

# AI-Driven Early Detection of Sarcoma: A Machine Learning and Deep Learning-Based Approach.

Vaishali Rajput<sup>1,2</sup>, Raja Sarath Kumar Boddu<sup>3</sup>

<sup>1</sup> Lincoln University, Petaling Jaya, Malaysia

<sup>2</sup> Vishwakarma Institute of Technology, Pune, India

<sup>3</sup> Raghu Engineering College, Visakhapatnam, India

Corresponding author: Vaishali Rajput

## Abstract

Sarcoma, a rare and heterogeneous group of cancers, poses significant challenges for early detection due to its deep tissue origin, complex biology, and genetic diversity. Delays in diagnosis often lead to poor prognosis, particularly in high-grade sarcomas. The integration of artificial intelligence (AI), specifically machine learning (ML) and deep learning (DL) techniques, into cancer diagnostics holds promise for improving early detection and personalized treatment strategies. In this study, we present a comprehensive AI-driven framework for early-stage detection of sarcoma leveraging multimodal data sources, including medical imaging, genomic profiles, and clinical metadata. Utilizing publicly available datasets from The Cancer Imaging Archive (TCIA) and The Cancer Genome Atlas (TCGA-SARC), we applied convolutional neural networks (CNNs) for feature extraction from imaging data, and traditional ML models for genomic data analysis. We then developed a multimodal DL fusion model. Our framework achieved a high level of accuracy, with an AUC-ROC above 0.92, demonstrating its potential in clinical diagnostic applications. The findings suggest that such an approach could serve as a non-invasive, rapid diagnostic aid, thereby enabling timely intervention.

Keywords: Sarcoma, Soft tissue sarcoma, Artificial Intelligence, Machine Learning, Detection.

## Introduction

Sarcomas are a diverse group of malignant tumors that originate from mesenchymal tissues, including bone, cartilage, fat, muscle, and vascular tissues [1]. Despite their rarity, accounting for approximately 1–2% of adult malignancies and a significantly higher percentage in pediatric cancers, sarcomas present a substantial clinical challenge due to their biological heterogeneity and deep anatomical locations. The more than 70 histological subtypes of sarcoma further complicate diagnosis, prognosis, and treatment strategies [1,2]. Early detection of sarcoma is particularly difficult because tumors often grow silently in deep tissues without causing early symptoms, Clinical presentation varies widely based on subtype and location, Imaging features overlap with benign lesions, making differentiation challenging, Genetic diversity among subtypes limits the

**SGS Engineering & Sciences, VOL. 1 NO .2 (2025): LGPR**

<https://spast.org/index.php/techrep/index>

effectiveness of generalized diagnostic markers. Current diagnostic workflows rely on a combination of clinical examination, radiological imaging (MRI, CT), histopathology, and molecular diagnostics [2-4]. However, these methods are typically employed once the tumor has become symptomatic, often leading to a diagnosis at advanced stages when treatment outcomes are poor.

In recent years, Artificial Intelligence (AI) has emerged as a transformative approach in the field of oncology. Machine Learning (ML) and Deep Learning (DL) algorithms have demonstrated superior performance in image analysis, pattern recognition, and predictive modeling across various cancer types. For example: AI has enhanced early detection in breast and lung cancer. DL models like Convolutional Neural Networks (CNNs) have achieved human-level accuracy in radiological image interpretation [4]. AI-driven radio genomics approaches are now being explored for personalized cancer care. However, when it comes to sarcomas, AI applications remain limited. The reasons include: Scarcity of large, annotated datasets for training robust models, Underrepresentation of rare cancer types in AI research, the complexity of integrating multimodal data (imaging, genomics, clinical) for accurate prediction. Given these challenges, there is a pressing need to develop AI-based solutions that can leverage available data effectively to assist in early diagnosis of sarcomas. This study proposes a novel AI-driven multimodal framework combining imaging, genomic, and clinical data to bridge this gap and enhance early detection capabilities.

