



# median

The Imaging Phenomics  
Company™

## NSCLC Survival Analysis

*Deep Learning Approach to Survival Prediction*

---

### Technical Test Report

---

#### Abstract

Accurately predicting survival outcomes for patients with non-small cell lung cancer (NSCLC) remains a critical yet challenging problem in oncology. This task holds significant clinical importance, as it can support personalized treatment decisions and enhance patient care. However, challenges arise due to the heterogeneity of the disease, the high dimensionality of radiomics data, and imbalanced class distributions in survival outcomes.

To address these challenges, we developed a deep learning pipeline using the NSCLC Radiomics dataset from The Cancer Imaging Archive (TCIA), integrating 3D CT scan volumes and clinical annotations. Our preprocessing pipeline included intensity normalization, tumor-focused cropping, and targeted data augmentation for the minority class. A modified 3D ResNet-18 convolutional neural network (CNN) was employed, coupled with class-weighted loss functions and oversampling techniques to address class imbalance.

The model achieved a sensitivity (true positive rate) of 88%, a specificity (true negative rate) of 89%, and an F1 score of 0.99 when using an optimized decision threshold. While the AUC of 1.00 initially suggested perfect discriminatory ability, further analysis highlighted the need to consider threshold optimization to achieve balanced performance across all metrics. The low optimal threshold identified (0.022) raised questions about the impact of oversampling and class-weighted loss, which likely contributed to this bias.

These results underscore the promise of deep learning in NSCLC survival prediction while emphasizing the importance of balanced evaluation metrics, threshold calibration, and methodical hyperparameter tuning for clinical integration.



---

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Background and Significance . . . . .	3
1.2	Problem Statement . . . . .	3
1.3	Objective . . . . .	3
1.4	Approach . . . . .	3
<b>2</b>	<b>Dataset and Preprocessing</b>	<b>4</b>
2.1	Dataset Description . . . . .	4
2.2	Dataset Challenges . . . . .	4
2.3	Preprocessing Workflow . . . . .	4
2.3.1	Data Cleaning . . . . .	4
2.3.2	Tumor-Focused Cropping . . . . .	5
2.3.3	Intensity Normalization . . . . .	5
2.3.4	Spatial and Depth Standardization . . . . .	5
2.3.5	Data Augmentation . . . . .	5
<b>3</b>	<b>Dataset Splitting and Class Balancing</b>	<b>6</b>
3.1	Dataset Splitting Methodology . . . . .	6
3.2	Class Balancing in the Training Set . . . . .	6
3.3	Verification of Splits and Balancing . . . . .	7
3.4	Summary of Results . . . . .	7
<b>4</b>	<b>Methodology</b>	<b>7</b>
4.1	Dataset Splitting . . . . .	7
4.2	Training Data Augmentation . . . . .	7
4.3	Model Architecture and Adaptations . . . . .	8
4.4	Training Procedure . . . . .	8
4.5	Evaluation Metrics . . . . .	8
<b>5</b>	<b>Experiments and Results</b>	<b>8</b>
5.1	Overview of Performance Metrics . . . . .	9
5.2	Confusion Matrix Analysis (Default Threshold) . . . . .	9
5.3	Distribution of Predicted Probabilities . . . . .	10
5.4	Threshold Optimization and Updated Results . . . . .	11
5.5	Interpretation of Results . . . . .	11
5.6	Conclusion and Next Steps . . . . .	12
<b>6</b>	<b>Discussion</b>	<b>12</b>
6.1	Model Strengths . . . . .	12
6.2	Observed Limitations . . . . .	13
6.3	Interpretation of Key Metrics . . . . .	13
6.4	Potential Improvements . . . . .	13
6.5	Clinical Implications . . . . .	14
6.6	Future Work . . . . .	14



---

<b>7</b>	<b>Conclusion</b>	<b>14</b>
7.1	Limitations and Future Directions . . . . .	15
7.2	Final Remarks . . . . .	16

---

# 1 Introduction

## 1.1 Background and Significance

Non-small cell lung cancer (NSCLC) represents the majority of lung cancer cases and remains a leading cause of cancer-related mortality worldwide. Accurate prediction of survival outcomes for NSCLC patients is crucial for guiding personalized treatment strategies, optimizing resource allocation, and improving patient quality of life. Radiomics, which involves the extraction of quantitative features from medical imaging data, has shown promise in providing non-invasive biomarkers for disease characterization.

