



End-to-End NIR-DOT Reconstruction with Hybrid CNN-Transformer

MSc Dissertation

Max Andrew Hart

School of Computer Science

College of Engineering and Physical Sciences

University of Birmingham

2024-25

This dissertation presents a novel hybrid CNN-Transformer architecture for enhanced Near-Infrared Diffuse Optical Tomography (NIR-DOT) reconstruction. NIR-DOT is a non-invasive imaging technique with significant potential for brain imaging and cancer detection, but faces challenges in reconstruction accuracy due to the ill-posed nature of the inverse problem.

We introduce a two-stage approach combining the spatial feature learning capabilities of CNNs with the contextual processing power of Transformers, further enhanced by the integration of local tissue context information. Using a comprehensive dataset of 5,000 synthetic phantoms, we demonstrate that our approach significantly improves reconstruction quality over traditional methods and standard deep learning approaches.

Our key innovation—tissue context integration—provides physiologically relevant priors that guide the reconstruction process. Experimental results show improvements in RMSE, SSIM, and PSNR metrics, with particularly notable enhancements in tumor region reconstruction. The proposed method represents a significant advance in NIR-DOT reconstruction technology with potential clinical applications in brain imaging and breast cancer detection.

Acknowledgements

Abbreviations

ACB

Apple Banana Carrot

Abstract	ii
Acknowledgements	iii
Abbreviations	iv
List of Figures	viii
List of Tables	ix
1 Introduction	1
1.1 Background and Motivation	1
1.2 Near-Infrared Diffuse Optical Tomography	2
1.3 Challenges in NIR-DOT Reconstruction	2
1.4 Research Objectives and Contributions	3
1.5 Dissertation Structure	3
2 Literature Review	5
2.1 Physics of NIR Light Propagation	5
2.2 Classical NIR-DOT Reconstruction Methods	5
2.3 Deep Learning for Medical Image Reconstruction	5
2.4 CNN Approaches in NIR-DOT	5
2.5 Transformers in Medical Imaging	5
2.6 Hybrid Architectures and Multi-Stage Training	5
2.7 Contextual Integration in Medical Imaging	5
2.8 Research Gap and Opportunity	5
3 Synthetic Phantom Data Generation	6
3.1 Physics-Based Forward Modeling	6
3.2 Geometric Phantom Construction	6
3.3 Optical Property Assignment	6

3.4	Surface Extraction and Probe Placement	6
3.5	Frequency-Domain Measurement Simulation	6
3.6	Dataset Composition and Statistics	6
3.7	Validation of Simulated Data	6
4	Hybrid CNN-Transformer Architecture	7
4.1	Architectural Overview	7
4.2	CNN Autoencoder Design	7
4.2.1	Encoder Architecture	7
4.2.2	Decoder Architecture	7
4.2.3	Residual Connections and Feature Extraction	7
4.3	NIR Measurement Processing	7
4.3.1	Spatial Awareness in Measurement Processing	7
4.3.2	Tissue Patch Extraction and Encoding	7
4.4	Transformer Encoder Design	7
4.4.1	Self-Attention Mechanism	7
4.4.2	Token-Type Embeddings	7
4.4.3	Feature Enhancement	7
4.5	Integration Strategy	7
4.5.1	Information Flow Between Components	7
4.5.2	Tissue Context Integration	7
5	Two-Stage Training Methodology	8
5.1	Training Strategy Overview	8
5.2	Stage 1: CNN Autoencoder Pre-Training	8
5.2.1	Identity Mapping Objective	8
5.2.2	AdamW with OneCycleLR Optimization	8
5.2.3	Hyperparameter Selection	8
5.3	Stage 2: Transformer Enhancement Training	8
5.3.1	Frozen Decoder Approach	8
5.3.2	Linear Warmup with Cosine Decay	8
5.3.3	Differential Weight Decay	8
5.4	Data Augmentation Strategy	8
5.5	Implementation Details	8
5.5.1	Hardware Optimization	8
5.5.2	Experiment Tracking	8
6	Experimental Results and Analysis	9
6.1	Experimental Setup	9
6.1.1	Dataset Preparation	9
6.1.2	Evaluation Metrics	9
6.1.3	Baseline Methods	9
6.2	Stage 1 Results: CNN Autoencoder Performance	9

6.2.1	Reconstruction Quality Metrics	9
6.2.2	Feature Learning Analysis	9
6.3	Stage 2 Results: Transformer Enhancement	9
6.3.1	Baseline Mode Performance	10
6.3.2	Enhanced Mode with Tissue Context	10
6.3.3	Comparative Analysis	10
6.4	Ablation Studies	10
6.4.1	Impact of Tissue Patch Size	11
6.4.2	Effect of Transformer Layers and Heads	11
6.4.3	Sensitivity to Tissue Context Quality	11
6.5	Visualization and Interpretation	11
6.5.1	Attention Map Analysis	11
6.5.2	Reconstruction Quality Visualization	11
6.5.3	Error Distribution Analysis	11
7	Discussion	12
7.1	Key Findings and Insights	12
7.2	Clinical Implications	12
7.3	Limitations of Current Approach	12
7.4	Computational Efficiency Considerations	12
8	Conclusion and Future Work	13
8.1	Summary of Contributions	13
8.2	Future Research Directions	13
8.3	Potential Clinical Applications	13
A	Implementation Details	15
B	Additional Experimental Results	16
C	Mathematical Derivations	17

