumetric Defect Segmentation.....	152
5.1	Introduction	152
5.2	Data	156
5.2.1	Mask Generation	157
5.2.2	Augmentation.....	158
5.3	Segmentation Methods.....	159
5.3.1	Model: Architecture and Training.....	159

5.3.2	Reference Sizing Metric: 6 dB Drop.....	161
5.4	Results and Discussion.....	161
5.4.1	Localisation.....	162
5.4.2	Sizing.....	167
5.5	Out-of-Distribution Testing.....	172
5.6	Conclusion.....	175
6	Self-Supervised Learning Segmentation.....	178
6.1	Introduction.....	178
6.2	Data.....	181
6.3	Pretext Learning.....	183
6.4	Inference.....	191
6.4.1	Methodology.....	191
6.4.2	Results.....	197
6.5	Industrial Demonstration.....	203
6.6	Conclusion.....	208
7	Summary and Future Work.....	213
7.1	Summary.....	213
7.2	Future Work and Final Remarks.....	217
8	Bibliography.....	221
	Appendix.....	242

1	i. List of Figures	
2	Figure 1: Increasing trend of CFRP use in commercial aircraft. Taken from [1].	18
3	Figure 2: Suggested automated data interpretation pipeline.	21
4	Figure 3: Illustration of reflected and transmitted waves for the interaction of an	
5	incident wave normal to the boundary of different mediums.	28
6	Figure 4: Illustration of reflected and transmitted waves for the interaction of an	
7	incident wave at an angle to the boundary of different mediums.	28
8	Figure 5: Example of pulse-echo inspection for a defect response (a) and defect free	
9	response (b).	30
10	Figure 6: Example of through-transmission inspection for a defect response (a) and	
11	defect free response (b).	31
12	Figure 7: Illustrations of standard array scanning methods: (a) Plane wave inspection,	
13	(b) steered inspection, (c) focused inspection.	33
14	Figure 8: Demonstration of how individual elements construct a linear phased array	
15	to produce B-scan and C-scan images when inspecting a component.	34
16	Figure 9: a) Representation of how A-scans are stacked to form B-scans. b) How B-	
17	scans are stacked to create a full UT volume.	35
18	Figure 10: Illustration of the nested relationship between AI, ML and DL.	46
19	Figure 11: Illustration of a perceptron.	48
20	Figure 12: Classical computer vision problems. (a) Image level classification, (b)	
21	Object detection, (c) pixel wise semantic segmentation, (d) instance level semantic	
22	segmentation. [64].	51
23	Figure 13: Illustration of a convolutional filter.	52
24	Figure 14: The structure of a generic CNN, consisting of convolutional, pooling and	
25	fully connected layers [65].	52
26	Figure 15: A breakdown showing the data type and material application of a	
27	selection of ML in NDE literature (54 journal papers).	Error! Bookmark not
28	defined.	
29	Figure 16: Yearly distribution of each data type analysed for a selection of ML in	
30	NDE literature (54 journal papers).	Error! Bookmark not defined.
31	Figure 17: Diagram of roller probe assembly. Adapted from [99].	70

1	Figure 18: Block Diagram illustration of the experimental data acquisition setup....	72
2	Figure 19: a) Overview of the experimental setup of KUKA KR90 and ultrasonic	
3	roller probe used for data acquisition. b) Close-up image of the experimental setup	
4	showing the assembly of the roller-probe and Force-Torque sensor as the robot end	
5	effector.	72
6	Figure 20: Overview of the experimental scan setup of KUKA KR90, Force-Torque	
7	sensor, and ultrasonic roller probe used for data acquisition.	73
8	Figure 21: The composite test sample showing 25 FBHs.....	74
9	Figure 22: a) Example of relative amplitude response from simulations, normalised	
10	signal, and Hilbert transform, applied to the original signal. b) Demonstration of how	
11	individual A-scans are time shifted to the front wall response.	77
12	Figure 23: a) Volumetric data with Hilbert transform applied only. b) Volumetric	
13	data with time shifting to the central response of the front wall peak. Both figures	
14	have been thresholded to remove the lowest 10% of amplitudes to aid in visual	
15	clarity.....	78
16	Figure 24: example of simulated (a) and experimental (b) C-scan responses of a 9	
17	mm diameter FBHs.	88
18	Figure 25: Flow diagram showing the process used for HPO of the CNN architecture	
19	and the use of the optimal architecture for classification evaluation.	89
20	Figure 26: CNN architecture example with a convolutional channel ratio of 2.	91
21	Figure 27: Example images of initial CycleGAN outputs.	96
22	Figure 28: Diagram showing how an example mid-cycle activation map loss is	
23	generated.	98
24	Figure 29: a) The model contains two mapping functions $G_{\text{Experimental}}: \text{Simulated} \rightarrow$	
25	Experimental and $G_{\text{Simulated}}: \text{Experimental} \rightarrow \text{Simulated}$, to transfer between the	
26	respective domains and the associated adversarial discriminators $D_{\text{Experimental}}$ and	
27	$D_{\text{Simulated}}$. b) When completing the full cycle from the simulated domain, the mid-	
28	cycle loss is added along with the cycle loss. c) When completing a cycle beginning	
29	in the experimental domain, the cycle loss is solely used as the mid-cycle loss is not	
30	calculable for the simulated domain.	99
31	Figure 30: Example of synthetic generated images from their corresponding	
32	simulated defect input, along with real experimental images for comparison.....	101

1	Figure 31: Example images showing the combination of real noise and simulated	
2	defect responses.	102
3	Figure 32: Density histogram showing the distribution of data from the clean sample.	
4	104
5	Figure 33: Example images showing the combination of C-scan simulated noise and	
6	simulated defect responses.....	105
7	Figure 34: Density histogram showing the random noise distribution from the total	
8	A-scans.....	108
9	Figure 35: Density histogram showing the distribution of deviation for structural	
10	noise from the mean structural noise pattern.	108
11	Figure 36: a) An example of how a structural noise profile is generated from the	
12	mean. b) A cleaner example of the final generated noise profile.....	109
13	Figure 37: An example of how structural and random noise profiles are combined at	
14	a B-scan level.	110
15	Figure 38: Example images showing the combination of A-scan simulated noise and	
16	simulated defect responses.....	110
17	Figure 39: Comparison of different real and synthetically generated C-scan image	
18	examples.....	112
19	Figure 40: Comparison of classification results for each dataset.....	112
20	Figure 41: Example of Grad-CAM visualisation of models trained on different	
21	datasets.	117
22	Figure 42: Overview of the pipeline for automated volumetric UT classification. .	126
23	Figure 43: a) A frame of 64 simulated A-scans for a simulated defect response. b) the	
24	corresponding A-scans with synthetically added noise for the same defect response.	
25	129
26	Figure 44: a) Complete ultrasonic volume of simulated A-scans for a defect response.	
27	b) the corresponding synthetically noised volume for the same defective response.	
28	Both figures have been thresholded to remove the lowest 10% of amplitudes to aid in	
29	visual clarity.....	130
30	Figure 45: a) Example of how scaling augmentation is done on an individual A-scan.	
31	b) Example of how dilation augmentation and padding is completed for an individual	
32	A-scan.	132

1	Figure 46: The VoxNet architecture. Where Conv (f,d,s) indicates the number of	
2	filters f, filter size d, and stride s, of the convolutional layer.....	136
3	Figure 47: Network architecture for the CustomNet.	138
4	Figure 48: Representation of the ResNet style searched space.....	139
5	Figure 49: Diagram of the searched residual block.	139
6	Figure 50: Overview of the process for NAS implementation.	141
7	Figure 51: The overall structure of the discovered architecture.	142
8	Figure 52: The details of each discovered residual block.....	143
9	Figure 53: Demonstration of the 6 dB drop method for defect sizing. a) Finding	
10	maximum defect response. b) Using the 6 dB loss in maximum amplitude to locate	
11	one edge of the defect. c) The corresponding defect edge detected using the 6 dB	
12	drop to determine the defect length.....	154
13	Figure 54: (a) The ground truth segmentation mask and (b) the corresponding	
14	simulated defect response. The overlay of both the mask and response is shown in	
15	(c). Colour mapping and axes are given in Figure 23.	157
16	Figure 55: Architecture diagram for the 3D U-Net. Blue boxes depict the feature	
17	maps. The number of channels and dimensions of the data (probe × scan × time) are	
18	denoted above and to the side of each feature map respectively.	160
19	Figure 56: Example 9 and 3 mm defects respectively; their experimental ultrasonic	
20	volumetric responses, thresholded responses (amplitudes >10% of maximum	
21	response for defect visualisation), and their corresponding predicted segmentations.	
22	Colour mapping and axes are given in Figure 23.	162
23	Figure 57: Depth localisation results.	164
24	Figure 58: In plane localisation results compared to 6 dB drop with reference to	
25	defect diameters (a) and the expanded (b), which shows the reference to the array	
26	pitch more clearly.....	165
27	Figure 59: Sizing results for the 6 dB drop method and U-Net predictions.	167
28	Figure 60: Comparison of sizing for synthetic and experimental data for all defect	
29	sizes (a) and defect diameters above 4 mm (b).	170
30	Figure 61: Corrected sizing results for the U-Net predictions.	172
31	Figure 62: Diagram illustrating an example of a through-component series, as	
32	examined for the pretext learning task.	184

1	Figure 63: Example of predicted distribution from an input series.	186
2	Figure 64: Probabilistic CNN architecture.....	187
3	Figure 65: Demonstration of the impact of stride (8 and 64) when sampling the	
4	training data.....	189
5	Figure 66: Test set mean Log-Likelihood for varying sampling strides.....	190
6	Figure 67: Flowchart of sequence prediction for inference.	192
7	Figure 68: Flowchart of the methodology overview for complete volumetric	
8	segmentation.	195
9	Figure 69: Demonstration of the impacts of different post-processing steps for defect	
10	sample 2 and threshold of 0.9999999.	195
11	Figure 70: Example of edge artefacts for a single sweep when using zero padding	
12	compared to edge padding (sample: defect 2, threshold: 0.9999999).....	196
13	Figure 71: Defect detection accuracy for each threshold and processing step.	198
14	Figure 72: Example B-scan across multiple raster passes showing the voxels	
15	highlighted as defective and the corresponding true defect size in white. This is	
16	shown for defect sample 1 with a threshold of 0.9999999.	200
17	Figure 73: Visualisations of volumetric ultrasonic responses, their corresponding	
18	overlayed segmentations, and component drawing.	203
19	Figure 74: Impact on data processing for SSL method on Learjet data using the	
20	SEARCH lab acquisition system.	206
21	Figure 75: Amplitude C-scan of the Learjet component section acquired using the	
22	Tecnatom setup. Both manual defect annotation and SSL defect segmentations are	
23	shown.	207
24	Figure 76: Screenshot of the graphical user interface developed to interact with DL	
25	model outputs from UT scans (left) and demonstration of the system integrated into a	
26	flexible robotic scanning system (right).....	220
27		

1 **ii. List of Tables**

2 Table 1: Appropriate UT imaging method for a selection of different composite
3 defects, where the number of tics corresponds to the increasing level of applicability
4 from none to high. Modified from [29]..... 36
5 Table 2: Summary of roller probe parameters 70
6 Table 3: Summary of samples and their defects. 74
7 Table 4: Summary of the datasets produced. 87
8 Table 5: HPO variables and their range of values..... 92
9 Table 6: Optimised hyperparameters used for CNN..... 93
10 Table 7: Average confusion matrix across 100 training iterations for a CNN trained
11 on experimental data. 93
12 Table 8: Average confusion matrix across 100 training iterations for a CNN trained
13 on simulated data..... 94
14 Table 9: Average confusion matrix across 100 training iterations for a CNN trained
15 on GAN generated synthetic data. 102
16 Table 10: Average confusion matrix across 100 training iterations for a CNN trained
17 on real noise data..... 103
18 Table 11: Average confusion matrix across 100 training iterations for CNN trained
19 on simulated C-scan noise data. 106
20 Table 12: Average confusion matrix across 100 training iterations for a CNN trained
21 on simulated A-scan noise data..... 111
22 Table 13: Summary of classification results for each dataset. 113
23 Table 14: Summary of the datasets produced. 130
24 Table 15: Average confusion matrices for VoxNet, CustomNet and the NAS
25 discovered architecture..... 144
26 Table 16: Comparison of classification results across the different architectures. The
27 means and standard deviations are presented as mean \pm std..... 144
28 Table 17: Comparison of the effects of data augmentation on the NAS discovered
29 architecture. The means and standard deviations are presented as mean \pm std..... 145
30 Table 18: Comparison of model sizes and inference time for each architecture. 145
31 Table 19: Summary of the datasets produced. 157
32 Table 20: Complete model centroid deviation results from the 6 dB drop..... 166

1	Table 21: Comparison in sizing of experimental and synthetic responses with the 6	
2	dB drop method. Where R is the ratio of synthetic response to experimental response.	
3	R gives the mean average of R for a given defect diameter.....	171
4	Table 22: Sizing and in-plane localisation results for out of distribution test defects.	
5	173
6	Table 23: Summary of samples used.	182
7	Table 24: MAE for different thresholds.....	200
8	Table 25: Localisation results.	202
9	Table 26: Summary of defects scanned in the Learjet 85 component.	204
10	Table 27: Comparison of correct defect detections and false indications for the	
11	different pre-processing and acquisition systems.	207
12	Table 28: Detection accuracy across thresholds and processing steps for each sample.	
13	242
14		

1 iii. Abbreviations

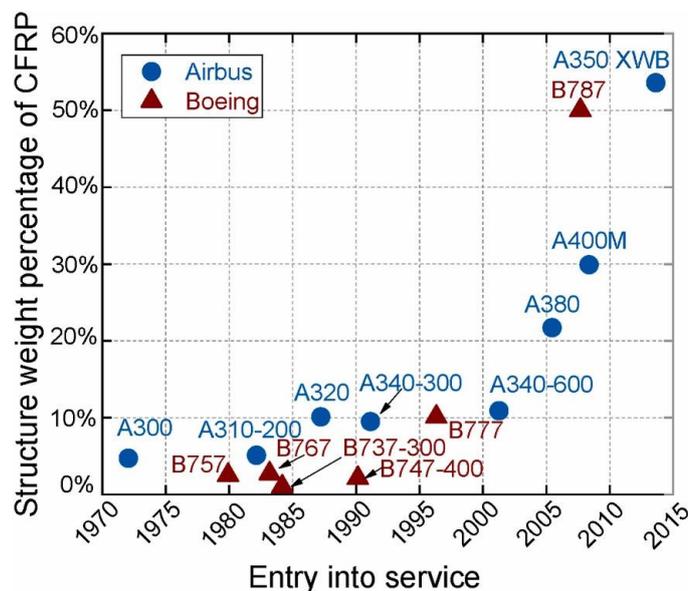
1D	-	1 Dimensional
2D	-	2 Dimensional
3D	-	3 Dimensional
AI	-	Artificial Intelligence
CAM	-	Class Activation Map
CFRP	-	Carbon Fiber Reinforced Polymer
CNN	-	Convolutional Neural Network
CV	-	Computer Vision
DL	-	Deep Learning
FBH	-	Flat-Bottom Holes
FEA	-	Finite Element Analysis
FMC	-	Full Matrix Capture
FOI	-	Foreign Object Inclusion
GAN	-	Generative Adversarial Network
GELU	-	Gaussian Error Linear Units
HPO	-	Hyperparameter Optimisation
MAE	-	Mean Absolute Error
ML	-	Machine Learning
MPL	-	Multi-Layer Perceptron
NAS	-	Neural Architecture search
NDE/T	-	Non-Destructive Evaluation/Testing
NLL	-	Negative Log-Likelihood
PNN	-	Probabilistic Neural Network
PTFE	-	Polytetrafluoroethylene
RE	-	Regularised Evolution
ReLU	-	Rectified Linear Unit
RNN	-	Recurrent Neural Networks
SHM	-	Structural Health Monitoring
SOTA	-	State of the Art
SSL	-	Self-Supervised Learning

TCG	-	Time-Compensate Gain
TFM	-	Total Focusing Method
ToF	-	Time of Flight
UT	-	Ultrasonic Testing
ViT	-	Vision Transformer
wt%	-	Percentage by Weight

1 Introduction:

2 1.1 Industrial Motivation and Research Context

3 Composite materials find extensive application within the aerospace, marine and civil
4 engineering sectors, with Carbon Fibre Reinforced Polymers (CFRP) being one of the
5 most prominent. The global demand for CFRP is expected to reach 285 kt in 2025,
6 rising from 181 kt in 2021 [1]. Figure 1 demonstrates the strong increasing trend of
7 CFRP use in commercial aircraft, with improved performance and efficiency in
8 operation driven by environmental pressures. Due to their high specific strength,
9 stiffness and corrosion resistance, composites are widely used for critical aerospace
10 components such as wings and fuselages [2], [3]. As the use of CFRP grows, the need
11 for effective testing of these safety critical components also increases. Ultrasonic Non-
12 Destructive Evaluation (NDE) is the most applied method for the inspection of
13 aerospace composites during manufacturing [4], [5]. Ultrasonic inspection is often a
14 manual task which can be time-intensive, challenging to scale, and exposed to human
15 factors (such as cognitive, physical or experience) which can lead to error [6].



16

17 *Figure 1: Increasing trend of CFRP use in commercial aircraft. Taken from [1].*

1 The integration of robotics into NDE has enabled the efficient automation of sensor
2 deployment for large-scale inspection processes [7]. However, despite the increased
3 flexibility of robotic scanning and the drastic reduction in scan time seen by
4 mechanised scanning compared to manual scanning (high degree of freedom robotics
5 are three times faster than gantry systems [8]), the interpretation of the results in
6 industry remains a challenging and time intensive task that requires highly trained and
7 qualified operators to interpret results according to existing standards [9], [10], [11],
8 [12], [13], [14]. Despite the significant improvements brought about by robotic NDE,
9 the need for expert human interpretation of results persists. This highlights the need
10 for further research and development of automated data interpretation techniques that
11 can supplement or even replace human interpretation, to improve the efficiency and
12 reliability of NDE in various industries. By reducing the dependence on human
13 interpretation, automation can enhance the consistency, repeatability, and traceability
14 of the NDE processes, while reducing inspection time and costs.

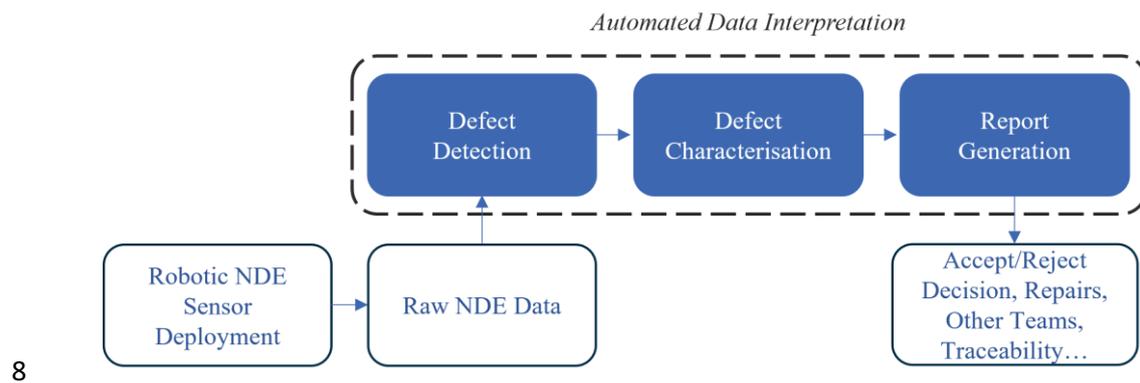
15 The interpretation of UT scan results by human operators presents two significant
16 drawbacks, namely, poor time efficiency and the risk of human error [12]. Low levels
17 of automation for data interpretation are feasible for mass-produced parts with
18 precisely known geometries, but this approach typically relies on hard-coded features
19 such as predefined time-gating, filtering, and amplitude thresholding, which may not
20 be adequate for complex tasks, with changes in manufacturing conditions, variations
21 in geometry, or defect characteristics [14]. DL has been identified as a key requirement
22 for transitioning from low to high levels of industrial automation [14]. Therefore, if a
23 Deep Learning (DL) approach could be created to automate the interpretation of
24 complex results and work alongside the robotic inspection, the required inspection

1 time and quality of large components could be improved significantly, allowing for
2 shorter signal interpretation time and a faster uptake of UT automation in aerospace
3 and other industries. Furthermore, increased automation could lead to better defect
4 detection capability, whilst improving consistency, traceability, and repeatability as
5 discussed by Cantero et al. [14].

6 DL approaches have been identified as one of the most important enabling factors for
7 the transition from low to high automation levels [14]. Despite the clear potential
8 benefits of applying DL techniques to ultrasonic signal analysis for composite
9 components, its uptake has been limited [14]. Shortage of training data is one of the
10 main challenges that hinders research developments in this area. This shortage,
11 combined with industrial concerns about the interpretability and compliance with
12 standards of DL models, has presented challenges for the effective use of DL
13 techniques. As a result, the adoption of DL in UT signal analysis for composite
14 components has been slow, despite its promising potential to enhance the accuracy and
15 efficiency of defect detection and characterisation.

1 1.2 Aims and Objectives

2 The objective of this thesis is to investigate the application of DL in the ultrasonic
3 inspection of carbon fibre composites, with the aim of enhancing the automation levels
4 observed in industrial aerospace ultrasonic testing. Considering typical tasks required
5 by NDE operators, it is conceivable to envision an automated data interpretation
6 pipeline, as depicted in Figure 2. The work presented in this thesis aims to breakdown
7 and tackle different stages of the proposed pipeline.



9 *Figure 2: Suggested automated data interpretation pipeline.*

10 1.3 Outline of Thesis Structure

11 The following chapters of the thesis are structured as follows:

- 12 • Chapter 2: Provides an overview of ultrasonic testing, composites, AI, along
13 with the experimental data collection, signal processing and samples used in
14 the following chapters.
- 15 • Chapter 3: Presents a solution to a lack of training data in NDE through the
16 exploration of different synthetic data generation methods.
- 17 • Chapter 4: Presents work on volumetric defect detection, which utilises and
18 expands on the work from the previous chapter to work with volumetric
19 datasets.

- 1 • Chapter 5: Explores a supervised deep learning approach for defect
2 segmentation in volumetric data, trying to tackle defect characterisation and
3 report generation.
- 4 • Chapter 6: Demonstrates a self-supervised method for volumetric defect
5 segmentation, including a case study showing industrial application.
- 6 • Chapter 7: Conclusion of the main findings of the thesis and discusses potential
7 future work.

8 1.4 Contributions to knowledge

9 This body of work investigates DL solutions for the automated analysis of UT data for
10 CFRP aerospace components. This thesis presents several unique and novel
11 contributions to the field of ML in NDE. They are summarised below.

- 12 • Providing a solution to the lack of training data for NDE ML models by
13 bridging the simulation to experimental domain gap with the development of
14 four different synthetic data generation techniques. The evaluation and
15 comparison of simulated and synthetic training data on experimental
16 classification performance. Additionally, modifications to CycleGAN,
17 including the introduction of an additional loss term and weighting the
18 synthetic data generator more than the simulated generator, were implemented
19 to improve the quality of synthetic data generation.
- 20 • The introduction of volumetric analysis for defect detection from UT data
21 using DL. A new architecture was discovered which greatly improved on
22 existing models. This was evaluated against an established VoxNet and hand-
23 crafted architecture. Domain specific augmentation methods were also

1 developed and evaluated to increase model performance.

- 2 • A 3D U-Net variant was employed for volumetric segmentation to achieve
3 accurate defect sizing and localization. The model was trained using supervised
4 learning, with synthetic training data and defect simulation parameters
5 providing the ground truth segmentation masks necessary for training.
- 6 • The development of a self-supervised model for volumetric segmentation of
7 UT data. The model learns from defect-free data series along the scan direction,
8 allowing it to establish a baseline representation. This approach enables
9 segmentation through anomaly detection, eliminating the need for labelled
10 defective training data.

11 1.4.1 Lead Author Journal Publications

- 12 1. **S. McKnight**, G. Pierce, E. Mohseni, C. MacKinnon, C. N. MacLeod, T.
13 O'Hare, and C. Loukas, '*A Comparison of Methods for Generating Synthetic*
14 *Training Data for Domain Adaption of Deep Learning Models in Ultrasonic*
15 *Non-Destructive Evaluation*'. NDT & E International vol 141, no. 102978, doi:
16 10.1016/j.ndteint.2023.102978.
- 17 2. **S. McKnight**, C. MacKinnon, G. Pierce, E. Mohseni, V. Tunukovic, C. N.
18 MacLeod, R. Vithanage, and T. O'Hare. '*3-Dimensional Residual Neural*
19 *Architecture Search for Ultrasonic Defect Detection*', IEEE Transactions on
20 Ultrasonics, Ferroelectrics, and Frequency Control, doi:
21 /10.1109/TUFFC.2024.3353408.
- 22 3. **S. McKnight**, V. Tunukovic, G. Pierce, E. Mohseni, R. Pyle, C. N. MacLeod,
23 and T. O'Hare. '*Advancing Carbon Fiber Composite Inspection: Deep*
24 *Learning-Enabled Defect Localization and Sizing via 3-Dimensional U-Net*

- 1 *Segmentation of Ultrasonic Data*'. IEEE Transactions on Ultrasonics,
2 Ferroelectrics, and Frequency Control, doi: 10.1109/TUFFC.2024.3408314.
3 4. **S. McKnight**, V. Tunukovic, A. Hifi, G. Pierce, E. Mohseni, C. N. MacLeod
4 and T. O'Hare, '*3-DUSSS: 3-Dimensional Ultrasonic Self Supervised*
5 *Segmentation*' [Under Review].

6 *For full research output, please see: [Google Scholar Profile](#).*

7

1 2 Literature and Background:

2 NDE constitutes a diverse array of methodologies employed for the inspection of
3 components without inducing damage. Prominent techniques utilised in this domain
4 are Radiography, Thermography, Electromagnetic methods, and Ultrasound. These
5 methods provide inspection capabilities suitable for components of varying
6 complexities and sizes. Each of these methods exhibits distinctive strengths and
7 weaknesses, necessitating a detailed decision-making process when selecting a
8 suitable approach, which often involves compromise to find the most applicable
9 method. Factors influencing this decision encompass a multitude of considerations,
10 including material properties, component geometry, safety considerations, resolution,
11 implementation feasibility, operational constraints, and defect typology, among others.
12 The intricate nature of this decision-making process underscores the complexity
13 inherent in NDE method selection. The application of appropriate NDE techniques can
14 significantly enhance the reliability and safety of structures and components across
15 diverse industrial sectors.

16 The work conducted in this thesis focuses on the use of Ultrasonic testing. For an in-
17 depth exploration of the various NDE methods pertinent to composite materials,
18 readers are directed to the comprehensive investigation conducted by S. Gholizadeth
19 [15]. This work provides valuable insights into the application of a diverse range of
20 NDE techniques tailored specifically to the assessment of composite materials,
21 providing invaluable insights into their respective merits and limitations.

1 2.1.1 Ultrasonic Testing

2 UT employs high-frequency acoustic waves, typically exceeding 20 kHz, to assess the
3 integrity of components. This technique offers versatility, allowing inspection of
4 materials ranging from metals to composites, and relies on the transmission,
5 propagation, and reception of ultrasonic waves. UT stands out as one of the most
6 prevalent NDE method owing to its capability for volumetric inspection, ability to
7 detect a diverse range of defects, along with its adaptability, user-friendliness, and
8 safety profile. Various ultrasonic-based inspection methods exist, tailored to different
9 applications such as utilising bulk waves for volumetric detection [16] or guided waves
10 for increased inspection ranges [17]. This thesis concentrates on bulk ultrasonic waves,
11 primarily aimed at sub-surface defect detection.

12 UT has gained extensive adoption and standardisation for volumetric inspection within
13 the aerospace industry, primarily due to its comparatively straightforward and safe
14 implementation in contrast with radiography, alongside its capacity to identify a broad
15 spectrum of volumetric defects [9], [15], [18], [19]. In the aerospace sector, composite
16 UT predominantly utilises bulk wave inspection, where sound waves are excited on
17 the surface of a component, and the reflected/scattered wave from internal scatterers
18 can provide valuable information about the volumetric discontinuities or properties of
19 the component.

20 *2.1.1.1 Bulk Wave Propagation in Isotropic Medium*

21 As waves propagate through real materials they lose energy, this effect is known as
22 attenuation. Attenuation of a wave is due to various mechanisms such as absorption or
23 scattering, but generally waves with longer wavelengths have lower attenuation than

1 waves with short wavelengths, allowing them to propagate further into a component.
2 However, higher frequency waves with shorter wavelengths interact with smaller
3 features allowing for higher resolution inspection. There is therefore a balance between
4 penetration depth and inspection resolution. This balance is constrained by the
5 relationship between operating frequency and wavelengths, as given in equation (1),
6 where f is the operating frequency and λ is the wavelength:

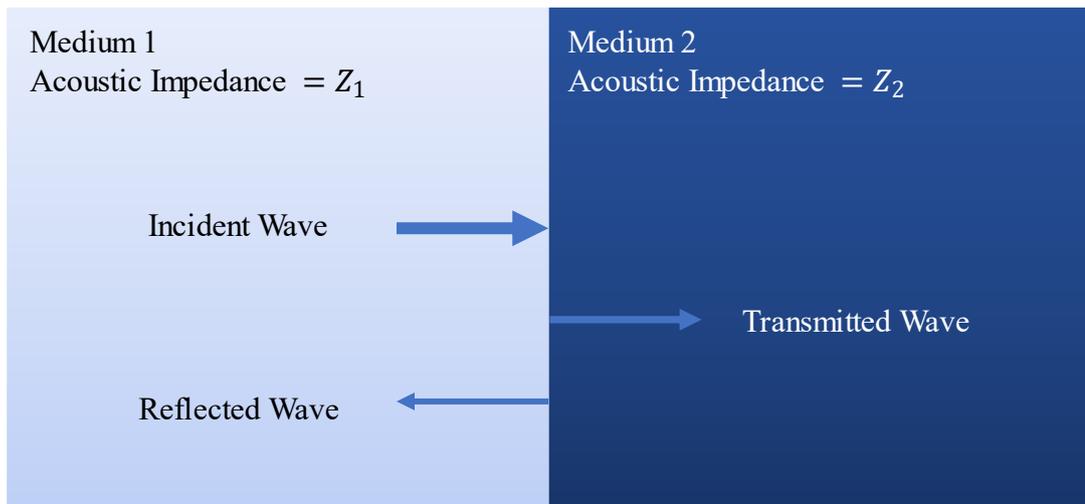
$$f = \frac{v}{\lambda} \quad (1)$$

7 *2.1.1.2 Reflection Refraction and Mode Conversion*

8 When an ultrasonic wave meets the boundary of two media with different acoustic
9 impedances it will undergo changes such as reflection, refraction, or mode conversion.
10 When a wave encounters a boundary perpendicular to the interface surface, typically
11 a portion of the wave is transmitted into the new medium, while another portion is
12 reflected (Figure 3). This is given by T and R respectively in equations (2) and (3).
13 Here, Z denotes the acoustic impedance of a given material. Acoustic impedance refers
14 to a materials resistance to the propagation of ultrasound. This is given by equation
15 (4), where ρ is the material density. Consequently, when interacting with materials
16 exhibiting a substantial disparity in acoustic impedance, the wave will undergo a
17 greater degree of reflection.

$$R = \frac{Z_2 - Z_1}{Z_1 + Z_2} \quad (2)$$

$$T = \frac{2Z_2}{Z_1 + Z_2} \quad (3)$$

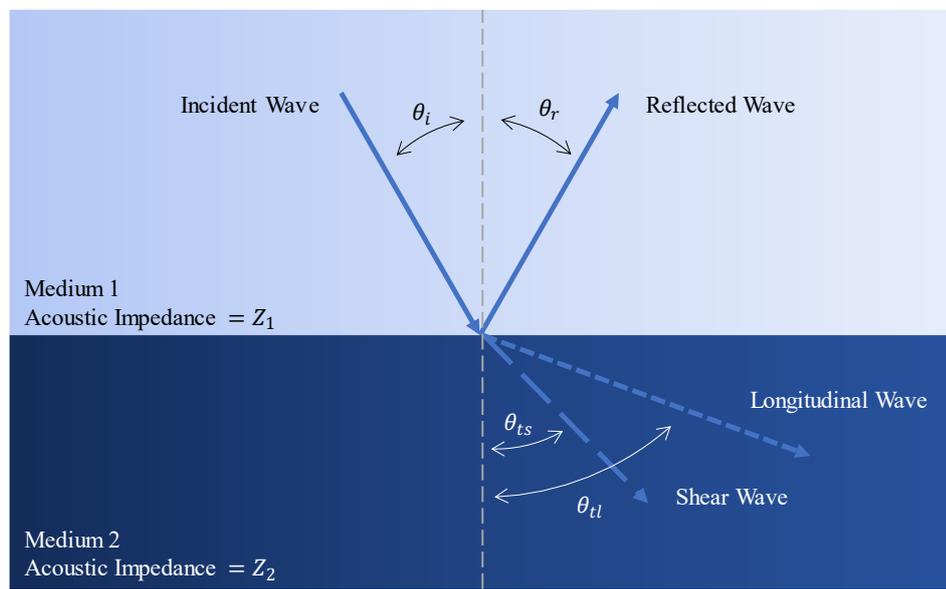


1

2 *Figure 3: Illustration of reflected and transmitted waves for the interaction of an incident wave*
 3 *normal to the boundary of different mediums.*

$$Z = \rho v \tag{4}$$

4 When the acoustic wave does not meet the interface boundary at a normal angle, a
 5 change of direction is observed (Figure 4). This is known as refraction and is defined
 6 by Snells law, given in equation (5).



7

8 *Figure 4: Illustration of reflected and transmitted waves for the interaction of an incident wave at an*
 9 *angle to the boundary of different mediums.*

10

1

$$\frac{\sin \theta_i}{v_i} = \frac{\sin \theta_{rl}}{v_{rl}} = \frac{\sin \theta_{ts}}{v_{ts}} = \frac{\sin \theta_{tl}}{v_{tl}} \quad (5)$$

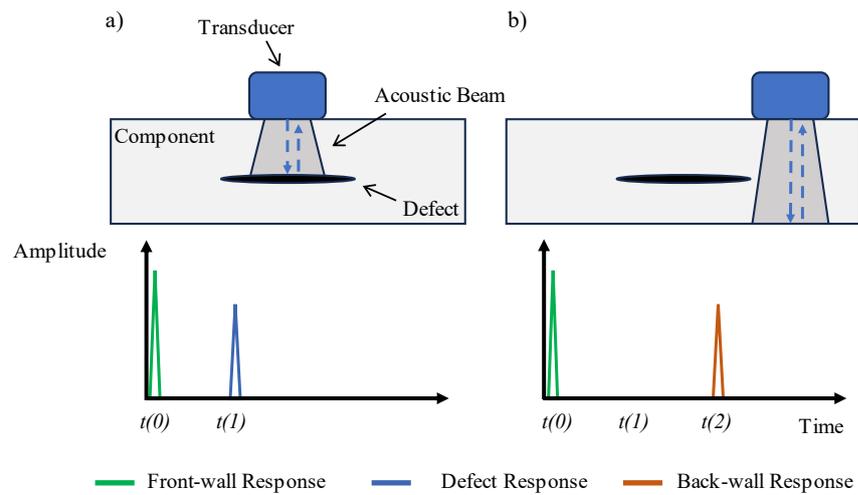
2 When the refracted angle $\theta_{tl} < 90^\circ$ (equation (5), Figure 4), known as the critical
3 angle, mode conversion occurs. Mode conversion is the propagation of a different
4 wave type e.g. longitudinal waves create shear waves. At the critical angle ($\theta_{tl} =$
5 90°), the incident longitudinal wave is converted to a surface following longitudinal
6 wave. Beyond the first critical angle ($\theta_{tl} > 90^\circ$) the longitudinal wave is totally
7 internally reflected, meaning it does not pass into the second medium but is instead
8 reflected back into the first medium, and only the shear wave is refracted into the
9 material.

10 While numerous other wave types and factors influence ultrasound, this section offers
11 an overview of the primary macro components contributing to bulk waves, which
12 constitute the majority of ultrasonic inspection [20].

13 *2.1.1.3 Single Element Ultrasound*

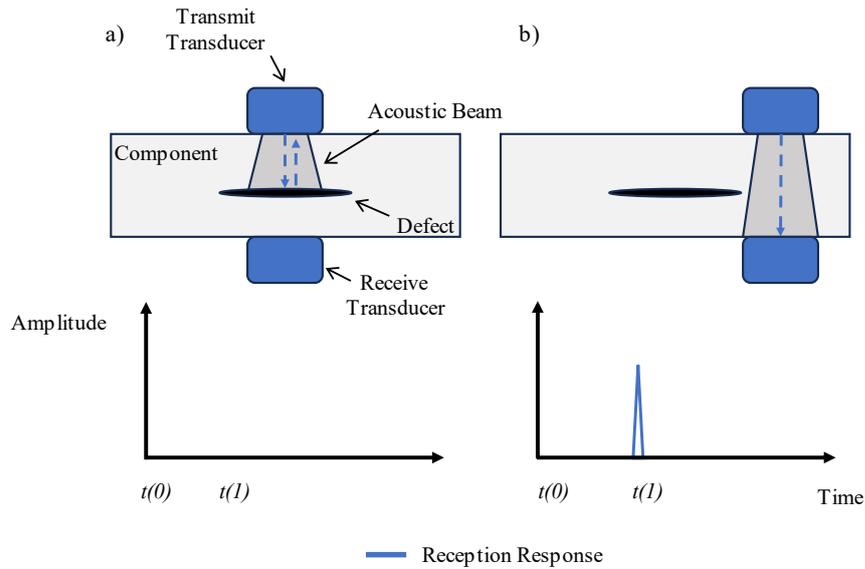
14 A single piezoelectric ultrasonic transducer is the simplest means to convert electrical
15 current to ultrasonic waves and vice versa. A single transducer can be used to inspect
16 the internal volume of a component along with other tests such as thickness
17 measurements using pulse-echo. In pulse-echo, the transducer is attached to the surface
18 of a component, often with a coupling medium (or couplant), which aids in removing
19 the air surrounding the components' surface, facilitating acoustic energy transfer from
20 the transducer to the test component. An acoustic wave is then propagated through the
21 surface of the component, traversing the thickness and reflecting off the backwall
22 surface in a healthy component (this can be used for thickness measurements if the

1 speed of sound in the component is known) or off a discontinuity in a defective
2 component (Figure 5). The reflection at the surface of the component and any later
3 reflections are recorded using the transducer which converts the pressure into a
4 voltage. This is often visualised in the form of an amplitude versus time (A-scan) plot.
5 Pulse-echo inspection sees wide applicability to industrial settings where access is
6 limited to a single side of a component.



7
8 *Figure 5: Example of pulse-echo inspection for a defect response (a) and defect free response (b).*

9 Through-transmission is a common alternative to pulse-echo which makes use of two
10 transducers, one for transmission and one for reception, on either side of the
11 component (Figure 6). However, this method requires access to both sides of the
12 component and requires good alignment of the two probes. This can make it more
13 challenging to implement in industrial settings than the pulse-echo method.



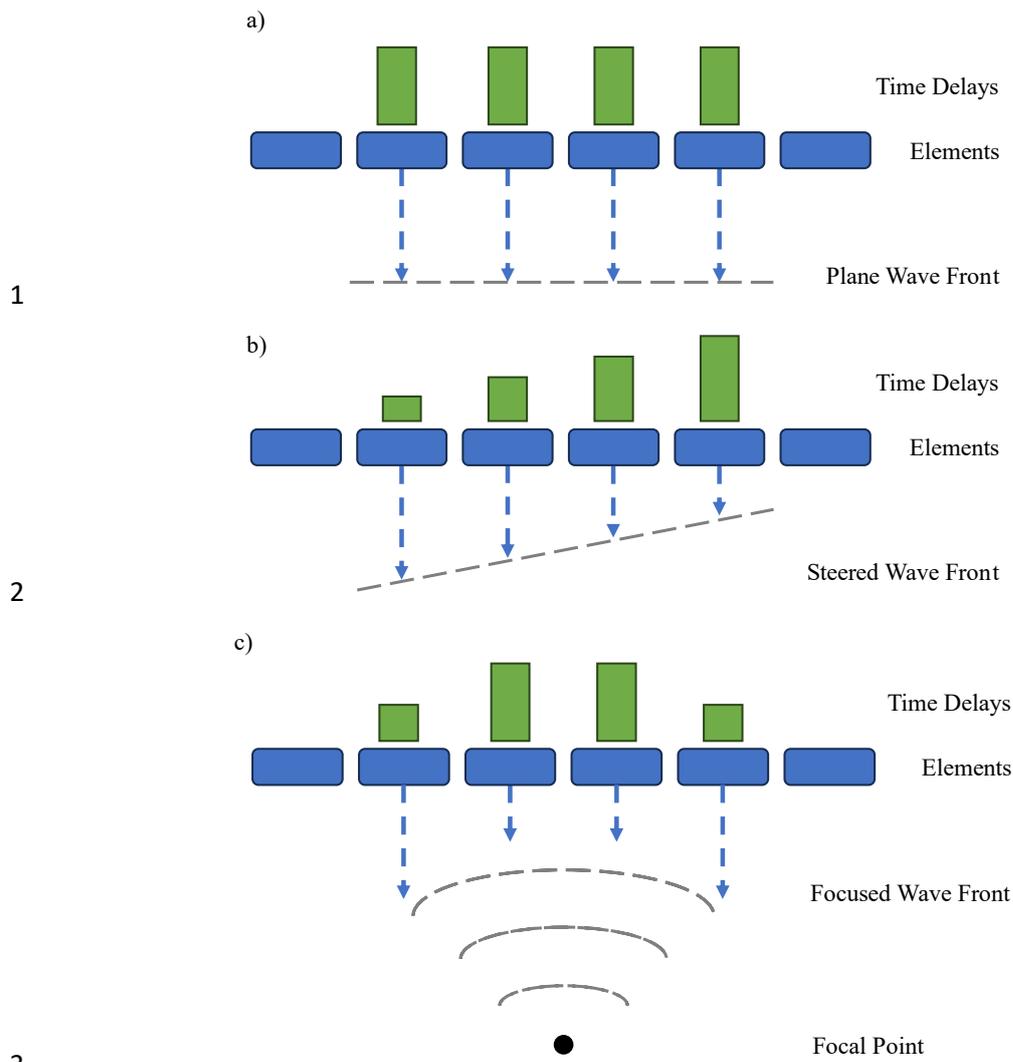
1

2 *Figure 6: Example of through-transmission inspection for a defect response (a) and defect free*
 3 *response (b).*

4 *2.1.1.4 Phased Arrays, Mechanised Scanning and Data Representation*

5 The use of Phased Array Ultrasonic Testing (PAUT) has become increasingly popular
 6 to generate and receive ultrasonic sound waves owing to their operational flexibility.
 7 Phased arrays employ independently controllable UT elements that enable the
 8 collection of richer datasets through more complex electronic scanning and imaging
 9 methods, such as beam steering, dynamic depth focusing, and variable sub-apertures
 10 [21]. Examples of how firing delays can generate different beams can be seen in Figure
 11 7. Recently Full Matrix Capture (FMC) has become increasingly popular as it collects
 12 the full set of transmit/receive element combinations which allows for different
 13 imaging methods to be applied post-acquisition such as the Total Focusing Method
 14 (TFM) [22]. When constructing a TFM image a delay and sum approach is used, where
 15 the known speed of sound in the material is used to give the expected time index of
 16 each FMC transmission/reception pair. For every pixel in the image the sum of each

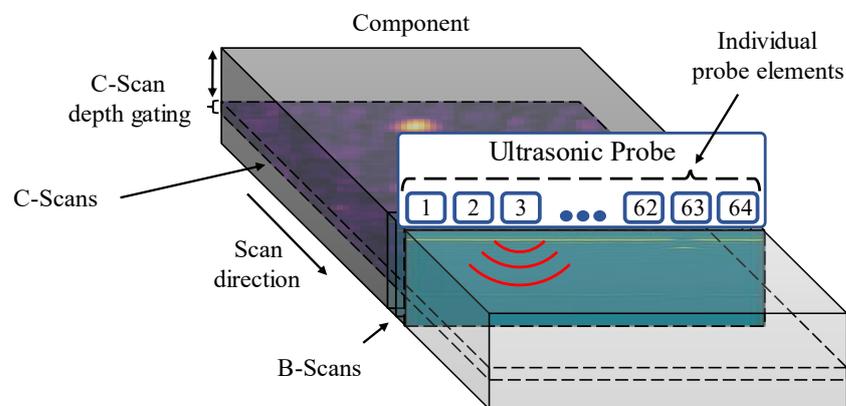
1 pairs corresponding time indexed amplitude gives the total magnitude of the response.
2 This leads to an increased signal to noise ratio and improved focusing within the
3 image. Such imaging methods work very well for isotropic materials with consistent
4 and well-defined acoustic velocities. However, for anisotropic medium with variable
5 speeds of sound, the calculation of appropriate time of flights becomes more
6 challenging and can be intractable for certain scenarios without the use of simulations
7 [23], making these advanced imaging methods often inappropriate for industrial
8 inspections, where efficiency is a key driver.



4 *Figure 7: Illustrations of standard array scanning methods: (a) Plane wave inspection, (b) steered*
5 *inspection, (c) focused inspection.*

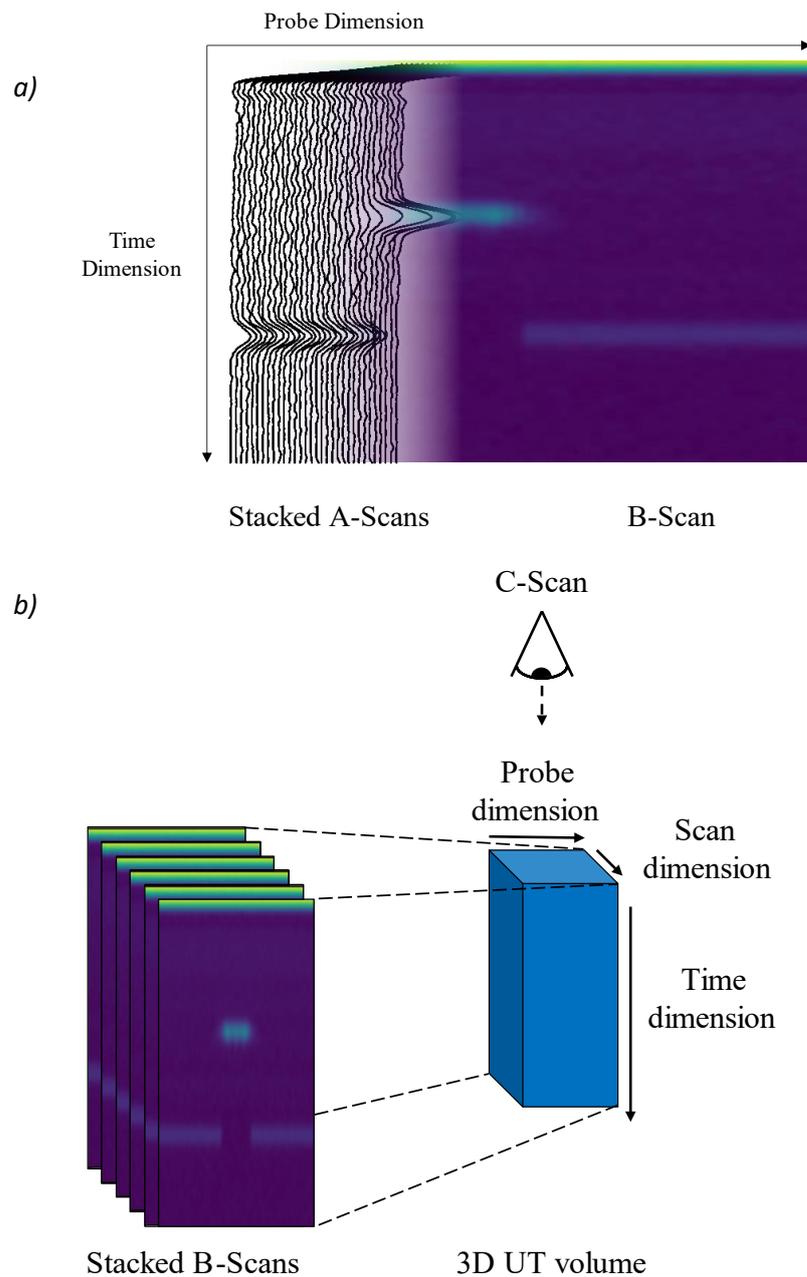
6 By controlling each individual element (or sub-aperture of elements) of a linear phased
7 array, depth-wise sectional images (B-scans) can be created in a single scan (Figure 8,
8 Figure 9). When combined with mechanised scanning perpendicular to the length of a
9 linear phased array, complete 3-dimensional (3D) volumetric scan data of components
10 can be generated by stacking multiple individual B-scans together at known positions
11 (Figure 9 (b)). Mechanised scanning was first introduced with gantry or bridge-based
12 systems which aimed to automate sensor delivery for simple or fixed geometries[24].

1 The adoption of more complex high axis industrial robots with greater degrees of
2 freedom has led to increased flexibility [8], [24], [25], [26], with previous systems
3 demonstrating impressive scanning rates of 25.3 square meters per hour [8]. The
4 integration of robotics and PAUT into NDE has revolutionised large-scale inspection
5 processes by enabling efficient automated inspection of large, complex components
6 [7]. It allows for increased flexibility, repeatability and the drastic reduction in scan
7 time seen compared to previous manual scanning approaches and has significant
8 positive implications for the reliability and safety of aerospace structures.



9

10 *Figure 8: Demonstration of how individual elements construct a linear phased array to produce B-*
11 *scan and C-scan images when inspecting a component.*



1 *Figure 9: a) Representation of how A-scans are stacked to form B-scans. b) How B-scans are stacked*
 2 *to create a full UT volume.*

3 UT data is commonly visualised as images, either by selecting a B-scan directly, or as
 4 an amplitude or time of flight C-scan; where either the maximum response amplitude
 5 (examples can be seen in Figure 22) or the time index of the maximum response
 6 amplitude within the volume is imaged to produce a top-down section view across the

1 sample. B-scan imaging offers the most comprehensive understanding of the
 2 volumetric response of a sample, yet it sacrifices comparative spatial information
 3 along the scan axis, which can make it impractical for defect sizing and inefficient for
 4 the complete analysis of large components. Conversely, C-scans maintain comparative
 5 spatial information but at the expense of compressing the temporal dimension and
 6 losing depth-wise information. As highlighted by Smith et al. [27], who extracted
 7 additional information such as ply orientation for composite inspection from C-scans,
 8 which is not readily available from B-scans. Typically, practitioners employ both
 9 imaging methods concurrently to evaluate components effectively and efficiently, as
 10 certain imaging methods are appropriate for different defects, as given in Table 1.

11 *Table 1: Appropriate UT imaging method for a selection of different composite defects, with*
 12 *applicability levels ranked in ascending order as: none, low, medium, high, v. (very) high. Modified*
 13 *from [28].*

Defect	A-Scan	B-Scan	Amplitude C-Scan	Time-of- Flight C-Scan
Delamination (<10 mm)	Medium	High	High	High
Delamination (>10 mm)	High	V. High	V. High	V. High
Disbond	High	V. High	V. High	V. High
Void	Medium	Medium	V. High	Medium
Impact	V. High	V. High	V. High	V. High
Porosity	Medium	Medium	High	Medium
Inclusion	High	High	High	Medium
Fiber Waviness	Low	High	High	None
Fiber misalignment	None	Low	High	None

14

15 2.1.2 Composites in Aerospace

16 Composites are versatile, often anisotropic materials that are widely used in many
 17 industries due to their favourable properties such as corrosion resistance, high specific
 18 strength, and high specific stiffness. Layered composite structures are generally
 19 anisotropic. This anisotropy allows for precise engineering to meet specific structural

1 requirements, making these composites ideal for high-performance applications [9],
2 [18], [19], [21], [29], [30], [31], [32], [33]. CFRP is a widely used layered composite
3 in the aerospace industry making up over 50 percent by weight (wt%) for the two most
4 recent long-range aircraft, the Airbus A350 and the Boeing 787, and up to 70-80 wt%
5 for private jets and helicopters [1]. CFRP is manufactured by layering multiple carbon
6 ply sheets, which are subsequently cured following the addition of a thermoset
7 polymer.

8 It is necessary for enhanced automation in composite manufacturing to improve
9 efficiency and mitigate costs. However, the manufacturing process of composite
10 components is complex which can introduce defects, compromising their integrity and
11 performance [15], [18], [21], [29], [31], [33], [34]. Given that rejected components
12 result in heightened expenses due to rework and repairs, automated testing also
13 emerges as a pivotal requirement. Defects can range from delamination and cracks to
14 foreign object inclusions, fibre distortions, and porosity [19], [34] (for further details
15 please refer to section 2.1.2.1). These defects represent an increased risk as they are
16 often not detectable by the naked eye and can affect strength and fatigue behaviour
17 [33]. Cyclic stresses from operation can cause these defects to grow to a critical level
18 where they may cause catastrophic failure of the component [31]. As the use of
19 composites in safety-critical parts continues to rise, the detection, characterisation, and
20 quantification of defects become increasingly important [31]. Moreover, as the
21 automated production of aerospace components continues to expand, so does the
22 requirement for automated testing. [33]. Automated NDE techniques are therefore well
23 placed to inspect these components as they can provide information about the integrity
24 of the structure below the surface at the scales required.

1 Intensive NDE is widely performed on aerospace components during manufacturing
2 and assembly [19]. However, some of the physical properties that make composites so
3 advantageous for use in components, also add significant challenges for NDE
4 inspection [19], [29], [30]. The anisotropic and inhomogeneous nature of composites
5 makes NDE far more challenging than compared to more traditional materials, such as
6 metals [21], [29]. For example, in ultrasonic NDE, wave propagation in anisotropic
7 composite structures is complex, with variability in scattering due to “macro”
8 structural features such as material lay-up and “micro” non-structural features such as
9 local anisotropy. This can lead to high attenuation of ultrasonic waves, reducing the
10 signal to noise ratio and in turn reducing the probability of detection [21], [32], [35].
11 In certain cases, this can pose challenges in distinguishing whether irregularities
12 detected during testing represent defects or variations in the base material [34].
13 Furthermore, accurately determining the ultrasound velocity within the sample adds
14 another layer of complexity. This parameter varies depending on specific factors such
15 as the resin type, fabric stacking sequences, and fibre orientation. Additionally, the
16 customisable anisotropic nature of composites, which can be tailored to meet design
17 specifications, further complicates matters. These intricacies introduce varying levels
18 of inconsistency across different components, posing significant complexities for the
19 NDE of composite materials.

20 *2.1.2.1 Defects*

21 The work conducted in this thesis focuses on defects occurring during manufacturing,
22 where a wide range of defects can be introduced [19]. However, damage can also be
23 caused during the service-life of a component from impacts, loading cycles etc. This
24 damage can often lead to cracking in cured samples which results in potential failure

1 modes that are challenging to predict due to the materials' anisotropic characteristics.
2 Further details on the specific defect types commonly seen in manufacturing are
3 outlined in this section.

4 *2.1.2.2 Voids and Porosities*

5 Voids and porosities are significant and common defects commonly encountered in
6 composite materials, which present significant challenges to their structural integrity
7 and performance [36]. Voids refer to isolated trapped air pockets or gas bubbles within
8 the material matrix, while porosities are microscopic pores which are often dispersed
9 throughout the composite structure. These defects can arise during the manufacturing
10 process due to incomplete resin impregnation, improper curing conditions, or
11 inadequate vacuum or pressure application during fabrication. Additionally, voids and
12 porosities can also result from the presence of contaminants or moisture in the
13 composite materials. Voids can often be detected ultrasonically since the significant
14 mismatch in acoustic properties between air and the composite material results in a
15 significant amplitude response. The size of porosities makes them far more challenging
16 to detect than voids, it is also very challenging to eliminate them completely from
17 manufacturing and most parts will have an allowable porosity percentage, for primary
18 aerospace structures this is often less than 2% [37]. Porosities therefore rarely give a
19 strong amplitude response ultrasonically and have to be detected due to an increased
20 attenuation, often determined as a lack or reduction in back wall response. The
21 presence of voids and porosities can compromise the mechanical properties of
22 composites, including strength, stiffness, and fatigue resistance, as they create stress
23 concentration points and reduce the effective load-bearing capacity of the material.

1 2.1.2.3 *Delaminations*

2 Delaminations represent a critical concern in composite materials, characterised by the
3 separation or splitting of layers within the laminate structure. They are one of the more
4 common defects found in composites [38]. These interfacial defects can arise during
5 the manufacturing process due to insufficient bonding between layers, voids, or resin-
6 rich regions. Delaminations can also occur because of mechanical loading, impact
7 damage, or environmental factors such as moisture absorption. Delaminations can
8 range in size and as a result have a ranging impact on a component from negligible to
9 severe. Left undetected or untreated, delaminations can propagate, leading to
10 significant reductions in structural integrity and mechanical performance. As
11 delamination's run parallel to the ply orientation and subsequently the surface of the
12 material, they are often well detected by ultrasonic testing due to their significant
13 amplitude response. However, in certain exceptional cases, plies may maintain contact
14 without the ability to transfer strain, a phenomenon commonly known as a "kissing
15 bond" [39]. Owing to the absence of separation between layers, detecting these specific
16 delamination's can prove particularly challenging.

