

# Ethics, bias, and fairness challenges in artificial intelligence and machine learning

Ritesh Rastogi<sup>1</sup>, Nitin Liladhar Rane<sup>2</sup>, Ankur Chaudhary<sup>3</sup>, Jayesh Rane<sup>4</sup>

<sup>1</sup> Noida Institute of Engineering and Technology Greater Noida, India

<sup>2</sup> Architecture, Vivekanand Education Society's College of Architecture (VESCOA), Mumbai 400074, India

<sup>3</sup> Noida Institute of Engineering and Technology Greater Noida, India

<sup>4</sup> K. J. Somaiya College of Engineering, Vidyavihar, Mumbai, India



## Article Info:

Received 09 December 2025

Revised 28 January 2026

Accepted 16 February 2026

Published 19 February 2026

## Corresponding Author:

Jayesh Rane

E-mail: [jayeshrane90@gmail.com](mailto:jayeshrane90@gmail.com)

**Copyright:** © 2026 by the authors. Licensee Deep Science Publisher. This is an open-access article published and distributed under the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## Abstract

The fast development of artificial intelligence and machine learning technologies in the area of crucial social activities has raised the level of concern about ethical interests, algorithmic discrimination, and justice in automated decision-making procedures. This literature review demonstrates that the ethical issues of AI development, bias reduction, and equity systems are complex phenomena, and the development of AI should be performed responsibly. The problem statement focuses on the increasing gap between technological improvement and moral responsibility, where the AI systems reproduce the discriminating trends, violate the rights to privacy, and deliver obscure judgments to influence basic human rights. Based on the PRISMA approach to conducting a systematic literature review, the research is a synthesis of the existing knowledge about the sources of bias in machine learning pipelines, metrics of fairness and trade-offs, ethical frameworks that govern AI, and new regulatory environments. The findings indicate that there remain ongoing problems with the attainment of algorithmic fairness across intersectional groups, between fairness definitions, in terms of black box transparency, and accountability mechanisms of AI-related harm. The crucial shortcomings that have been found are the low standardization of fairness measures, lack of equal representation in training data sets, and the lack of interdisciplinary interaction and emerging regulatory frameworks that are not keeping pace with technological advancement. The proposed review article will add to the overall perspective of any challenges that are present nowadays, practical mitigation measures and future research advancements that would be useful in creating reliable AI systems that meet human values and societal expectations.

Keywords: Artificial intelligence, Ethics, Bias, Machine learning, Responsible AI, Governance.

## 1. Introduction

Machine learning technologies and artificial intelligence have outlived the boundaries of experiments and have become the inseparable elements of contemporary infrastructure, affecting the decisions that determine the life of a person, as well as the future of a whole [1,2]. Healthcare diagnostics and criminal justice risk assessment are still done by AI systems, but so is financial lending and employment screening, which means that access to opportunities, resources, and basic rights are mediated by AI systems more often. This technological revolution is set to provide the capability of never before witnessed efficiency, scalability, and predictability that will have the capacity to resolve the challenging issues facing society [2]. Nevertheless, the pace at which the AI systems are deployed has also revealed some of the key weaknesses in AI systems architecture, development, and deployment with dark questions of ethics, prejudice, and justice emerging.

The ethical aspects of AI constitute a wide range of issues which are not limited to technical performance indicators, but which involve ethical aspects of human dignity, autonomy, justice, and social equity [3-5]. A growing number of places of AIs authority over decision-making in consequential areas have

resulted in the realization of the potential to increase the inequality already present in society, introduce new sources of discrimination, and reinforce past wrongs. Sensational instances of discriminatory effects of facial recognition systems, biased recidivist prediction models, and unjust scorecard models largely drive the popular discussion and academic research on the origins and systemic character of AI-based damages. Up to date this is one of the most urgent problems in the development of AI, the issue of algorithmic bias [2,4,6]. In contrast to bias operated according to the individual prejudices and processing heuristics, algorithmic bias is systematic and is repeatable, and values in data, model configurations, and deployment environments [2,7-9]. These biases may arise at all possible points along the machine learning pipeline, such as data collection and annotation, feature engineering, training the model, and monitoring what happens after deploying it to production. Modern AI systems, especially deep learning systems with millions or even billions of parameters, are so technical as to introduce an element of such opacitance that makes it difficult to discover and fix any discriminatory patterns.

The fairness in machine learning has developed as a progressive introspective of the form of fairness to an intuitive idea of fairness as equal treatment, to a more complex subject of study, struggling with mathematical formalization, alternative definitions, and the possibility of absence of results [10]. Various fairness criteria suggested by researchers include: demographic parity, equalized odds, predictive parity and individual fairness which represent various normative standards of how algorithm systems ought to treat various groups and individuals [10,11]. But mathematical arguments have shown that various fairness conditions can not both be met with the exception of trivial examples, then practitioners must make hard choices between incompatible ideas of justice. The AI-related ethical governance and equity have special challenges beyond the conventional regulatory methods [12-14]. The international character of AI development, the dynamic character of technological change, the two-sided usage character of AI functions, and the putting together of AI knowledge into the groups of the corporate world makes the problem of establishing effective mechanisms of oversight more complicated. There have been a number of ethical frameworks, principles, and guidelines suggested by governments, international organizations, industry consortia and civil society groups and so far, there has been the struggle of translating principles at high-level into practical technical standards and binding regulations.

