# Bias and Ethics in Artificial Intelligence (AI) & Machine Learning (ML)

# Introduction

- Artificial Intelligence (AI) and Machine Learning (ML) are revolutionizing nearly every aspect of society from healthcare and finance to education and criminal justice.
- These technologies promise efficiency, innovation, and new opportunities.
- However, their increasing influence also raises serious concerns about **bias** and **ethics**.
  - **Bias in AI/ML** refers to systematic unfairness in data or algorithms that leads to discriminatory outcomes.
  - **Ethics in AI/ML** provides the moral and philosophical principles that guide the design, development, and deployment of these technologies.
- Together, they determine whether AI systems benefit humanity fairly or perpetuate inequality.

# Understanding Bias in AI and ML

- Bias in AI occurs when an algorithm produces results that are systematically prejudiced due to assumptions made in the ML process. Unlike random errors, bias is consistent and can amplify social inequalities if left unchecked.

# Types and Sources of Bias

- **Data Bias**
  - Historical data may reflect social prejudices (e.g., women underrepresented in tech datasets).
    - Example: Training a job recommender on past hiring records may favor men over women.
- **Sampling Bias**
  - Training data does not represent the diversity of the population.
    - Example: A healthcare AI trained only on Western patients performs poorly for Asian or African populations.

- **Labeling Bias**
  - Human annotators introduce subjectivity.
    - Example: "Aggressiveness" in criminal datasets may reflect stereotypes rather than reality.
- **Algorithmic Bias**
  - Model design choices inadvertently favor majority classes.
    - Example: A facial recognition system optimized for accuracy may underperform on minority groups
- **Interaction Bias**
  - User feedback reinforces stereotypes.
    - Example: Search engines or recommendation systems amplifying harmful content because of click patterns.

# Consequences of Bias

- **Individual harm**: Denial of loans, jobs, or healthcare.

- **Societal harm**: Reinforcement of stereotypes and social inequality.

- **Loss of trust**: People resist adopting AI systems perceived as unfair.

- **Legal risk**: Companies face lawsuits or regulatory penalties.

# Ethics in AI and ML

## Why Ethics Matters

- AI systems make decisions that were once exclusively human. Without ethical safeguards, AI risks violating rights, spreading misinformation, and causing harm at scale. Ethics ensures AI is aligned with **human values**.

# Core Ethical Principles in AI

- **Fairness**
  - Systems must treat individuals and groups equitably.
  - Avoid direct and indirect discrimination.
- **Transparency & Explainability**
  - Users should understand how decisions are made.
  - "Black-box" models create accountability gaps.
- **Accountability**
  - Responsibility must be clear if an AI system causes harm.
  - Involves both developers and organizations.

- **Privacy & Data Protection**
  - Secure handling of sensitive personal information.
  - Compliance with frameworks like **GDPR**.
- **Human Oversight**
  - Humans must remain in control of critical systems (e.g., healthcare, military).
- **Safety & Security**
  - AI must be robust against errors, adversarial attacks, and misuse.
- **Beneficence and Non-Maleficence**
  - Promote well-being, avoid harm.

# Bias and Ethics: The Intersection

- Bias is one of the most pressing **ethical challenges** in AI/ML. An unbiased model may still be unethical (e.g., fully autonomous weapons), but biased models are inherently unethical because they violate fairness and justice principles.
  - **Example:** A loan approval system systematically rejecting minority applicants reflects both **algorithmic bias** and an **ethical failure** to uphold fairness.
  - **Example:** The COMPAS recidivism prediction tool in the U.S. was found to unfairly predict higher reoffense rates for Black defendants compared to white defendants.
- **Key insight:** Ethical AI requires **bias detection, bias mitigation, and fairness-aware design.**