List of Figures

List of Tables

1.1 Background and Motivation

Diffuse Optical Tomography (DOT) is a non-invasive imaging modality that uses near-infrared (NIR) light to probe tissue optical properties, particularly absorption (μ_a) and reduced scattering coefficients (μ'_s). By measuring how NIR photons migrate through tissue, DOT can provide valuable structural and functional information about biological tissue, including haemodynamic changes and tumour localisation [?, ?]. In recent years, Near-Infrared Diffuse Optical Tomography (NIR-DOT) has gained prominence in biomedical imaging due to its non-ionising nature, low cost, and ability to provide unique functional contrast compared to conventional modalities such as MRI or CT.

However, the potential of NIR-DOT has not yet been fully realised in clinical applications. One of the central challenges is the reconstruction of three-dimensional tissue optical property maps from boundary NIR measurements. This inverse problem is severely ill-posed: many different internal tissue configurations can give rise to similar surface measurements. The result is that reconstructions are often noisy, blurred, and highly sensitive to measurement error or modelling assumptions. Overcoming these challenges requires methodological innovation in both data-driven learning and physics-aware modelling.

Recent advances in deep learning, particularly convolutional neural networks (CNNs) and transformers, have revolutionised inverse imaging tasks across modalities. Inspired by these developments, our research leverages a hybrid CNN–transformer architecture for NIR-DOT reconstruction, building upon and extending the work of Robin Dale and colleagues [?]. Dale’s work demonstrated the feasibility of combining convolutional feature extraction with transformer-based long-range dependency modelling to improve optical imaging reconstruction. Our research pushes this line further by developing a two-stage, student–teacher training paradigm that integrates synthetic phantom generation, robust standardisation, and physics-aware evaluation to achieve state-of-the-art reconstructions in raw optical units.

1.2 Near-Infrared Diffuse Optical Tomography

NIR-DOT operates on the principle of photon migration in highly scattering media. When NIR light in the range of 650–950 nm is introduced into tissue, photons undergo multiple scattering events and absorption before being detected at surface detectors. By modelling this transport, typically with the diffusion approximation to the radiative transfer equation, it is possible to infer the internal distribution of absorption and scattering coefficients.

The forward problem in NIR-DOT consists of solving the diffusion equation given known tissue optical properties and boundary conditions, producing boundary fluxes or detector measurements. The inverse problem, however, is the reconstruction of tissue properties from these measurements. This inversion is complicated by the exponential attenuation of light, the limited number of source–detector pairs, and the inherent ambiguity of mapping low-dimensional measurements back to high-dimensional volumes. Consequently, the problem is both underdetermined and ill-conditioned.

Traditional reconstruction approaches rely on iterative optimisation with regularisation, often requiring prior assumptions about smoothness or sparsity. These methods are computationally expensive and can yield unstable solutions. Data-driven methods, by contrast, can learn mappings directly from simulated or empirical data, enabling faster and more robust reconstructions. In particular, CNNs are adept at learning spatial hierarchies from volumetric data, while transformers excel at capturing long-range dependencies across measurements. Combining these paradigms offers a powerful new direction for NIR-DOT.

1.3 Challenges in NIR-DOT Reconstruction

The central challenge in NIR-DOT is the ill-posed nature of the inverse problem. Specifically:

- **Non-uniqueness:** Different tissue structures can yield nearly identical boundary measurements.
- **Instability:** Small perturbations in measurements can lead to disproportionately large errors in reconstructions.
- **High dimensionality:** Reconstructing 64^3 voxel volumes from relatively few measurements (typically 1000 per phantom) represents an extreme underdetermined problem.
- **Noise sensitivity:** Instrumental noise, modelling errors, and physiological variability further degrade reconstruction accuracy.

From a computational standpoint, another challenge lies in balancing fidelity to physical principles with the flexibility of data-driven learning. Purely physics-based models are slow and fragile, while purely data-driven models risk producing reconstructions that violate known physics. Hybrid approaches must therefore enforce physical consistency while still benefiting from the representational power of deep networks.

1.4 Research Objectives and Contributions

The goal of this dissertation is to design, implement, and evaluate a hybrid CNN–transformer pipeline for NIR-DOT reconstruction that addresses the above challenges through a principled two-stage training framework. The specific objectives are:

1. **Data Simulation:** Develop a high-fidelity synthetic phantom generation and forward-modelling pipeline to create realistic training datasets with controllable optical properties.
2. **Stage 1 Training:** Pre-train a 3D CNN autoencoder on ground truth volumes to learn robust spatial representations of tissue optical properties.
3. **Stage 2 Training:** Introduce a transformer encoder to process NIR measurements, integrating spatial embeddings and sequence modelling to produce latent codes compatible with the Stage 1 decoder.
4. **Student–Teacher Learning:** Employ a teacher–student paradigm where the frozen CNN encoder provides latent targets to guide the transformer-based student network.
5. **Physics-Aware Evaluation:** Evaluate reconstructions in raw physical units (mm^{-1}) with metrics such as Dice coefficient, RMSE, and contrast ratio, ensuring interpretability and clinical relevance.
6. **End-to-End System Integration:** Build a modular, extensible codebase with clean separation of data generation, training, and evaluation, enabling reproducibility and future extensions.

