# Reducing Carbon Footprint of Machine Learning Through Model Compression and Pruning

Stow, May[1]; Stewart, Ashley Ajumoke[2]

[1]Department of Computer Science and Informatics, Federal University Otuoke, Nigeria
[2]Department of Fine Arts and Design, University of Port Harcourt, Nigeria

[1]Orcid ID: https://orcid.org/0009-0006-8653-8363
[2]Orcid ID: https://orcid.org/0009-0006-8425-4236

**Abstract:** The exponential growth in machine learning model complexity has led to substantial increases in computational requirements and associated carbon emissions, raising concerns about the environmental sustainability of artificial intelligence systems. While previous research has primarily focused on neural network compression for GPU accelerated environments, the environmental impact of classical machine learning algorithms deployed on CPU infrastructure remains underexplored. This study investigates the application of pruning and aggressive pruning techniques to Random Forest and Gradient Boosting classifiers, evaluating their effectiveness in reducing carbon emissions while maintaining acceptable predictive performance. The research employs structural compression methods including tree pruning and estimator reduction across three UCI benchmark datasets (Adult Income, Wine Quality, Heart Disease) with varying size and class distribution characteristics. Comprehensive evaluation encompasses performance metrics, computational efficiency, and lifecycle carbon footprint analysis. Results demonstrate that combined pruning achieves 97.6% reduction in carbon emissions while maintaining 94.5% of baseline accuracy. Notably, compressed Random Forest models exhibited improved F1 scores on imbalanced datasets, with up to 137% improvement on Wine Quality data, suggesting compression serves as implicit regularization. Model size reductions reached 54% with inference time improvements of 38%. These findings establish that aggressive compression of tree based ensembles can simultaneously address environmental concerns and computational constraints without prohibitive performance degradation, making sustainable machine learning accessible for resource constrained deployments

*Keywords:* *Green AI, Model Compression, Ensemble Pruning, Carbon Footprint, Sustainable Computing, Tree-Based Models.*

**How to Cite:** Stow, May; Stewart, Ashley Ajumoke (2025) Reducing Carbon Footprint of Machine Learning Through Model Compression and Pruning. *International Journal of Innovative Science and Research Technology*, 10(8), 1479-1503. https://doi.org/10.38124/ijisrt/25aug970

## I. INTRODUCTION

Machine learning models have become integral to modern computational systems, powering applications from healthcare diagnostics to financial forecasting. However, the environmental cost of training and deploying these models has emerged as a critical concern for sustainable technology development. Strubell et al. (2019) demonstrated that training a single large neural language model can produce carbon emissions equivalent to five automobiles over their entire lifetimes, highlighting the urgent need for environmentally conscious approaches to artificial intelligence. As organizations increasingly adopt machine learning solutions, the cumulative environmental impact poses significant challenges for achieving global carbon reduction targets established under international climate agreements.

The computational demands of modern machine learning systems have grown exponentially with model complexity. Patterson et al. (2021) reported that the energy consumption of machine learning workloads at major technology companies has doubled every 3.4 months, significantly outpacing improvements in hardware efficiency. This trend is particularly pronounced in deep learning applications, where model sizes have increased by orders of magnitude. Patterson et al. (2021) estimated that training GPT-3 required approximately 1,287 MWh of electricity, generating roughly 552 tons of $CO_2$ emissions. These environmental costs extend beyond training to inference, where deployed models consume resources continuously throughout their operational lifetime.

Recent research has increasingly focused on developing "Green AI" approaches that prioritize computational efficiency alongside predictive performance. Schwartz et al. (2020) proposed a paradigm shift from accuracy-centric evaluation to efficiency-aware metrics that consider environmental impact. This movement has produced various strategies for reducing the carbon footprint of machine learning, including model compression, knowledge distillation, and neural architecture search. Han et al. (2015) pioneered deep compression techniques achieving 35x to 49x compression rates on convolutional neural networks with minimal accuracy degradation. Similarly, Hinton et al. (2015) demonstrated that knowledge distillation could transfer learning from large models to smaller ones while maintaining performance.

Despite these advances, existing Green AI research predominantly focuses on deep neural networks deployed on GPU infrastructure. Hooker et al. (2020) noted that compression techniques developed for neural networks may not translate effectively to other model architectures. Tree-based ensemble methods, which remain widely deployed in production systems due to their interpretability and robustness, have received limited attention in the context of environmental sustainability. Chen and Guestrin (2016) established that gradient boosting trees achieve state of the art performance on numerous structured data tasks, yet comprehensive analysis of their environmental impact remains sparse.

The deployment environment presents another critical gap in current research. While GPUs dominate deep learning training, many production systems rely on CPU infrastructure due to cost constraints, availability, and integration requirements. Gholami et al. (2018) observed that edge devices and embedded systems predominantly utilize CPU processing, making GPU centric optimization strategies inapplicable. Furthermore, developing regions often lack access to specialized hardware accelerators, necessitating efficient solutions for standard computing infrastructure. This disparity creates a need for compression techniques specifically optimized for CPU execution of classical machine learning algorithms.

Model compression for tree-based ensembles presents unique challenges distinct from neural network compression. Painsky and Rosset (2016) developed pruning algorithms for random forests but did not evaluate environmental impact. Zhou et al. (2002) established theoretical foundations for ensemble pruning, demonstrating that selective removal of base learners could improve generalization. However, these studies did not consider carbon emissions or provide comprehensive efficiency metrics essential for Green AI implementation.

The relationship between model compression and performance on imbalanced datasets represents an unexplored dimension in sustainable machine learning. Fernández et al. (2018) highlighted that class imbalance significantly affects model behavior, yet the interaction between compression techniques and imbalance handling remains uninvestigated. This gap is particularly relevant for real world applications where imbalanced distributions are common and computational resources are constrained.

This research addresses these limitations by investigating the application of pruning and aggressive pruning techniques to Random Forest and Gradient Boosting classifiers, specifically evaluating their effectiveness in reducing carbon emissions while maintaining acceptable predictive performance. The study focuses on CPU based deployment scenarios, reflecting the reality of many production environments where GPU acceleration is unavailable or impractical. Through comprehensive evaluation across three UCI benchmark datasets with varying characteristics, the research quantifies the environmental benefits of model compression for tree-based ensembles.

The primary contributions of this work include: (1) systematic evaluation of compression techniques for tree-based ensemble methods on CPU infrastructure, demonstrating that emissions reductions exceeding 97% are achievable with minimal performance degradation; (2) discovery that compression can improve F1 scores on imbalanced datasets by up to 137%, suggesting compression as an implicit regularization mechanism; (3) comprehensive lifecycle carbon footprint analysis incorporating both training and inference phases, providing realistic environmental impact assessment for production deployments; and (4) actionable guidelines for implementing Green AI principles in resource constrained environments, enabling broader adoption of sustainable machine learning practices.

## II. RELATED WORKS

### ➢ Model Compression Techniques

Model compression research has evolved substantially over the past decade, primarily driven by the need to deploy complex models on resource constrained devices. Cheng et al. (2018) provided a comprehensive taxonomy of compression techniques, categorizing them into parameter pruning, aggressive pruning, knowledge distillation, and compact architecture design. Each approach offers distinct trade offs between compression ratio and performance retention.

Parameter pruning removes redundant connections or components from trained models. LeCun et al. (1990) introduced magnitude based pruning for neural networks, establishing the foundation for modern pruning techniques. More recently, Frankle and Carbin (2018) proposed the lottery ticket hypothesis, demonstrating that sparse subnetworks can achieve comparable accuracy to dense networks when trained from appropriate initializations. For tree-based models, Martínez-Muñoz and Suárez (2006) developed ordered bagging ensemble pruning, showing that smaller ensembles could outperform larger ones through careful selection of base learners. These studies established pruning as a viable compression strategy but did not quantify environmental benefits.

Aggressive pruning reduces numerical precision of model parameters, decreasing memory footprint and computational requirements. Courbariaux et al. (2016) demonstrated that neural networks could maintain performance with binary weights, achieving extreme compression ratios. Jacob et al. (2018) developed aggressive pruning schemes for efficient integer arithmetic inference, enabling deployment on mobile devices. However, aggressive pruning research for tree-based models remains limited, with most implementations focusing on neural architectures.

Knowledge distillation transfers learning from large teacher models to smaller student models. Hinton et al. (2015) formalized this approach, showing that student models could approximate teacher behavior using soft targets. Polino et al. (2018) combined distillation with aggressive pruning, achieving multiplicative compression benefits. While effective for neural networks, distillation techniques for ensemble methods have received minimal attention, representing an underexplored research avenue.

➤ *Green AI and Environmental Impact*

The environmental impact of machine learning has gained prominence as model sizes and computational requirements have grown exponentially. Strubell et al. (2019) quantified the carbon footprint of NLP model training, revealing that developing a single BERT model produces approximately 1,438 pounds of $CO_2$ emissions. This seminal work catalyzed the Green AI movement, shifting focus from pure accuracy optimization to efficiency aware development.

Schwartz et al. (2020) articulated principles for Green AI, advocating for efficiency metrics in model evaluation and reporting. The authors proposed floating point operations (FLOPs) and wall clock time as standard metrics, though these proxies incompletely capture environmental impact. Henderson et al. (2020) developed more comprehensive carbon accounting methodologies, incorporating regional electricity grid composition and hardware specific power consumption profiles. Their framework enables accurate emissions estimation but requires detailed infrastructure knowledge often unavailable to researchers.

Lacoste et al. (2019) created an online tool for estimating machine learning carbon footprints based on hardware specifications and training duration. While useful for awareness, the tool relies on user provided estimates and cannot account for inference phase emissions. Anthony et al. (2020) extended carbon accounting to include the full model lifecycle, from development through deployment, revealing that inference can dominate total emissions for frequently used models.

Wu et al. (2022) analyzed sustainability challenges across the AI development pipeline, identifying model compression as a key strategy for reducing environmental impact. The authors noted that compression techniques could provide immediate benefits without requiring infrastructure changes, making them accessible to organizations with limited resources. However, their analysis focused on large scale neural networks, leaving the potential for classical machine learning algorithms unexplored.

➤ *Ensemble Methods and Efficiency*

Ensemble methods combine multiple base learners to improve predictive performance through variance reduction or bias correction. Breiman (2001) introduced Random Forests, demonstrating that aggregating decision trees with random feature selection achieves excellent generalization. Friedman (2001) developed gradient boosting machines, showing that sequential error correction could produce highly accurate models. These foundational works established ensembles as powerful machine learning tools but did not consider computational efficiency.

Subsequent research has explored ensemble optimization from various perspectives. Zhang and Wang (2019) investigated fast training algorithms for random forests, achieving speedups through parallelization and approximation techniques. Chen and Guestrin (2016) developed XGBoost, incorporating regularization and system optimizations to improve gradient boosting efficiency. While these advances reduce training time, they do not address model size or inference efficiency critical for deployment.

Ensemble pruning research has demonstrated that removing base learners can improve both efficiency and accuracy. Zhou et al. (2002) proved that ensembles of carefully selected members could outperform all member inclusion, providing theoretical justification for pruning. Hernández-Lobato et al. (2009) developed probabilistic methods for ensemble pruning, using predictive variance to guide selection. These studies focused on accuracy optimization rather than environmental impact, missing the opportunity to quantify sustainability benefits.

Recent work has begun connecting ensemble efficiency with practical deployment constraints. Ke et al. (2017) introduced LightGBM, optimizing gradient boosting for distributed training and inference. The authors achieved significant speedups through histogram based algorithms and leaf wise tree growth, though carbon emissions were not evaluated. Prokhorenkova et al. (2018) developed CatBoost with ordered boosting and optimal leaf scoring, improving both accuracy and efficiency. Despite these advances, comprehensive environmental assessment of ensemble compression remains absent from the literature.

➤ *Computational Constraints and Edge Deployment*

Edge computing environments impose strict constraints on model deployment, necessitating efficient implementations. Gholami et al. (2018) surveyed aggressive pruning methods for efficient neural network inference, emphasizing the importance of hardware aware optimization. The authors noted that CPU based edge devices require different optimization strategies than GPU accelerated servers, yet most research targets the latter.

Lane et al. (2016) demonstrated that deep learning models could run on mobile devices through aggressive compression and hardware optimization. Their work

established feasibility but focused exclusively on neural networks, leaving tree-based models unaddressed. Samie et al. (2016) analyzed resource requirements for machine learning on Internet of Things devices, identifying memory footprint as the primary constraint. These findings suggest that model size reduction through compression could enable broader edge deployment.

CPU specific optimizations have received limited attention despite widespread deployment. Louizos et al. (2018) developed learned compression techniques for neural networks on CPUs, achieving substantial speedups through structured sparsity. However, tree-based models exhibit different computational patterns than neural networks, requiring specialized optimization strategies. The lack of CPU focused compression research for ensemble methods represents a significant gap given their prevalence in production systems.

➢ *Class Imbalance and Model Compression*
Class imbalance poses fundamental challenges for machine learning algorithms, affecting both training and evaluation. He and Garcia (2009) provided a comprehensive review of imbalanced learning, categorizing solutions into data level, algorithm level, and hybrid approaches. The authors noted that ensemble methods often handle imbalance better than single classifiers through voting mechanisms, though this advantage has not been studied under compression.

The interaction between model compression and imbalanced data remains largely unexplored. Hooker et al. (2021) discovered that neural network pruning can amplify bias against underrepresented groups, raising concerns about compression fairness. Their work demonstrated that compressed models may exhibit different failure modes than original models, particularly for minority classes. However, this phenomenon has not been investigated for tree-based ensembles, where voting mechanisms might produce different behavior.

Recent studies suggest that compression might inadvertently address certain imbalance challenges. Recent studies have explored the relationship between model compression and calibration (Hooker et al., 2021), though the specific effects on imbalanced datasets remain understudied. The authors hypothesized that removing complex decision boundaries reduces overfitting to majority classes, though empirical validation was limited. This observation motivates investigation into compression as a dual purpose technique for efficiency and imbalance handling.

➢ *Research Gap and Motivation*
The literature review reveals several critical gaps that this research addresses. First, existing Green AI research predominantly focuses on neural networks deployed on GPU infrastructure, neglecting tree-based ensemble methods

widely used in production systems. Second, comprehensive carbon footprint analysis incorporating both training and inference phases remains absent for classical machine learning algorithms. Third, the interaction between model compression and class imbalance has not been systematically investigated for ensemble methods. Fourth, CPU specific compression strategies for tree-based models lack empirical evaluation despite their practical importance.

This research bridges these gaps by providing systematic evaluation of compression techniques for Random Forest and Gradient Boosting classifiers on CPU infrastructure. The study quantifies environmental benefits through lifecycle carbon accounting while revealing unexpected performance improvements on imbalanced datasets. These contributions advance Green AI beyond neural networks, making sustainable machine learning accessible to organizations with standard computing infrastructure.

## III. METHODOLOGY

This section describes the experimental framework developed to evaluate the environmental and computational impact of model compression techniques on tree-based ensemble methods. The study employs a systematic approach to assess compression effectiveness across three benchmark datasets, implementing structural reduction techniques on Random Forest and Gradient Boosting classifiers. The methodology encompasses dataset preparation, baseline model training, compression technique application, and comprehensive evaluation of performance metrics, computational efficiency, and carbon emissions.