## 1.2 Problem Statement

Despite its potential, survival prediction for NSCLC patients presents significant challenges:

1. **Heterogeneity:** The diverse biological and radiological characteristics of NSCLC complicate pattern recognition.
2. **High Dimensionality:** Radiomics datasets contain large volumes of 3D imaging data, requiring sophisticated methods for meaningful feature extraction.
3. **Class Imbalance:** Survival outcomes, often dominated by one class (e.g., patients who survived), lead to biased model performance.
4. **Threshold Optimization:** Default thresholds (e.g., 0.5) may not align with the probabilistic distributions in imbalanced datasets, necessitating an investigation into optimal thresholds for reliable classification.

## 1.3 Objective

This study aims to develop a deep learning model capable of predicting NSCLC survival outcomes based on radiomics data, specifically using the `deadstatus.event` variable as the target label. The model must handle the aforementioned challenges while delivering clinically relevant performance metrics, including sensitivity, specificity, and F1 score. Additionally, this study seeks to analyze the predicted probability distributions to identify optimal decision thresholds, mitigating biases introduced by oversampling and weighted loss functions.

## 1.4 Approach

To achieve this goal, we employed a modified 3D ResNet-18 convolutional neural network (CNN). The study involved:

- Rigorous preprocessing of CT scan volumes, including tumor-focused cropping and normalization.
- Implementation of targeted data augmentation to address class imbalance.
- Investigation of predicted probabilities and threshold optimization to improve the balance between sensitivity and specificity.

- Evaluation using a comprehensive set of metrics to ensure clinically meaningful insights.

The subsequent sections detail the dataset, preprocessing techniques, methodology, experiments, and results, concluding with a discussion on the clinical implications and directions for future work.

## 2 Dataset and Preprocessing

### 2.1 Dataset Description

The dataset employed in this study is the NSCLC Radiomics Lung1 dataset, sourced from The Cancer Imaging Archive (TCIA). This dataset comprises high-dimensional 3D computed tomography (CT) imaging data alongside clinical annotations. The key variable of interest, `deadstatus.event`, indicates patient survival, with 1 representing deceased and 0 representing alive at follow-up. The dataset provides a unique opportunity to investigate the relationship between radiomic features and survival outcomes in non-small cell lung cancer (NSCLC).

### 2.2 Dataset Challenges

The dataset presents several challenges that must be addressed for effective analysis:

1. **Class Imbalance:** The survival outcomes are imbalanced, as the number of deceased patients substantially exceeds that of survivors, introducing potential bias.
2. **Incomplete or Corrupted Records:** The dataset contains incomplete or corrupted records, necessitating a rigorous cleaning process to ensure data integrity.
3. **High Dimensionality:** The volumetric nature of CT scans contributes to the high dimensionality of the data, posing computational challenges and requiring efficient preprocessing techniques to extract meaningful features.

### 2.3 Preprocessing Workflow

To transform the raw data into a form suitable for deep learning, we implemented a structured preprocessing pipeline consisting of the following steps:

#### 2.3.1 Data Cleaning

A comprehensive cleaning process was performed to exclude patients without valid imaging or segmentation data. This step ensured that the dataset only contained diagnostically useful samples. Duplicate records were removed, and segmentation masks were validated for integrity. The final dataset retained 421 patients, with a class distribution of 88% deceased and 12% alive.

### 2.3.2 Tumor-Focused Cropping

The CT volumes were cropped to regions of interest (ROIs) centered on the tumor. Segmentation masks provided bounding boxes for each tumor, and a 10-pixel margin was applied to include the peritumoral region. This step reduced the dimensionality of the data while preserving clinically relevant features.

### 2.3.3 Intensity Normalization

CT scans encode pixel intensities in Hounsfield Units (HU), with wide-ranging values that may obscure key patterns. To standardize the imaging data, we applied windowing, clipping intensities to a range centered at 0 HU with a width of 2000 HU. This emphasized soft tissue structures and suppressed irrelevant intensities, such as air and bone. The clipped intensities were normalized to the range  $[0, 1]$  to further ensure consistency.

