# AI Ethics Project: Comprehensive Report

**Author**: AI Ethics Analysis

**Date**: November 30, 2025

**Project**: Complete AI Ethics Assessment and Audit

---

## Table of Contents

---

# Part 1: Theoretical Understanding

## 1. Short Answer Questions

### Q1: Define algorithmic bias and provide two examples of how it manifests in AI systems.

**Algorithmic bias** is when systematic errors in machine learning algorithms produce unfair and discriminatory outcomes that disadvantage certain groups of people based on characteristics like race, gender, age, or socioeconomic status.

**Two Examples:**

1. **Amazon's Hiring Tool**: The AI recruiting system was trained on historical resume data that predominantly came from male candidates. As a result, the algorithm learned to penalize resumes containing words like "women's" (e.g., "women's chess club captain") and downgraded graduates from all-women's colleges. This manifested as gender bias where qualified female candidates were systematically ranked lower than male counterparts.

2. **Facial Recognition Systems**: AI-powered facial recognition has been shown to have significantly higher error rates when identifying people with darker skin tones, particularly Black women. This bias stems from training data that overrepresents lighter-skinned individuals, causing the system to misidentify minorities at rates up to 35% higher than for white individuals. This has led to wrongful arrests and discriminatory surveillance.

# Q2: Explain the difference between transparency and explainability in AI. Why are both important?

**Transparency** refers to openness about an AI system's design, development, and operation. It means making visible:

- What data was used to train the model
- How the model was built (architecture, algorithms)
- Who developed it and for what purpose
- What the system's limitations are
- How it's being deployed and monitored

**Explainability** (or interpretability) refers to the ability to understand and articulate why an AI system made a specific decision or prediction. It answers:

- Which factors influenced this particular output?
- How did the model reach this conclusion?
- What would need to change for a different outcome?

**Why Both Matter:**

- **Transparency without explainability**: You might know that a neural network with 10 million parameters was trained on customer data, but still have no idea why it denied a specific person's loan application. This enables accountability at the system level but doesn't help individuals understand decisions affecting them.

- **Explainability without transparency**: You might receive an explanation for why you were rejected ("low credit score weighted heavily"), but if the underlying training data and model design are secret, you can't audit for bias or challenge systematic unfairness.

**Together, they enable**:

- **Accountability**: Developers and organizations can be held responsible for outcomes
- **Trust**: Users can have confidence in systems they understand
- **Debugging**: Identifying and fixing biases or errors
- **Legal Compliance**: Meeting requirements like GDPR's "right to explanation"
- **Fairness**: Detecting and addressing discriminatory patterns

# Q3: How does GDPR (General Data Protection Regulation) impact AI development in the EU?

The **GDPR** is a comprehensive data protection law in the European Union that significantly impacts AI development through several key provisions:

**1. Consent Requirements**:

- All data collection must be consensual—companies cannot use data to train AI models without explicit user permission
- Consent must be freely given, specific, informed, and unambiguous
- Users can withdraw consent at any time

**2. Right to Explanation (Article 22)**:

- Individuals have the right to not be subject to fully automated decisions with legal or significant effects without human involvement
- Users can demand explanations for algorithmic decisions affecting them
- This pushes developers toward more interpretable AI models

**3. Data Minimization**:

- Companies can only collect data strictly necessary for the stated purpose
- This limits the availability of large datasets often used for training AI
- Developers must justify why each data point is needed

**4. Right to Erasure ("Right to be Forgotten")**:

- Users can request deletion of their personal data
- This complicates AI systems that have already been trained on that data
- Models may need retraining when data is removed

**5. Privacy by Design**:

- Privacy protections must be built into systems from the beginning, not added later
- AI developers must conduct Data Protection Impact Assessments (DPIAs) for high-risk applications
- Requires techniques like differential privacy, federated learning, or anonymization

**Impact on AI Development**:

- **Slower Development**: Obtaining proper consent and ensuring compliance slows data collection and model training
- **Limited Data Access**: Stricter rules mean smaller, potentially less diverse training datasets
- **Higher Costs**: Compliance requires legal review, technical privacy enhancements, and ongoing monitoring
- **Innovation in Privacy-Preserving AI**: Drives development of techniques like federated learning, homomorphic encryption, and synthetic data generation
- **Massive Penalties**: Non-compliance can result in fines up to €20 million or 4% of global annual revenue, creating strong incentives for ethical practices
- **Competitive Disadvantage**: EU companies may have access to less data than competitors in regions with weaker privacy laws

However, GDPR also promotes **responsible AI development** by forcing developers to consider privacy and fairness from the outset, potentially leading to more trustworthy and ethically sound systems.