17 2.1.2.4 *Foreign Object Inclusions*

18 Foreign Object Inclusions (FOI) encompass any unintended foreign materials
19 embedded within composite structures during manufacturing or service, and can
20 include particles, fibres, or debris from processing equipment [40]. These inclusions
21 may lead to localised stress concentrations, delamination, or initiation of cracks, which
22 diminish the overall strength and durability of the composite material. The detection
23 of FOIs within composite materials is influenced by the diverse properties of the

1 inclusion materials. For instance, fibrous inclusions may pose challenges for ultrasonic
2 detection because their acoustic properties closely resemble those of the surrounding
3 medium after curing. In contrast, FOIs with markedly distinct acoustic properties from
4 the composite matrix can be readily identified. This variability in detectability
5 underscores the importance of understanding the acoustic characteristics of different
6 inclusion materials and employing appropriate inspection techniques tailored to their
7 specific properties. In practice, the manufacturing of high-value components often
8 takes place within controlled environments such as clean rooms to minimise the
9 occurrence of FOIs. These sterile settings aim to reduce the introduction of
10 contaminants and foreign materials during the manufacturing process. However,
11 despite stringent measures, it is impossible to completely eliminate the risk of FOIs.

12 *2.1.2.5 Fibre distortions (or marcols)*

13 Fibre distortions can occur during the manufacturing process, where fibres may
14 experience misalignment, waviness, or kinking, compromising the intended
15 mechanical properties of the composite [41], [42]. Fiber distortions can result from
16 various factors, including improper handling, resin flow issues, or inadequate
17 consolidation during curing. These distortions can significantly impact the strength,
18 stiffness, and fatigue resistance of the composite, as they introduce weak points and
19 stress concentrations along the fibre-matrix interface. Additionally, fibre distortions
20 can affect the uniform distribution of load-bearing capabilities within the composite
21 structure, leading to non-uniform mechanical performance.

1 2.1.2.6 *Resin Rich/Starved Areas*

2 Resin-rich regions occur when excess resin accumulates within the composite, leading
3 to non-uniformity and potential weakening of the material. Conversely, resin-starved
4 areas occur when there is insufficient resin to fully saturate the reinforcing fibres,
5 resulting in inadequate bonding and reduced load-bearing capacity [43]. These issues
6 can arise during the manufacturing process due to improper resin infusion, resin flow,
7 or vacuum pressure. Resin-rich areas may lead to increased weight, reduced stiffness,
8 and potential delamination, while resin-starved regions can result in reduced strength,
9 increased susceptibility to cracking, and compromised durability. Quality control
10 measures, such as optimised resin infusion techniques and careful monitoring of resin-
11 to-fibre ratios, can help mitigate these issues during manufacturing. Changes in
12 ultrasonic attenuation can be used to evaluate resin rich or starved areas, as these
13 impact the components fibre-volume-fraction which has a direct impact on acoustic
14 impedance.

15 2.1.3 *Ultrasonic Testing of Composites*

16 UT has been widely adopted and standardised for testing in the aerospace industry due
17 to its ease of implementation and ability to detect a wide variety of defects [9], [15],
18 [18], [19]. The typical use for UT is with normal incident longitudinal waves so that
19 wave propagation is independent of ply orientation [19]. The interaction of
20 longitudinal waves and ply thickness is normally weak due to the wavelengths at
21 typical testing range of single megahertz frequencies being much greater than the
22 typical ply thicknesses [19]. The National Composites Network released a best practice
23 guide for applying NDE to composites [44]. This includes what inspection methods
24 are best suited for different types of defects. At the time of writing, UT pulse-echo had

1 been proven to detect the greatest number of flaws. Hsu et al. and Bossi and Georgeson
2 [19], [45] discussed specifically what defects UT can detect in aerospace composites,
3 and the differences in ultrasonic flaw interactions of each defect type. They also
4 suggested the most appropriate methods for detection of different defects. When
5 inspecting composite UT data inspectors will often use both time of flight and
6 amplitude C-scans simultaneously to aid with inspection [45].

7 Whilst many research papers rely on typical manufactured defects, Kokorov et al. [32],
8 aimed to simulate typical manufacturing defects more accurately in composites, where
9 the structural defects have similar physical and chemical properties as the binding
10 material. Their results demonstrated that UT can efficiently detect defects which have
11 similar physical and chemical properties to the composite binding material, helping to
12 bridge the gap between the detection capability of naturally occurring and
13 manufactured defects.

14 2.1.4 Artificial Intelligence

15 2.1.4.1 *Fundamentals*

16 Artificial Intelligence (AI) serves as an overarching concept encompassing machines
17 capable of leveraging knowledge and addressing diverse problem-solving tasks that
18 typically require human intelligence [46], [47]. There has been significant historical
19 research into different approaches for achieving AI. Early efforts in AI, particularly
20 between the 1950s and 1990s, focused heavily on symbolic reasoning and rule-based
21 systems. These paradigms formed the foundation of symbolic AI, a branch of AI where
22 knowledge was explicitly encoded into structured rules and logical relationships. This
23 approach aimed to model intelligence through explicit representation of knowledge,

1 enabling reasoning and decision-making in domains with clearly defined parameters
2 [48]. These systems demonstrated significant promise in early AI research for well-
3 defined tasks, with notable examples such as IBM's Deep Blue [49], a chess-playing
4 program, and ELIZA [50], the first chatbot. However, while rule-based approaches
5 excelled in structured environments, they often fail when encountering situations
6 which do not match their heuristics, which can often be the case for the complexity
7 and variability of real-world scenarios [51]. Model-based approaches attempt to
8 address these problems by basing solutions upon a theoretical model of a component
9 or system [52]. They are less dependent on expert opinion and offer higher flexibility
10 and scope for expansion. However, these systems can become expensive and complex
11 to construct. Symbolic AI has faced varying challenges such as the manual effort
12 required for knowledge acquisition, the rigidity of rule-based frameworks, challenges
13 with scalability, and its inability to effectively handle perception tasks like image and
14 speech recognition.

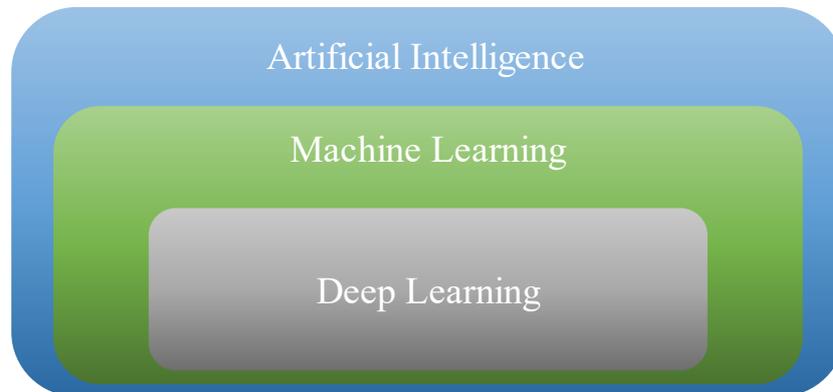
15 To overcome these challenges, the late 1980s and early 1990s marked a pivotal shift
16 toward connectionist approaches [53], which emphasized data-driven learning through
17 artificial neural networks. Unlike symbolic systems, which relied on predefined rules,
18 connectionist methods allowed machines to learn patterns and representations directly
19 from data, making them more adaptable and robust in handling noisy or incomplete
20 information. This shift was aided by advances in algorithms, particularly the
21 development of backpropagation [54], which facilitated the training of multi-layer
22 networks. Connectionist models proved especially adept at tasks requiring perception
23 and pattern recognition, areas where symbolic AI struggled. Concurrently, advances
24 in computing power and the growth of the global internet provided the computational

1 resources and vast datasets needed to train more complex and capable models, further
2 accelerating the adoption of connectionist paradigms.

3 While connectionist approaches gained significant traction, they did not entirely
4 replace other AI methodologies. Instead, the field expanded to incorporate a range of
5 techniques that remain relevant today. These include fuzzy logic for reasoning under
6 uncertainty, rule-based systems for applications requiring interpretability, model-
7 based reasoning for leveraging structured system models, evolutionary algorithms
8 inspired by natural selection, hybrid systems, and data-driven machine learning (ML)
9 methods [55], [56], [57].

10 ML constitutes the broad subset within the field of AI, emphasising the development
11 of algorithms and models capable of learning from data (such as connectionist
12 approaches) (Figure 10). This often allows for solving complex problems where rules-
13 based approaches would be intractable. ML problems can be categorised into either
14 classification or regression tasks. In classification tasks, the objective is to categorise
15 data into distinct classes or groups, whereas in regression tasks, the aim is to predict
16 continuous numerical values. DL is a subset of ML (Figure 10) as brought about by
17 the connectionist movement which, loosely inspired by the human brain, makes use of
18 neural networks with at least one hidden layer. In general, DL is employed to tackle
19 more intricate and demanding tasks compared to alternative ML methods. However,
20 achieving success in DL often necessitates access to substantial amounts of data. DL
21 is widely used and has been successfully applied to many challenging tasks such as
22 health care (diagnostic assistance, drug discovery, virtual healthcare etc. [58]), natural

1 language processing (translation, summarisation etc. [59]), or Computer Vision (CV)
2 (pose estimation, depth estimation, autonomous navigation etc. [60]).



3
4 *Figure 10: Illustration of the nested relationship between AI, ML and DL.*

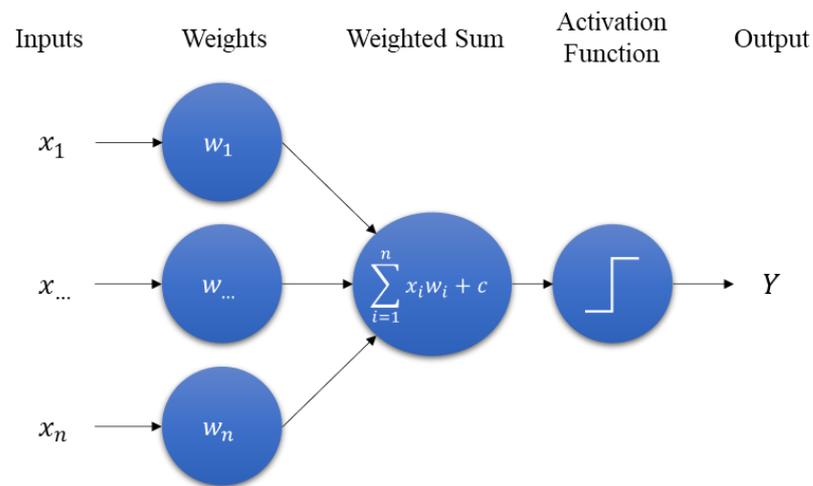
5 When trying to implement a ML solution there are many factors to consider. These
6 can be broken down into distinct areas: 1) the model which provides the solution, 2)
7 the data which allows the model to improve its performance, 3) a training regime to
8 effectively transfer the knowledge from the data onto the model. There is not a one-
9 size-fits-all method for every task and finding the right approach for each domain is
10 not straightforward. Variations in one aspect often influence other areas, making the
11 process complex and non-trivial.

12 The massive increase in data availability and compute resources over the last couple
13 of decades has led to an explosion in the use of DL and has provided state-of-the-art
14 (SOTA) solutions to tasks previously thought unsolvable, such as complex perception
15 tasks like image processing. In the context of NDE, DL offers distinct advantages over
16 traditional AI approaches like symbolic AI. While symbolic AI relies on predefined
17 rules and logical reasoning, it can struggle to handle unstructured data, such as images,
18 signals, and sensor outputs, which are common in NDE applications. Symbolic
19 systems are also less flexible when faced with the variability and noise inherent in real-

1 world data. On the other hand, DL models, can learn features from raw data, enabling
2 them to adapt to complex, noisy, and high-dimensional inputs without the need for
3 exhaustive manual feature engineering which can be highly application specific and
4 require extensive expert intervention. This makes DL particularly effective for defect
5 detection and classification in materials, where patterns can be subtle and
6 unpredictable. Given these advantages, this thesis primarily concentrates on DL
7 solutions. Consequently, the following sections offer an overview of DL to provide
8 necessary context and understanding.

9 *2.1.4.2 Deep Learning*

10 The perceptron, introduced by McCulloch and Pitts in 1943 [61] and first demonstrated
11 by Frank Rosenblatt in 1957 [62], is the basic building block for the original neural
12 networks and was inspired by workings of biological neurons in the brain [63]. The
13 functioning mechanism involves taking inputs, summing them following a linear
14 transformation using weights along with the addition of a constant term, known as a
15 bias. Subsequently a non-linear transformation is applied which facilitates
16 backpropagation and allows for stacking of multiple layers of neurons, this is often
17 referred to as an activation function (Figure 11). By combining multiple perceptron's
18 together to produce Multi-Layer Perceptron's (MLPs) it is possible to approximate any
19 function [64]. Advances in neural network research has introduced modifications to
20 the perceptron, but the underlying principles are the foundations for most neural
21 networks.



1
2 *Figure 11: Illustration of a perceptron.*

3 DL encompasses a range of different neural network architectures which all have at
4 least one hidden layer [65]. Traditionally, these architectures faced constraints on the
5 number of layers due to challenges in information propagation through deep networks,
6 a phenomenon known as the vanishing gradient problem. [66]. However, as research
7 into deeper networks continued solutions to the vanishing gradient problem emerged
8 [67]. This facilitated the construction of significantly deeper and larger neural
9 networks, empowering them to learn far more intricate features and undertake more
10 complex tasks.

11 Many different network architectures exist, and they are often task dependent, with the
12 construction of the best network often a challenging part in the process in applying
13 DL. However, despite the near infinite array of different architecture combinations,
14 there are a few key principles that underpin the training of these models. These are
15 forward pass, network loss, and back propagation.

16 Forward pass is the process of passing data through a model, crucially this can then be
17 evaluated using some metric known as loss. The network loss is the measure of how

1 well our model fits to the training data. A model with a high loss hasn't learnt much
2 information from the training data whilst a model with a low loss is capable of mapping
3 to the training data very well. There are multiple different types of loss functions used
4 in training DL models, for details on specific loss functions please refer to the relevant
5 literature [68], [69].

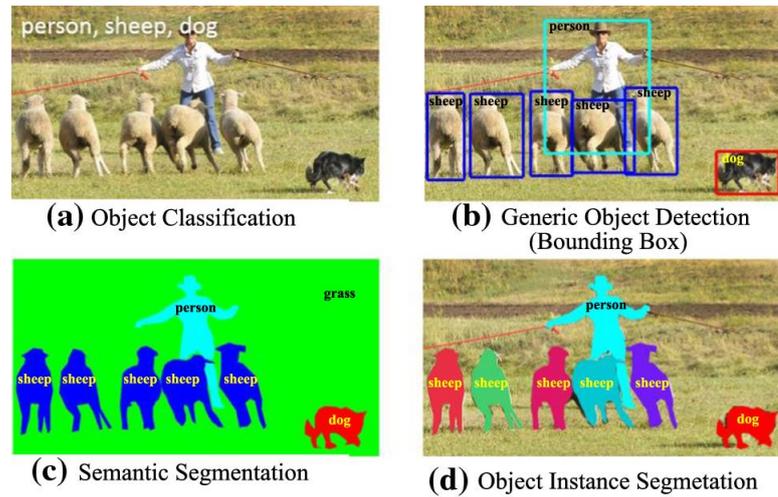
6 After the model's performance is assessed using the forward pass and a loss value is
7 calculated, backpropagation is then used to adjust the model to better fit to the training
8 data. This generally involves adjusting the model parameters using one of many
9 optimisation algorithms based on the gradients computed during backpropagation. The
10 field of optimisation methods for DL models is a large one with many different
11 approaches available. The Adam optimiser [70] has proven to be highly effective in
12 many cases [67], [71]. Nonetheless, numerous alternative options exist, and once
13 again, the choice of optimiser is often contingent on the specific application. For
14 additional insights on model optimisers, please consult the available literature. [72],
15 [73].

16 This iterative learning process is generally repeated multiple times during training until
17 a model is produced that has learnt some representation from the training data.
18 Typically, a model is evaluated during training against a separate dataset known as the
19 validation set. The validation set comprises data that is from the same distribution as
20 the training data, but the model has not been exposed to during training. This
21 evaluation allows for monitoring the model's performance on unseen data, providing
22 insights into its generalisation capabilities. By assessing performance on the validation
23 set, training can be stopped at an appropriate point, preventing overfitting to the

1 training data, and ensuring that the model performs well on new, unseen data. This
2 approach helps to ensure that the model learns meaningful patterns from the data,
3 rather than simply memorising the training examples. Once training is completed final
4 model performance is evaluated on a separate, unseen test set. This ensures there is
5 unbiased reporting of the final model performance.

6 *2.1.4.3 Computer Vision and Convolutional Neural Networks*

7 CV is the field of using computers to interpret and draw information from visual data.
8 Over the last couple of decades DL has revolutionised the field of CV [74], [75]. There
9 are many different CV applications, but they can generally be categorised into
10 classification, object detection, and segmentation tasks. In classification tasks, models
11 aim to extract features from images or video data and categorise them into distinct
12 groups based on the visual information. For object detection the goal is to localise
13 specific objects within an image or video, generally by drawing bounding boxes
14 around them. This process integrates classification within an image and combines it
15 with approximate instance localisation. For segmentation type tasks the objective is to
16 classify each pixel within an image, thus not only achieving classification but also
17 highly accurate localisation (Figure 12). This can be extended to instance segmentation
18 where individual objects are distinguished from one another, enabling precise
19 delineation and identification of each object instance within the image. Considering
20 that ultrasonic data is often visualised as images, comparing these tasks with the
21 automated ultrasonic inspection pipeline (Figure 2) underscores a notable alignment.
22 This implies that computer vision methodologies could offer potential solutions to
23 many of the challenges in automated UT data analysis.

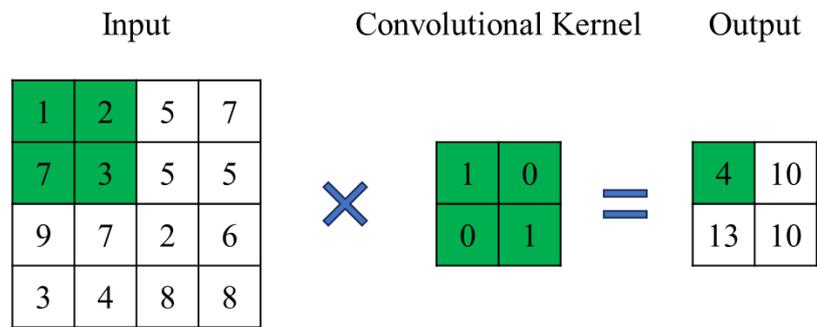


1

2 *Figure 12: Classical computer vision problems. (a) Image level classification, (b) Object detection,*
 3 *(c) pixel wise semantic segmentation, (d) instance level semantic segmentation. [76]*

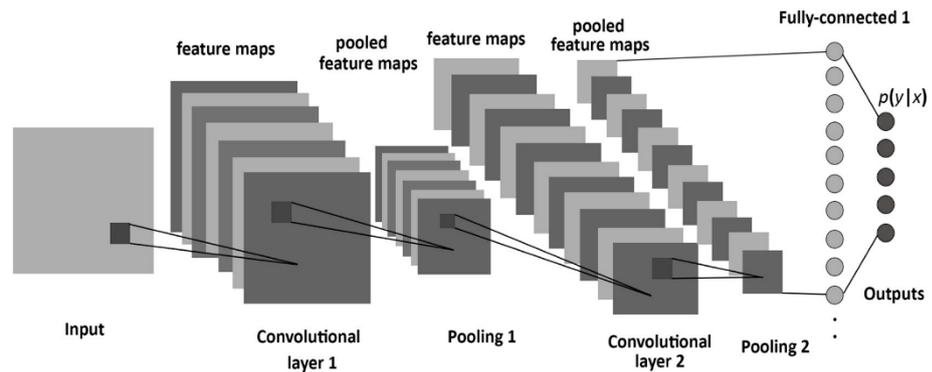
4 As previously discussed, DL encompasses a wide array of different model
 5 architectures, but in CV applications Convolutional Neural Networks (CNNs) have
 6 repeatedly demonstrated wide scale success in image classification and other CV tasks
 7 [74].

8 CNNs make use of convolutional operations, which are built into convolutional layers
 9 to extract information from images (Figure 13). Convolutional layers are often
 10 combined with down sampling layers for dimensionality reduction and activation
 11 functions to construct complete network architectures (Figure 14).



1

2 *Figure 13: Illustration of a convolutional filter.*



3

4 *Figure 14: The structure of a generic CNN, consisting of convolutional, pooling and fully connected*
 5 *layers [77].*

6 Over the past decade CNNs have evolved significantly with changes in architectures,
 7 training methodologies, data availability, and increased computational power all
 8 contributing to improved results. 2012 marked a significant milestone in the evolution
 9 of CNNs with the emergence of AlexNet [78]. AlexNet utilised deep convolutional
 10 layers, overlap pooling, Rectified Linear Unit (ReLU) activation functions, and
 11 dropout for regularisation. This was the first demonstration of the remarkable potential
 12 of CNNs at scale, significantly improving on previous results for the ImageNet Large
 13 Scale Visual Recognition Challenge [79]. This achievement underscored the efficacy
 14 of CNN methodologies in the domain of CV, triggering a surge of interest and
 15 subsequent advancements in CNN architectures and applications. After the
 16 introduction of AlexNet, the field witnessed notable progressions, exemplified by the

1 advent of VGG in 2014 [80]. VGG introduced novel architectural features such as
2 smaller convolutional filters coupled with deeper network structures, thereby
3 enhancing the capacity of CNNs to extract intricate features from visual data. VGG's
4 simplicity in using uniform layers made it easier to understand and implement,
5 although it came with increased computational cost due to its depth and large number
6 of parameters. In 2015, the introduction of ResNet [67] presented a breakthrough by
7 introducing skip connections, which helped mitigate the vanishing gradient problem.
8 Prior to the introduction of ResNet the vanishing gradient problem was limiting the
9 depth of CNNs being trained, which in turn limited their ability to learn complex data
10 representations. The vanishing gradient problem refers to the phenomenon in deep
11 neural networks where gradients calculated during backpropagation diminish
12 exponentially as they propagate backward through multiple layers of the network. This
13 diminishing gradient effect can hinder the training process by making it difficult to
14 update the weights of early layers effectively, leading to slower convergence or even
15 stagnation in learning. Skip connections enabled the gradient to flow directly through
16 the network, making it feasible to train exceedingly deep networks comprising
17 hundreds of layers and thus capture more complex features.

18 CNNs continued to advance, with competition from other architectures such as
19 transformers providing an alternative method and often achieving near state-of-the-art
20 results in large visual recognition tasks [81]. Vision Transformers (ViTs) represent a
21 departure from traditional CNNs by leveraging transformer architecture, originally
22 popularised in natural language processing. ViTs process images as sequences of
23 patches, using the self-attention mechanisms to capture global and local dependencies
24 within these patches. This approach has shown promising results in various computer

1 vision tasks, demonstrating its potential to compete with and sometimes outperform
2 CNNs, particularly in tasks requiring long-range dependencies and global context
3 understanding. However, ViTs generally demand larger amounts of training data
4 compared to CNNs due to their dependence on self-attention mechanisms, which
5 require sufficient diverse examples to effectively learn both global and local
6 dependencies within image patches. In contrast, CNNs traditionally rely on local
7 receptive fields and hierarchical feature extraction, which can sometimes generalise
8 well even with smaller datasets. Therefore, while ViTs have shown promising results,
9 particularly in tasks demanding global context understanding, their performance often
10 hinges on the availability and quality of extensive training data.

11 Despite the impressive performance and popularity of ViTs, in 2022 Liu et al.
12 introduced ConvNext [82], reaffirming CNNs' ongoing effectiveness in computer
13 vision tasks. ConvNext demonstrated how different architectural modernisations could
14 improve a standard ResNet to compete with Swin Transformers (a form of ViT) [83].
15 This work brought together and analysed the impact of several strategies for
16 modernising CNNs (often inspired by ViTs). This included the integration of skip-
17 connections, separate down sampling layers, the exploration of alternative and fewer
18 activation functions such as Gaussian Error Linear Units (GELU), a reduction in
19 normalisations and the use of layer normalisation instead of batch normalisation, and
20 experimentation with varying kernel sizes. This work resulted in a significant
21 advancement in the capabilities of CNNs, showcasing their adaptability and
22 competitiveness against transformer-based architecture.

1 It is evident that, akin to the broader landscape of DL, the architecture of modern CNNs
2 encompasses a diverse array of configurations and design choices, often tailored to the
3 specific requirements of individual tasks [84], [85]. This diversity stems from
4 fundamental differences in the design principles and components of these
5 architectures, which significantly impact their performance and suitability for various
6 computer vision challenges.

7 *2.1.4.4 Data*

8 ML systems depend on data to learn from during training. Consequently, a substantial
9 amount of high-quality training data is crucial in the development of ML, especially
10 DL models as they require large amounts of data to effectively update numerous
11 parameters and uncover complex solutions to tasks. The availability of large open-
12 source datasets such as CIFAR-10 [86] and ImageNet [79] has enabled massive growth
13 in CV and has contributed significantly to the enhanced accuracy of models.

14 Zhu et al. [87] evaluated how the size of a training datasets impacts the quality of
15 object detection models. They concluded that simple models become saturated
16 surprisingly early and are unable to benefit from very large datasets. However, more
17 complex models can make use of greater features and benefit from significant datasets
18 (potentially up to 10^{12}). They go on to suggest that the relationship between model
19 complexity and saturation of datasets may be logarithmic; suggesting that as model
20 complexity increases datasets will need to increase exponentially before the model
21 becomes saturated. However, the returns from super large datasets quickly diminish
22 due to the logarithmic nature of the performance increase.

1 Data augmentation is a valuable technique to try and get more out of limited training
2 datasets by modifying the original dataset. For images, there exist a wide range of
3 augmentation methods from geometric transformations, such as flipping, rotations,
4 translation, cropping or colour space adjustments, to image mixing or random erasing
5 [88]. Whilst, many different augmentation methods are available, the selection of the
6 appropriate method/methods is dependent on the specific application area, as not all
7 methods will always be suitable. For example, random erasing is not always a label-
8 preserving transformation and may make the image unrecognisable. In an NDE
9 context, if you were to erase the defect signature from an image, the image would not
10 represent the original defect classification. Therefore, this method may not be
11 applicable or would require adjustment/manual intervention such as object-aware
12 random erasing.

13 *2.1.4.5 Summary*

14 AI, ML, and DL are broad and complex fields with a wealth of research and
15 applications that cannot be fully covered in this thesis. This section introduces key
16 concepts and terms relevant to the scope of this work. AI aims to enable machines to
17 replicate tasks typically requiring human intelligence. Early AI approaches often
18 focused on rule-based systems to encode knowledge and reasoning capabilities. These
19 systems excelled in structured environments however, they faced challenges in
20 adapting to the complexity of real-world data, particularly in areas of perception. To
21 address these limitations, the late 20th century saw a shift toward connectionist
22 approaches, particularly artificial neural networks, which emphasized data-driven
23 learning. Advances like backpropagation, increasing computational power, and the
24 proliferation of large datasets catalyzed the adoption of these methods, laying the
25 groundwork for modern DL.

1 While DL demands substantial data and computational resources, it has revolutionized
2 various fields, particularly in image analysis. Unlike traditional symbolic AI, DL
3 models automatically extract features from data, making them particularly effective
4 for handling data variability. CNNs are particularly pivotal in computer vision,
5 excelling in tasks like image classification, object detection, and segmentation.
6 Architectural milestones, including AlexNet, VGG, and ResNet, have significantly
7 enhanced the capability of CNNs to process and interpret visual data. More recently,
8 ViTs and ConvNext have introduced advancements, highlighting the ongoing
9 evolution of neural network architectures.

10 The success of DL relies heavily on the availability of high-quality data for training.
11 Large datasets like ImageNet have been instrumental in advancing DL models,
12 although data augmentation techniques remain essential for enhancing model
13 performance when datasets are limited. In the context of NDE, careful consideration
14 is required to ensure augmentation methods preserve the integrity of defect-related
15 data. As DL continues to evolve, its adaptability and effectiveness position it as a
16 cornerstone technology for solving complex challenges in fields like NDE, particularly
17 in automated analysis of ultrasonic data.

18 2.1.5 Artificial Intelligence in Non-Destructive Evaluation

19 The use of AI systems for NDE have been explored as an alternative solution to manual
20 data processing. Historically, simple rule-based algorithms, such as thresholding or
21 comparison to reference images were used [89], [90]. However, they were often unable
22 to replicate human level performance due to a lack of reliable results as a result of
23 noise and variations from operational environments [91]. As a result, human operators
24 remain more trusted as they can adapt to changes in the data.

25 Despite these limitations, rule-based systems have demonstrated notable successes.
26 For instance, commercial software tools such as “NDT kit” offers some automated
27 detection capabilities based on pre-defined data processing steps such as thresholding
28 or subtraction from a reference scan. This acts as an operator assistance tool, however
29 it still requires significant operator interaction, for example to narrow down the area

1 of interest [92]. Despite not being a fully automated solution, analysis times have been
2 greatly improved with such tools. Similarly, “ADA” represents an effort to encode
3 inspection procedures into an automated data analysis pipeline [93]. While promising,
4 ADA faces challenges in rejecting noise artifacts caused by scanning processes—
5 artifacts that an experienced human operator would intuitively dismiss as non-defects.
6 This approach also depends on precise parameterization of acceptance criteria,
7 typically calibrated to reference samples. Such characteristics make it well-suited to
8 aligning with regulatory requirements. However, these systems have demonstrated
9 efficacy only in scenarios involving simplified geometries and have yet to prove their
10 robustness in inspecting complex industrial components or challenging test conditions.
11 For dealing with complexities in geometries Guo et al. suggested the use of a model-
12 based approach, which relied on adjusting gating parameters to align with the digital
13 representation of the component [94]. Whilst this demonstrated promise, it requires an
14 accurate model of the component and correctly aligning or registering this model to
15 the real-world component.

16 Rule-based approaches have shown utility for processing signals with well-defined
17 characteristics. They benefit from explainability but are limited in their ability to
18 handle variability and complexity caused by environmental changes, manufacturing
19 inconsistencies, and intricate geometries or materials. These limitations have
20 prevented their widespread adoption in complex industrial applications. Data-driven
21 ML methods, by contrast, offer greater adaptability, as they are not bound by rigid
22 predefined parameters. A meta-analysis by Sergio Cantero-Chinchilla et al. [14]
23 underscored this gap and highlighted DL as a promising solution, owing to its
24 capability to extract complex features from data.

25 *2.1.5.1 Machine Learning in Non-Destructive Evaluation*

26 The application of ML in NDE is still a relatively young field, with the earliest found
27 example published in 1992 [95], however there are plenty of examples of researchers
28 demonstrating the effectiveness of different methods for specific problems.

1 As the problem space in NDE is highly variable depending on material, geometry,
2 defects of interest, operating conditions, and sensing modality, comparing different
3 ML applications in NDE can be challenging and expecting models to transfer between
4 all NDE scenarios is at this current moment not possible. This therefore makes direct
5 comparisons in literature challenging. One method which helps to break down and
6 compare ML in NDE literature is by grouping based upon data analysis type and
7 material application. Such as, ML applied to timeseries data (e.g. [96], [97], [98]) like
8 ultrasonic waveforms or to image level analysis (e.g. [99], [100], [101]), typically
9 constructed from multiple timeseries sources. The benefit of this is that it allows for
10 comparison of more closely aligned ML methods, as for example image level
11 ultrasonic analysis may have more in common with image level thermography analysis
12 than time-series A-scan analysis.

13 While grouping ML approaches in NDE by data type and material application aids in
14 facilitating comparisons between similar approaches, defining a true (SOTA) remains
15 challenging. The NDE field encompasses a wide range of inspection tasks, from defect
16 detection to material classification, each with unique requirements influenced by
17 component geometry, material properties, data availability, defect characteristics, and
18 sensing modalities. This diversity complicates the development of universally
19 applicable models and creates difficulty in establishing reliable benchmarks. Unlike
20 fields like computer vision, where standardized datasets such as ImageNet [79] enable
21 direct performance comparisons, NDE lacks widely accepted public datasets. This
22 absence limits cross-study validation and makes it difficult for the community to
23 identify and agree upon SOTA methods or approaches.

1 By contrast, to other fields NDE data is often restricted by industrial data-sharing and
2 is highly specialized with formats that vary widely, from time-series data to volumetric
3 imaging. Without accepted datasets, researchers rely primarily on proprietary data,
4 limiting the ability to replicate findings and preventing a cohesive comparison across
5 different studies. This reliance on specific, often inaccessible datasets restricts the
6 establishment of a widely recognized SOTA model within the NDE community.

7 Despite these challenges, certain ML architectures or approaches that have achieved
8 success elsewhere can be applied in NDE, though they typically require adaptation to
9 meet the distinct requirements of each application. For instance, U-Net based
10 architectures have been widely accepted as the SOTA for image segmentation [102].
11 The modified approaches can then be compared to the widely established SOTA for
12 similar tasks, as was done with DefectDet [103], which was introduced as a SOTA
13 approach for defect detection in B-scans of metal samples and compared to other
14 commonly used object detection architectures, but its results remain unvalidated
15 beyond the original study and application due to the inaccessibility of the dataset. The
16 need for customization of techniques to NDE highlights the difficulty of establishing
17 a one-size-fits-all SOTA model in NDE, as models must frequently be tailored to suit
18 highly specific conditions. These task-specific modifications, while effective for
19 individual cases, further complicate the cross-comparison of ML models across
20 studies.

21 Some researchers have taken steps to address data availability issues by making their
22 datasets and code publicly accessible. For example, a dataset generated from fabricated
23 steel plates with manufactured flaws was released as an open-source ultrasonic
24 imaging dataset USimgAIST [104], containing over 7,000 images of steel plates with

1 and without flaws. While this initiative represents an important step toward more open
2 NDE ML research, such contributions remain rare and highly specific, underscoring
3 the need for broader efforts to build shared benchmarks and encourage collaboration
4 across the field.

5 Another trend in ML in NDE literature, is that older works focused on analysing time
6 series data. These often used a large amount of hand crafted or extracted features to
7 reduce the dimensionality of the data [96]. Whilst time-series analysis is still popular,
8 thanks to improvements in the CV field and computing performance, a greater
9 percentage of recent works have focused on image-based methods [99]. While ML has
10 been applied across various NDE sensing methodologies, the literature often
11 emphasises metal welds as a primary focus, with limited exploration of composites.
12 Composites, which are frequently imaged differently and possess distinct properties
13 from metals, remain relatively underexplored in this context.

14 For each task, regardless of data type, most applications of ML in NDE typically fall
15 into one of three categories. Gardner et al. [7] suggested a hierarchy for automated
16 NDE problems, summarising and explaining these problems. This was based on the
17 Rytter's Hierarchy previously used in Structural Health Monitoring (SHM) [105]:

- 18 1. DETECTION: the method gives a qualitative indication that damage might be
19 present in the structure.
- 20 2. CLASSIFICATION: the method gives information about the type of damage.
- 21 3. ASSESSMENT: the method gives an estimate of the extent of the damage.

22 The hierarchy provides a clear picture of both the importance of each task and
23 increasing complexity. It is also important to recognise that different industrial

1 scenarios will have different requirements and whilst full defect assessment such as
2 type, size and position may be required for some settings, others may just require
3 limited detection capability. The hierarchy is in strong agreement with the automated
4 data analysis pipeline proposed in Figure 2, which groups Classification and
5 Assessment under “Characterisation”, whilst demonstrating their integration into a
6 broader automation framework.

7 As previously discussed, ML methods rely on training data, with increasing quantities
8 of data required often for more complex methods. The availability of training data is
9 therefore a primary concern when applying ML to NDE. NDE encounters challenges
10 in generating large datasets of real samples and defects, primarily due to the limited
11 availability of such samples, particularly of defective examples, and the significant
12 time investment required to collect the data. To address this challenge, researchers
13 often resort to employing transformations and other techniques to augment the
14 datasets. Alternatively, other papers have proposed using simulated data with
15 encouraging results [100], [106]. However, for this to be effective it is crucial that the
16 simulated data accurately reflects the problem. There is often a challenge when using
17 synthetic data that when testing the model trained on synthetic data one must be careful
18 not to commit an ‘inverse crime’ [107], by simply validating the results from features
19 extracted from a synthetic model with synthetic data generated from the same synthetic
20 model. This would produce a result that does not generalise well to the real problem.
21 To negate this, it is therefore preferable to test models using real experimentally
22 acquired datasets.

1 As highlighted by Cantero-Chinchilla et al. [14], in order for deep learning to see
2 adoption in NDE by both industry and regulators the explainability of models is a key
3 requirement which will help to build trust in the results. Explainable AI (XAI) is a
4 significant, ongoing area of research which involves many different approaches to
5 understand the reasoning behind AI systems [108], [109]. These approaches can range
6 from model-specific to model-agnostic explainability approaches, but there is
7 generally a trade-off between model accuracy and interpretability. Models which give
8 a level of uncertainty in their own results is one method of helping to build trust in
9 their outputs, as it will help to de-mystify the “black-box” of deep learning. Abdar et
10 al. [110], have performed a comprehensive review of uncertainty quantification
11 methods in both traditional ML and DL applications. They use the medical field as a
12 case study for the need for uncertainty quantification. There is significant overlap
13 between the medical and NDE industries when applying DL, for example the
14 challenges in acquiring enough data and the conservative nature of the industries. An
15 alternative method to explainability is using symbolic representations of expert
16 knowledge and formalising these representations as a rule-based system. Young et al.
17 [111] explored how these representations could be both manually and automatically
18 captured and formalised for industrial fault detection systems.

19 For a detailed review of DL data analysis in NDE please refer to Cantero-Chinchilla
20 et al. work [14]. This review paper provides a comprehensive summary of DL,
21 automated ultrasonic methods, and how these relate to different levels of industrial
22 automation. Whilst this work focuses on ultrasonic data, many of the conclusions and
23 challenges identified are transferable to different sensing modalities. For example, the
24 authors suggested different classifications for levels of automation; based upon

1 automation levels published by the European Union Aviation Safety Agency [112],
2 [113]. This can be broadly applied to different NDE methodologies:

3 Level 0 – Classical NDE

4 Level 1 – Operator assistance

5 Level 2 – Partial automation

6 Level 3 – Operational automation

7 Level 4 – Full automation

8 2.1.6 Summary of Key Challenges

9 The automation of NDE data analysis remains a significant challenge. While efforts
10 have been made to develop automated systems using rule-based logic, their adoption
11 in industry has been limited. A key reason for this is their lack of flexibility and limited
12 applicability to real-world inspection scenarios, where human interpretation often
13 plays a critical role in interpreting data from complex geometries or handling
14 variations from environmental changes or anomalies. ML methods offer a promising
15 alternative by enabling systems to "learn" the flexibility required to handle these
16 complexities—something that is difficult to encode explicitly. DL, as the SOTA in
17 many perception tasks, including CV, presents a compelling avenue for exploration in
18 this context, despite inherent challenges such as the need for greater explainability.

19 DL research in NDE is at a relatively early stage, but it has still seen varied application
20 to different inspection scenarios. For UT this is mainly applied to metallic welds, with
21 a distinct lack of research in the application of DL to composites. From the current
22 literature, it is evident that there are key shared challenges when applying DL to NDE.
23 These are outlined in the remainder of section 2.1.6.

1 *2.1.6.1 Lack of Training Data*

2 Lack of training data may be the biggest challenge when applying DL to NDE. This
3 holds especially true for naturally occurring, real defects, which have not been
4 artificially manufactured, as they provide responses that more accurately represent
5 real-world responses. However, gathering sufficient examples of naturally occurring
6 defects for training is challenging, and accurately establishing and labeling the ground
7 truth for these defects presents an additional difficulty. Meanwhile, training on
8 manufactured defects can lead to models which do not perform well on real defective
9 data. There are two main approaches to overcome the problem of overfitting due to
10 small datasets:

- 11 1. **Experimentally increase the dataset:** This is challenging as NDE real-
12 flaw data is not readily available and using manufactured defects can lead
13 to a lack of generalisation (as previously discussed).
- 14 2. **Augment the available data:** This method has been widely and
15 successfully adopted in other ML applications, such as computer vision
16 [88]. It has exhibited use in NDE as well in different forms. However,
17 successfully augmenting a small dataset to cover the wide application of
18 real-world flaws is a challenging task, which requires a well thought out
19 methodology based on good understanding of the NDE modality and the
20 targeted defects. This is currently the most widely adopted approach in
21 NDE literature. Where researchers will generally create datasets from a
22 handful of manufactured signals, then augment these to produce datasets
23 ranging in sizes anywhere from 100's to over 1 million (typical papers use
24 an augmented datasets on the scale of four order of magnitude) [114], [115],

1 [116].

2 In NDE scenarios, there are various methods to augment training data:

- 3 • Simple linear data augmentation or sub-sampling techniques, akin to those
4 commonly employed in other CV applications, are often utilised for images
5 [115].
- 6 • Virtual flaws can be introduced, where flawed signals are implanted into noisy
7 background data to augment the number of flawed signals [116], [117].
- 8 • Different types of transfer learning can be applied to mitigate overfitting of
9 networks when dealing with small datasets [118].
- 10 • Simulated data, such as finite element analysis (FEA) simulations, or a
11 combination of simulated and real data, can be employed to expand the size of
12 datasets [99].
- 13 • Generative Adversarial Networks (GANs) can be leveraged to generate entirely
14 new synthetic datasets [119], [120].

15 Regardless of the technique employed, the crucial takeaway is the importance of
16 enhancing the networks' ability to generalise by exposing it to a greater distribution of
17 samples during training. There are a wide range of other approaches for reducing
18 model overfitting include adjusting hyperparameters or implementing dropout etc.
19 [121], [122]. However, it is widely acknowledged that if the initial data is of poor
20 quality or not representative of the target domain or distribution, the model
21 performance will invariably be subpar.

1 *2.1.6.2 Defining the Problem Scope*

2 When implementing a DL methodology, it is crucial to have a well-defined problem.
3 This ensures that effective training data, representative of the problem, can be
4 established. When applying DL to NDE, the problem can often be broken down into;
5 what is the task (detection, classification, assessment etc), and what is the target
6 domain (signals, scans, materials, defects etc.). For example, characterising a defect in
7 steel welds is different to detecting a defect in composites, which is a different task to
8 segmenting cracks in concrete, and the training data will need to reflect this. When
9 trying to detect defects it is important to not only understand what types of defects are
10 likely to occur, but also to understand what types of signals and scans are most
11 appropriate. This dictates what will serve as the model input. All these factors
12 influence the applicability of different types of models. Depending on the specific
13 scenario, it may be appropriate to construct a custom network, while in other cases,
14 utilising an off-the-shelf model and fine-tuning it to meet the task requirements may
15 be more suitable.

16 *2.1.6.3 Model Evaluation and Explainability*

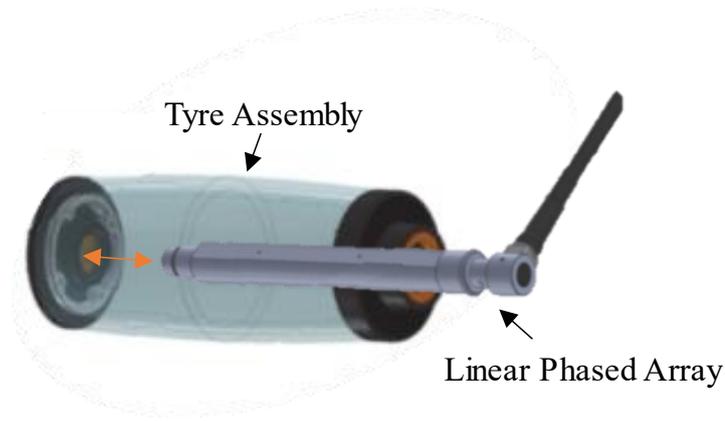
17 NDE is frequently employed in highly safety-critical industries. However, the
18 adoption of new technologies in these sectors tends to be conservative, given the
19 paramount importance placed on safety and reliability. For adoptions of DL/ML
20 methods into these industries, it is important to effectively communicate and evaluate
21 the performance of these models. A great way of doing this is by comparing them to
22 human operators and POD scores, as well as traditional ML statistical evaluation
23 methods. The testing of models should also be done very carefully with a large enough,

1 unseen, test set that is representative of the target domain. This should also be pushed
2 past the target domain to determine the models' limitations. Furthermore, it is essential
3 to demystify the 'black box' nature of deep learning DL, particularly in highly regulated
4 industries. To achieve this, XAI methods should be employed to make it easier to
5 understand how models arrive at their predictions and what influences specific
6 decisions. XAI can enhance transparency and is likely to be crucial for meeting
7 regulatory standards and ethical requirements, helping to build trust and allowing DL
8 models to be safely integrated into critical applications such as aerospace. Researchers
9 should also be clear about the limitations of models even if these are a result of the
10 NDE method itself.

1 2.2 Experimental Data Collection

2 Much of the work in this thesis was conducted using UT data collected experimentally
3 with a robotic system from composite samples. The methodology for data acquisition
4 is largely consistent across the chapters of this thesis. Hence, this section provides
5 background information about the data collection process and the samples used.

6 The ultrasonic data was acquired at room temperature using a robotically deployed
7 unfocused linear phased array. The array used was an Olympus Inspection Solutions
8 RollerFORM-5L64 [123] (Figure 15), which had a central frequency of 5 MHz and
9 was made up of 64-elements with a pitch of 0.8 mm and elevation of 6.4 mm. The
10 roller probe was initially filled with deionised water as per the manufacturer's
11 guidelines; however, this was later switched to non-corrosive glycol. Glycol has
12 similar acoustic properties to water and is an appropriate alternative as per the
13 manufacturer's guidance. The benefits of using glycol were two-fold. Firstly, when
14 using water, the array was required to be removed and cleaned regularly, as constant
15 submersion led to an increased risk of delamination of the array's protective layers.
16 Secondly, the increased viscosity of glycol meant that the formation of microbubbles
17 on the surface of the array was less common which minimised any acoustic barriers.
18 The elements were driven at 100 V with a receiver gain of 22.5 dB to maximise the
19 front wall amplitude without signal saturation. The sample rate was 100 MHz. A Peak
20 NDT Ltd. MicroPulse 6 or LPTA [124] (both supporting at least 64 transmit/receive
21 channels) was used for ultrasonic control with a digital band pass filter applied on
22 reception to filter out frequencies between 2 and 6 MHz. A summary of roller probe
23 parameters is given in Table 2.



1

2 *Figure 15: Diagram of roller probe assembly. Adapted from [124].*

3 *Table 2: Summary of roller probe parameters*

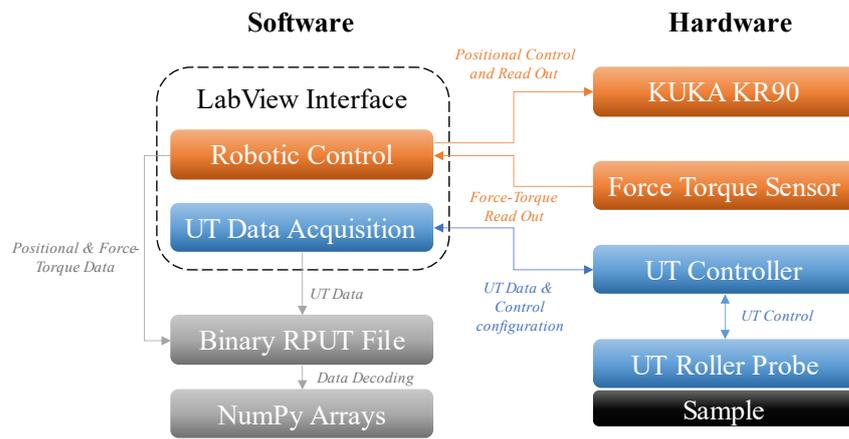
Manufacturer	Olympus Inspection Solutions
Model	RollerFORM-5L-64
Central operating frequency	5 MHz
Number of elements	64
Pitch	0.8 mm
Elevation	6.4 mm
Driving Voltage	100 V
Receiver Gain	22.5 dB

4

5 The pulse repetition frequency was set to collect a 4-element aperture, unfocused B-
 6 scan every 0.8 mm with A-scan speed of 10 mm/s which was controlled using a fully
 7 automated robotic system built around a KUKA KR 90 R3100 extra HA industrial
 8 robot (Figure 17) [125]. Robotic scanning enabled the concatenation of encoded B-
 9 scans to form volumetric datasets. Robotic scanning was essential in accurate
 10 concatenation of scans. This process was highly repeatable thanks to a ± 0.04 mm pose
 11 repeatability. To ensure a steady coupling of the roller-probe to the surface of the
 12 component and consistent transfer of acoustic wave energy into the sample at different
 13 scanning positions, Force-Torque compensation was used to control the contact force
 14 on the samples surface with feedback from the force axis perpendicular to the sample.
 15 This was accomplished with integration of a Schunk GmbH & Co. FTN-GAMMA-

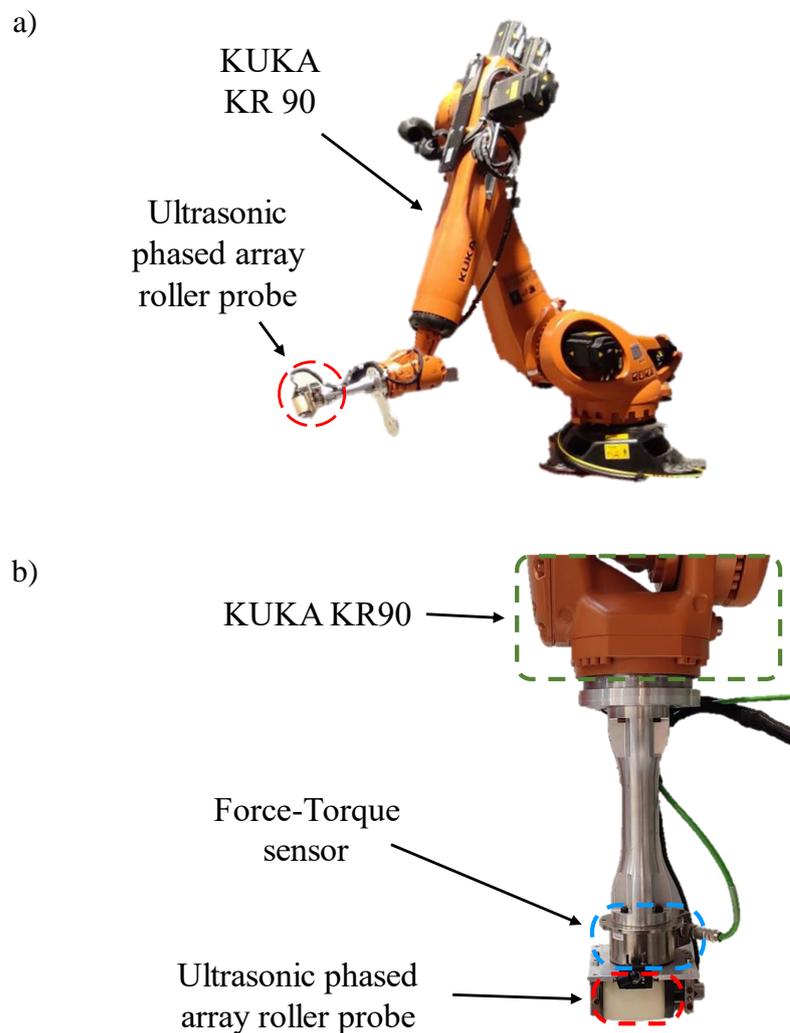
1 IP65 SI-130-10 Force-Torque sensor, mounted between the robot's flange and the
2 roller-probe. Initial experiments used a contact force of 35 N (Chapter 3), however
3 latter experiments increased the force to 70N which increased the consistency of
4 acoustic transmission by ensuring consistent tyre compression throughout the scan.
5 The Force-Torque sensor had a programmed limit of 130N in the primary force axis.
6 This was set to avoid any potential damage caused by errors during the inspection.
7 Water was used as an acoustic couplant in the scanning process. Similar data
8 acquisition setups are used in industry and has been employed for data collection on
9 large composite aerospace components previously [26].

10 Communication between the robotic interface and data acquisition was handled
11 through a custom LabVIEW interface. A custom binary file format (.rput) was
12 developed for this use case to efficiently store ultrasonic, robotic, and Force-Torque
13 data for each B-scan. The development of the binary file format resulted in an
14 approximate 2 times reduction in memory compared to previously employed CSV file
15 formats. The robotic information included the position at each frame which allowed
16 for accurate rasterisation of multiple scan passes. An accompanying Python script
17 decoded the custom binary file into usable NumPy arrays which allowed for efficient
18 post processing and ML development in Python. A block diagram of the setup is
19 illustrated in Figure 16.

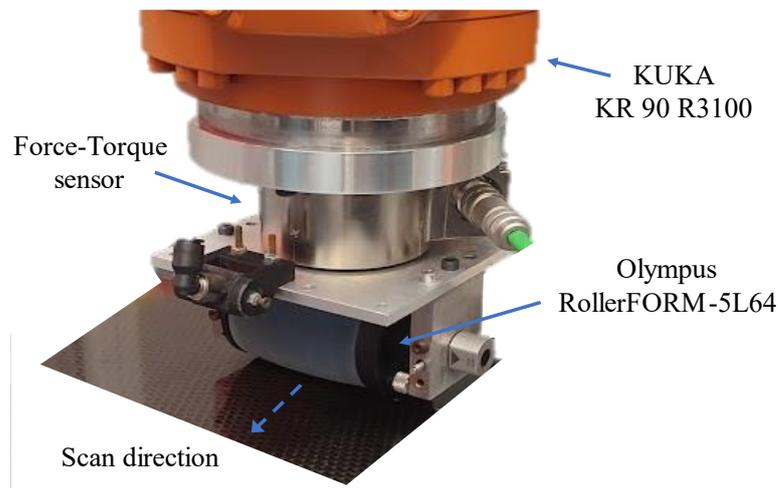


1

2 *Figure 16: Block Diagram illustration of the experimental data acquisition setup.*



3 *Figure 17: a) Overview of the experimental setup of KUKA KR90 and ultrasonic roller probe used for*
 4 *data acquisition. b) Close-up image of the experimental setup showing the assembly of the roller-*
 5 *probe and Force-Torque sensor as the robot end effector.*



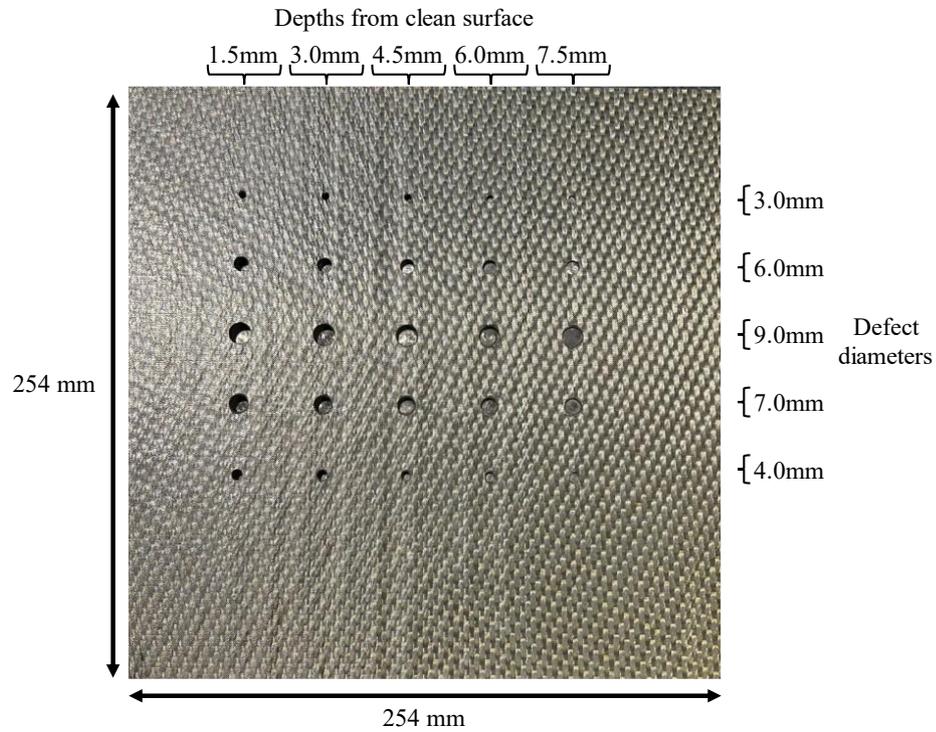
1

2 *Figure 18: Overview of the experimental scan setup of KUKA KR90, Force-Torque sensor, and*
 3 *ultrasonic roller probe used for data acquisition.*

4 Experimental ultrasonic data was acquired from CFRP samples, both with and without
 5 artificially introduced defects, to serve as test data. To imitate delamination defects,
 6 which are one of the most common defects in composites [101] and a significant life-
 7 limiting failure mode [126], Flat-Bottom Holes (FBHs) were drilled from the backside
 8 of the samples. Such defects are simple to produce post-cure and give similar responses
 9 to delamination's. In addition to this, their consistent known geometry allows for them
 10 to be used as references for defect sizing.