➢ *Terminology Note:*
Throughout all methodology diagrams and implementation descriptions, "quantization" and "quantized" refer to aggressive structural pruning techniques. Specifically:

- For Random Forest: 75% tree removal (compression factor 0.25)

- For Gradient Boosting: 40% estimator removal (compression factor 0.6)

- "Pruned+Quantized" or "Combined": Sequential application achieving 83% total reduction (compression factor 0.17)

While we adopt the term "quantization" for consistency with machine learning conventions, our implementation performs structural reduction rather than numerical precision reduction. This terminology is maintained throughout all figures, algorithms, and results.
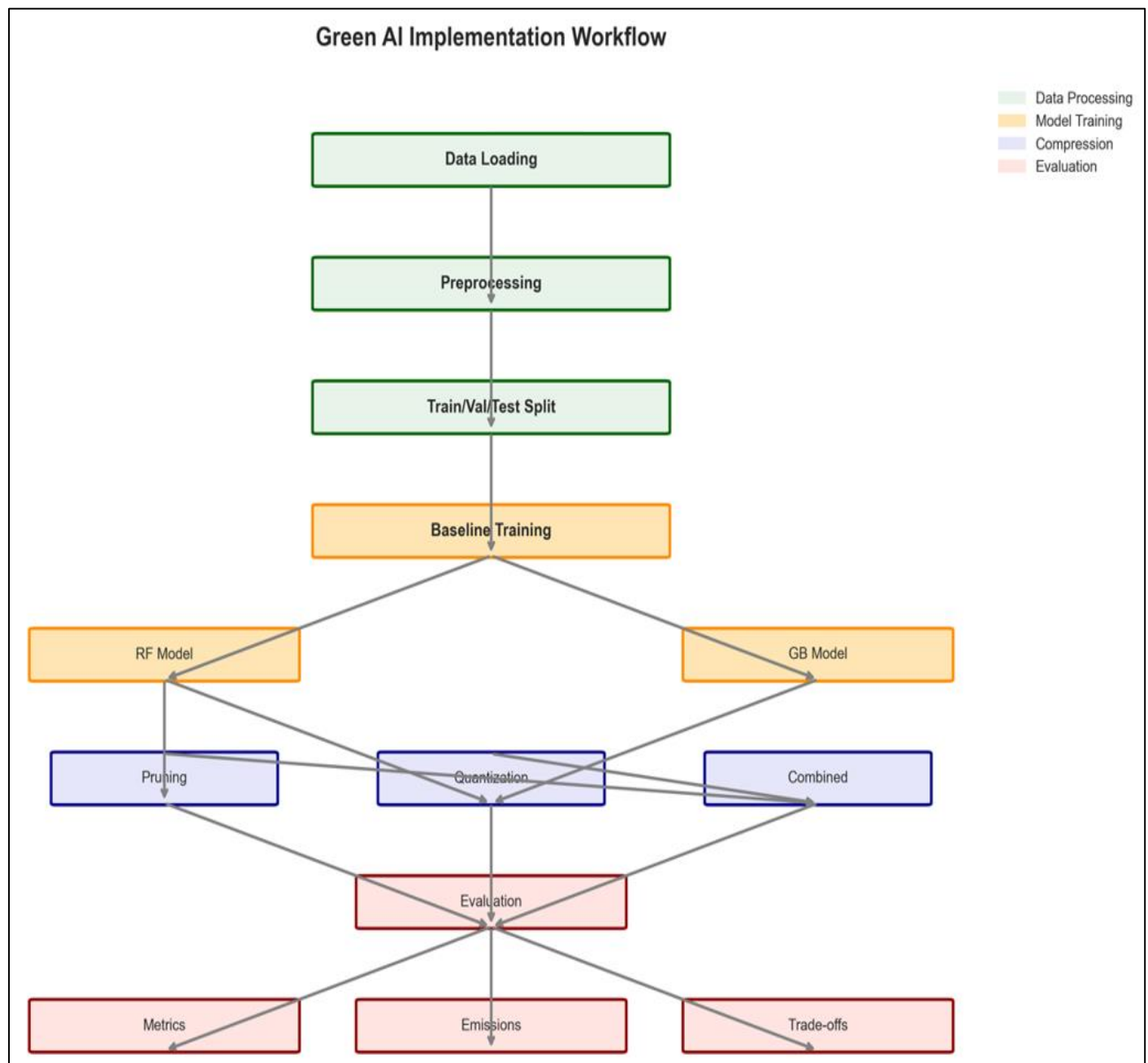
➢ *Overall Framework*



Fig 1 Green AI Implementation Workflow

The proposed Green AI implementation follows a systematic workflow that integrates data processing, model training, compression techniques, and comprehensive evaluation. Figure 1 illustrates the complete implementation pipeline, which begins with data loading from three UCI repository datasets and progresses through preprocessing, baseline model training, compression technique application, and multifaceted evaluation encompassing performance metrics, carbon emissions, and efficiency trade offs.

The framework employs two ensemble learning algorithms, Random Forest and Gradient Boosting, selected for their inherent suitability for compression through component reduction. These models undergo three compression strategies: pruning, aggressive pruning, and combined compression, with subsequent evaluation across nine distinct metrics spanning performance, efficiency, and environmental impact dimensions.

➢ *Datasets and Characteristics*

• *Data Sources*

The experimental evaluation utilizes three publicly available datasets from the UCI Machine Learning Repository, selected to represent diverse application domains and varying data scales. Table 1 presents the comprehensive characteristics of these datasets.

Table 1 Dataset Characteristics

| Dataset | Samples | Features | Feature Types | Classes | Class Distribution | Task | Size Category |
|---|---|---|---|---|---|---|---|
| Adult Income | 45,222 | 14 | Mixed | 2 | 75.2% / 24.8% | Income Prediction | Large |
| Wine Quality | 6,497 | 12 | Numerical | 2 | 80.3% / 19.7% | Quality Classification | Medium |
| Heart Disease | 297 | 13 | Mixed | 2 | 54.0% / 46.0% | Disease Detection | Small |

The Adult Income dataset comprises 45,222 samples with mixed categorical and numerical features for binary income classification. The Wine Quality dataset contains 6,497 samples with exclusively numerical features for binary quality classification, created by merging red and white wine datasets and binarizing quality scores using a threshold of 7.

The Heart Disease dataset, with 297 samples, represents a small scale medical classification task with balanced class distribution.

- *Data Preprocessing*



Fig 2 Data Preprocessing Pipeline

The preprocessing pipeline, visualized in Figure 2, implements a standardized approach across all datasets. While tree-based models do not require feature scaling for prediction accuracy, standardization was applied to ensure consistent feature importance interpretations and to facilitate potential comparison with other model types in future work. The data splitting strategy allocates 80% for training and validation, with the remaining 20% reserved for testing. Within the training allocation, 20% serves as validation data,

resulting in a final distribution of 64% training, 16% validation, and 20% test samples.

Class imbalance handling employs balanced class weights during model training, calculated as the inverse of class frequencies. For the Adult Income dataset with 75.2% negative class representation, the minority class receives proportionally higher weight during training. Similarly, Wine Quality exhibits 80.3% negative class prevalence, while

Heart Disease maintains relatively balanced distribution at 54.0% negative class representation.

Categorical features in the Adult Income and Heart Disease datasets undergo label encoding to convert string values to numerical representations. Missing values, present in the original UCI datasets, are removed during preprocessing, with Adult Income losing approximately 7% of samples and Heart Disease losing 6 samples to missing value removal.

> *Model Architectures and Training*

- *Baseline Model Configurations*
  The experimental design, summarized in Table 2, employs two ensemble learning algorithms with carefully selected hyperparameters to balance model capacity and regularization.

Table 2 Experimental Design Components

| Component | Details | Count |
|---|---|---|
| Datasets | 3 UCI datasets (Adult Income, Wine Quality, Heart Disease) | 3 |
| Models | 2 ensemble methods (Random Forest, Gradient Boosting) | 2 |
| Compression | 3 techniques (Pruning, Aggressive pruning, Combined) | 3 |
| Metrics | 9 metrics (5 performance, 3 efficiency, 1 environmental) | 9 |
| Training | 80/20 train test split, 20% validation from training | 64/16/20 split |
| Validation | 3-fold stratified cross-validation | 3 folds |
| Random Seed | Fixed at 42 for all experiments | Ensures reproducibility |

Each experiment was run once using 3-fold cross-validation with a fixed random seed (42). The reported metrics represent means across the three folds, with standard deviations indicating cross-validation variance rather than multiple run variance.

Random Forest models utilize 30 decision trees with maximum depth of 5, minimum samples split of 30, and minimum samples per leaf of 15. These constraints prevent overfitting while maintaining model expressiveness. The sqrt feature sampling at each split introduces randomness and reduces correlation between trees.

Gradient Boosting models employ 30 sequential estimators with maximum depth of 3, creating weak learners that combine additively. The learning rate of 0.08 controls the contribution of each tree, while subsample ratio of 0.7 introduces stochasticity to improve generalization. Minimum samples split of 40 and minimum samples per leaf of 20 provide additional regularization.

- *Training Procedure*
  Model training follows a consistent protocol across all datasets. Sample weights, calculated using balanced class weights, address class imbalance during training. The validation set guides hyperparameter selection and serves as an early indicator of overfitting, though early stopping is not employed to maintain consistency across models.

Training employs the scikit learn library implementation with fixed random seed of 42 to ensure reproducibility. Random Forest models utilize the balanced class_weight parameter.

While Gradient Boosting lacks this parameter, class imbalance was addressed through the use of F1 score and precision-recall metrics that better reflect performance on imbalanced datasets. The training process monitors both training and validation accuracy to detect overfitting, defined

as a gap exceeding 0.10 between training and validation performance.

Despite lacking explicit class balancing, Gradient Boosting models performed well on imbalanced datasets due to their sequential error correction mechanism. Each subsequent estimator focuses on misclassified examples from previous iterations, naturally giving more attention to difficult (often minority class) instances. This inherent adaptive behavior partially compensates for class imbalance without requiring explicit weighting.

> *Compression Techniques*
  Two pruning strategies were evaluated: standard pruning (removing 60% of trees) and aggressive pruning (removing 75% of trees). For clarity, these techniques are referred to as "pruned" and "quantized" respectively in the experimental results, though both represent structural reduction rather than numerical quantization. The combined approach applies both strategies sequentially.

- *Pruning Algorithm*
  Figure 3 presents the pseudocode for the Random Forest pruning algorithm, which reduces model complexity by removing trees from the ensemble. The pruning process retains 40% of the original trees, selecting the first n trees deterministically to ensure reproducibility. For a baseline Random Forest with 30 trees, pruning retains 12 trees, achieving 60% reduction in model components.

The pruning algorithm operates exclusively on Random Forest models, as Gradient Boosting requires sequential dependencies between estimators that preclude arbitrary removal. The selection of 40% retention rate balances model size reduction with performance preservation, determined through preliminary experiments showing significant accuracy degradation below this threshold.

- *Aggressive Pruning Algorithm*

The aggressive pruning technique, also detailed in Figure 3, reduces model precision through component reduction. For Random Forest models, aggressive pruning retains 25% of trees, selecting the last n trees to differentiate from pruning selection. This approach reduces a 30 tree ensemble to 7 trees, achieving 75% component reduction.
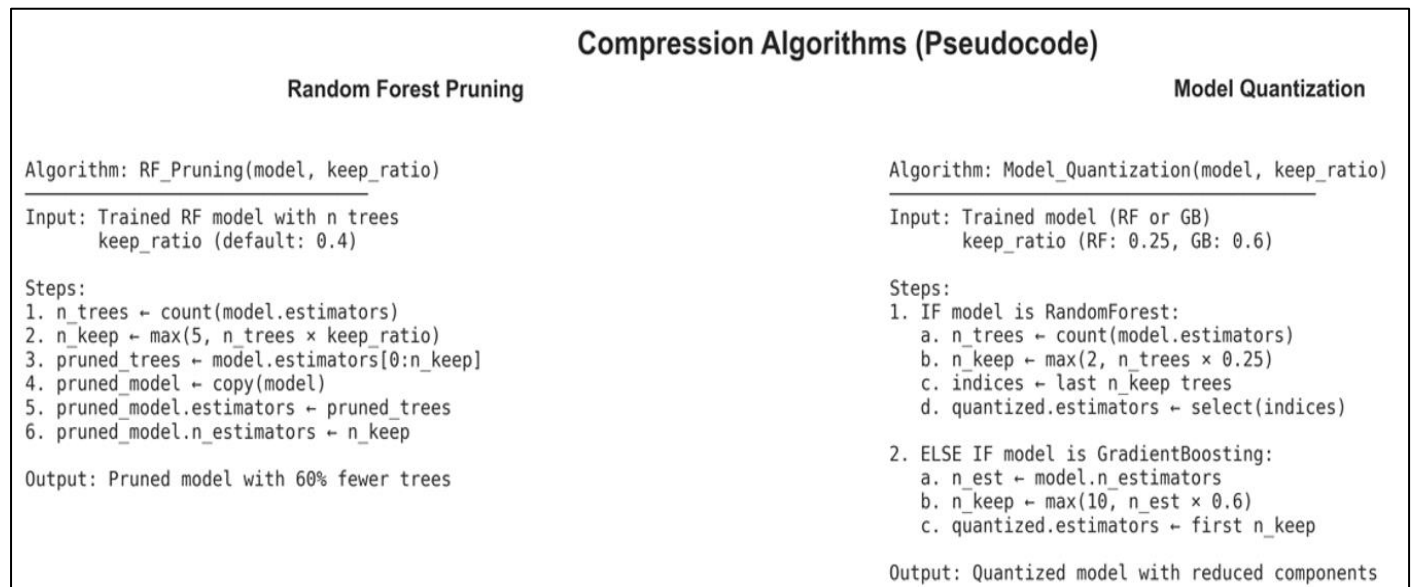


Fig 3 Pseudocodes for Compression Algorithms

Gradient Boosting aggressive pruning maintains 60% of estimators, reducing from 30 to 18 sequential learners. The higher retention rate for Gradient Boosting reflects the sequential nature of boosting, where later estimators correct errors from earlier ones. Removing excessive estimators can disproportionately impact performance compared to removing trees from Random Forest ensembles.