# 2. Ethical Principles Matching

**A) Justice** - Fair distribution of AI benefits and risks.

**B) Non-maleficence** - Ensuring AI does not harm individuals or society.

**C) Autonomy** - Respecting users' right to control their data and decisions.

**D) Sustainability** - Designing AI to be environmentally friendly.

---

# Part 2: Case Study Analysis

# Case 1: Biased Hiring Tool - Amazon's AI Recruiting System

## Background

Amazon's AI recruiting tool was designed to automate resume screening but was discovered to systematically penalize female candidates, particularly for technical positions.

## 1. Source of Bias

The primary source of bias stemmed from the **training data**:

- **Historical Data Bias**: The model was trained on resumes submitted to Amazon over a 10-year period, predominantly from male candidates. Since the tech industry historically has had male-dominated hiring patterns, the AI learned to favor patterns associated with male resumes.

- **Feature Engineering Issues**: The model identified gendered language patterns as predictive features. For example, resumes containing words like "women's" (e.g., "women's chess club captain") were penalized, and the system downgraded graduates from all-women's colleges.

- **Proxy Variables**: The model inadvertently used gender as a proxy through correlated features such as sports teams, clubs, and educational institutions that are gender-associated.

## 2. Three Fixes to Make the Tool Fairer

### Fix 1: Diverse and Balanced Training Data

- Reconstruct the training dataset to include equal representation of successful male and female candidates
- Use synthetic data augmentation to balance gender representation

- Include successful candidates from diverse backgrounds and institutions
- Remove historical hiring data that reflects past discriminatory practices

### Fix 2: Fairness-Aware Machine Learning Techniques

- Implement adversarial debiasing where a secondary model attempts to predict gender from the hiring model's internal representations, and the primary model is trained to prevent this
- Apply reweighing techniques to give more importance to underrepresented groups during training
- Use fairness constraints during optimization (e.g., demographic parity or equalized odds constraints)
- Remove explicitly gendered terms and gender-proxy features from the feature set

### Fix 3: Human-in-the-Loop Validation and Ongoing Monitoring

- Implement mandatory human review for all AI recommendations, especially for borderline cases
- Create diverse hiring panels to review AI-flagged candidates
- Establish regular bias audits with gender-disaggregated performance metrics
- Set up feedback loops where human recruiters can flag biased decisions to retrain the model
- Implement explainability tools to understand why candidates are ranked certain ways

# 3. Metrics to Evaluate Fairness Post-Correction

### Demographic Parity Metrics:

- **Selection Rate Parity**: The percentage of male vs. female candidates recommended should be proportional to the qualified applicant pool
  - Formula: $P(\hat{Y}=1|Gender=Male) \approx P(\hat{Y}=1|Gender=Female)$

### Equalized Odds Metrics:

- **True Positive Rate (TPR) Parity**: The rate at which qualified male and female candidates are correctly identified should be equal
  - $TPR\_male \approx TPR\_female$
- **False Positive Rate (FPR) Parity**: The rate at which unqualified candidates are incorrectly recommended should be equal across genders
  - $FPR\_male \approx FPR\_female$

### Predictive Parity Metrics:

- **Positive Predictive Value (PPV) Parity**: Among candidates recommended by the AI, the success rate should be equal across genders
  - $PPV\_male \approx PPV\_female$

### Calibration Metrics:

- **Calibration Across Groups**: For candidates assigned similar scores, actual success rates should be similar regardless of gender
  - For score S: $P(Success|Score=S, Gender=Male) \approx P(Success|Score=S, Gender=Female)$

- **Disparate Impact Ratio**: Ratio of selection rates between protected and unprotected groups should be ≥ 0.8 (80% rule)
    - DIR = (Female Selection Rate) / (Male Selection Rate) should be between 0.8 and 1.25
- **Statistical Parity Difference**: Absolute difference in selection rates should be minimal (ideally < 0.1)

**Process Metrics:**

- Time-to-hire equity across gender groups
- Interview invitation rates by gender
- Offer acceptance rates by gender
- Long-term success metrics (performance reviews, retention) by gender

---

# Case 2: Facial Recognition in Policing

## Background

Facial recognition systems deployed in law enforcement have been shown to misidentify minorities, particularly Black and Asian individuals, at significantly higher rates than white individuals.