11 Prior to introducing defects, clean scans of each sample were taken to form defect-free
 12 datasets. The use of the same CFRP base sample ensured that the trained models
 13 learned defect-specific responses rather than the underlying properties of different
 14 composite samples. Two 254 x 254 x 8.6 mm (WxDxH) composite samples were
 15 provided by Spirit AeroSystems. The samples were all manufactured to the BAPS 260
 16 specification using a Resin Transfer Infusion Process, made using non-crimp fabric
 17 and Cycom 890 resin. The ply layers had a repeating lay-up pattern of 0, 45, -45, and
 18 90 degrees and a density of 1440 kg/m³. In the first sample 15 FBHs were drilled from

1 the backside to simulate defects. The defects were 3.0, 6.0 and 9.0 mm in diameter,
 2 with each individual defect diameter drilled to depths of 1.5, 3.0, 4.5, 6.0, 7.5 mm from
 3 the front surface. The different defect diameters were spaced 30 mm apart with
 4 different depth defects separated by 35 mm. In the second sample, 25 FBH were
 5 drilled to the same depths as the first sample but with additional defect sizes of 4.0 and
 6 7.0 mm as shown in Figure 19. All defects were manufactured to tolerances in depth
 7 of +/- 0.3 mm, and diameter of +/- 0.2 mm. A summary of the samples is provided in
 8 Table 3.



9
 10 *Figure 19: The composite test sample showing 25 FBHs.*

11 *Table 3: Summary of samples and their defects.*

Sample	Number of Defects	Diameters (mm +/- 0.2)
1	15 Flat-Bottom Holes	3.0, 6.0, 9.0
2	25 Flat-Bottom Holes	3.0, 4.0, 6.0, 7.0, 9.0

12

1 2.3 Signal Processing

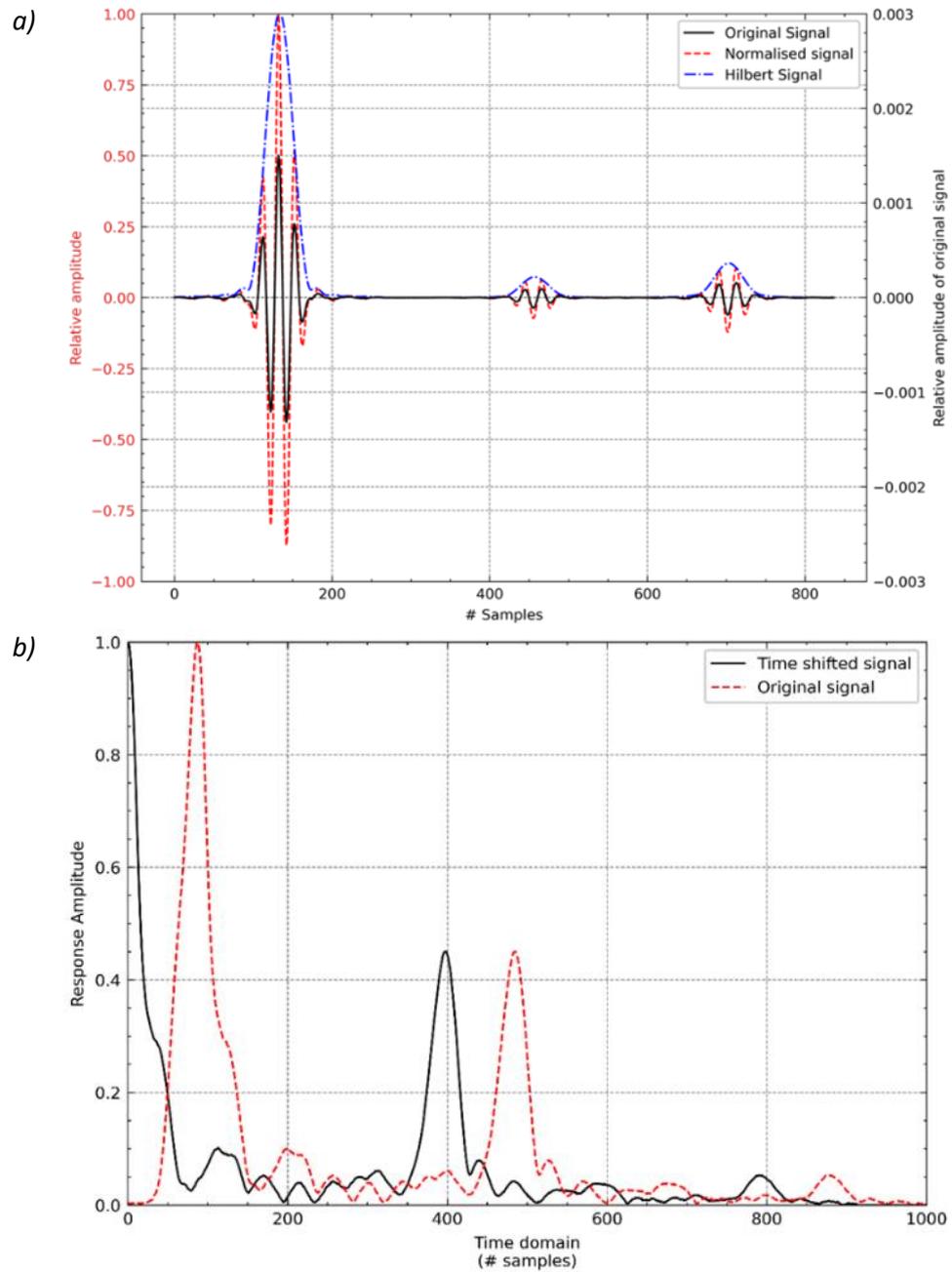
2 There are many different signal processing techniques that can be applied used to
3 increase the interpretability of the raw radio-frequency time series data collected from
4 each transducer (or group of transducers). Commonly when imaging, the envelope of
5 the signal is taken, this is often done using the Hilbert transform. The Hilbert transform
6 is used to obtain the analytical signal, the absolute of which gives an enveloped signal
7 which is useful for calculating the instantaneous response of a time series. This is
8 beneficial when imaging as it not only gives a positive response, but it also smooths
9 the signal, unlike simply rectifying the signal, where the absolute value of response is
10 taken. This approach is a standard signal processing technique used when generating
11 C-scan images from time series ultrasonic data [127]. The theory and mathematics
12 behind the Hilbert Transform are well documented, for further details please refer to
13 the relevant literature [128].

14 In many UT scenarios, the signals obtained yield strong responses from geometric
15 features, such as the front or back wall response. Including these features in the
16 imaging process can often mask out other useful information due to their high
17 amplitude response. To address this issue, gating is frequently employed in the time
18 domain to selectively remove these geometric features and establish a window that
19 focuses the imaging specifically on the area of interest.

20 Many different signal processing methods are used in UT. Examples of these include
21 the Wavelet transform; where a signal is decomposed into shifted and scaled versions
22 of a base function (or wavelet), the Fourier transform; where the signal is decomposed
23 into sine waves of specific frequencies, signal filtering; where a portion of the signal

1 is removed, this is often done in the frequency domain but can also be done in the time
2 domain to remove features not of interest, or enveloping; where the boundary of the
3 oscillating signal is extracted. The use of different pre-processing techniques is often
4 application dependent, for specific details on these please refer to the relevant literature
5 [127], [129], [130]. The work conducted in this thesis primarily involves the use of the
6 methods detailed in the following section.

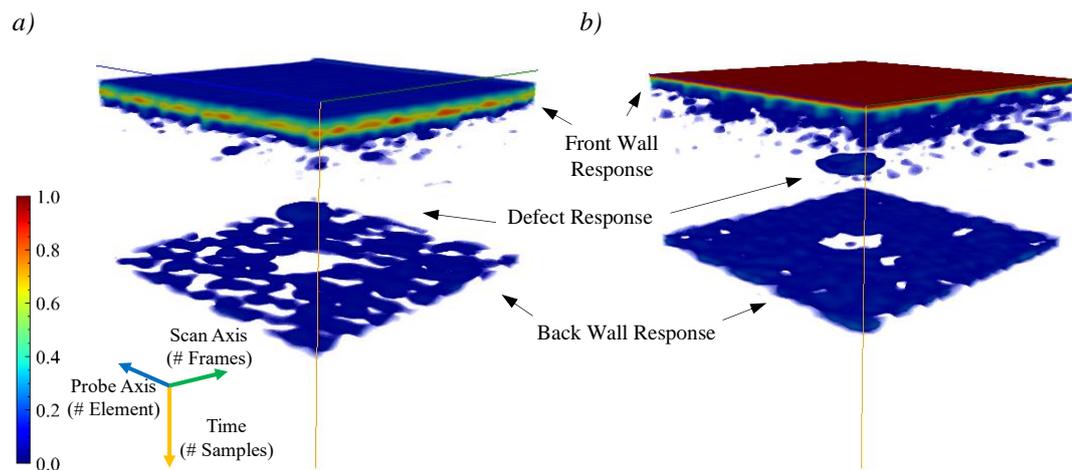
7 For the work conducted in this thesis, data collected (either experimentally or through
8 simulations) was in the form of radio frequency A-scans, also known as amplitude
9 scans. To make the raw radio frequency data more useful, different pre-processing
10 steps were used. Initially, the A-scans were centred at zero mean amplitude and
11 enveloped using the Hilbert transform previously discussed, as shown in Figure 20 (a).
12 Subsequently, each dataset was normalised between 0 and 1 by dividing by its
13 maximum absolute peak amplitude. Normalisation is not only a beneficial step in data
14 processing for ML training stability, but it also allowed for direct comparison of the
15 different datasets as amplitudes from the simulations are relative to each other and are
16 not reflective of experimental voltage values. In addition to facilitating dataset
17 comparability, the data preprocessing steps enabled the incorporation of domain-
18 knowledge into the analysis pipeline. For instance, prior understanding dictates that
19 the front wall response corresponds to the surface of the component and, under
20 standard inspection conditions, should consistently appear at the same point in time.
21 While machine learning models could theoretically learn this information given a
22 sufficiently large and representative dataset, embedding this knowledge during
23 preprocessing helps to eliminate noise in the data and assists the models in learning
24 effectively.



2 *Figure 20: a) Example of relative amplitude response from simulations, normalised signal, and*
 3 *Hilbert transform, applied to the original signal. b) Demonstration of how individual A-scans are time*
 4 *shifted to the front wall response.*

5 For chapter 3, which is focused on the interpretation of C-scan images, the front and
 6 back wall were gated out. For chapters 4 and 5, which focus on volumetric analysis,
 7 once the data was normalised, the offset in the time domain was compensated for by

1 aligning the peak front wall response to the origin. This made sure that features were
2 correctly aligned in the time domain and helped to account for any variability in the
3 acoustic path length between individual transducers and the surface of the sample.
4 Figure 20 (b) Shows how the time shifting was done for an individual A-scan with the
5 Hilbert transform applied. Figure 21 shows the effect of this on a complete ultrasonic
6 experimental volume.



7 *Figure 21: a) Volumetric data with Hilbert transform applied only. b) Volumetric data with time*
8 *shifting to the central response of the front wall peak. Both figures have been thresholded to remove*
9 *the lowest 10% of amplitudes to aid in visual clarity.*

1 2.4 Dataset Introduction

2 The resolution of the UT data in the array dimension was constrained by the element
3 pitch, and the scan width was restricted by the number of elements in the array. This
4 limited the inspection data to 64 voxels in the array dimension. To match this, 64 B-
5 scans were concatenated in the scan dimension to create cuboidal datasets, as seen in
6 Figure 21. As the element pitch was 0.8 mm, and the robotic-scanning speed was
7 regulated with the pulse repetition frequency to ensure a B-scan offset of 0.8 mm. This
8 enabled the generation of volumes with square voxels in the spatial domains, along
9 both the probe and scan directions. By utilising this approach, the work was able to
10 achieve a standardised volumetric resolution that was consistent throughout the
11 datasets. Since ultrasound values are just echo amplitude responses received from
12 within the inspected component by the array and presented in levels of voltage
13 response, the images/volumes were kept in single channel grayscale as any colours did
14 not have any physical significance.

15 2.5 Conclusion

16 This chapter provided an overview of ultrasonic technology, covering topics from
17 basic acoustic propagation to phased arrays and mechanized scanning. It also
18 underscored the growing use of composite materials in aerospace applications and
19 safety-critical components, which heightens the need for effective NDE to identify a
20 variety of potential manufacturing defects. UT is widely adopted in the industry for
21 inspecting aerospace composites due to its ease of use and reliable detection
22 capabilities across various defect types.

23

1 The chapter also introduced foundational concepts in AI, with an emphasis on data-
2 driven approaches like DL. A summary of the applications of these methods to NDE
3 was presented, along with a discussion of the specific challenges encountered in
4 integrating these techniques within the field.

5 Finally, the chapter outlined the experimental data collection methodology, and the
6 samples used throughout this thesis, detailing aspects such as signal processing and
7 dataset construction.

1 3 Synthetic Data: A Solution to Training Data Scarcity

2 3.1 Introduction

3 Despite the clear opportunity, ML has seen limited uptake in UT signal analysis,
4 particularly for composite components, which present a more challenging case with
5 additional structural noise compared to isotropic and homogeneous materials. A clear
6 barrier to research developments is the lack of training data [14]. This combined with
7 industrial questions over interpretability and compliance with standards has presented
8 challenges for the use of ML, and DL particularly. Modern manufacturing processes
9 aim to reduce the production of defects, meaning large volumes of real defect
10 responses are simply not available; especially ones that represent the full distribution
11 of defect classes and wide variability within these classes that are present from
12 inconsistencies in manufacturing. Furthermore, the manufacturing volumes of
13 aerospace components can be small, and stringent protocols for data protection of civil
14 and military components compounds the issue of data scarcity. Most commonly,
15 previous works have aimed to experimentally increase their datasets using
16 manufactured defects [32], [131], [132]. However, whilst these approaches can
17 demonstrate research concepts, they are unlikely to give UT responses that accurately
18 represent real-world responses especially not at the same variability seen within real
19 defects. Other authors have demonstrated success using simulated data developed
20 using FEA software to model defects and ray-based models to create Plane Wave
21 Capture, which uses a physics-based understanding of the wave propagation to
22 produce accurate responses based on bulk material properties [106]. However, this is
23 typically done for isotropic and homogenous steel samples which have very low
24 attenuation and noise, and have less modelling complexity compared to composites,

1 which are acoustically anisotropic and produce large amounts of UT wave attenuation
2 and scattering noise. Furthermore, this noise is often produced structurally from the
3 internal ply/fibre bundle interfaces of the composite material and is not random.
4 Therefore, neglecting this structural noise component and merely addressing the
5 random noise through the addition of randomly distributed noise to the datasets may
6 give unrealistic images or obscure defect responses. It is therefore important to
7 understand what gives rise to the complete noise, how this can be modelled, and how
8 this impacts our DL models. Most modern FEA software can account for ply
9 interactions, but it needs intensive material acoustic property characterisations,
10 modelling effort, and very long time-transient simulations. Therefore, composites are
11 often modelled using average bulk properties and not done at the individual ply level.
12 As an alternative to full FEA software, semi-analytical physics-based software has
13 been shown to produce experimentally accurate defect responses [133], [134]. This
14 software is much less computationally expensive than full FEA and can be used for
15 simulating composite responses based on bulk material properties [135]. This provides
16 a great opportunity to simulate vast amounts of defect responses with low
17 computational cost however, it does lack the complexities of structural noise response.
18 Synthetic datasets are widely used in ML to augment small training datasets [136] and
19 they offer a potential solution to the lack of defect data in UT. This chapter looks at
20 different novel methods of generating synthetic datasets from simulated data for
21 composite UT. These novel synthetic data generation methods are comparatively
22 evaluated on their experimental classification performance when used for training a
23 CNN. Hyperparameter Optimisation (HPO) is used to select an appropriate CNN
24 architecture that can represent the solution space for our task. GAN [137] are one of

1 the approaches investigated and have seen success in generating and augmenting
2 training data.

3 GANs are a class of generative ML models which rely on the relationship between two
4 networks: a generator and a discriminator. During training the two networks learn
5 together; the generator tries to produce data that the discriminator is unable to
6 distinguish from the ground truth, whilst the discriminator learns to better identify real
7 and generated data. GANs are often used to create or fill in images and as such have
8 been explored as a means for augmenting limited data sets for CNNs. To augment data,
9 GANs are often used to generate additional samples to further populate a distribution
10 of a particular target case, relying on the variability within the GAN to provide a
11 greater variability in training examples [119], [138], [139]. An additional benefit of
12 using synthetic data generated from GANs could be in the anonymisation of the
13 original dataset to address security or privacy concerns [140]. This could help to
14 alleviate concerns over data sharing of real NDE defects from industry.

15 The specific GAN used in this work to tackle a data shortage challenge for the first
16 time in the NDT domain is CycleGAN, which is a conditional GAN that has
17 demonstrated good results in unpaired image-to-image translation tasks [141]. This
18 GAN approach aims to combine NDT data generated from physics-based simulations
19 with GAN augmentation to create a dataset based upon physically accurate defect
20 responses that better resemble experimental data. The approach uses a modified
21 CycleGAN architecture to learn the mapping from simulated UT data to experimental
22 UT data. Specific, novel modifications, integrally an additional loss function, help to
23 encourage accurate defect signal reproduction whilst allowing for the addition of

1 experimental noise. With this approach, large quantities of highly varied simulated
2 defects can be produced in a relatively short time as compared to experiments or FEA,
3 and using the GANs mapping, produce large quantities of experimentally
4 representative synthetic data. The overall goal of the work in this chapter is to identify
5 the best methods for generating synthetic datasets in UT of composites to help unlock
6 the potential of DL in NDT applications.

7 This chapter provides details on how experimental and simulated UT testing is
8 gathered and processed into defective and non-defective image datasets (section 3.3).
9 In section 3.4, information is provided on the use of a CNN architecture for evaluation
10 of classification performance and details on the HPO method used for architecture
11 selection. Comparison is made with the experimental classification performance
12 between simulated and experimental data in section 3.5. The different methods of
13 synthetic data generation are then explored in section 3.6 with the effects on
14 classification performance evaluated. Section 3.7.1 provides a summary of the
15 classification results. Finally, section 3.7.2 introduces Grad-CAM as a method to help
16 with model interpretability when comparing synthetic data to experimental data and
17 discusses the full results of this work.

18 3.2 Simulated Data

19 A simulated dataset of the experimental test sample discussed in section 2.2 was
20 constructed using a semi-analytical, physics-based, commercial NDT simulation
21 software – CIVA [142]. Flaw interaction in CIVA is made up of three computation
22 stages: incoming transient ultrasonic field arising on the defect, field-to-flaw
23 interaction according to the Kirchhoff approximation, and prediction of the sensitivity

1 at reception using Auld's reciprocity theorem. The Kirchhoff approximation assumes
2 the wave does not propagate into the defect which is appropriate for the FBHs
3 modelled here [133], [143]. As the software adopts a semi-analytical approach, it
4 allowed for simulations to be completed with significantly reduced computational cost
5 (with an approximate 60x reduction in time) compared to FEA methods. Since the
6 focus of this work was the opportunity to produce large datasets for UT, this was a
7 significant benefit of the semi-analytical software approach.

8 CIVA simulation software is physics based, and has been widely used for commercial
9 UT simulation work, and experimentally validated for UT [133]. Therefore, there was
10 reasonable confidence that the modelling of wave propagation and its interaction with
11 defects were representative, producing reliable defect responses as well as being
12 computationally efficient. In addition, the simulated defect dimensions and positions
13 were readily controlled, allowing duplication of the exact experimental setup. This
14 allowed for efficient, complete annotations of the dataset to be generated at the point
15 of simulation, which opens further opportunities beyond classification, such as
16 segmentation where each pixel within the image is individually classified which could
17 be beneficial for defect sizing etc. A significant downside of using a semi-analytical
18 software as opposed to FEA is that the software was unable to model each distinct
19 composite layer response leading to differences between the simulations and measured
20 experimental responses, such as the lack of coherent scatter from ply interfaces.
21 However, in creation of the model, the individual layers were still constructed but were
22 only used to estimate equivalent homogeneous (and anisotropic) material properties.
23 A single ply layer was constructed and alternated with 0, 45, -45, and 90 degrees to
24 match the experimental sample as closely as possible. The resulting multilayer

1 structure was homogenised so that it was consistent with a homogeneous medium
2 having mechanical properties equivalent to those of the multi-ply composite. The
3 homogenisation gives good bulk propagation characteristics but removes structural
4 noise due to reflections from ply boundaries. This is a limitation of this simulation
5 method but is a necessary trade off against the computational cost of FEA when
6 producing large datasets. The fibre density was also set to 50 % to give the density
7 which best matched the experimental sample value of 1440 kg/m³. A parametric study
8 simulation was setup which used the composite bulk properties previously calculated
9 and varied the diameter and depth of defects. The study matched the experimental
10 setup with 3.0, 6.0 and 9.0 mm defects at depths of 1.5, 3.0, 4.5, 6.0, and 7.5 mm from
11 the surface. Both the front and backwall surface reflections were included in the model.
12 The full simulations for the complete dataset took less than 6 hours on a desktop
13 computer with a 24-Core 3.79 GHz CPU and 128 Gb of memory.

14 3.3 Image Generation

15 Once the data was normalised and the signal processing steps outlined in section 2.3
16 were applied. The data was truncated to remove the front and back wall echoes across
17 the full dataset. Then the maximum amplitudes were taken at varying depths of 5
18 samples in the time domain to produce C-scans (sampling rate of 100 MHz), refer to
19 Figure 8 (*C-scan depth gating*) for visualisation of the image extraction process. This
20 enabled direct comparison for multiple different response images to be generated for
21 each defect. From these C-scans, the images which represented a defect response were
22 collected. For the experimental samples, data was also collected from the reference
23 sample to obtain defect free C-scan images. In total this produced 334 defective
24 images from the experimental training sample, 150 defective images from the

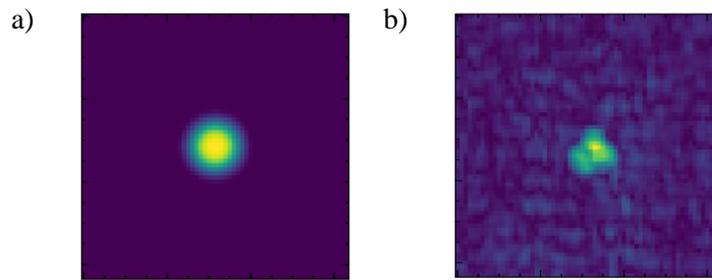
1 experimental test sample and 640 defect free images from the reference sample. This
 2 was split into 334 clean training images and the rest were used for testing. From the
 3 simulated dataset, 154 defective images were produced. Figure 22 shows how the
 4 simulated responses were significantly different from the experimental data. The
 5 simulated responses have far greater signal to noise ratio than the experimental
 6 responses and lacked the background response that is typically seen in experimental
 7 scans from the composite ply interactions, with a mean signal to noise ratio of over
 8 400 times the simulated defective datasets compared to the defective test dataset. A
 9 summary of the datasets generated from the experimental and simulated data is given
 10 in Table 4. For details on experimental samples and data acquisition, please refer to in
 11 section 2.2.

12 *Table 4: Summary of the datasets produced.*

Data source	Dataset	Number of images
Experimental test sample (15 Flat-Bottom Holes)	Defective test	150
Experimental train sample (25 Flat-Bottom Holes)	Defective train	334
Experimental reference sample	Clean test	148
	Clean train	334
Simulated experimental test sample (15 Flat-Bottom Holes)	Simulated defective	154

13

14



1 *Figure 22: example of simulated (a) and experimental (b) C-scan responses of a 9 mm diameter*
2 *FBHs.*

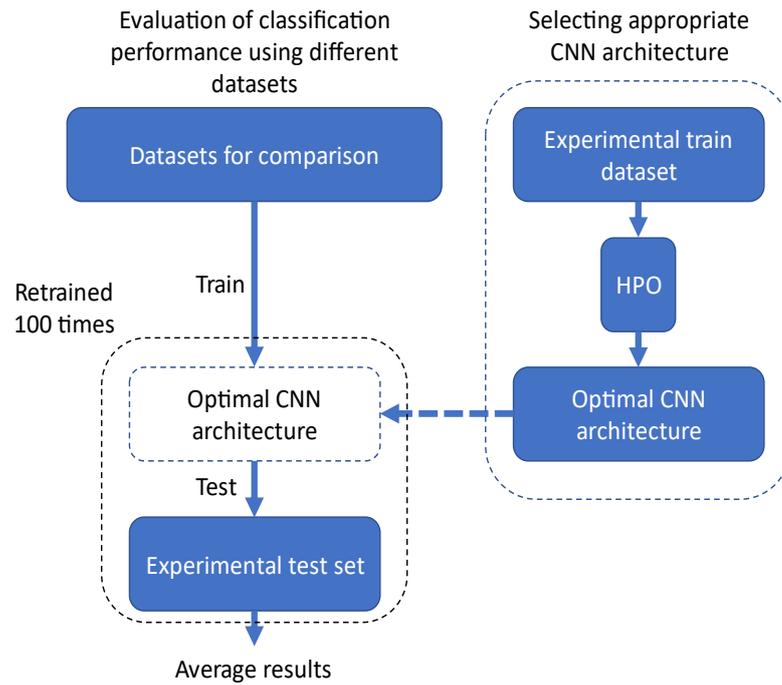
3 3.4 CNN Classification and Evaluation

4 3.4.1 Classification Evaluation with CNNs

5 The aim of this work is to evaluate different methods of modifying simulated data to
6 make them more effective at training Deep Learning models for experimental
7 classification tasks. It is therefore important that we evaluate our synthetic datasets
8 with respect to a classification metric. A CNN was used to evaluate and compare the
9 classification performance of different synthetic and experimental datasets. CNNs
10 have repeatedly demonstrated wide scale success in image classification and are
11 appropriate for this task [74].

12 Since the focus of this work was to compare synthetic datasets and not on optimal
13 classification accuracy, the CNN was kept constant for each dataset. Whilst the CNN
14 should be kept lightweight to reduce the computational cost of testing each synthetic
15 dataset, it was also important that the CNN had adequate complexity to learn the task.
16 To make sure the CNN had enough complexity to represent the solution space, a
17 genetic algorithm was deployed for hyperparameter optimisation (HPO) of a CNN
18 when trained on experimental data. A genetic algorithm is a heuristic search method
19 that mimics natural selection seen in biological evolution.

1 As the datasets used in the study are small, there was a degree of variability in the
 2 classification results. To negate this, when training the classifier, the CNN was re-
 3 trained for each synthetic dataset with a fresh initialisation 100 times and the average
 4 results were taken. Each CNN was evaluated on the same experimental dataset of 298
 5 images, made up from the experimental clean and defective test dataset. Figure 23
 6 shows the methodology used for classification evaluation.



7 *Figure 23: Flow diagram showing the process used for HPO of the CNN architecture and the use of*
 8 *the optimal architecture for classification evaluation.*

9 To quantitatively assess the performance of the classifiers, confusion matrices were
 10 generated, and precision, recall and F1 scores were calculated according to (7), (8) and
 11 (9).

$$Precision = \frac{TP}{(TP + FP)} \quad (6)$$

$$Recall = \frac{TP}{(TP + FN)} \quad (7)$$

$$F1 = \frac{(2 \times Precision \times Recall)}{(Precision + Recall)} \quad (8)$$

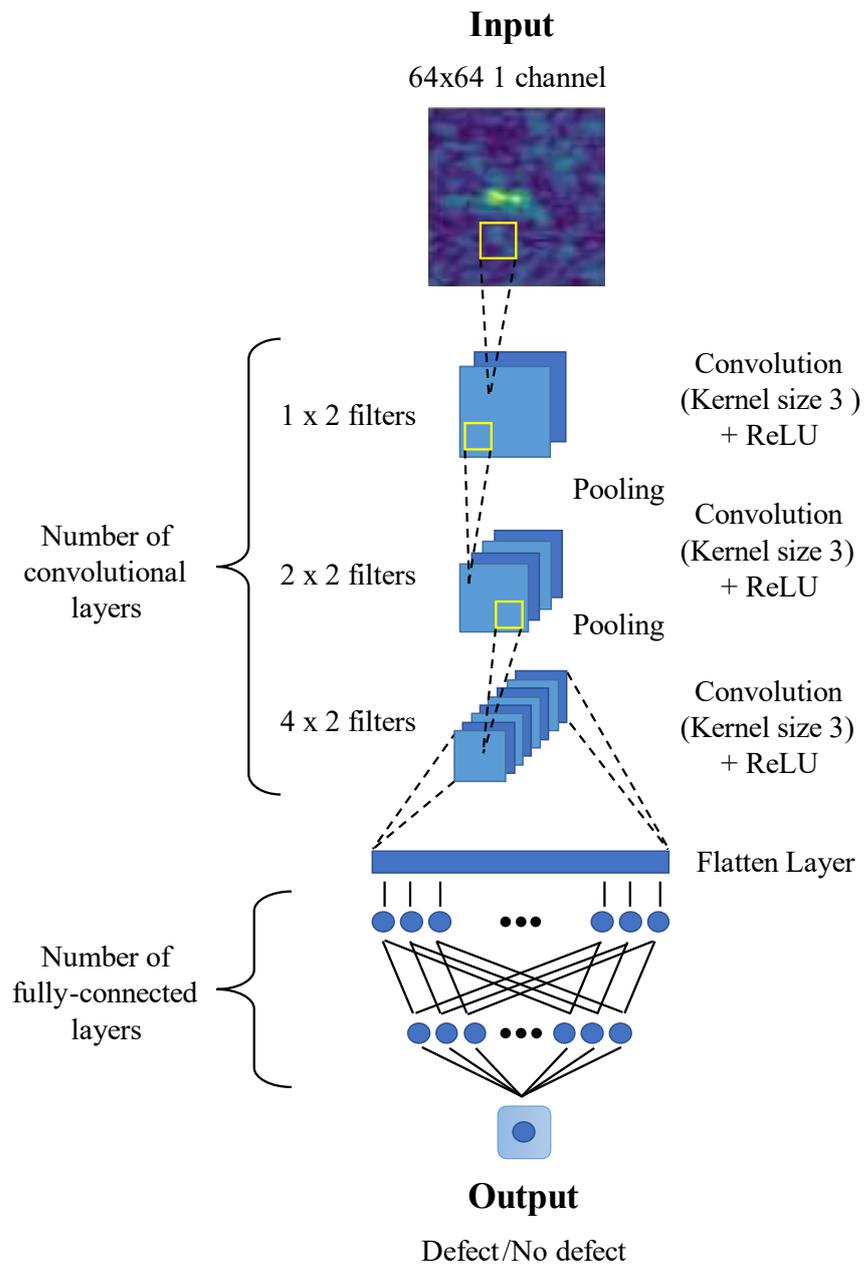
1

2 Where TP is true positive, FP is false positive, and FN is false negative, with positives
3 being the presence of a defect. Each result was individually averaged using a simple
4 mean across the 100 training cycles.

5 3.4.2 Hyperparameter Optimisation on Experimental Data

6 A genetic algorithm was used to perform HPO on the experimental training (defective
7 and clean) dataset to determine the parameters for the CNN. The model had at least 1
8 convolutional layer. Each convolutional layer had a fixed kernel size of 3 and used a
9 ReLU as the activation function [144] followed by max pooling with a kernel size of
10 2. The number of convolutional layers was parameterised with the number of filters
11 given by a constant out-channel ratio and the number of out channels from the previous
12 layer. The out-channel ratio was also parameterised. The network always had at least
13 one fully connected layer, from the flattened layer to the single output node, with a
14 sigmoid activation function for binary classification. There were a variable number of
15 fully connected layers and each hidden fully connected layer used ReLU activation.
16 The number of nodes on each hidden layer was equally distributed by dividing the
17 number of nodes in the flattened layer by the total number of layers and removing this
18 from the previous hidden layer each time. The optimised hyperparameters also
19 included batch size, early stop, learning rate, momentum, and number of epochs. The

- 1 values for the HPO variables are given in Table 5. Figure 24 shows an example of the
- 2 network with three convolutional layers, 2 hidden layers, and an out-filter ratio of 2.



3 *Figure 24: CNN architecture example with a convolutional channel ratio of 2.*

4

1 *Table 5: HPO variables and their range of values.*

Variable Parameter	Range
Number of fully connected layers	1 - 6
Number of convolutional layers	1 - 6
Channel ratio for convolutional layer filters	1 - 3
Batch size	16, 32, 64, 128, 256
Early stop	0 - 5
Learning rate	0.00001 – 0.5 (log scale)
Momentum	0 - 1
Number of epochs	100 - 500

2 The HPO was performed using the experimental train dataset, made up of 334 defect
3 images and the same number of defect free images from the clean train dataset. The
4 genetic algorithm used was a variant of Regularised Evolution (RE) [145], which was
5 adapted for continuous and integer valued hyperparameters. The algorithm was
6 initialised with a Population (P) of 128 configurations generated via a random search.
7 At each iteration RE sampled 5 configurations from the population, the model with the
8 highest evaluation score within this sample was selected and a new child configuration
9 was generated by mutating one of the parents hyperparameters. This child model is
10 then trained and prepended to the population with the ‘oldest’ model discarded. This
11 assisted in avoiding the system becoming trapped in local minima, as high performing
12 models relative to the population will be exploited for P iterations before being
13 discarded and allowing the process to explore new areas of the search space. This
14 method was run for 512 iterations. During each model evaluation, the dataset was
15 randomly subsampled without replacement with 80% of the dataset used for training
16 and 20% used for testing. The F1 score was calculated over 10 iterations of training
17 and testing data samples with the average F1 score used as the evaluation metric. The
18 optimum final network had an average F1 score of 0.978. The optimum

1 hyperparameters are outlined in Table 6. The network was implemented using the
 2 PyTorch framework [146].

3 *Table 6: Optimised hyperparameters used for CNN.*

Variable Parameter	Optimal value
Number of fully connected layers	1
Number of convolutional layers	3
Channel ratio for convolutional layer filters	3
Batch size	16
Early stop	1
Learning rate	0.014
Momentum	0.176
Number of epochs	264

4 3.5 Experimental and Simulated Data Classification

5 Performance

6 3.5.1 Experimental Results

7 For comparison to the synthetic datasets, a model was trained on the experimental test
 8 dataset and the same number of clean images sampled from the clean test dataset with
 9 a train/test split of 80% and 20% respectively. This gave a total of 60 test images. The
 10 averaged results across 100 training iterations, gave a mean model accuracy ($\pm \sigma$) of
 11 $89.8 \pm 9.8\%$, with average F1, precision and recall scores of 0.887 ± 0.112 , $0.974 \pm$
 12 0.135 and 0.826 ± 0.119 , respectively. The average confusion matrix for the
 13 experimentally trained model is given in Table 7.

14 *Table 7: Average confusion matrix across 100 training iterations for a CNN trained on experimental*
 15 *data.*

		Predicted	
		Defect	No defect
True	Defect	30.0	1.0
	No defect	5.1	23.9

1 3.5.2 Simulated Results

2 A model was also trained on the simulated, unmodified defect response data and the
3 same real defect free images generated from the defective test sample which were used
4 for the experimental results. This was made up of 154 simulated defect images and
5 154 real defect free images sampled from the clean train dataset. After 100 training
6 iterations, the model gave an average accuracy ($\pm \sigma$) of $62.8 \pm 4.9\%$, with average F1,
7 precision and recall scores of 0.394 ± 0.109 , 1.00 ± 0.00 and 0.252 ± 0.100 ,
8 respectively. The average confusion matrix for the model trained on simulated data is
9 given in Table 8.

10 *Table 8: Average confusion matrix across 100 training iterations for a CNN trained on simulated*
11 *data.*

		Predicted	
		Defect	No defect
True	Defect	150.0	0.0
	No defect	119.7	37.3

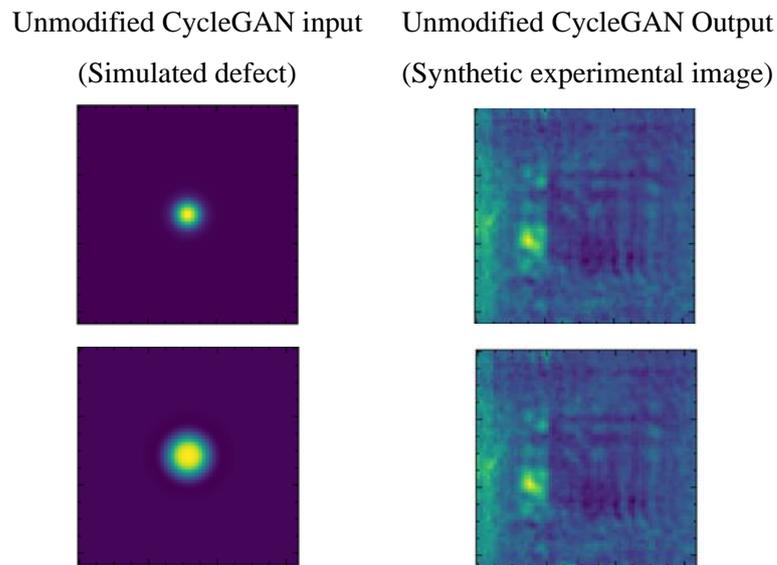
12 3.6 Methods of Synthetic Noise Generation

13 In this work four separate methods were explored to map simulated data to more
14 experimentally representative synthetic datasets by adding noise. The first approach
15 uses a modified CycleGAN to learn the mapping between simulated and experimental
16 data. The second approach aims to utilise the fact that clean ultrasonic images are
17 comparatively much more available than defect data, by combining both real clean
18 images and defect simulations. The final two approaches studied the noise profiles
19 seen in the experimental data and attempted to simulate these at both the C-scan image
20 level and the individual A-scan level.

1 3.6.1 Approach 1: Simulated to Experimental Domain Mapping with 2 CycleGAN

3 To learn the mapping between simulated and experimental data, an image-to-image
4 translation GAN was used. CycleGAN was chosen as it has shown promising results
5 in unpaired image-to-image translation, and works particularly well for style transfer
6 tasks which this application is similar to [141]. Unlike conventional GANs that rely
7 on adversarial loss alone, CycleGAN introduces an additional cyclic loss, which
8 enforces reconstruction for a full cycle; from source domain to target domain and
9 crucially back to source domain. The addition of cyclic loss removes the necessity for
10 paired training data, allowing for unsupervised domain transfer. Not requiring paired
11 images in training was a significant advantage as it provided greater freedom in the
12 images used in training. Furthermore, from an NDT perspective, if this approach was
13 extended to naturally occurring defects, it would be impossible to accurately simulate
14 the complexity of naturally occurring experimental defect responses to produce a
15 completely paired dataset.

16 Implementing the standard CycleGAN directly with the parameters given in the
17 original paper[141], was unable to accurately reproduce ultrasonic images with the
18 simulated defect responses present. Furthermore, the generated images suffered from
19 significant mode collapse. Mode collapse occurs when the generator repeatedly
20 outputs a single type of image, due to finding one image that is successful in fooling
21 the discriminator. Figure 25 shows an example of this, where different input simulated
22 defect responses produce the same output. The original implementation was done in
23 Pytorch and was trained for 200 epochs, with a batch size of 4, 6 residual blocks, and
24 an identity loss of 5 (half the cycle consistency loss).



1 *Figure 25: Example images of initial CycleGAN outputs.*

2 *3.6.1.1 CycleGAN Modifications– Mid-Cycle Activation Map*

3 It has been demonstrated that adjusting the loss function of CycleGAN can improve
 4 performance for specific tasks [147]. To improve the performance of the original
 5 CycleGAN [141] for this task a variety of adjustments were made, with the most
 6 significant being the introduction of a mid-cycle activation map loss.

7 My model contains two mapping functions $G_{\text{Experimental}}$ (G_{exp}): Simulated \rightarrow
 8 Experimental and $G_{\text{Simulated}}$ (G_{sim}): Experimental \rightarrow Simulated and associated
 9 adversarial discriminators $D_{\text{Experimental}}$ (D_{exp}) and $D_{\text{Simulated}}$ (D_{sim}). $D_{\text{Experimental}}$
 10 encourages $G_{\text{Experimental}}$ to translate experimental images into outputs indistinguishable
 11 from real experimental images, and vice versa for $D_{\text{Simulated}}$ and $G_{\text{Simulated}}$. Both cycles
 12 include the cycle consistency loss that was introduced in the original paper Figure 27
 13 (b, c). To further encourage accurate defect reproduction, I introduce a mid-cycle
 14 activation map loss for the simulated image cycle Figure 27 (b).

1 The mid-cycle activation map loss aimed to give the algorithm freedom to alter the
2 noise profile whilst retaining constraint over the original defect response. The need for
3 this was clear from the original implementation as the defect response can easily be
4 washed out (Figure 25). To do this, the simulated input image was used to generate an
5 activation map. This activation map was a normalised version of the original simulated
6 input image to a range of 0 and 1. The simulated responses allowed for this unique
7 implementation as the background responses were uniform. By normalising the
8 activation map, the effect of background response was zeroed, and only inaccurate
9 reconstructions of defect responses were punished, whilst maintaining even weak
10 defect responses. Next, a scale factor was calculated to allow for adjustments of defect
11 size. This was calculated by taking all non-zero values (defect response) from the
12 activation map and dividing by the total image area. The L1 unreduced absolute error
13 between the generated image and the simulated image was then calculated. The
14 activation map was then applied to focus the loss to the defect response and minimise
15 the loss from the noise. This new loss map was then divided by the scale factor
16 previously calculated from the activation map. This means that the loss function is
17 indiscriminate of defect size and does not punish larger defects more significantly than
18 smaller defects. Finally, the mean was taken to get the reduced value, which was fed
19 into the combined generator loss function given by (9). Figure 26 demonstrates this
20 process with an example image.

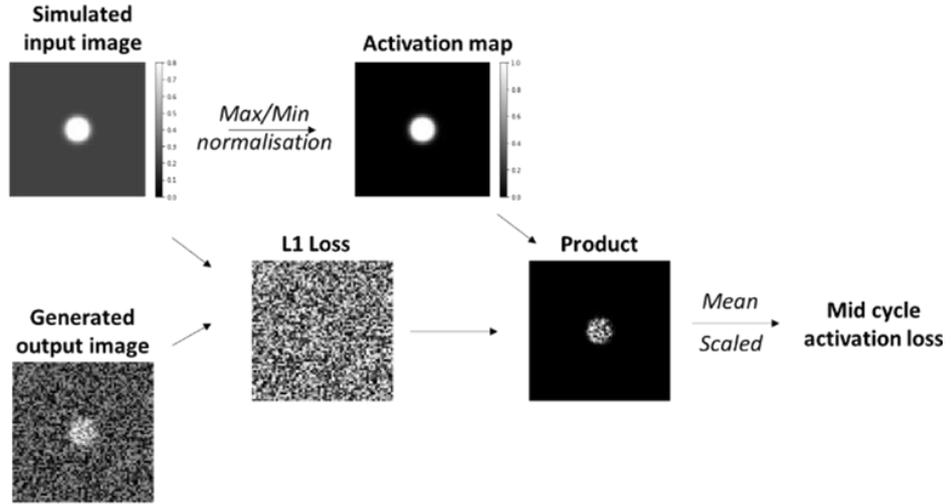
21

$$E_{sim \sim P(sim)} \left[\frac{L_{activmap}(G_{exp})}{K_{scale\ factor}} \times M_{activation\ map\ [0 \rightarrow 1]} \right] \quad (9)$$

1

$$L_{total}(G_{exp}, G_{sim}, D_{exp}, D_{sim}) = L_{GAN}(G_{exp}) + \frac{2}{3}L_{GAN}(D_{exp}) + L_{GAN}(G_{sim}) + \frac{1}{3}L_{GAN}(D_{sim}) + \frac{\lambda}{3}(2L_{cyc}(G_{exp}) + L_{cyc}(G_{sim})) + 2\lambda L_{activmap}(G_{exp}) \quad (10)$$

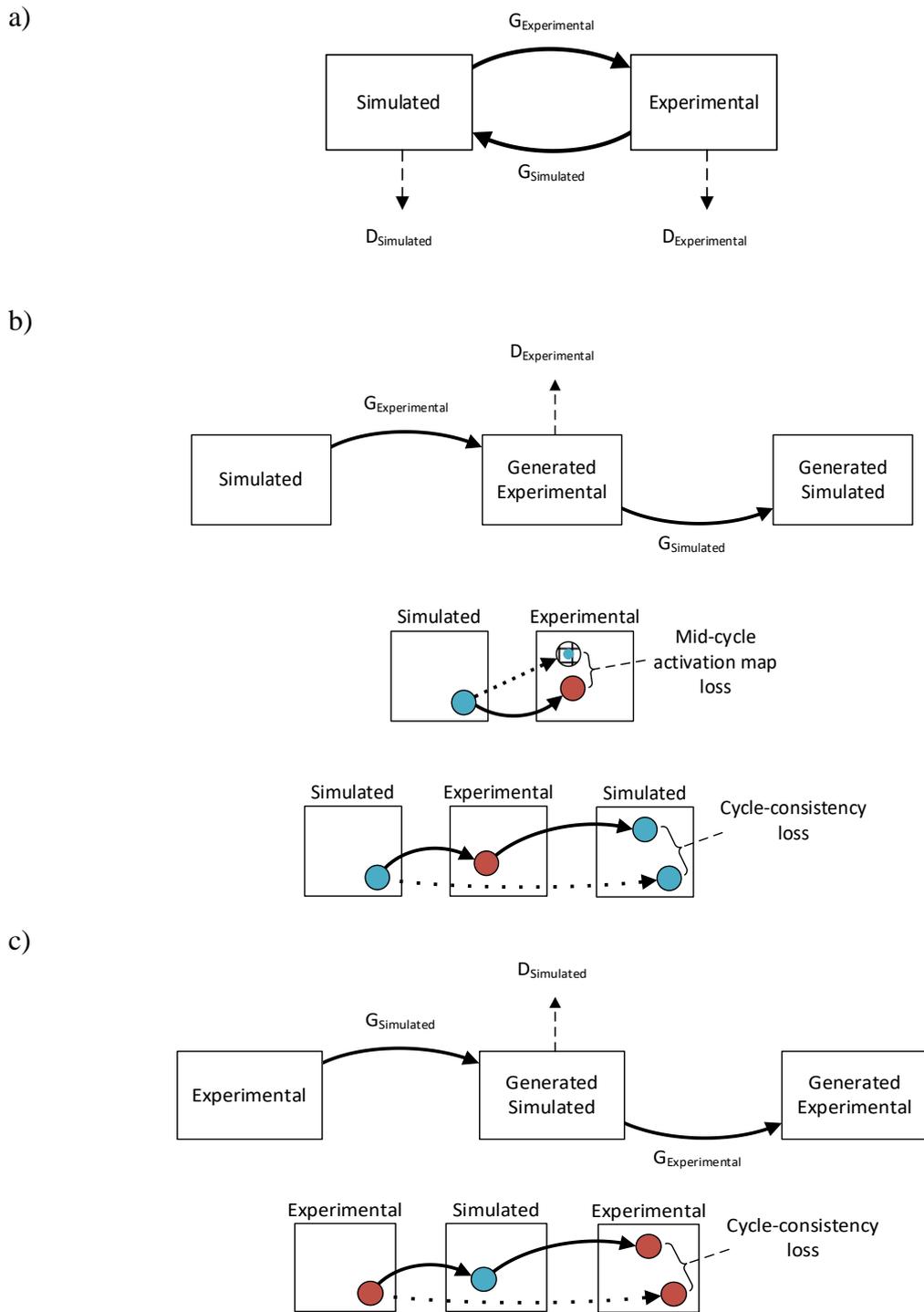
2 Where L_{total} is the total loss, L_{GAN} and L_{cyc} represent the adversarial loss and cyclic
 3 loss given in the original paper [141], $E_{sim \sim P(sim)}$ represents the expectation over the
 4 batch of simulated samples, $M_{activation\ map\ [0 \rightarrow 1]}$ is the normalised activation map, λ
 5 is a coefficient to balance the relative importance of each loss function during training.



6

7 Figure 26: Diagram showing how an example mid-cycle activation map loss is generated.

8



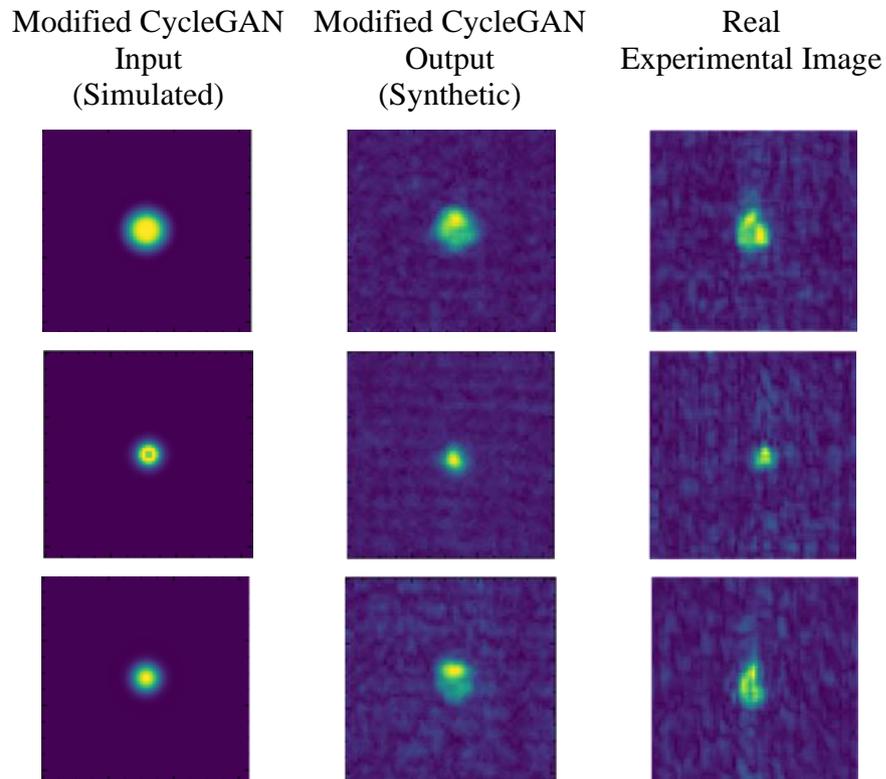
1 *Figure 27: a) The model contains two mapping functions $G_{\text{Experimental}}: \text{Simulated} \rightarrow \text{Experimental}$ and*
 2 *$G_{\text{Simulated}}: \text{Experimental} \rightarrow \text{Simulated}$, to transfer between the respective domains and the associated*
 3 *adversarial discriminators $D_{\text{Experimental}}$ and $D_{\text{Simulated}}$. b) When completing the full cycle from the*
 4 *simulated domain, the mid-cycle loss is added along with the cycle loss. c) When completing a cycle*
 5 *beginning in the experimental domain, the cycle loss is solely used as the mid-cycle loss is not*
 6 *calculable for the simulated domain.*

1 The mid-cycle activation map loss is only applied in the direction going from
2 simulated responses to generated experimental images, as it relies on the clean defect
3 response of simulated images. This is demonstrated by Figure 27 (b, c). For better
4 images in this task, the cycle loss was also adjusted to give twice the weighting of the
5 simulated input cycle compared to the experimental cycle, whilst the discriminator loss
6 for identification of experimental images was weighted twice as much as the
7 discriminator for simulated images. This was done to further remove restrictions on
8 noise generation and further encourage accurate defect response, whilst focusing on
9 generation of experimental images over simulated images. The cycle loss coefficient
10 (λ , equation (10)) was set to 100, with the mid cycle activation loss set to double the
11 cycle loss. To further improve the results, the CycleGAN model used was adjusted
12 from the original implementation [141] to perform better on the lower resolution 64x64
13 ultrasound images, by optimising the size of the first generator convolutional layers to
14 3x3 instead of 7x7, with 6 residual blocks used. The model was trained from scratch
15 with a learning rate of 0.0002 which decayed linearly after 100 epochs to zero for the
16 remaining training. For training, the GAN used the experimental defective train dataset
17 of 334 images, and the simulated defective dataset of 154 images. The GAN model
18 was trained over 2300 epochs using a batch size of 128 using an NVIDIA GeForce
19 RTX 3090 and took less than 8 hours to train. All other parameters were unmodified
20 from the original paper [141]. The GAN model was created using the Pytorch
21 framework.

22 Once trained, the learnt mapping from the GAN was used to convert the original 154
23 simulated images to a new synthetic dataset of defective images. The synthetic dataset
24 produced high quality ultrasonic amplitude images which are visually comparable to

1 experimentally obtained images, examples of images generated from their
2 corresponding simulated input are shown in Figure 28.

3



4 *Figure 28: Example of synthetic generated images from their corresponding simulated defect input,*
5 *along with real experimental images for comparison.*

6 *3.6.1.2 Classification Results*

7 Training the CNN with GAN generated synthetic dataset and an equal number of clean
8 images sampled from the clean train set, had a significant increase in classification
9 performance compared to unprocessed simulated data when tested on the experimental
10 clean and defective test datasets of 298 total images. After 100 training iterations, the
11 model gave an average accuracy ($\pm \sigma$) of $87.0 \pm 11.5\%$, (up from 62.8%) with average
12 F1, precision and recall scores of 0.837 ± 0.204 , 0.926 ± 0.231 and 0.775 ± 0.200
13 respectively. The average confusion matrix for the model is given in Table 9.

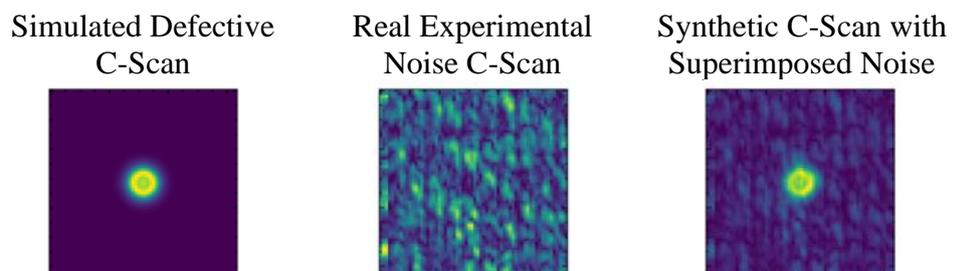
1 *Table 9: Average confusion matrix across 100 training iterations for a CNN trained on GAN*
 2 *generated synthetic data.*

		Predicted	
		Defect	No defect
True	Defect	144.6	5.4
	No defect	33.3	114.8

3

4 3.6.2 Approach 2: Experimental C-Scan Noise Superposition

5 Out of the 334 clean experimental C-scan images from the clean train dataset, 154
 6 were randomly sampled to match the size of the simulated dataset. The simulated
 7 defect images were then combined with the real noise images by summation at an
 8 individual pixel level. To not exceed the normalised upper value limit of 1, if a pixel
 9 value exceeded 1 due to the addition of noise, it was clipped to remain within the limit.
 10 This was done instead of re-normalising the dataset as this would have reduced the
 11 noise distribution from the experimental data. From the new dataset, the images where
 12 the noise was greater than the signal were removed. This left 83 final images. An
 13 example of this is demonstrated in Figure 29.



14 *Figure 29: Example images showing the combination of real noise and simulated defect responses.*

15 A considerable downside of the real noise approach is that it is not a fully simulated
 16 approach. This restricts its ability to scale as it requires an equal number of clean
 17 experimental images as simulated images. However, the experimental data required is
 18 from defect-free images which are more accessible and considerably easier to acquire

1 than real defect responses. The computational complexity of scaling this approach to
 2 a large number of images would be low. Therefore, if adequate clean images were
 3 available this technique could be used to produce a large dataset.

4 *3.6.2.1 Classification Results*

5 Training the CNN with the experimental noise synthetic dataset and an equal number
 6 of clean images sampled from the clean training set had a significant increase in
 7 classification performance when tested on the experimental clean and defective test
 8 datasets compared to the simulated data but was unable to match the results from the
 9 GAN generated dataset. After 100 training iterations, the model gave an average
 10 accuracy ($\pm \sigma$) of $77.4 \pm 7.8\%$, with average F1, precision and recall scores of $0.688 \pm$
 11 0.179 , 0.950 ± 0.218 and 0.545 ± 0.158 respectively. The average confusion matrix
 12 for the model is given in Table 10.