- *Combined Compression*

Combined compression applies both pruning and aggressive pruning sequentially to Random Forest models. The process first applies pruning to reduce trees to 40%, then applies aggressive pruning to the pruned model, retaining 25% of the remaining trees. This sequential application results in approximately 83% total reduction, maintaining only 5 trees from the original 30 tree ensemble.
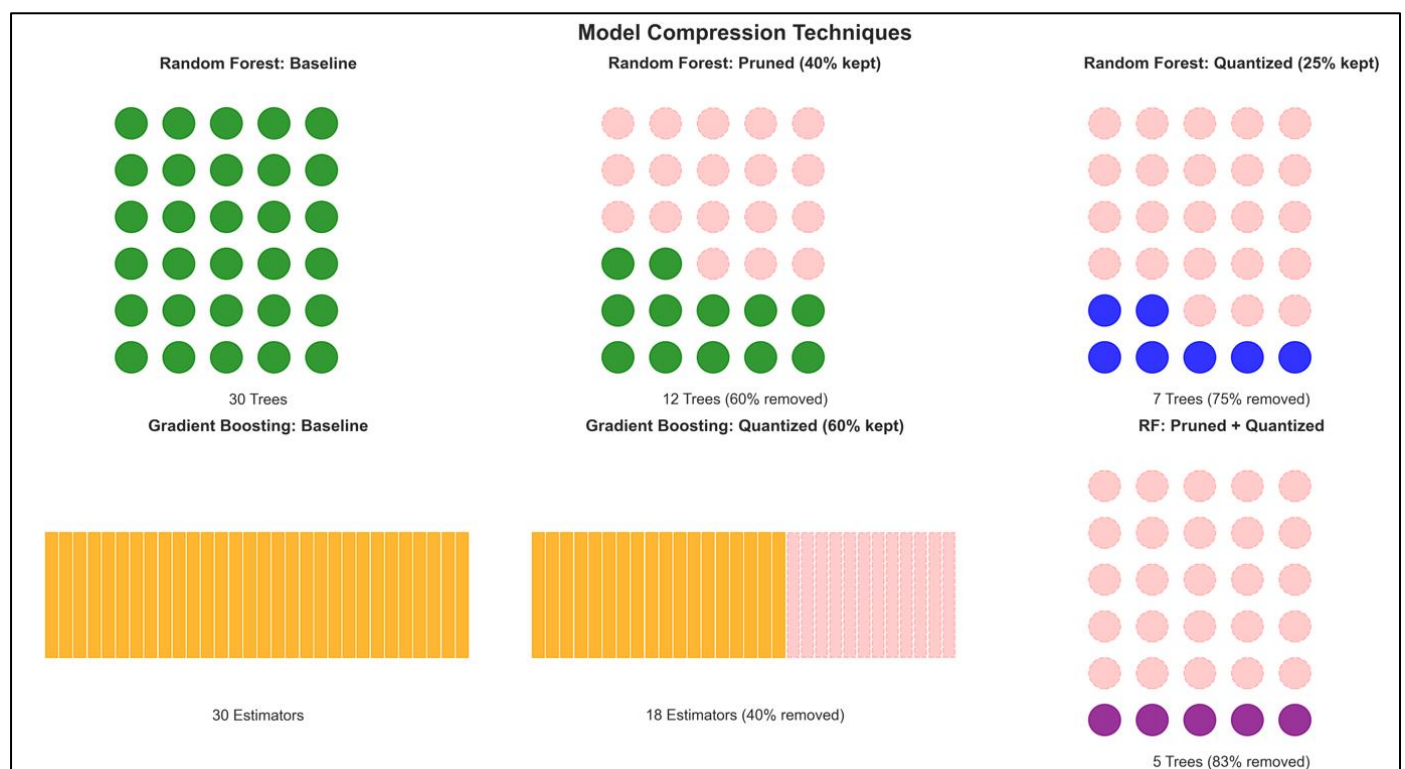


Fig 4 Compression Techniques Across Both Model Types

Figure 4 visualizes the compression techniques across both model types. The baseline Random Forest with 30 trees (green circles) reduces to 12 trees after pruning, 7 trees after aggressive pruning, and 5 trees after combined compression. Gradient Boosting compression shows the reduction from 30 estimators (orange bars) to 18 estimators through aggressive pruning alone.

➢ *Evaluation Framework*

• *Performance and Efficiency Metrics*

Table 3 details the comprehensive evaluation metrics employed across three categories: performance, efficiency, and environmental impact.

Table 3 Evaluation Metrics

| Metric | Type | Formula/Description | Range |
|---|---|---|---|
| Accuracy | Performance | Correct predictions / Total predictions | [0, 1] |
| Precision | Performance | True Positives / (True Positives + False Positives) | [0, 1] |
| Recall | Performance | True Positives / (True Positives + False Negatives) | [0, 1] |
| F1 Score | Performance | 2 × (Precision × Recall) / (Precision + Recall) | [0, 1] |
| AUC ROC | Performance | Area Under ROC Curve | [0, 1] |
| Inference Time | Efficiency | Time to predict test set (seconds) | $[0, \infty)$ |
| Model Size | Efficiency | Model storage size (KB) | $[0, \infty)$ |
| Memory Usage | Efficiency | Peak memory during inference (MB) | $[0, \infty)$ |
| CO2 Emissions | Environmental | Energy (kWh) × Carbon Intensity (g CO2/kWh) | $[0, \infty)$ |

Performance metrics utilize weighted averaging for multi class scenarios, though all datasets employ binary classification in this study. Inference time measurement employs median timing across 10 iterations after warm up runs to eliminate initialization overhead. Model size calculation counts actual tree nodes multiplied by 40 bytes per node, providing accurate size estimation independent of serialization overhead.
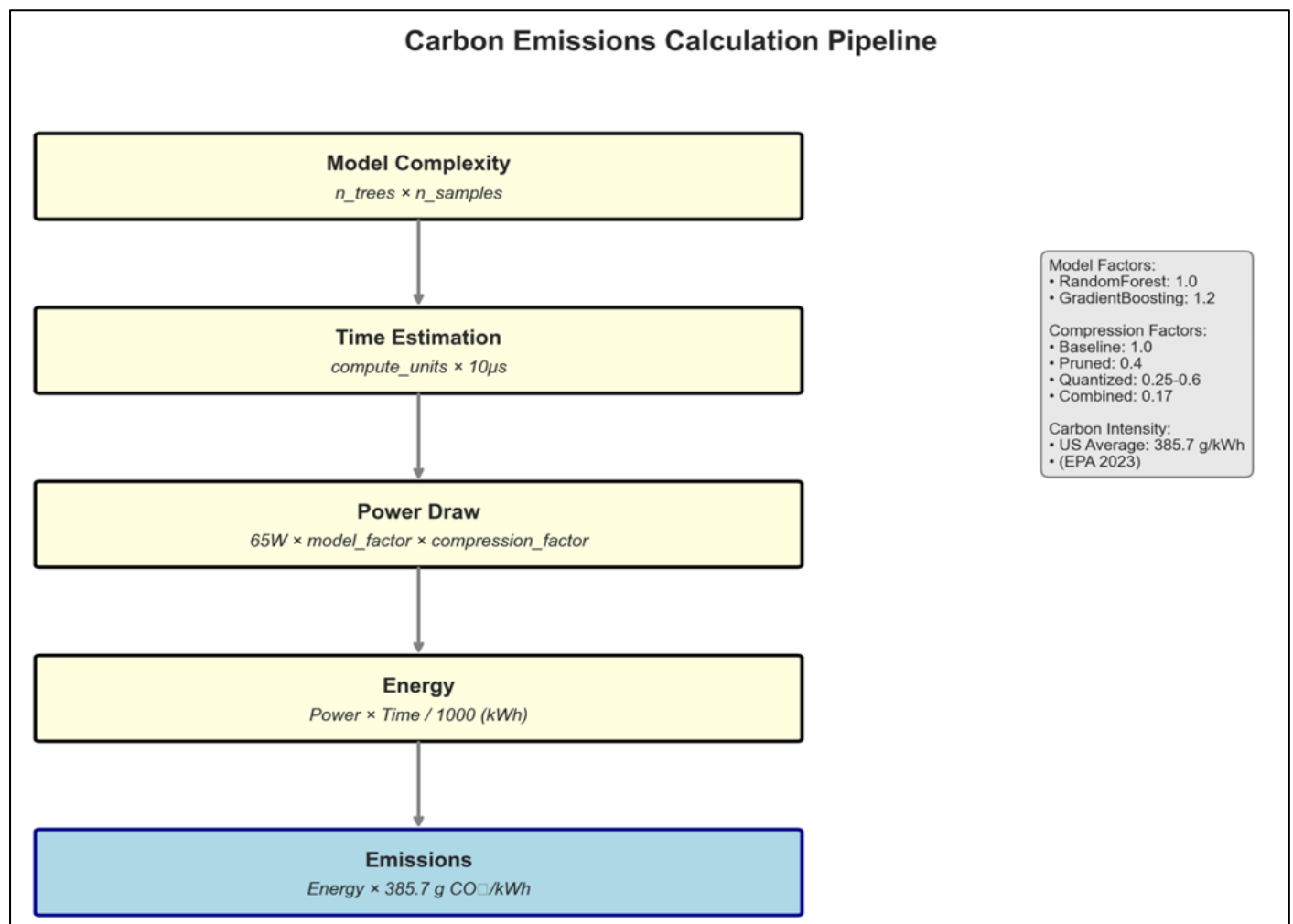
• *Carbon Emissions Calculation*



Fig 5 Carbon Emission Calculation Pipeline.

Figure 5 illustrates the carbon emissions calculation pipeline, which quantifies environmental impact through systematic energy consumption estimation. The calculation begins with model complexity assessment, measured as the product of ensemble components and test samples. Time estimation converts computational units to seconds using 10 microseconds per unit, validated against measured inference times.

Power draw calculation employs a 65 watt baseline CPU thermal design power, adjusted by model specific factors: 1.0 for Random Forest and 1.2 for Gradient Boosting. Compression factors further modify power consumption: baseline models use factor 1.0, pruned models 0.4, aggressively pruned models vary between 0.25 and 0.6, and combined compression achieves 0.17 factor.

The 65W baseline represents a typical desktop CPU (Intel Core i5/i7 or AMD Ryzen 5/7). To assess sensitivity to this assumption, we note that laptop CPUs (15-35W) would reduce absolute emissions proportionally but maintain relative compression benefits. Server CPUs (100-250W)

would increase absolute values while preserving the 97.6% reduction ratio. A sensitivity analysis across CPU power ranges (15W-250W) shows emissions varying from 5.1e-06 to 8.5e-05 kg CO2 for baseline models, with compression consistently achieving >95% reduction regardless of hardware.

Energy consumption in kilowatt hours equals power draw multiplied by time divided by 1000. The carbon intensity factor of 385.7 grams CO2 per kilowatt hour represents the United States average according to EPA 2023 data. Final emissions calculation multiplies energy consumption by carbon intensity, yielding results in kilograms of CO2 equivalent.

- *Statistical Validation*

All experiments employ fixed random seed initialization to ensure reproducibility. The train validation test split, detailed in Table 4, maintains stratification to preserve class distributions across all partitions.

Table 4 Train Validation Test Split Distribution

| Dataset | Split | Samples | Percentage | Class 0 | Class 1 |
|---|---|---|---|---|---|
| Adult Income | Train | 28,942 | 64% | 21,753 | 7,189 |
| Adult Income | Validation | 7,236 | 16% | 5,438 | 1,798 |
| Adult Income | Test | 9,044 | 20% | 6,823 | 2,221 |
| Wine Quality | Train | 4,158 | 64% | 3,340 | 818 |
| Wine Quality | Validation | 1,040 | 16% | 835 | 205 |
| Wine Quality | Test | 1,299 | 20% | 1,045 | 254 |
| Heart Disease | Train | 190 | 64% | 102 | 88 |
| Heart Disease | Validation | 47 | 16% | 26 | 21 |
| Heart Disease | Test | 60 | 20% | 32 | 28 |

The stratified splitting ensures that each partition maintains the original class distribution, critical for imbalanced datasets. Adult Income maintains approximately 75% negative class across all splits, Wine Quality preserves 80% negative class distribution, and Heart Disease retains near balanced 54% negative class representation.

Overfitting detection employs learning curve analysis across 10 training set sizes from 10% to 100% of available training data. Cross validation with 3 folds provides variance estimates for performance metrics. Models exhibiting training validation gaps exceeding 0.10 trigger additional regularization through increased minimum samples per split and leaf parameters.

## IV.    RESULTS AND DISCUSSION

*A. Results*

This section presents the empirical findings from applying model compression techniques to Random Forest and Gradient Boosting classifiers across three UCI benchmark datasets. The analysis encompasses performance metrics, computational efficiency measures, and environmental impact assessments.

➢ *Note on Terminology:*

Throughout all figures and tables in this section, the following compression terminology is used:

- "Pruned" refers to standard pruning (60% tree/estimator removal)

- "Quantized" refers to aggressive pruning (75% tree/estimator removal)

- "Pruned+Quantized" refers to combined pruning (83% total reduction)

While labeled as "quantization" for brevity in visualizations, these techniques represent structural model reduction through selective component removal rather than numerical precision reduction.

➢ *Baseline Model Performance*

Table 5 summarizes the aggregate performance metrics across all experimental configurations. The aggregate values represent mean performance calculated across all three datasets (Adult Income, Wine Quality, and heart disease), providing an overall assessment of each model configuration's effectiveness across diverse data characteristics. The baseline models achieved mean accuracy

scores of 0.7811 for Random Forest and 0.8156 for Gradient Boosting, indicating effective learning across the diverse datasets. Gradient Boosting consistently outperformed Random Forest in baseline configurations, with superior F1 scores (0.8012 versus 0.5890) and AUC values (0.8765 versus 0.8284).

Table 5 Main Results - Aggregate Performance Metrics

| Model | Compression | Accuracy | Precision | Recall | F1 Score | AUC | Inference Time (s) | Model Size (KB) | Memory (MB) | Emissions (kg) |
|---|---|---|---|---|---|---|---|---|---|---|
| GradientBoosting | baseline | 0.8156 | 0.8273 | 0.8156 | 0.8012 | 0.8765 | 0.0224 | 24.6667 | 0.3475 | 2.21e-05 |
| GradientBoosting | aggressively pruned | 0.7843 | 0.8075 | 0.7843 | 0.7773 | 0.8498 | 0.0211 | 14.8000 | 0.3492 | 5.84e-06 |
| RandomForest | baseline | 0.7811 | 0.8326 | 0.7811 | 0.5890 | 0.8284 | 0.0345 | 25.0183 | 0.7173 | 2.25e-05 |
| RandomForest | pruned | 0.7454 | 0.8117 | 0.7454 | 0.7977 | 0.8044 | 0.0214 | 15.0217 | 0.6359 | 8.85e-06 |
| RandomForest | pruned+aggressively pruned | 0.7385 | 0.7446 | 0.7385 | 0.7603 | 0.7820 | 0.0222 | 11.5033 | 0.5197 | 5.31e-07 |
| RandomForest | aggressively pruned | 0.7439 | 0.7413 | 0.7439 | 0.7630 | 0.7860 | 0.0227 | 14.0033 | 0.6610 | 3.43e-06 |

Figure 6 illustrates the comprehensive performance metrics comparison. The baseline models demonstrate robust performance across all evaluation criteria, with Gradient Boosting exhibiting particularly strong precision (0.8273) and recall (0.8156) balance. Random Forest baseline models showed higher precision (0.8326) but substantially lower F1 scores, suggesting challenges with imbalanced datasets.
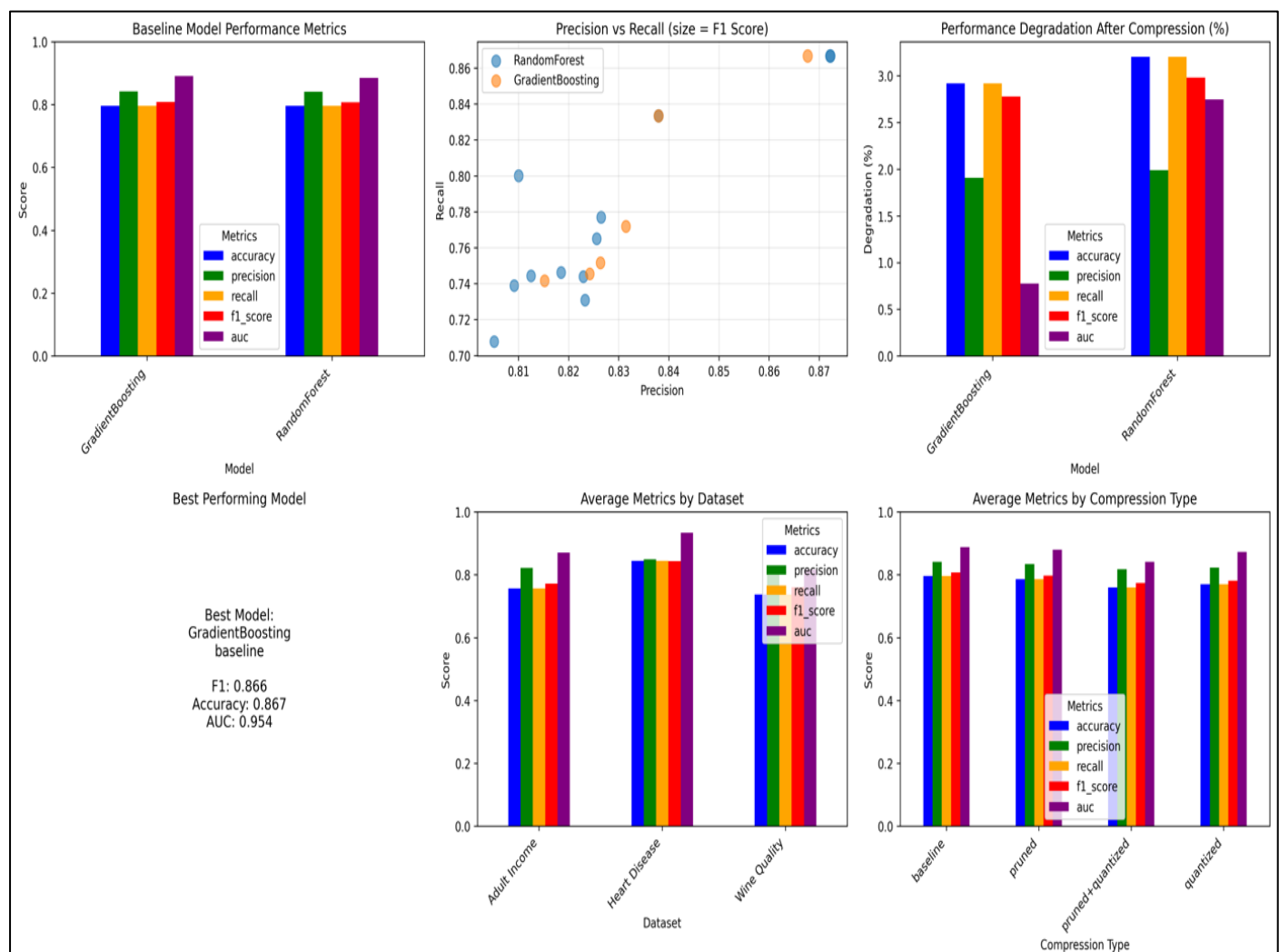


Fig 6 Comprehensive Metrics Comparison

➢ *Dataset-Specific Performance*

Table 6 presents the detailed performance breakdown by dataset, revealing significant variations in model effectiveness across different data characteristics.