The main contributions of this work are as follows:

- A novel hybrid CNN–transformer architecture for NIR-DOT reconstruction.
- A two-stage training pipeline combining autoencoder pre-training with transformer-based measurement encoding.
- Integration of synthetic phantom simulation with physics-consistent data preprocessing and standardisation.
- A teacher–student learning paradigm for latent-space alignment across stages.
- A comprehensive evaluation framework in raw optical units, ensuring clinically meaningful assessment of reconstructions.

1.5 Dissertation Structure

The remainder of this dissertation is organised as follows:

- **Chapter 2: Literature Review** surveys prior work in NIR-DOT reconstruction, deep learning for inverse problems, and hybrid architectures combining CNNs and transformers.

- **Chapter 3: Methods** details the data simulation pipeline, model architectures, training strategies, and evaluation metrics used in this research.
- **Chapter 4: Experiments** presents experimental design, training protocols, and validation strategies.
- **Chapter 5: Results** reports quantitative and qualitative results for both Stage 1 and Stage 2 training, including ablation studies and comparison with baseline methods.
- **Chapter 6: Discussion** interprets the findings, analyses limitations, and explores implications for clinical translation.
- **Chapter 7: Conclusion** summarises the contributions, highlights potential future research directions, and reflects on the broader impact of this work.

- 2.1 Physics of NIR Light Propagation
- 2.2 Classical NIR-DOT Reconstruction Methods
- 2.3 Deep Learning for Medical Image Reconstruction
- 2.4 CNN Approaches in NIR-DOT
- 2.5 Transformers in Medical Imaging
- 2.6 Hybrid Architectures and Multi-Stage Training
- 2.7 Contextual Integration in Medical Imaging
- 2.8 Research Gap and Opportunity

- 3.1 Physics-Based Forward Modeling**
- 3.2 Geometric Phantom Construction**
- 3.3 Optical Property Assignment**
- 3.4 Surface Extraction and Probe Placement**
- 3.5 Frequency-Domain Measurement Simulation**
- 3.6 Dataset Composition and Statistics**
- 3.7 Validation of Simulated Data**

4.1 Architectural Overview

4.2 CNN Autoencoder Design

4.2.1 Encoder Architecture

4.2.2 Decoder Architecture

4.2.3 Residual Connections and Feature Extraction

4.3 NIR Measurement Processing

4.3.1 Spatial Awareness in Measurement Processing

4.3.2 Tissue Patch Extraction and Encoding

4.4 Transformer Encoder Design

4.4.1 Self-Attention Mechanism

4.4.2 Token-Type Embeddings

4.4.3 Feature Enhancement

4.5 Integration Strategy

4.5.1 Information Flow Between Components

4.5.2 Tissue Context Integration

5.1 Training Strategy Overview

5.2 Stage 1: CNN Autoencoder Pre-Training

5.2.1 Identity Mapping Objective

5.2.2 AdamW with OneCycleLR Optimization

5.2.3 Hyperparameter Selection

5.3 Stage 2: Transformer Enhancement Training

5.3.1 Frozen Decoder Approach

5.3.2 Linear Warmup with Cosine Decay

5.3.3 Differential Weight Decay

5.4 Data Augmentation Strategy

5.5 Implementation Details

5.5.1 Hardware Optimization

5.5.2 Experiment Tracking

6.1 Experimental Setup

6.1.1 Dataset Preparation

6.1.2 Evaluation Metrics

6.1.3 Baseline Methods

6.2 Stage 1 Results: CNN Autoencoder Performance

6.2.1 Reconstruction Quality Metrics

6.2.2 Feature Learning Analysis

6.3 Stage 2 Results: Transformer Enhancement

6.3.1 Baseline Mode Performance

6.3.2 Enhanced Mode with Tissue Context

6.3.3 Comparative Analysis

6.4 Ablation Studies

6.4.1 Impact of Tissue Patch Size

6.4.2 Effect of Transformer Layers and Heads

6.4.3 Sensitivity to Tissue Context Quality

6.5 Visualization and Interpretation

6.5.1 Attention Map Analysis

6.5.2 Reconstruction Quality Visualization

6.5.3 Error Distribution Analysis

7.1 Key Findings and Insights

7.2 Clinical Implications

7.3 Limitations of Current Approach

7.4 Computational Efficiency Considerations

Conclusion and Future Work

8.1 Summary of Contributions

8.2 Future Research Directions

8.3 Potential Clinical Applications

Bibliography

APPENDIX A

Implementation Details

APPENDIX B

Additional Experimental Results

APPENDIX C

Mathematical Derivations