13 *Table 10: Average confusion matrix across 100 training iterations for a CNN trained on real noise*
 14 *data.*

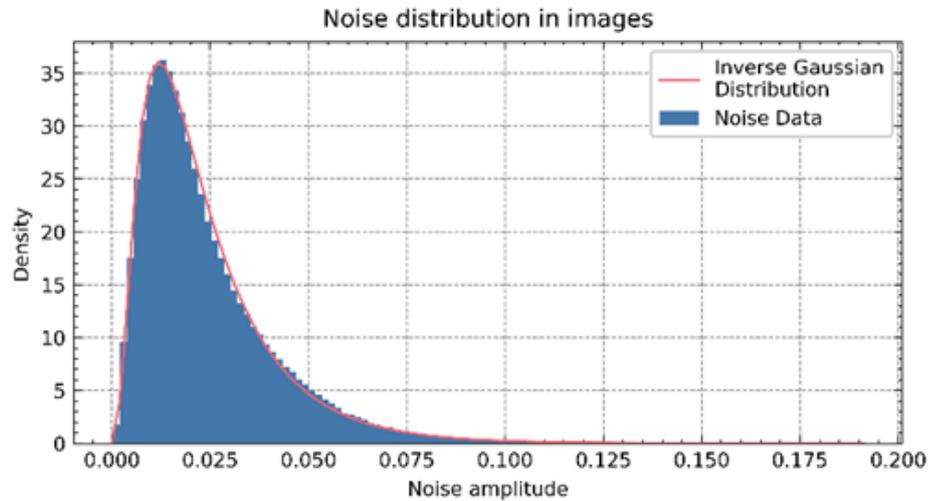
		Predicted	
		Defect	No defect
True	Defect	150.0	0.0
	No defect	67.3	80.7

15

16 *3.6.3 Approach 3: Simulated C-Scan Noise*

17 To reduce the experimental demand of the real noise superposition approach requiring
 18 a unique experimental image for each simulation, a study was conducted to understand
 19 if it was possible to fully simulate the experimental noise profile. To do this, the noise
 20 distribution from the clean experimental C-scan images of the defect free sample were
 21 analysed by plotting a histogram. It can be seen from Figure 30 that this noise profile

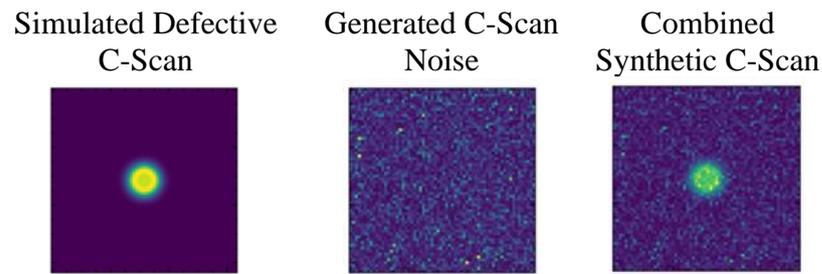
1 is well aligned with an inverse gaussian distribution given by μ 0.410, loc -0.003 and
2 scale of 0.066, the distribution was implemented using SciPy, with details of the
3 distribution given by [148].



4 *Figure 30: Density histogram showing the distribution of data from the clean sample.*

5 The simulated defect images were then combined with a noise pattern which was
6 randomly generated for each image from an inverse gaussian distribution, Figure 30,
7 with the previously determined parameters. The images were combined by summation
8 at an individual pixel level. As per the real noise method, to not exceed the normalised
9 upper value limit of 1, if a pixel value exceeded 1 it was clipped to remain within the
10 limit. From the new synthetic dataset, the images where the noise was greater than the
11 signal were removed, and we were left with 80 C-scan final images. An example of
12 this is demonstrated in image Figure 31.

13



1 *Figure 31: Example images showing the combination of C-scan simulated noise and simulated defect*
 2 *responses.*

3 The implementation of C-scan noise at scale would be considerably easier than the real
 4 noise approach. This is as fully simulating the noise profile from an appropriate
 5 experimental distribution requires little additional experimental data acquisition after
 6 a suitable population has been sampled. Furthermore, the computational complexity
 7 of this implementation is as efficient as the real noise approach and could scale well
 8 to produce a large dataset. Whilst it benefits from simplicity, the approach could be
 9 extended to account for local correlations between pixels, to more accurately simulate
 10 local relationships resulting from the composite ply structure.

11 *3.6.3.1 Classification results*

12 Training the CNN with the C-scan noise synthetic dataset and an equal number of
 13 clean images sampled from the clean training set produced poorer results than the
 14 superimposed real noise dataset but still improved significantly in classification
 15 performance over the simulated dataset when tested on the experimental clean and
 16 defective test datasets. After 100 training iterations, the model gave an average
 17 accuracy ($\pm \sigma$) of $74.3 \pm 8.1\%$, with average F1, precision and recall scores of $0.629 \pm$
 18 0.195 , 0.930 ± 0.255 and 0.482 ± 0.164 respectively. The average confusion matrix
 19 for the model is given in Table 11.

20

1 *Table 11: Average confusion matrix across 100 training iterations for CNN trained on simulated C-*
 2 *scan noise data.*

		Predicted	
		Defect	No defect
True	Defect	150.0	0.0
	No defect	76.7	71.3

3

4 3.6.4 Approach 4: Simulated Ultrasonic A-Scan Noise

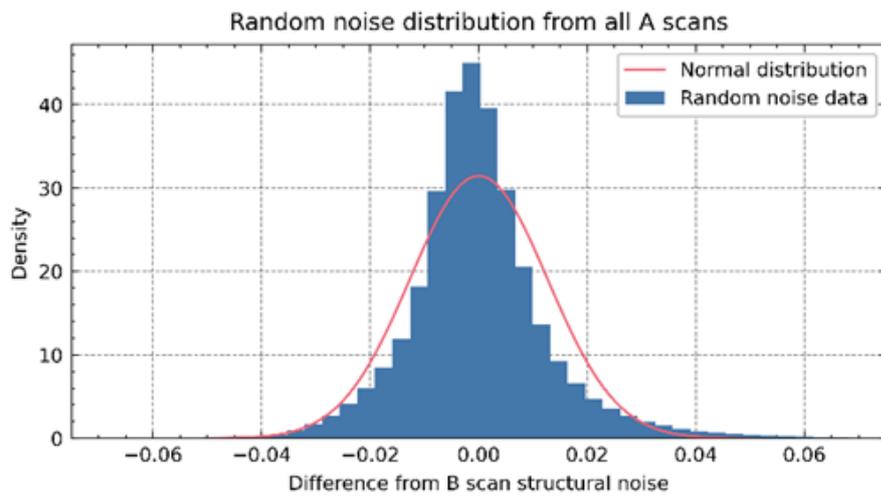
5 An approach of fully generating a simulated noise profile at an A-scan level was also
 6 investigated which is better aligned to how noise occurs from the physical response of
 7 ultrasonic signals. For each individual time trace signal, the complete noise profile is
 8 composed of both structured noise and random noise. Structured noise consists of
 9 physically accurate responses, just not from a known feature. These are likely due to
 10 the interaction of different composite plies and the component geometry with the
 11 propagated ultrasonic waves. Whereas random noise is independent of the samples
 12 structure and could be due to random electrical noise for example.

13 It was assumed that for a given B-scan, the structural noise profile will remain
 14 constant, as for a given B-scan the ultrasonic wave and ply layer interactions and
 15 therefore backscattering noise should be similar. Therefore, at a B-scan level, it is
 16 possible to remove most of the random noise by mean averaging the individual A-
 17 scans together at each sample interval leaving the structural noise component. For each
 18 A-scan in each B-scan, it is then possible to work out the random noise component
 19 from the differences between each A-scan and the structural noise component on a per
 20 sample basis. These combined differences can be plotted on a histogram to represent
 21 the random noise population of a B-scan. This process was completed for each
 22 individual B-scan. The random noise profiles were combined to give a greater number

1 of samples for the distribution. From Figure 32, it can be seen that this distribution is
2 approximated by a normal distribution with 0.000 mean and a standard deviation of
3 0.013.

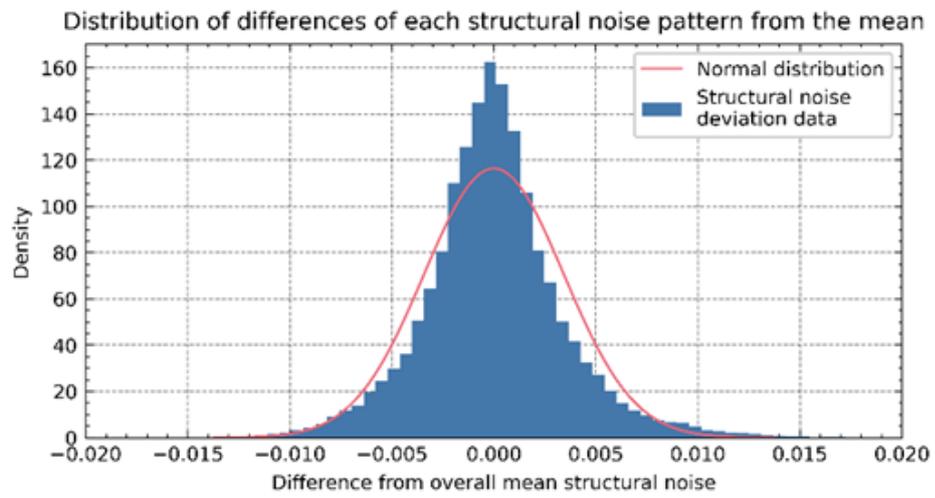
4 To learn the variation of the structural noise components across B-scans, the average
5 B-scan structural noise was first calculated by averaging each individual B-scan noise
6 profile on a mean sample basis. The difference between the mean and each individual
7 B-scan structural noise profile was calculated on a per sample basis and again plotted
8 on a histogram (Figure 33). This can be approximated by a normal distribution with
9 mean 0.000 and standard deviation 0.003.

10 To generate a new noise pattern for a B-scan, a new structural noise pattern was
11 generated by taking the overall mean structural noise pattern and adding variation
12 based on the normal distribution previously calculated. To make this signal more
13 representative of the Hilbert transformed A-scan data, a Savitzky–Golay filter [149]
14 was applied to smooth the data (Figure 34). Afterwards, a random noise profile was
15 added to the generated A-scan baseline signal, following the previously determined
16 normal distribution for each A-scan in Figure 32. Figure 35 helps to illustrate this
17 process at A-scan and B-scan levels. The simulated responses were then combined
18 with the generated combined noise profiles using a per sample summation. As per
19 previous methods, to not exceed the normalised upper amplitude value limit of 1, pixel
20 values exceeding 1 were clipped to remain within the limit. From the new dataset, the
21 images where the noise was greater than the signal was removed resulting in 126 C-
22 scan final images. An example of the final images is demonstrated in Figure 36.



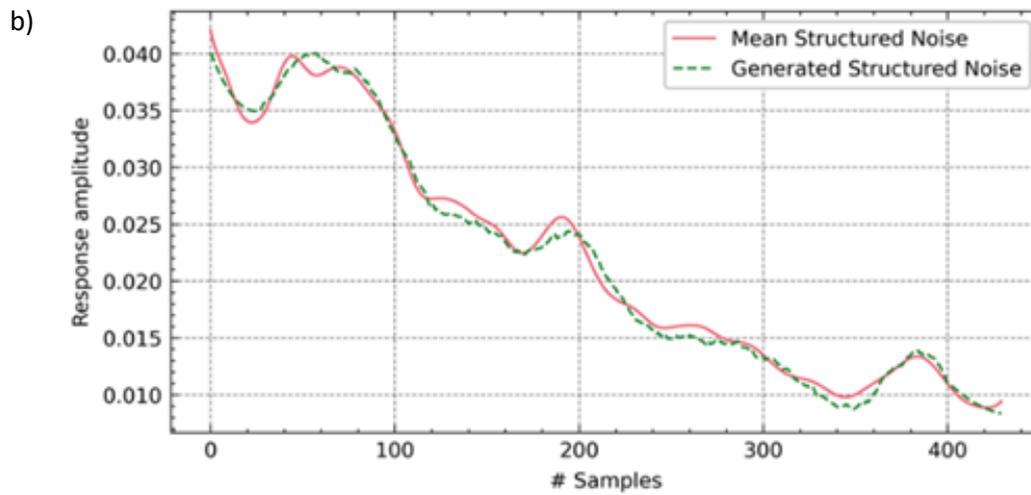
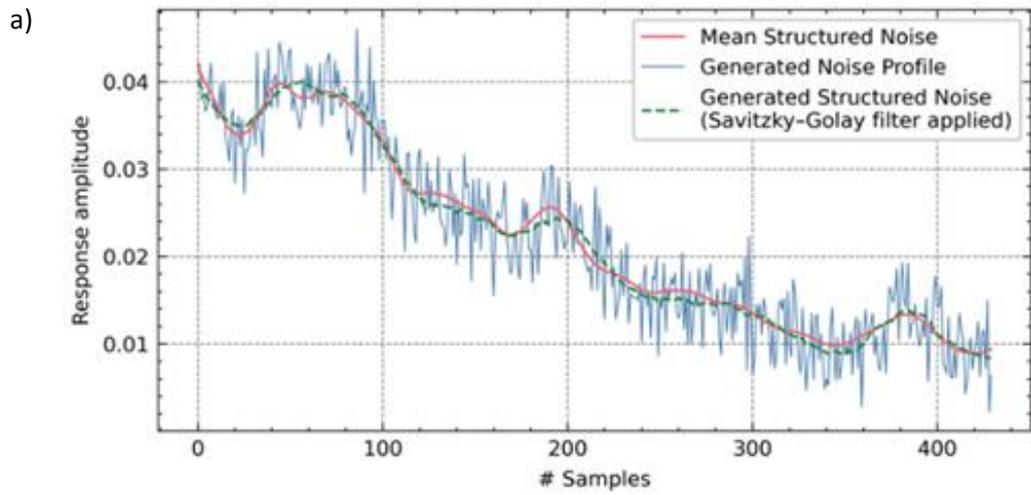
1

2 *Figure 32: Density histogram showing the random noise distribution from the total A-scans.*



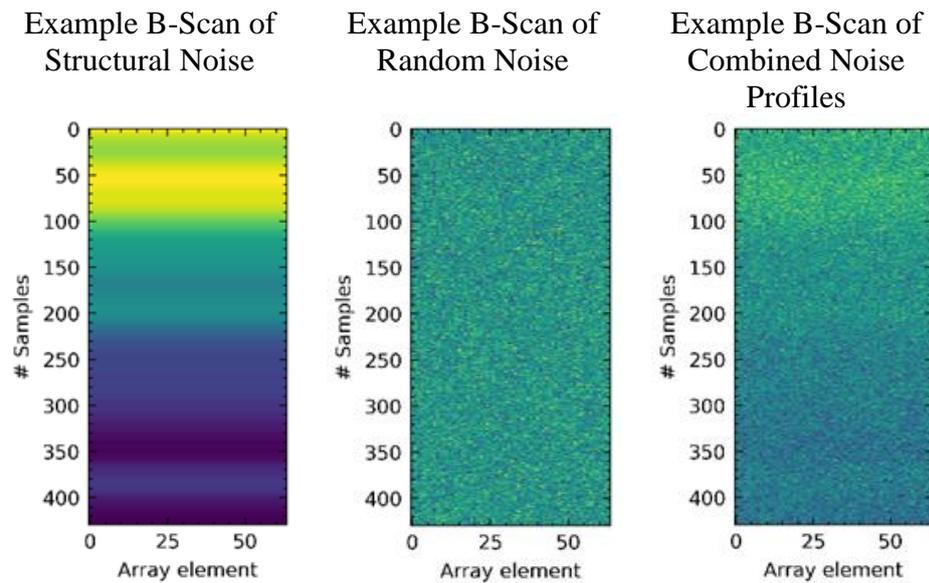
3

4 *Figure 33: Density histogram showing the distribution of deviation for structural noise from the mean*
 5 *structural noise pattern.*

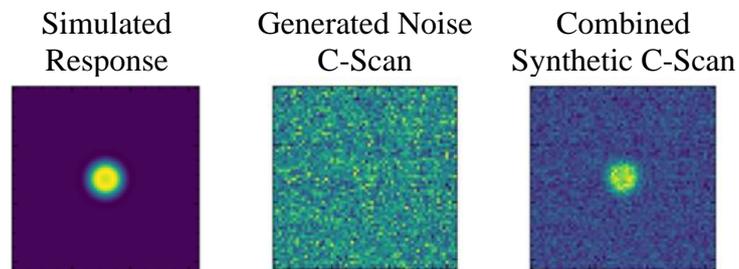


1 *Figure 34: a) An example of how a structural noise profile is generated from the mean. b) A cleaner*
 2 *example of the final generated noise profile.*

3



1 *Figure 35: An example of how structural and random noise profiles are combined at a B-scan level.*



2 *Figure 36: Example images showing the combination of A-scan simulated noise and simulated defect*
 3 *responses.*

4 Whilst implementing the A-scan noise profile does require experimental analysis and
 5 characterisation, the application to simulated data is a fully simulated approach. In
 6 addition, by adding noise at an A-scan level, it allows for the potential of three-
 7 dimensional volumetric analysis, or analysis of B-scan images, which is not possible
 8 with any of the other methods. However, it requires a greater level of analysis
 9 compared to the C-scan level noise method before implementation. Furthermore, as
 10 the generation of the noise pattern is required on a per B-scan level, an additional
 11 computational step is required to cover the number of B-scans. This is therefore less
 12 computationally efficient than both the real noise and C-scan noise implementation.

1 3.6.4.1 Classification results

2 Training the CNN with the A-scan noise synthetic dataset and an equal number of
3 clean images sampled from the clean training set, gave the second-best classification
4 performance when tested on the experimental test datasets after the GAN generated
5 dataset. After 100 training iterations, the model gave an average accuracy ($\pm \sigma$) of 80.0
6 $\pm 6.2\%$, with average F1, precision and recall scores of 0.738 ± 0.141 , 0.970 ± 0.171
7 and 0.598 ± 0.124 respectively. The average confusion matrix for the model is given
8 in Table 12.

9 *Table 12: Average confusion matrix across 100 training iterations for a CNN trained on simulated A-*
10 *scan noise data.*

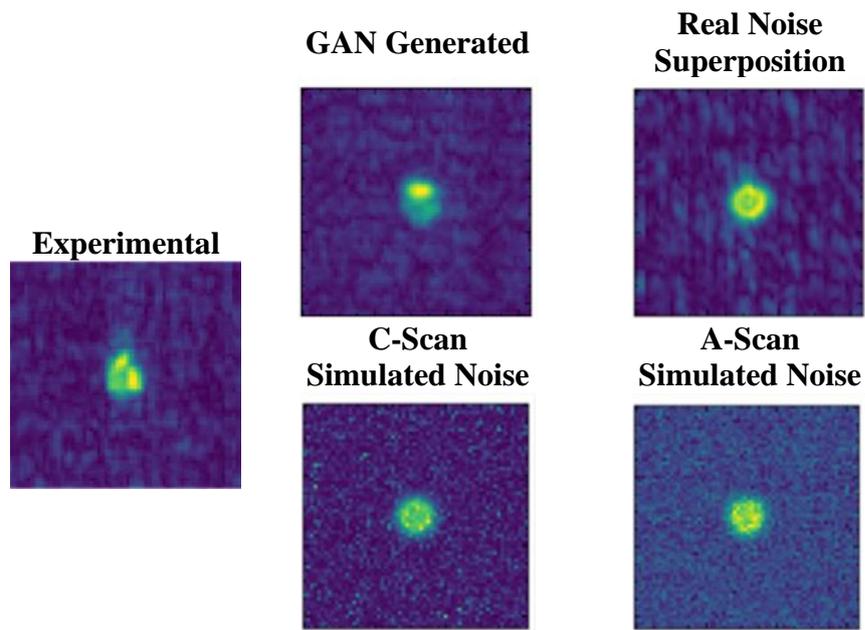
		Predicted	
		Defect	No defect
True	Defect	150.0	0.0
	No defect	59.5	88.5

11 3.7 Discussion

12 3.7.1 Comparison of Classification Results

13 Figure 37 shows examples of C-scan images produced by the different synthetic data
14 generation methods. The classification results are summarised in Figure 38 and Table
15 13, which show the mean (μ) and standard deviation (σ) accuracy and F1 scores, and
16 full evaluation metrics respectively for each dataset investigated.

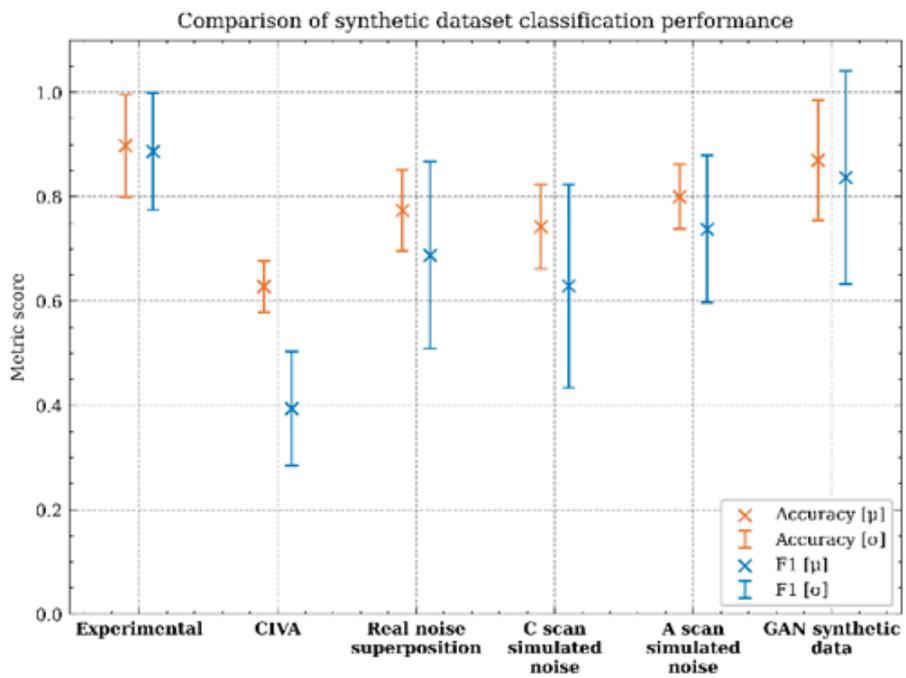
1



2 *Figure 37: Comparison of different real and synthetically generated C-scan image examples.*

3

4



5

6 *Figure 38: Comparison of classification results for each dataset.*

1 *Table 13: Summary of classification results for each dataset.*

Training dataset	Evaluation metric			
	Accuracy	F1	Precision	Recall
Experimental	89.8%	0.887	0.974	0.826
CIVA	62.8%	0.394	1.00	0.252
Modified CycleGAN	87.0%	0.837	0.925	0.775
Real Noise Superposition	77.4%	0.688	0.950	0.545
C-Scan Simulated Noise	74.3%	0.629	0.930	0.482
A-Scan Simulated Noise	80.0%	0.738	0.970	0.598

2

3 Simulated UT data of defect responses in composites lacks the complexity of
 4 experimental noise. In this work, it was demonstrated that when CNN classifiers are
 5 trained on purely simulated data and tested on real experimental data a significant
 6 adverse impact on classification performance is observed, with an average F1 score of
 7 0.39. However, four novel strategies were proposed and explored in this research for
 8 creating synthetic composite UT datasets to reduce this effect with the aim to better
 9 simulate real experimental data. According to the results of this study, all four methods
 10 showed significant increases in classification performance compared to the original
 11 simulated dataset. Among these, the modified CycleGAN generated synthetic dataset
 12 produced significantly better classification results than the other methods, with an
 13 average F1 score of 0.84. This neared the classifier trained on a subset of the
 14 experimental dataset, but due to the reduction in available experimental training and
 15 test data due to the train/test split this should not be considered a direct comparison.
 16 For direct comparison an additional experimental dataset would have been required
 17 for training and testing on the complete test set.

1 Superimposed experimental noise, simulated C-scan noise, and simulated A-scan
2 noise produced similar mean accuracy results, but the simulated A-scan noise synthetic
3 dataset produced the best average F1 score of the three, with 0.74. It is interesting that
4 the simulated A-scan noise dataset outperformed the real noise synthetic dataset. This
5 may be due to the fact the real noise obscures the defect response features too much.
6 Alongside the ability to accurately simulate noise response, a further reason for
7 improved classification results for GAN and A-scan synthetic datasets may be their
8 ability to account for depth wise signal attenuation and adjust the noise levels with
9 respect to depth and signal response. This produces more appropriate noise levels for
10 deeper and weaker defect responses and allows for the preservation of many more
11 simulated responses. Unlike simulated C-scan and real noise approaches which are
12 defect depth agnostic and therefore result in the rejection of more images due to the
13 concealment of low-level responses with noise profiles that are not depth matched.
14 These methods could be extended to include a finer depth wise noise implementation,
15 but this is outside the scope of this work and is left for future investigation.

16 These results demonstrate that in scenarios where noisy experimental environments
17 can cause real data to vary greatly from simulated data, synthetic methodologies for
18 noising data provide an opportunity for generating more effective training data. This
19 is particularly beneficial as we retain the accuracy and fully labelled nature of physics-
20 based simulations, which allow us to fully control the simulation of different defect
21 class types and the variability within them.

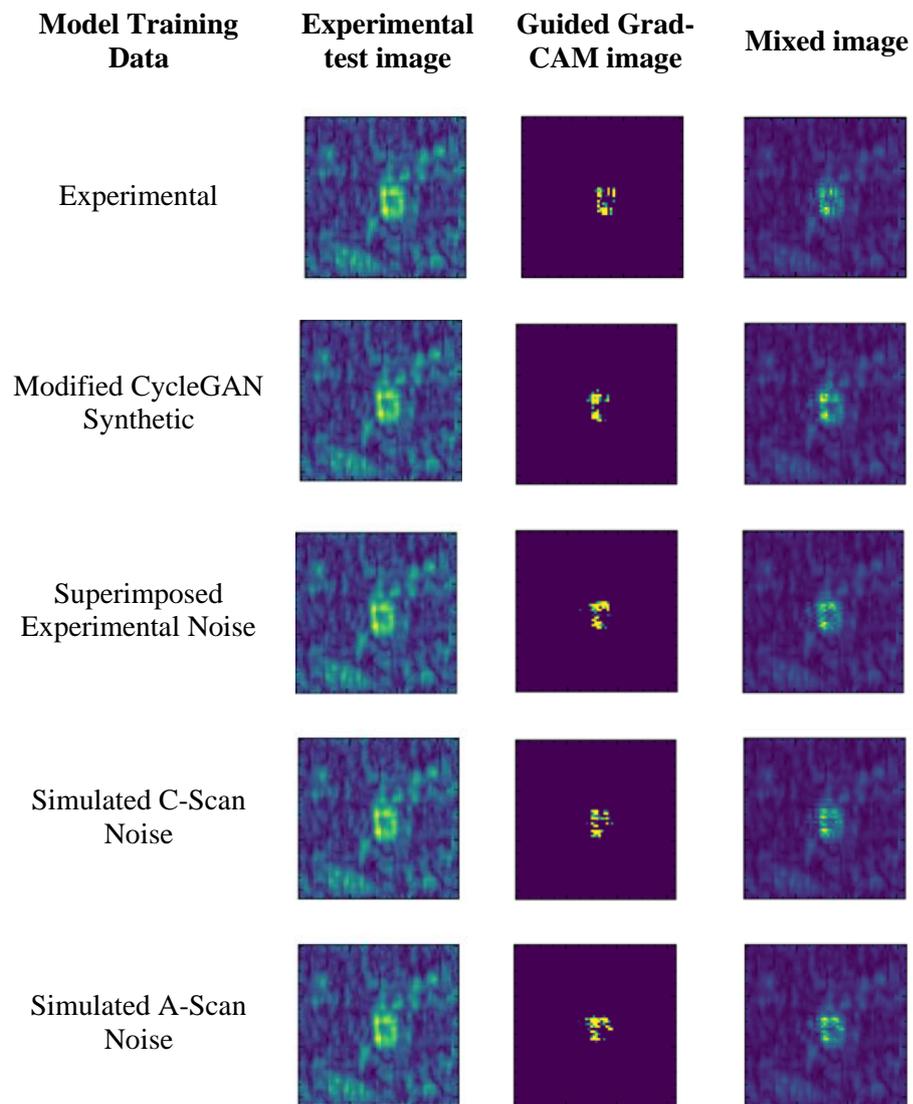
22 When considering the broader aim of generating large synthetic datasets that could be
23 used to create a database of realistic training examples, it is important to consider the

1 ease and robustness of synthetic data generation. Training of the CycleGAN is a
2 delicate process and whilst it has been able to produce realistic images for FBHs, it
3 may struggle to generalise to other defects without significantly broader examples of
4 defects in training. This would largely defeat the point of the synthetic data generation
5 in this instance. Furthermore, the training of an effective GAN model is still extremely
6 challenging and the process of hyperparameter selection is not robust. It is therefore
7 favourable to consider an approach that is robust to different defect types and can be
8 scaled. For scalability, a fully simulated method is preferable over a method which still
9 requires significant collection of experimental data. Therefore, the real noise approach
10 is superseded by both the A-scan and C-scan synthetic approaches. The C-scan noise
11 approach is slightly easier to implement than the A-scan as it requires less
12 experimental data analysis and can be done at the C-scan image level instead of the A-
13 scan level. However, the A-scan noise approach allows for noising of the full
14 volumetric data, which could provide opportunities in three-dimensional data analysis.
15 Further work could be done to explore the distribution of C-scan noise at different
16 depths to enable maintenance of a larger number of simulated responses in a simpler
17 way than the complex A-scan noise simulation method. This could potentially combine
18 some of the benefits of both the A-scan and C-scan noise approaches but would remove
19 the opportunity for volumetric data analysis if done at an image level. In certain
20 scenarios, gathering clean experimental data may not be a limiting factor and in this
21 case, it could be beneficial to expand the real noise superposition method to align the
22 images depth wise between the experimental and simulated domains. This would help
23 to better account for localised structural noise, and likely give improved classification
24 performance similar to the A-scan noise approach.

1 3.7.2 Model Interpretability with Guided Grad-CAM

2 A key barrier to the uptake of Machine Learning in NDT is a lack of model
3 interpretability [14] and the use of synthetic data has the potential to further mystify
4 this process. To help tackle this issue, Guided Gradient-weighted Class Activation
5 Mapping (Guided Grad-CAM) was implemented for a randomly selected model
6 trained from each dataset and evaluated on experimental data. Guided Grad-CAM is a
7 technique for producing ‘visual explanations’ of CNNs with the goal of making them
8 more transparent and explainable [150]. Guided Grad-CAM gives a visual indication
9 of what inputs are used for positive class prediction. Whilst it does not give any
10 information about how or why the inputs are used for the prediction within the model
11 it has been shown to help users place greater trust in a model. The method combines
12 Guided backpropagation and Class Activation Maps (CAM) to create visualisations
13 which indicate relevant image regions for class-discriminative predictions. Guided
14 Grad-CAM is not a complete solution for model interpretability; however, the goal is
15 to visually compare if the models trained on synthetic data are using similar inputs for
16 prediction compared to models trained on experimental data, in the hope that this
17 provides trust in the viability of using synthetic datasets. Figure 39 shows the defective
18 experimental test image, and both the associated Guided Grad-CAM image which
19 gives a visual indication of significant regions contributing to defective predictions,
20 and a mixed image which combines the Guided Grad-CAM and the input image with
21 a respective weighting of 1.5.

22



1 *Figure 39: Example of Grad-CAM visualisation of models trained on different datasets.*

2 It has been identified in literature that model interpretability is a key limiting factor in
3 the uptake of DL in NDT. Guided Grad-CAM was implemented to try and minimise
4 the obscurity that using synthetic data could produce. Whilst model interpretability is
5 a complex field of research and interpretability is challenging to quantify, it is hoped
6 that the Guided Grad-CAM results at least indicate that models trained on synthetic
7 data are learning similar features compared to models trained on purely experimental
8 data.

1 The Grad-CAM images consistently and accurately highlight defect pixels in the
2 context of defect detection. The substantial similarity between Grad-CAM results from
3 models trained on synthetic data and those trained on experimental data suggests that
4 both types of models focus on the same relevant features. This consistency across
5 different training data sources enhances confidence in the models' predictions and
6 supports the use of synthetic data as a viable alternative to experimental data for
7 training.

8 By providing clear visual evidence of accurate defect region identification, Grad-CAM
9 helps to address concerns about the reliability and validity of using synthetic data. This
10 assurance increases confidence in deploying DL models trained on synthetic data for
11 practical NDT applications.

12 3.8 Conclusion

13 Deep learning provides an attractive solution for helping to automate the interpretation
14 of ultrasonic testing NDT data results in quality assurance processes. A barrier to
15 implementation is that DL approaches typically demand large quantities of
16 representative training data to allow accurate and reliable predications to be
17 established. However, since modern manufacturing processes strive to reduce the
18 incidence of defect formation, there is a paucity of real-world defect data available for
19 ML training. Despite this work focusing on FBHs, by employing physics-based
20 simulations of ultrasonic response to defects, it is possible to generate large sets of
21 defect data, with variability in defect types, sizes, and orientations. A drawback in such
22 simulation is in replicating the same noise distributions encountered in experimental
23 measurements, and this is challenging without increasing model complexity to the

1 point of computational intractability. In this study, 4 techniques to map the noise
2 distribution of experimental data onto our simulated data were presented to improve
3 the performance of subsequent ML based classification of defects. A generative
4 network was used to learn the mapping between simulated and experimental images,
5 this resulted in a mean F1 score of 0.843. A method of combining clean experimental
6 images with simulated images was introduced which resulted in a mean F1 score of
7 0.688. To remove the requirement for clean experimental images, two methods of fully
8 generating synthetic noise profiles, C-scan and A-scan noise, were presented; the latter
9 being based on a closer physical representation of how noise is produced
10 experimentally. These methods produced mean F1 classification results on
11 experimental data of 0.629 and 0.738, respectively. Whilst each method produced a
12 significant improvement in classification over the purely simulated data, with a
13 modified loss function to encourage accurate defect response, CycleGAN showed the
14 greatest improvement in classification performance, allowing us to maintain the utility
15 of simulating data from physics-based models and convert them to more
16 experimentally realistic synthetic datasets. However, it was identified that other
17 synthetic data generation methods may be more appropriate for generating large
18 datasets, such as A-scan noise due to their greater robustness.

19 Model interpretability is a significant challenge for the uptake in use of Deep Learning
20 in UT, with the use of synthetic data likely to further add ambiguity. To help minimise
21 this, Guided Grad-CAM was implemented which visually indicated that models
22 trained on synthetic data were learning similar features to models trained on
23 experimental data for classification. This aids in providing confidence that the methods
24 of generating synthetic data are appropriate for training experimental classifiers.

1 Whilst classification results for individual synthetic datasets had room for
2 improvement, this work demonstrates that the synthetic data generation methods were
3 able to successfully transfer the simulation domain closer to the experimental domain.
4 This demonstrates a viable approach to training DL models when experimental data is
5 unavailable, as with many NDT applications.

6 Future work will look to maximise the classification accuracy of specific models. This
7 could be done by combining this work with additional domain adaption techniques,
8 which have shown promise in previous literature [99]. Further investigation will also
9 be conducted to optimise individual model classification accuracy by performing HPO
10 directly using synthetic datasets. This would demonstrate the effects of performing
11 HPO on a model trained on a synthetic dataset and whether this improves its
12 classification in the experimental domain. This would also eliminate the need for
13 experimental data entirely when training a DL classifier as both the parameter
14 optimisation and training could be conducted in the fully synthetic domain. This would
15 require only a small amount of experimental data for testing. Additionally, the next
16 steps in this work will look to see if the style transfer can be extended across the full
17 range of defect types and tested on naturally occurring experimental defects. It would
18 also be beneficial to identify if it is possible to detect more challenging defects such
19 as superficial defects using similar methods or if the UT approach would require
20 modification. If successful, large, fully annotated, synthetic datasets could be
21 efficiently produced, opening the potential for further use of DL in NDT.

22

1 4 Volumetric Detection

2 4.1 Introduction

3 Synthetic datasets are widely used in ML to augment small training datasets [136] and
4 they have been successfully implemented for UT of composites with encouraging
5 results for 2-dimensional (2D) classification of C-scan images in the previous chapter.
6 Part of this work builds upon the work in the previous chapter to extend one of the
7 synthetic data generation methods to make it applicable for full 3D volumes. The
8 synthetic datasets are based on simulations from semi-analytical physics based
9 software that has been shown to produce experimentally accurate defect responses
10 [133], [134]. This software offers a less computationally expensive alternative to FEA,
11 allowing for the simulation of composite responses based on bulk material properties
12 [135].

13 When ML is used to interpret UT NDE data in literature, it is typically applied to
14 interpret A-scan time traces or 2D images constructed from A-scans [97], [98], [101],
15 [104], [106], [131], [132]. Compared to B-scans, A-scans lack all spatial information
16 and nowadays, they are rarely used alone to characterise defects by human operators
17 since the introduction of phased arrays. C-scans preserve detailed spatial information,
18 however constructing the 2D image from the volumetric data necessitates the
19 compression of temporal information. Whilst C-scans excel in capturing intricate
20 spatial details, their need for temporal compression results in minimal representation
21 of through-depth features. Compression of A-Scans to C-Scans often removes useful
22 features such as the backwall response, which can be important when detecting defects
23 with a low reflective index such as porosity [151]. Furthermore, to produce C-scan

1 images appropriate gating must be applied to remove the front wall surface response.
2 This can be challenging when trying to detect near-surface defects. In the aerospace
3 industry, operators typically start with a C-scan to gain a complete picture of defect
4 responses and then move to analysis of B-scans for further information about the
5 nature of the responses [152], [153]. Whilst current ML approaches in literature make
6 use of data in formats that are easily interpreted by humans (images or time-traces),
7 ML algorithms are not limited to image-level analysis and have proved very capable
8 at interpreting 3D volumetric data [154], [155]. By implementing algorithms capable
9 of volumetric interpretation, we retain all spatial and depth information, this gives the
10 algorithms more relevant features to learn from and removes the need for image pre-
11 processing and gating.

12 CNNs have been used effectively for decades in a wide variety of image and
13 volumetric analysis tasks with models such as ResNet typically having tens of millions
14 of parameters [156], and are still widely used as backbones or standalone architectures
15 [74]. However, these networks are typically applied to data of similar dimensions, or
16 data which has been scaled to give even dimensionality of each axis. UT data has
17 extreme aspect ratios due to the difference in requirements of sample rate in the spatial
18 and time dimensions. The upper limits of the temporal and spatial resolution are
19 determined by the sample thickness and by factors such as the number of elements in
20 the ultrasonic probe or the scan length, respectively. Compressing the data in the time
21 dimension to match the spatial dimension, normally dictated by the sub-aperture pitch
22 and the scan acquisition rate, would result in a substantial loss of depth information.
23 Alternatively, the spatial dimensions could be upscaled to match the number of
24 samples in the time dimension, but this is highly inefficient, creating data instances

1 that would require large amounts of memory, and would make training intractable.
2 Therefore, retaining the original dimensionality and aspect ratio of the UT data is
3 highly preferable. Using CNNs to interpret images with high dimensionality is not new
4 and the use of rectangular kernels instead of square kernels in CNNs has given positive
5 results for classification of speech signals, which have high aspect ratios when
6 represented as spectrogram images [157]. The work in this chapter makes use of a
7 similar approach for volumetric data.

8 Network architecture design is a key component of effectively leveraging machine
9 learning techniques. Traditionally, network design heuristics and 'rules of thumb'
10 would be used, in tandem with domain expert knowledge to construct a specific
11 architecture. Automatic architecture design or Neural Architecture Search (NAS) is a
12 development on this approach where a practitioner can leverage compute to aid the
13 process of architecture selection. This process, which can be considered a subset of
14 hyperparameter optimisation, generally involves an iterative process of selecting,
15 training, and evaluating architectures. In its simplest form, a 'Random Search'
16 involves repeating the above process until some threshold or limit in terms of
17 performance or computation budget is reached. More complex approaches to NAS
18 often focus on efficient model evaluations, making use of proxy evaluation methods
19 [158], [159] or efficient sampling algorithms [160], [161] attempting to make the
20 largest improvement with each evaluation.

21 This chapter presents a comparative analysis of the performance achieved from three
22 separate architectures for defect detection in volumetric ultrasonic data. The first,
23 VoxNet [162], is prevalent in the literature for volumetric classification problems, the

1 second architecture presents modifications to VoxNet for this task using a traditional
2 network design approach, and finally a discovered architecture from NAS.

3 VoxNet is a 3D CNN initially proposed for classification of LiDAR, RGBD and CAD
4 data. It has since been used as a backbone for different volumetric classification tasks
5 [163]. Additionally, notable contributions of this study to knowledge in the field
6 encompass the introduction of domain-specific augmentations, which exert a
7 substantial impact on the classification performance. Furthermore, synthetic data
8 generation techniques are leveraged from prior 2D work to generate 3D UT datasets
9 from semi-analytical simulations, effectively addressing one of the prominent
10 challenges encountered in the application of deep learning for NDE: the scarcity of
11 effective training data.

12 This work presents a novel DL architecture designed to process volumetric UT data.
13 In contrast to prior methods relying on time-series data or 2D image-based approaches,
14 which diminish spatial or temporal features, whilst often requiring additional
15 processing. The main contributions of this work are:

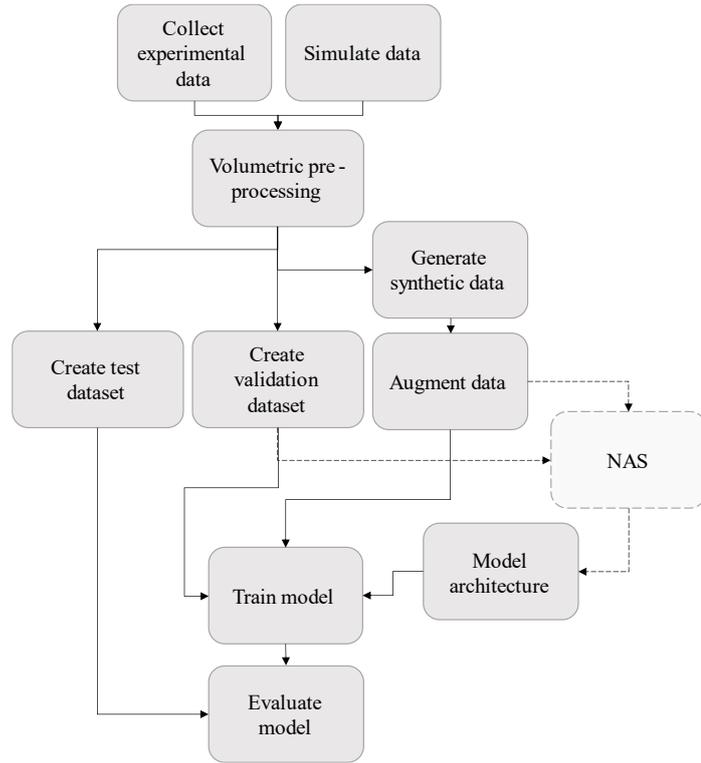
- 16 • Interpretation of volumetric UT data, instead of images or time signals. This
17 reduces preprocessing requirements and allows the model to learn from greater
18 features.
- 19 • Introduction of two domain specific methods for data augmentation, helping
20 with the domain transfer from synthetic to experimental data.
- 21 • Discovery of a novel 3D CNN architecture through NAS.

1 4.1.1 Procedure

2 In this work, the automated data interpretation is simplified by inspecting the complete
3 volumetric data, eliminating image processing steps like gating to remove front and
4 back wall responses, while preserving all spatial and temporal information. Whilst the
5 models are trained on synthetic data, they are tested using experimentally collected UT
6 data from samples with manufactured defects that aim to mimic delaminations.
7 Manufactured defects are commonly used in literature to act as test cases and qualify
8 NDE techniques and operators where naturally occurring defects are not always
9 available [32], [131], [132]. An overview of the simulation and deep learning pipeline
10 is presented in Figure 40. Figure 40 also shows how NAS can fit into this process, with
11 Figure 48 providing a more detailed overview of the NAS pipeline.

12

1



2

3 *Figure 40: Overview of the pipeline for automated volumetric UT classification.*

4 4.2 Data

5 4.2.1 CIVA Simulations

6 Due to the lack of available experimental training data, a simulated dataset was
7 constructed for training. This was done using CIVA, a semi-analytical physics-based
8 commercial NDE simulation software [142]. CIVA has the ability to model wave
9 propagation and interactions with defects. It has been validated, showing good
10 agreement to experimental results for different UT scenarios [133], [134].
11 Additionally, the software is computationally efficient when compared to other
12 alternatives such as Finite Element Analysis (FEA). The full control of the simulated
13 domain enabled the modelling of similar defects and material properties to the

1 experimental domain. However, the use of semi-analytical software instead of FEA
2 had limitations in that the software was unable to model responses from ply
3 interactions and lacks noise seen in experimental data. As a result, differences existed
4 between the simulations and measured experimental responses, leading to the use of
5 the synthetic data generation steps discussed in Section 4.2.2 to reduce the differences
6 between simulated and experimental domains.

7 To set up the simulation, the individual layers of composite were constructed and used
8 to generate equivalent homogeneous material properties of the experimental CFRP
9 samples. A single ply layer was constructed and alternated repeatedly with 8 layers at
10 orientations of 0, 45, -45, and 90 degrees to match the experimental sample as closely
11 as possible (as given in section 2.2). The resulting multilayer structure was
12 homogenised to be consistent with a homogeneous medium having mechanical
13 properties equivalent to those of the multi-ply composite, with the fibre density set to
14 50% best match the experimental sample's density of 1440 kg/m^3 . To simulate the
15 waveform, a sinusoidal wave of 5 MHz was employed, accompanied by a Hanning
16 filter that provided a bandwidth of 66% at 12 dBs.

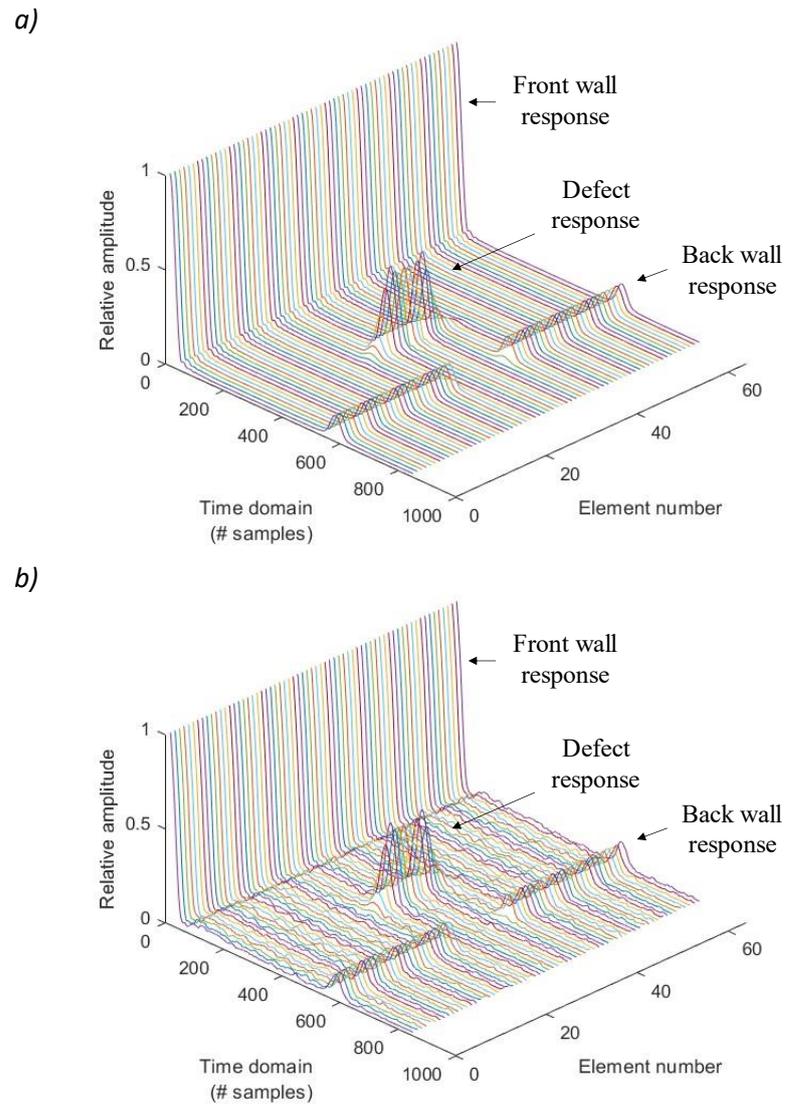
17 For running multiple, sequential simulations, a parametric study was set up, using the
18 composite bulk properties previously calculated and varying the diameter and depth of
19 defects. FBH defects were simulated with diameters from 3.0 mm to 15.0 mm,
20 increasing every 0.5 mm, with varying depths from 1.5 mm to 7.0 mm from the
21 surface, in increments of 1.5 mm. A defect-free simulation was also run to provide the
22 basis for defect-free synthetic data. Both the front and back wall surface reflections
23 were included in the model. The full simulations took less than 15 hours on a desktop

1 computer with a 24-Core 3.79 GHz CPU and 128 GB of memory. An example of a
2 simulated defect is shown in Figure 42 (a).

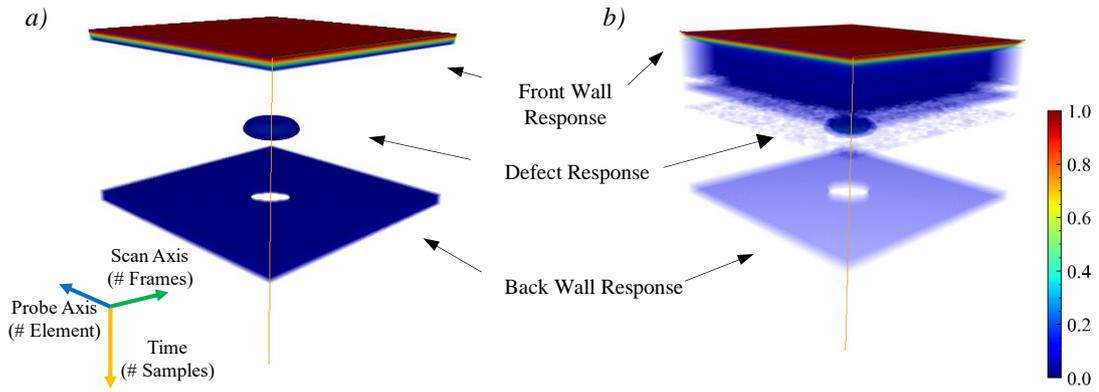
3 4.2.2 Synthetic Data Generation

4 Work conducted in chapter 3 has shown that semi-analytical simulated data alone is
5 not representative enough of the experimental domain [164]. Therefore, there is a need
6 for methods of translating the simulated domain closer to the experimental domain.
7 Fully statistical methods of generating noise are advantageous as they can be re-
8 sampled continuously to keep generating unique noise profiles which are in line with
9 experimental data. In this work, we extend previous work in generating 2D synthetic
10 images and propose a new approach for adding noise to complete volumetric UT data.
11 The previous study concluded [164] that A-scan level noise was the best fully
12 generative statistical method for adding noise. Additionally, all the other approaches,
13 except for the simulated A-scan noise, introduced noise at an image level, which is
14 intractable for volumetric data. To adapt the methodology described in the previous
15 chapter for the analysis of full volumetric data, unique noise profiles for each A-scan
16 were generated and subsequently summed with the simulated responses past the front
17 wall.

18 Figure 41 shows an example of the addition of noise on simulated data at an A-scan
19 level and Figure 42 demonstrates this for a complete ultrasonic volume. The statistical
20 noise distributions of the A-scans were calculated from a separate hold out sample with
21 the same layup and thickness as the test samples. For further details on building up the
22 noise profiles, please refer to the previous work [164].



1 *Figure 41: a) A frame of 64 simulated A-scans for a simulated defect response. b) the corresponding*
 2 *A-scans with synthetically added noise for the same defect response.*



1 *Figure 42: a) Complete ultrasonic volume of simulated A-scans for a defect response. b) the*
 2 *corresponding synthetically noised volume for the same defective response. Both figures have been*
 3 *thresholded to remove the lowest 10% of amplitudes to aid in visual clarity.*

4 A summary of the datasets generated from the experimental and synthetic data is given
 5 in Table 14.

6 *Table 14: Summary of the datasets produced.*

Data source	Dataset	Number of datapoints
Simulated defect responses (300 Flat-Bottom Holes)	Synthetic defective train	300
Simulated defect free response	Synthetic defect free train	300
Experimental defect reference sample (15+25 Flat-Bottom Holes)	Defect test (70%)	25
	Defect validation (30%)	15
Experimental defect free reference sample	Defect free test (70%)	25
	Defect free validation (30%)	15

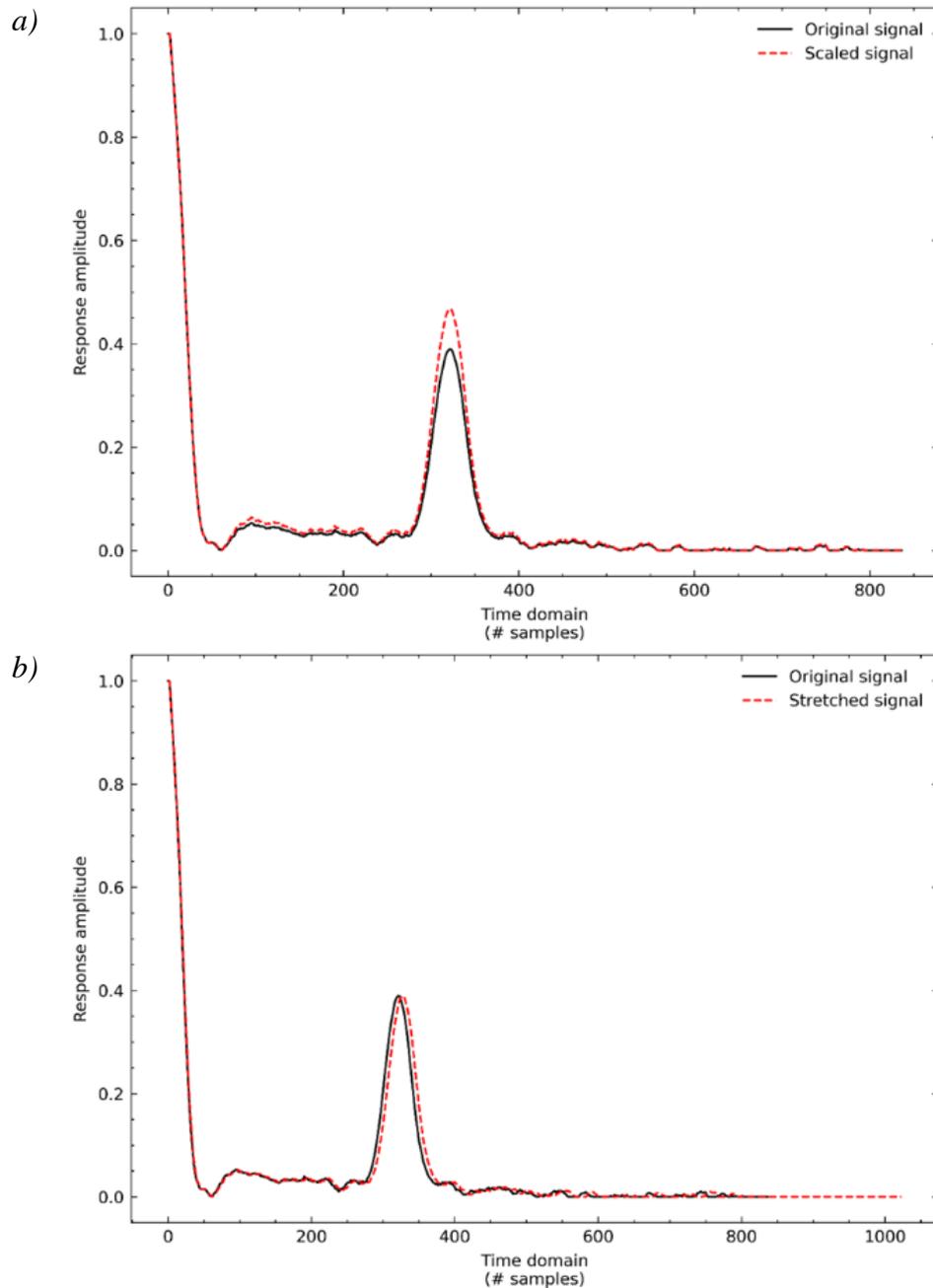
7 4.2.3 Augmentation

8 The generalisability of ML models is a critical aspect of their performance. One
 9 approach to improve generalisability is to augment the training data. Augmenting the
 10 training data makes the task more challenging by adding noise at the training stage,
 11 reducing the likelihood of overfitting, and often improves performance in the target
 12 domain. This is particularly important when the target (experimental) domain is
 13 different from the training (synthetic) domain.

1 As demonstrated in Figure 42 (a), there is little variation between simulated A-Scans.
2 However, this is not the case for experimentally acquired data. Received amplitudes
3 are affected by the sensitivity of individual elements, variations of couplant on the
4 surface of the sample, and roughness of surface finishes (particularly for manufactured
5 defects). The anisotropy of CFRP can also result in variations in attenuation, which
6 impact effect received amplitudes. Surface roughness and local changes in fibre
7 density due to the materials inherent anisotropy produce small changes in time of
8 flight. Traditional augmentation methods such as those used for images (e.g. crop, mix-
9 up, flipping etc.) do not model these variations well and can produce unrealistic
10 examples.

11 Therefore, in this study, we introduce two types of augmentation that were generated
12 online for each minibatch during training. These augmentations aim to mimic the inter-
13 element response variability observed within the UT probes used for data collection.