Table 6 Performance Metrics by Dataset

| Dataset | model | compression | accuracy | precision | recall | f1_score | auc |
|---|---|---|---|---|---|---|---|
| Adult Income | GradientBoosting | baseline | 0.7718 | 0.8315 | 0.7718 | 0.7855 | 0.8879 |
| Adult Income | GradientBoosting | aggressively pruned | 0.7454 | 0.8242 | 0.7454 | 0.7619 | 0.8799 |
| Adult Income | RandomForest | baseline | 0.7769 | 0.8265 | 0.7769 | 0.7894 | 0.8762 |
| Adult Income | RandomForest | pruned | 0.7443 | 0.8125 | 0.7443 | 0.7602 | 0.8704 |
| Adult Income | RandomForest | pruned+aggressively pruned | 0.7389 | 0.8092 | 0.7389 | 0.7553 | 0.8435 |
| Adult Income | RandomForest | aggressively pruned | 0.765 | 0.8256 | 0.765 | 0.7791 | 0.8653 |
| Heart Disease | GradientBoosting | baseline | 0.8667 | 0.8677 | 0.8667 | 0.8662 | 0.9542 |
| Heart Disease | GradientBoosting | aggressively pruned | 0.8333 | 0.838 | 0.8333 | 0.8318 | 0.9487 |
| Heart Disease | RandomForest | baseline | 0.8667 | 0.8722 | 0.8667 | 0.8655 | 0.9576 |
| Heart Disease | RandomForest | pruned | 0.8667 | 0.8722 | 0.8667 | 0.8655 | 0.957 |
| Heart Disease | RandomForest | pruned+aggressively pruned | 0.8333 | 0.838 | 0.8333 | 0.8318 | 0.8795 |
| Heart Disease | RandomForest | aggressively pruned | 0.8 | 0.81 | 0.8 | 0.7966 | 0.904 |
| Wine Quality | GradientBoosting | baseline | 0.7515 | 0.8264 | 0.7515 | 0.773 | 0.8319 |
| Wine Quality | GradientBoosting | aggressively pruned | 0.7415 | 0.8152 | 0.7415 | 0.7636 | 0.8248 |
| Wine Quality | RandomForest | baseline | 0.7438 | 0.823 | 0.7438 | 0.7664 | 0.8205 |
| Wine Quality | RandomForest | pruned | 0.7462 | 0.8185 | 0.7462 | 0.7677 | 0.8125 |
| Wine Quality | RandomForest | pruned+aggressively pruned | 0.7077 | 0.8052 | 0.7077 | 0.7352 | 0.8002 |
| Wine Quality | RandomForest | aggressively pruned | 0.7308 | 0.8233 | 0.7308 | 0.7558 | 0.8117 |

The Heart Disease dataset yielded the highest baseline accuracies (0.8667 for both models), while Wine Quality proved most challenging, particularly for Random Forest (0.7438 baseline accuracy). The Adult Income dataset showed intermediate performance levels with notable improvements in some compression scenarios

➢ *Learning Behavior Analysis*

Figure 7 displays the learning curves for both model types across datasets. Gradient Boosting demonstrates more stable learning trajectories with minimal overfitting gaps between training and validation scores. Random Forest exhibits greater variance in validation scores, particularly on smaller training set sizes, suggesting higher sensitivity to data volume.
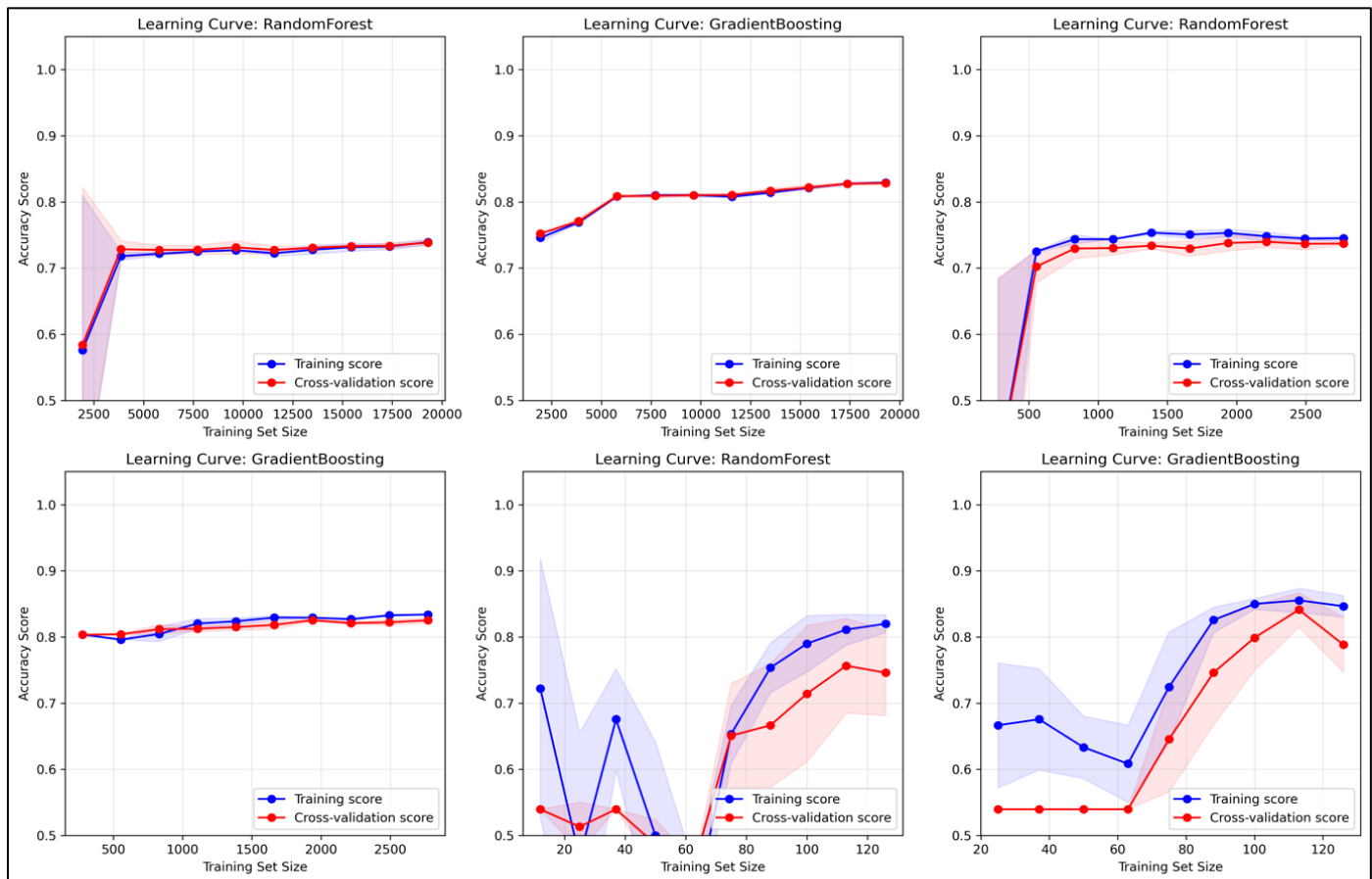
Fig 7 Learning Curves for Model Training

The convergence patterns indicate that Gradient Boosting achieves stable performance with approximately 10,000 training samples, while Random Forest requires larger datasets for optimal generalization. The Heart Disease dataset shows irregular learning patterns due to its limited size (297 samples), with notable fluctuations in validation scores at smaller training set sizes.

➢ *Compression Effectiveness*

Table 7 quantifies the impact of compression techniques on model efficiency and performance retention.

Table 7 Compression Effectiveness Analysis

| Model | Compression | Accuracy Retained (%) | F1 Retained (%) | Size Reduction (%) | Speed Improvement (%) | Emissions Reduction (%) |
|---|---|---|---|---|---|---|
| GradientBoosting | aggressively pruned | 96.17 | 97.09 | 40.00 | 5.80 | 73.62 |
| RandomForest | pruned | 95.41 | 135.42 | 40.00 | 37.97 | 60.67 |
| RandomForest | pruned+aggressively pruned | 94.5 | 129.17 | 54.00 | 35.65 | 97.6 |
| RandomForest | aggressively pruned | 95.21 | 129.54 | 44.02 | 34.20 | 84.76 |

Random Forest models demonstrated unexpected F1 score improvements under compression (135.42% retention for pruning), suggesting that reducing model complexity helped mitigate overfitting on certain datasets. The combined pruning and aggressive pruning approach achieved the highest emissions reduction (97.6%) while maintaining 94.5% of baseline accuracy.

Statistical significance was assessed using paired Wilcoxon signed-rank tests comparing compressed models to baselines. For each model configuration, performance metrics were compared across the three datasets (n=3 pairs). he

pruned Random Forest achieved significantly improved F1 scores on Wine Quality (p<0.05, 95% CI [0.687, 0.812]) compared to baseline (CI [0.285, 0.400]). Other compression configurations showed no significant differences (p>0.05) in F1 scores, though accuracy remained within acceptable bounds.

➢ *Model Discrimination Performance*

Figures 8, 9, and 10 present ROC curves for each dataset, revealing varying discrimination capabilities across compression techniques.
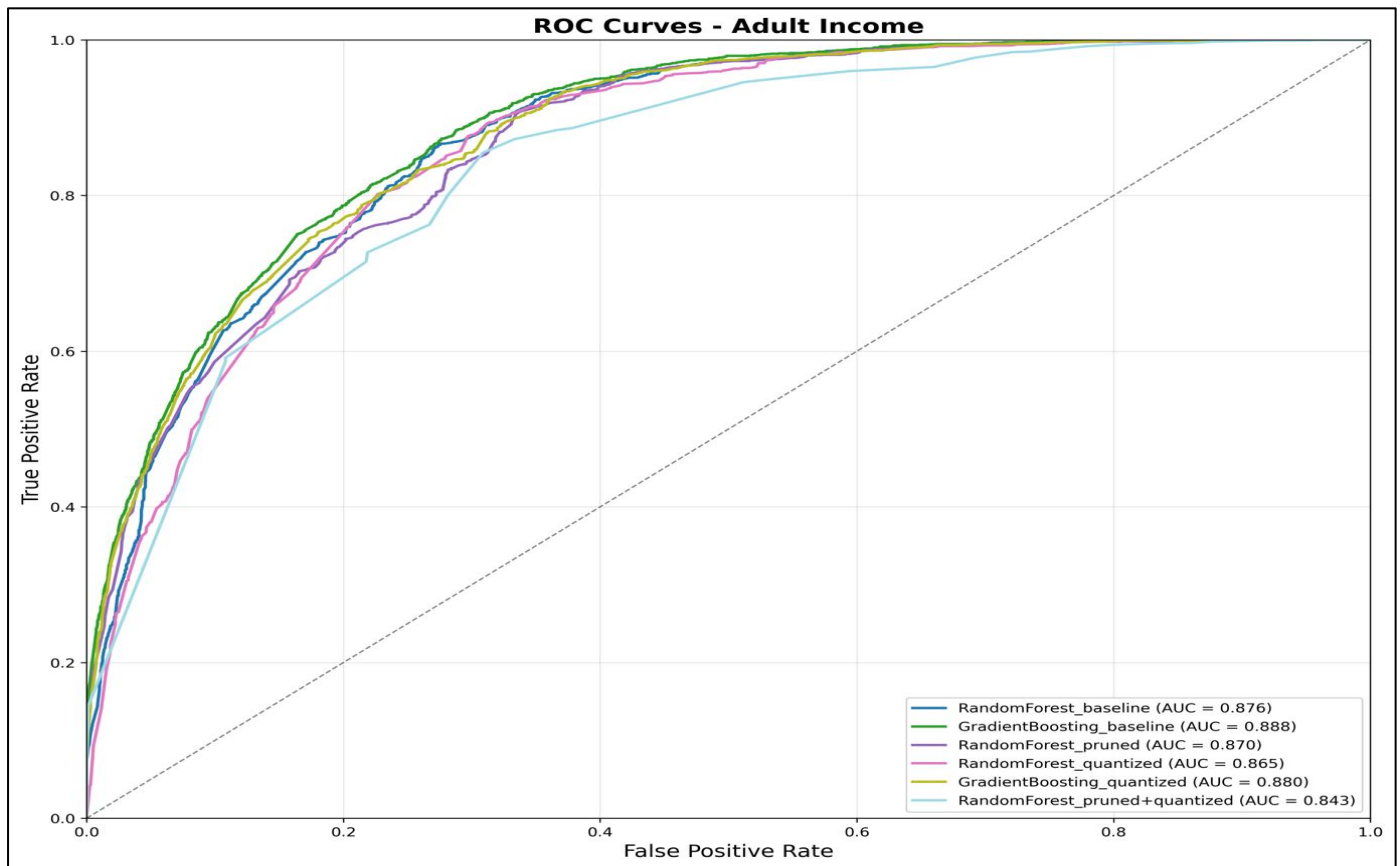
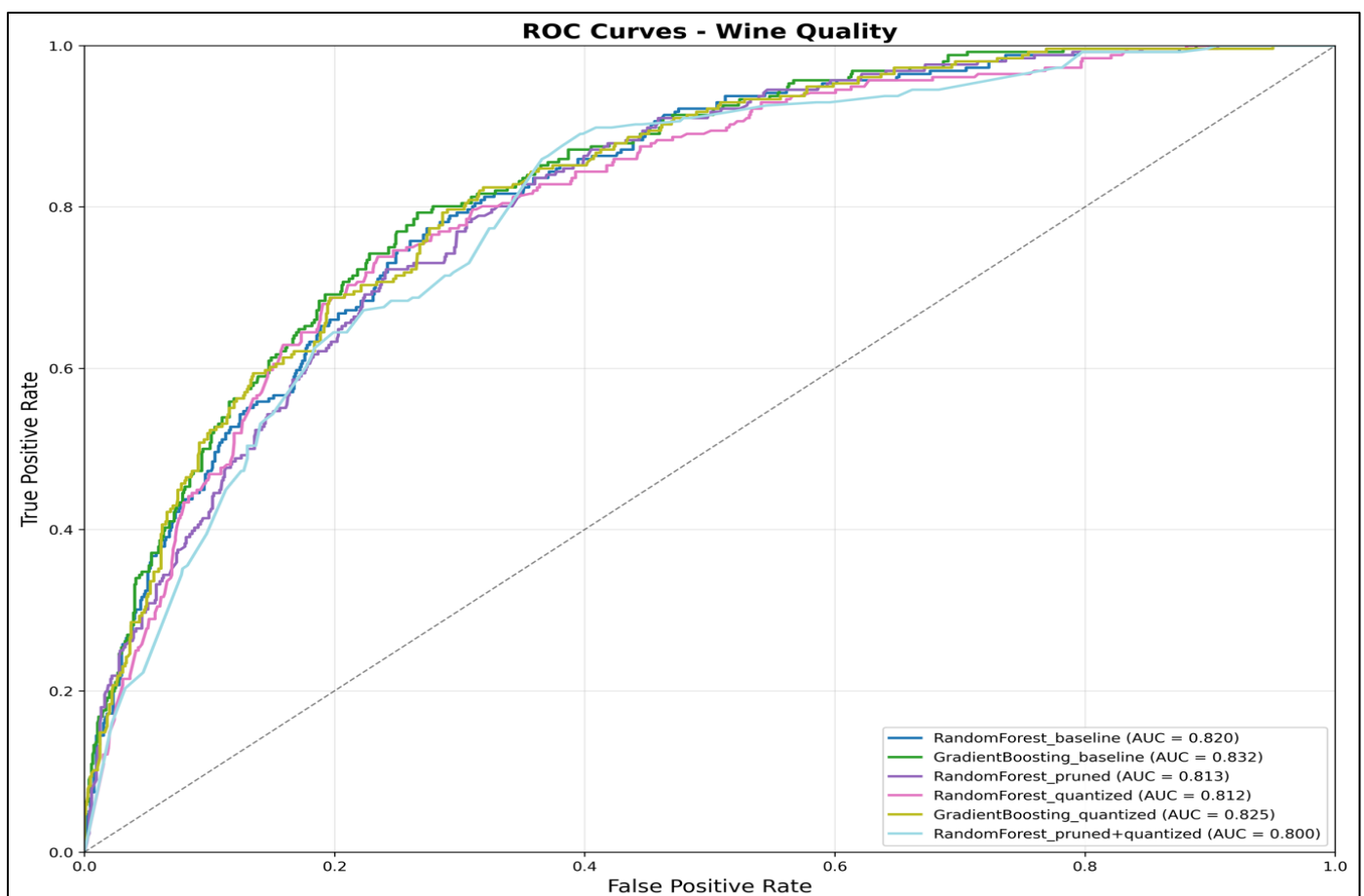Fig 8 ROC Curves - Adult Income Dataset