# 1. Ethical Risks

### Wrongful Arrests and Criminalization

- **False Positives Leading to Wrongful Detention**: Higher misidentification rates for minorities mean innocent people are more likely to be wrongly arrested, detained, and potentially prosecuted
- **Compounding Bias**: These systems can amplify existing racial biases in policing, leading to disproportionate targeting of minority communities
- **Psychological and Social Harm**: Wrongful arrests cause trauma, job loss, reputational damage, and erosion of trust in law enforcement
- **Due Process Violations**: Individuals may be detained based on algorithmic decisions without proper probable cause

### Privacy Violations

- **Mass Surveillance**: Facial recognition enables warrantless, continuous surveillance of public spaces, disproportionately affecting minority neighborhoods that are often over-policed
- **Chilling Effects**: Knowledge of constant monitoring can suppress freedom of movement, assembly, and expression, particularly in communities already marginalized
- **Data Collection Without Consent**: Individuals' biometric data is collected and stored without knowledge or consent
- **Scope Creep**: Data collected for one purpose (e.g., identifying suspects) can be repurposed for broader surveillance

### Erosion of Civil Liberties

- **Presumption of Guilt**: Being flagged by facial recognition can create a presumption of guilt rather than innocence
- **Lack of Transparency**: Individuals often don't know when they've been scanned or why they were flagged
- **Unequal Protection Under Law**: Higher error rates for minorities mean unequal application of surveillance and enforcement

### Reinforcement of Systemic Discrimination

- **Feedback Loops**: If minorities are arrested more due to false positives, this creates more "criminal records" that justify further surveillance
- **Historical Bias Amplification**: Training data often reflects historical patterns of discriminatory policing
- **Resource Misallocation**: Over-policing of minority communities based on flawed data perpetuates inequity

### Accountability Gaps

- **Diffusion of Responsibility**: Errors are blamed on "the algorithm," making it difficult to hold anyone accountable
- **Opacity**: Proprietary systems lack transparency, preventing independent audits
- **Limited Recourse**: Victims of misidentification often have limited legal remedies

# 2. Policies for Responsible Deployment

### Pre-Deployment Requirements

### Policy 1: Mandatory Bias Testing and Certification

- Require independent, third-party testing of facial recognition systems before deployment
- Systems must demonstrate accuracy rates above 99% across all demographic groups (race, gender, age)
- Publish disaggregated performance metrics publicly
- Require annual re-certification with updated test data

### Policy 2: Strict Use Case Limitations

- Prohibit use of facial recognition for:
    - Mass surveillance in public spaces
    - Identification at protests or political gatherings
    - Real-time identification without prior warrant
- Limit use to specific, serious crimes (e.g., violent felonies, missing persons)
- Require judicial warrant before deployment in most cases

### Deployment Safeguards

### Policy 3: Human Review and Verification Requirements

- Mandate that facial recognition can NEVER be the sole basis for arrest or detention
- Require human expert review of all matches before any law enforcement action
- Implement minimum confidence thresholds (e.g., 95%+) before human review

- Require corroborating evidence beyond facial recognition match

### Policy 4: Transparency and Explainability Standards

- Require open-source algorithms or source code escrow for judicial review
- Law enforcement must disclose use of facial recognition in arrest affidavits
- Defendants must receive all facial recognition data and match scores
- Public registry of when and where systems are deployed

### Accountability Mechanisms

### Policy 5: Data Protection and Privacy Safeguards

- Limit database sources to criminal mugshots only (no DMV, social media, or commercial databases)
- Mandatory data minimization: delete scans that don't match within 24 hours
- Prohibition on sharing data with third parties or other agencies without court order
- Right to know if you've been scanned and matched

### Policy 6: Community Oversight and Consent

- Require community input before deployment through public hearings
- Create civilian oversight boards with subpoena power to audit usage
- Allow communities to opt-out of facial recognition deployment
- Regular public reporting on usage statistics and error rates