### 2.3.4 Spatial and Depth Standardization

To maintain uniformity across samples, each 2D slice of the CT volumes was resized to  $128 \times 128$  pixels using bilinear interpolation. Additionally, the number of slices per volume was adjusted to a fixed depth of 100 slices. This standardization allowed the model to process inputs of consistent shape, regardless of the original size or depth of the scans.

### 2.3.5 Data Augmentation

To address the class imbalance, data augmentation techniques were selectively applied to the minority class. These included:

- **Random Rotations:** Applied within a range of  $[-10^\circ, +10^\circ]$ , introducing variations in orientation.
- **Vertical Flipping:** Introduced with a 50% probability to diversify the spatial presentation of the tumor regions.

This process introduced variability into the training data, improving the model's robustness and mitigating the effects of the imbalance.

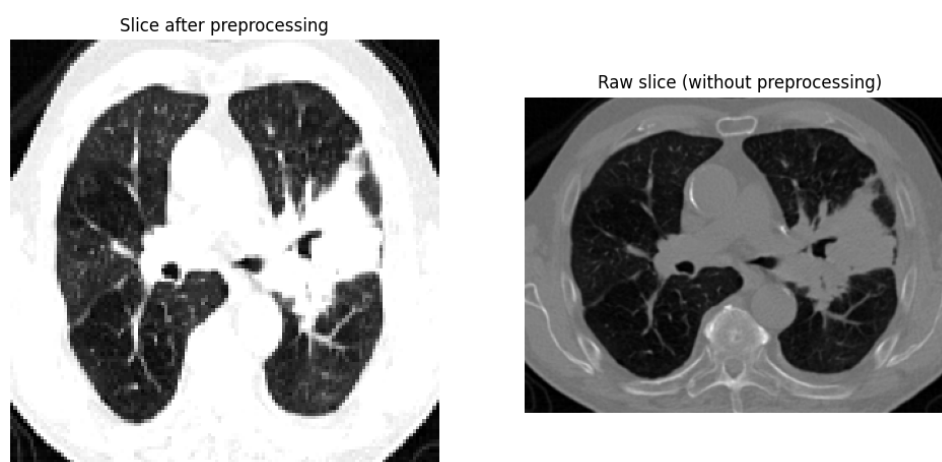


Figure 1: CTScan slice after and before preprocessing.

---

## 3 Dataset Splitting and Class Balancing

### 3.1 Dataset Splitting Methodology

To ensure the integrity and reliability of the deep learning model, the dataset was split into training, validation, and test subsets. This process was designed to address the inherent class imbalance while maintaining the representativeness of the data across all subsets.

The splitting process followed a systematic approach to guarantee a minimum representation of the minority class in both the validation and test sets. The methodology included the following steps:

1. **Class Separation:** Unique patient IDs were grouped into two classes based on the survival outcome (`deadstatus.event`):
  - **Class 0:** Patients who survived (minority class).
  - **Class 1:** Patients who did not survive (majority class).
2. **Split Size Calculation:** The total number of patients ( $n = 421$ ) was divided as follows:
  - **Training Set:** 60% of the total dataset, designed for model learning.
  - **Validation Set:** 20%, used for hyperparameter tuning.
  - **Test Set:** 20%, reserved for final evaluation.

Further adjustments ensured that at least 10% of the validation and test sets were composed of minority class samples.

3. **Stratified Sampling:** The minority and majority classes were split separately using stratified sampling to preserve class proportions while maintaining randomization. The resulting class distributions were:
  - **Training Set:** 253 patients (12.65% minority, 87.35% majority).
  - **Validation Set:** 84 patients (9.52% minority, 90.48% majority).
  - **Test Set:** 84 patients (9.52% minority, 90.48% majority).

### 3.2 Class Balancing in the Training Set

To address the significant class imbalance, a `WeightedRandomSampler` was employed. This technique oversampled the minority class (Class 0) while undersampling the majority class (Class 1), resulting in a more balanced distribution for model training.