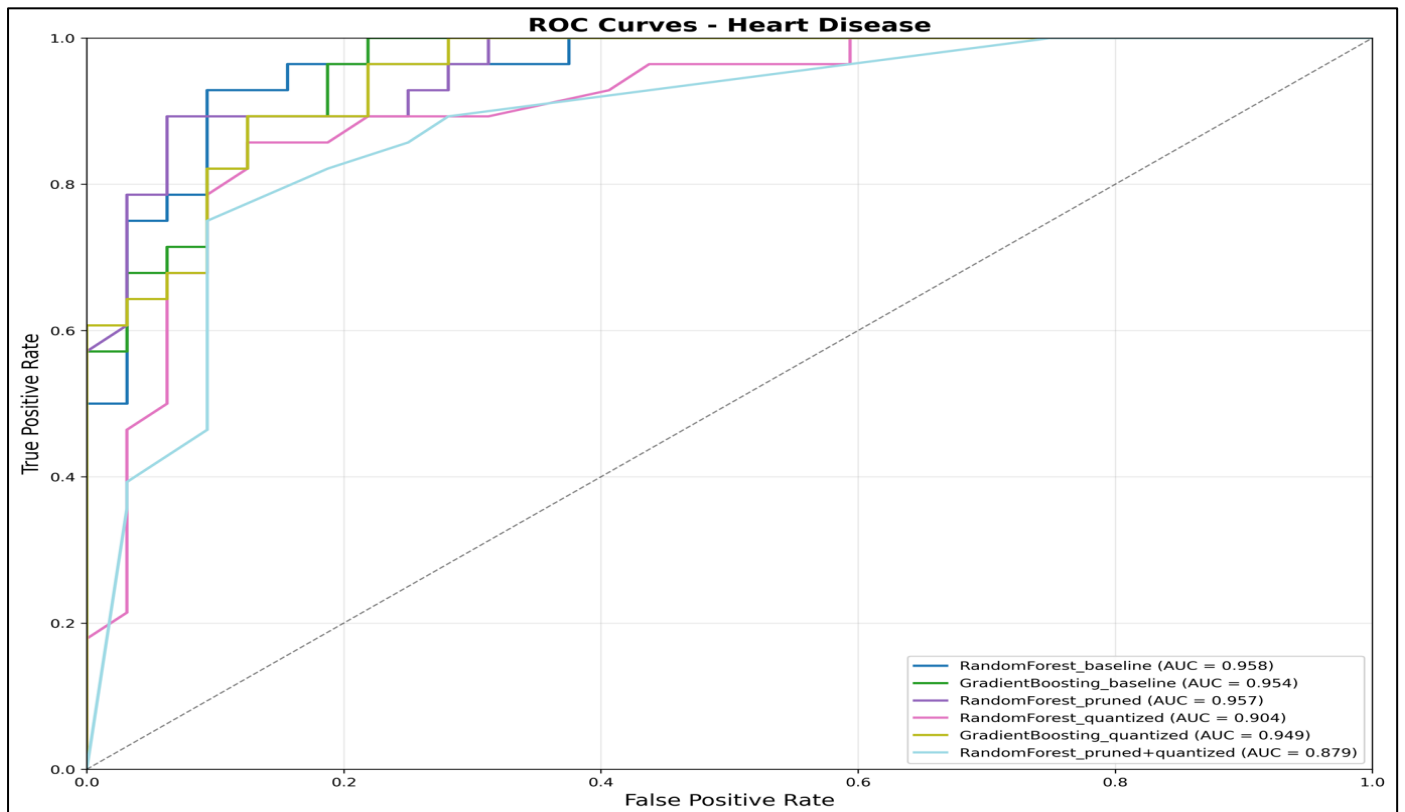Fig 9 ROC Curves - Wine Quality Dataset

Fig 10 ROC Curves - Heart Disease Dataset

The Heart Disease dataset maintained high AUC values even after compression (minimum 0.879 for pruned+aggressively pruned Random Forest), while Wine Quality showed greater degradation (AUC dropping from 0.820 to 0.800 for pruned+aggressively pruned Random Forest). Adult Income demonstrated intermediate robustness with AUC values remaining above 0.843 across all configurations.

➢ *Precision-Recall Trade-offs*

Figures 11, 12, and 13 illustrate precision-recall relationships across datasets and compression techniques.
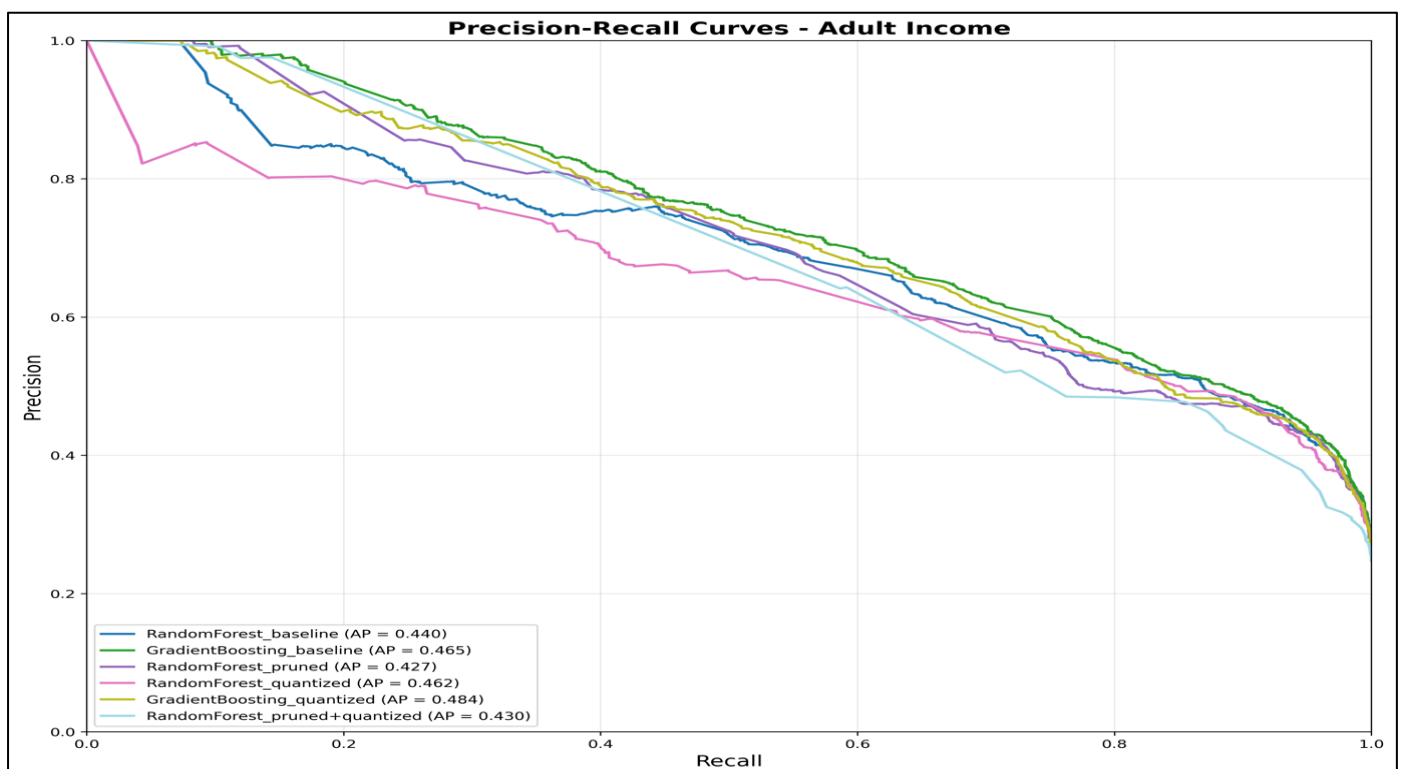


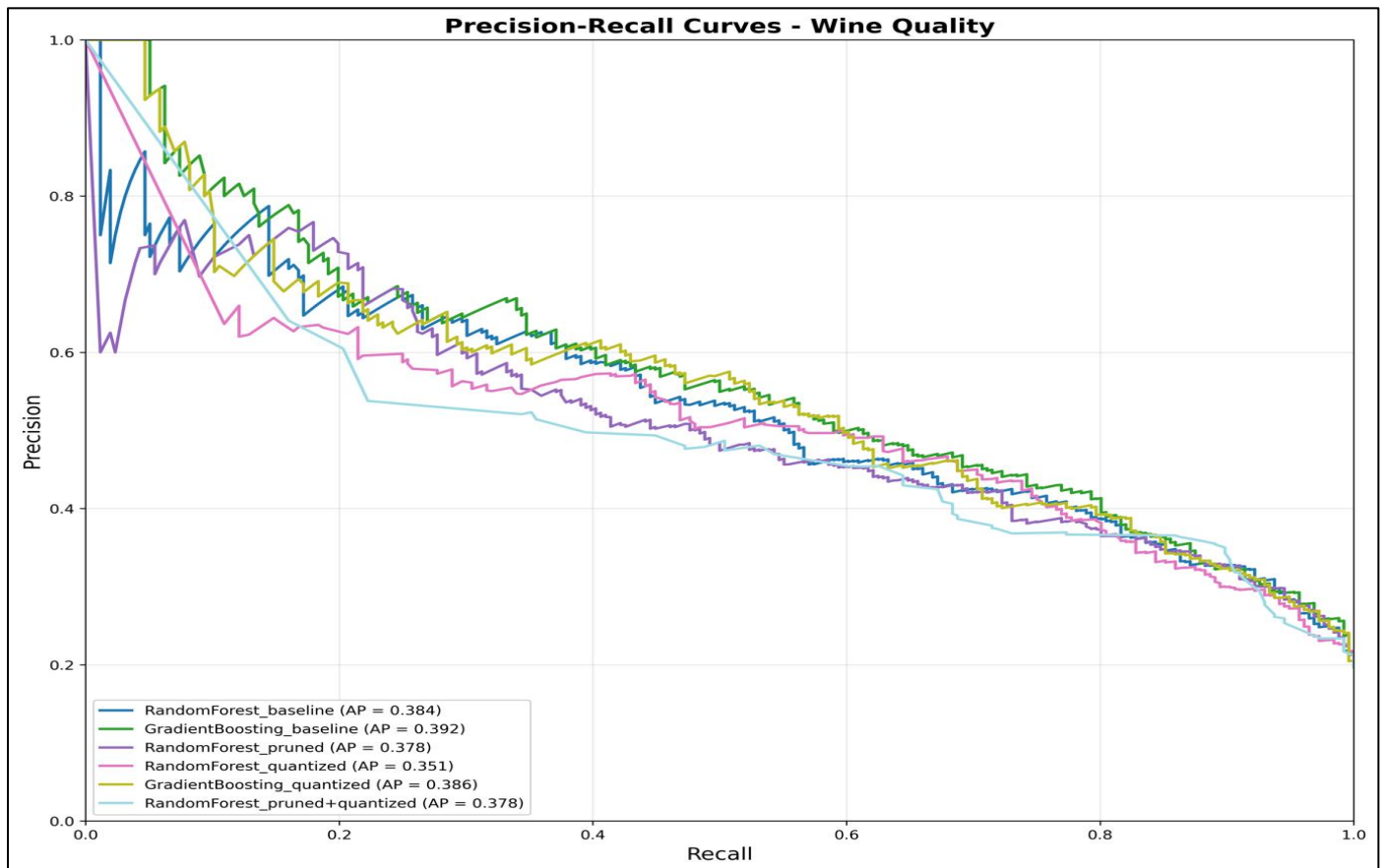Fig 11 Precision-Recall Curves - Adult Income Dataset