### Policy 7: Legal Liability and Redress

- Hold vendors liable for discriminatory performance through contract terms
- Create civil cause of action for victims of misidentification
- Automatic expungement of records for false positive arrests
- Compensation fund for victims of wrongful detention

### Policy 8: Training and Expertise Requirements

- Require specialized training for officers using facial recognition
- Training must include bias awareness and system limitations
- Designate certified specialists for system operation
- Document all usage with justification and outcomes

### Policy 9: Sunset Provisions and Ongoing Evaluation

- Implement 2-year sunset clauses requiring re-authorization
- Mandatory annual third-party audits of deployment outcomes
- Disaggregated analysis of impact on different demographic groups
- Authority to suspend deployment if disparate impact is identified

### Policy 10: Alternative Investigation Methods

- Require documentation that traditional investigative methods were considered

- Invest in less invasive alternatives (e.g., tip lines, witness interviews)
- Prioritize community policing over technological surveillance

---

# Part 3: Practical Audit

## COMPAS Recidivism Dataset Bias Audit

## Executive Summary

This practical audit analyzed the COMPAS (Correctional Offender Management Profiling for Alternative Sanctions) recidivism risk assessment dataset using Python and statistical analysis. The audit reveals significant racial disparities in how the algorithm assigns risk scores, with African-American defendants being disproportionately classified as high-risk compared to Caucasian defendants.

## Methodology

**Tools Used**:

- Python 3.x with pandas, numpy, matplotlib, seaborn
- Statistical analysis libraries (sklearn)
- AI Fairness 360 (IBM's toolkit) for fairness metrics
- COMPAS dataset from ProPublica's investigation

**Analysis Approach**:

1. Load and preprocess COMPAS dataset
2. Calculate demographic parity metrics
3. Analyze false positive and false negative rate disparities
4. Generate visualizations of bias patterns
5. Develop remediation recommendations

## Key Findings

### Finding 1: Disparate High-Risk Classification Rates

- African-American high-risk rate: ~60-65%
- Caucasian high-risk rate: ~35-40%
- Disparity: Approximately 20-25 percentage points difference

### Finding 2: False Positive Rate Disparity *(Most Critical)*

- African-American FPR: ~45-50%
- Caucasian FPR: ~23-25%

- **Impact**: African-Americans are approximately **2x more likely** to be wrongly labeled as high-risk

### Finding 3: False Negative Rate Disparity *(Inverse Pattern)*

- Caucasian FNR: ~45-50%
- African-American FNR: ~25-30%
- **Impact**: Caucasian defendants are more likely to receive lower risk scores despite similar or higher likelihood of recidivism

### Finding 4: Predictive Parity vs. Error Rate Parity Trade-off

- The system cannot simultaneously achieve equal FPR AND FNR across racial groups while maintaining overall accuracy
- Current configuration prioritizes overall accuracy over equal treatment

# Visualizations Generated

The audit code produces six key visualizations:

1. Risk score distribution by race
2. High-risk classification rate comparison
3. False positive rate comparison
4. False negative rate comparison
5. Decile score distribution
6. Actual recidivism rate by predicted risk score

# Root Causes

**Data Bias**: Training data reflects decades of racially biased policing patterns

**Structural Factors**: Systemic inequalities in criminal justice create feedback loops

**Algorithmic Design**: Optimization for accuracy rather than fairness across groups

# Remediation Steps

**Immediate Actions**:

1. Transparency mandate for all risk assessment tools
2. Mandatory bias testing before deployment
3. Human override capabilities
4. Right to explanation for defendants

**Technical Interventions**: 5. Fairness-aware training with explicit constraints 6. Feature auditing to remove race proxies 7. Data rebalancing techniques 8. Adversarial debiasing

**Policy Reforms**: 9. Limited use cases for risk assessments 10. Regular reauditing with public reporting 11. Accountability mechanisms and compensation 12. Community oversight boards

# Code Implementation

The complete audit code (`compas_audit.py`) includes:

- Automated data loading from local Excel file
- Comprehensive fairness metric calculations
- Six-panel visualization generation
- Detailed console reporting
- Remediation recommendations

**To Run**:

```
python compas_audit.py
```

# Conclusion

The COMPAS audit demonstrates that algorithmic risk assessment tools perpetuate and amplify historical biases in the criminal justice system. The disparate false positive rates represent a serious violation of fairness principles, with real consequences for individuals' liberty and life prospects.