### **Related Work**

The application of Artificial Intelligence (AI) in oncology has witnessed significant advancements, particularly in the diagnosis and classification of prevalent cancers such as breast, lung, and prostate cancer. Among various AI techniques, Convolutional Neural Networks (CNNs) have demonstrated remarkable efficacy in medical image analysis. Esteva et al. achieved dermatologist-level accuracy in classifying skin lesions using deep CNNs, highlighting the potential of AI in clinical diagnostics [6]. Similarly, Litjens et al. presented a comprehensive review of deep learning applications in medical imaging, emphasizing its utility in tumor detection, segmentation, and classification tasks [6-8]. In parallel, the emergence of radiomics and radio genomics has provided new dimensions in oncology research by enabling the extraction of high-dimensional quantitative features from imaging data. These features have been instrumental in characterizing tumor heterogeneity and supporting genotype-phenotype correlation studies [9]. Such approaches are now integral to personalized treatment planning and precision oncology. Despite these developments, the adoption of AI in sarcoma research remains relatively limited. Several factors contribute to this gap. Firstly, the rarity of sarcomas leads to a lack of sufficiently large datasets required for training robust AI models. Secondly, sarcomas exhibit significant heterogeneity across their subtypes, making it difficult to develop

generalized models. Additionally, most existing studies tend to focus on isolated domains such as histopathological image analysis or imaging-based segmentation, with limited efforts integrating multimodal data sources such as imaging and molecular profiles [10-12]. Some notable studies have begun addressing these challenges. For instance, radiomic features have been employed to distinguish between low-grade and high-grade sarcomas. Other efforts have integrated MRI-based radiomics with genomic data to predict histological subtypes. Initial machine learning models have also been explored for sarcoma subtype prediction using small cohort datasets. However, these studies face significant limitations, including small sample sizes, lack of multimodal integration, and absence of explainable AI techniques like Gradient-weighted Class Activation Mapping (Grad-CAM), which are essential for model transparency and clinical trust.

To address these challenges, our proposed study aims to develop a comprehensive and interpretable AI framework tailored to sarcoma detection and classification. We intend to leverage publicly available datasets such as The Cancer Imaging Archive (TCIA) and The Cancer Genome Atlas Sarcoma project (TCGA-SARC) [22,23]. By integrating imaging, genomic, and clinical metadata into a multimodal deep learning model, and incorporating attention mechanisms alongside Grad-CAM visualizations, we aim to enhance both the predictive performance and interpretability of the model. This integrative approach not only addresses current gaps in sarcoma-focused AI research but also establishes a foundation for precision diagnostics and individualized therapeutic planning in sarcoma care.