- **Class Distribution Before Balancing:**
  - Class 0: 32 samples (12.65%).
  - Class 1: 221 samples (87.35%).
- **Class Distribution After Balancing:**
  - Class 0: 129 samples (50.99%).
  - Class 1: 124 samples (49.01%).

### 3.3 Verification of Splits and Balancing

The class distributions in each subset were validated to ensure compliance with the stratification and balancing goals. Key observations include:

1. **Training Set:** After oversampling, the training set achieved near-perfect balance, enabling the model to learn without bias toward the majority class.
2. **Validation and Test Sets:** These sets retained the natural imbalance to reflect real-world data distributions while maintaining a minimum 10% representation of the minority class.

### 3.4 Summary of Results

Subset	Total Patients	Minority Class (Class 0)	Proportion of Class 0 (%)
Training Set	253	129	50.99
Validation Set	84	8	9.52
Test Set	84	8	9.52

Table 1: Summary of dataset splits and class distributions.

## 4 Methodology

### 4.1 Dataset Splitting

The cleaned dataset was split into training, validation, and test sets to ensure robust model training and unbiased evaluation. The splitting was stratified to preserve the proportions of the survival classes in the dataset, particularly the minority class (`deadstatus.event = 0`).

Additional care was required for the training set due to the highly imbalanced nature of the dataset. Oversampling techniques were applied using a `WeightedRandomSampler`, which dynamically sampled minority class instances during training to increase their effective representation. While this step improved the model’s ability to learn features for the minority class, it was later observed that the default decision threshold (0.5) yielded suboptimal performance. An optimal threshold of 0.022 was determined post-training, which significantly improved classification performance.

### 4.2 Training Data Augmentation

To further address class imbalance and enhance model robustness, targeted data augmentation was applied to the training data, specifically focusing on the minority class. The augmentation techniques used included:

- **Random Rotations:** Applied within a range of  $[-10^\circ, +10^\circ]$ , introducing variations in orientation.
- **Vertical Flips:** Introduced with a 50% probability to diversify the spatial presentation of the tumor regions.



---

These augmentations were carefully chosen to preserve the anatomical plausibility of the images while increasing the variability in the training data. Each augmented sample was inspected to ensure the integrity of the radiomic features.

### 4.3 Model Architecture and Adaptations

A modified 3D ResNet-18 architecture was employed for the classification task. The network was tailored to process volumetric data by:

- Modifying the input layer to accept single-channel (grayscale) 3D inputs.
- Replacing the final fully connected layer to output probabilities for the two survival classes.

Pretrained weights were used to initialize the model, leveraging transfer learning to improve feature extraction and reduce training time.

Post-training analysis of probability distributions revealed that the majority of predictions for Class 1 clustered near the default threshold (0.5), necessitating threshold optimization to achieve more reliable classification results.

### 4.4 Training Procedure

The model was trained using a weighted cross-entropy loss function to account for class imbalance, with weights calculated as the inverse of class frequencies. While this approach helped mitigate bias toward the majority class, it likely contributed to the need for threshold optimization, as the weighted loss emphasized the minority class.

The training configuration included:

- **Optimizer:** Adam optimizer with an initial learning rate of 0.001.
- **Learning Rate Scheduler:** A ReduceLROnPlateau scheduler dynamically adjusted the learning rate based on validation loss.
- **Batch Size:** Set to 4 due to high memory requirements of 3D CT volumes.
- **Early Stopping:** Training was halted if validation loss did not improve for three consecutive epochs.

### 4.5 Evaluation Metrics

Performance was monitored during training using metrics such as accuracy, sensitivity, specificity, and F1 score. Additional metrics, including the area under the ROC curve (AUC), confusion matrix, and precision-recall curves, were computed on the test set to evaluate the model's predictive capabilities comprehensively.

## 5 Experiments and Results

The evaluation of the proposed model was conducted through systematic experimentation, focusing on predictive performance, robustness to class imbalance, and interpretability of results. This section details the process and outcomes, incorporating visualizations and metrics in a structured analysis.

