

Part 1: Theoretical Understanding - Expanded Analysis

1. Short Answer Questions - Detailed Explanations

Q1: Algorithmic Bias - Deeper Examination

Definition Refinement: Algorithmic bias refers to systematic and repeatable errors in AI systems that create unfair outcomes, typically disadvantaging certain groups while privileging others. This bias stems from three primary sources:

1. Data Bias - When training datasets are unrepresentative or reflect historical prejudices
2. Model Bias - When algorithm design choices inadvertently favor certain outcomes
3. Deployment Bias - When the context of use differs from development conditions

Additional Examples:

3. Criminal Justice: COMPAS risk assessment software showed racial bias by falsely flagging Black defendants as higher risk at nearly twice the rate of white defendants
4. Healthcare: Pulse oximeters (using AI algorithms) provided less accurate readings for patients with darker skin tones, potentially delaying critical care

Mitigation Strategies:

- Diverse dataset collection and auditing
- Bias detection tools like IBM's AI Fairness 360
- Inclusive development teams to identify potential blind spots

Q2: Transparency vs. Explainability - Extended Analysis

Comparative Framework:

Aspect	Transparency	Explainability
Focus	System-level understanding	Decision-level understanding
Audience	Developers, regulators	End-users, affected parties
Implementation	Open-source models, documentation	Simplified explanations, visualizations
Challenges	Intellectual property concerns	Technical complexity of explanations

Industry Applications:

- Healthcare: Explainability crucial for diagnostic AI (e.g., showing which image features suggested a tumor)
- Finance: Transparency needed for regulatory compliance in algorithmic trading systems

Emerging Solutions:

- "Explainable AI" (XAI) techniques like LIME and SHAP values
- Regulatory frameworks requiring both (e.g., EU AI Act)

Q3: GDPR Impact - Comprehensive Overview

Key Provisions Affecting AI:

1. Article 22 - Right to human review of automated decisions
2. Articles 13-15 - Right to meaningful information about automated processing
3. Article 35 - Data Protection Impact Assessments for high-risk systems

Case Studies:

1. Automated Recruitment: A French company was fined €200,000 for covertly analyzing candidates' personality traits via AI without proper disclosure

2. Credit Scoring: Dutch authorities banned an AI system that used 170+ variables (including social connections) to assess creditworthiness

Ongoing Challenges:

- Balancing innovation with compliance costs (estimated at €20B annually for EU businesses)
- Addressing "black box" AI systems that resist explanation
- Cross-border data flows in global AI development

2. Ethical Principles - Expanded Framework

Extended Definitions and Applications

A) Justice in AI

- Definition: Equitable distribution of benefits and burdens across all demographic groups
- Applications:
 - Fair allocation of healthcare resources via predictive algorithms
 - Equitable access to AI-powered educational tools
- Challenges:
 - Defining fairness metrics (statistical parity vs. equal opportunity)
 - Addressing historical inequities embedded in data

B) Non-maleficence in Practice

- Implementation Strategies:
 - Harm prediction models during development
 - Continuous monitoring for unintended consequences
- Case Example: Facebook's suicide prevention AI reduced response time from 11 hours to 10 minutes while minimizing false positives

C) Autonomy in Digital Age

- Emerging Issues:
 - Dark patterns in AI interfaces that manipulate choices
 - Informed consent for data collection in IoT devices
- Solutions:
 - Granular privacy controls
 - "Right to be forgotten" implementations

D) Sustainability Metrics

- Environmental Impact:
 - Training large language models can emit 300,000 kg of CO₂ (equivalent to 125 round-trip flights from NY to Beijing)
 - Energy-efficient AI chips (e.g., Google's TPUs reduce energy use by 83%)
- Circular Economy Applications:
 - AI-optimized recycling systems
 - Predictive maintenance reducing industrial waste

3. Additional Ethical Considerations

Emerging AI Ethics Challenges

1. Epistemic Injustice

- When AI systems systematically disregard certain ways of knowing
- Example: Voice assistants failing to understand non-native accents

2. Moral Deskillng

- Over-reliance on AI eroding human decision-making capabilities
- Example: Doctors potentially losing diagnostic skills due to AI dependence

3. Value Alignment Problem

- Difficulty encoding human values into AI systems
- Case Study: Autonomous vehicles' ethical decision-making in accident scenarios

Comparative Analysis of Ethical Frameworks

Framework	Focus Area	AI Application Example	Limitations
Utilitarian	Greatest good	Pandemic response algorithms	May sacrifice minorities

Deontologic al	Rule-based ethics	Privacy-preserving AI design	Inflexible in crises
Virtue Ethics	Moral character	AI personality design	Hard to quantify

4. Future Directions in AI Ethics

Technological Solutions

- Differential privacy techniques
- Federated learning for decentralized data
- Blockchain for audit trails

Policy Developments

- Proposed EU AI Act risk classifications
- US Algorithmic Accountability Act
- China's generative AI regulations

Industry Initiatives

- Partnership on AI (Apple, Google, Facebook)
- IEEE's Ethically Aligned Design standards
- OpenAI's charter for beneficial AGI