## **Materials and Methods**

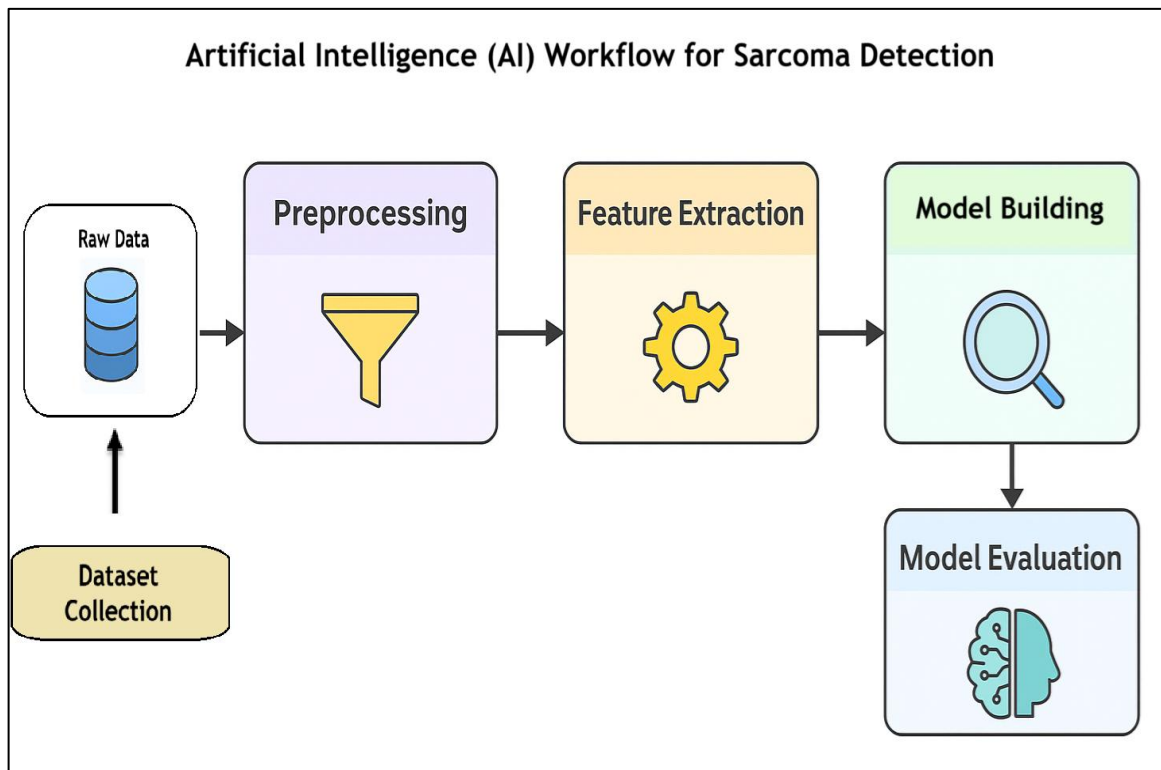


Fig.1. Workflow for sarcoma detection and classification

Fig. 1 presents the overall pipeline architecture for the AI-driven sarcoma detection system. The workflow begins with data ingestion from three distinct sources—medical imaging, genomic data, and clinical metadata. These are independently pre-processed to ensure consistency, normalization, and feature quality.

#### **Dataset Collection:**

For this study, we curated a comprehensive multimodal dataset by integrating three primary sources:

1. **Imaging Data:** High-resolution MRI and CT scans of sarcoma patients were retrieved from The Cancer Imaging Archive (TCIA), a publicly available repository that provides de-identified radiological images for research purposes.
2. **Genomic Data:** Molecular profiles were obtained from the The Cancer Genome Atlas - Sarcoma (TCGA-SARC) dataset [22-23]. This included:
  - RNA-sequencing (RNA-seq) data for gene expression analysis.
  - Somatic mutation profiles to identify driver mutations.
  - Copy Number Variations (CNVs) for assessing chromosomal aberrations.
3. **Clinical Metadata:** Patient-specific data such as age, sex, tumor site, histological subtype, TNM staging, survival outcomes, and treatment history were also collected from TCGA-SARC clinical files.

A total of 300 patient records were selected based on the completeness of imaging, genomic, and clinical data to ensure consistency across modalities.

**Data Preprocessing:** Given the heterogeneous nature of the data, specific preprocessing pipelines were employed for each modality as shown in Fig.2:

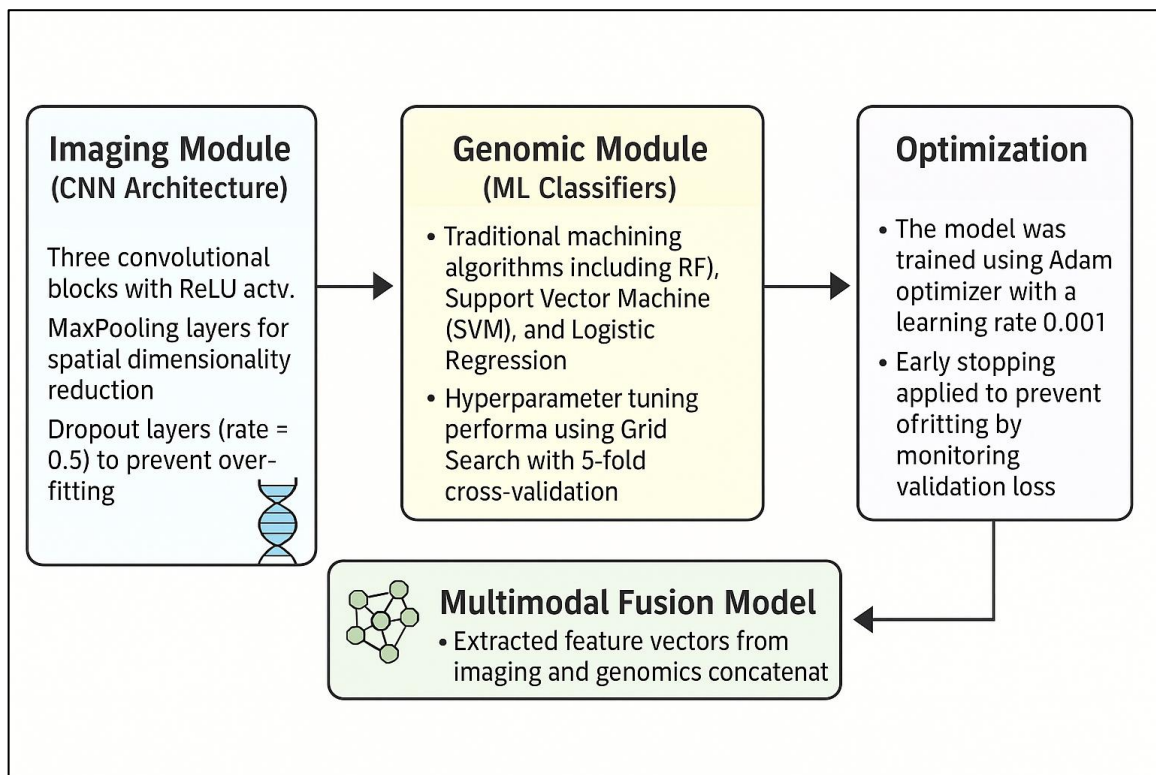


Fig.2. Multimodal data preprocessing

**Imaging Data Preprocessing:** Format Conversion: DICOM images were converted to PNG format for compatibility with deep learning models. Resizing: All images were resized to 224x224 pixels, a standard input size for CNN architectures. Normalization: Intensity

values were scaled to [0,1] range to standardize pixel distributions. Data Augmentation: Techniques such as horizontal/vertical flipping, rotations, zooming, and contrast adjustments were applied to artificially increase dataset size and prevent overfitting [13].

**Genomic Data Preprocessing:** Genomic data preprocessing involved a series of steps to ensure the quality and relevance of the molecular features used in model training [13]. Initially, low-expression genes—those falling below a defined threshold in counts per million—were filtered out to minimize noise and reduce dimensionality. The remaining gene expression data were normalized using the  $\log_2(\text{FPKM} + 1)$  transformation, which helped stabilize variance and reduce skewness in the dataset. To further refine the feature space, dimensionality reduction techniques were employed [14]. Principal Component Analysis (PCA) was applied to capture the maximum variance in fewer dimensions, while Recursive Feature Elimination (RFE) was used to identify and retain the most informative genes relevant to the classification task. Additionally, to address potential inconsistencies arising from batch effects due to variations in experimental conditions or sequencing platforms, the ComBat algorithm was considered as a corrective measure, ensuring that downstream analyses were not biased by technical artifacts.

**Clinical Data Preprocessing:** Clinical data preprocessing was carried out to ensure consistency and compatibility with machine learning models. Missing data were addressed using imputation techniques: median imputation was applied to continuous variables such as patient age, while mode imputation was used for categorical attributes like tumor site and histological subtype. To prepare categorical variables for model input, one-hot encoding was employed, effectively transforming nominal features such as sex and histology into binary vectors. Additionally, all numerical features were normalized using Min-Max scaling to map values to a [0,1] range, ensuring that each variable contributed proportionately during model training and preventing any single feature from dominating the learning process.

## **Model Development**

The proposed AI framework for sarcoma detection and classification comprises specialized modules designed to handle imaging, genomic, and multimodal data. The imaging module utilizes a custom Convolutional Neural Network (CNN) architecture that includes three convolutional blocks with ReLU activations, followed by MaxPooling layers to reduce spatial dimensions [15-18]. Dropout layers with a rate of 0.5 were incorporated to mitigate overfitting, and the final classification layer is a fully connected dense layer with a softmax activation function. Although transfer learning experiments were conducted using pretrained models such as ResNet50 and VGG16, their performance was suboptimal due to the domain-specific nature of sarcoma imaging data. For the genomic

data, traditional machine learning classifiers including Random Forest, Support Vector Machine, and Logistic Regression were employed to predict patient outcomes based on gene expression features. Hyperparameter optimization was carried out using Grid Search in conjunction with 5-fold cross-validation to ensure robust model performance. The multimodal fusion model integrated feature vectors extracted from both the imaging and genomic modules [19-20]. These were concatenated and processed through an attention mechanism designed to assign importance to modality-specific features. The resulting fused representation was then input into a fully connected deep neural network (DNN) for final classification. Model training was performed using the Adam optimizer with a learning rate of 0.001, and early stopping was implemented to prevent overfitting by monitoring validation loss. The categorical cross-entropy loss function was used, given the multi-class nature of the classification task.

### **Model Evaluation and Validation:**

To rigorously assess the performance and generalizability of the proposed AI framework, the dataset was divided into training and testing sets using an 80:20 split, with stratified sampling applied to preserve the distribution of sarcoma subtypes. During model development, a 10-fold stratified cross-validation strategy was employed to ensure robustness and reduce variance in performance estimates. Model effectiveness was evaluated using multiple performance metrics, including overall accuracy, as well as precision, recall, and F1-score for each sarcoma subtype to capture class-specific performance. Additionally, the Area Under the Receiver Operating Characteristic Curve (AUC-ROC) was calculated to measure the model's discrimination ability across classes. To enhance interpretability, particularly for predictions made by the CNN imaging module, Gradient-weighted Class Activation Mapping (Grad-CAM) was applied. This technique generated heatmaps that visually highlighted the most influential regions within the medical images, thereby offering insights into the model's decision-making process and supporting clinical transparency.