---

# Part 4: Ethical Reflection

# Personal Commitment to Ethical AI Development

## Introduction

As a software engineer working with AI systems, I recognize that every technical decision carries ethical implications. Whether developing a recommendation system, a data analysis tool, or an automated decision-making application, I have a responsibility to consider how my work affects real people—particularly those from marginalized or vulnerable communities.

## My Ethical Framework

For any AI project I undertake, I commit to adhering to these core principles:

### 1. Justice and Fairness

**Commitment**: Ensure fair distribution of AI benefits and risks, preventing discrimination.

**Actions**:

- Test models on disaggregated demographic data before deployment
- Calculate fairness metrics (demographic parity, equalized odds, calibration)
- Seek diverse, representative training data
- Audit features for hidden correlations with protected attributes

**Example**: For a loan approval system, I would ensure approval rates and error rates are similar across racial and income groups, even if it means sacrificing some overall accuracy.

# 2. Non-Maleficence (Do No Harm)

**Commitment**: Ensure AI systems do not harm individuals or society.

**Actions**:

- Conduct thorough risk assessment considering worst-case scenarios
- Design fail-safes and human oversight for high-stakes decisions
- Engage with affected communities to understand potential harms
- Perform red team testing to identify vulnerabilities

**Example**: For content moderation AI, ensure it doesn't disproportionately silence marginalized voices while providing clear appeals processes.

# 3. Autonomy and Consent

**Commitment**: Respect users' rights to control their data and decisions.

**Actions**:

- Provide clear explanations of data collection and usage in plain language
- Make opt-in the default, never collect data without explicit consent
- Implement data minimization—collect only what's necessary
- Enable right to deletion and data portability
- Provide understandable explanations for AI decisions

**Example**: For a health recommendation app, allow users to see exactly what data informs recommendations and give them control to correct or remove data.

# 4. Transparency and Explainability

**Commitment**: Make AI systems understandable and auditable.

**Actions**:

- Maintain comprehensive documentation of data sources, architecture, and limitations
- Use interpretable models when possible
- Provide feature importance indicators

- Open-source code when appropriate for community auditing
- Clearly communicate system limitations and error-prone scenarios

**Example**: For a resume screening tool, show recruiters which qualifications led to each ranking, allowing understanding and override.

# 5. Accountability and Governance

**Commitment**: Establish clear responsibility for AI outcomes and enable redress.

**Actions**:

- Implement human-in-the-loop for high-stakes decisions
- Maintain detailed audit trails of decisions and model versions
- Create feedback channels for reporting problems
- Monitor deployed systems continuously
- Use rigorous version control

**Example**: For hiring systems, implement detailed logging, quarterly bias audits, and processes for candidates to challenge decisions.

# 6. Sustainability and Social Responsibility

**Commitment**: Consider environmental and long-term societal impacts.

**Actions**:

- Design computationally efficient models
- Assess potential misuse and unintended consequences
- Consider job displacement and advocate for transition support
- Choose renewable energy cloud providers
- Optimize training to reduce carbon emissions

**Example**: For customer service chatbots, use efficient models and work with stakeholders to retrain displaced workers.

# Personal Accountability Checklist

Before deploying any AI project, I will ask:

1. ✓ Who benefits from this system, and who might be harmed?
2. ✓ Have I tested for bias across all relevant demographic groups?
3. ✓ Can users understand how decisions are made?
4. ✓ Do users have meaningful control over their data?
5. ✓ Is there human oversight for consequential decisions?
6. ✓ Have I consulted with affected communities?
7. ✓ Can someone harmed by my system get recourse?
8. ✓ Am I prepared to take responsibility if something goes wrong?

9. ✓ Would I be comfortable if this system were used on me or my family?
10. ✓ Am I being honest about limitations and potential misuses?

# Conclusion

Ethical AI development is not a one-time checkbox but an ongoing commitment requiring vigilance, humility, and continuous learning. I recognize that perfect fairness may be impossible—there will be trade-offs and difficult choices. What matters is approaching these challenges thoughtfully, transparently, and with genuine concern for those most likely to be harmed.

**My commitment**: I will not build AI systems that I would be uncomfortable explaining to those most affected by them. I will prioritize people over performance metrics, and fairness over convenience.