14 The first type of augmentation is related to the magnitude of response measured by the
15 UT elements, which varies due to many factors not included in the simulation, such as
16 manufacturing tolerances of the sample and the UT array probe, wear and tear of the
17 probe and electrical wires/connections, or inter-layer multiple scattering of the sound
18 waves. To mimic these whilst preserving the correct normalisation, each A-scan was
19 scaled by a constant past the front wall. The scale factor was sampled from a uniform
20 distribution to give a scale factor between 80-120%. An example of this is given in
21 Figure 43 (a).



1 *Figure 43: a) Example of how scaling augmentation is done on an individual A-scan. b) Example of*
 2 *how dilation augmentation and padding is completed for an individual A-scan.*

3 The second type of augmentation mimics any changes in ultrasonic travel time seen
 4 by different elements. This can be caused by a variety of factors, such as variations in
 5 component sound speed due to the anisotropic nature of composites, departure from
 6 central frequency for certain elements, etc. To simulate this 1-D interpolation was used

1 to randomly stretch or compress the signal in the time domain. The dilation amount
2 was randomly sampled from a uniform distribution for each A-scan up to ± 15
3 samples. An example of this is given in Figure 43 (b).

4 By introducing these augmentation methods, we aim to improve the generalisability
5 of the models to the experimental domain. The online nature of these augmentations
6 means that they can be easily incorporated into the training process without the need
7 for additional data collection or pre-processing steps. To ensure consistent length of
8 data in the time domain, each A-scan was padded with zeroes to a length of 1024
9 samples during training. Further investigation could be conducted to identify
10 additional domain specific augmentation methods which could help to bridge the
11 simulation to experimental domain gaps. This could include variations to frequency
12 and bandwidth.

13 4.3 Network Architectures

14 In this work we investigated the performance of three different 3D CNN architectures
15 for binary classification of 3D defect and defect free UT data with extreme aspect
16 ratios.

17 The first 3D CNN, VoxNet, was designed for similar volumetric classification tasks
18 and acts as a baseline architecture. For low aspect ratios CNNs (such as VoxNet)
19 typically make use of square or cuboidal kernels which are appropriate for their equal
20 (or near equal) aspect ratios. The use of CNNs on data with more extreme aspect ratios
21 is less common and is particularly extreme for UT data between the time and the spatial
22 domains, with an aspect ratio of 16.

1 To overcome this challenge, a task-specific architecture was hand-crafted by adapting
2 VoxNet in a manner that follows the traditional approach to architecture design. This
3 custom network is specifically designed to tackle the extreme aspect ratio problem and
4 enhance overall classification performance.

5 As an alternative to traditional architecture design, neural architecture search was
6 employed to develop a third architecture for comparison. For each model Adam
7 optimiser [70] was used with a constant learning rate of 0.001, β_1 of 0.9 and β_2 of
8 0.999. A batch size of 8 was utilised in the training process. The chosen loss function
9 for this model was binary cross-entropy, with a sigmoid activation function applied to
10 the final layer to facilitate classification.

11 Due to the small amounts of experimental test data, there was a likelihood of noisy
12 results during both training and testing phases. To mitigate this, each model was
13 trained ten times with varying random initialisations, and their individual results were
14 averaged across the performance metrics. This gives a better representation of the
15 model's performance by averaging out any noisy results due to the small datasets.

16 During the training phase, a fixed validation set comprising 30% of the total test data
17 was randomly selected from each class of experimental data. This set was used to
18 monitor the model's performance and minimise the risk of overfitting. The models
19 were trained with a patience of 10 epochs, where the training process monitored binary
20 cross entropy loss on the validation data, for improvement. If there was no
21 enhancement for a consecutive period of 10 epochs, the training process was halted.
22 The model parameters with the lowest validation loss were used to evaluate the
23 classification performance on the test set. This approach ensured that the final model's

1 defect detection performance was evaluated using the parameters that had the best
2 ability to generalise to the target domain, as opposed to the model that had overfit to
3 the synthetic domain.

4 4.3.1 Evaluation metrics

5 To quantitatively assess the binary classification performance of each network,
6 average mean accuracy, precision, recall and F1 scores were calculated according to
7 Equations 11-14.

$$Accuracy = (TP + TN) / (TP + TN + FP + FN) \quad (11)$$

$$Precision = TP / (TP + FP) \quad (12)$$

$$Recall = TP / (TP + FN) \quad (13)$$

$$F1 = (2 \times Precision \times Recall) / (Precision + Recall) \quad (14)$$

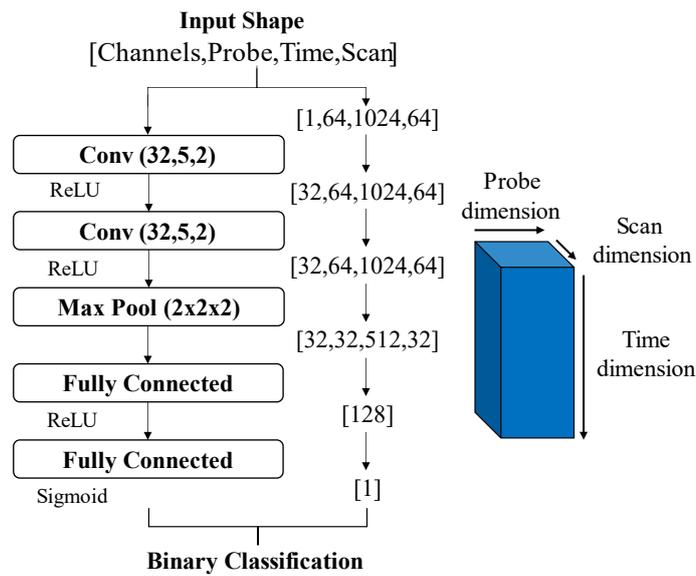
8 Where TP is true positive, TN is true negative, FP is false positive, and FN is false
9 negative, with positives being the presence of a defect within an inspected volume.
10 Each result was individually averaged using a simple mean across the 10 training
11 cycles.

12 4.3.2 VoxNet: Baseline Architecture

13 Introduced by Maturana and Scherer, VoxNet [162] is a 3D CNN designed to tackle
14 classification problems of 3D data that can be represented as voxels to form an
15 occupancy grid. Originally tested on LiDAR, RGBD and CAD data it has since been
16 used as the backbone for methods tested on ModelNet40 [163].

17 While the data from UT for this task differs from the datasets previously employed
18 with VoxNet, the process of converting data into voxel-based format within the

1 VoxNet pipeline is well aligned to the 3D representation of UT data. As a result,
 2 VoxNet was employed to establish baseline model performance metrics for this task.
 3 VoxNet is constructed using two 3D convolutional layers with cuboidal kernels,
 4 followed by a pooling layer and two fully connected layers (Figure 44). For further
 5 details on the model please refer to the original paper. VoxNets total number of
 6 parameters is 235M.



7
 8 *Figure 44: The VoxNet architecture. Where Conv (f,d,s) indicates the number of filters f, filter size d,*
 9 *and stride s, of the convolutional layer.*

10 4.3.3 Hand Designed Architecture

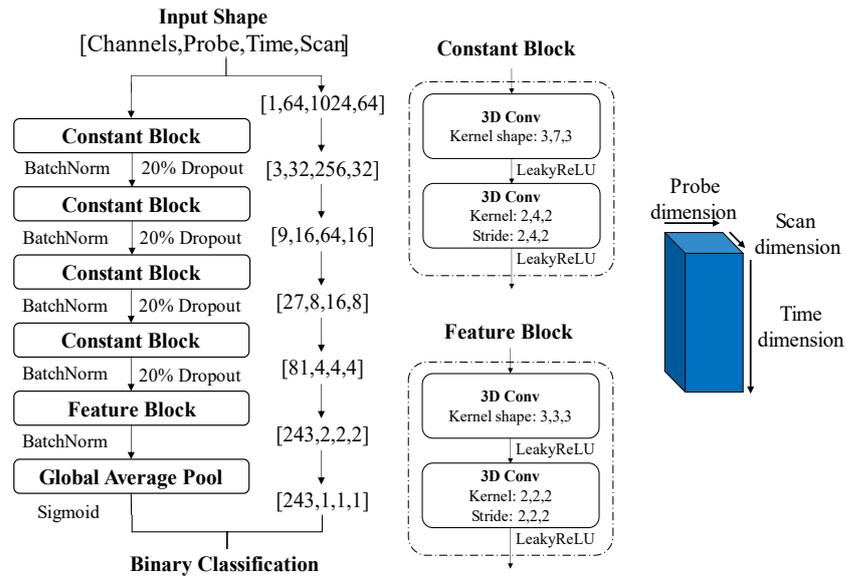
11 The second architecture, referred to as CustomNet, demonstrates a conventional
 12 approach to architectural design. In this context, adaptations to VoxNet have been
 13 implemented to contemporise and enhance its performance specifically for the given
 14 task.

15 The UT dataset stands out to previous VoxNet datasets due to its higher dimensionality
 16 coupled with notable differences in spatial and temporal dimensions. To effectively

1 handle these unique attributes, adjustments were made to the models' architecture.
2 Specifically, the number of convolutional layers was increased to enhance the
3 extraction of meaningful features from the complex dataset. Additionally, cuboidal
4 kernels with non-uniform dimensions were employed in the initial four blocks (refer
5 to constant blocks in Figure 45) of the model. This approach aimed to address the
6 uneven dimensionality inherent in the data, ultimately equalising the dimensions and
7 contributing to a more robust feature representation throughout the network (Figure
8 45). After this a feature block with cube kernels of equal dimensionality could be used
9 (refer to feature block in Figure 45).

10 In the process of updating VoxNet, we incorporated convolutional layers for pooling
11 instead of the previously employed max-pooling layers. Additionally, ReLU was
12 substituted with LeakyReLU, and batch normalisation was introduced. To mitigate
13 overfitting, dropout and global average pooling were employed to reduce the number
14 of features for classification, avoiding the use of large fully connected layers. These
15 modifications are geared towards improving the model's performance by incorporating
16 contemporary practices that have shown substantial performance benefits, as
17 highlighted in previous studies [82].

18 The final architecture is given by the diagram in Figure 45. The total parameter size of
19 the network was estimated to be 1.28 M parameters.



1

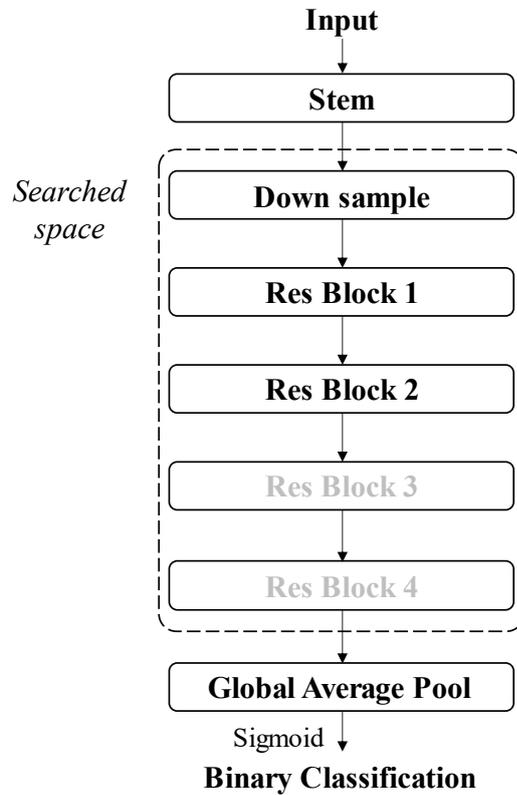
2 *Figure 45: Network architecture for the CustomNet.*

3 4.3.4 NAS Discovered: 3D ResNet based Neural Architecture Search

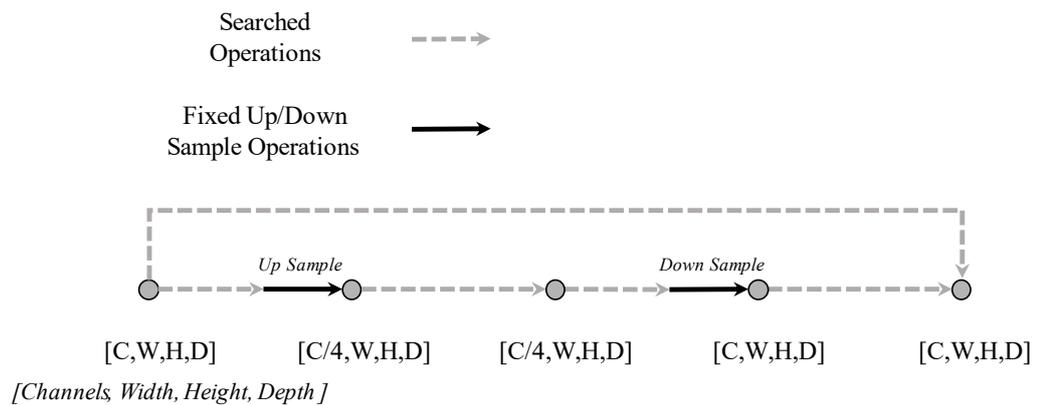
4 4.3.4.1 Neural Architecture Search

5 The final architecture was developed through NAS of a modified ResNet search space
6 to account for 3D convolutions and operations. One of the challenges in applying NAS
7 to a new domain task is the design of the search space. For this task, a new search-
8 space framework which utilises a novel search space based on a ResNet-like structure
9 is introduced. A fixed stem was used to down sample the data by a factor of 4 in the
10 spatial dimensions and a factor of 8 in the time dimension whilst aiming to retain
11 information through increasing the channels to 64. A further down sample block with
12 average pooling followed by two to four residual blocks were all searched individually.
13 An overview of the structure can be seen in Figure 46. The residual blocks and
14 bottleneck features of the ResNet architecture are retained, whilst searching operations
15 for each edge within the residual block. This provided a large diversity of architectures,
16 which is key to attaining good performance in a novel application, whilst also ensuring

1 that many networks conformed to successful design principles. Each residual block
 2 contained two fixed point-wise convolutions used to down and up sample the number
 3 of channels. Figure 47 shows an example of a residual block denoting the searched
 4 and fixed operations.



5
 6 *Figure 46: Representation of the ResNet style searched space.*



7
 8 *Figure 47: Diagram of the searched residual block.*

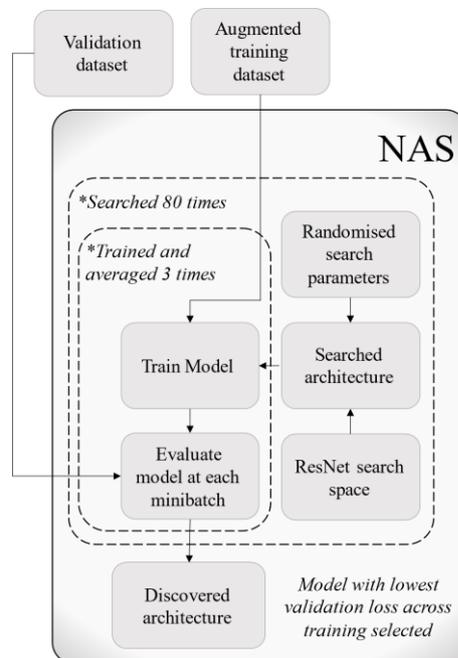
1 These blocks were then stacked in groups, with the resolution down sampled between
2 groups. Equation 15 gives the probability of a new group being created for each
3 residual block, otherwise they were added to the current group. This makes groups
4 unlikely to be extremely long or short.

$$P(\text{newResBlockGroup}) = \frac{1}{\text{currentResBlockGroupSize}} \quad (15)$$

5 The primitive operations of a search space are the list of operations which are assigned
6 to the edges of a network architecture. The implemented approach incorporated a
7 standard set of operations commonly found in the NAS literature. These operations
8 comprised of convolutions, pooling, and skip connections, which are widely
9 recognised and utilised within the field. These operations were all 3D due to the
10 dimensionality of the data. In contrast to standard practice, which makes use of
11 separable convolutions, the approach presented in this study deployed both depth-wise
12 and point-wise convolutions as the fundamental convolutions within the search space.
13 This significantly reduced the number of parameters in each operation of the
14 architecture, greatly reducing the computational cost. Specifically, the depth-wise
15 convolutions were applied with equidimensional cube kernels, of size 3, 5, or 7,
16 coupled with dilation values that ranged from 1 to 4. Skip connections, point-wise
17 convolutions, as well as average and max pooling operations were also searched for.
18 For the pooling operations, equidimensional cube kernels of size 3, 5, or 7, with a
19 dilation value of one were employed. The search encompassed the exploration of
20 Gaussian Error Linear Unit (GELU) activation function and batch normalisation, as
21 well as the absence of activation and normalisation operations. This allowed for
22 architectures with fewer activation and normalisation function which has been shown

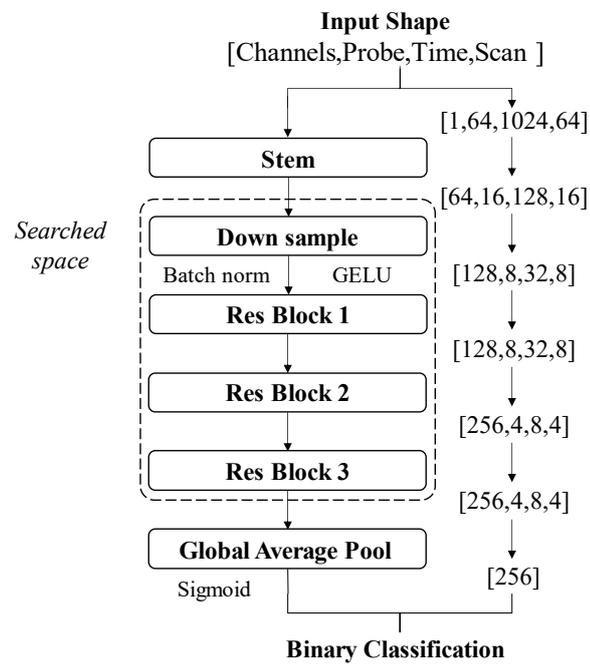
1 to be beneficial [82]. The searched down sample operation had a fixed kernel size with
 2 two in the spatial dimensions a four in the temporal dimension with a dilation of one.
 3 Throughout the relevant operations, a stride of one was employed.

4 A simple random search was applied to this search space for 80 iterations. Each model
 5 was evaluated using the validation dataset, with the lowest loss on validation across
 6 the training taken as the evaluation metric. For each searched architecture, a model
 7 was retrained with new initialisations three times and the mean evaluation metrics were
 8 used when selecting the discovered architecture, this ensured a more accurate estimate
 9 of model performance. Cross validation was unable to be used as the combination of
 10 NAS and domain transfer would have resulted in data leak between the NAS stage and
 11 the final model test evaluation stage. Figure 48 provides an overview of the NAS
 12 process and demonstrates how separation of the validation and test set were maintained
 13 in context of the complete model pipeline, given in Figure 40.



14
 15 *Figure 48: Overview of the process for NAS implementation.*

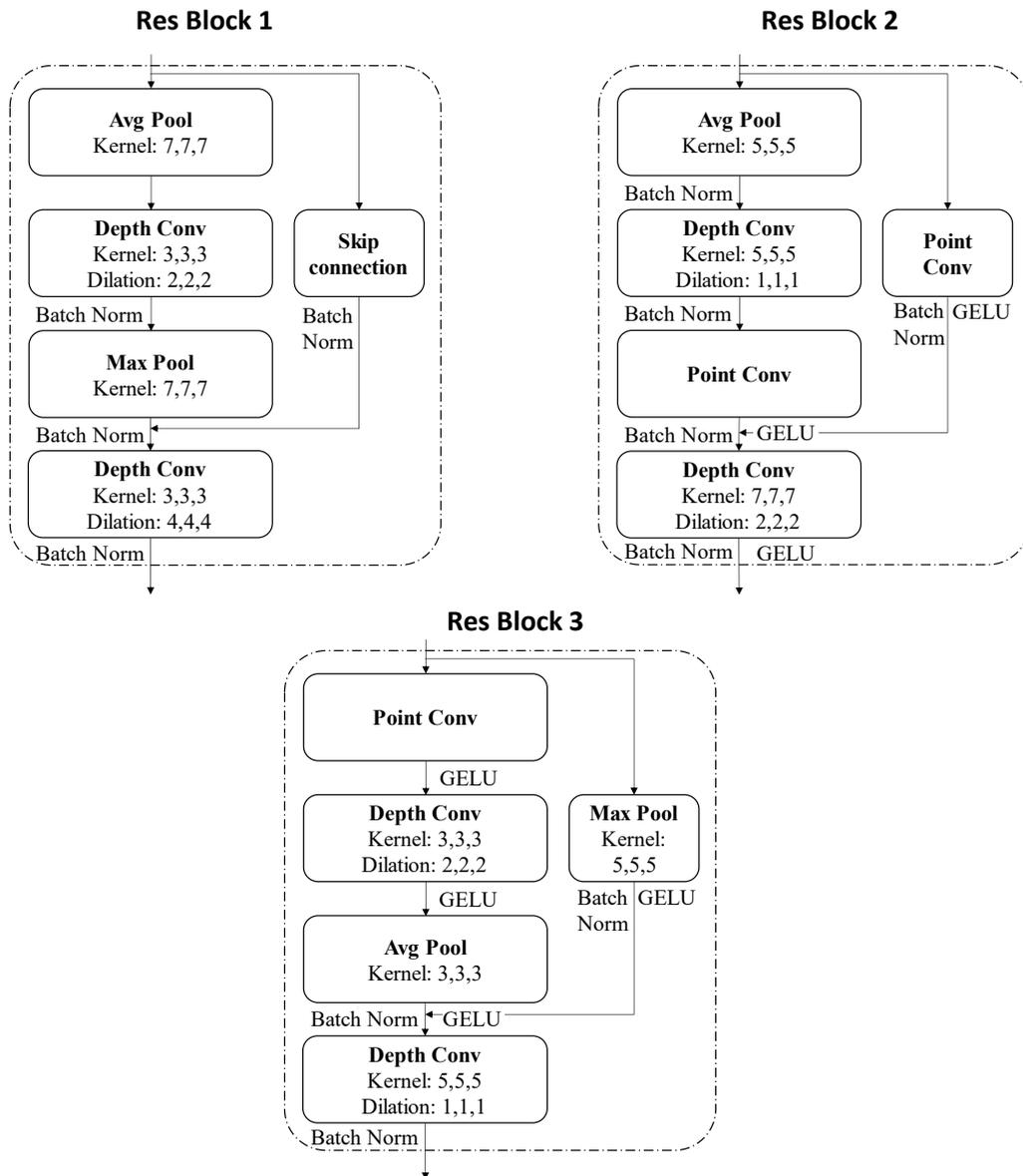
- 1 The final discovered architecture had 1.03 M parameters and is given in Figure 49,
- 2 with the details of the residual blocks given in Figure 50.



3

4 *Figure 49: The overall structure of the discovered architecture.*

5



1 *Figure 50: The details of each discovered residual block.*

2

1 Results

2 *Table 15: Average confusion matrices for VoxNet, CustomNet and the NAS discovered architecture.*

	True	Predicted	
		No Defect	Defect
<i>VoxNet</i>	No Defect	15.8	9.2
	Defect	2.1	22.9
<i>CustomNet</i>	No Defect	24.3	0.7
	Defect	3.2	21.8
<i>NAS Discovered</i>	No Defect	25.0	0.0
	Defect	0.0	25.0

3

4 *Table 16: Comparison of classification results across the different architectures. The means and*
5 *standard deviations are presented as mean \pm std.*

Model	VoxNet	CustomNet	NAS
Accuracy	0.774 \pm 0.184	0.922 \pm 0.095	1.00 \pm 0.00
F1	0.825 \pm 0.114	0.904 \pm 0.130	1.00 \pm 0.00
Precision	0.793 \pm 0.219	0.975 \pm 0.050	1.00 \pm 0.00
Recall	0.916 \pm 0.073	0.872 \pm 0.197	1.00 \pm 0.00

6 Table 15 provides average confusion matrixes for the test results of the VoxNet,
7 CustomNet and NAS discovered models. Table 16 presents a summary of each
8 method's performance, displaying the mean and standard deviation across various
9 performance metrics.

10 The architecture discovered by NAS consistently produced ideal results when trained
11 using data augmentation, with a mean classification accuracy of 1.00 and a standard
12 deviation of 0.00 across the 10 separate training iterations, this demonstrated high
13 confidence in the model's conclusions and robust design for the target domain.

14

1 *Table 17: Comparison of the effects of data augmentation on the NAS discovered architecture. The*
 2 *means and standard deviations are presented as mean \pm std.*

Augmentation	None	Scaling	Both
Accuracy	0.776 \pm 0.178	0.806 \pm 0.227	1.00 \pm 0.00
F1	0.830 \pm 0.123	0.842 \pm 0.168	1.00 \pm 0.00
Precision	0.745 \pm 0.19	0.846 \pm 0.227	1.00 \pm 0.00
Recall	0.972 \pm 0.044	0.916 \pm 0.182	1.00 \pm 0.00

3 *Table 18: Comparison of model sizes and inference time for each architecture.*

Model	VoxNet	CustomNet	NAS
Total Parameters (M)	235	1.28	1.03
Total Size (MB)	1779	557	93
Inference Time (seconds)	0.37	0.03	0.40

4 Table 17 demonstrates the impact of data augmentation on the best performing NAS
 5 model. Discarding data augmentation completely during training had a significant
 6 impact on the classification performance with a 22.4% drop in mean accuracy, along
 7 with a standard deviation increase in 17.8%, which demonstrated a significant
 8 reduction in statistical confidence. Whilst the addition of amplitude scaling
 9 augmentation improved the mean accuracy, it was only by 3%. This demonstrates the
 10 importance of using both augmentation methods in parallel for increased
 11 generalisability to the experimental domain.

12 Table 18 provides a summary of model sizes, and inference times for a single batch of
 13 test data. The NAS discovered architecture has a total size 16.7% and 5.2% smaller
 14 than CustomNet and VoxNet respectively. Whilst the CustomNet was 12 times faster
 15 at inference than the next closest, VoxNet.

1 4.4 Discussion

2 VoxNet demonstrated it was able to learn features from synthetic data and performed
3 reasonably well on experimental data, with a mean F1 score of 0.825. However, its
4 significant standard deviation in accuracies between training instances demonstrates
5 that the architecture was not well optimised for the problem. The CustomNet improved
6 on the accuracy of VoxNet substantially by 14.8%, whilst also reducing the standard
7 deviation of results by 8.9% which indicated an increase in consistent generalisability
8 to the experimental domain. This illustrates the benefits of tailoring architectural
9 modifications to address the needs of specialised tasks. The experimental results
10 demonstrated that the architecture discovered from NAS greatly outperformed the
11 other two in terms of classification accuracy. Whilst all the models used in this work
12 are not large and are considerably smaller than typical sizes for 2D ResNet's and other
13 CNNs [156], the NAS model was able to achieve the highest performance with a
14 significantly lower model size, at only 5.2% the memory requirement of VoxNet. The
15 black box nature of DL makes it difficult to specify which features lead to this
16 improvement in performance. This complexity is in fact a large motivator for NAS as
17 the design space is too large for a human to efficiently find an optimal network
18 architecture. It is anticipated that the addition of skip connections and the ability to
19 vary operations at different depths added by the NAS has a significant positive impact
20 on performance. This demonstrated the importance of utilising neural architecture
21 search to optimise CNNs.

22 Due to the large fully connected layers, VoxNet results in a far greater number of total
23 parameters than the other two networks. This results in a model which occupies far

1 more memory. Whilst CustomNet and the NAS discovered architecture have a
2 comparable number of parameters the discovered network is far smaller. This is a
3 result of many of its operations being far more efficient, such as the separation of point
4 and depth wise convolutions. Whilst the discovered architecture is the smallest, its
5 inference time takes the longest due to the greater architectural complexity of the
6 model and its operations. This said, all models have acceptable inference time and can
7 process 8 samples in under half a second. However, CustomNet is notably twelve times
8 faster at inference than the second fastest network, VoxNet, which could be an
9 advantage in some industrial settings.

10 When trained without data augmentation the NAS model performed significantly
11 worse. Furthermore, the performance was only slightly improved by adding amplitude
12 scaling augmentation alone. For best performance, both augmentation methods were
13 needed in combination. This indicates that despite accurate synthetic data generation,
14 data augmentation still has a significant role in producing generalisable models to the
15 experimental domain.

16 Whilst ideal classification was achieved consistently for the discovered architecture
17 when trained with data augmentation, this was tested on detection of manufactured
18 defects only. Specifically, back drilled holes which are perpendicular to the
19 propagating sound wave and act as ideal reflectors. This makes them comparably
20 easier to detect than other defects. Whilst samples with naturally occurring defects are
21 challenging to get access to, future work would benefit by expanding the simulation
22 scope and testing the models on naturally occurring defects which will likely prove
23 more challenging to detect. For more challenging detection and characterisation tasks

1 a more sophisticated search optimisation algorithm could be employed to discover
2 architectures more efficiently.

3 The achieved classification results suggest that the synthetic data generation process
4 is a viable approach for producing fully synthetic 3D UT volumetric datasets that
5 closely map to the experimental domain and enable the development of effective
6 classifiers. However, due to the substantial improvement in classification performance
7 achieved through the implementation of data augmentation methods, it is important to
8 acknowledge that disparities between the synthetic and experimental domains persist.
9 This observation underscores the necessity for augmentation techniques to further
10 enhance the generalisability of the model. Nonetheless, it is worth noting that the data
11 augmentation methods employed in this study proved to be highly effective in aiding
12 not only generalisability but also in facilitating the transfer of knowledge across
13 domains.

14 The key benefits for analysing the complete 3-D volumetric data instead of processed
15 images were the ability to learn from greater features, the reduction in pre-processing
16 requirements, and the potential reduction in inference time by analysing the complete
17 volume all at once. The impact on inference time is challenging to quantify, however
18 if comparing the compute required to process 64 B-scan images (the equivalent spatial
19 scan data), without parallelisation for equivalent 2D classifiers, there is the potential
20 for up to 64 times saving in inference time for the same scan area. Despite these
21 advantages there are still potential benefits to analysing UT data as images. One of
22 these is the many opportunities for detection of a single defect in multiple B-scans. It
23 is likely that defects will span multiple B-scan images, and as such by analysing each

1 B-scan there are multiple chances to detect an individual defect. This means an
2 individual defect can still be detected even if individual defective images are
3 incorrectly classified. However, the opportunities for characterisation and localisation
4 of defects are far greater when retaining the volumetric spatial information and this
5 work opens future prospects for 3D classification and segmentation which would be
6 much more challenging if using C-scans or B-scans alone.

7 The research outcomes demonstrated the considerable potential of employing 3D-
8 CNNs in conjunction with well-designed data augmentation techniques and optimised
9 architecture search spaces to address challenging 3D classification tasks characterised
10 by extreme aspect ratios, as observed in the context of UT. Insufficient utilisation of
11 data augmentation severely hampered the model's ability to generalise to experimental
12 datasets, leading to suboptimal classification performance. Likewise, choosing an
13 unsuitable model architecture could result in the failure to capture crucial features
14 necessary for accurate classification. Consequently, it is imperative to thoroughly
15 consider both aspects during the design of a classification model for 3D UT data to
16 ensure optimal performance.

17 4.5 Conclusion

18 Deep learning has demonstrated prior success in ultrasonic non-destructive evaluation
19 when applied to either time series or image data. However, analysing only time series
20 or image data can result in a significant loss of information in either the temporal or
21 spatial domains. This work proposes the use of 3D convolutional neural networks to
22 classify complete volumetric ultrasound data without compression, retaining all spatial
23 and temporal information. This approach not only reduced the need for accurate gating
24 when constructing C-scan images but also decreased the amount of signal processing

1 required. To train the models, synthetic data was generated from semi-analytical
2 simulations, while experimentally collected ultrasonic responses from manufactured
3 defects were used for testing. Two forms of data augmentation were implemented
4 based on physical variations seen in experimental ultrasonic responses to improve the
5 model's classification performance in the experimental domain. Furthermore, the
6 performance of three different architectures; one existing in the literature, one hand-
7 designed based on current practices, and one designed by NAS from a ResNet search
8 space modified for 3D, were compared.

9 The first architecture, VoxNet, performed reasonably well on experimental data,
10 achieving a mean F1 score of 0.825. However, its notable standard deviation in
11 accuracies during training suggests suboptimal architecture optimisation for this task.
12 CustomNet's greatly improved on VoxNet with an accuracy increase of 14.8%, whilst
13 reducing the standard deviations in accuracy by 8.9%, hence demonstrating an
14 architecture better optimised for this task.

15 The third architecture, designed by NAS, when trained with data augmentation, gave
16 the best results, providing 100% classification accuracy. The impact of online domain
17 specific augmentation was notable, leading to a 22.4% decrease in mean accuracy for
18 the NAS model when augmentation was omitted.

19 Overall, this work demonstrated that it is possible to train successful DL models to
20 classify full volumetric ultrasonic data for NDE. The issue of a lack of data in most
21 NDE situations was addressed by successfully implementing synthetic data generation
22 in 3D. The work highlighted the importance of appropriate architecture selection and

1 effective data augmentation when translating between synthetic and experimental
2 domains, with both factors essential in achieving high classification accuracy.

3 The focus of this work was on the use of volumetric datasets, and whilst 100%
4 classification accuracy was achieved through effective NAS, it is recognised that FBHs
5 are generally simple defects to detect by human operators.

6 Future work aims to increase the complexity of the task by detecting a wider range of
7 more challenging defects and expanding the simulation scope to better cover naturally
8 occurring defects, where performance can be measured against human operators in a
9 more realistic industrial scenario. Further work also aims to extend the problem to
10 defect classification and sizing.

1 5 Supervised Volumetric Defect Segmentation

2 5.1 Introduction

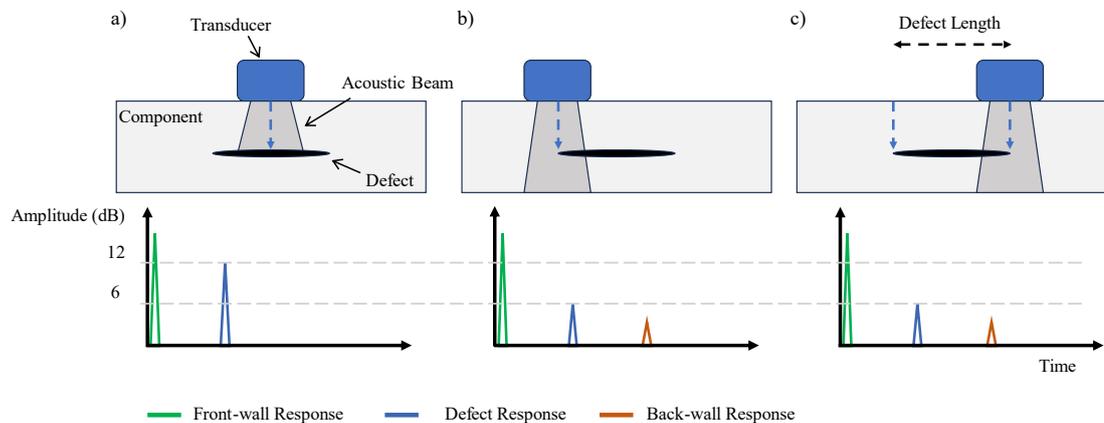
3 Automated NDE data interpretation and reporting can be broken down into distinct
4 sections which in many cases will flow sequentially, as presented in Figure 9. In most
5 settings, defect detection is the primary task to be completed. Once defects are detected
6 it is important to evaluate their size against acceptance criteria. In some applications
7 the acceptance criteria will vary depending on the defect type, orientation, and
8 location, so the identification of different defect types and the location in the geometry
9 is also important. In an industrial setting it is crucial to report on the NDE findings for
10 traceability and to allow for downstream testing, design decisions, and re-work. It is
11 therefore not acceptable to just detect a defect. Defects must be evaluated to extract as
12 much information about the defect as possible, such as sizing, type, location etc [165].
13 This is covered by the characterisation and quantification section of the pipeline
14 (Figure 9). Whilst there is potential for a single end-to-end model/system to complete
15 the whole interpretation process; by breaking the tasks up it will allow for greater
16 model testing and improved understanding and comparison to human operators.
17 Breaking down the data interpretation process into separate tasks not only enables
18 comprehensive model testing but also provides valuable insights into how errors
19 propagate through the system. This decomposition allows for a more nuanced
20 evaluation, shedding light on the strengths and limitations of each component of the
21 analysis. Additionally, the isolation of specific tasks enhances the ability to discern
22 how errors manifest at different stages, making it possible to facilitate a more effective
23 comparison and agreement with human operators. Their expertise can be leveraged to
24 refine and optimise each element. This stepwise methodology ensures a thorough

1 assessment of the entire system and establishes a foundation for collaborative decision-
2 making between automated systems and human operators. Moreover, this modular
3 approach enhances the system's adaptability to different requirements, allowing for
4 customisation and optimisation based on specific application needs.

5 Much of NDE automation research focuses on defect detection, often overlooking the
6 significant burden placed on NDE operators for report generation. By analysing
7 complete ultrasonic volumes of a sample, this work focuses on defect sizing and 3D
8 localisation, which can be used directly when evaluating defects against accept/reject
9 criteria and provide accurate defect positional information, which is useful in
10 supporting re-work etc. This is accomplished through segmentation of ultrasonic
11 volumes, which allows for automated generation of computer-aided design files to
12 further reduce the burden on NDE reporting and can be used for building digital twins
13 of components for testing [166], [167].

14 Accurate defect sizing is a key metric in determining if a component is safe against
15 standardised acceptance criteria. The 6 dB drop method is a widely accepted method
16 for defect sizing and is commonly used in industrial standards [16]. The technique
17 relies on the utilisation of a single transducer and the peak amplitude from the defect
18 response to determine the boundaries of a defect response by detecting the point at
19 which the transducer is directly over the edge of the defect as determined by a 50%
20 energy dissipation from the reflector, manifesting as a 6 dB reduction in amplitude as
21 indicated in Figure 51. The method benefits from being based on physical properties
22 and is fully explainable. The 6 dB drop is often extended and applied to amplitude C-
23 Scans and phased arrays for thresholding defect areas [92], [168], [169]. Whilst the 6

1 dB drop method is widely established it does have limitations. Primarily, the defect
 2 must be larger than the acoustic beam to get an accurate value for the peak response
 3 amplitude [168]. In addition, real defect responses generally do not follow the ideal
 4 defect response curve, often leading to under-sizing defects [170]. To combat this,
 5 alternative amplitude drop thresholds are used in different industrial settings [16],
 6 [171]. These are often component specific and require experimental determination. Li
 7 et al. proposed an alternative method which utilised a generalised regression neural
 8 network and took additional features into account to provide dynamic thresholding of
 9 a C-scan image for more accurate defect sizing than the 6 dB drop method [172].



10

11 *Figure 51: Demonstration of the 6 dB drop method for defect sizing. a) Finding maximum defect*
 12 *response. b) Using the 6 dB loss in maximum amplitude to locate one edge of the defect. c) The*
 13 *corresponding defect edge detected using the 6 dB drop to determine the defect length.*

14 Whilst the 6 dB drop can be used for in-plane defect localisation, depth-wise
 15 localisation requires information from the time trace signal. Cheng et al. showed a
 16 promising method for depth localisation of defects in CFRP panels using different DL
 17 approaches [173] with A-scan signals. They reported a minimum depth relative error
 18 of 9% for the hybrid CNN-LSTM, reducing the error by 96% compared to relying on
 19 the peak-to-peak time-of-flight measurement alone. DL models present an avenue for

1 advancing the automation of NDE data interpretation, and are becoming more
2 prevalent in the literature, especially in the context of defect detection when dealing
3 with images or A-scans [14]. DL is particularly well suited to addressing challenging
4 automation tasks, where a traditional method may not be available, such as defect
5 characterisation and quantification. There are several examples of DL models
6 demonstrating the ability to exceed human performance in certain situations [100],
7 [174].

8 The previous chapter demonstrated the effectiveness of using DL to detect defects
9 from volumetric ultrasonic data. This chapter presents an alternative method for defect
10 sizing using 3D U-Net for volumetric segmentation of ultrasonic data, evaluated
11 against the established 6 dB drop method. U-Net is a DL model, introduced in 2015,
12 which proposed an architecture for medical 2D image segmentation [175]. The authors
13 highlighted the problem that for localisation and classification, pixel level annotation
14 is required which is much more intensive. This makes large, labelled, datasets
15 unreachable. They highlight that the sliding window approach which is often used
16 suffers from two key drawbacks: 1. It is inefficient to run the network for each window,
17 and 2. There is a trade-off between localisation accuracy and use of context. Despite
18 advances in computer vision the U-Net architecture is still widely popular and shows
19 impressive results in many different segmentation tasks [102]. Çiçek et al. [176]
20 extended the original U-Net paper for 3D segmentation of highly variable kidney
21 volumes, giving 3D U-Net, which showed impressive results. By incorporating
22 algorithms with the capability to interpret volumetric data, it is ensured that all spatial
23 and depth wise information is preserved. This approach provides the model with more
24 pertinent features to learn from and eliminates the necessity for image pre-processing

1 and gating. Furthermore, it enables comprehensive 3D defect localisation, a key focus
2 of this study.

3 Section 5.2 provides information on the generation of synthetic training data, and the
4 creation of ground truth segmentation masks. It also covers augmentation applied
5 during training. In Section 5.3, the network architecture, training specifics, and the
6 reference 6 dB drop method for sizing are detailed. The results and discussion are
7 presented in Section 5.4, which is divided into localisation and sizing. To the best of
8 the authors' knowledge, this marks the first utilisation of a 3D U-Net for sizing and
9 localisation of defects in volumetric ultrasonic testing data, offering several
10 advantages:

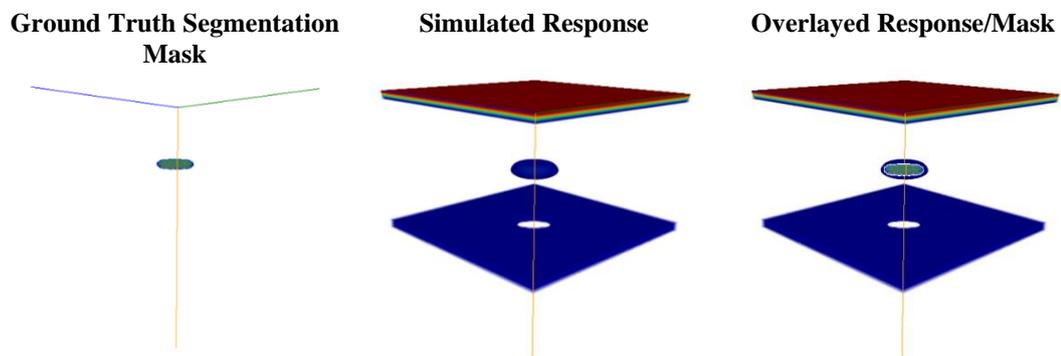
- 11 • Reduced pre-processing times as no thresholding, gating, or generation of
12 images is needed.
- 13 • Results of the developed model trained exclusively on synthetic datasets,
14 outperform industry standard 6 dB drop method for sizing by 35% on
15 experimental test data.
- 16 • Complete localisation of defects within 3D space. Which enables easy
17 extraction of data for downstream processes i.e. testing with Finite Element
18 Analysis.

19 5.2 Data

20 The previous chapter utilised realistic synthetic volumetric data for training. The same
21 synthetic training data was used to train the 3D U-Net model in a fully supervised
22 manner without the need for any experimental training data. The trained model's
23 performance was evaluated against a fully experimental test dataset.

1 5.2.1 Mask Generation

2 An important advantage of employing simulated data as the foundation for training
 3 datasets lies in the capacity to have full control over the simulation input parameters.
 4 This control can be harnessed for the automated generation of ground truth masks
 5 during the training process, a task that would pose more significant challenges when
 6 training models on experimental data. In the context of this study, defect diameter and
 7 depths were utilised to create segmentation masks for defects with a nominal thickness.
 8 Figure 52 provides an illustration of a ground truth defect mask alongside its
 9 corresponding simulation.



10 *Figure 52: (a) The ground truth segmentation mask and (b) the corresponding simulated defect*
 11 *response. The overlay of both the mask and response is shown in (c). Colour mapping and axes are*
 12 *given in Figure 21.*

13 Table 19 provides a summarised description of the datasets created from both the
 14 experimental and synthetic data sources.

15 *Table 19: Summary of the datasets produced.*

Data source	Dataset	Number of datapoints
Simulated defect responses and segmentation masks (300 Flat-Bottom Holes)	Synthetic defective <i>Train (80%)</i>	240
	Synthetic defective <i>Validation (20%)</i>	60
Experimental defect sample (25+15 Flat-Bottom Holes)	Sample 1 ~ Diameters: 3, 4, 6, 7, 9 mm <i>Test</i>	25
	Sample 2 ~ Diameters: 3, 6, 9 mm <i>Test</i>	15

1 5.2.2 Augmentation

2 The generalisability of ML models is a critical aspect of their performance and
3 accounts for differences in the source and target domain. Augmenting the training data
4 improves model generalisability by adding noise at the training stage, reducing the
5 likelihood of overfitting. This often improves performance in the target domain
6 particularly when the target (experimental) domain is different from the source
7 (synthetic) domain.

8 In this study, two domain-specific augmentation techniques that have shown their
9 effectiveness in the previous chapter for augmenting volumetric ultrasonic responses
10 were employed during training. Standard computer vision augmentation methods (e.g.
11 mix-up, cut mix etc.) do not translate directly to UT data as they would impact the
12 underlying signal response. The first type of augmentation is concerned with response
13 magnitude. Magnitude can vary due to various factors unaccounted for in the
14 simulation, such as manufacturing variances in the sample and the UT array probe,
15 wear on the probe and its electrical connections, and the complexities of sound wave
16 scattering between layers. To replicate these variations while maintaining the
17 appropriate data normalisation, the amplitude of each A-scan was adjusted by a
18 constant factor beyond the front wall. This factor was randomly selected from a
19 uniform distribution, resulting in a scaling factor ranging from 80% to 120%.

20 The second augmentation method aims to replicate phase aberration - the variations in
21 ultrasonic travel time between elements [177]. These variations can result from a range
22 of factors, including fluctuations in the sound speed of composite materials due to their
23 anisotropic properties, and deviations from the central frequency for specific elements.

1 To simulate phase aberration, a 1-D interpolation technique was employed to randomly
2 stretch or compress the signal in the time domain. The extent of dilation was randomly
3 determined from a uniform distribution for each A-scan, allowing for dilation of up to
4 ± 300 ns.

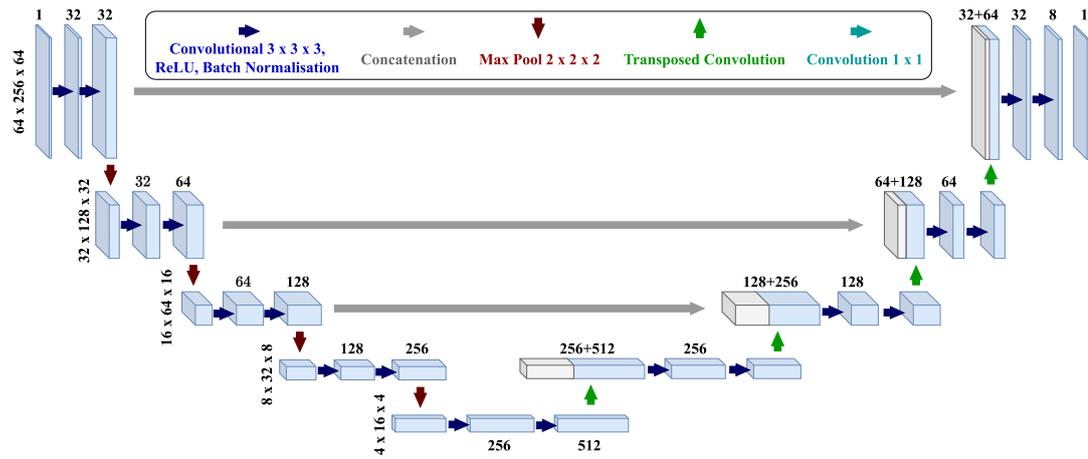
5 The objective in implementing these augmentation methods is to enhance the models'
6 ability to generalise effectively within the experimental domain. The convenience of
7 real-time augmentations allows for their integration into the training process,
8 eliminating the need for additional data collection or preprocessing steps. To maintain
9 consistent data length in the time domain, each A-scan was extended by zero padding,
10 resulting in a length of 1024 samples during training. Subsequently, to mitigate
11 computational demands, each volume was down sampled in the time domain by a
12 factor of 4.

13 5.3 Segmentation Methods

14 5.3.1 Model: Architecture and Training

15 In this chapter, volumetric segmentation was carried out through the training and
16 deployment of a customised 3D U-Net architecture. The design of the architecture
17 drew inspiration from [176], but extended to five convolutional blocks with a sigmoid
18 layer applied to the output. Models of varying number of blocks were tested, and the
19 inclusion of an additional convolutional block resulted in a 40.8% reduction in
20 validation loss. Further optimisation of hyperparameters and architecture may result
21 in a performance increase, however this was outside the scope of this work. A
22 graphical representation of the overarching architectural design is presented in Figure

1 53.



1 5.3.2 Reference Sizing Metric: 6 dB Drop

2 The 6 dB drop criterion represents the prevailing industrial methodology for defect
3 sizing, as documented in the literature [92], [169]. The underlying principle of this
4 method relies on the utilisation of a single transducer to pinpoint the edge of a defect
5 response by detecting the moment when exactly 50% of the energy is reflected by the
6 defect, corresponding to a 6 dB decrease in amplitude [16]. The technique is repeated
7 on the opposing boundary of the defect, and the resultant displacement results in the
8 measured defect's length Figure 51. This fundamental principle can be further
9 extended to encompass the sizing of defects from amplitude C-scan images produced
10 from employing phased array transducers [92], [168], [169].

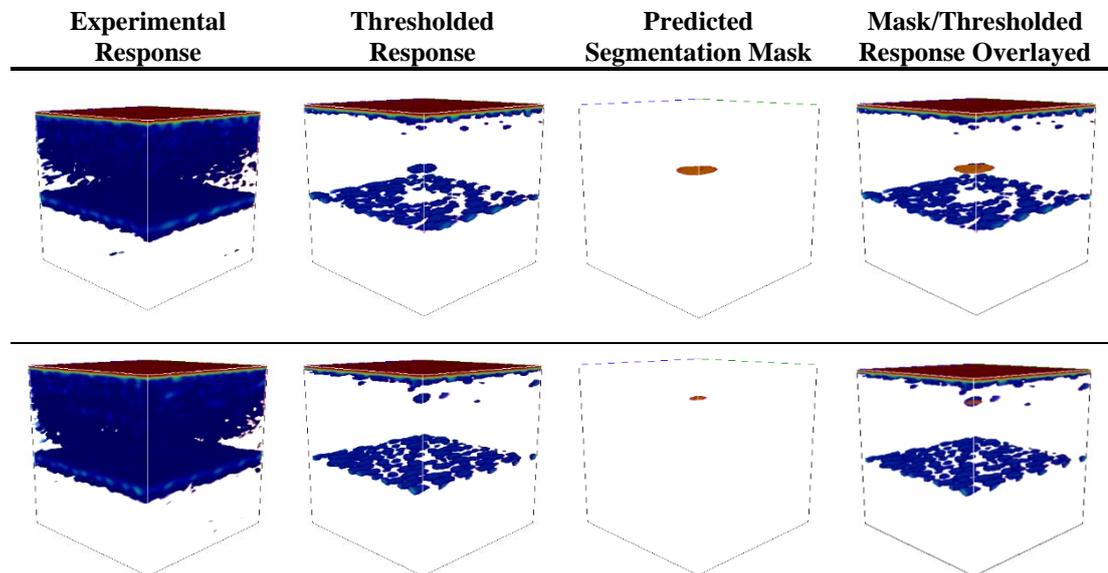
11 In this research, sizing in the fibre plane using the segmentations predicted by the U-
12 Net model are compared with the established 6 dB drop. Given the prior knowledge
13 of the manufactured defects present in the reference sample, recognised as circular
14 FBH, and to mitigate variations in diameters through the component thickness, the
15 defect diameters were computed based on the maximum segmented area through the
16 depth of the sample, based on pixel summation, using the formula outlined in equation
17 16.

$$Diameter = 2 \sqrt{\frac{\sum_{pixels} \left(\frac{\text{Max}_{through\ depth} (Segmented\ Volume)}{\pi} \right)}{}} \quad 16$$

18 5.4 Results and Discussion

19 Figure 54 presents examples of experimental defect responses and their corresponding
20 segmentation masks as generated from the 3D U-Net. Along with demonstrating
21 results, these visualisations could be used by human operators to sense-check the

1 models' predictions and reject inaccurate model predictions easily, providing the
2 possibility for the method to contribute to a human-in-the-loop semi-autonomous NDE
3 system.



4 *Figure 54: Example 9 and 3 mm defects respectively; their experimental ultrasonic volumetric*
5 *responses, thresholded responses (amplitudes >10% of maximum response for defect visualisation),*
6 *and their corresponding predicted segmentations. Colour mapping and axes are given in Figure 21.*

7 5.4.1 Localisation

8 The localisation of defects in 3D can be deconstructed into two primary components:
9 in-plane localisation and through-thickness depth-wise localisation. It is important to
10 note that the widely adopted 6 dB drop criterion only addresses in-plane localisation
11 and it does not provide information regarding depth-wise localisation.

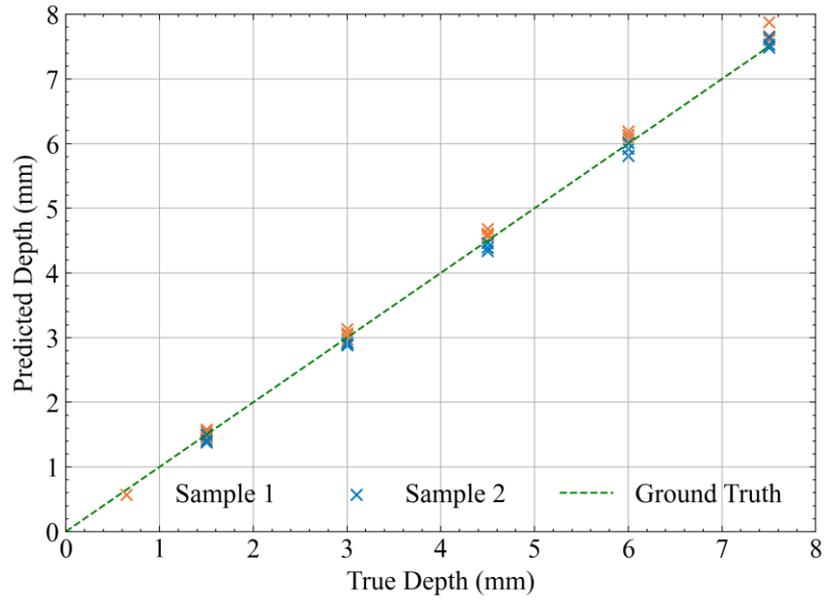
12 The 6 dB drop method can produce inaccuracies in defect sizing as discussed in section
13 5.4.2. However, the circular shape of the test defects ensures that any errors in sizing,
14 which might cause changes in diameter, will have minimal impact on the position of
15 the defect's centroid. Obtaining an accurate ground truth for in-plane localisation less
16 than 1.0 mm is infeasible due to the cumulative positional errors introduced throughout

1 the experimental setup. Therefore, this research uses the 6 dB drop criterion as a
2 reference standard, to validate the agreement between the 6 dB and U-Net's in-plane
3 localisation. By comparison, experimental through-thickness depth measurements for
4 defects are considerably easier to acquire, which allows for a direct assessment.

5 *5.4.1.1 Depth*

6 The segmentation of volumetric ultrasonic data offers a distinct advantage compared
7 to the 6 dB drop method due to its capacity for depth-wise localisation. This eliminates
8 the need to employ multiple data types for characterisation, such as amplitude and
9 time-of-flight C-scans. The depth-wise position determined from the segmented
10 volume is represented by the mean segmented depth. In Figure 55 the predicted defect
11 depth is compared to the true measured depth of the reference defects. The segmented
12 volumes demonstrate a excellent level of accuracy in depth-wise localisation, as
13 evidenced by a Mean Absolute Error (MAE) of 0.08 mm. This precision can be
14 attributed to the substantially higher sampling rate in the temporal domain in contrast
15 to the spatial domains, resulting in significantly superior temporal resolution when
16 compared to spatial resolution.