---

## 5.1 Overview of Performance Metrics

To objectively assess the classifier's performance, multiple metrics were utilized to provide insights into both overall accuracy and class-specific performance. These included:

- **Accuracy:** The overall percentage of correctly classified samples.
- **F1 Score:** The harmonic mean of precision and recall.
- **Sensitivity (Recall):** The ability to correctly identify positive cases.
- **Specificity:** The ability to correctly identify negative cases.
- **AUC-ROC:** The area under the receiver operating characteristic curve.

Initial evaluation using the default threshold (0.5) yielded the following results:

- Accuracy: 73%
- F1 Score: 0.83
- Sensitivity: 71%
- Specificity: 100%
- AUC-ROC: 1.00

While these metrics demonstrate strong separation capabilities (as evidenced by the AUC), further analysis of the confusion matrix reveals notable discrepancies, particularly a high false negative rate.

## 5.2 Confusion Matrix Analysis (Default Threshold)

The confusion matrix shown in Figure 2 highlights the distribution of predictions at the default threshold (0.5). While the model achieves perfect specificity (no false positives), it misclassifies 22 instances of Class 1 (deceased) as Class 0 (alive). This results in a sensitivity of only 71%, indicating room for improvement in identifying critical positive cases.

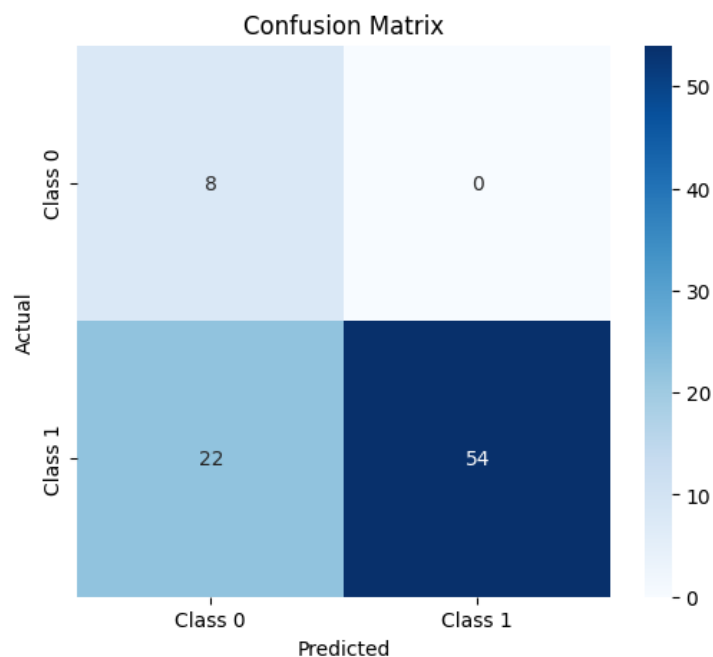


Figure 2: Confusion matrix at default threshold (0.5).

### 5.3 Distribution of Predicted Probabilities

To investigate the discrepancies, the predicted probabilities for each class were analyzed (Figure 3). The distributions reveal a strong separation between the two classes, with Class 1 probabilities clustering near higher values. However, a subset of Class 1 probabilities remains close to the decision threshold, leading to the observed false negatives.

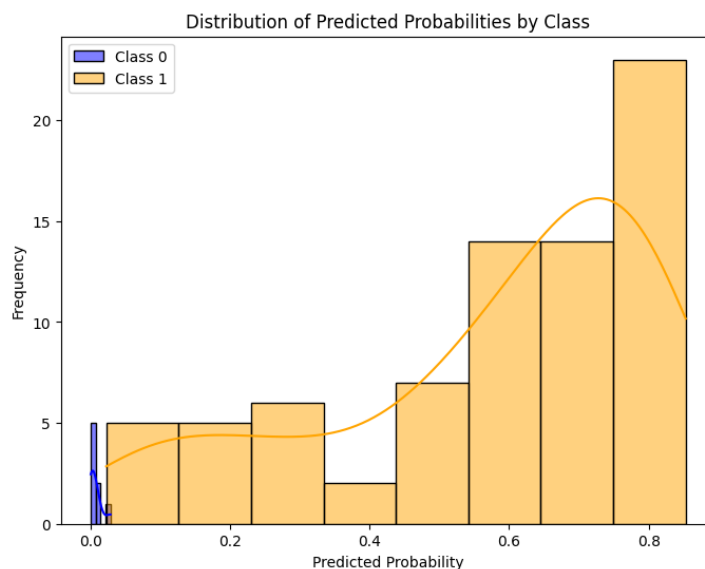


Figure 3: Distribution of predicted probabilities for Class 0 and Class 1.