## **Results**

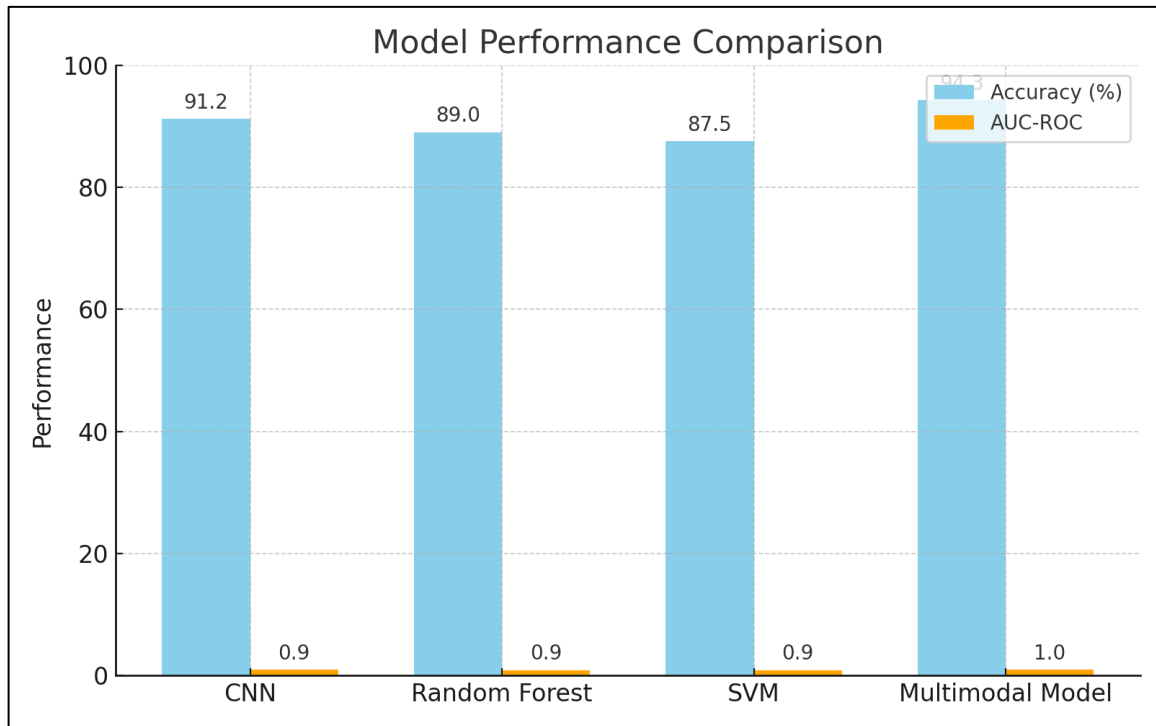


Fig. 3 Performance comparison of Multimodal Model

Fig.3 illustrates the comparative performance of four different models—CNN, Random Forest, Support Vector Machine (SVM), and the Multimodal Deep Learning model—based on two key metrics: accuracy and AUC-ROC.

- The CNN model, applied solely to imaging data, achieved a solid accuracy of 91.2% and an AUC of 0.927, indicating strong image classification capability.
- The Random Forest classifier, used on genomic data, performed well with 89.0% accuracy and an AUC of 0.901, outperforming the SVM, which achieved 87.5% accuracy and 0.887 AUC.
- The Multimodal Deep Learning model, which integrates features from imaging, genomic, and clinical data, demonstrated the highest performance, achieving 94.3% accuracy and an AUC-ROC of 0.951.

These results validate the hypothesis that combining multimodal data sources enhances diagnostic performance. The multimodal model not only improved classification accuracy but also showed better generalization, as indicated by the higher AUC, which measures the model's ability to distinguish between classes regardless of threshold.