1

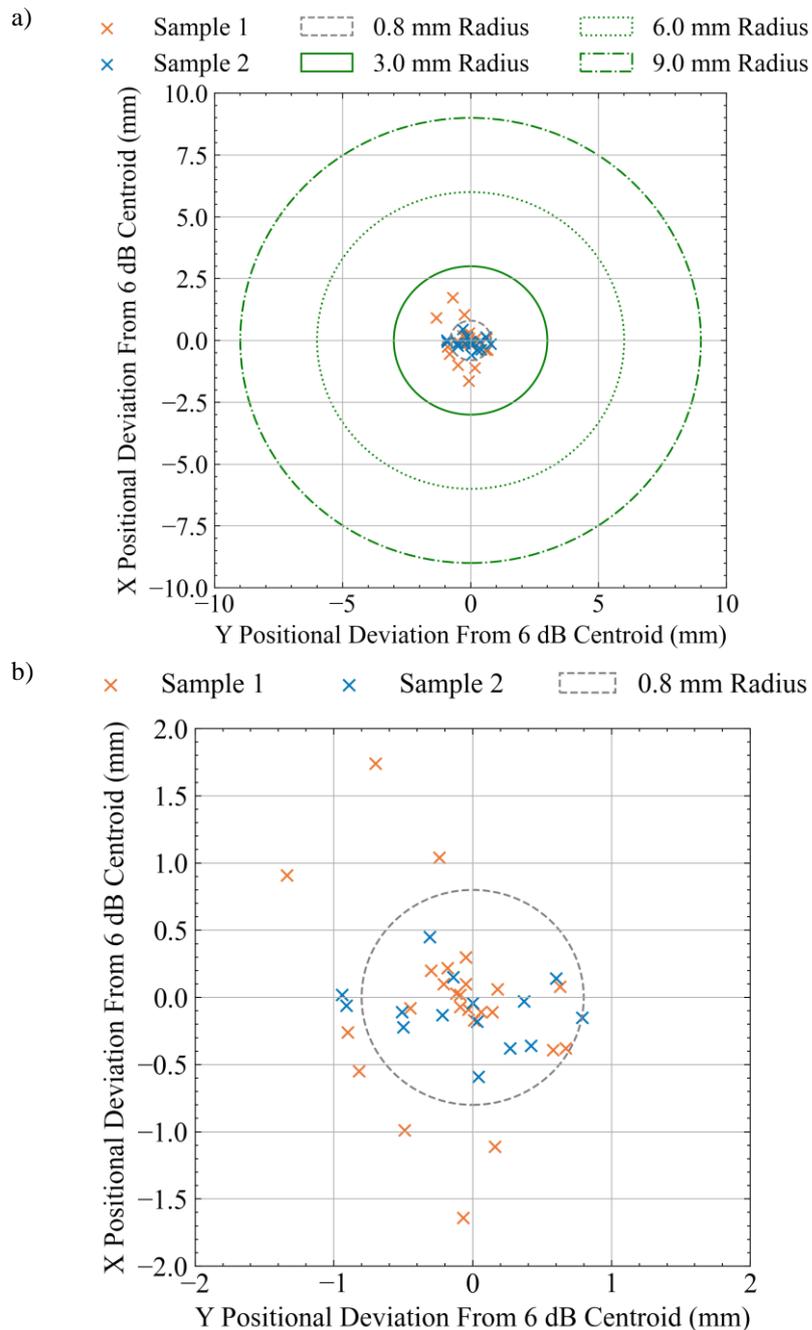


2

3 *Figure 55: Depth localisation results.*

4 *5.4.1.2 In plane*

5 The performance of the models' in-plane localisation is quantified by measuring their
6 deviation from the centroid area compared to the 6 dB drop. Figure 56 visually presents
7 the centroid deviation with reference to defect sizes and the pitch of array elements.
8 As depicted in the figure, 75% of the variations (30 out of 40 defects) are below the
9 0.8 mm array pitch.



1 *Figure 56: In plane localisation results compared to 6 dB drop with reference to defect diameters (a)*
 2 *and the expanded (b), which shows the reference to the array pitch more clearly.*

3 Table 20 provides a comprehensive overview of the in-plane localisation outcomes,
 4 including the absolute distance between the 6 dB criterion and the centroid determined
 5 by the model. The Mean Absolute Error (MAE) of 0.57 mm demonstrates a substantial
 6 concordance with the established industrial benchmark represented by the 6 dB drop

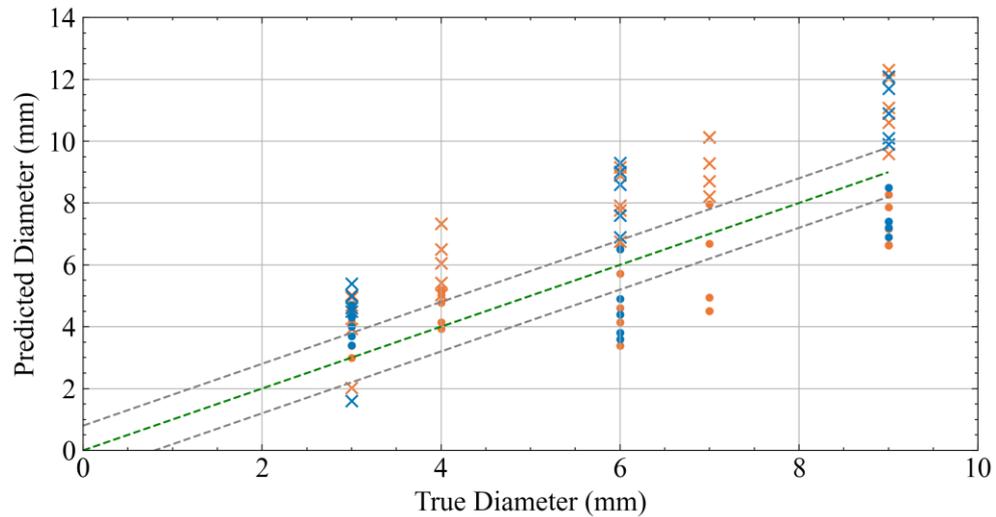
1 criterion for in-plane localisation. Notably, the MAE, being less than the 0.8 mm array
 2 pitch, which establishes the spatial resolution, underscores the robust agreement
 3 between the model-based in-plane localisation and the standard reference.

4 *Table 20: Complete model centroid deviation results from the 6 dB drop.*

Sample 1				Sample 2			
Defect Diameter (mm)	Deviation from X (mm)	Deviation from Y (mm)	Absolute displacement (mm)	Defect Diameter (mm)	Deviation from X (mm)	Deviation from Y (mm)	Absolute displacement (mm)
9	-0.05	0.10	0.11	9	0.42	-0.36	0.55
	0.58	-0.39	0.70		-0.50	-0.22	0.54
	-0.49	-0.99	1.10		-0.94	0.02	0.94
	0.18	0.06	0.19		-0.31	0.45	0.55
	0.01	-0.17	0.18		-0.14	0.15	0.21
7	0.63	0.08	0.63	6	0.37	-0.03	0.37
	-0.12	0.03	0.13		-0.51	-0.11	0.52
	-0.07	-1.64	1.64		-0.91	-0.06	0.91
	-0.18	0.22	0.29		0.00	-0.04	0.04
	-0.09	0.02	0.09		-0.22	-0.13	0.25
6	-0.03	-0.09	0.09	3	0.27	-0.38	0.47
	0.16	-1.11	1.12		0.04	-0.59	0.59
	0.67	-0.38	0.77		0.79	-0.15	0.80
	-0.24	1.04	1.06		0.03	-0.18	0.18
	-0.09	-0.07	0.12		0.60	0.14	0.61
4	-1.34	0.91	1.62	3	-0.82	-0.55	0.98
	-0.45	-0.08	0.46		0.14	-0.11	0.18
	-0.05	0.30	0.30		-0.21	0.10	0.23
	-0.70	1.74	1.88		-0.30	0.20	0.36
	-0.90	-0.26	0.94		0.06	-0.11	0.13
Mean Average Error (MAE)			0.61				0.50
Total MAE							0.57

1 5.4.2 Sizing

× Sample 1 (UNet) × Sample 2 (UNet) - - - Ground Truth
• Sample 1 (6 dB) • Sample 2 (6 dB) - - - Array Pitch (0.8 mm)



2

3 *Figure 57: Sizing results for the 6 dB drop method and U-Net predictions.*

4 Figure 57 presents a summary of the defect diameter predictions from the original
5 segmented from the 6 dB drop and U-Net areas for each defect compared to the known
6 ground truth area.

7 5.4.2.1 6 dB approach

8 With a MAE of 1.35 mm, our findings demonstrate a reasonable degree of accuracy,
9 which, when coupled with suitable safety factors, is likely to be deemed adequate in
10 industrial settings. Nevertheless, it is crucial to acknowledge that real-world responses
11 often deviate significantly from the ideal defect response, leading the 6 dB drop
12 approach to systematically underestimate defect sizes [170].

13 This limitation has prompted the utilisation of alternative amplitude drop methods for
14 sizing in industry [16], [171], wherein the threshold values are frequently determined
15 through experimental calibration. The experimental data collected for this research
16 corroborates this tendency for under sizing defect responses (Figure 57), particularly

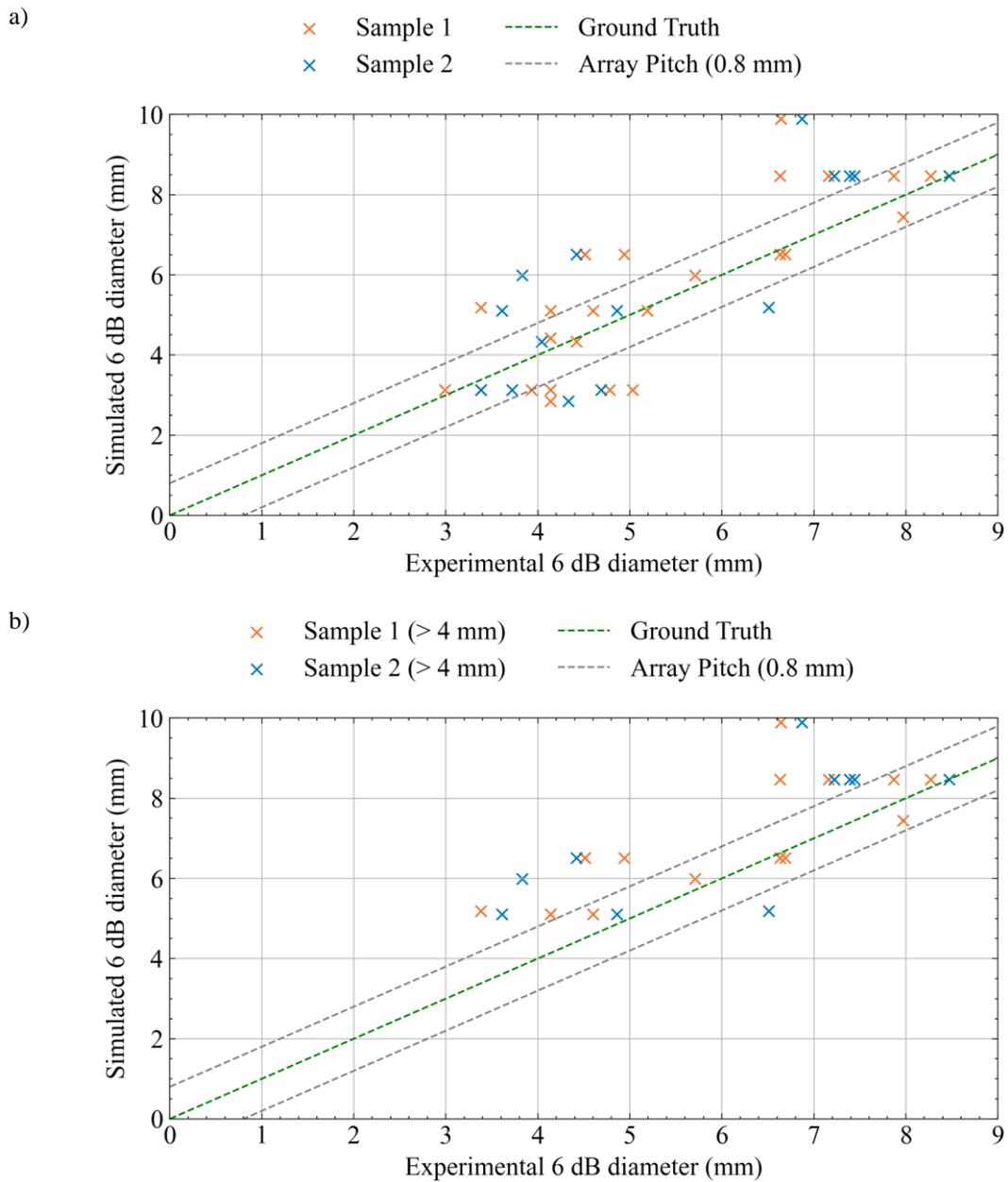
1 for defects exceeding 4.0 mm in diameter. Which exhibit a mean undersize of 1.37
2 mm. Conversely, our results reveal a tendency to oversize defects of 3.0 and 4.0 mm
3 in diameter.

4 It is essential to recognise that any sizing method relying on maximum amplitude
5 necessitates the defect to be substantially larger than the acoustic beam to accurately
6 ascertain the maximum acoustic response. The experimental setup employed in this
7 work, which utilised a 4-element sub-aperture with an element pitch of 0.8 mm,
8 resulted in an effective transducer width of 3.2 mm. Considering this in combination
9 with the spatial resolution limitations imposed by the fixed 0.8 mm pitch for each beam
10 step, along with accounting for any beam spread, it was determined that the
11 experimental setup was inadequate for the precise sizing of defects measuring 4 mm
12 or less when employing an amplitude drop method. This leads to the average
13 oversizing of 3.0 and 4.0 mm defects of 0.82 mm, counter to the expectation of the 6
14 dB amplitude drop under sizing defects. The under-sizing of defects larger than the
15 diameter of the acoustic beam can likely be attributed to the curved edges of the
16 defects. Since these defects do not maintain orthogonality to the sizing axis, the result
17 is a diminished reflector when accounting for the three-dimensional nature of the
18 response. This leads to a 6 dB decrease in acoustic energy closer to the centre of the
19 defect, rather than at the true defect edge. The inconsistency of this method and the
20 need for varying amplitude drop thresholds adds to the complexity of consistent defect
21 sizing in industry and could be a concern for safety critical parts.

1 5.4.2.2 *U-Net*

2 The initial segmentation of U-Net masks yielded a MAE of 2.09 mm, which represents
3 a 55% increase in error compared to the 6 dB drop method. As depicted in Figure 57,
4 there was a consistent tendency to overestimate defect sizes across the range of
5 diameters. This observation underscores the U-Net approach's reduced reliance on
6 absolute peak amplitudes and its ability to deliver more consistent performance across
7 a range of defect sizes, even when the defect size is not greater than the width of the
8 acoustic beam. Moreover, it is worth noting that in numerous industrial applications
9 for safety-critical components, it is preferable to overestimate rather than
10 underestimate defect sizes.

11 It is imperative to delve into the reasons for the model's consistent trend of oversizing
12 defects. Given that the model exhibited convergence during training on synthetic
13 datasets, the overestimation observed during testing hints at a domain disparity
14 between training on synthetic data and testing on experimental data. To elucidate this
15 distinction, a comparative analysis using the 6 dB drop method was conducted between
16 the responses derived from the experimental data and those generated by synthetic data
17 for corresponding defect sizes and depths, as illustrated in Figure 58 (a). As previously
18 noted, defects of 3 and 4 mm in diameter were too small to be accurately sized using
19 this experimental setup and the 6 dB drop method, and thus, were excluded from this
20 analysis (Figure 58 (b)).



1 *Figure 58: Comparison of sizing for synthetic and experimental data for all defect sizes (a) and defect*
 2 *diameters above 4 mm (b).*

3 The comparison reveals that synthetic responses tend to yield larger defect sizes than
 4 experimental responses when employing the 6 dB drop method. Since our model's
 5 ground truth during training was based on synthetic response masks, it becomes
 6 apparent why there exists a propensity to overestimate defect sizes in our model;
 7 simulated responses tend to produce spatially larger defect responses than

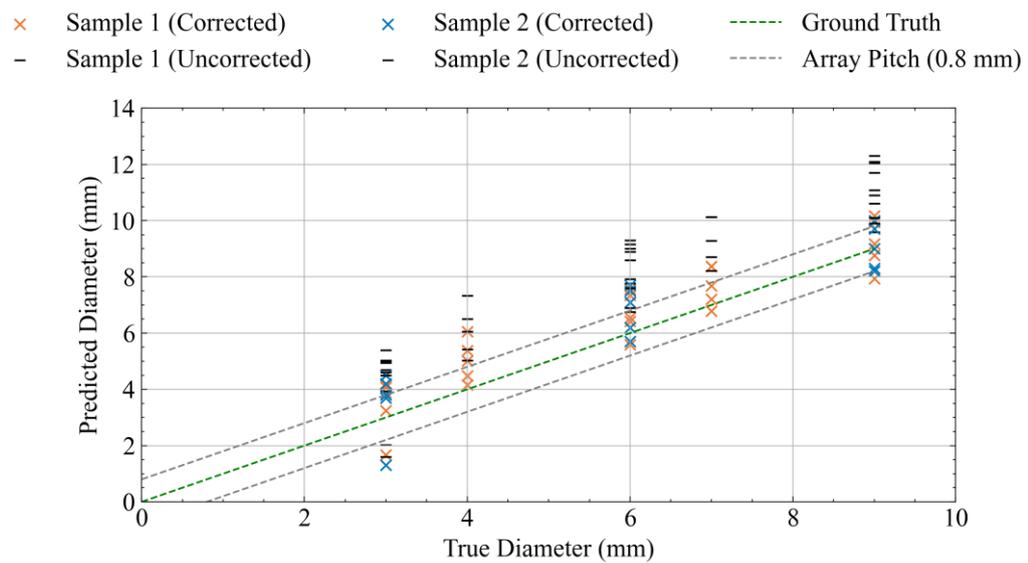
1 experimental. This disparity can be attributed to a combination of factors, CFRP acting
2 as an acoustic collimator is likely the most impactful cause. Due to the anisotropic
3 nature of CFRP, attenuation increases significantly as the propagation angle to normal
4 increases [178]. This effect was not captured in the semi-analytical simulations, as it
5 did not account for the anisotropic acoustic attenuation. Consequently, these
6 simulations resulted in increased acoustic beam spread and larger defect responses in
7 the simulated domain compared to the experimental domain.

8 *Table 21: Comparison in sizing of experimental and synthetic responses with the 6 dB drop method.*
9 *Where R is the ratio of synthetic response to experimental response. \bar{R} gives the mean average of R*
10 *for a given defect diameter.*

Sample 1					Sample 2				
Defect Diameter (mm)	6 dB drop diameter		$R = \frac{synth}{exp}$	\bar{R}	Defect Diameter (mm)	6 dB drop diameter		$R = \frac{synth}{exp}$	\bar{R}
	Experimental response	Synthetic response				Experimental response	Synthetic response		
9	6.63	8.47	1.28	1.21	9	8.47	8.47	1.00	1.18
	7.87	8.47	1.08						
	7.16	8.47	1.18						
	8.27	8.47	1.02						
	6.64	9.89	1.49						
7	4.51	6.51	1.44	1.22	6	6.51	5.19	0.80	1.26
	4.94	6.51	1.32						
	4.51	6.51	1.44						
	6.69	6.51	0.97						
	7.97	7.44	0.93						
6	3.38	5.19	1.53	1.18	3	3.38	3.13	0.93	0.83
	4.14	5.11	1.23						
	4.60	5.11	1.11						
	5.71	5.99	1.05						
	6.63	6.51	0.98						
4	4.78	3.13	0.65	0.82		4.04	4.33	1.07	
	3.93	3.13	0.80						
	5.03	3.13	0.62						
	4.14	4.42	1.07						
	5.19	5.11	0.98						
3	3.38	3.13	0.93	0.88		3.38	3.13	0.93	
	4.14	3.13	0.76						
	2.99	3.13	1.05						
	4.14	2.85	0.69						
	4.42	4.33	0.98						

11 To rectify this domain disparity, a constant R (given as the ratio between synthetic and
12 experimental response diameters) can be computed based on the disparities in the 6 dB
13 drop method for each defect, as detailed in Table 21. The correction factor is
14 determined by the mean R calculated across defects with diameters exceeding 4 mm.

1 The resulting correction factor is 1.21 (as previously noted, defects with diameters ≤ 4
 2 mm were deemed too small to be accurately sized with this experimental setup and the
 3 6 dB drop method). When this correction factor is applied to the defect sizes
 4 determined from the model predictions, a significantly improved agreement with the
 5 known defect sizes is achieved, as shown in the Figure 59. This correction results in a
 6 reduction of the MAE by 58%, bringing it down to 0.88 mm. Consequently, defect
 7 sizing exceeds the accuracy of the 6 dB drop method significantly by 35%.



8
 9 *Figure 59: Corrected sizing results for the U-Net predictions.*

10 5.5 Out-of-Distribution Testing

11 The model utilised in this study underwent supervised training, a process commonly
 12 employed in ML where the model learns patterns and relationships from labelled
 13 training data. Typically, supervised models excel in predicting examples that fall
 14 within the distribution of the training data. In this research, the training data was
 15 generated from simulated data, designed to replicate the geometry and characteristics
 16 of FBHs observed during testing. However, it's important to note that naturally
 17 occurring defects may exhibit significant variations in geometry. While defects in

1 composite materials often manifest in-plane, their characteristics can vary widely.
2 Training a model to generalise across such diverse conditions would likely necessitate
3 a substantially larger and more varied training dataset.

4 To give an insight into the model's generalisability to different defect types and
5 geometries, its sizing and in-plane localisation performance was also evaluated on 15
6 square, 6 mm wide, Polytetrafluoroethylene (PTFE) inserts from a different sample.
7 These different defect types can be considered a semantic distribution shift from the
8 FBHs seen in training, and therefore categorised as “near Out-of-Distribution” [179].
9 5 of the defects were located within 2 plies of the front wall surface, 5 were in the
10 middle of the component and 5 were located within 2 plies of the back surface. This
11 highlighted a limitation of the model, in that defects very near to front or backwall
12 surfaces proved too challenging to segment. To account for this the front and back wall
13 responses were removed. The results of this evaluation are detailed in Table 22,
14 providing valuable insights into the model's performance for defects out of its training
15 distribution.

16 *Table 22: Sizing and in-plane localisation results for out of distribution test defects.*

True	Width (mm)				Deviation in X (mm)	Deviation in Y (mm)	Absolute Displacement (mm)
	Predicted	Predicted Error	Corrected	Corrected Error			
6	5.18	-0.82	4.28	-1.72	-0.18	-1.6	1.61
	6.69	0.69	5.53	-0.47	-0.24	-0.14	0.28
	8.98	2.98	7.42	1.42	0.08	0.47	0.48
	8.98	2.98	7.42	1.42	0.06	-0.45	0.45
	6.65	0.65	5.49	-0.51	-0.09	-0.18	0.20
	6.60	0.60	5.45	-0.55	-0.15	-0.21	0.26
	1.96	-4.04	1.62	-4.38	2.22	1.86	2.90
	7.11	1.11	5.88	-0.12	-0.44	0.07	0.45
	8.80	2.80	7.27	1.27	-0.03	0.51	0.51
	4.23	-1.77	3.50	-2.50	-0.42	-1.23	1.30
	5.99	-0.01	4.95	-1.05	-1.44	-0.69	1.60
	3.10	-2.90	2.56	-3.44	-0.82	-0.6	1.02
	8.80	2.80	7.27	1.27	0.09	0.1	0.13
	6.50	0.50	5.37	-0.63	-0.5	-0.07	0.50
	5.06	-0.94	4.18	-1.82	-1.04	-0.85	1.34
MAE		1.71		1.50			0.86

1 Whilst this does not serve to test the wide variety of defect geometries which occur
2 naturally, the analysis provides insights into the model's efficacy to generalise to
3 defects seen outside of the training distribution. A MAE of 1.50 mm for defect width
4 was observed after the correction factor was applied. This represents a 71% increase
5 in sizing error compared to defects within the training distribution. In-plane
6 localisation performed much better with an MAE of 0.86 mm, a 51% increase in
7 localisation error however, this is still in line with the element pitch. As these inserts
8 were embedded pre-cure it is not possible to extract a true measured ground truth of
9 through-thickness localisation as is the case for FBHs, this analysis has therefore been
10 omitted. These results highlight a limitation of the proposed model and training
11 regime. Synthetic data effectively addresses the challenge of acquiring labelled
12 training data and can train an effective model for expected defects. However, for
13 defects outside of the training distribution, there is a significant drop in performance.
14 Furthermore, for edge cases such as defects near geometric features, additional pre-
15 processing steps may be required; further limiting the generalisability of the model.
16 Most DL work applied to NDE has a very specific and controlled application, and
17 there's a notable challenge in finding models that have demonstrated effective
18 generalisability across various materials, defect types, and component geometries. In
19 future work the authors hope to further expand the synthetic training data to encompass
20 a far wider range of defect types and geometries whilst also simulating different
21 component geometries, accounting for edge cases such as near front and back wall
22 responses. With a much larger synthetic training set it is hoped that a far more
23 generalisable model can be trained.

1 5.6 Conclusion

2 This chapter proposes the use of a 3D U-Net to size and localise defects in CFRP by
3 segmenting volumetric ultrasound data. Defect sizing is a crucial piece of information
4 for evaluating defects against standards and acceptance criteria whilst accurate
5 localisation is beneficial if re-work is required. A key benefit of the approach is that
6 the use of volumetric ultrasound data allowed for through-thickness and in-plane
7 localisation whilst removing the requirement for gating and reducing pre-processing.
8 Such gating and amplitude threshold selection is often performed manually by the
9 NDE operator leading to significant data errors if gates and thresholds are incorrectly
10 set.

11 Simulations were used to generate synthetic data and ground truth segmentation masks
12 for training. This was a key requirement and allowed for the training of a segmentation
13 model in a fully supervised manor. Experimentally collected ultrasonic responses from
14 manufactured reference defects were used for testing. Sizing and in-plane localisation
15 were evaluated against the widely accepted 6 dB drop standard, and through thickness
16 localisation was compared to the measured ground truth.

17 The models' depth-wise localisation showed excellent results with a MAE of 0.08 mm.
18 In-plane localisation had good agreement with the accepted 6 dB drop standard with a
19 MAE of 0.57 mm. The significant resolution differences in the spatial and temporal
20 domains resulted in differences of error scales for in-plane and depth wise localisation.
21 This is a limitation of using a fixed pitch array, but the errors in-plane are reasonable
22 when compared to the array pitch of 0.8 mm, which is the limiting factor for spatial
23 resolution.

1 The 6 dB drop consistently undersized experimental defects greater than 4 mm in
2 diameter whilst the U-Net produced segmentation masks that consistently oversized
3 defects. The U-Net's oversizing was consistent across defect sizes, showing that it was
4 less reliant in maximum amplitude, making it more robust to sizing defects smaller
5 than the width of the acoustic beam, which proved inconsistent for amplitude drop
6 methods. In industrial settings other factors such as parallelism of defects to the
7 inspection surface would also impact the maximum signal response, introducing
8 further inaccuracies to amplitude-based sizing methods. Upon investigation of the
9 synthetic and experimental data domains it was evident that the experimental
10 responses gave rise to smaller defect responses. By correcting for this disparity
11 between the source and target domain using a single correction factor it became
12 possible to reduce the MAE for defect sizing from 2.09 mm to 0.88 mm. The corrected
13 defect sizing from U-Net gave a 35% reduction in MAE sizing compared to the
14 commonly accepted 6 dB drop method. Despite this the 6 dB drop method is based on
15 physical understanding of defect responses and whilst it has limitations the results are
16 directly explainable. Whilst the generation of a segmentation map aides in providing
17 some explanation for defect sizing and is more interrogatable than a regression model,
18 it still relies on a deep learning approach which is less explainable than an amplitude-
19 based threshold. Even so, a clear benefit of the volumetric segmentation is that it can
20 be translated directly into a computer-aided design file which could dramatically
21 increase the efficiency of subsequent report generation and simulation-based testing
22 of components.

23 While the study demonstrated promising results in defect sizing and localisation, it's
24 essential to acknowledge a current limitation: the absence of real defects from

1 industrial manufacturing processes for testing. Although out-of-distribution cases were
2 examined using PTFE inserts, evaluating performance on naturally occurring defects
3 would be advantageous. The out-of-distribution testing underscored the necessity for
4 a more extensive distribution of training data. Future work aims to address this by
5 expanding the dataset, incorporating different probe frequencies, and testing on
6 naturally occurring defects with irregular shapes. Furthermore, the objective is to
7 integrate this research with previous work on defect detection to develop an end-to-
8 end system for automated NDE data processing in industrial manufacturing
9 environments.

1 6 Self-Supervised Learning Segmentation

2 6.1 Introduction

3 While the 6 dB drop method remains a cornerstone in defect sizing, advancements in
4 DL techniques offer promising avenues for enhancing defect localisation and sizing
5 accuracy. By leveraging the capabilities of fully supervised volumetric methods, the
6 previous chapter has made significant strides in addressing the complexities associated
7 with defect characterisation and quantification. Analysing volumetric data offers a
8 wealth of information for defect characterisation that surpasses what can be achieved
9 through individual image analysis alone. Additionally, this approach proves
10 advantageous in reducing manual preprocessing tasks, such as gating out structural
11 responses, which are often labour-intensive. Moreover, volumetric segmentation
12 masks of components open avenues for various downstream tasks, such as building
13 digital twins [166], performing FEA, and alleviating the burden of report generation.

14 However, the previous fully supervised method faces limitations due to the necessity
15 of hard to acquire labelled training data, a common requirement for any fully
16 supervised training approach. In many NDE applications, obtaining accurately
17 annotated labelled datasets of real defects is challenging and often not possible, and is
18 one of the main barriers for applying DL to NDE [14]. To address this scarcity issue,
19 work in previous chapters utilised synthetic training data generated from simulations.
20 By controlling simulation parameters, accurate labelling of segmentation ground truth
21 could be achieved automatically. However, accurately simulating the full distribution
22 of defects and their variations is computationally demanding, and it is challenging to
23 ensure fidelity between the simulated and experimental domains [180]. As a result, if

1 defect segmentation could be performed without the reliance on large positively
2 labelled datasets, it would offer significant advantages.

3 Self-Supervised Learning (SSL) [181], [182] is a method for training DL models in a
4 supervised manner without the need for labelled training data. Generally, labels are
5 generated through auxiliary tasks or by leveraging inherent structures in the data itself,
6 enabling the model to learn meaningful representations without explicit human
7 annotation. By training in a supervised manner, models are often able to learn more
8 detailed feature representations than unsupervised approaches. SSL introduces the
9 ability to leverage large amounts of unlabelled data, reducing the need for costly and
10 time-consuming annotation. Its versatility has meant that it has been applied broadly
11 from computer vision to natural language processing and has demonstrated impressive
12 performance in many notable DL tasks, such as with large language models [183],
13 [184] and large vision models [185], [186]. In general, SSL works by formulating
14 pretext tasks that require the model to predict certain aspects of the input data based
15 solely on the input itself. These pretext tasks are designed to be easily computable from
16 the raw data without the need for external annotations. By solving these pretext tasks,
17 the model learns to extract meaningful features and representations from the data,
18 which can then be transferred to downstream tasks. For example, in natural language
19 processing, the model may be tasked with predicting missing words in a sentence (e.g.,
20 masked language modelling [187]) or predicting the next word in a sequence (e.g.,
21 language modelling [183], [184]). Similarly, in computer vision, SSL tasks may
22 involve predicting the rotation, colorization, or spatial arrangement of patches within
23 an image [188], [189], [190], [191]. Another prevalent strategy in SSL is contrastive
24 learning, where the model learns to differentiate between positive and negative pairs

1 of data samples. By maximising the similarity between positive pairs (e.g., different
2 augmentations of the same image) while minimising the similarity between negative
3 pairs (e.g., augmentations of different images), such as Siamese networks [192].

4 This chapter presents a novel approach aimed at advancing the detection, localisation,
5 and segmentation of defects within ultrasonic volumes. The approach leverages the
6 capabilities of SSL coupled with a 1D CNN to achieve 3D segmentation of defects
7 from volumetric ultrasonic testing data of composite components. We employ pretext
8 learning to predict distributions of amplitudes from ultrasonic series. During the
9 inference stage, the pre-trained 1D CNN is deployed to flag any regions within the
10 component that exhibit anomalous behaviour. This process capitalises on the insights
11 gleaned from the pretext learning task, where the model familiarises itself with clean
12 samples. Unlike traditional approaches, which often necessitate extensive positive
13 training examples, this methodology operates on the principle of anomaly detection
14 rather than classification into specific defect classes. By reframing the problem as one
15 of anomaly detection, rather than attempting to categorise defects into predefined
16 classes, we circumvent the significant challenge of acquiring an extensive set of
17 positive training examples. Moreover, this approach offers the advantage of being
18 defect-agnostic, thereby mitigating concerns regarding the generalisability of the
19 model to novel defects. This characteristic alleviates the need for meticulous fine-
20 tuning and ensures that the model remains robust across a range of defect types, as
21 long as they produce an anomalous response, eliminating the burden of adapting the
22 system for each new defect type encountered. Overall, by adopting this innovative
23 methodology, critical limitations in the application of DL to NDE are addressed,
24 paving the way for more efficient and robust defect detection in ultrasonic inspection

1 processes whilst giving information on not just detection but also complete volumetric
2 localisation and segmentation which has only been previously achieved with fully
3 supervised training [193].

4 Section 6.2 of this chapter outlines the data used and any pre-processing requirements.
5 In Section 6.3, the pretext learning is detailed. The results and discussion are presented
6 in Section 6.4.2. This, for the first time, introduces an SSL method for volumetric
7 defect segmentation of ultrasonic testing data, offering several advantages:

- 8 • No defective training data is required.
- 9 • Minimal preprocessing.
- 10 • Geometric features are retained.
- 11 • Enhanced generalisability for defect detection is achieved by reframing the
12 problem as anomaly detection.

13 6.2 Data

14 Additional CFRP samples of varying thicknesses, supplied by Spirit AeroSystems,
15 were employed in this study. For pretext learning, samples verified as defect-free
16 through ultrasonic inspection, analysed by an NDE operator, were selected. These
17 samples were segregated into distinct datasets for training, validation, and testing
18 purposes. During the inference phase, the previously introduced defective samples
19 with FBHs were used for testing (see section 2.2). Additionally, as used for the out-of-
20 distribution test for the supervised segmentation method, a final stepped sample with
21 square 6 mm wide PTFE inserts presents a more challenging geometry and defect
22 responses, offering a representation closer to naturally occurring defects. A summary
23 of the samples and dataset characteristics is presented in Table 23.

1 Data acquisition was conducted using an automated ultrasonic system centred around
2 a 64 element, 5MHz Olympus linear phased array roller probe, with an element pitch
3 of 0.8 mm. Further details regarding the experimental setup and composite samples
4 can be found in section 2.2. The segmentation methodology employed in this study
5 was intentionally designed to be versatile, thereby maximising its applicability across
6 different scenarios. Consequently, only minimal generic data preprocessing was
7 applied. The data preprocessing involved enveloping the signal using the Hilbert
8 Transform (see section 2.3). Notably, techniques such as TCG, peak-alignment, and
9 gating out of geometric features were deliberately omitted. This helps to ensure more
10 consistent and reproducible results across different experimental setups, as it does not
11 require adjusting for specific setups or components, leading to less calibration and
12 setup time. It also reduces the risk of human error and biases that can be introduced
13 during these (potentially) manual processes. Consequently, this leads to more robust
14 and reliable data analysis.

15 *Table 23: Summary of samples used.*

	Sample	Thickness (mm)	Dataset Size [Probe Time Frames]	Details
<i>Pretext Learning</i>	Clean 1	2.75	122 350 260	Training
	Clean 2	4.25	122 450 260	
	Clean 3	4.25	122 450 260	
	Clean 4	6.00	122 600 260	
	Clean 5	6.00	122 600 260	Validation
	Clean 6	8.60	122 700 260	Test
<i>Inference</i>	Defective 1	8.60	183 700 260	15 FBH
	Defective 2	8.60	305 700 230	25 FBH
	Defective 3	7.50,9.60,11.80	488 1050 112	15 Inserts

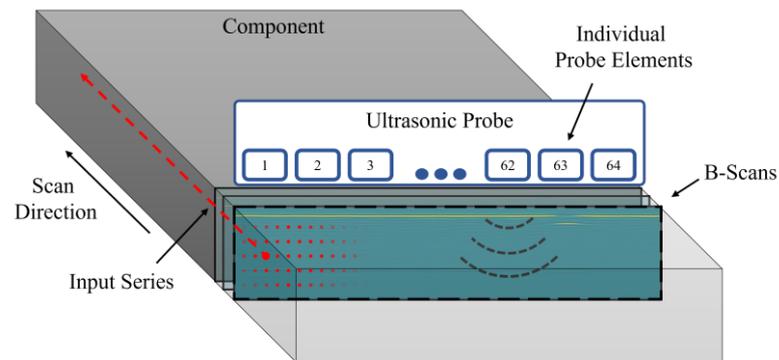
16

1 6.3 Pretext Learning

2 Typically, ultrasonic signals are considered as A-scans, which is often appropriate.
3 However, when inspecting large components C-scan images often provide the most
4 impactful and clear information about the structure of the component. C-scans
5 compress information in the time-domain either through amplitude imaging or
6 amplitude indexing (Time-of-Flight). This loss in temporal information is done as a
7 result of trying to maximise comparative spatial information, which can often be the
8 most helpful when evaluating components of known geometry, particularly for
9 composite components where structural noise can show significant variability in the
10 time-trace. Operators will often produce C-scan images at different thickness gates to
11 compare different through-thickness slices. Intuitively therefore it can be deduced that
12 considering comparative spatial information for given time-windows is as, if not more,
13 important than direct temporal comparisons along a single time trace. Ideally spatial
14 comparisons would be conducted at every depth but for manual interpretation this is
15 intractable due to the substantial amount of data produced.

16 Mechanised linear phased array scanning is commonly employed in the inspection of
17 large-scale industrial components [7], [8]. Arrays operate within acceptable tolerances
18 of element sensitivity. Whilst compensation can be done to account for inter-element
19 variations within arrays through calibration, some level of variation between elements
20 is likely to remain during scanning. Considering this and the importance of spatial
21 comparison discussed earlier, it follows that it is appropriate to analyse a singular
22 series through a component at a particular time step (Figure 60), and is the basis for
23 this work. This approach minimises variations from array elements by concentrating

1 on the same reception element or sub-aperture, while facilitating spatial comparison at
2 a given depth. Furthermore, this methodology can be easily conceptualised as a self-
3 supervised learning task. By treating it as a 1D task, we not only reduce model size
4 and complexity but also benefit from a substantial increase in training data abundance
5 compared to using 2D or 3D training sets from the same sample availability.

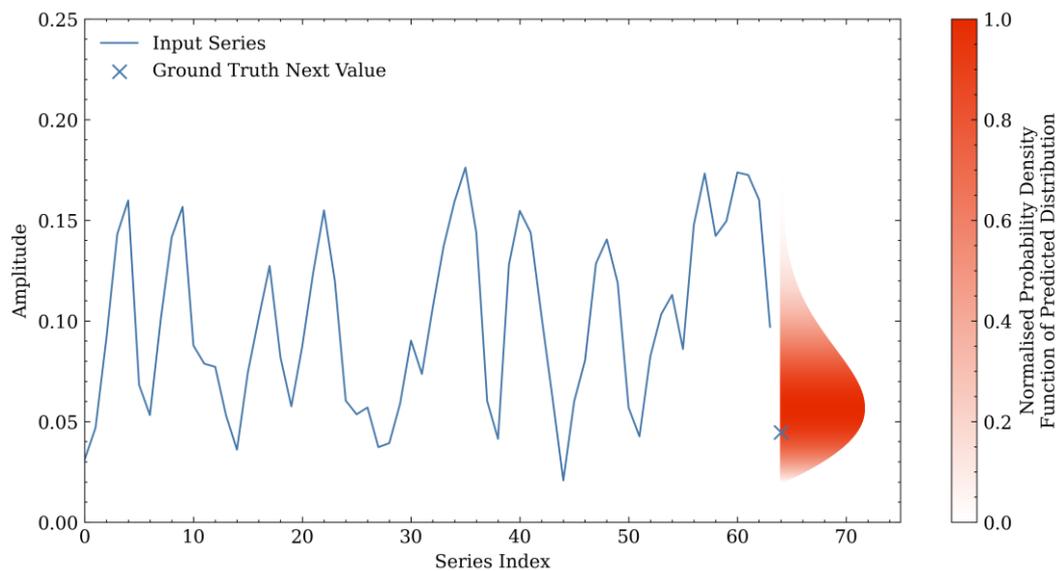


6
7 *Figure 60: Diagram illustrating an example of a through-component series, as examined for the*
8 *pretext learning task.*

9 For the pretext training task in this study, inspiration is drawn from language modelling
10 (where a model tries to predict the next token in a sequence), aiming to predict the next
11 value in the series for a clean sample. While alternative methods such as contrastive
12 learning or generative learning could be utilised, the task of predicting the next value
13 lends itself well to in-process inference and, more broadly, to the downstream
14 inference task of volumetric segmentation.

15 Preliminary trials, focused on single-value prediction, which demonstrated good
16 performance on the pretext task but performed poorly on the inference task due to the
17 lack of accountability for variability within sequences. This made establishing
18 appropriate thresholds challenging and often lead to noisy inference results. For
19 instance, while high-amplitude areas of the scan (such as the front wall) might yield

1 relatively accurate predictions, the significant variability posed an issue. Although
2 variable or percentage-based thresholds provided some benefit, the underlying
3 problem persisted. To enhance this method, a more sophisticated Probabilistic Neural
4 Network (PNN) approach was employed. In this approach, the model attempts to
5 predict the distribution that corresponds to the likelihood of the expected next value in
6 the scan sequence, as depicted in Figure 61. The distribution shows the expected
7 probabilities of different values based on the model's learnt knowledge and the input
8 series, where the mode represents the most likely expected value. This allows the
9 model to account for areas of prior variability or lack of variability by widening or
10 tightening the distribution. Thresholds can then be set as confidence intervals against
11 these distributions, automatically accommodating series variability. We employed a
12 sequence length of 64 values for prediction, aiming to strike a balance between having
13 a sufficiently large receptive field to learn about the distribution and patterns of
14 amplitude responses within the component, while also ensuring that the sequence
15 length is not excessively long, which could limit inference to large components and
16 minimise access to training data.

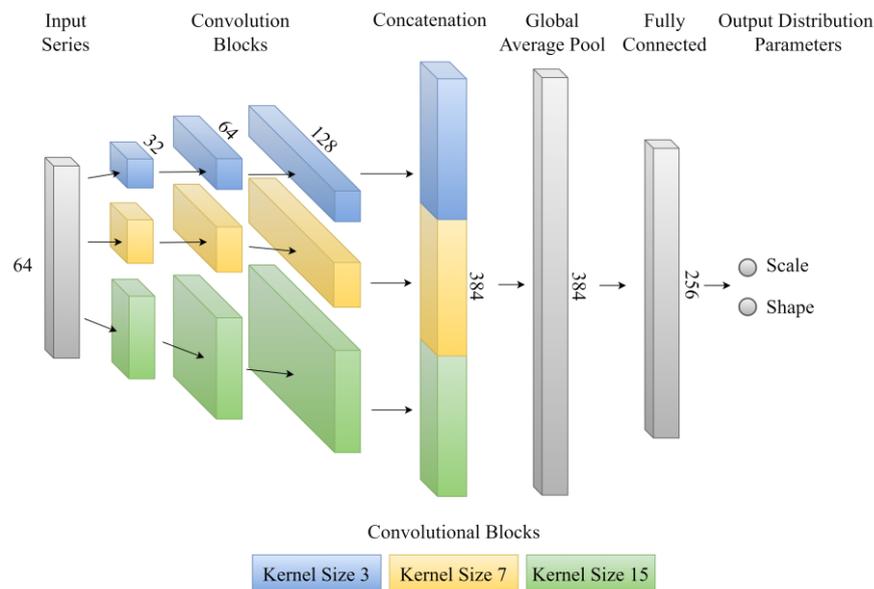


1
2 *Figure 61: Example of the models' predicted distribution for the next value from a given input series*
3 *and comparison to the true value.*

4 Various models for series forecasting exist, ranging from DL methods such as
5 Recurrent Neural Networks (RNNs) and Transformers to probabilistic approaches like
6 Gaussian Processes [194], which can be adapted for probabilistic prediction based on
7 specific requirements. In this study, a 1D multi-head CNN, inspired by InceptionTime
8 [195], was adopted for the architecture, with short series (<128) shallow networks with
9 a lower number of shorter filter lengths being deemed sufficient. This has enabled the
10 1D CNN to be a lightweight model, with 486,242 trainable parameters occupying just
11 1.94 MB of memory.

12 The model probabilistic approach is characterised by predicting the scale and
13 concentration parameters of a two-parameter Weibull distribution as outputs. The
14 choice of a two-parameter Weibull distribution enables the modelling of various
15 distribution shapes for continuous positive values (as consistent with enveloped
16 amplitude values), accommodating non-symmetric distributions through both left and
17 right skewed data. While this model and method proved effective for the downstream

1 segmentation task (please see section 6.4.2), the exploration and evaluation of
 2 alternative probabilistic regression methods and architectures which could lead to the
 3 prediction of better fitting distributions has been left for future work. The model
 4 architecture is depicted in Figure 62, where each convolutional block consists of a
 5 convolutional layer with the specified kernel size and stride of 1, followed by a
 6 convolutional down-sampling layer with a stride and kernel size of 2. A LeakyReLU
 7 activation function is employed between each convolutional and fully connected layer.



8

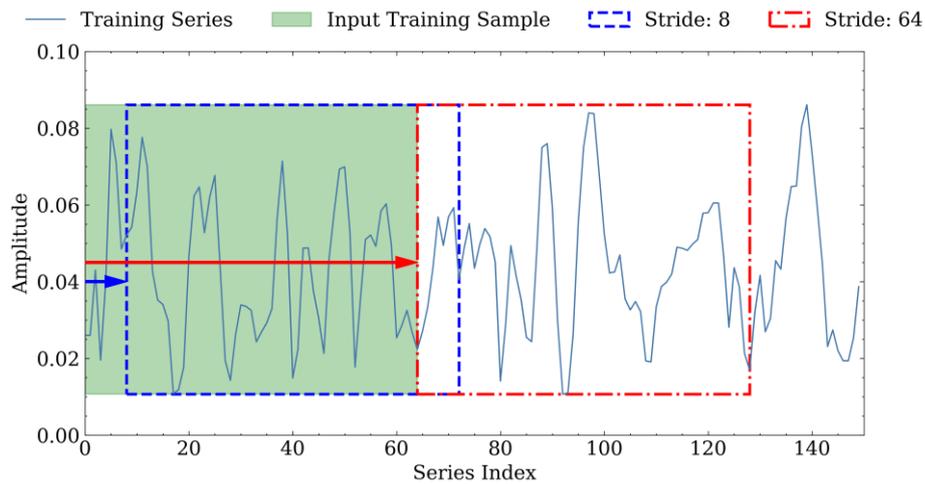
9 *Figure 62: Probabilistic CNN architecture.*

10 During pretext learning, clean samples 1-4 were used for training with a batch size of
 11 65536, whilst clean samples 5 and 6 were kept as holdouts for the validation and test
 12 sets respectively (Table 23). Adam optimiser [70], with a learning rate of 1×10^{-6} is
 13 used to minimise the Negative Log-Likelihood (NLL) loss for the Weibull distribution,
 14 as given by equation (17), where $f(a, b | x_i)$ is the Weibull probability density
 15 function parameterised by Scale (a) and Shape (b). To minimise overfitting, a patience
 16 of 3 epochs was used when evaluating the validation set to determine early stopping.

$$\textit{Weibull NLL Loss} = -\log \prod_{i=1} f(a, b | x_i) = -\sum_{i=0}^n \log f(a, b | x_i) \quad (17)$$

1 During training the data is down sampled in the time domain by every 5 samples
2 (approximately 30 μm in depth within the CFRP samples). This was done as points
3 next to each other in the time domain exhibit minimal variation because of the high
4 ultrasonic sampling rate; offering limited additional information to be learnt and an
5 increased computational cost during training.

6 During training, a hyperparameter which arises for this problem is the stride for
7 sampling data during training. Consider a single full-length scan series: the rate at
8 which this full length is sampled for new training samples is determined by the stride
9 of the window applied, as demonstrated in Figure 63. In DL, it is generally
10 advantageous to maximise the amount of training data available. Using a stride of 1
11 achieves this by providing maximum available training data. However, whilst each
12 sample is different by a single point, there exists significant overlap between
13 neighbouring series, exposing the model to very similar data during training. In such
14 cases, the model may not gain additional information from the additional samples and
15 could become prone to overfitting the training data.

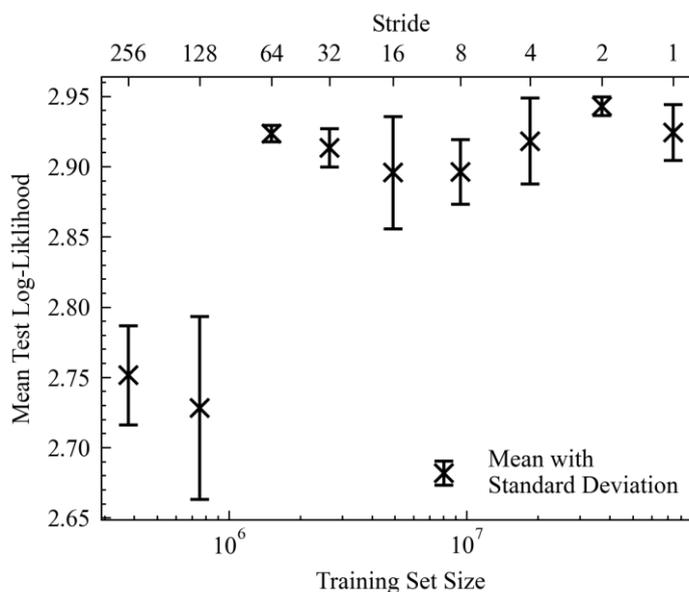


1

2 *Figure 63: Demonstration of the impact of stride (8 and 64) when sampling the training data.*

3 Alternatively, a stride length of 64 can be employed, where each training example
 4 represents a distinct series without information sharing due to overlap with other
 5 training examples. While this reduces overlap, it also significantly reduces the
 6 available training data. Picking an appropriate stride for sampling training series is
 7 therefore a trade-off between training set size and overlap in the training data,
 8 potentially leading to overfitting. To explore this relationship, we tested the following
 9 stride values: 1, 2, 4, 8, 16, 32, 64, 128, 256 which corresponded to dataset sizes of
 10 74.12 M, 37.06 M, 18.53 M, 9.46 M, 4.92 M, 2.65 M, 1.51 M, 0.76 M, 0.38 M
 11 respectively. It is worth highlighting that an advantage of using series in this way is
 12 the production of very large training datasets compared to typical image-based ML
 13 research in NDE. For this test, the stride of the validation set was fixed at 64, and the
 14 testing stride was set to 1 to maximise the test set size (which also matches the
 15 requirement during inference). The model was trained three times for each stride. For
 16 each training iteration, the mean log-likelihood was recorded across the test set, and
 17 the mean and standard deviation of the results from the three training runs are reported
 18 in Figure 64. A higher log-likelihood between the predicted distribution of the

1 expected value and the ground-truth measured value indicates that the model's
2 prediction is better aligned with the measured value. This suggests that the model is
3 better at capturing the underlying patterns or relationships in the data, leading to more
4 accurate and precise predictions.



5

6 *Figure 64: Test set mean Log-Likelihood for varying sampling strides.*

7 As depicted Figure 64, a stepped improvement is observed when the training datasets
8 exceed 1 million samples (stride < 128). However, beyond this threshold, the model
9 appears to exhibit diminishing returns from additional training data. Similar trends to
10 this have been documented in the existing literature [196]. Perhaps surprisingly, there
11 is not a consistent increase in performance for a reduction in stride. Whilst this could
12 indicate model saturation, there could also be a detailed relationship between training
13 data and overfitting at play. Notably, strides of 64 and 2 exhibited lower deviations in
14 results, coupled with good performance, suggesting the emergence of stable solutions.
15 It's important to note that different pretext learning models and training datasets may
16 exhibit varying relationships with stride length. Nonetheless, this analysis underscores

1 the significance of tuning this parameter and provides insight into the required training
2 dataset sizes. In future work, it is hoped to explore this further and assess its impact on
3 inference performance.

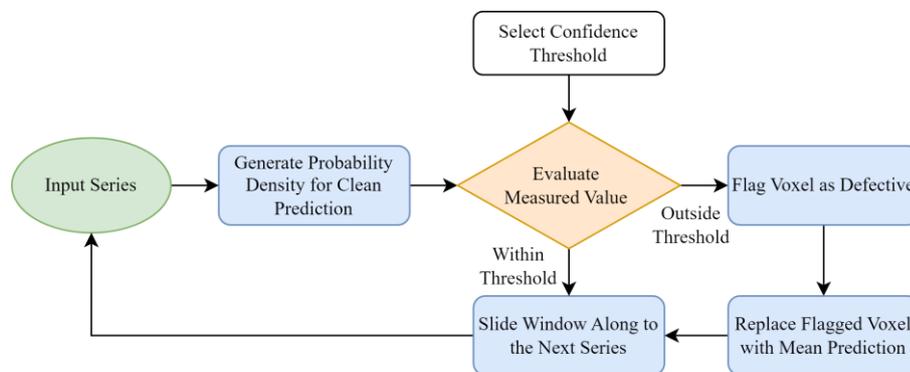
4 6.4 Inference

5 6.4.1 Methodology

6 By training on defect free data, the pretext learning phase yields a model which
7 predicts the expected distributions of the next defect free value for a given time index
8 and reception element. By leveraging this pre-learnt information, the 1D model is
9 utilised by applying it to each preceding 64 series of volumetric data to forecast the
10 distribution for each point in the subsequent B-scan. This enables anticipation of the
11 distribution of data points across the entire next B-scan under the assumption of a
12 defect-free component.

13 During inference each volumetric frame (B-scan) is processed in a sequential manner,
14 comparing the predicted clean B-scan to the measured experimental result at every
15 point in time for each receptive element or sub-aperture of elements. Since the model
16 is trained only on clean data to predict the next clean value in the series, if there is
17 agreement between the values in the predicted frame and the measured experimental
18 frame they are marked as defect free. The inference window is then moved to the next
19 frame in the data and the process is continued. This acts to update the predictive model
20 with the most current information and keeps good alignment with variations seen in
21 specific samples/scans.

1 However, if the next B-scan contains defective voxels this will result in an increased
 2 amplitude response around the defect which is outside the distribution of clean
 3 sequences seen in the prior learning. As a result, when comparing the predicted (clean)
 4 frame to the measured (defective) frame there will be poor agreement locally around
 5 the defect response. These voxels can therefore be marked as defective within the
 6 volume; locally segmenting the defect. For the series used for the subsequent frame,
 7 prior predictions not marked as defective are treated as normal with the series window
 8 sliding along from the experimental series. However, for defective voxels, the mean is
 9 taken from the probabilistic output as an a priori estimate of a clean response and is
 10 used when evaluating future sequences. Doing this, ensures that the model is always
 11 predicting expected distributions of clean responses based on prior clean sequences or
 12 the closest estimate to clean values, and errors as a result of defective responses do not
 13 impact future predictions, as consistent with the pretext learning stage. This process is
 14 outlined in Figure 65.



15
 16 *Figure 65: Flowchart of sequence prediction for inference.*

17 Given the probabilistic output, a threshold must be set to evaluate the model's clean
 18 prediction against the measured value to determine if a voxel is defective. This
 19 threshold can be set as a confidence threshold on the predicted clean output

1 expectation, making it adaptive to variations in series amplitude responses. There are
2 different ways to determine an appropriate confidence threshold such as experimental
3 calibration. An alternative is to base it off an allowable false-call rate. Since this
4 detection analysis is for the entire volume and due to the resolution difference in the
5 spatial and time domain, the number of voxels is far larger (700 times for test samples
6 1 and 2) than for image level analysis. For the same expected number of absolute false
7 calls for image analysis the false call rate for the volume would therefore need to be
8 far lower (due to the increased number of voxels). To account for this a confidence
9 threshold is chosen based on a much lower allowable false-call rate than would be
10 expected for image analysis. It is important to note that this false-call rate is done on a
11 per voxel basis and not per defect basis. In this study results are presented for false-
12 call rates ranging from 1% to 0.00001%. Even at the aggressive lower false-call rate
13 of 0.00001, sample 2 which has 5.6 M voxels would still be expected to produce
14 approximately 6 false-calls. Whilst for these samples this is a low number, for larger
15 parts this would scale cubically for any increase in the number of B-scans as a result
16 of a longer scan.

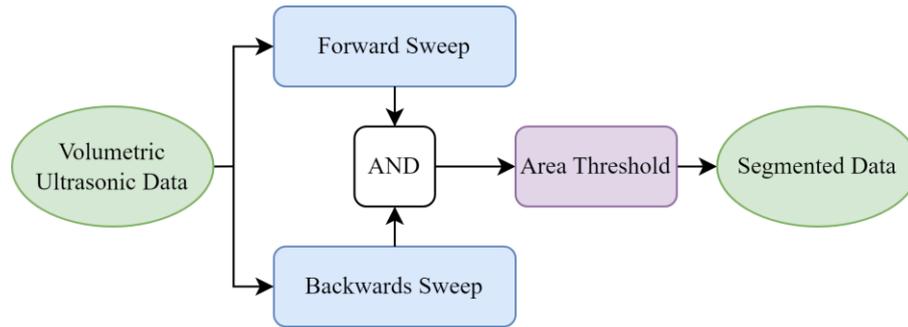
17 For volumetric inspection, the inference process can be completed just following the
18 forward scan (as most applicable to in-process inspection), herein denoted as the
19 forward sweep, or post scan the analysis can be completed in reverse – simulating the
20 scan from the other direction, herein denoted as the backwards sweep. Taking the
21 logical AND of both passes as the final segmented volume acts to increase confidence
22 in defective predictions by having them detected twice. This helps to further remove
23 false positive indications and cleans up the segmented response as seen in Figure 67.