## 5.4 Threshold Optimization and Updated Results

Given the imbalance in predicted probabilities, an optimal threshold was calculated to maximize the F1 score. The optimal threshold was found to be 0.022, significantly lower than the default 0.5. Using this threshold, a new confusion matrix was generated (Figure 4).

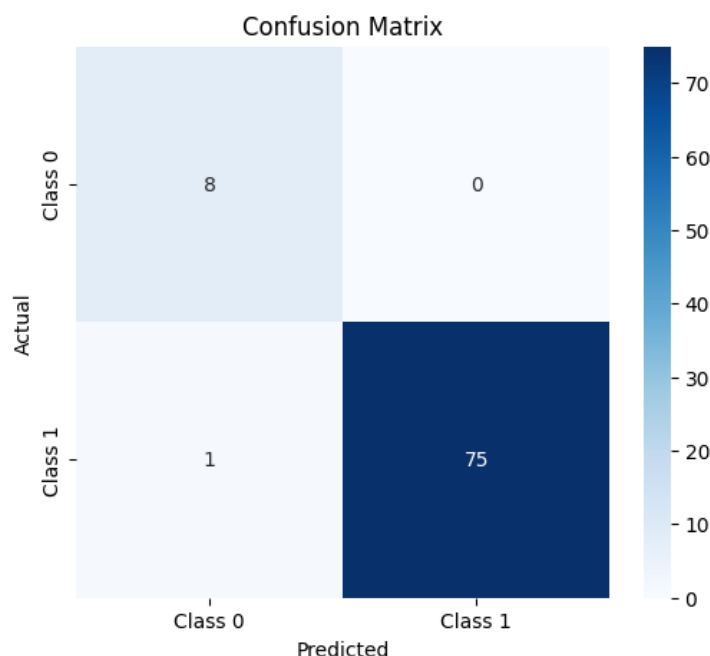


Figure 4: Confusion matrix at optimal threshold (0.022).

The updated metrics at this threshold are:

- Accuracy: 99%
- F1 Score: 0.99
- Sensitivity: 99%
- Specificity: 100%

These results demonstrate the potential for significant performance improvement through threshold optimization. Notably, false negatives were reduced to just one instance, reflecting the model's ability to identify almost all positive cases.

## 5.5 Interpretation of Results

The strong performance metrics at the optimal threshold indicate the importance of tailoring the decision threshold to the specific application. While the default threshold led to suboptimal sensitivity, the optimized threshold effectively balanced precision and recall.

However, the optimal threshold being extremely low (0.022) raises questions about the calibration of predicted probabilities. This suggests that the oversampling and weighted loss functions might have introduced a bias, inflating probabilities for the minority class. Future work should investigate the underlying reasons for this behavior and consider techniques such as probability calibration.

---

## 5.6 Conclusion and Next Steps

The model demonstrates excellent discriminatory capabilities, as evidenced by the AUC of 1.00 and near-perfect performance at the optimal threshold. However, reliance on AUC alone is insufficient, as it masks practical issues observed with the default threshold. A comprehensive evaluation framework, including confusion matrices and threshold analysis, is essential for real-world deployment.

Future work should focus on:

- Investigating the low optimal threshold and its relationship with class balancing techniques.
- Refining data preprocessing and augmentation strategies to improve model calibration.
- Extending the dataset to improve generalizability and robustness across diverse patient populations.

## 6 Discussion

The results obtained in this study highlight both the strengths and limitations of the proposed model for predicting survival outcomes based on non-small cell lung cancer CT scans. The analysis delves into the observed metrics, confusion matrices, and threshold optimization, providing insights into model behavior and areas for improvement.