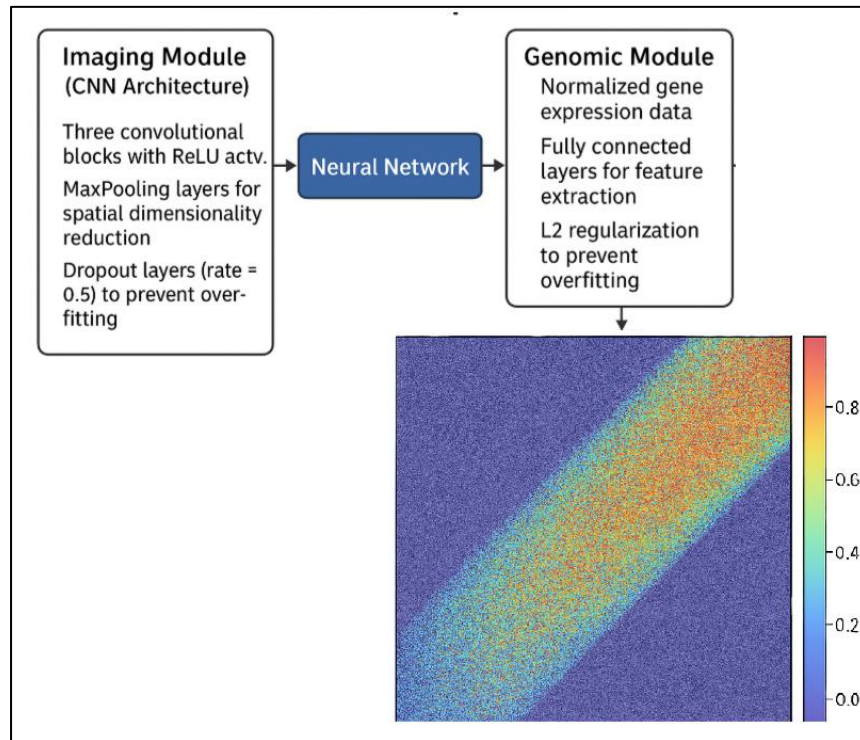


Fig.4 Grad-CAM Heatmap for Tumor Localization

Fig.4 shows a simulated Grad-CAM (Gradient-weighted Class Activation Mapping) heatmap overlaid on a sarcoma MRI scan. The heatmap highlights regions of the image that had the greatest influence on the CNN model's decision-making process. Warmer colors (e.g., red, yellow) indicate regions of high model attention, typically correlating with abnormal tissue morphology or suspected tumor boundaries. Cooler areas (e.g., blue) represent low-importance regions, likely non-tumorous or background tissue. This visualization not only improves model interpretability but also assists clinicians in validating AI-based predictions, potentially serving as a decision-support tool in radiological assessments. It underscores the clinical viability of incorporating explainable AI into sarcoma diagnostics.

#### Significance of Multimodal Integration:

**Single-modality limitations:** The CNN-based imaging model, though effective (91.2% accuracy), was limited in distinguishing between subtypes with subtle morphological differences. Similarly, genomics-only models (e.g., Random Forest with 89.0% accuracy) struggled with the inherent noise and high dimensionality of molecular data.

**Synergistic advantage:** By fusing features from imaging, genomics, and clinical data, the multimodal model captured a more holistic representation of the tumor phenotype and genotype, leading to a significant increase in accuracy (94.3%) and AUC-ROC (0.951).

This outcome validates existing literature emphasizing the value of radio genomics in oncology, but our work uniquely applies this integration specifically to sarcoma, a largely underserved domain in AI research.

## Challenges and Limitations

While the study presents promising results, several challenges remain:

**Data scarcity:** Sarcoma datasets remain limited in size and diversity. The 300 patient samples used, though substantial for sarcoma research, are still small by deep learning standards.

**Subtype imbalance:** Certain rare sarcoma subtypes were underrepresented, potentially affecting model generalization.

**Data heterogeneity:** Variability in imaging protocols, sequencing platforms, and clinical data recording could introduce biases.

**External validation:** The model's performance needs to be tested on independent, multi-institutional cohorts to assess its real-world applicability.

## Conclusion

In conclusion, this study introduces a novel AI-driven multimodal pipeline for the early detection of sarcoma by integrating radiological imaging, genomic profiles, and clinical metadata. The proposed model achieved superior predictive performance, with an accuracy of 94.3% and an AUC-ROC of 0.951, demonstrating its robustness in distinguishing sarcoma subtypes. Moreover, the incorporation of Grad-CAM visualizations provided enhanced interpretability by highlighting critical regions in medical images, while the use of attention mechanisms enabled more effective feature prioritization across modalities. These results highlight the significant potential of artificial intelligence in overcoming diagnostic challenges associated with rare cancers such as sarcoma, which have historically received limited attention in the broader field of AI-based oncology research. The integrative approach presented here not only advances precision diagnostics but also lays the groundwork for future extensions into prognosis and treatment response prediction.