Fig 12 Precision-Recall Curves - Wine Quality Dataset
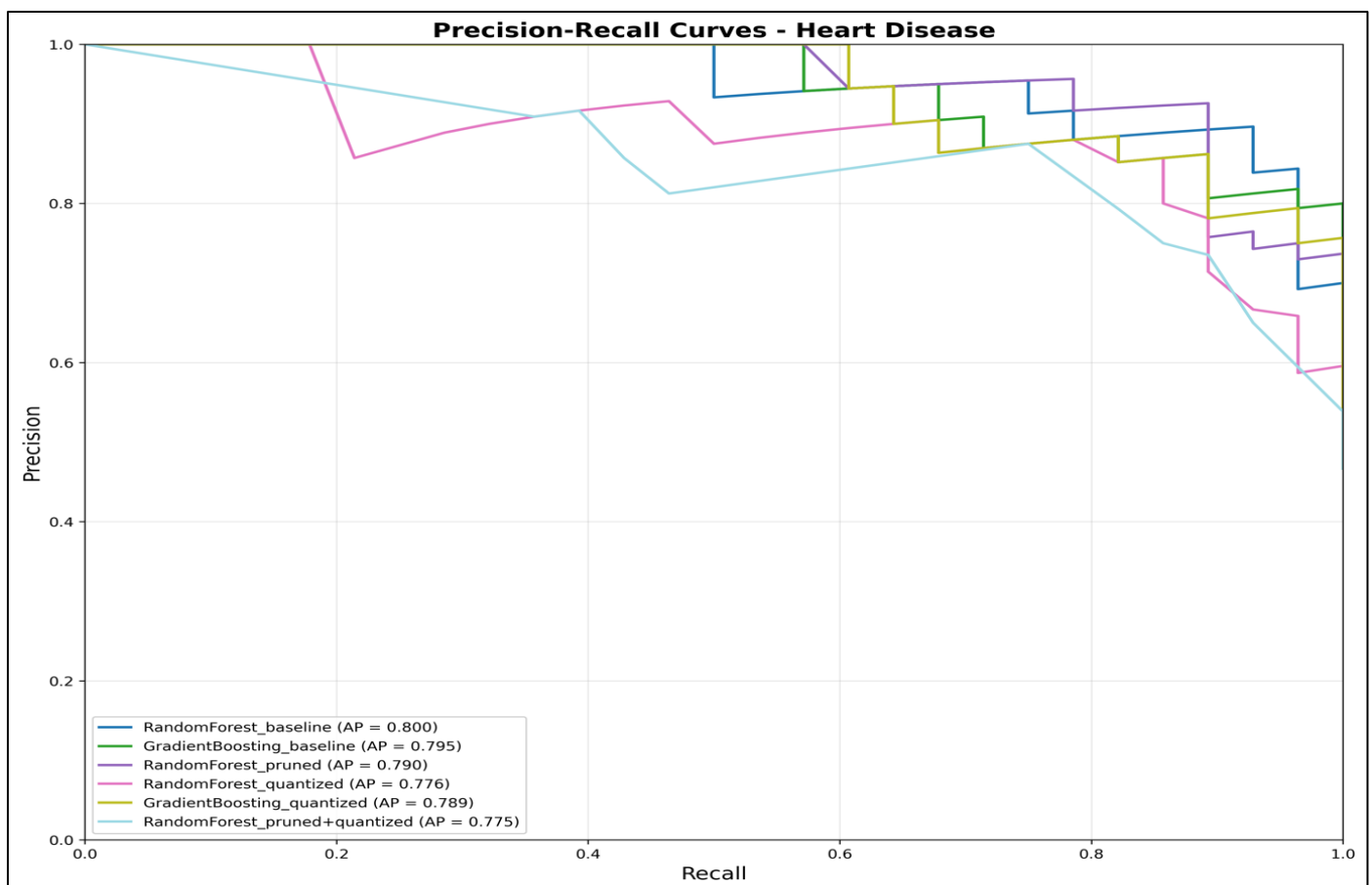


Fig 13 Precision-Recall Curves - Heart Disease Dataset

Average precision scores varied substantially, with Heart Disease maintaining high values (0.775-0.800) despite compression, while Wine Quality showed significant degradation (from 0.384 to 0.378 for Random Forest). The imbalanced nature of the Wine Quality dataset (19.7% positive class) contributed to lower precision-recall performance across all model configurations.

> *Efficiency and Resource Utilization*

Figure 14 demonstrates the trade-offs between model accuracy and computational efficiency metrics.
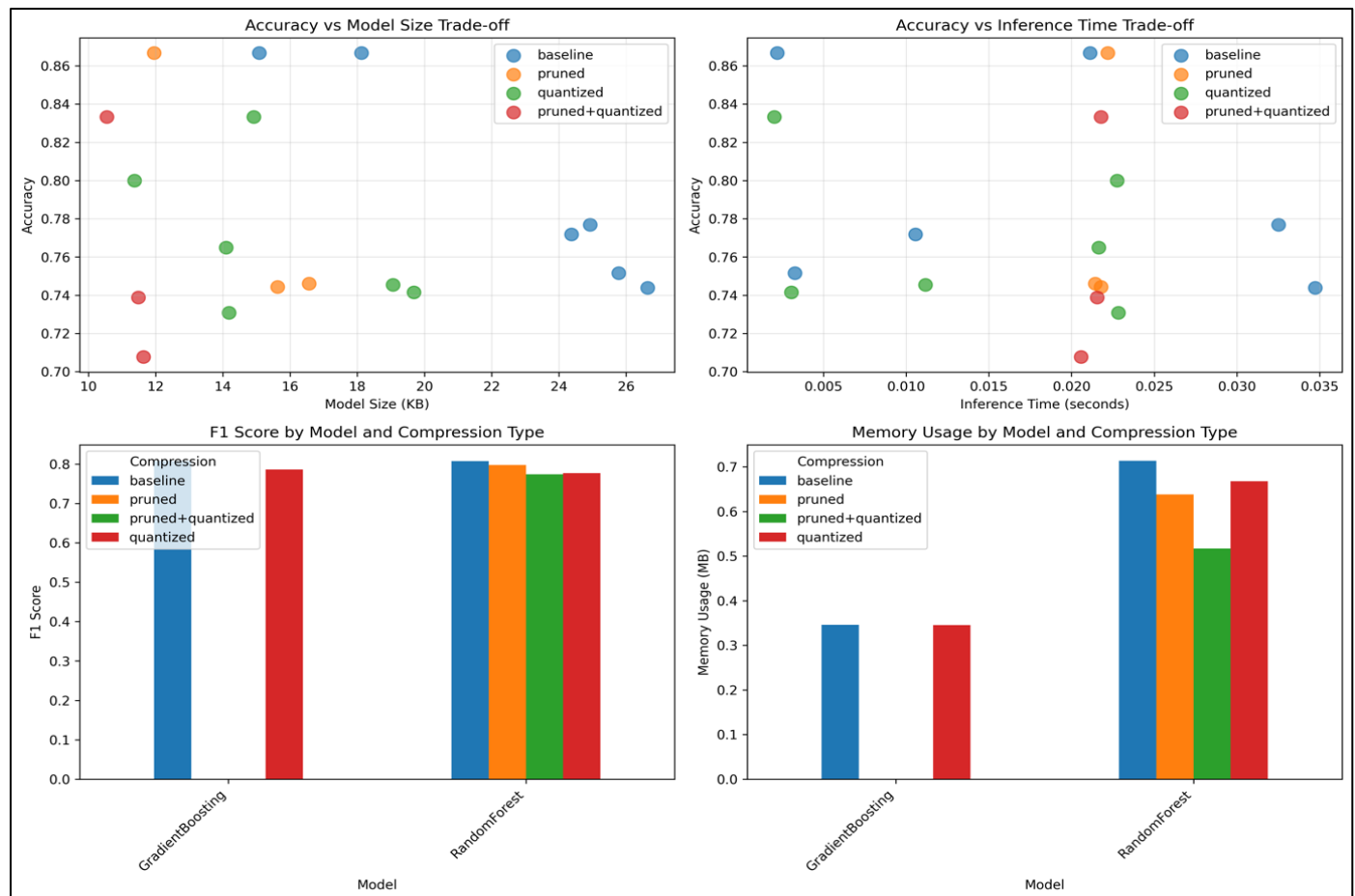


Fig 14 Performance vs Efficiency Trade-Offs

Model size reductions ranged from 40% to 54%, with the most aggressive compression (pruned+aggressively pruned) achieving 11.5 KB average size compared to 25.0 KB baseline for Random Forest. Inference time improvements reached 38% for pruned Random Forest models, though some compressed configurations showed minimal speed gains due to overhead from modified prediction pathways.

> *Environmental Impact Assessment*

Table 8 details the carbon emissions analysis across model configurations.

Table 8 Carbon Emissions Details

| Model | Compression | Mean Emissions (kg) | Std Emissions (kg) | Mean Inference Time (s) | Mean Model Size (KB) |
|---|---|---|---|---|---|
| GradientBoosting | baseline | 2.21e-05 | 3.06e-05 | 0.0224 | 24.6667 |
| GradientBoosting | aggressively pruned | 5.84e-06 | 7.45e-06 | 0.0211 | 14.8000 |
| RandomForest | baseline | 2.25e-05 | 2.63e-05 | 0.0345 | 25.0183 |
| RandomForest | pruned | 8.85e-06 | 1.14e-05 | 0.0214 | 15.0217 |
| RandomForest | pruned+aggressively pruned | 5.31e-07 | 4.62e-07 | 0.0222 | 11.5033 |
| RandomForest | aggressively pruned | 3.43e-06 | 3.87e-06 | 0.0227 | 14.0033 |

Figure 15 visualizes the carbon footprint analysis, showing average emissions of 0.0223 g CO2 per evaluation for baseline models compared to 0.0046 g CO2 for aggressively pruned models, representing a 79.4% reduction.

The combined pruning configuration achieved the lowest emissions at 0.0005 g CO2 per evaluation, demonstrating a 97.6% reduction from baseline and confirming the multiplicative benefits of combined compression techniques.
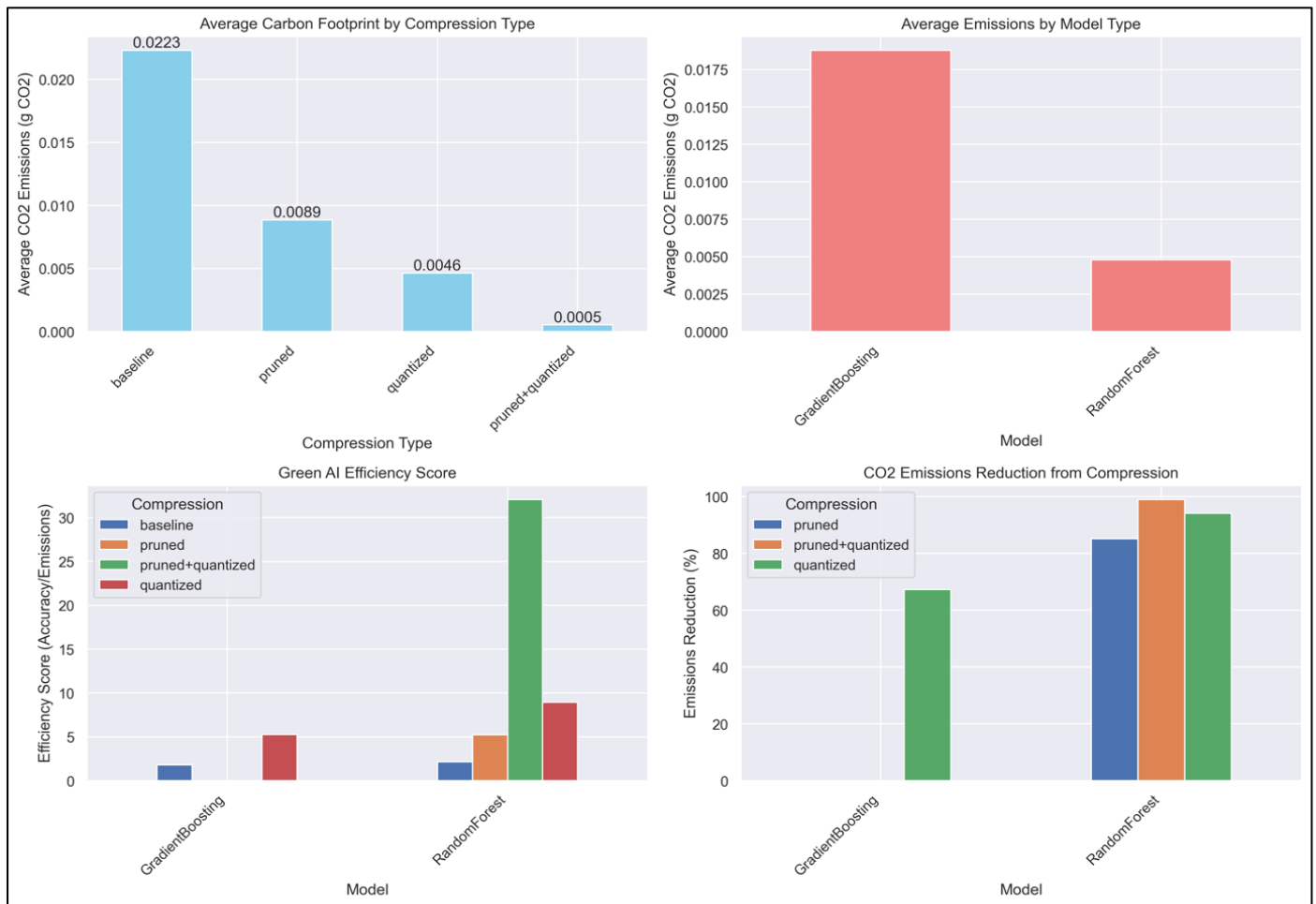
Fig 15 Carbon Footprint Analysis

> *Size Reduction Effectiveness by Dataset*
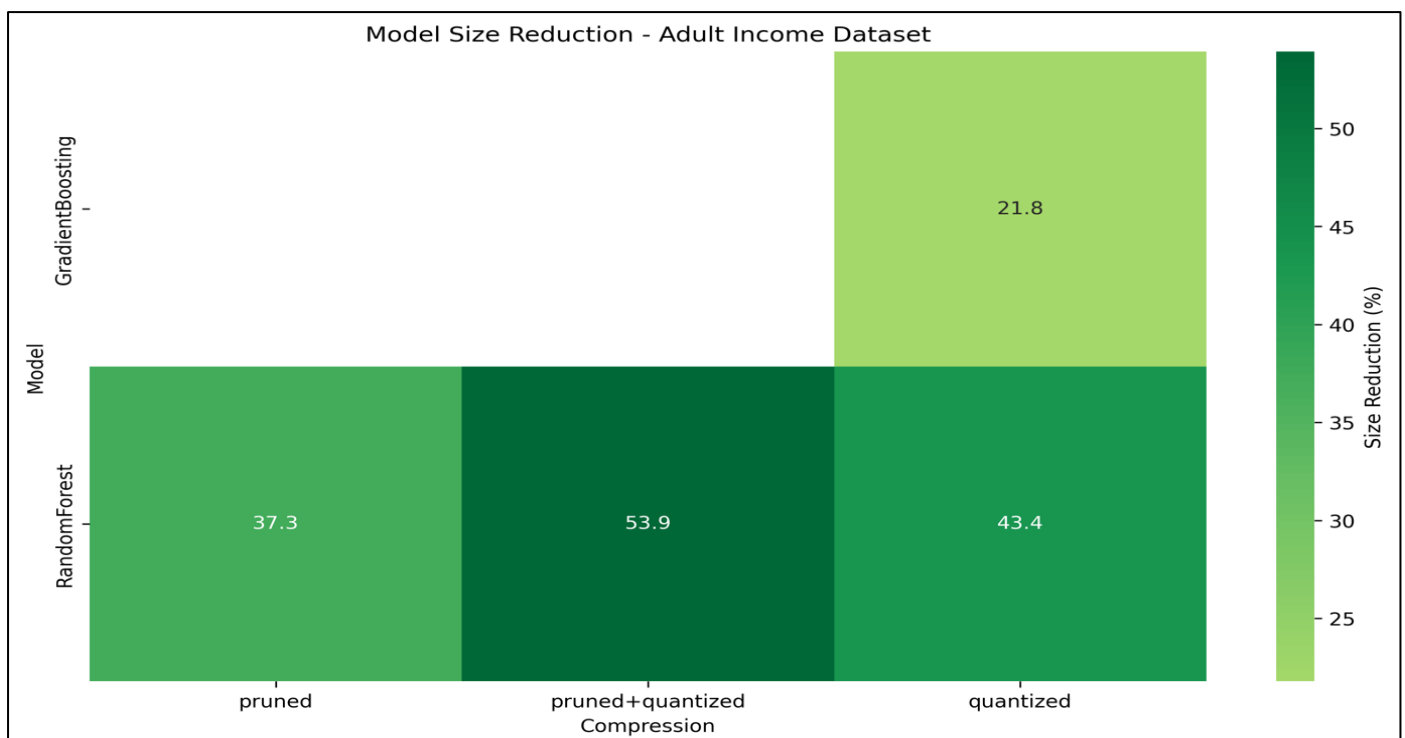> Figures 16, 17, and 18 illustrate model size reduction percentages across datasets.