1 To further mitigate the impact of false-calls, a minimum area threshold is applied to
2 detect only defects above a critical size and limit low-area voxel clusters resulting from
3 the probabilistic nature of predictions. For CFRP applications, defects generally align
4 parallel to the ply orientation. It is this area which is typically used to assess critical
5 defect size. Area thresholding is calculated using equation (18), where *Area Opening*
6 refers to the *skimage.morphology.area_opening* [197] which, removes all connected
7 components smaller than the *Filter* for the volume. The connectivity determines what
8 neighbours are considered connected components, for this application all 8
9 neighbouring voxels are considered for a 2D slice through the volume. For other
10 applications and materials, where the defects are not primarily in-plane, 3D connected
11 components may be more appropriate.

$$\textit{Thresholded Volume} = \textit{Area Opening}(\textit{Volume}[:, \textit{depth}, :], \textit{Filter}, \textit{Connectivity}) \quad (18)$$

12 The minimum defect size threshold can be tailored based on the specific application
13 requirements. For testing purposes, where the minimum defect size was 3.0 mm in
14 diameter, the area threshold was adjusted to exclude any indications smaller than this.
15 For applications where thin crack-type defects are typical, such as in metals, it may be
16 beneficial to consider a more complex shape-based thresholding method, such as
17 evaluating the defect aspect ratio to account for extreme geometries. However, for
18 composites, which generally exhibit in-plane defects, area thresholding is deemed
19 appropriate and is used by Spirit AeroSystems to determine defect criticality. Area
20 segmentation is completed as the final processing stage. The complete workflow for
21 volumetric instance segmentation is illustrated in Figure 66. Details on the impact of
22 each processing step for different thresholds can be found the in the following Results
23 section (Figure 69) and Table 28 of the Appendix.

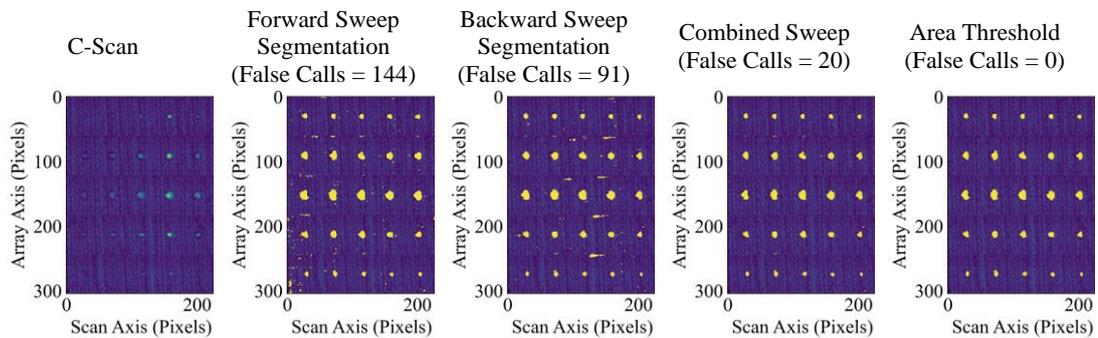
1



2

3 *Figure 66: Flowchart of the methodology overview for complete volumetric segmentation.*

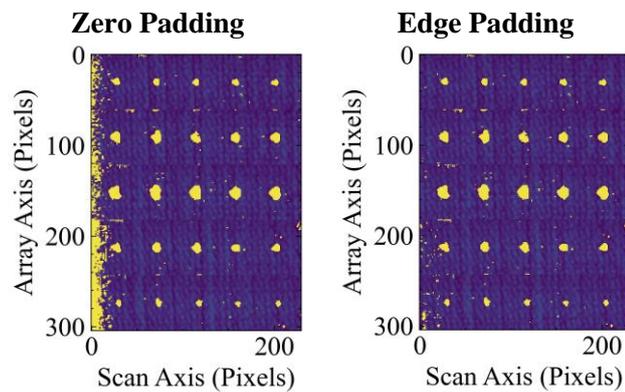
4



5 *Figure 67: Demonstration of the impacts of different post-processing steps for defect sample 2 and*
6 *threshold of 0.9999999.*

7 To process the first frames where there are not 64 prior frames for prediction, padding
 8 is used to populate the remaining missing information for each sequence. Whilst zero
 9 padding, constant edge padding or reflect padding can all work as the combined
 10 sweeps clear up any errors as a result of not having enough prior information for
 11 predictions. They can produce artefacts for a single sweep at the edge of the scan due
 12 to a lack of prior information for the model to give an accurate prediction. Reflect and
 13 edge padding give the best results, reducing the reliance on sweeps to clean up edge
 14 artefacts. Figure 68 shows an example of edge artefacts from using zero padding
 15 compared to edge padding.

16



1 *Figure 68: Example of edge artefacts for a single sweep when using zero padding compared to edge*
 2 *padding (sample: defect 2, threshold: 0.9999999).*

3 Different defects exhibit varying characteristics. This method reframes the problem of
 4 positive defect identification as one of anomaly prediction, which is beneficial for
 5 generalisability to a wide range of defects. However, as the model has not learnt
 6 specific positive features of defect responses, any change in amplitude response,
 7 whether from a defect or geometric feature, can lead to an anomalous prediction. Real
 8 components often feature geometric elements such as step changes in thickness. While
 9 inference can still be conducted in such cases, analysis must be conducted in parallel
 10 with any geometric alterations so that they are present in the series used for prediction.
 11 Therefore, if there is a part with very complex geometry it may have to be sectioned
 12 for inference. To evaluate the method's performance under these conditions and to
 13 assess its effectiveness with different defect types, we applied the approach to a
 14 stepped sample with PTFE inserts (defective sample 3). These defects are inserted pre-
 15 cure and are more representative of naturally occurring defects compared to FBHs.
 16 For inference the model used was the best performing during pretext learning.
 17 Inference is a sequential process, but the computational cost can be easily minimised
 18 by batching the predictions for a complete frame; due to the lightweight 1D CNN. To
 19 reduce the computational cost during inference, the time domain was down sampled

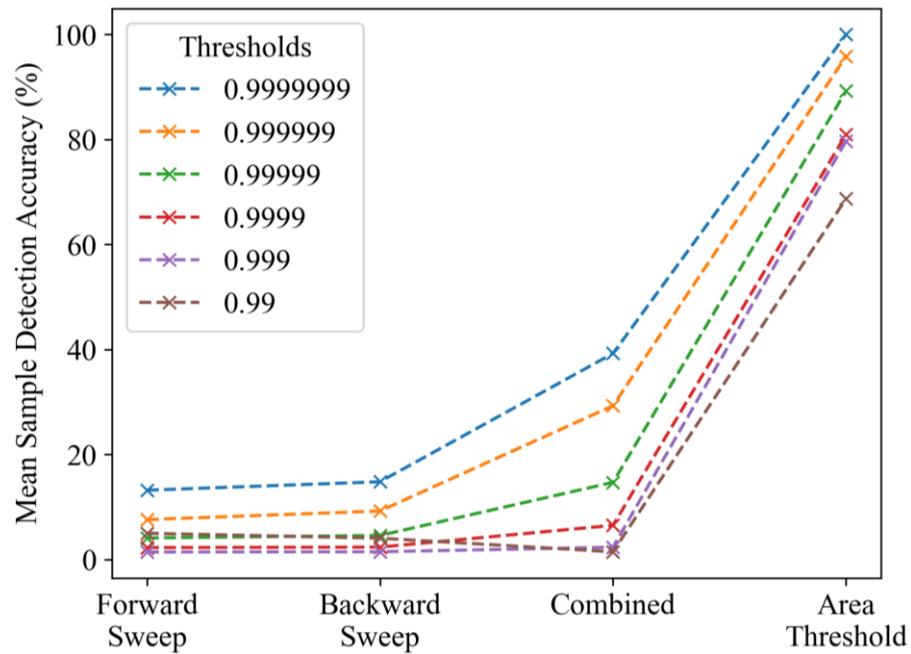
1 by a factor of 10 for testing. Previous work demonstrated that it is still possible to get
2 very high levels of depth-wise resolution with this sampling [193]. Inference for a
3 single frame took approximately 0.05s when batched. To complete the full processing
4 pipeline of forwards and backwards sweeps and area thresholding took less than 35s
5 for each sample. Testing was conducted on a workstation equipped with a NVIDIA
6 GeForce RTX 3090 running the Pytorch [146] framework. The following section
7 presents results for defect detection, sizing, localisation, and visualisations for
8 qualitative examples of volumetric segmentation.

9 6.4.2 Results

10 6.4.2.1 Detection

11 Detection performance of the model was evaluated for thresholds ranging from 0.99
12 to 0.9999999, with results reported for each processing stage. Evaluated was
13 conducted on amplitude C-scans of the segmented volumes to avoid the impact of
14 repeat echoes stemming from defect indications. The ground truth defect mask was
15 generated using manual identification of defects, with the 6 dB drop used for locating
16 the centroid of each defect. The true defect sizes were then used around the centroids
17 to mark out the defect areas. A defect is considered detected if the predicted mask
18 overlaps with the ground truth mask. If there is no overlap the prediction is considered
19 a false-positive. Across all thresholds and processing stages, all defects were
20 successfully identified, with no missed detections; however, a notable decline in
21 accuracy is observed as a result false positive indication. Whilst the absolute number
22 of false positives can be large, they are often very small in size (as demonstrated in
23 Figure 67) and can therefore be effectively removed with area thresholding, which can

1 be derived from critical defect size (as dictated by industrial testing standards). The
 2 results for detection accuracy are presented in Figure 69 (with a full breakdown
 3 available in Table 28 in the appendix). It is feasible to completely minimise false
 4 positives, thus achieving a detection accuracy rate of 100%, using the processing steps
 5 suggested in Figure 66 and a threshold of 0.9999999.



6
 7 *Figure 69: Defect detection accuracy for each threshold and processing step.*

8 6.4.2.2 Defect Sizing

9 The impact of defect sizing has been assessed across the same range thresholds as used
 10 for detection, with the results summarised in Table 23. While the deviation of absolute
 11 error remains relatively stable, there exists a strong negative correlation between
 12 increasing the confidence threshold and a reduction in MAE for defect sizing. Notably,
 13 at a threshold of 0.9999999, the MAE of 1.41 mm, aligns closely with the findings of
 14 a previous study which reported errors in sizing inaccuracy for the 6dB drop method
 15 of 1.35 mm [193]. Moreover, this performance surpasses that of the fully supervised

1 3D U-Net method previously presented, prior to adjustments for differences between
2 synthetic and experimental domains.

3 The method has been shown to only oversize defects. This stems from the underlying
4 mechanism of detecting anomalous voxels, which does not always correspond
5 perfectly in a one-to-one manner to actual defect size; due to impacts of ultrasonic
6 imaging such as beam spread etc. The impact of this is visually demonstrated through
7 an example B-scan in Figure 70. Consistency in oversizing defects yields two main
8 advantages:

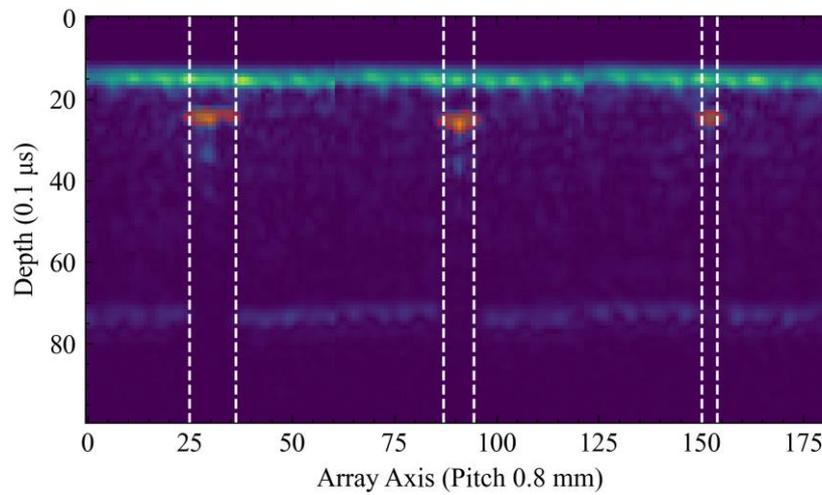
9 **Safety Considerations:** Oversizing defects, while potentially sacrificing some
10 precision, inherently reduces the risk of overlooking critical flaws. In safety-critical
11 industries such as aerospace or structural engineering, the consequences of under
12 sizing defects can be severe, leading to structural failures or operational hazards. By
13 consistently erring on the side of caution and oversizing defects, the method provides
14 a safety buffer, ensuring that potential weaknesses are identified and addressed before
15 they escalate into safety incidents.

16 **Calibration Potential:** The method's consistent error pattern facilitates calibration
17 procedures, akin to those commonly employed in traditional amplitude drop sizing
18 methods. Calibration involves establishing a correlation between measured responses
19 and actual defect sizes, thereby refining the method's accuracy and reliability. The
20 method's predictable oversizing behaviour lends itself well to calibration processes,
21 allowing for adjustments that compensate for systematic errors. For instance, for a
22 threshold of 0.9999999, when utilising defect sample 1 as a calibration reference, the
23 MAE for Defect Samples 2 and 3 decreases to 0.58 mm, a 57% reduction in MAE.

1 *Table 24: MAE for different thresholds.*

Sample	Defect Width (mm)	MAE for given threshold					
		0.99	0.999	0.9999	0.99999	0.999999	0.9999999
Defective 1	9	2.75	2.25	1.41	1.27	1.13	1.02
	6	3.39	2.55	2.05	1.67	1.44	1.27
	3	3.78	2.91	2.60	2.21	1.83	1.75
Defective 2	9	3.14	2.38	1.86	1.49	1.20	0.93
	7	3.26	2.61	2.19	1.83	1.50	1.23
	6	3.18	2.62	2.10	1.74	1.37	1.23
	4	3.67	2.97	2.45	2.17	1.93	1.79
	3	3.93	3.39	2.94	2.58	2.32	2.06
Defective 3	6	3.06	2.32	1.73	1.46	1.29	1.21
Mean		3.37	2.71	2.20	1.87	1.59	1.41
Standard Deviation		0.66	0.61	0.69	0.67	0.66	0.68

2



3

4 *Figure 70: Example B-scan across multiple raster passes showing the voxels highlighted as defective*
 5 *and the corresponding true defect size in white. This is shown for defect sample 1 with a threshold of*
 6 *0.9999999.*

7 *6.4.2.3 Localisation*

8 For linear scanning of composite materials, achieving accurate volumetric localisation
 9 necessitates an understanding of both in-plane and depth-wise localisations. For depth-
 10 wise localisation, the evaluation focuses on FBHs, while excluding inserts. Inserts,
 11 positioned between layers pre-cure, pose a challenge in obtaining true ground truth
 12 depth measurements post-curing, unlike FBHs which can be accurately measured post

1 cure to obtain a ground truth depth measurement. Through-thickness measurements
2 are derived from the mean thickness through the centroid of the segmented defect.

3 In contrast, acquiring a ground truth measurement for in-plane localisation presents its
4 own challenges. Cumulative positional errors inherent in the experimental setup render
5 the attainment of accurate ground truth measurements (<1.0 mm) infeasible.
6 Consequently, this research adopts the 6 dB drop criterion as a reference standard for
7 validating the agreement between segmented in-plane localisation and the 6 dB drop
8 criterion. By doing this the in-plane localisation can be compared against an industry
9 standard method as a validation benchmark. To assess in-plane localisation, the
10 Euclidean distance between the centroid of the 2D projected segmentation mask and
11 the 6 dB drop masks is computed.

12 The integration of both through-thickness and in-plane localisation provides a
13 comprehensive understanding of the method's ability to locate defects within the
14 volume. Table 25 presents the mean and standard deviation of results. The MAE for
15 depth localisation surpassed that of in-plane localisation, with comparable standard
16 deviations for both. Although this method showed improvement for in-plane
17 localisation compared to previous fully supervised volumetric segmentation methods,
18 it fell short of achieving the high accuracy levels seen in through-thickness depth
19 localisation. This discrepancy might be partly attributed to the additional pre-
20 processing steps used for the fully supervised method, such as peak alignment.
21 Enhancing through-thickness accuracy could potentially be achieved by increasing the
22 temporal sampling rate during inference, although for most applications, the current
23 level of accuracy is likely sufficient. Both metrics demonstrate considerable

1 localisation performance, with MAEs well below 0.5 mm. This level of precision is
 2 likely more than sufficient for typical industrial rework scenarios, as it aligns with the
 3 accuracy required for precise tooling operations. Additionally, the in-plane MAE is far
 4 below the element pitch, which is the limiting factor for in-plane imaging resolution.
 5 All in all, the model performs well in volumetric defect localisation.

6 *Table 25: Localisation results.*

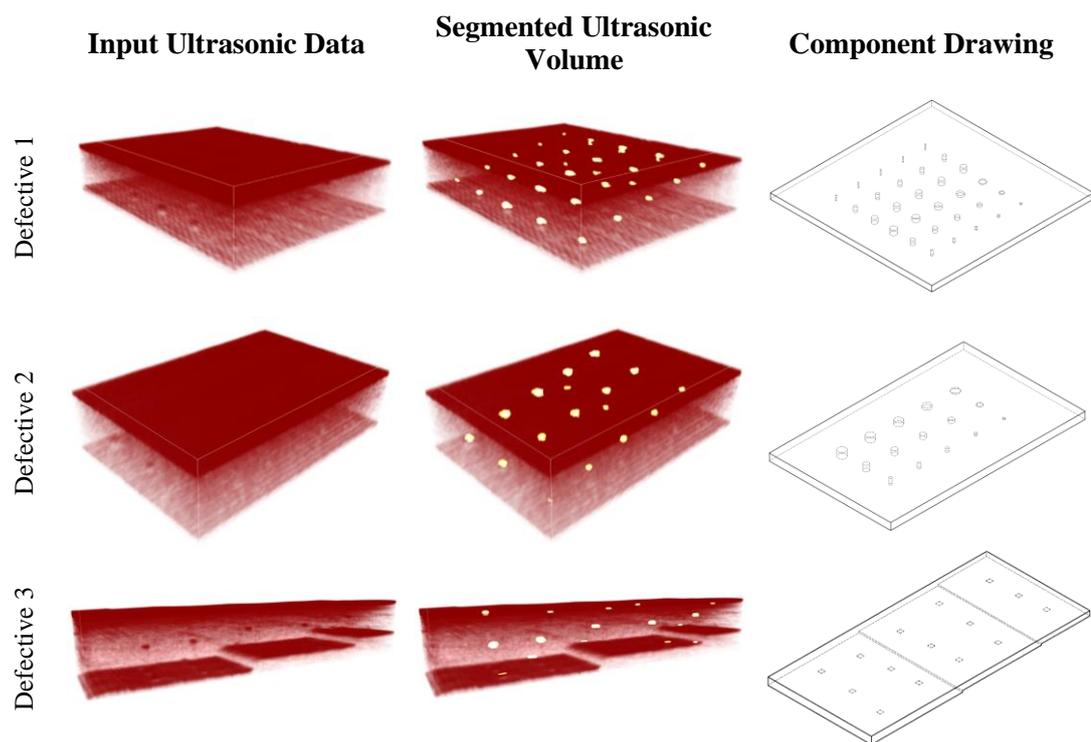
Sample	Defect Width (mm)	Depth		In-Plane Distance	
		MAE	Standard Deviation	MAE	Standard Deviation
Defective 1	9	0.34	0.29	0.41	0.11
	6	0.10	0.05	0.50	0.28
	3	0.16	0.05	0.26	0.16
Defective 2	9	0.33	0.15	0.28	0.15
	7	0.30	0.08	0.32	0.10
	6	0.23	0.17	0.44	0.14
	4	0.23	0.15	0.31	0.08
	3	0.36	0.06	0.40	0.16
Defective 3	6	-	-	0.27	0.18
Total		0.26	0.17	0.37	0.18

7

8 *6.4.2.4 Visualisations*

9 While the previous section quantitatively demonstrated the method's ability to localise
 10 defects in 3D, this capability is best appreciated through visualisations. Volumetric
 11 segmentation offers a significant advantage over image-based segmentation by
 12 providing comprehensive localisation information about defects. This enhanced
 13 localisation allows for more thorough reconstructions, which are particularly
 14 beneficial for constructing digital twins of components for testing or reporting
 15 purposes. To illustrate this benefit, Figure 71 presents 3D visual examples of the
 16 volumetric segmentation results, showcasing the model's performance across various
 17 defects and stepped samples. Despite the absence of TCG and significant variations in
 18 defect response levels at different thicknesses due to attenuation, the method manages

1 to deliver relatively consistent segmentation masks at different depths for defects of
2 the same size by considering local temporal samples during inference. This resilience
3 to variations in thickness and attenuation levels enhances the reliability and
4 applicability of the segmentation method across a range of scenarios. After visualising
5 the 3D segmented volumes, the opportunity of constructing digital twins becomes
6 evident and highlights a clear benefit of the method.



7 *Figure 71: Visualisations of volumetric ultrasonic responses, their corresponding overlaid*
8 *segmentations, and component drawing.*

9 6.5 Industrial Demonstration

10 Given that this PhD research was sponsored by Spirit AeroSystems, it was important
11 to investigate the applicability and efficacy of the proposed method when applied to
12 industrial data and components. This sub-section gives an initial look into how the SSL
13 method can be applied to data acquired from an industrial component and a different
14 industrial acquisition setup.

1 An industrial sample was supplied by Spirit AeroSystems to test the performance of
 2 the SSL volumetric segmentation method. The sample is a stiffened panel,
 3 manufactured to BAPS 260-007, from the Learjet 85 with stringers which give rise to
 4 variable geometry. There are a range of PTFE insert which simulate
 5 inclusion/delamination defects. A 240 mm by 490 mm section of the component was
 6 scanned which contained 12 defects at a range of locations. A breakdown of the
 7 positions of the defects are given in Table 26.

8 *Table 26: Summary of defects scanned in the Learjet 85 component.*

Defect Position	Dimensions (mm) [Count]	Total
Skin (Near Tool/Top Surface)	20×10 [2], 10×10 [2], 5×5 [2]	3
Stringer (Near Tool/Top Surface)	20×10 [1], 10×10 [1], 5×5 [1]	3
Stringer (Near Bag/Bottom Surface)	20×10 [2], 10×10 [2], 5×5 [2]	6

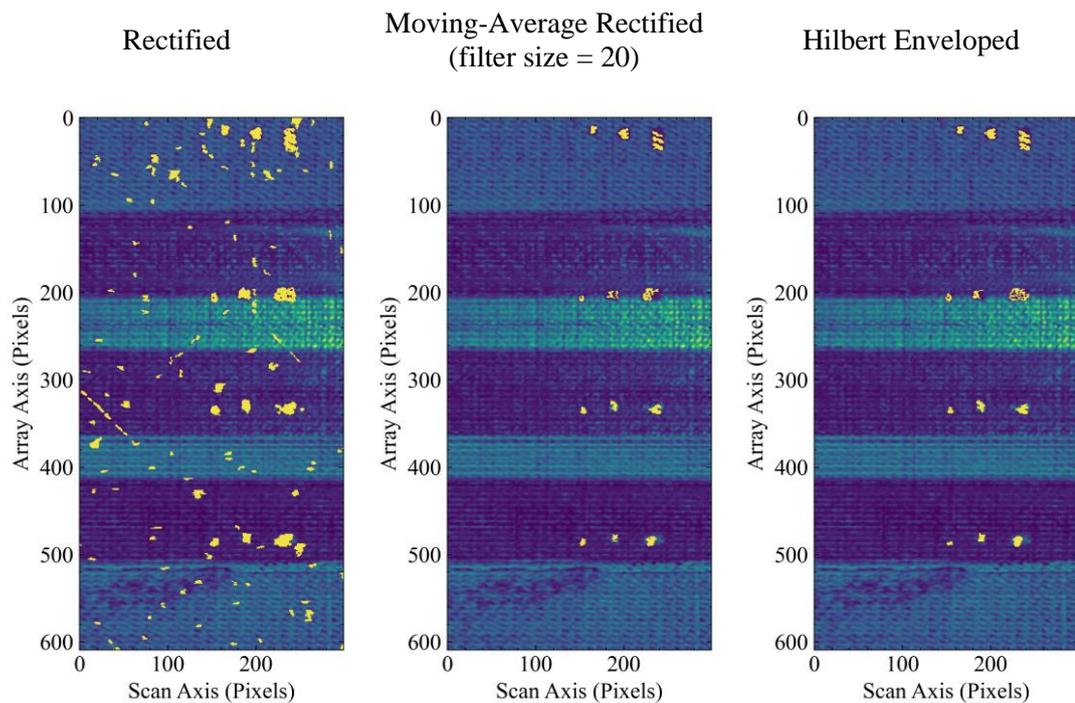
9 Typical industrial data interpretation times for the Learjet sample can total 140
 10 minutes; 80 minutes if the component is defect free, with a potential additional 60
 11 minutes if the component contains defects. The component was scanned by Spirits’
 12 Belfast site using their Tecnomat phased array ultrasonic acquisition setup. A C-scan
 13 produced from the Tecnomat setup of the analysed section is given in Figure 73.

14 Data acquired at Spirits’ Belfast site is usually analysed and stored post rectification,
 15 where the absolute value of the full waveform data is used to removed negative values.
 16 This rectification process compresses the data and removes phase information, similar
 17 to taking the envelope of the signal. However, unlike the enveloped signal the rectified
 18 signal gives an unsmooth response with multiple local maxima and minima for a given
 19 instantaneous response.

1 The SSL segmentation method is designed to be applied to smooth enveloped data and
2 can highlight anomalous responses for given temporal signals. However, the method
3 fails when applied to rectified data (as demonstrated in Figure 72). The presence of
4 multiple maxima and minima for a given response makes it impractical to compare
5 indexes over time, as minor temporal shifts can lead to significant changes in
6 amplitude within a response. To address this issue the use of a moving average filter
7 was explored. This can be applied to smooth the rectified signal, minimising the impact
8 of multiple minima within a response. This approach has limitations such as requiring
9 manually selecting the filter size to balance adequate smoothing and preserving the
10 signal's response. However, it does provide a potential solution for the analysis of
11 industrial data where the full waveform or enveloped signal is not available.

12 The results of applying the SSL method to the Learjet sample, based on data collected
13 using the setup presented in Chapter 2, are presented herein (Figure 72). The
14 evaluation focuses on the effects of different data processing techniques to ensure
15 alignment with industrial data collection, storage, and interpretation. Specifically, the
16 analysis includes the use of rectified data, moving-average rectified data, and Hilbert
17 enveloped data. For this study, a minimum defect size filter was set to remove any
18 defects below 3 mm in width, as consistent with the previous tests and a median
19 threshold of 0.99999 was selected.

20

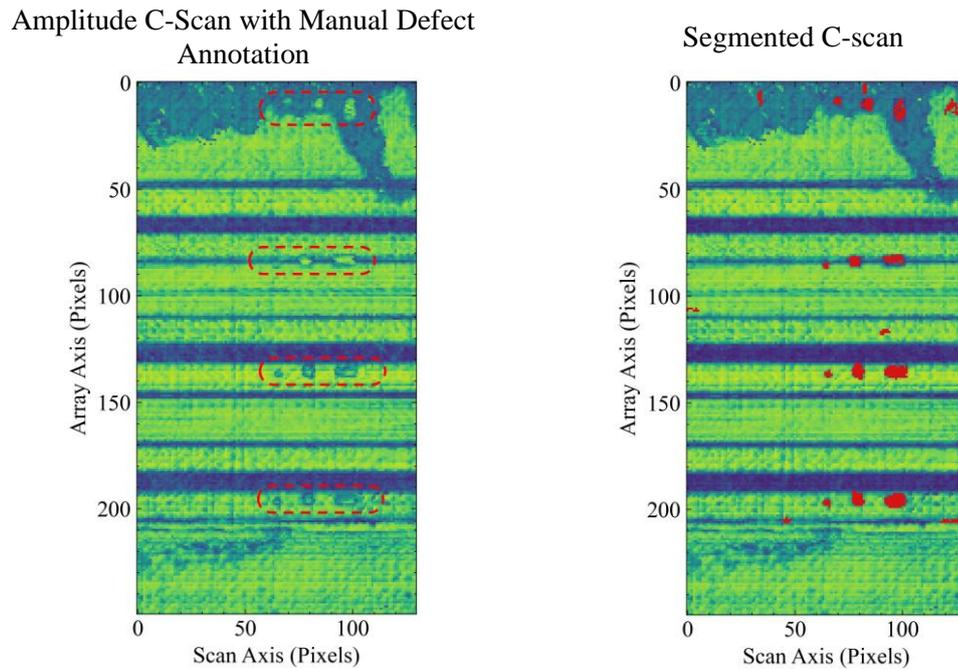


1 *Figure 72: Impact on data processing for SSL method on Learjet data using the SEARCH lab*
 2 *acquisition system.*

3 As qualitatively demonstrated though Figure 72, applying the SSL method directly to
 4 rectified data produces a significant number of false positives as expected. However,
 5 the moving average technique serves as an effective solution, enabling the method to
 6 be applied to industrial data where full waveform or enveloped data is not available.

7 To demonstrate transferability of the method to the acquisition system used by Spirit
 8 at their Belfast site, the method is applied to the moving average rectified data taken
 9 using the Tecnatom setup. The segmented C-scan for this is shown in Figure 73. Whilst
 10 there are additional false positives compared to the data acquired using the roller probe
 11 (on the Strathclyde SEARCH acquisition set up, Figure 72), each defect is clearly
 12 segmented. Despite the SSL model not having seen data from this acquisition setup,
 13 the results show promise for the method to be applied to both industrial parts and
 14 industrial acquisition systems. The results also demonstrate that in instances where the

1 Hilbert enveloped data is unavailable, applying a moving average to the rectified data
 2 can act as a practical workaround for the application of this methodology. It would be
 3 beneficial to conduct a wider study in future to assess the impact of different thresholds
 4 and moving average filter sizes. However, these results give promise to the application
 5 of this method to both industrial parts and acquisition setups.



6 *Figure 73: Amplitude C-scan of the Learjet component section acquired using the Tecnatom setup.*
 7 *Both manual defect annotation and SSL defect segmentations are shown.*

8 *Table 27: Comparison of correct defect detections and false indications for the different pre-*
 9 *processing and acquisition systems.*

Acquisition System	Pre-processing	Correct Defect Detections	False Indications
SEARCH Lab Roller Probe	Rectified	12	101
	Moving Average	12	0
	Enveloped	12	0
Spirit Tecnatom	Moving Average	12	7

10

1 Table 27 highlights that although defects were consistently detected, different data pre-
2 processing techniques significantly impacted the level of false indications. It also
3 demonstrates that an appropriate moving average filter can match the performance of
4 the enveloped data in this case. The method has shown promising transferability to
5 Sprints' Tecnatom acquisition system, with the moving average filter clearly detecting
6 all defects. Despite seven false positive indications, which could be due to a domain
7 shift or due to the moving average filter not being specifically tuned for this acquisition
8 system. This underscores a significant disadvantage of using this method with an
9 additional pre-processing step: the introduction of a tuneable parameter due to the lack
10 of full-waveform data. While a longer moving average filter kernel might help reduce
11 false positives, it would also diminish defect response amplitudes (and potentially lead
12 to missed defects). Therefore, it is preferable to apply the method to the enveloped
13 data from the raw-waveform data as originally designed. In the future, acquiring full-
14 waveform data using the Tecnatom acquisition system would be beneficial to
15 determine if the additional false positives could be eliminated by using the enveloped
16 data. It would also be beneficial to study the impact of changes in moving average
17 filter size for when enveloped data is unavailable.

18 6.6 Conclusion

19 This chapter introduces a new approach for ultrasonic volumetric defect segmentation
20 using SSL to address the need for defective training data. The method has been
21 demonstrated for CFRP composites over different samples, defects, and geometries.

22 Volumetric segmentation provides information on defect sizing and localisation and
23 allows for advanced visualisations which can facilitate the creation of digital twins for

1 reporting or testing. One of the biggest challenges when applying DL methods to NDE
2 is the requirement for training data. This is compounded for volumetric segmentation
3 tasks as volumetric training data is far less available and segmentation models typically
4 require highly intensive labelling at a per voxel level.

5 This method employs a 1D CNN to learn the expected distributions for clean data
6 along a given time index for a receptive element. This learned information is then used
7 to generate a volumetric defect segmentation mask. Several advantages stem from this
8 approach. Firstly, it circumvents the need for expensive and challenging-to-obtain
9 labelled training data. By simplifying the pretext task to a 1D series, the accessibility
10 of training data significantly increases; with millions of training samples derived from
11 just four samples. Since this training data is defect-free, the limitation of obtaining
12 large quantities of defective components no longer applies. Moreover, the size of the
13 model used was only 1.94 MB. This means that although batching during inference is
14 preferred for time-saving purposes, it is possible to reduce batch sizes or eliminate
15 them altogether to allow the model to run on machines with more limited hardware.

16 Additionally, by framing the problem as one of anomaly prediction rather than positive
17 defect identification, the method becomes more generalisable and robust to various
18 types of defects, as long as they exhibit an amplitude response that deviates from the
19 expected. This was demonstrated through testing on more challenging to detect PTFE
20 inserts which were more accurately sized than the FBH. Furthermore, by treating the
21 task as one of probabilistic sequence prediction, the segmentation approach becomes
22 more interpretable compared to using a 3D end-to-end segmentation model. This aids

1 in demystifying some of the "black box" nature of DL predictions, which is a
2 significant challenge in the NDE industry.

3 All defects were reliably detected across thresholds and processing steps. False
4 positive indications were successfully eliminated for a 100% detection accuracy with
5 the complete processing pipeline and a threshold of 0.9999999. The results showcased
6 good in-plane and through-thickness localisation, exhibiting improvements over the
7 previously presented fully supervised model [193] for in-plane localisation, with
8 MAEs of 0.37 and 0.26 mm, respectively. Although through-thickness localisation
9 performed slightly worse than previously demonstrated 3D U-Net model, it remains
10 suitable for most applications, and it still falls below the in-plane localisation error.
11 MAEs for defect sizing were also presented, revealing a negative correlation between
12 sizing error and threshold. For a threshold of 0.9999999, the MAE aligned with the 6
13 dB drop at 1.41 mm. Due to the method's nature of detecting anomalous voxels in the
14 ultrasonic domain, which may not always correspond directly to sizes in the physical
15 domain, achieving accurate sizing remains challenging. However, the consistent over-
16 sizing of defects enabled consistent calibration, which lead to a 57% reduction in sizing
17 error.

18 The industrial case study demonstrated how the methodology could be successfully
19 applied to an industrial part. It also highlighted a key challenge on how industrial data
20 storage and representation can impact the applicability of the technique. However, the
21 use of a moving average filter proved as a suitable solution to scenarios where the full
22 waveform data may not be available to envelope using the Hilbert transform. Whilst
23 further investigation would be beneficial regarding optimal moving average filter

1 lengths, along with testing on other components, the findings proved promising for the
2 transferability of the approach to industrial components and data collection methods.

3 However, there are limitations to this method. Firstly, the method requires one
4 consistent geometric axis to be applied to the data. While we demonstrated its
5 robustness on stepped samples, more complex shapes may prove challenging, and the
6 method would likely require modifications to accommodate such complexities.

7 Secondly, most DL models shift the computational load of understanding a task to
8 training, which can be done offline, prior to inference. This typically results in rapid
9 inference results – a key benefit over other interpretation methods. However, the
10 sequential nature of inference for this method means that there is an increased
11 computational cost and time during interpretation. Although the full interpretation time
12 was minimised to a maximum of 35 seconds using series batching for each B-scan, for
13 less powerful machines, this may not always be feasible. Additionally, for larger parts,
14 the inference time will increase as a result of longer scans. Despite this inference is
15 still likely to be far quicker than human inspection in real-world applications.

16 In future work it would be interesting to evaluate this method against other, (perhaps
17 simpler) methods of anomaly detection. It would also be beneficial to explore different
18 models for the pre-text learning task and conduct a detailed analysis of how well the
19 model can predict the distribution for the next consecutive value. Whilst the PNN gave
20 results that performed well for segmentation, series forecasting is a significant, active
21 area of research, and it is likely that a better approach exists to predict more accurate
22 distributions of expected data points. A robust exploration of different approaches and
23 how their accuracy impacts segmentation performance would be beneficial. It would

- 1 also be worthwhile exploring if this approach can be applied to different materials and
- 2 scanning methodologies.

1 7 Summary and Future Work

2 7.1 Summary

3 The primary objective of this research was to develop DL solutions aimed at enhancing
4 automation in the ultrasonic inspection of aerospace composites. To this end, the four
5 chapters that form the core of the research presented in this thesis each contribute to a
6 distinct aspect of applying DL to UT of composites. Supported closely by Spirit
7 AeroSystems, this research holds significant potential for advancing automation
8 within industrial settings.

9 Chapter 1 began by outlining the industrial motivation that contextualised this
10 research. It underscored the growing significance of CFRP in the aerospace industry,
11 which in turn highlighted the importance of NDE for ensuring component integrity.
12 UT was identified as the most widely used NDE technique for inspecting aerospace
13 composites. Advancements and the increasing uptake of robotics for sensor
14 deployment has greatly progressed NDE automation. Despite this, the interpretation
15 of UT results remained a largely manual process. This manual interpretation not only
16 limits NDE automation, but also creates a bottleneck in the manufacturing process and
17 introduces the potential for human error.

18 Chapter 2 provided an overview of background knowledge relevant to the thesis. It
19 began with an introduction to the range of NDE techniques commonly used for
20 composites, focusing on UT and its application in the aerospace industry. The use of
21 phased arrays and mechanised scanning was also discussed, demonstrating how large-
22 scale aerospace components can be inspected in industrial settings. The chapter also

1 explored the application of composites in aerospace and identified typical defects that
2 can occur. Background information on AI, with an emphasis on CNNs was provided
3 due to their importance in this research. Finally, the current use of ML in NDE was
4 investigated, addressing key challenges such as data scarcity, problem definition, and
5 model evaluation.

6 Chapter 3 addressed what is arguably the biggest challenge in applying ML to NDE:
7 the scarcity of training data. To evaluate this, a CNN classifier was employed and four
8 different methods for generating synthetic UT C-scans were investigated. Each method
9 aimed to bridge the sim-to-real domain gap between simulated and real defect
10 responses. Among these methods, one utilised a generative DL model, specifically
11 CycleGAN. A modification to the model's loss function was introduced to constrain
12 significant changes to defect responses while allowing the model to freely modify the
13 noise response. Although the modified CycleGAN produced the most effective
14 synthetic data, it was deduced that other methods performed nearly as well and might
15 be preferable due to their robustness and ease of implementation. Guided Grad-CAM
16 was employed to compare models trained on synthetic and experimental data. While
17 this provided only a qualitative analysis, it suggested that models trained on synthetic
18 data use similar features for classification as those trained on experimental data. This
19 finding reinforces the potential of synthetic data to mitigate the issue of training data
20 scarcity and helps to limit the additional interpretability concern which may occur from
21 the use of synthetic training data. This work proved fundamental to the subsequent
22 chapters 4 and 5, as it established a robust framework for generating synthetic training
23 data, which was crucial for the development and training of more advanced machine
24 learning models.

1 Chapter 4 delved into the potential of harnessing the complete 3D ultrasonic volume
2 for defect detection. While human operators typically visualise this volume as images
3 for easier interpretation, ML systems are not confined to such 2D representations. The
4 underlying intuition was that by refraining from compressing the time or spatial
5 domain during imaging, the ML model could glean richer information to assess defect
6 responses. This endeavour necessitated the generation of synthetic volumetric datasets
7 for training, building upon the A-scan noise generation method introduced in previous
8 chapters. This method, initially designed for generating synthetic C-scan images, was
9 extended to cater to volumetric datasets. A NAS was conducted, leading to the
10 discovery of an optimal architecture that was then evaluated against two other models.
11 To further bridge the gap between synthetic and experimental domains, two domain-
12 specific augmentation methods were introduced. These methods significantly
13 enhanced the classification performance. The study conclusively demonstrated the
14 feasibility of training a DL model to effectively detect defects using the complete
15 volumetric ultrasonic data.

16 Chapter 5 advanced beyond the scope of the previous chapter by focusing on defect
17 segmentation within the ultrasonic volume. This expanded approach provided valuable
18 additional information for defect characterisation, such as defect sizing along with in-
19 plane and through-thickness localisation. These aspects are crucial for tasks such as
20 conducting repairs or constructing digital twins, along with determining the criticality
21 of defects. To achieve this, a modified 3D U-Net was trained in a fully supervised
22 manner using synthetic defect responses and their corresponding defect masks.
23 Evaluation of the localisation performance was conducted against known defect depth
24 and 6 dB centroid, yielding a MAE of 0.08 mm and 0.57 mm, respectively. Sizing

1 accuracy was compared to the 6 dB drop criterion. Initially, there was a 55% increase
2 in error compared to the 6 dB drop criterion. However, after applying a correction
3 factor to account for the disparity between the experimental and synthetic domains, a
4 35% reduction in error was achieved over the 6 dB drop method in defect sizing.
5 Furthermore, the model was subjected to testing on out-of-distribution defect
6 responses from PTFE inserts. This testing revealed a 71% increase in defect sizing
7 error compared to in-distribution defects. This underscores a notable limitation of
8 supervised training, which necessitates adequate coverage of the test distribution of
9 defects during training. Despite the utility of synthetic data in mitigating data scarcity,
10 modelling the full range of defects and the inherent variability within defect types, as
11 necessary for wide generalisability, remains a challenge in NDE.

12 Chapter 6 furthered the research into segmentation of the full ultrasonic volume by
13 developing a self-supervised approach for training a segmentation model. In contrast
14 to chapter 5, which employed supervised learning, this method did not require
15 examples of defect responses during training and only used defect free data. This
16 represents a significant advantage as experimental defect-free data is much more
17 readily available. As noted in the previous chapter, fully labelled experimental
18 responses are rarely available, and accurately simulating the range of defect responses
19 required for training a supervised model is a substantial challenge. The SSL model
20 demonstrated good results for in-plane and through-thickness localisation, with MAEs
21 of 0.26 mm and 0.37 mm, respectively. Since the SSL method is not trained to identify
22 true defect size in the physical domain from the ultrasonic domain, the initial MAE for
23 defect sizing was greater than that achieved with fully supervised training. However,
24 due to the consistency of the method, it is possible to calibrate to physical defect sizes.

1 This calibration resulted in an MAE of less than the element pitch, specifically 0.58
2 mm. Furthermore, by treating the segmentation problem as one of anomaly detection
3 the method generalises well to different defect responses as demonstrated with testing
4 on PTFE inserts as well as FBHs. The effectiveness of the method was demonstrated
5 on an industrial component. Where industrial data collection limits the availability of
6 enveloped data processing, it was evidenced that an appropriate moving average filter
7 could be applied to allows the applicability of this method. This approach
8 simultaneously solves the problem of needing defect responses during training, allows
9 for full 3D segmentation, and demonstrates impressive generalisability.

10 7.2 Future Work and Final Remarks

11 There remain challenges in applying DL to NDE. Many valuable lessons have been
12 learnt during this research and are useful for future research or industrial DL NDE
13 tasks. The author believes that, in general, detection problems should be framed as
14 issues of anomaly detection and addressed using unsupervised or self-supervised
15 methods. This is primarily due to the lack of large training datasets available for
16 supervised training and improved generalisability. It has been demonstrated that even
17 defect segmentation, which traditionally relies on supervised learning, can be reframed
18 in this manner with significant benefits. In the future, large open-source datasets may
19 become available, facilitating the development of generalisable supervised models.
20 However, until this happens, alternative forms of training are likely to be more robust
21 to the variability encountered in NDE and even if these were to become available may
22 still be inferior to supervised training methods.

23 There is still scope for further investigation into the research presented in this thesis.

24 The most noteworthy future work would be the incorporation of a wide range of

1 naturally occurring defects into the test sets to evaluate model generalisability. As
2 observed in Chapter 5, model performance can decrease when exposed to samples
3 outside of the training distribution. Before deploying DL-based solutions in industry,
4 it is crucial to robustly test them on a variety of naturally occurring defects. If such
5 testing were conducted, one would likely see a decrease in performance for the
6 presented models trained in a supervised manner, due to the lack of extensive training
7 data.

8 Chapters 3-5 demonstrated that synthetic data can be effective for training models.
9 However, to produce generalisable models for industrial applications, it would likely
10 require generating a much larger synthetic training set that encompasses a broader
11 range of variability seen in naturally occurring defects. This would necessitate the
12 accurate simulation of a greater variety of defect responses, encompassing different
13 classes of defects as well as greater intra-class variability. In addition to variability in
14 defects, the base sample should also be diversified, to account for differences in
15 components such as ply layup and geometric changes. Should these datasets become
16 available, much of the work conducted in Chapters 4-5 could be revisited, with
17 retraining and evaluation to cover a wider spectrum of defects and component
18 variability.

19 Chapters 5 and 6 demonstrated how volumetric segmentation could be used for defect
20 characterisation. However, for complete characterisation, a key piece of information
21 is missing: the class of defect. Multi-class classification was not explored in this work
22 primarily due to the lack of available test data. The volumetric defect detection model
23 presented in Chapter 4 could be easily modified and retrained to provide multi-class

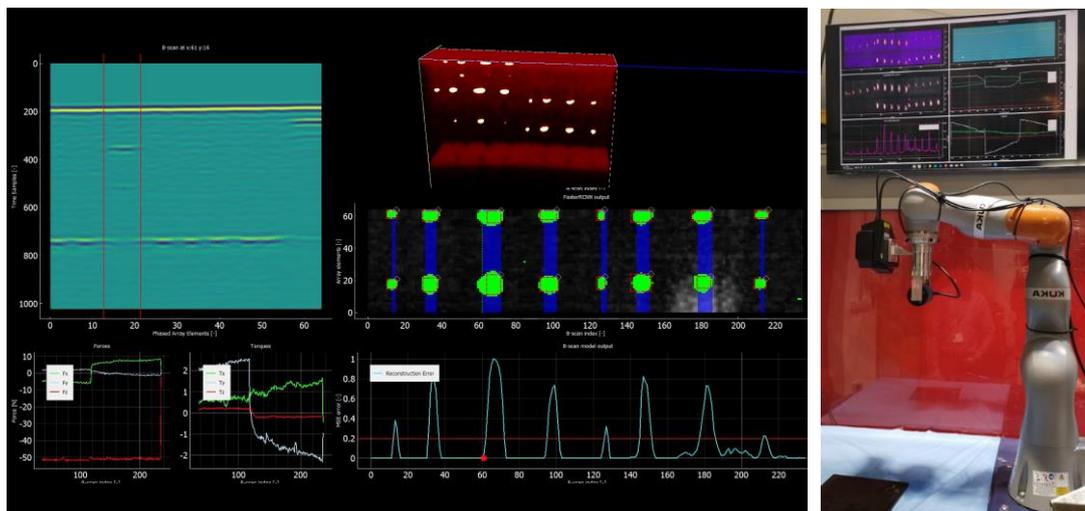
1 predictions. Here, the additional information available by analysing the full volume
2 may prove beneficial in improving classification accuracy over 2D image analysis. To
3 train such a model, an expanded synthetic (or experimental) training dataset with
4 associated class labels would be required. Alternatively, SSL approaches such as
5 Siamese networks [198] could be effective if large labelled training datasets are
6 unavailable. However, any approach would still necessitate a test set with multiple
7 defect classes for evaluation.

8 Whilst chapter 6 demonstrated the potential of the methodology on industrial
9 components, further exploration would be beneficial. Industrial scans are likely to vary
10 in quality and sensing methodologies, which may necessitate further research. Real
11 components can vary significantly in geometry, making automated detection more
12 challenging, as geometrical responses can often be indistinguishable from defect
13 responses without additional information. While Chapter 6 demonstrated the
14 effectiveness of the SSL model on stepped samples and an industrial sample with
15 stringers, the model necessitated inference parallel to any change in thickness.
16 Industrial components with more complex geometries could limit its applicability or
17 require modification of the methodology.

18 In addition to testing on a greater variety of industrial components, there is also a
19 regulatory factor to consider for industrial adoption. Further work would be beneficial
20 to assess how these tools and systems could be used in industry and what is required
21 to approve them for use. Regulatory approval may become the biggest barrier to the
22 uptake and use of these techniques, necessitating comprehensive evaluations to ensure
23 compliance with industry standards and regulations. Addressing these challenges will

1 be crucial for the successful deployment of these advanced methodologies in real-
2 world industrial settings.

3 Initial work has been undertaken to build on this DL research to produce automated
4 UT software solutions with robotic integration. Figure 74 shows a screenshot of the
5 developed software and its integration with a collaborative robot-based sensor delivery
6 platform. This system was recently presented at the BINDT Aerospace Conference
7 [199]. However, there remains significant potential for further development. This
8 could involve collaborating with NDE operators to create more specialised and
9 effective human-in-the-loop automation tools, designed with operator feedback. The
10 integration with robotic systems paves the way for future advancements in automation,
11 including the ability to automatically re-scan or perform rework on components. By
12 taking a holistic approach, incorporating operator insights, and advancing robotic
13 integration, the system could become more adaptive and efficient, leading to improved
14 defect detection and repair processes in industrial applications.



15 *Figure 74: Screenshot of the graphical user interface developed to interact with DL model outputs*
16 *from UT scans (left) and demonstration of the system integrated into a flexible robotic scanning*
17 *system (right).*

8 Bibliography

- [1] J. Zhang, G. Lin, U. Vaidya, and H. Wang, ‘Past, present and future prospective of global carbon fibre composite developments and applications’, *Composites Part B: Engineering*, vol. 250, p. 110463, Feb. 2023, doi: 10.1016/j.compositesb.2022.110463.
- [2] J. Skoczylas, S. Samborski, and M. Kłonica, ‘THE APPLICATION OF COMPOSITE MATERIALS IN THE AEROSPACE INDUSTRY’, *Journal of Technology and Exploitation in Mechanical Engineering*, vol. 5, Jan. 2019, doi: 10.35784/jteme.73.
- [3] B. Parveez, M. I. Kittur, I. A. Badruddin, S. Kamangar, M. Hussien, and M. A. Umarfarooq, ‘Scientific Advancements in Composite Materials for Aircraft Applications: A Review’, *Polymers (Basel)*, vol. 14, no. 22, p. 5007, Nov. 2022, doi: 10.3390/polym14225007.
- [4] J. Jodhani, A. Handa, A. Gautam, Ashwni, and R. Rana, ‘Ultrasonic non-destructive evaluation of composites: A review’, *Materials Today: Proceedings*, vol. 78, pp. 627–632, Jan. 2023, doi: 10.1016/j.matpr.2022.12.055.
- [5] R. Henrich and U. Schnars, ‘Applications of NDT Methods on Composite Structures in Aerospace Industry’, *e-Journal of Nondestructive Testing*, vol. 11, no. 12, Dec. 2006, Accessed: Jun. 19, 2024. [Online]. Available: <https://www.ndt.net/search/docs.php3?id=4180&msgID=0&rootID=0>
- [6] S. Cumblidge, A. D’Agostino, S. Morrow, C. Franklin, and N. Nughes, ‘Review of Human Factors Research in Nondestructive Examination.’, in *US Nuclear Regulatory Commission-Pacific Northwest National Laboratory 7th European-American Workshop on Reliability of NDE.*, Feb. 2017.
- [7] P. Gardner *et al.*, ‘Machine learning at the interface of structural health monitoring and non-destructive evaluation’, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 378, no. 2182, p. 20190581, Oct. 2020, doi: 10.1098/rsta.2019.0581.
- [8] C. Mineo *et al.*, ‘Robotic Geometric and Volumetric Inspection of High Value and Large Scale Aircraft Wings’, in *2019 IEEE 5th International Workshop on Metrology for AeroSpace (MetroAeroSpace)*, Jun. 2019, pp. 82–86. doi:

- 1 10.1109/MetroAeroSpace.2019.8869667.
- 2 [9] L. Séguin-Charbonneau, J. Walter, L.-D. Thérout, L. Scheed, A. Beausoleil, and
3 B. Masson, ‘Automated defect detection for ultrasonic inspection of CFRP
4 aircraft components’, *NDT & E International*, vol. 122, p. 102478, Sep. 2021,
5 doi: 10.1016/j.ndteint.2021.102478.
- 6 [10] ‘Introduction to non-destructive testing’, Aerospace Testing International.
7 Accessed: Nov. 17, 2021. [Online]. Available:
8 [https://www.aerospacetestinginternational.com/features/introduction-to-non-](https://www.aerospacetestinginternational.com/features/introduction-to-non-destructive-testing.html)
9 [destructive-testing.html](https://www.aerospacetestinginternational.com/features/introduction-to-non-destructive-testing.html)
- 10 [11] F. W. Margrave, K. Rigas, D. A. Bradley, and P. Barrowcliffe, ‘The use of neural
11 networks in ultrasonic flaw detection’, *Measurement*, vol. 25, no. 2, pp. 143–154,
12 Mar. 1999, doi: 10.1016/S0263-2241(98)00075-X.
- 13 [12] J. Ye, S. Ito, and N. Toyama, ‘Computerized Ultrasonic Imaging Inspection: From
14 Shallow to Deep Learning’, *Sensors (Basel)*, vol. 18, no. 11, p. 3820, Nov. 2018,
15 doi: 10.3390/s18113820.
- 16 [13] B. Valeske, A. Osman, F. Römer, and R. Tschuncky, ‘Next Generation NDE
17 Sensor Systems as IIoT Elements of Industry 4.0’, *Research in Nondestructive*
18 *Evaluation*, vol. 31, no. 5–6, pp. 340–369, Nov. 2020, doi:
19 10.1080/09349847.2020.1841862.
- 20 [14] S. Cantero-Chinchilla, P. D. Wilcox, and A. J. Croxford, ‘Deep learning in
21 automated ultrasonic NDE -- developments, axioms and opportunities’,
22 *arXiv:2112.06650 [eess]*, Dec. 2021, Accessed: Jan. 12, 2022. [Online].
23 Available: <http://arxiv.org/abs/2112.06650>
- 24 [15] S. Gholizadeh, ‘A review of non-destructive testing methods of composite
25 materials’, *Procedia Structural Integrity*, vol. 1, pp. 50–57, 2016, doi:
26 10.1016/j.prostr.2016.02.008.
- 27 [16] M. V. Felice and Z. Fan, ‘Sizing of flaws using ultrasonic bulk wave testing: A
28 review’, *Ultrasonics*, vol. 88, pp. 26–42, Aug. 2018, doi:
29 10.1016/j.ultras.2018.03.003.
- 30 [17] S. C. Olisa, M. A. Khan, and A. Starr, ‘Review of Current Guided Wave
31 Ultrasonic Testing (GWUT) Limitations and Future Directions’, *Sensors (Basel)*,
32 vol. 21, no. 3, p. 811, Jan. 2021, doi: 10.3390/s21030811.