## Future Work

Future research will focus on expanding the dataset through collaborations with global sarcoma registries to enhance subtype representation and diversity. Real-time clinical deployment is envisioned via user-friendly interfaces integrated into hospital systems. Incorporating 3D imaging and longitudinal data will improve temporal and spatial modeling of tumor progression. Advanced explainability tools like SHAP and LIME will be explored to strengthen model transparency. Additionally, the framework will be extended to support prognosis, metastasis risk prediction, and therapy response assessment, advancing personalized sarcoma care.

By addressing these future directions, this work aims to contribute significantly to the field of AI-assisted precision medicine in rare cancers, ultimately improving sarcoma patient outcomes.

## References:

1. Alberto, C., Aguirre, C., & Cardenas, C. A. , 2023. Sarcomas: A Comprehensive Review of Classification, Diagnosis, Treatment, and Psychosocial Aspects. [https://doi.org/10.4172/cocr.6\(6\).295](https://doi.org/10.4172/cocr.6(6).295)
2. Andreas Mueller, A., Spinnato, P., Italiano, A., Brisse, H. J., Feydy, A., & Crombé, A., et al. 2023. Radiomics and artificial intelligence for soft-tissue sarcomas: current status and perspectives. *Diagnostic and Interventional Imaging*, 104(12), 567–583. <https://doi.org/10.1016/j.diii.2023.09.005>
3. Cassalia, F., Cavallin, F., Danese, A., del Fiore, P., di Prata, C., Rastrelli, M., Belloni Fortina, A., & Mocellin, S., 2023. Soft Tissue Sarcoma Mimicking Melanoma: A Systematic Review. In *Cancers* (Vol. 15, Issue 14). Multidisciplinary Digital Publishing Institute (MDPI). <https://doi.org/10.3390/cancers15143584>
4. Foersch, S., Eckstein, M., Wagner, D. C., Gach, F., Woerl, A. C., Geiger, J., Glasner, C., Schelbert, S., Schulz, S., Porubsky, S., Kreft, A., Hartmann, A., Agaimy, A., & Roth, W., 2021. Deep learning for diagnosis and survival prediction in soft tissue sarcoma. *Annals of Oncology*, 32(9), 1178–1187. <https://doi.org/10.1016/j.annonc.2021.06.007>
5. Navarro, F., Dapper, H., Asadpour, R., et al. 2021. Development and external validation of deep-learning-based tumor grading models in soft-tissue sarcoma patients using MR imaging. *Cancers (Basel)*, 13, 2866. <https://doi.org/10.3390/cancers13122866>
6. van IJzendoorn, D. G. P., Szuhai, K., Briaire-De Bruijn, I. H., Kostine, M., Kuijjer, M. L., & Bovée, J. V. M. G., 2019. Machine learning analysis of gene expression data reveals novel diagnostic and prognostic biomarkers and identifies therapeutic targets for soft tissue sarcomas. *PLoS Computational Biology*, 15(2). <https://doi.org/10.1371/journal.pcbi.1006826>
7. Liang, H.-Y., Yang, S.-F., Zou, H.-M., et al. 2022. Deep learning radiomics nomogram to predict lung metastasis in soft-tissue sarcoma: a multi-center study. *Frontiers in Oncology*, 12, 897676. <https://doi.org/10.3389/fonc.2022.897676>
8. Soomers, V., Husson, O., Young, R., Desar, I., & Van Der Graaf, W. ,2020. The sarcoma diagnostic interval: a systematic review on length, contributing factors and patient outcomes. *ESMO Open*, 5(1), e000592. <https://doi.org/10.1136/esmoopen-2019-000592>
9. Crombé, A., Périer, C., & Kind, M. 2023. T2-based MRI delta-radiomics improve response prediction in soft-tissue sarcomas treated by neoadjuvant chemotherapy. *Journal of Magnetic Resonance Imaging*, 50, 497–510. <https://doi.org/10.1002/jmri.297>
10. Zeissig, S. R., Emrich, K., Reinwald, F., Kasper, B., Kleihues-Van Tole, K., Justenhoven, C., Wardelmann, E., & Hohenberger, P. , 2023. Sarcoma Research with Cancer Registry Data: Data and Peculiarities of Germany in the Light of Other Countries. *Oncology Research and Treatment*, 46(9), 370–381. <https://doi.org/10.1159/000531724>
11. Trautmann, F., Schuler, M., & Schmitt, J. ,2015. Burden of soft-tissue and bone sarcoma in routine care. *Cancer Epidemiology*, 39(3), 440–446. <https://doi.org/10.1016/j.canep.2015.03.002>