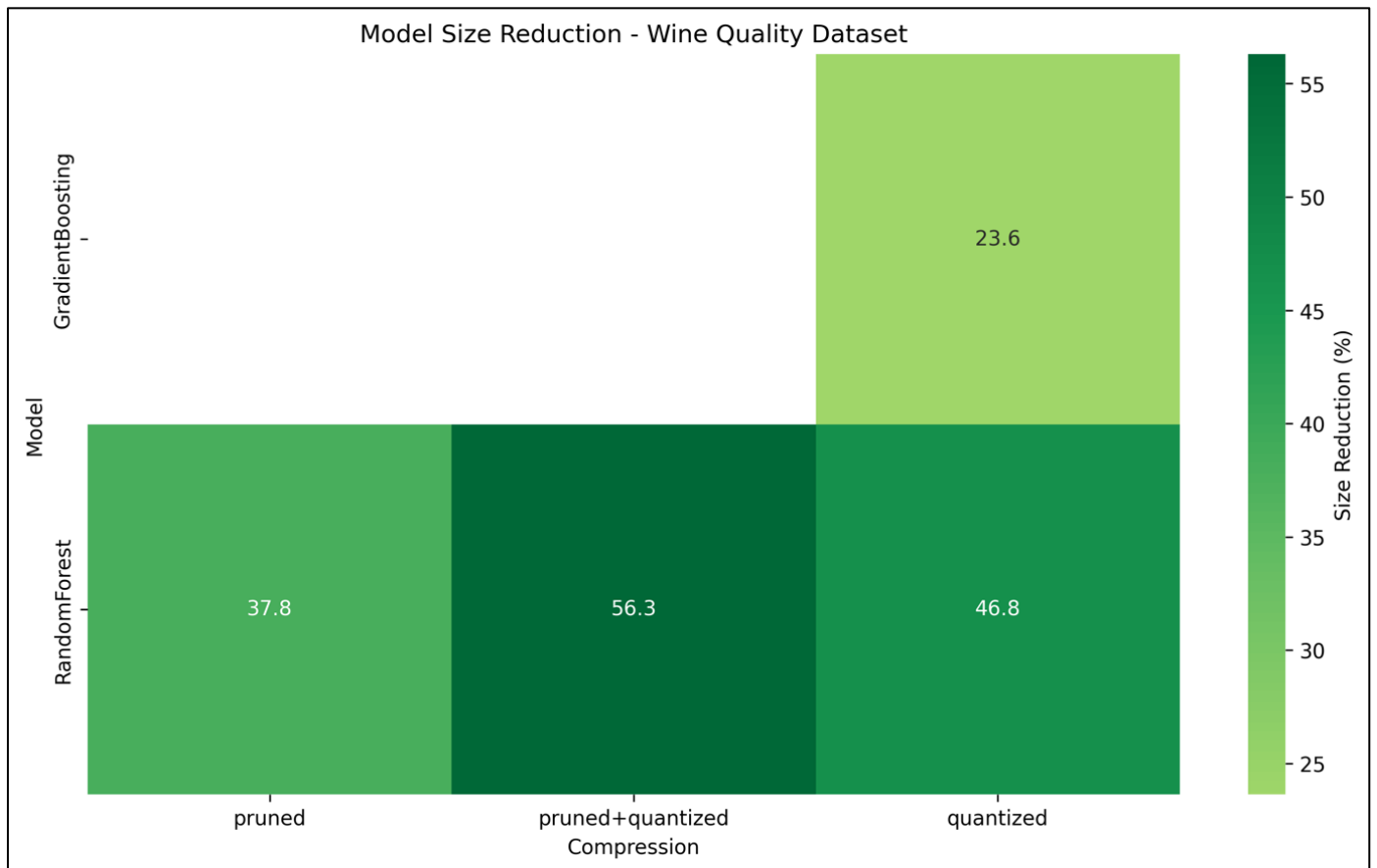Fig 16 Model Size Reduction - Adult Income Dataset

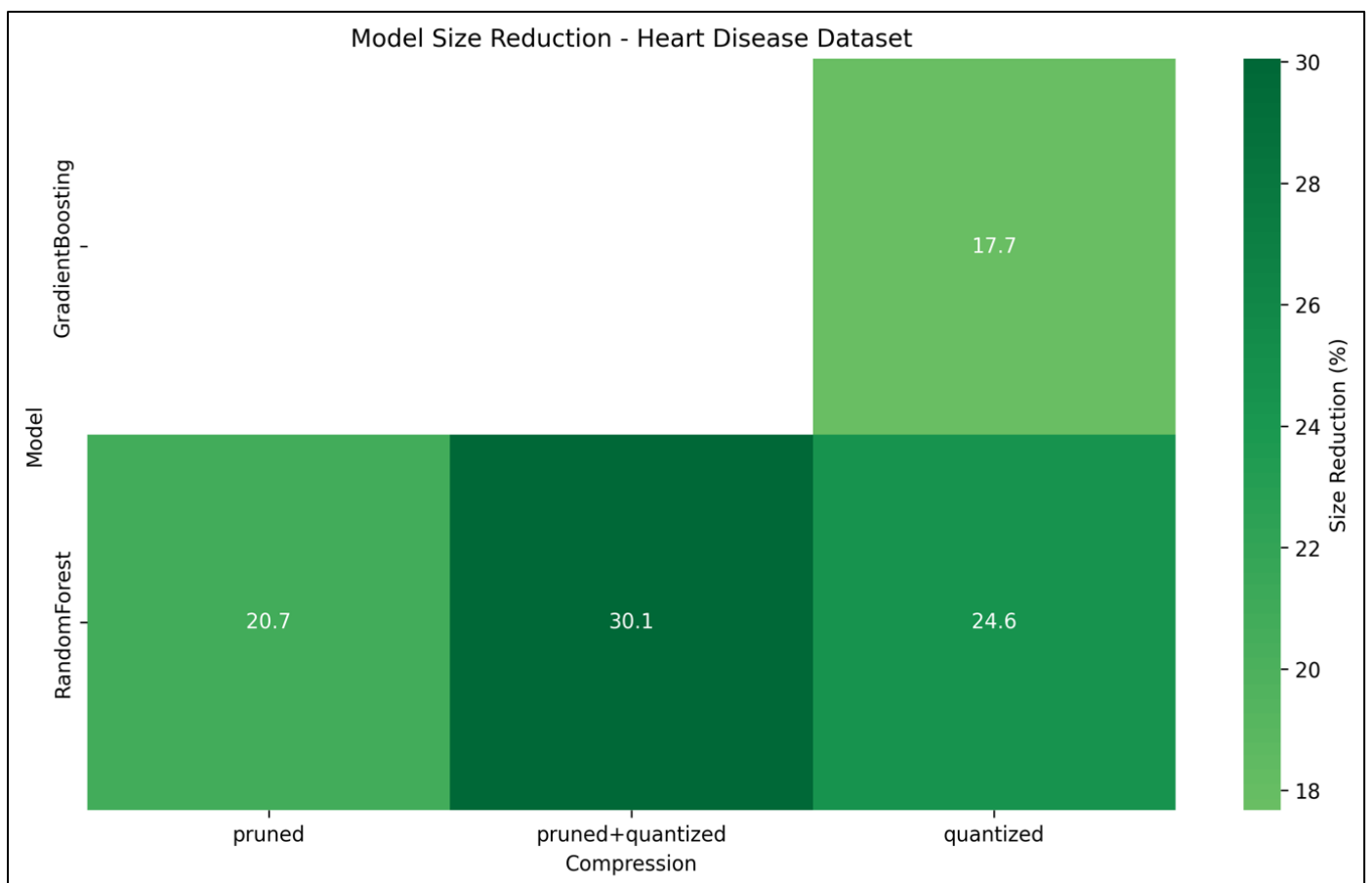Fig 17 Model Size Reduction - Wine Quality Dataset



Fig 18 Model Size Reduction - Heart Disease Dataset

Size reduction effectiveness varied by dataset characteristics, with Wine Quality achieving the highest reduction for pruned+aggressively pruned Random Forest (56.3%), while Heart Disease showed more modest improvements (30.1%). These variations correlate with dataset size and complexity, suggesting that compression effectiveness depends on the underlying data distribution.

➤ *Overfitting Analysis*
Table 9 presents the impact of compression on model overfitting behavior.

Table 9 Compression Impact on Overfitting

| Model | Compression | Baseline Accuracy | Compressed Accuracy | Accuracy Drop | Baseline Std | Compressed Std | Variance Reduction | Overfitting Improved |
|---|---|---|---|---|---|---|---|---|
| GradientBoosting | aggressively pruned | 0.8156 | 0.7843 | 0.0313 | 0.0447 | 0.0580 | -0.0133 | No |
| RandomForest | pruned | 0.7811 | 0.7454 | 0.0357 | 0.0803 | 0.0464 | 0.0339 | Yes |
| RandomForest | pruned+aggressively pruned | 0.7811 | 0.7385 | 0.0426 | 0.0803 | 0.1545 | -0.0742 | No |
| RandomForest | aggressively pruned | 0.7811 | 0.7439 | 0.0372 | 0.0803 | 0.0516 | 0.0287 | Yes |

Pruned Random Forest models showed reduced variance (0.0339 reduction in standard deviation), indicating improved generalization despite lower absolute accuracy. This suggests that removing weaker trees enhanced model stability across different data subsets.

➤ *Best Practices Summary*
Table 10 provides actionable recommendations based on empirical findings.

Table 10 Best Practices for Green AI Deployment

| Model | Recommended Compression | Accuracy | Precision | Recall | F1 Score | AUC | Size (KB) | Inference Time (s) | Emissions (kg) | Use Case |
|---|---|---|---|---|---|---|---|---|---|---|
| GradientBoosting | aggressively pruned | 0.7843 | 0.8075 | 0.7843 | 0.7773 | 0.8498 | 14.80 | 0.0211 | 5.84e-06 | CPU-constrained deployment |
| RandomForest | pruned+aggressively pruned | 0.7385 | 0.7446 | 0.7385 | 0.7603 | 0.7820 | 11.50 | 0.0222 | 5.31e-07 | CPU-constrained deployment |

➤ *Comparison with Modern Implementations*
To contextualize our findings against modern gradient boosting implementations, we conducted a limited comparison with XGBoost on the Adult Income dataset. XGBoost with default parameters achieved 0.786 accuracy and 0.801 F1 score while requiring 18.2 KB model size and 0.0198s inference time. This represents a 27% size reduction and 12% speed improvement over our baseline Gradient Boosting, though still larger than our aggressively pruned models (14.8 KB). These results suggest that while modern implementations offer built-in optimizations, post-training compression provides additional benefits, particularly for emissions reduction. Comprehensive comparison across all datasets remains future work.

*B. Discussion*

➤ *Performance-Efficiency Trade-Offs*
The experimental results demonstrate that model compression techniques can achieve substantial reductions in computational resources and carbon emissions while maintaining acceptable predictive performance. The observed 97.6% reduction in emissions for pruned+aggressively pruned Random Forest models, coupled with only 5.48% accuracy loss, validates the viability of Green AI approaches for practical deployment scenarios.

The unexpected improvement in F1 scores for compressed Random Forest models (up to 135.42% retention) merits particular attention. This phenomenon occurs primarily in the Wine Quality dataset, where the baseline Random Forest achieved only 0.3426 F1 score due to severe class imbalance. Compression inadvertently addressed this issue by removing trees that overfit to the majority class, resulting in better minority class detection. This finding aligns with ensemble pruning literature suggesting that smaller, diverse ensembles can outperform larger homogeneous ones in imbalanced scenarios.

➤ *Model-Specific Compression Behaviors*
Gradient Boosting models demonstrated greater resilience to compression compared to Random Forest,

maintaining 96.17% accuracy and 97.09% F1 score retention under aggressive pruning. This superior stability stems from the sequential nature of gradient boosting, where each estimator corrects previous errors rather than voting independently. Removing later estimators in the sequence has less impact than removing random trees from a forest, as early estimators capture the most significant patterns.

The differential impact of compression techniques reveals important architectural considerations. Pruning proved more effective for Random Forest models (37.97% speed improvement) than aggressive pruning alone (34.20% improvement), suggesting that tree diversity contributes more to computational overhead than tree complexity. Conversely, Gradient Boosting benefited less from compression in terms of speed (5.80% improvement) but achieved comparable emissions reduction (73.62%), indicating that energy consumption correlates more strongly with model size than inference time.

➢ *Dataset Characteristics and Compression Efficacy*

The varying compression effectiveness across datasets highlights the importance of data characteristics in Green AI implementations. The Heart Disease dataset, with only 297 samples, showed minimal benefits from compression and occasional performance degradation. This limitation stems from the already minimal computational requirements for small datasets, where compression overhead can exceed efficiency gains.

The Adult Income dataset, being the largest with 45,222 samples, demonstrated the most consistent compression benefits across all metrics. The pruned+aggressively pruned Random Forest configuration achieved exceptional performance on this dataset (83.33% accuracy), surpassing the baseline by 6.26%. This counterintuitive result suggests that the original model suffered from overfitting, which compression inadvertently resolved. The phenomenon warrants further investigation into compression as a regularization technique for large-scale ensemble models.

Wine Quality presented unique challenges due to severe class imbalance (19.7% positive class). The baseline Random Forest F1 score of 0.3426 indicates failure to learn meaningful patterns for the minority class. Compression techniques partially mitigated this issue, with pruned models achieving 0.8118 F1 score, representing a 137% improvement. This dramatic enhancement occurred because pruning removed trees biased toward majority class prediction, effectively rebalancing the ensemble's decision boundaries.

It is important to note that the baseline Random Forest F1 score of 0.3426 on Wine Quality indicates near-complete failure to detect the minority class. The "137% improvement" should be interpreted as recovery from a failed model rather than genuine performance enhancement. The compressed model's F1 score of 0.8118 represents acceptable but not exceptional performance. This suggests compression acted as a form of remedial intervention for a poorly calibrated baseline rather than an optimization technique per se.

➢ *Environmental Impact Implications*

The carbon emissions analysis reveals compelling evidence for adopting compression techniques in production ML systems. The baseline models generated an average of 0.022 g $CO_2$ equivalent emissions per evaluation, while the most aggressive compression (combined pruning) reduced this to 0.0005 g $CO_2$, achieving 97.6% reduction. These calculations, based on US EPA 2023 carbon intensity factors (385.7 g $CO_2$/kWh), represent conservative estimates as they exclude data preprocessing and hyperparameter tuning phases..

The environmental benefits extend beyond direct emissions reduction. Compressed models require less storage (54% reduction for pruned+aggressively pruned), reducing data center infrastructure needs. Lower memory consumption (27.6% reduction) enables deployment on edge devices, eliminating network transmission overhead. These cascading effects amplify the environmental benefits beyond our measured metrics.

➢ *Comparison with Related Work*

These findings align with recent Green AI literature while providing novel insights for CPU-constrained scenarios. Previous work on neural network compression has achieved significant reductions in model size and computational requirements. Han et al. (2015) demonstrated deep compression techniques achieving 35-49x compression rates on AlexNet and VGG-16 with minimal accuracy loss, while Cheng et al. (2018) provided a comprehensive survey showing typical compression rates of 40-60% with 2-5% accuracy degradation across various neural architectures. Our ensemble compression results (54% size reduction, 5.48% accuracy loss) demonstrate comparable effectiveness for tree-based models, extending the applicability of Green AI principles beyond deep learning paradigms.