- 1 [18] C. Meola, S. Boccardi, G. M. Carlomagno, N. D. Boffa, E. Monaco, and F. Ricci,
2 'Nondestructive evaluation of carbon fibre reinforced composites with infrared
3 thermography and ultrasonics', *Composite Structures*, vol. 134, pp. 845–853,
4 Dec. 2015, doi: 10.1016/j.compstruct.2015.08.119.
- 5 [19] D. K. Hsu, '15 - Non-destructive evaluation (NDE) of aerospace composites:
6 ultrasonic techniques', in *Non-Destructive Evaluation (NDE) of Polymer Matrix*
7 *Composites*, V. M. Karbhari, Ed., in Woodhead Publishing Series in Composites
8 Science and Engineering. , Woodhead Publishing, 2013, pp. 397–422. doi:
9 10.1533/9780857093554.3.397.
- 10 [20] 'Ultrasonic Test - an overview | ScienceDirect Topics'. Accessed: Jun. 14, 2024.
11 [Online]. Available: [https://www.sciencedirect.com/topics/physics-and-](https://www.sciencedirect.com/topics/physics-and-astronomy/ultrasonic-test)
12 [astronomy/ultrasonic-test](https://www.sciencedirect.com/topics/physics-and-astronomy/ultrasonic-test)
- 13 [21] I. Papa, V. Lopresto, and A. Langella, 'Ultrasonic inspection of composites
14 materials: Application to detect impact damage', *International Journal of*
15 *Lightweight Materials and Manufacture*, vol. 4, no. 1, pp. 37–42, Mar. 2021, doi:
16 10.1016/j.ijlmm.2020.04.002.
- 17 [22] C. Holmes, B. W. Drinkwater, and P. D. Wilcox, 'Post-processing of the full
18 matrix of ultrasonic transmit–receive array data for non-destructive evaluation',
19 *NDT & E International*, vol. 38, no. 8, pp. 701–711, Dec. 2005, doi:
20 10.1016/j.ndteint.2005.04.002.
- 21 [23] S. Moon *et al.*, 'FEA-Based Ultrasonic Focusing Method in Anisotropic Media
22 for Phased Array Systems', *Applied Sciences*, vol. 11, no. 19, Art. no. 19, Jan.
23 2021, doi: 10.3390/app11198888.
- 24 [24] C. Aguado, 'Ultrasonic Techniques and Industrial Robots: Natural Evolution of
25 Inspection Systems'.
- 26 [25] D. Garnier, D. Garnier, P. Louviot, and A. Tachattahte, 'Robotised UT
27 Transmission NDT of Composite Complex Shaped Parts'.
- 28 [26] C. Mineo *et al.*, 'Flexible integration of robotics, ultrasonics and metrology for
29 the inspection of aerospace components', *AIP Conference Proceedings*, vol.
30 1806, no. 1, p. 020026, Feb. 2017, doi: 10.1063/1.4974567.
- 31 [27] R. A. Smith and B. Clarke, 'Ultrasonic C-scan determination of ply stacking
32 sequence in carbon-fibre composites.', *Insight - Non-Destructive Testing and*

- 1 *Condition Monitoring*, vol. 36, no. 10, pp. 741–747, 1994.
- 2 [28] P. Vaara and J. Leinonen, ‘Technology Survey on NDT of Carbon-fiber
3 Composites’, 2012. Accessed: Feb. 26, 2024. [Online]. Available:
4 [https://www.semanticscholar.org/paper/Technology-Survey-on-NDT-of-Carbon-](https://www.semanticscholar.org/paper/Technology-Survey-on-NDT-of-Carbon-fiber-Composites-Vaara-Leinonen/7d8b35f344c13149b177995ed210431f5d36dff0)
5 [fiber-Composites-Vaara-](https://www.semanticscholar.org/paper/Technology-Survey-on-NDT-of-Carbon-fiber-Composites-Vaara-Leinonen/7d8b35f344c13149b177995ed210431f5d36dff0)
6 [Leinonen/7d8b35f344c13149b177995ed210431f5d36dff0](https://www.semanticscholar.org/paper/Technology-Survey-on-NDT-of-Carbon-fiber-Composites-Vaara-Leinonen/7d8b35f344c13149b177995ed210431f5d36dff0)
- 7 [29] B. Djordjevic, ‘Non Destructive Test Technology for the Composite’, p. 7, Jan.
8 2009.
- 9 [30] D. Moore and A. M.-E. Dorado, ‘Composite Material Characterization using
10 Acoustic Wave Speed Measurements’, presented at the ASNT Annual
11 Conference, 2015, p. 5. [Online]. Available: <https://www.osti.gov/biblio/1326353>
- 12 [31] Ley, O. and V. Godinez, ‘Non-destructive evaluation (NDE) of aerospace
13 composites: application of infrared (IR) thermography’, doi:
14 10.1533/9780857093554.3.309.
- 15 [32] A. Kokurov and D. Subbotin, ‘Ultrasonic detection of manufacturing defects in
16 multilayer composite structures’, *IOP Conference Series: Materials Science and*
17 *Engineering*, vol. 1023, p. 012013, Jan. 2021, doi: 10.1088/1757-
18 899X/1023/1/012013.
- 19 [33] F. Heinecke and C. Willberg, ‘Manufacturing-Induced Imperfections in
20 Composite Parts Manufactured via Automated Fiber Placement’, *J. Compos. Sci.*,
21 vol. 3, no. 2, p. 56, Jun. 2019, doi: 10.3390/jcs3020056.
- 22 [34] M. Jolly *et al.*, ‘Review of Non-destructive Testing (NDT) Techniques and their
23 Applicability to Thick Walled Composites’, *Procedia CIRP*, vol. 38, pp. 129–
24 136, Jan. 2015, doi: 10.1016/j.procir.2015.07.043.
- 25 [35] H. Taheri and A. A. Hassen, ‘Nondestructive Ultrasonic Inspection of Composite
26 Materials: A Comparative Advantage of Phased Array Ultrasonic’, *Applied*
27 *Sciences*, vol. 9, no. 8, Art. no. 8, Jan. 2019, doi: 10.3390/app9081628.
- 28 [36] S. Bayat, A. Jamzad, N. Zobeiry, A. Poursartip, P. Mousavi, and P. Abolmaesumi,
29 ‘Temporal enhanced Ultrasound: A new method for detection of porosity defects
30 in composites’, *Composites Part A: Applied Science and Manufacturing*, vol.
31 164, p. 107259, Jan. 2023, doi: 10.1016/j.compositesa.2022.107259.
- 32 [37] G. Fernlund, J. Wells, L. Fahrang, J. Kay, and A. Poursartip, ‘Causes and

- 1 remedies for porosity in composite manufacturing’, *IOP Conf. Ser.: Mater. Sci.*
2 *Eng.*, vol. 139, p. 012002, Jul. 2016, doi: 10.1088/1757-899X/139/1/012002.
- 3 [38] K. Senthil, A. Arockiarajan, R. Palaninathan, B. Santhosh, and K. M. Usha,
4 ‘Defects in composite structures: Its effects and prediction methods – A
5 comprehensive review’, *Composite Structures*, vol. 106, pp. 139–149, Dec. 2013,
6 doi: 10.1016/j.compstruct.2013.06.008.
- 7 [39] R. Telford, A. O’Carroll, R. S. Pierce, and T. M. Young, ‘A novel method to
8 produce kiss-bonds in composites components for NDI and characterisation
9 purposes’, *Composites Part B: Engineering*, vol. 173, p. 106926, Sep. 2019, doi:
10 10.1016/j.compositesb.2019.106926.
- 11 [40] A. Poudel, S. S. Shrestha, J. S. Sandhu, T. P. Chu, and C. G. Pergantis,
12 ‘Comparison and analysis of Acoustography with other NDE techniques for
13 foreign object inclusion detection in graphite epoxy composites’, *Composites*
14 *Part B: Engineering*, vol. 78, pp. 86–94, Sep. 2015, doi:
15 10.1016/j.compositesb.2015.03.048.
- 16 [41] H. M. Hsiao and I. M. Daniel, ‘Effect of fiber waviness on stiffness and strength
17 reduction of unidirectional composites under compressive loading’, *Composites*
18 *Science and Technology*, vol. 56, no. 5, pp. 581–593, Jan. 1996, doi:
19 10.1016/0266-3538(96)00045-0.
- 20 [42] G. Karami and M. Garnich, ‘Effective moduli and failure considerations for
21 composites with periodic fiber waviness’, *Composite Structures*, vol. 67, no. 4,
22 pp. 461–475, Mar. 2005, doi: 10.1016/j.compstruct.2004.02.005.
- 23 [43] C. Shah, S. Bosse, and A. von Hehl, ‘Taxonomy of Damage Patterns in
24 Composite Materials, Measuring Signals, and Methods for Automated Damage
25 Diagnostics’, *Materials (Basel)*, vol. 15, no. 13, p. 4645, Jul. 2022, doi:
26 10.3390/ma15134645.
- 27 [44] K. Ajay, ‘Best Practice Guide Non-Destructive Testing of Composites’. National
28 Composites Network, 2013. [Online]. Available:
29 <https://avaloncs1.files.wordpress.com/2013/01/ncn-best-practice-ndt.pdf>
- 30 [45] R. H. Bossi and G. E. Georgeson, ‘16 - Nondestructive testing of aerospace
31 composites’, in *Polymer Composites in the Aerospace Industry (Second Edition)*,
32 P. Irving and C. Soutis, Eds., in Woodhead Publishing Series in Composites

- 1 Science and Engineering. , Woodhead Publishing, 2020, pp. 461–489. doi:
2 10.1016/B978-0-08-102679-3.00016-2.
- 3 [46] W. Ertel, *Introduction to Artificial Intelligence*. in Undergraduate Topics in
4 Computer Science. Cham: Springer International Publishing, 2017. doi:
5 10.1007/978-3-319-58487-4.
- 6 [47] A. Toosi, A. G. Bottino, B. Saboury, E. Siegel, and A. Rahmim, ‘A Brief History
7 of AI: How to Prevent Another Winter (A Critical Review)’, *PET Clin*, vol. 16,
8 no. 4, pp. 449–469, Oct. 2021, doi: 10.1016/j.cpet.2021.07.001.
- 9 [48] G. F. Luger and W. A. Stubblefield, *Artificial Intelligence and the Design of*
10 *Expert Systems*. Benjamin/Cummings Publishing Company, 1989.
- 11 [49] M. Campbell, A. J. Hoane, and F. Hsu, ‘Deep Blue’, *Artificial Intelligence*, vol.
12 134, no. 1, pp. 57–83, Jan. 2002, doi: 10.1016/S0004-3702(01)00129-1.
- 13 [50] J. Weizenbaum, ‘ELIZA—a computer program for the study of natural language
14 communication between man and machine’, *Commun. ACM*, vol. 9, no. 1, pp.
15 36–45, Jan. 1966, doi: 10.1145/365153.365168.
- 16 [51] J. M. Mira, ‘Symbols versus connections: 50 years of artificial intelligence’,
17 *Neurocomputing*, vol. 71, no. 4, pp. 671–680, Jan. 2008, doi:
18 10.1016/j.neucom.2007.06.009.
- 19 [52] Y. Shang, ‘5 - Expert Systems’, in *The Electrical Engineering Handbook*, W.-K.
20 Chen, Ed., Burlington: Academic Press, 2005, pp. 367–377. doi: 10.1016/B978-
21 012170960-0/50031-1.
- 22 [53] P. Smolensky, ‘Connectionist AI, symbolic AI, and the brain’, *Artif Intell Rev*,
23 vol. 1, no. 2, pp. 95–109, Jun. 1987, doi: 10.1007/BF00130011.
- 24 [54] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, ‘Learning representations by
25 back-propagating errors’, *Nature*, vol. 323, no. 6088, pp. 533–536, Oct. 1986,
26 doi: 10.1038/323533a0.
- 27 [55] I. H. Sarker, ‘AI-Based Modeling: Techniques, Applications and Research Issues
28 Towards Automation, Intelligent and Smart Systems’, *Sn Computer Science*, vol.
29 3, no. 2, p. 158, Feb. 2022, doi: 10.1007/s42979-022-01043-x.
- 30 [56] M. Elahi, S. O. Afolaranmi, J. L. Martinez Lastra, and J. A. Perez Garcia, ‘A
31 comprehensive literature review of the applications of AI techniques through the
32 lifecycle of industrial equipment’, *Discov Artif Intell*, vol. 3, no. 1, p. 43, Dec.

- 1 2023, doi: 10.1007/s44163-023-00089-x.
- 2 [57] K. Ofosu-Ampong, ‘Artificial intelligence research: A review on dominant
3 themes, methods, frameworks and future research directions’, *Telematics and*
4 *Informatics Reports*, vol. 14, p. 100127, Jun. 2024, doi:
5 10.1016/j.teler.2024.100127.
- 6 [58] S. A. Alowais *et al.*, ‘Revolutionizing healthcare: the role of artificial intelligence
7 in clinical practice’, *BMC Medical Education*, vol. 23, no. 1, p. 689, Sep. 2023,
8 doi: 10.1186/s12909-023-04698-z.
- 9 [59] D. Khurana, A. Koli, K. Khatter, and S. Singh, ‘Natural language processing:
10 state of the art, current trends and challenges’, *Multimed Tools Appl*, vol. 82, no.
11 3, pp. 3713–3744, Jan. 2023, doi: 10.1007/s11042-022-13428-4.
- 12 [60] R. Szeliski, *Computer Vision: Algorithms and Applications*. Springer Nature,
13 2022.
- 14 [61] W. S. McCulloch and W. Pitts, ‘A logical calculus of the ideas immanent in
15 nervous activity’, *Bulletin of Mathematical Biophysics*, vol. 5, no. 4, pp. 115–
16 133, Dec. 1943, doi: 10.1007/BF02478259.
- 17 [62] F. Rosenblatt, *The Perceptron, a Perceiving and Recognizing Automaton Project*
18 *Para*. Cornell Aeronautical Laboratory, 1957.
- 19 [63] F. Rosenblatt, ‘The perceptron: A probabilistic model for information storage and
20 organization in the brain’, *Psychological Review*, vol. 65, no. 6, pp. 386–408,
21 1958, doi: 10.1037/h0042519.
- 22 [64] K. Hornik, M. Stinchcombe, and H. White, ‘Multilayer feedforward networks are
23 universal approximators’, *Neural Networks*, vol. 2, no. 5, pp. 359–366, Jan. 1989,
24 doi: 10.1016/0893-6080(89)90020-8.
- 25 [65] Y. LeCun, Y. Bengio, and G. Hinton, ‘Deep learning’, *Nature*, vol. 521, no. 7553,
26 pp. 436–444, May 2015, doi: 10.1038/nature14539.
- 27 [66] Y. Bengio, P. Simard, and P. Frasconi, ‘Learning long-term dependencies with
28 gradient descent is difficult’, *IEEE Transactions on Neural Networks*, vol. 5, no.
29 2, pp. 157–166, Mar. 1994, doi: 10.1109/72.279181.
- 30 [67] K. He, X. Zhang, S. Ren, and J. Sun, ‘Deep Residual Learning for Image
31 Recognition’, Dec. 10, 2015, *arXiv*: arXiv:1512.03385. doi:
32 10.48550/arXiv.1512.03385.

- 1 [68] L. Ciampiconi, A. Elwood, M. Leonardi, A. Mohamed, and A. Rozza, ‘A survey
2 and taxonomy of loss functions in machine learning’, Jan. 13, 2023, *arXiv*:
3 arXiv:2301.05579. doi: 10.48550/arXiv.2301.05579.
- 4 [69] Q. Wang, Y. Ma, K. Zhao, and Y. Tian, ‘A Comprehensive Survey of Loss
5 Functions in Machine Learning’, *Ann. Data. Sci.*, vol. 9, no. 2, pp. 187–212, Apr.
6 2022, doi: 10.1007/s40745-020-00253-5.
- 7 [70] D. P. Kingma and J. Ba, ‘Adam: A Method for Stochastic Optimization’, Jan. 29,
8 2017, *arXiv*: arXiv:1412.6980. doi: 10.48550/arXiv.1412.6980.
- 9 [71] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, ‘Densely Connected
10 Convolutional Networks’, *arXiv:1608.06993 [cs]*, Jan. 2018, Accessed: Oct. 04,
11 2021. [Online]. Available: <http://arxiv.org/abs/1608.06993>
- 12 [72] R. Abdulkadimov, P. Lyakhov, and N. Nagornov, ‘Survey of Optimization
13 Algorithms in Modern Neural Networks’, *Mathematics*, vol. 11, no. 11, Art. no.
14 11, Jan. 2023, doi: 10.3390/math11112466.
- 15 [73] R. Sun, ‘Optimization for deep learning: theory and algorithms’, Dec. 18, 2019,
16 *arXiv*: arXiv:1912.08957. doi: 10.48550/arXiv.1912.08957.
- 17 [74] J. Chai, H. Zeng, A. Li, and E. W. T. Ngai, ‘Deep learning in computer vision: A
18 critical review of emerging techniques and application scenarios’, *Machine*
19 *Learning with Applications*, vol. 6, p. 100134, Dec. 2021, doi:
20 10.1016/j.mlwa.2021.100134.
- 21 [75] N. O’Mahony *et al.*, ‘Deep Learning vs. Traditional Computer Vision’, in
22 *Advances in Computer Vision*, vol. 943, K. Arai and S. Kapoor, Eds., in *Advances*
23 *in Intelligent Systems and Computing*, vol. 943. , Cham: Springer International
24 Publishing, 2020, pp. 128–144. doi: 10.1007/978-3-030-17795-9_10.
- 25 [76] L. Liu *et al.*, ‘Deep Learning for Generic Object Detection: A Survey’, *Int J*
26 *Comput Vis*, vol. 128, no. 2, pp. 261–318, Feb. 2020, doi: 10.1007/s11263-019-
27 01247-4.
- 28 [77] S. Albelwi and A. Mahmood, ‘A Framework for Designing the Architectures of
29 Deep Convolutional Neural Networks’, *Entropy*, vol. 19, no. 6, Art. no. 6, Jun.
30 2017, doi: 10.3390/e19060242.
- 31 [78] A. Krizhevsky, I. Sutskever, and G. E. Hinton, ‘ImageNet Classification with
32 Deep Convolutional Neural Networks’, in *Advances in Neural Information*

- 1 *Processing Systems*, Curran Associates, Inc., 2012. Accessed: Feb. 27, 2024.
2 [Online]. Available:
3 [https://proceedings.neurips.cc/paper_files/paper/2012/hash/c399862d3b9d6b76](https://proceedings.neurips.cc/paper_files/paper/2012/hash/c399862d3b9d6b76c8436e924a68c45b-Abstract.html)
4 [c8436e924a68c45b-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2012/hash/c399862d3b9d6b76c8436e924a68c45b-Abstract.html)
- 5 [79] ‘ImageNet’. Accessed: Jan. 29, 2022. [Online]. Available: <https://image-net.org/>
- 6 [80] K. Simonyan and A. Zisserman, ‘Very Deep Convolutional Networks for Large-
7 Scale Image Recognition’, Apr. 10, 2015, *arXiv*: arXiv:1409.1556. doi:
8 10.48550/arXiv.1409.1556.
- 9 [81] A. Dosovitskiy *et al.*, ‘An Image is Worth 16x16 Words: Transformers for Image
10 Recognition at Scale’, Jun. 03, 2021, *arXiv*: arXiv:2010.11929. doi:
11 10.48550/arXiv.2010.11929.
- 12 [82] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, ‘A ConvNet
13 for the 2020s’, Mar. 02, 2022, *arXiv*: arXiv:2201.03545. Accessed: Jun. 16, 2023.
14 [Online]. Available: <http://arxiv.org/abs/2201.03545>
- 15 [83] Z. Liu *et al.*, ‘Swin Transformer: Hierarchical Vision Transformer using Shifted
16 Windows’, Aug. 17, 2021, *arXiv*: arXiv:2103.14030. doi:
17 10.48550/arXiv.2103.14030.
- 18 [84] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, P. Martinez-
19 Gonzalez, and J. Garcia-Rodriguez, ‘A survey on deep learning techniques for
20 image and video semantic segmentation’, *Applied Soft Computing*, vol. 70, pp.
21 41–65, Sep. 2018, doi: 10.1016/j.asoc.2018.05.018.
- 22 [85] A. Ghosh, N. D. Jana, S. Das, and R. Mallipeddi, ‘Two-Phase Evolutionary
23 Convolutional Neural Network Architecture Search for Medical Image
24 Classification’, *IEEE Access*, vol. 11, pp. 115280–115305, 2023, doi:
25 10.1109/ACCESS.2023.3323705.
- 26 [86] ‘CIFAR-10 and CIFAR-100 datasets’. Accessed: Jan. 29, 2022. [Online].
27 Available: <http://www.cs.toronto.edu/~kriz/cifar.html>
- 28 [87] X. Zhu, C. Vondrick, C. Fowlkes, and D. Ramanan, ‘Do We Need More Training
29 Data?’, *Int J Comput Vis*, vol. 119, no. 1, pp. 76–92, Aug. 2016, doi:
30 10.1007/s11263-015-0812-2.
- 31 [88] C. Shorten and T. M. Khoshgoftaar, ‘A survey on Image Data Augmentation for
32 Deep Learning’, *Journal of Big Data*, vol. 6, no. 1, p. 60, Jul. 2019, doi:

- 1 10.1186/s40537-019-0197-0.
- 2 [89] Y. Sun, P. Bai, H. Sun, and P. Zhou, ‘Real-time automatic detection of weld
3 defects in steel pipe’, *NDT & E International*, vol. 38, no. 7, pp. 522–528, Oct.
4 2005, doi: 10.1016/j.ndteint.2005.01.011.
- 5 [90] A. A. Markov, D. A. Shpagin, and M. N. Shilov, ‘Ultrasonic Multichannel Flaw
6 Detector for Testing Rails with Signal Recording’, *Russian Journal of*
7 *Nondestructive Testing*, vol. 39, no. 2, pp. 105–114, Feb. 2003, doi:
8 10.1023/B:RUNT.0000008386.07674.97.
- 9 [91] A. Osman, Y. Duan, and V. Kaftandjian, ‘Applied Artificial Intelligence in NDE’,
10 in *Handbook of Nondestructive Evaluation 4.0*, N. Meyendorf, N. Ida, R. Singh,
11 and J. Vrana, Eds., Cham: Springer International Publishing, 2021, pp. 1–35. doi:
12 10.1007/978-3-030-48200-8_49-1.
- 13 [92] S. Barut, V. Bissauge, G. Ithurralde, and W. Claassens, ‘Computer-aided analysis
14 of ultrasound data to speed-up the release of aerospace CFRP components’, *e-*
15 *Journal of Nondestructive Testing*, vol. 17, no. 07, Jul. 2012, [Online]. Available:
16 <https://www.ndt.net/search/docs.php3?id=12429&msgID=0&rootID=0>
- 17 [93] J. Aldrin, C. Coughlin, D. Forsyth, and J. Welter, ‘Progress on the Development
18 of Automated Data Analysis Algorithms and Software for Ultrasonic Inspection
19 of Composites’, *AIP Conference Proceedings*, vol. 1581, Jan. 2014, doi:
20 10.1063/1.4865058.
- 21 [94] D. Guo, G. Jiang, X. Lin, and Y. Wu, ‘Automated ultrasonic testing for 3D laser-
22 rapid prototyping blisk blades’, in *2016 7th International Conference on*
23 *Mechanical and Aerospace Engineering (ICMAE)*, Jul. 2016, pp. 214–218. doi:
24 10.1109/ICMAE.2016.7549537.
- 25 [95] S. J. Song and L. W. Schmerr, ‘Ultrasonic flaw classification in weldments using
26 probabilistic neural networks’, *J Nondestruct Eval*, vol. 11, no. 2, pp. 69–77, Jun.
27 1992, doi: 10.1007/BF00568290.
- 28 [96] S.-J. Song, H.-J. Kim, and H. Cho, ‘Development of an intelligent system for
29 ultrasonic flaw classification in weldments’, *Nuclear Engineering and Design*,
30 vol. 212, no. 1, pp. 307–320, Mar. 2002, doi: 10.1016/S0029-5493(01)00495-2.
- 31 [97] N. Munir, H.-J. Kim, J. Park, S.-J. Song, and S.-S. Kang, ‘Convolutional neural
32 network for ultrasonic weldment flaw classification in noisy conditions’,

- 1 *Ultrasonics*, vol. 94, pp. 74–81, Apr. 2019, doi: 10.1016/j.ultras.2018.12.001.
- 2 [98] P. Zacharis, G. West, G. Dobie, T. Lardner, and A. Gachagan, ‘Data-Driven
3 Analysis of Ultrasonic Inspection Data of Pressure Tubes’, *Nuclear Technology*,
4 vol. 202, no. 2–3, pp. 153–160, Jun. 2018, doi:
5 10.1080/00295450.2017.1421803.
- 6 [99] R. J. Pyle, R. L. T. Bevan, R. R. Hughes, A. A. S. Ali, and P. D. Wilcox, ‘Domain
7 Adapted Deep-Learning for Improved Ultrasonic Crack Characterization Using
8 Limited Experimental Data’, *IEEE Transactions on Ultrasonics, Ferroelectrics,*
9 *and Frequency Control*, vol. 69, no. 4, pp. 1485–1496, Apr. 2022, doi:
10 10.1109/TUFFC.2022.3151397.
- 11 [100] I. Virkkunen, T. Koskinen, O. Jessen-Juhler, and J. Rinta-aho, ‘Augmented
12 Ultrasonic Data for Machine Learning’, *J Nondestruct Eval*, vol. 40, no. 1, p. 4,
13 Jan. 2021, doi: 10.1007/s10921-020-00739-5.
- 14 [101] D. Medak, L. Posilovic, M. Subasic, M. Budimir, and S. Loncaric, ‘Automated
15 Defect Detection From Ultrasonic Images Using Deep Learning’, *IEEE Trans.*
16 *Ultrason., Ferroelect., Freq. Contr.*, vol. 68, no. 10, pp. 3126–3134, Oct. 2021,
17 doi: 10.1109/TUFFC.2021.3081750.
- 18 [102] F. Isensee, P. F. Jaeger, S. A. A. Kohl, J. Petersen, and K. H. Maier-Hein, ‘nnU-
19 Net: a self-configuring method for deep learning-based biomedical image
20 segmentation’, *Nat Methods*, vol. 18, no. 2, Art. no. 2, Feb. 2021, doi:
21 10.1038/s41592-020-01008-z.
- 22 [103] D. Medak, L. Posilović, M. Subašić, M. Budimir, and S. Lončarić, ‘DefectDet:
23 A deep learning architecture for detection of defects with extreme aspect ratios
24 in ultrasonic images’, *Neurocomputing*, vol. 473, pp. 107–115, Feb. 2022, doi:
25 10.1016/j.neucom.2021.12.008.
- 26 [104] J. Ye and N. Toyama, ‘Benchmarking Deep Learning Models for Automatic
27 Ultrasonic Imaging Inspection’, *IEEE Access*, vol. 9, pp. 36986–36994, 2021,
28 doi: 10.1109/ACCESS.2021.3062860.
- 29 [105] A. Rytter, ‘Vibrational Based Inspection of Civil Engineering Structures’,
30 1993.
- 31 [106] R. J. Pyle, R. L. T. Bevan, R. R. Hughes, R. K. Rachev, A. A. S. Ali, and P. D.
32 Wilcox, ‘Deep Learning for Ultrasonic Crack Characterization in NDE’, *IEEE*

- 1 *Trans. Ultrason., Ferroelect., Freq. Contr.*, vol. 68, no. 5, pp. 1854–1865, May
2 2021, doi: 10.1109/TUFFC.2020.3045847.
- 3 [107] A. Wirgin, ‘The inverse crime’, *arXiv:math-ph/0401050*, Jan. 2004, Accessed:
4 Jan. 06, 2022. [Online]. Available: <http://arxiv.org/abs/math-ph/0401050>
- 5 [108] S. Ali *et al.*, ‘Explainable Artificial Intelligence (XAI): What we know and
6 what is left to attain Trustworthy Artificial Intelligence’, *Information Fusion*, vol.
7 99, p. 101805, Nov. 2023, doi: 10.1016/j.inffus.2023.101805.
- 8 [109] S. A. and S. R., ‘A systematic review of Explainable Artificial Intelligence
9 models and applications: Recent developments and future trends’, *Decision*
10 *Analytics Journal*, vol. 7, p. 100230, Jun. 2023, doi:
11 10.1016/j.dajour.2023.100230.
- 12 [110] M. Abdar *et al.*, ‘A review of uncertainty quantification in deep learning:
13 Techniques, applications and challenges’, *Information Fusion*, vol. 76, pp. 243–
14 297, Dec. 2021, doi: 10.1016/j.inffus.2021.05.008.
- 15 [111] A. Young *et al.*, ‘Capturing Symbolic Expert Knowledge for the Development of
16 Industrial Fault Detection Systems: Manual and Automated Approaches’,
17 *International journal of COMADEM*.
- 18 [112] ‘EASA Artificial Intelligence Roadmap 1.0’, EASA. Accessed: Jan. 12, 2022.
19 [Online]. Available: [https://www.easa.europa.eu/document-library/general-](https://www.easa.europa.eu/document-library/general-publications/easa-artificial-intelligence-roadmap-10)
20 publications/easa-artificial-intelligence-roadmap-10
- 21 [113] ‘EASA releases its Concept Paper “First usable guidance for Level 1 machine
22 learning applications”’, EASA. Accessed: Jan. 12, 2022. [Online]. Available:
23 [https://www.easa.europa.eu/newsroom-and-events/news/easa-releases-its-](https://www.easa.europa.eu/newsroom-and-events/news/easa-releases-its-concept-paper-first-usable-guidance-level-1-machine-0)
24 concept-paper-first-usable-guidance-level-1-machine-0
- 25 [114] H. Huang, Q. Li, and D. Zhang, ‘Deep learning based image recognition for
26 crack and leakage defects of metro shield tunnel’, *Tunnelling and Underground*
27 *Space Technology*, vol. 77, pp. 166–176, Jul. 2018, doi:
28 10.1016/j.tust.2018.04.002.
- 29 [115] D. Weimer, B. Scholz-Reiter, and M. Shpitalni, ‘Design of deep convolutional
30 neural network architectures for automated feature extraction in industrial
31 inspection’, *CIRP Annals*, vol. 65, no. 1, pp. 417–420, Jan. 2016, doi:
32 10.1016/j.cirp.2016.04.072.

- 1 [116] O. Siljama, T. Koskinen, O. Jessen-Juhler, and I. Virkkunen, ‘Automated Flaw
2 Detection in Multi-channel Phased Array Ultrasonic Data Using Machine
3 Learning’, *J Nondestruct Eval*, vol. 40, no. 3, p. 67, Sep. 2021, doi:
4 10.1007/s10921-021-00796-4.
- 5 [117] I. Virkkunen, K. Miettinen, and T. Packalén, ‘Virtual flaws for NDE training
6 and qualification’, in *e-Journal of Nondestructive Testing*, Prague, Oct. 2014, p.
7 8. [Online]. Available: <https://www.ndt.net/?id=16779>
- 8 [118] O. Janssens, R. Van de Walle, M. Loccufier, and S. Van Hoecke, ‘Deep learning
9 for infrared thermal image based machine health monitoring’, *IEEE-ASME*
10 *TRANSACTIONS ON MECHATRONICS*, vol. 23, no. 1, Art. no. 1, 2018, doi:
11 10.1109/TMECH.2017.2722479.
- 12 [119] M. Frid-Adar, I. Diamant, E. Klang, M. Amitai, J. Goldberger, and H.
13 Greenspan, ‘GAN-based Synthetic Medical Image Augmentation for increased
14 CNN Performance in Liver Lesion Classification’, *Neurocomputing*, vol. 321, pp.
15 321–331, Dec. 2018, doi: 10.1016/j.neucom.2018.09.013.
- 16 [120] S. Sundaram and N. Hulkund, ‘GAN-based Data Augmentation for Chest X-
17 ray Classification’, Jul. 06, 2021, *arXiv*: arXiv:2107.02970. Accessed: Oct. 10,
18 2022. [Online]. Available: <http://arxiv.org/abs/2107.02970>
- 19 [121] M. M. Bejani and M. Ghatee, ‘A systematic review on overfitting control in
20 shallow and deep neural networks’, *Artif Intell Rev*, vol. 54, no. 8, pp. 6391–
21 6438, Dec. 2021, doi: 10.1007/s10462-021-09975-1.
- 22 [122] H. K. Dishar and L. A. Muhammed, ‘A Review of the Overfitting Problem in
23 Convolution Neural Network and Remedy Approaches’, *J. Al-Qadisiyah Comp.*
24 *Sci. Math.*, vol. 15, no. 2, Sep. 2023, doi: 10.29304/jqcm.2023.15.2.1240.
- 25 [123] ‘RollerFORM: Phased Array Wheel Probe’. Accessed: Jun. 06, 2023. [Online].
26 Available: <https://www.olympus-ims.com/en/rollerform/>
- 27 [124] ‘Ultrasonic Inspection Solutions for NDT | Peak NDT’. Accessed: Feb. 28,
28 2024. [Online]. Available: <https://www.peakndt.com/>
- 29 [125] M. Vasilev *et al.*, ‘Sensor-Enabled Multi-Robot System for Automated
30 Welding and In-Process Ultrasonic NDE’, *Sensors*, vol. 21, no. 15, Art. no. 15,
31 Jan. 2021, doi: 10.3390/s21155077.
- 32 [126] Y. Fu and X. Yao, ‘A review on manufacturing defects and their detection of

- 1 fiber reinforced resin matrix composites’, *Composites Part C: Open Access*, vol.
2 8, p. 100276, Jul. 2022, doi: 10.1016/j.jcomc.2022.100276.
- 3 [127] R. Draï, F. Sellidj, M. Khelil, and A. Benchaala, ‘Elaboration of some signal
4 processing algorithms in ultrasonic techniques: application to materials NDT’,
5 *Ultrasonics*, vol. 38, no. 1, pp. 503–507, Mar. 2000, doi: 10.1016/S0041-
6 624X(99)00082-7.
- 7 [128] A. V Oppenheim and R. W. Schafer, ‘Discrete-Time Signal Processing’.
8 Accessed: Feb. 29, 2024. [Online]. Available: [https://www.pearson.com/en-](https://www.pearson.com/en-us/subject-catalog/p/discrete-time-signal-processing/P200000003226/9780137549771)
9 [us/subject-catalog/p/discrete-time-signal-](https://www.pearson.com/en-us/subject-catalog/p/discrete-time-signal-processing/P200000003226/9780137549771)
10 [processing/P200000003226/9780137549771](https://www.pearson.com/en-us/subject-catalog/p/discrete-time-signal-processing/P200000003226/9780137549771)
- 11 [129] P. Angam, K. Vijayarekha, S. Abraham, and V. Balasubramaniam, ‘Fourier
12 Analysis of Ultrasonic TOFD Signals for Defect Detection in Austenitic Stainless
13 Steel Welds’, *International Journal of Computer Applications*, vol. 71, pp. 14–
14 17, Jun. 2013, doi: 10.5120/12385-8737.
- 15 [130] J. C. Lazaro, ‘Noise reduction in ultrasonic NDT using discrete wavelet
16 transform processing’, in *2002 IEEE Ultrasonics Symposium, 2002.*
17 *Proceedings.*, Oct. 2002, pp. 777–780 vol.1. doi:
18 10.1109/ULTSYM.2002.1193514.
- 19 [131] M. Meng, Y. J. Chua, E. Wouterson, and C. P. K. Ong, ‘Ultrasonic signal
20 classification and imaging system for composite materials via deep convolutional
21 neural networks’, *Neurocomputing*, vol. 257, pp. 128–135, Sep. 2017, doi:
22 10.1016/j.neucom.2016.11.066.
- 23 [132] B. Wang, Y. Li, Y. Luo, X. Li, and T. Freiheit, ‘Early Event Detection in a Deep-
24 learning Driven Quality Prediction Model for Ultrasonic Welding’, *Journal of*
25 *Manufacturing Systems*, vol. 2021, pp. 325–336, Jun. 2021, doi:
26 10.1016/j.jmsy.2021.06.009.
- 27 [133] S. Lonné, L. D. Roumilly, L. L. Ber, S. Mahaut, and G. Cattiaux,
28 ‘EXPERIMENTAL VALIDATION OF CIVA ULTRASONIC SIMULATIONS’,
29 2006, [Online]. Available:
30 [https://www.semanticscholar.org/paper/EXPERIMENTAL-VALIDATION-OF-](https://www.semanticscholar.org/paper/EXPERIMENTAL-VALIDATION-OF-CIVA-ULTRASONIC-Lonn%C3%A9-Roumilly/16b85af3b6a96d4657c9902ca8652fbbf93cbf2e)
31 [CIVA-ULTRASONIC-Lonn%C3%A9-](https://www.semanticscholar.org/paper/EXPERIMENTAL-VALIDATION-OF-CIVA-ULTRASONIC-Lonn%C3%A9-Roumilly/16b85af3b6a96d4657c9902ca8652fbbf93cbf2e)
32 [Roumilly/16b85af3b6a96d4657c9902ca8652fbbf93cbf2e](https://www.semanticscholar.org/paper/EXPERIMENTAL-VALIDATION-OF-CIVA-ULTRASONIC-Lonn%C3%A9-Roumilly/16b85af3b6a96d4657c9902ca8652fbbf93cbf2e)

- 1 [134] M. Darmon *et al.*, ‘VALIDATION OF AN ULTRASONIC
2 CHARACTERIZATION TECHNIQUE FOR ANISOTROPIC MATERIALS:
3 COMPARISON OF EXPERIMENTS WITH BEAM PROPAGATION
4 MODELLING’, presented at the 2019 INTERNATIONAL CONGRESS ON
5 ULTRASONICS, Bruges Belgium.
- 6 [135] K. Jezzine, D. Segur, R. Ecault, and N. Dominguez, ‘Simulation of ultrasonic
7 inspections of composite structures in the CIVA software platform’, *e-Journal of*
8 *Nondestructive Testing*, vol. 21, no. 07, Jul. 2016, [Online]. Available:
9 <https://www.ndt.net/search/docs.php3?id=19438&msgID=0&rootID=0>
- 10 [136] A. Figueira and B. Vaz, ‘Survey on Synthetic Data Generation, Evaluation
11 Methods and GANs’, *Mathematics*, vol. 10, no. 15, Art. no. 15, Jan. 2022, doi:
12 10.3390/math10152733.
- 13 [137] I. J. Goodfellow *et al.*, ‘Generative Adversarial Networks’, Jun. 10, 2014,
14 *arXiv*: arXiv:1406.2661. Accessed: Aug. 14, 2023. [Online]. Available:
15 <http://arxiv.org/abs/1406.2661>
- 16 [138] A. Antoniou, A. Storkey, and H. Edwards, ‘Data Augmentation Generative
17 Adversarial Networks’, *arXiv:1711.04340 [cs, stat]*, Mar. 2018, Accessed: Feb.
18 03, 2022. [Online]. Available: <http://arxiv.org/abs/1711.04340>
- 19 [139] S. Motamed, P. Rogalla, and F. Khalvati, ‘Data augmentation using Generative
20 Adversarial Networks (GANs) for GAN-based detection of Pneumonia and
21 COVID-19 in chest X-ray images’, *Informatics in Medicine Unlocked*, vol. 27,
22 p. 100779, Jan. 2021, doi: 10.1016/j.imu.2021.100779.
- 23 [140] N. Park, M. Mohammadi, K. Gorde, S. Jajodia, H. Park, and Y. Kim, ‘Data
24 Synthesis based on Generative Adversarial Networks’, *Proc. VLDB Endow.*, vol.
25 11, no. 10, pp. 1071–1083, Jun. 2018, doi: 10.14778/3231751.3231757.
- 26 [141] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, ‘Unpaired Image-to-Image
27 Translation using Cycle-Consistent Adversarial Networks’, Aug. 24, 2020, *arXiv*:
28 arXiv:1703.10593. Accessed: May 24, 2022. [Online]. Available:
29 <http://arxiv.org/abs/1703.10593>
- 30 [142] ‘EXTENDE, Experts in Non Destructive Testing Simulation with CIVA
31 Software’. Accessed: Nov. 07, 2022. [Online]. Available:
32 <https://www.extende.com/>

- 1 [143] Extende, ‘Civa 2020 User Manual’. Nov. 22, 2019.
- 2 [144] V. Nair and G. E. Hinton, ‘Rectified Linear Units Improve Restricted
3 Boltzmann Machines’, p. 8.
- 4 [145] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, ‘Regularized Evolution for
5 Image Classifier Architecture Search’, Feb. 16, 2019, *arXiv*: arXiv:1802.01548.
6 doi: 10.48550/arXiv.1802.01548.
- 7 [146] ‘PyTorch’. Accessed: Jul. 26, 2023. [Online]. Available:
8 <https://www.pytorch.org>
- 9 [147] C. D. Walsh, J. Edwards, and R. H. Insall, ‘Ensuring accurate stain
10 reproduction in deep generative networks for virtual immunohistochemistry’,
11 Apr. 14, 2022, *arXiv*: arXiv:2204.06849. Accessed: Aug. 25, 2022. [Online].
12 Available: <http://arxiv.org/abs/2204.06849>
- 13 [148] ‘scipy.stats.invgauss — SciPy v1.11.1 Manual’. Accessed: Aug. 14, 2023.
14 [Online]. Available:
15 <https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.invgauss.html>
- 16 [149] ‘scipy.signal.savgol_filter — SciPy v1.11.1 Manual’. Accessed: Aug. 14, 2023.
17 [Online]. Available:
18 [https://docs.scipy.org/doc/scipy/reference/generated/scipy.signal.savgol_filter.ht](https://docs.scipy.org/doc/scipy/reference/generated/scipy.signal.savgol_filter.html)
19 [ml](https://docs.scipy.org/doc/scipy/reference/generated/scipy.signal.savgol_filter.html)
- 20 [150] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra,
21 ‘Grad-CAM: Visual Explanations from Deep Networks via Gradient-based
22 Localization’, *Int J Comput Vis*, vol. 128, no. 2, pp. 336–359, Feb. 2020, doi:
23 10.1007/s11263-019-01228-7.
- 24 [151] T. B. Hudson, P. J. Follis, J. J. Pinakidis, T. Sreekantamurthy, and F. L.
25 Palmieri, ‘Porosity detection and localization during composite cure inside an
26 autoclave using ultrasonic inspection’, *Composites Part A: Applied Science and*
27 *Manufacturing*, vol. 147, p. 106337, Aug. 2021, doi:
28 10.1016/j.compositesa.2021.106337.
- 29 [152] Institute of Fundamentals of Machinery Design, Silesian University of
30 Technology, Gliwice, 44-100, Poland, A. Wronkowicz, A. Katunin, and K.
31 Dragan, ‘Ultrasonic C-Scan Image Processing Using Multilevel Thresholding for
32 Damage Evaluation in Aircraft Vertical Stabilizer’, *IJIGSP*, vol. 7, no. 11, pp. 1–

- 1 8, Oct. 2015, doi: 10.5815/ijigsp.2015.11.01.
- 2 [153] Z. Zhang, M. Liu, Q. Li, and Y. Ang, ‘Visualized characterization of diversified
3 defects in thick aerospace composites using ultrasonic B-scan’, *Composites*
4 *Communications*, vol. 22, p. 100435, Dec. 2020, doi:
5 10.1016/j.coco.2020.100435.
- 6 [154] Y. Zhou *et al.*, ‘Multi-task learning for segmentation and classification of
7 tumors in 3D automated breast ultrasound images’, *Medical Image Analysis*, vol.
8 70, p. 101918, May 2021, doi: 10.1016/j.media.2020.101918.
- 9 [155] Y. Liu, ‘3D Image Segmentation of MRI Prostate Based on a Pytorch
10 Implementation of V-Net’, *J. Phys.: Conf. Ser.*, vol. 1549, no. 4, p. 042074, Jun.
11 2020, doi: 10.1088/1742-6596/1549/4/042074.
- 12 [156] M. Leong, D. Prasad, Y. T. Lee, and F. Lin, ‘Semi-CNN Architecture for
13 Effective Spatio-Temporal Learning in Action Recognition’, *Applied Sciences*,
14 vol. 10, p. 557, Jan. 2020, doi: 10.3390/app10020557.
- 15 [157] A. M. Badshah *et al.*, ‘Deep features-based speech emotion recognition for
16 smart affective services’, *Multimed Tools Appl*, vol. 78, no. 5, pp. 5571–5589,
17 Mar. 2019, doi: 10.1007/s11042-017-5292-7.
- 18 [158] H. Liu, K. Simonyan, and Y. Yang, ‘DARTS: Differentiable Architecture
19 Search’, Apr. 23, 2019, *arXiv*: arXiv:1806.09055. doi:
20 10.48550/arXiv.1806.09055.
- 21 [159] J. Mellor, J. Turner, A. Storkey, and E. J. Crowley, ‘Neural Architecture Search
22 without Training’, Jun. 11, 2021, *arXiv*: arXiv:2006.04647. doi:
23 10.48550/arXiv.2006.04647.
- 24 [160] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, ‘Algorithms for Hyper-
25 Parameter Optimization’, in *Advances in Neural Information Processing*
26 *Systems*, Curran Associates, Inc., 2011. Accessed: Jun. 16, 2023. [Online].
27 Available:
28 [https://papers.nips.cc/paper_files/paper/2011/hash/86e8f7ab32cfd12577bc2619](https://papers.nips.cc/paper_files/paper/2011/hash/86e8f7ab32cfd12577bc2619bc635690-Abstract.html)
29 [bc635690-Abstract.html](https://papers.nips.cc/paper_files/paper/2011/hash/86e8f7ab32cfd12577bc2619bc635690-Abstract.html)
- 30 [161] S. Schrodi, D. Stoll, B. Ru, R. Sukthanker, T. Brox, and F. Hutter, ‘Construction
31 of Hierarchical Neural Architecture Search Spaces based on Context-free
32 Grammars’, Jun. 05, 2023, *arXiv*: arXiv:2211.01842. Accessed: Jun. 16, 2023.

- 1 [Online]. Available: <http://arxiv.org/abs/2211.01842>
- 2 [162] D. Maturana and S. Scherer, 'VoxNet: A 3D Convolutional Neural Network for
3 real-time object recognition', in *2015 IEEE/RSJ International Conference on*
4 *Intelligent Robots and Systems (IROS)*, Sep. 2015, pp. 922–928. doi:
5 10.1109/IROS.2015.7353481.
- 6 [163] Z. Liu *et al.*, 'VB-Net: Voxel-Based Broad Learning Network for 3D Object
7 Classification', *Applied Sciences*, vol. 10, no. 19, Art. no. 19, Jan. 2020, doi:
8 10.3390/app10196735.
- 9 [164] S. McKnight *et al.*, 'GANs and alternative methods of synthetic noise
10 generation for domain adaption of defect classification of Non-destructive
11 ultrasonic testing', Jun. 02, 2023, *arXiv*: arXiv:2306.01469. doi:
12 10.48550/arXiv.2306.01469.
- 13 [165] M. Gower, G. Sims, R. Lee, S. Frost, M. Stone, and M. Wall, 'Measurement
14 Good Practice Guide', no. 78.
- 15 [166] Y. Wang, F. Tao, Y. Zuo, M. Zhang, and Q. Qi, 'Digital-Twin-Enhanced Quality
16 Prediction for the Composite Materials', *Engineering*, vol. 22, pp. 23–33, Mar.
17 2023, doi: 10.1016/j.eng.2022.08.019.
- 18 [167] D. G. Puttaraju and H. G. Hanumantharaju, 'Finite element analysis and
19 validation of tensile properties of carbon fiber reinforced polymer matrix
20 composites', *Materials Today: Proceedings*, vol. 62, pp. 2800–2807, Jan. 2022,
21 doi: 10.1016/j.matpr.2022.02.188.
- 22 [168] A. Hauffe, F. Hähnel, and K. Wolf, 'Comparison of algorithms to quantify the
23 damaged area in CFRP ultrasonic scans', *Composite Structures*, vol. 235, p.
24 111791, Mar. 2020, doi: 10.1016/j.compstruct.2019.111791.
- 25 [169] S. Barut and N. Dominguez, 'NDT Diagnosis Automation: a Key to Efficient
26 Production in the Aeronautic Industry', *e-Journal of Nondestructive Testing*, vol.
27 21, no. 07, Jul. 2016, [Online]. Available:
28 <https://www.ndt.net/search/docs.php3?id=19184&msgID=0&rootID=0>
- 29 [170] S. Kumaran and S. Rani, 'Application of 6db Drop Technique to Estimate the
30 Width of Sub Assembly Ring Top Using Pulse Echo Ultrasonic Technique',
31 *International Journal of Engineering and Technology*, vol. 5, pp. 4771–4775,
32 Jan. 2013.

- 1 [171] P. Ciorau, ‘Comparison Between -6 DB and -12 DB Amplitude Drop
2 Techniques for Length Sizing’.
- 3 [172] X. Li, Y. Wang, P. Ni, H. Hu, and Y. Song, ‘Flaw sizing using ultrasonic C-scan
4 imaging with dynamic thresholds’, *insight*, vol. 59, no. 11, pp. 603–608, Nov.
5 2017, doi: 10.1784/insi.2017.59.11.603.
- 6 [173] X. Cheng, G. Ma, Z. Wu, H. Zu, and X. Hu, ‘Automatic defect depth estimation
7 for ultrasonic testing in carbon fiber reinforced composites using deep learning’,
8 *NDT & E International*, vol. 135, p. 102804, Apr. 2023, doi:
9 10.1016/j.ndteint.2023.102804.
- 10 [174] P. M. Cheng and H. S. Malhi, ‘Transfer Learning with Convolutional Neural
11 Networks for Classification of Abdominal Ultrasound Images’, *J Digit Imaging*,
12 vol. 30, no. 2, pp. 234–243, Apr. 2017, doi: 10.1007/s10278-016-9929-2.
- 13 [175] O. Ronneberger, P. Fischer, and T. Brox, ‘U-Net: Convolutional Networks for
14 Biomedical Image Segmentation’, in *Medical Image Computing and Computer-
15 Assisted Intervention – MICCAI 2015*, N. Navab, J. Hornegger, W. M. Wells, and
16 A. F. Frangi, Eds., in Lecture Notes in Computer Science. Cham: Springer
17 International Publishing, 2015, pp. 234–241. doi: 10.1007/978-3-319-24574-
18 4_28.
- 19 [176] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, ‘3D
20 U-Net: Learning Dense Volumetric Segmentation from Sparse Annotation’, Jun.
21 21, 2016, *arXiv*: arXiv:1606.06650. doi: 10.48550/arXiv.1606.06650.
- 22 [177] M. Sharifzadeh, H. Benali, and H. Rivaz, ‘Phase Aberration Correction: A
23 Convolutional Neural Network Approach’, *IEEE Access*, vol. 8, pp. 162252–
24 162260, 2020, doi: 10.1109/ACCESS.2020.3021685.
- 25 [178] K. Ono and A. Gallego, ‘Attenuation of lamb waves in CFRP plates’, *Journal
26 of Acoustic Emission*, vol. 30, pp. 109–124, Jan. 2012.
- 27 [179] J. Yang, K. Zhou, and Z. Liu, ‘Full-Spectrum Out-of-Distribution Detection’,
28 *Int J Comput Vis*, vol. 131, no. 10, pp. 2607–2622, Oct. 2023, doi:
29 10.1007/s11263-023-01811-z.
- 30 [180] S. McKnight *et al.*, ‘A comparison of methods for generating synthetic training
31 data for domain adaption of deep learning models in ultrasonic non-destructive
32 evaluation’, *NDT & E International*, vol. 141, p. 102978, Jan. 2024, doi:

- 1 10.1016/j.ndteint.2023.102978.
- 2 [181] V. Rani, S. T. Nabi, M. Kumar, A. Mittal, and K. Kumar, ‘Self-supervised
3 Learning: A Succinct Review’, *Arch Computat Methods Eng*, vol. 30, no. 4, pp.
4 2761–2775, May 2023, doi: 10.1007/s11831-023-09884-2.
- 5 [182] Z. Zhao, L. Alzubaidi, J. Zhang, Y. Duan, and Y. Gu, ‘A comparison review of
6 transfer learning and self-supervised learning: Definitions, applications,
7 advantages and limitations’, *Expert Systems with Applications*, vol. 242, p.
8 122807, May 2024, doi: 10.1016/j.eswa.2023.122807.
- 9 [183] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, ‘BERT: Pre-training of
10 Deep Bidirectional Transformers for Language Understanding’, May 24, 2019,
11 *arXiv*: arXiv:1810.04805. doi: 10.48550/arXiv.1810.04805.
- 12 [184] P. P. Ray, ‘ChatGPT: A comprehensive review on background, applications,
13 key challenges, bias, ethics, limitations and future scope’, *Internet of Things and
14 Cyber-Physical Systems*, vol. 3, pp. 121–154, Jan. 2023, doi:
15 10.1016/j.iotcps.2023.04.003.
- 16 [185] M. Oquab *et al.*, ‘DINOv2: Learning Robust Visual Features without
17 Supervision’, Feb. 02, 2024, *arXiv*: arXiv:2304.07193. doi:
18 10.48550/arXiv.2304.07193.
- 19 [186] ‘CLIP: Connecting text and images’. Accessed: Apr. 03, 2024. [Online].
20 Available: <https://openai.com/research/clip>
- 21 [187] R. He, A. Ravula, B. Kanagal, and J. Ainslie, ‘RealFormer: Transformer Likes
22 Residual Attention’, Sep. 10, 2021, *arXiv*: arXiv:2012.11747. doi:
23 10.48550/arXiv.2012.11747.
- 24 [188] S. Gidaris, P. Singh, and N. Komodakis, ‘Unsupervised Representation
25 Learning by Predicting Image Rotations’, Mar. 20, 2018, *arXiv*:
26 arXiv:1803.07728. doi: 10.48550/arXiv.1803.07728.
- 27 [189] M. Noroozi and P. Favaro, ‘Unsupervised Learning of Visual Representations
28 by Solving Jigsaw Puzzles’, Aug. 22, 2017, *arXiv*: arXiv:1603.09246. doi:
29 10.48550/arXiv.1603.09246.
- 30 [190] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, ‘Context
31 Encoders: Feature Learning by Inpainting’, Nov. 21, 2016, *arXiv*:
32 arXiv:1604.07379. doi: 10.48550/arXiv.1604.07379.

- 1 [191] R. Zhang, P. Isola, and A. A. Efros, ‘Colorful Image Colorization’, Oct. 05,
2 2016, *arXiv*: arXiv:1603.08511. doi: 10.48550/arXiv.1603.08511.
- 3 [192] G. Koch, R. Zemel, and R. Salakhutdinov, ‘Siamese Neural Networks for One-
4 shot Image Recognition’.
- 5 [193] S. McKnight *et al.*, ‘Advancing Carbon Fiber Composite Inspection: Deep
6 Learning-Enabled Defect Localization and Sizing via 3-Dimensional U-Net
7 Segmentation of Ultrasonic Data’, *IEEE Transactions on Ultrasonics,*
8 *Ferroelectrics, and Frequency Control*, pp. 1–1, 2024, doi:
9 10.1109/TUFFC.2024.3408314.
- 10 [194] B. Lim and S. Zohren, ‘Time-series forecasting with deep learning: a survey’,
11 *Philosophical Transactions of the Royal Society A: Mathematical, Physical and*
12 *Engineering Sciences*, vol. 379, no. 2194, p. 20200209, Feb. 2021, doi:
13 10.1098/rsta.2020.0209.
- 14 [195] H. I. Fawaz *et al.*, ‘InceptionTime: Finding AlexNet for Time Series
15 Classification’, *Data Min Knowl Disc*, vol. 34, no. 6, pp. 1936–1962, Nov. 2020,
16 doi: 10.1007/s10618-020-00710-y.
- 17 [196] J. P. E. Schouten, C. Matek, L. F. P. Jacobs, M. C. Buck, D. Bošnački, and C.
18 Marr, ‘Tens of images can suffice to train neural networks for malignant
19 leukocyte detection’, *Sci Rep*, vol. 11, no. 1, p. 7995, Apr. 2021, doi:
20 10.1038/s41598-021-86995-5.
- 21 [197] ‘scikit-image: Image processing in Python — scikit-image’. Accessed: Apr. 06,
22 2024. [Online]. Available: <https://scikit-image.org/>
- 23 [198] A. Mehmood, M. Maqsood, M. Bashir, and Y. Shuyuan, ‘A Deep Siamese
24 Convolution Neural Network for Multi-Class Classification of Alzheimer
25 Disease’, *Brain Sci*, vol. 10, no. 2, p. 84, Feb. 2020, doi:
26 10.3390/brainsci10020084.
- 27 [199] ‘Aerospace Event 2024’. Accessed: Jun. 04, 2024. [Online]. Available:
28 [https://www.bindt.org/events-and-awards/PastEventsandWebinars/aerospace-](https://www.bindt.org/events-and-awards/PastEventsandWebinars/aerospace-event-2024/?cookie-accept=1)
29 [event-2024/?cookie-accept=1](https://www.bindt.org/events-and-awards/PastEventsandWebinars/aerospace-event-2024/?cookie-accept=1)
- 30
31

1 Appendix

2 Table 28: Detection accuracy across thresholds and processing steps for each sample.

Threshold	Defective Sample	Detection Accuracy % (False positives)			
		Forward Sweep	Backward Sweep	Combined Sweep	Area Threshold
0.9999999	1	16.85 (74)	12.40 (106)	39.47 (23)	100.00 (0)
	2	14.79 (144)	21.55 (91)	55.56 (20)	100.00 (0)
	3	7.98 (173)	10.49 (128)	22.73 (51)	100.00 (0)
0.999999	1	9.20 (148)	7.54 (184)	28.85 (37)	93.75 (1)
	2	8.28 (277)	13.30 (163)	42.37 (34)	100.00 (0)
	3	5.34 (266)	6.91 (202)	16.67 (75)	93.75 (1)
0.99999	1	4.66 (307)	3.83 (377)	12.83 (102)	83.33 (3)
	2	4.28 (559)	5.94 (396)	20.49 (97)	96.15 (1)
	3	3.42 (395)	3.99 (337)	10.79 (124)	88.24 (2)
0.9999	1	2.17 (677)	2.02 (729)	5.68 (249)	71.43 (6)
	2	2.33 (1046)	2.76 (880)	7.99 (288)	92.59 (2)
	3	2.33 (586)	2.34 (584)	5.88 (240)	78.95 (4)
0.999	1	1.31 (1129)	1.36 (1085)	1.89 (780)	65.22 (8)
	2	1.33 (1704)	1.39 (1780)	2.32 (1053)	89.29 (3)
	3	1.68 (759)	1.71 (806)	2.85 (512)	84.62 (2)
0.99	1	4.95 (288)	4.79 (298)	1.51 (980)	50.00 (15)
	2	5.73 (411)	3.81 (632)	1.30 (1896)	96.15 (1)
	3	4.34 (331)	3.60 (402)	1.62 (913)	60.00 (10)

3