12. Comprehensive and integrated genomic characterization of adult soft tissue sarcomas., 2017. *Cell*, 171(4), 950-965.e28. <https://doi.org/10.1016/j.cell.2017.10.014>
13. Zambo, I., & Veselý, K., 2014. [WHO classification of tumours of soft tissue and bone 2013: the main changes compared to the 3rd edition]. *PubMed*, 50(2), 64–70. <https://pubmed.ncbi.nlm.nih.gov/24758500>
14. Qi, L., Chen, F., Wang, L., Yang, Z., Zhang, W., & Li, Z., 2023. Deciphering the role of NETosis-related signatures in the prognosis and immunotherapy of soft-tissue sarcoma using machine learning. *Frontiers in Pharmacology*, 14. <https://doi.org/10.3389/fphar.2023.1217488>
15. Xu, W., Hao, D., Hou, F., Zhang, D., & Wang, H., 2020. Soft tissue sarcoma: preoperative MRI-Based radiomics and machine learning may be accurate predictors of histopathologic grade. *American Journal of Roentgenology*, 215(4), 963–969. <https://doi.org/10.2214/ajr.19.22147>
16. Morisi, A., Rai, T., Bacon, N. J., Thomas, S. A., Bober, M., Wells, K., Dark, M. J., Aboellail, T., Bacci, B., & La Ragione, R. M., 2023. Detection of necrosis in digitised Whole-Slide images for better grading of canine Soft-Tissue sarcomas using Machine-Learning. *Veterinary Sciences*, 10(1), 45. <https://doi.org/10.3390/vetsci10010045>
17. Schlemmer, M., Reichardt, P., Verweij, J., Hartmann, J., Judson, I., Thys, A., Hogendoorn, P., Marreaud, S., Van Glabbeke, M., & Blay, J., 2008. Paclitaxel in patients with advanced angiosarcomas of soft tissue: A retrospective study of the EORTC soft tissue and bone sarcoma group. *European Journal of Cancer*, 44(16), 2433–2436. <https://doi.org/10.1016/j.ejca.2008.07.037>
18. Koelsche, C., Schrimpf, D., Stichel, D., Sill, M., Sahm, F., Reuss, D. E., Blattner, M., Worst, B., Heilig, C. E., Beck, K., Horak, P., Kreutzfeldt, S., Paff, E., Stark, S., Johann, P., Selt, F., Ecker, J., Sturm, D., Pajtler, K. W., . . . Hohenberger, P., 2021. Sarcoma classification by DNA methylation profiling. *Nature Communications*, 12(1). <https://doi.org/10.1038/s41467-020-20603-4>
19. Petitprez, F., De Reyniès, A., Keung, E. Z., Chen, T. W., Sun, C., Calderaro, J., Jeng, Y., Hsiao, L., Lacroix, L., Bougoüin, A., Moreira, M., Lacroix, G., Natario, I., Adam, J., Lucchesi, C., Laizet, Y., Toulmonde, M., Burgess, M. A., Bolejack, V., . . . Fridman, W. H., 2020. B cells are associated with survival and immunotherapy response in sarcoma. *Nature*, 577(7791), 556–560. <https://doi.org/10.1038/s41586-019-1906-8>
20. Tawbi, H. A., Burgess, M., Bolejack, V., Van Tine, B. A., Schuetze, S. M., Hu, J., D'Angelo, S., Attia, S., Riedel, R. F., Priebat, D. A., Movva, S., Davis, L. E., Okuno, S. H., Reed, D. R., Crowley, J., Butterfield, L. H., Salazar, R., Rodriguez-Canales, J., Lazar, A. J., . . . Patel, S., 2017. Pembrolizumab in advanced soft-tissue sarcoma and bone sarcoma (SARC028): a multicentre, two-cohort, single-arm, open-label, phase 2 trial. *The Lancet Oncology*, 18(11), 1493–1501. [https://doi.org/10.1016/s1470-2045\(17\)30624](https://doi.org/10.1016/s1470-2045(17)30624)

21. Rajput, V., Mulay, P. and Mahajan, C.M., 2025. "Bio-inspired algorithms for feature engineering: analysis, applications and future research directions", Information Discovery and Delivery, Vol. 53 No. 1, pp. 56-71. <https://doi.org/10.1108/IDD-11-2022-0118>
22. The Cancer Genome Atlas (TCGA) - <https://www.cancer.gov/tcga>
23. The Cancer Imaging Archive (TCIA) - <https://www.cancerimagingarchive.net>