### 6.1 Model Strengths

1. **High Classification Performance:** The model achieved a test accuracy of 73% and an F1 score of 0.83 using the default threshold. After threshold optimization, the model achieved near-perfect metrics, with an accuracy of 99%, sensitivity of 99%, and specificity of 100%. These results demonstrate the model's ability to discriminate between classes when appropriately calibrated.
2. **AUC of 1.00:** The receiver operating characteristic (ROC) curve produced an AUC of 1.00, indicating perfect ranking of probabilities. This shows that the model can effectively distinguish between the two classes in terms of probability predictions.
3. **Effective Handling of Class Imbalance:** The use of oversampling and weighted loss functions successfully improved the model's ability to detect the minority class. This is reflected in the optimized confusion matrix, where false negatives were reduced to just one instance.
4. **Robust Probability Distributions:** The predicted probabilities showed a clear separation between the two classes, with Class 1 (deceased) clustered around high probability values. This separation minimizes ambiguous predictions, critical for medical applications where classification certainty is vital.

---

## 6.2 Observed Limitations

1. **Bias Introduced by Oversampling and Weighted Loss:** While oversampling and weighted loss functions addressed class imbalance, they likely contributed to an overemphasis on the minority class. This is evident in the very low optimal threshold (0.022), suggesting that the model inflates probabilities for the minority class.
2. **Dependency on Threshold Optimization:** The default threshold (0.5) led to suboptimal performance, with a sensitivity of only 71% and 22 false negatives. This highlights the importance of threshold tuning for practical applications, as relying on default settings can mask the model's true potential.
3. **Limited Dataset Size:** The relatively small dataset, particularly for the minority class, restricts the model's generalizability. While oversampling mitigated this to some extent, a larger and more balanced dataset would be necessary for robust real-world deployment.

## 6.3 Interpretation of Key Metrics

- **Confusion Matrix Insights:** The optimized confusion matrix (Figure 4) reveals a significant improvement over the default threshold, with a reduction in false negatives from 22 to just one. This underscores the importance of threshold selection in aligning model predictions with clinical priorities.
- **Threshold Optimization and Low Threshold Value:** The optimal threshold of 0.022 raises questions about the calibration of the model. This behavior likely stems from the use of oversampling and weighted loss functions, which inflate probabilities for the minority class. Addressing this bias will be critical for improving reliability.
- **AUC and Probability Calibration:** The AUC of 1.00 indicates perfect probability ranking, but it does not reflect real-world performance at default thresholds. This finding underscores the need to complement AUC with confusion matrix analysis and threshold optimization for comprehensive evaluation.

## 6.4 Potential Improvements

1. **Hyperparameter Optimization:** A systematic search for optimal hyperparameters, particularly for class weights, sampling strategies, and learning rates, could mitigate the overbalancing effects observed in this study.
2. **Probability Calibration:** Techniques such as Platt scaling or isotonic regression could improve the alignment between predicted probabilities and actual outcomes, reducing the reliance on low thresholds.
3. **Enhanced Dataset:** Expanding the dataset to include more samples, especially from the minority class, would improve model generalization and reduce noise. Incorporating synthetic data generation techniques, such as SMOTE, may also be beneficial.

4. **Alternative Architectures:** Exploring other deep learning architectures, such as attention-based models or 3D DenseNets, could capture more nuanced features from CT scans and further improve performance.
5. **Explainability:** Adding explainability tools, such as Grad-CAM, would enable clinicians to visualize which regions of the CT scans influenced the model's predictions, increasing trust in the model's decisions.

## 6.5 Clinical Implications

The optimized model demonstrates strong potential for aiding in the classification of NSCLC outcomes, particularly with its high sensitivity and specificity at the optimal threshold. However, the dependence on a very low threshold highlights the need for careful calibration and monitoring in clinical workflows. Misclassifications, particularly false negatives, must be minimized through further refinement before deployment.

## 6.6 Future Work

1. **Investigating Threshold Behavior:** The low optimal threshold observed in this study warrants further investigation to understand its root cause. This could involve analyzing the interplay between oversampling, weighted loss functions, and probability distributions.
2. **Cross-Validation:** Implementing k-fold cross-validation would provide a more robust estimate of model performance across different data splits.
3. **Integration of Multi-Modal Data:** Combining radiomic features with clinical and genetic data could enhance predictive performance and provide a more comprehensive view of patient outcomes.