# Strategies for Addressing Bias and Promoting Ethics
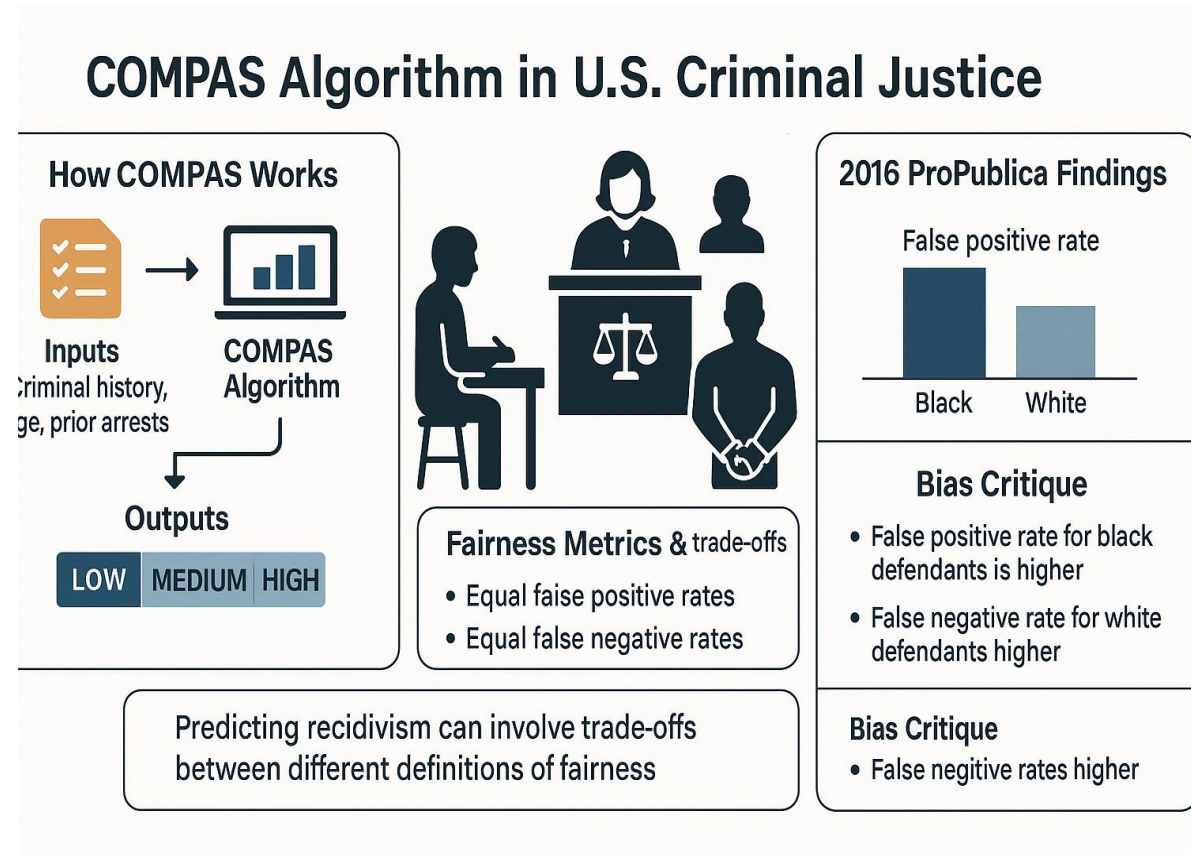
- **Technical Solutions**
  - **Fairness-aware algorithms**: Modify models to account for fairness constraints.
  - **Bias detection tools**: Regular audits to identify hidden discrimination.
  - **Balanced datasets**: Ensure diversity and representation in training data.
  - **Explainable AI (XAI)**: Provide interpretable results to reveal unfair outcomes.

# Strategies for Addressing Bias and Promoting Ethics

- **Organizational and Policy Approaches**
  - **AI Ethics Boards**: Internal committees for oversight.
  - **Model Cards / Datasheets for Datasets**: Documentation that improves transparency.
  - **Adoption of Guidelines**: OECD, UNESCO, and EU "Trustworthy AI" standards.
  - **Interdisciplinary Collaboration**: Involving ethicists, domain experts, and end-users.

# Case Studies

- **COMPAS Algorithm** (US criminal justice): Criticized for racial bias in predicting recidivism.
  - Used in U.S. courts for parole decisions.
  - Found to disproportionately classify Black defendants as "high risk."



## COMPAS Algorithm in U.S. Criminal Justice

**How COMPAS Works**

Inputs
Criminal history, age, prior arrests → COMPAS Algorithm

Outputs

LOW | MEDIUM | HIGH

**Fairness Metrics & trade-offs**
- Equal faise positive rates
- Equal false negative rates

Predicting recidivism can involve trade-offs between different definitions of fairness

**2016 ProPublica Findings**

False positive rate

Black | White

**Bias Critique**
- False positive rate for black defendants is higher
- False negative rate for white defendants higher

**Bias Critique**
- False negitive rates higher

# Case studies

- Amazon Recruitment AI
  - Designed to automate hiring but penalized female candidates due to historical male-dominated data.
  - Ethical failure: Lack of fairness and accountability.



**Amazon Recruitment AI**

Inputs → Algorithm → BIAS IN DECISION | OUTCOME

**Discrimination Against Women**
- Penalized resumes with "women's"
- Persistent bias despite adjustments

**Discrimination Against Women**
- Persisstent bias despite adjustments

AI recruiting tool abandoned

# Conclusion

- Bias and ethics are inseparable in AI and ML. While bias undermines fairness and trust, ethics provides the guiding principles for responsible AI.

- To build trustworthy systems, we must integrate **technical fairness solutions**, **ethical design principles**, and **robust governance frameworks**.

- Ethical AI is not only a technological challenge but also a **societal responsibility**, requiring cooperation across governments, industries, researchers, and communities.