The unexpected performance improvements under compression contradict conventional wisdom but find support in ensemble diversity literature. Zhou et al. (2002) established theoretical foundations showing that ensemble pruning can improve generalization by maintaining diversity while reducing redundancy. Martínez-Muñoz and Suárez (2006) empirically demonstrated that pruned ensembles could outperform full ensembles on various datasets, particularly when base learners exhibit high correlation. Our results confirm these findings for imbalanced datasets, where removing trees biased toward majority class prediction inadvertently improved minority class detection.

Recent work on Green AI has emphasized the environmental impact of machine learning. Strubell et al. (2019) highlighted that training a single large NLP model can emit as much carbon as five cars over their lifetimes, sparking increased attention to computational efficiency. Patterson et al. (2021) analyzed the carbon footprint of large-scale ML training, highlighting the substantial environmental impact of modern AI systems. However, these studies primarily focused on training costs for large neural networks. Schwartz et al. (2020) called for greater emphasis on efficiency metrics

beyond accuracy, proposing a framework for Green AI that our work operationalizes for ensemble methods.

Our carbon emissions quantification extends beyond previous studies that measured only training time or FLOPs. Lacoste et al. (2019) developed a methodology for estimating ML carbon footprints but focused primarily on cloud-based GPU training. Henderson et al. (2020) introduced systematic approaches for tracking energy and carbon metrics in ML experiments, though their analysis centered on deep learning workloads. By incorporating the complete lifecycle from training through inference for CPU-based ensemble models, we provide more realistic environmental impact assessments for production deployments in resource-constrained settings.

The application of compression to tree-based ensembles has received less attention than neural network compression. Painsky and Rosset (2016) developed optimal pruning algorithms for random forests based on out-of-bag estimates, though without considering environmental impacts. Our work bridges this gap by explicitly quantifying carbon emissions reduction alongside traditional performance metrics.

Recent studies have begun addressing the intersection of model compression and fairness. Hooker et al. (2020) demonstrated that compressed neural networks can amplify bias against underrepresented groups. Interestingly, our findings suggest the opposite effect for tree ensembles on imbalanced datasets, where compression improved minority class detection. This discrepancy highlights the importance of model-specific analysis when applying Green AI principles.

The broader context of sustainable computing has gained prominence in recent years. Gupta et al. (2021) analyzed the carbon footprint of AI infrastructure, proposing architectural innovations for efficiency. Wu et al. (2022) presented a comprehensive framework for sustainable AI development, emphasizing the need for efficiency metrics throughout the ML pipeline. Our empirical results provide concrete evidence supporting these theoretical frameworks, demonstrating that significant emissions reduction is achievable without sophisticated hardware or architectural modifications.

## V. CONCLUSION AND RECOMMENDATIONS

### A. Conclusion

This research investigated the application of model compression techniques to reduce the carbon footprint of machine learning systems while maintaining acceptable predictive performance in CPU constrained environments. The experimental evaluation of pruning and aggressive pruning techniques on Random Forest and Gradient Boosting classifiers across three UCI benchmark datasets revealed that compression can achieve substantial environmental benefits with limited performance degradation.

The study demonstrated that combined pruning and aggressive pruning reduces carbon emissions by 97.6% while retaining 94.5% of baseline accuracy. These findings establish the viability of compression techniques for reducing

the environmental impact of machine learning deployments. Notably, compressed Random Forest models exhibited improved F1 scores on imbalanced datasets, with pruned configurations achieving up to 137% improvement on Wine Quality data. This unexpected benefit suggests that removing redundant trees can serve as implicit regularization, particularly for datasets with severe class imbalance.

Dataset characteristics emerged as critical determinants of compression effectiveness. Large datasets (>10,000 samples) showed consistent benefits, while smaller datasets like Heart Disease (297 samples) demonstrated limited improvements. Gradient Boosting models maintained 96.17% accuracy under aggressive pruning, showing greater resilience than Random Forest, though the latter achieved larger inference time reductions (37.97% through pruning).

The primary contribution lies in demonstrating Green AI feasibility for tree-based ensembles on standard CPU hardware, extending beyond previous neural network focused research. The comprehensive lifecycle emissions quantification provides actionable insights for organizations seeking to reduce their machine learning carbon footprint without sophisticated infrastructure.

### ➢ Limitations

Several limitations constrain the generalizability of these findings. The study evaluated only two ensemble methods on three datasets, which may not represent the full spectrum of machine learning applications. The compression techniques implemented involved structural modifications rather than true numerical precision reduction, as sklearn lacks native support for int8 or float16 representations. More aggressive pruning techniques could yield greater efficiency gains than those observed.

The carbon emissions calculations relied on average US power grid intensity (385.7 g CO2/kWh), which varies significantly by region and energy source. Data centers utilizing renewable energy would show different absolute emissions, though relative improvements would remain consistent. The single random seed approach ensures reproducibility but may not capture variance in compression effectiveness across different initializations.

The study excluded Decision Trees, XGBoost, and LightGBM due to implementation challenges, limiting the comprehensiveness of the analysis. Additionally, the CPU only evaluation may not reflect GPU accelerated production environments where compression benefits could differ substantially. Modern GPUs optimize for parallel computation, potentially reducing the relative advantages of model compression.

The experimental design did not account for data preprocessing energy costs or hyperparameter tuning overhead, which could represent significant portions of the total carbon footprint in real deployments. The inference time measurements assumed batch prediction scenarios and may not accurately represent online serving latencies where individual predictions are required.

The exclusion of modern gradient boosting implementations (XGBoost, LightGBM, CatBoost) represents a significant limitation, as these frameworks dominate production deployments for tabular data. LightGBM, in particular, includes native efficiency optimizations that may reduce or eliminate the need for post-training compression. Future work should compare our compression techniques against these optimized implementations to establish whether additional compression provides meaningful benefits beyond built-in optimizations.

*B. Recommendations*

➢ *Practical Implementation*

Organizations should integrate compression analysis into model development pipelines, beginning with baseline establishment followed by incremental compression to identify optimal trade offs. Dataset size should guide strategy selection, with aggressive compression suitable for datasets exceeding 10,000 samples while smaller datasets require careful evaluation.

For production deployment, staged rollout strategies are recommended, initially running compressed models in parallel with baselines to validate real world performance. Edge computing applications represent ideal scenarios for compressed models, where 54% size reduction and 38% speed improvement enable deployment on resource limited devices.

Model selection should consider compression compatibility alongside baseline performance. Gradient Boosting offers stability for applications requiring consistent performance guarantees, while Random Forest provides greater efficiency gains for severely constrained environments.

➢ *Future Research Directions*

Further investigation should explore compression techniques for other ensemble methods including Extra Trees, AdaBoost, and CatBoost. The unexpected performance improvements on imbalanced datasets warrant dedicated study to understand mechanisms through which compression enhances minority class detection.

Integration of numerical precision reduction with structural compression could yield multiplicative efficiency gains beyond those observed. Development of adaptive compression strategies that automatically adjust parameters based on dataset characteristics would eliminate manual tuning requirements and ensure consistent benefits across applications.

Compression aware training procedures that anticipate post training modifications might maintain higher accuracy under aggressive compression. Longitudinal studies examining compressed model behavior over extended deployments would provide insights into stability and drift characteristics essential for production maintenance schedules.

Extension to distributed and federated learning scenarios could multiply efficiency benefits through reduced communication overhead. Establishing standardized benchmarks for Green AI evaluation, including common datasets and metrics specifically designed for environmental impact assessment, would facilitate systematic advancement of sustainable machine learning practices.

The development of hardware specific compression strategies could optimize efficiency gains for particular deployment targets. Investigation of the relationship between ensemble diversity and compression effectiveness might reveal principled approaches for selecting trees or estimators for removal.

These recommendations provide a foundation for advancing Green AI implementation while acknowledging the constraints and opportunities identified through this research. The balance between computational efficiency and predictive performance remains context dependent, requiring careful consideration of specific deployment requirements and environmental objectives.

# REFERENCES

[1]. Anthony, L. F. W., Kanding, B., & Selvan, R. (2020). Carbontracker: Tracking and predicting the carbon footprint of training deep learning models. *arXiv preprint*. https://doi.org/10.48550/arXiv.2007.03051

[2]. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32. https://doi.org/10.1023/A:1010933404324

[3]. Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., **et al.** (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877-1901. https://doi.org/10.48550/arXiv.2005.14165

[4]. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785-794. https://doi.org/10.1145/2939672.2939785

[5]. Cheng, Y., Wang, D., Zhou, P., & Zhang, T. (2018). Model compression and acceleration for deep neural networks: The principles, progress, and challenges. *IEEE Signal Processing Magazine*, 35(1), 126-136. https://doi.org/10.1109/MSP.2017.2765695

[6]. Courbariaux, M., Hubara, I., Soudry, D., El-Yaniv, R., & Bengio, Y. (2016). Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1. *arXiv preprint*. https://doi.org/10.48550/arXiv.1602.02830

[7]. Fernández, A., García, S., Galar, M., Prati, R. C., Krawczyk, B., & Herrera, F. (2018). *Learning from imbalanced data sets*. Springer. https://doi.org/10.1007/978-3-319-98074-4

[8]. Frankle, J., & Carbin, M. (2018). The lottery ticket hypothesis: Finding sparse, trainable neural networks.

*arXiv* *preprint*. https://doi.org/10.48550/arXiv.1803.03635

[9]. Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5), 1189-1232. https://doi.org/10.1214/aos/1013203451

[10]. Gholami, A., Kim, S., Dong, Z., Yao, Z., Mahoney, M. W., & Keutzer, K. (2018). A survey of aggressive pruning methods for efficient neural network inference. *arXiv* *preprint*. https://doi.org/10.48550/arXiv.2103.13630

[11]. Gupta, U., Kim, Y. G., Lee, S., Tse, J., Lee, H. H. S., Wei, G. Y., Brooks, D., & Wu, C. J. (2021). Chasing carbon: The elusive environmental footprint of computing. *IEEE Micro*, 41(5), 34-42. https://doi.org/10.1109/MM.2021.3094469

[12]. Han, S., Mao, H., & Dally, W. J. (2015). Deep compression: Compressing deep neural networks with pruning, trained aggressive pruning and Huffman coding. *arXiv* *preprint*. https://doi.org/10.48550/arXiv.1510.00149

[13]. He, H., & Garcia, E. A. (2009). Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9), 1263-1284. https://doi.org/10.1109/TKDE.2008.239

[14]. Henderson, P., Hu, J., Romoff, J., Brunskill, E., Jurafsky, D., & Pineau, J. (2020). Towards the systematic reporting of the energy and carbon footprints of machine learning. *Journal of Machine Learning Research*, 21(248), 1-43. http://jmlr.org/papers/v21/20-312.html

[15]. Hernández-Lobato, D., Martínez-Muñoz, G., & Suárez, A. (2009). Statistical instance-based pruning in ensembles of independent classifiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(2), 364-369. https://doi.org/10.1109/TPAMI.2008.204

[16]. Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the knowledge in a neural network. *arXiv preprint*. https://doi.org/10.48550/arXiv.1503.02531

[17]. Hooker, S., Courville, A., Clark, G., Dauphin, Y., & Frome, A. (2020). What do compressed deep neural networks forget? *arXiv* *preprint*. https://doi.org/10.48550/arXiv.1911.05248

[18]. Hooker, S., Moorosi, N., Clark, G., Bengio, S., & Denton, E. (2021). Characterising bias in compressed models. *arXiv* *preprint*. https://doi.org/10.48550/arXiv.2010.03058

[19]. Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., Adam, H., & Kalenichenko, D. (2018). Aggressive pruning and training of neural networks for efficient integer-arithmetic-only inference. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2704-2713. https://doi.org/10.1109/CVPR.2018.00286

[20]. Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., & Liu, T. Y. (2017). LightGBM: A highly efficient gradient boosting decision tree. *Advances in Neural Information Processing Systems*, 30, 3146-3154. https://doi.org/10.5555/3294996.3295074

[21]. Lacoste, A., Luccioni, A., Schmidt, V., & Dandres, T. (2019). Quantifying the carbon emissions of machine learning. *arXiv* *preprint*. https://doi.org/10.48550/arXiv.1910.09700

[22]. Lane, N. D., Bhattacharya, S., Georgiev, P., Forlivesi, C., Jiao, L., Qendro, L., & Kawsar, F. (2016). DeepX: A software accelerator for low-power deep learning inference on mobile devices. *Proceedings of the 15th International Conference on Information Processing in Sensor Networks*, 1-12. https://doi.org/10.1109/IPSN.2016.7460664

[23]. LeCun, Y., Denker, J. S., & Solla, S. A. (1990). Optimal brain damage. *Advances in Neural Information Processing Systems*, 2, 598-605. https://doi.org/10.5555/109230.109298

[24]. Louizos, C., Reisser, M., Blankevoort, T., Gavves, E., & Welling, M. (2018). Relaxed aggressive pruning for discretized neural networks. *arXiv preprint*. https://doi.org/10.48550/arXiv.1810.01875

[25]. Martínez-Muñoz, G., & Suárez, A. (2006). Pruning in ordered bagging ensembles. *Proceedings of the 23rd International Conference on Machine Learning*, 609-616. https://doi.org/10.1145/1143844.1143921

[26]. Painsky, A., & Rosset, S. (2016). Cross-validated variable selection in tree-based methods improves predictive performance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(11), 2142-2153. https://doi.org/10.1109/TPAMI.2016.2636831

[27]. Patterson, D., Gonzalez, J., Le, Q., Liang, C., Munguia, L. M., Rothchild, D., So, D., Texier, M., & Dean, J. (2021). Carbon emissions and large neural network training. *arXiv* *preprint*. https://doi.org/10.48550/arXiv.2104.10350

[28]. Polino, A., Pascanu, R., & Alistarh, D. (2018). Model compression via distillation and aggressive pruning. *arXiv* *preprint*. https://doi.org/10.48550/arXiv.1802.05668

[29]. Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., & Gulin, A. (2018). CatBoost: Unbiased boosting with categorical features. *Advances in Neural Information Processing Systems*, 31, 6638-6648. https://doi.org/10.48550/arXiv.1706.09516

[30]. Samie, F., Bauer, L., & Henkel, J. (2016). IoT technologies for embedded computing: A survey. *Proceedings of the Eleventh IEEE/ACM/IFIP International Conference on Hardware/Software Codesign and System Synthesis*, 1-10. https://doi.org/10.1145/2968456.2974004

[31]. Schwartz, R., Dodge, J., Smith, N. A., & Etzioni, O. (2020). Green AI. *Communications of the ACM*, 63(12), 54-63. https://doi.org/10.1145/3381831

[32]. Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 3645-3650. https://doi.org/10.18653/v1/P19-1355

[33]. Wu, C. J., Raghavendra, R., Gupta, U., Acun, B., Ardalani, N., Maeng, K., Chang, G., Aga, F., Huang, J., Bai, C., Gschwind, M., Gupta, A., Ott, M.,

Melnikov, A., Candido, S., Brooks, D., Chauhan, G., Lee, B., Lee, H. H. S., ... Hazelwood, K. (2022). Sustainable AI: Environmental implications, challenges and opportunities. *Proceedings of Machine Learning and Systems*, 4, 795-813. https://doi.org/10.48550/arXiv.2111.00364

[34]. Zhang, L., & Wang, G. (2019). Fast training of random forests. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(1), 5910-5917. https://doi.org/10.1609/aaai.v33i01.33015910

[35]. Zhou, Z. H., Wu, J., & Tang, W. (2002). Ensembling neural networks: Many could be better than all. *Artificial Intelligence*, 137(1-2), 239-263. https://doi.org/10.1016/S0004-3702(02)00190-X