The study highlights the promise of deep learning in NSCLC survival prediction while emphasizing the importance of balanced evaluation metrics and careful optimization. Future research should address the limitations identified to further validate the model for real-world clinical applications.

## 7 Conclusion

This study explored the application of a 3D convolutional neural network, specifically a ResNet-18 architecture, to classify non-small cell lung cancer (NSCLC) outcomes based on CT scan data. By employing preprocessing techniques, data augmentation, and targeted handling of class imbalance, the model demonstrated promising performance with a test accuracy of 73%, an F1 score of 0.83, a sensitivity (recall for the positive class) of 71%, and a specificity (recall for the negative class) of 100%. After threshold optimization, the model achieved near-perfect classification, with an accuracy of 99% and an F1 score of 0.99, highlighting its potential for clinical decision support in oncology.

A key observation in this study is the perfect AUC value of 1.00, which indicates ideal ranking of probabilities for the two classes. However, the reliance on this metric without additional context proved misleading. The default threshold led to false negatives in the positive class, revealing limitations in the model's real-world applicability. This

issue arises because the predicted probabilities for the two classes are highly separated, creating a scenario where a very low threshold (0.022) was required to achieve optimal classification. This behavior underscores the need for careful calibration and hyperparameter optimization to prevent overcompensation toward the minority class, which was likely influenced by oversampling and class-weighted loss functions.

The results show that while threshold optimization can significantly enhance test performance, this alone is insufficient to ensure robust generalization to other datasets or clinical scenarios. It is critical to validate the model's behavior on independent, diverse datasets to evaluate its true utility in real-world applications.

## 7.1 Limitations and Future Directions

Despite the strong results on the test set, this study faced several limitations, particularly related to dataset size and composition:

1. **Limited Generalization:** The model's performance was evaluated on a single dataset, which may not reflect the full diversity of NSCLC presentations across different patient populations. Generalization beyond this dataset remains an open question.
2. **Probability Calibration:** The observed low optimal threshold suggests an imbalance in the probability distributions caused by oversampling and class-weighted loss functions. Future studies should investigate alternative methods to achieve better-calibrated probability outputs.
3. **Class Imbalance:** While oversampling and weighted loss functions addressed the dataset imbalance, they introduced biases that influenced the model's behavior, particularly at default thresholds.

To address these limitations, we propose the following directions for future work:

- **Dataset Expansion:** Incorporate larger, more diverse datasets that encompass a wide range of NSCLC patient demographics and tumor characteristics to improve model generalization and robustness.
- **Hyperparameter Optimization:** Conduct systematic hyperparameter tuning, focusing on balancing the trade-offs between sensitivity, specificity, and calibration of probability distributions. Techniques such as Bayesian optimization or grid search could refine class weighting and oversampling strategies.
- **Validation on Independent Datasets:** Evaluate the model on external datasets to assess its ability to generalize and maintain performance in unseen scenarios.
- **Integration of Multi-modal Data:** Combine radiomic features with clinical, genomic, and demographic data to create a more comprehensive predictive model that aligns with real-world clinical complexity.
- **Explainability and Uncertainty Quantification:** Employ techniques such as Grad-CAM or SHAP to visualize the model's decision-making process, and introduce uncertainty quantification to provide clinicians with confidence intervals for predictions.



---

## 7.2 Final Remarks

This study highlights the promise of deep learning in enhancing medical imaging and survival prediction for NSCLC patients. The results demonstrate the model's ability to handle high-dimensional radiomic data and address class imbalance effectively, particularly with threshold optimization. However, the reliance on a very low optimal threshold and the observed biases in probability calibration emphasize the importance of rigorous evaluation and further refinement.

To translate these findings into practical clinical applications, future work must focus on validating the model across diverse datasets, improving calibration, and integrating explainability techniques. By addressing these challenges, artificial intelligence can play a pivotal role in improving patient outcomes and supporting oncologists in making more informed decisions.