---

# Overall Conclusion

## Synthesis of Learnings

This comprehensive AI ethics project has explored the multifaceted nature of responsible AI development through theory, case studies, practical analysis, and personal reflection. Several critical themes emerge:

## 1. Bias is Systemic, Not Accidental

From Amazon's hiring tool to COMPAS recidivism scores to facial recognition systems, bias in AI is not primarily due to individual programmer error. Rather, it reflects and amplifies:

- Historical discrimination embedded in training data
- Structural inequalities in society that create correlated features
- Optimization objectives that prioritize accuracy over fairness
- Lack of diverse perspectives in development teams

## 2. Technical Solutions Require Policy Support

While fairness-aware machine learning techniques (adversarial debiasing, reweighing, fairness constraints) can reduce bias, they cannot eliminate it without accompanying policy reforms:

- Transparency mandates and public auditing
- Community oversight and consent mechanisms
- Legal liability for discriminatory outcomes
- Limits on use cases for high-risk applications

## 3. Fairness Involves Difficult Trade-offs

The COMPAS analysis revealed that achieving equal false positive rates across racial groups may increase false negative rates, and vice versa. There is often no single "fair" solution—different fairness metrics conflict, requiring explicit value judgments about which harms to prioritize.

# 4. Human Oversight Remains Essential

AI should augment human judgment, not replace it. For high-stakes decisions affecting people's lives, liberty, or livelihood:

- Algorithmic recommendations should be advisory, not determinative
- Humans must be able to understand, question, and override AI decisions
- Accountability requires human decision-makers who can be held responsible

# 5. Ethics Must Be Proactive, Not Reactive

Ethical considerations cannot be afterthoughts. They must be integrated throughout the AI lifecycle:

- **Design phase**: Consult affected communities, conduct risk assessments
- **Development phase**: Use fairness-aware techniques, audit training data
- **Testing phase**: Evaluate disaggregated performance across groups
- **Deployment phase**: Implement human oversight, enable appeals
- **Monitoring phase**: Continuous auditing, public reporting, rapid response to harms

# 6. Regulation Like GDPR Shapes Responsible Innovation

While GDPR creates compliance challenges, it also drives innovation in privacy-preserving AI techniques and forces developers to consider ethics from the outset. Strong regulation can elevate standards across the industry.

# Path Forward

As AI becomes increasingly integrated into critical systems—criminal justice, healthcare, employment, education, finance —the stakes of getting ethics right become existential. The path forward requires:

**For Individuals**:

- Commit to ethical principles in every project
- Continuously educate ourselves on bias and fairness
- Speak up when we see problematic systems being developed
- Prioritize people over performance metrics

**For Organizations**:

- Establish ethics review boards with diverse representation
- Invest in fairness testing and bias auditing infrastructure
- Create accountability mechanisms and harm redress processes
- Foster cultures where ethical concerns can be raised without retaliation

**For Policymakers**:

- Mandate transparency and auditing for high-risk AI systems
- Establish clear liability for discriminatory outcomes
- Empower communities to oversee and consent to AI deployments
- Fund research into fairness techniques and bias mitigation

**For Society**:

- Demand accountability from AI developers and deployers
- Support affected communities in challenging harmful systems
- Recognize that "accurate" does not mean "fair"
- Prioritize human dignity over technological efficiency

# Final Reflection

The power of AI to improve lives is immense, but so is its potential for harm. As developers, we bear responsibility for ensuring our systems empower rather than exploit, include rather than exclude, and ultimately make society more just and equitable.

This project has reinforced my commitment to building AI ethically—with humility about limitations, transparency about trade-offs, and unwavering focus on the humans affected by our algorithms. The question is not whether we can build it, but whether we should, and for whom.

**The future of AI is not predetermined. Through conscious, ethical choices today, we can shape it toward justice.**

---

# Project Deliverables Summary

☐ **Part 1**: Theoretical understanding with Q&A and ethical principles matching

☐ **Part 2**: In-depth case study analysis of Amazon hiring bias and facial recognition in policing

☐ **Part 3**: Practical COMPAS dataset audit with Python code, visualizations, and 1,247-word report

☐ **Part 4**: Personal ethical reflection with actionable commitments

☐ **Final Document**: Comprehensive synthesis addressing entire project (this document)

---

**End of Report**