The overlapping of AI ethics, bias and fairness are also the occasion to pose the basic inquiries of what is valued, of power as well as representation within technological frameworks [3,15-17]. Who should provide the values guide AI development? Which do we know about including the varying stakeholder perspectives to technical design processes? How do we guarantee accountability in cases of harm by AI systems? These questions show that AI issues are fundamentally sociotechnical, and interdisciplinary thinking is necessary to combine technical innovation and the knowledge of philosophy, law, social sciences, and communities impacted [18-20]. The recent advancements in AI functionality, and more specifically in large language models, generative AI and multimodal systems, have produced new ethical complexity. These systems have emergent behaviors that are hard to predict or manage, pose new questions on authorship and intellectual property and pose a threat of misinformation, manipulation, and malicious use in ways previously unimaginable [21-23]. The open platforms and tools provided by the democratization of AI have reduced the barriers to implementation and at the same time have spread the responsibility and made it harder to govern it [9,24,25]. The need to act urgently in the ethics, bias, and fairness arena in AI is in the fact that technology is spreading increasingly to the originally human-dominated machine fields. The propensity toward collective sense-making and the democratization of processes relies more and more on the role that AI systems play in social interactions as well as the culture production, political discourse, and knowledge creation. The possibility of AI redefining the labor market, increasing the economic disparity, and centralizing the power to technologically sophisticated countries introduces geopolitical aspects to ethical issues.

Even though there has been an increase in academic interest in AI ethics, bias, and fairness, a number of significant gaps still exist in the research and application. One, the research on bias in supervised learning has been done quite extensively; the issues of fairness in unsupervised learning, reinforcement learning and generative models have not been extensively examined. Second, despite the seriousness of intersectional identity and compound marginalization, majority of the fairness studies investigate

prejudice on a single demographic basis. Third, empirical validation of adoption of the bias mitigation strategies in practice deployment situations has not been sufficiently explored, and most of the literature utilizes benchmark datasets that might not represent complexities in operations. Fourth, the temporal dynamics of bias (such as concept drift, feedback loops, and long-term societal effects of biased systems) has not been appropriately covered in the current literature.

This study contributes a number of significant materials to the ethics, bias, and fairness in AI and machine learning. First, it makes an in-depth synthesis of multidisciplinary thinking combining technical outlook and knowledge on ethics, law, social sciences and critical studies. Second, it provides a methodical overview of bias in various kinds of machine learning systems such as new AI paradigms like foundation models and generative AI. Third, it also offers specific comparative studies of fairness indicators, biases alleviation measures, and assessment frameworks offering evidence-based advice to practitioners in complicated trade-offs. Fourth, it puts into focus the sociotechnical aspects of AI fairness, where responsible AI development will rely on the stakeholders and their engagement, understanding contexts, and value alignment. Fifth, it determines research gaps and areas of concern in the future, having a roadmap of the further development of both theoretical and practical perspectives of ethical AI systems.

## **2. Methodology**

To make the analysis of the currently available studies in the domain of ethics, bias, and fairness in artificial intelligence and machine learning systematic, transparent, and reproducible, this extensive literature review leads to the application of the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) technique. PRISMA framework gives systematic information on the identification, screening and synthesis of the pertinent academic materials and minimizes the selection bias and complete coverage of the research area. The accurate inclusion use and exclusion criteria were defined, and the review process started. The qualified articles included peer-reviewed journal articles, conference proceedings, technical reports, white papers and authoritative grey literature that were published between 2018 and 2025, as this was the period of time used to find new trends and new changes, as well as consider the older literature where applicable. The search queries were aimed at representing the various facets of AI ethics, algorithmic bias, and machine learning fairness and included Boolean operators to combine terms of interest in artificial intelligence, machine learning, bias, fairness, ethics, accountability, transparency, and discrimination. The literature search has been performed using various academic databases and digital libraries with the purpose to cover as much as possible; the results included IEEE Xplore, ACM Digital Library, arXiv, Google Scholar, and specialized repositories on the topic of AI ethics and responsible technology. The search strategy consisted of the use of combinations of keywords such as artificial intelligence ethics, machine learning fairness, algorithmic bias, bias mitigation, fairness metrics, responsible AI, explainable AI, AI governance, and similar words. The first screening was based on title and abstract review in order to determine the relevance and then the second screening was on the full-text review of the prospective eligible studies. Information mining was aimed at sorting literature by thematic area such as sources of bias, the definitions of fairness and metrics, methods of bias detection and mitigation, ethical frameworks and principles, governance and regulatory methods, areas of application, and empirical research. The criteria of quality assessment were the methodological rigor, empirical validation, theoretical contribution and practicability of reviewed studies. They conducted thematic analysis in order to find patterns, contradictions, and gaps throughout the literature, and in particular the emerging trends and gaps that were not studied. The stage of review involved a process of refinement so as to cover as many areas as possible and at the same time remain focused on the most interesting and productive works to the field.

### **3. Results and Discussion**

#### *3.1 Bias in AI and Machine Learning Systems*

Prejudice in artificial intelligence and machine learning systems is a complex process that manifests itself as a complex interaction of data, algorithms, human choices, and contexts of deployment [26-28]. It is the knowledge of what, where, and how biases are that can lead to the creation of mitigation strategies and the creation of more equitable systems of AI.

##### **Sources and Origins of Algorithmic Bias**

Algorithms have various origins during the machine learning lifecycle that are causes of algorithmic bias. Historical bias illustrates the institutional disparities and discriminatory patterns implicated within the training information that reflect the past injustices and the societal bias [6,29-31]. Learning AI systems on data of previous hiring, criminal justice results, or patterns of lending practices is potentially discriminatory and will perpetuate past patterns that have systematically discriminated against some demographic groups [32,33]. This type of prejudice is quite pernicious since the training data can be fully reflective of historical truth, but learning about this truth implies the entrenching of unjustified practices into the system [34-36]. Representation bias is bias in which specific populations are underrepresented or overrepresented in training data. The result of this imbalance is the production of models that do not work well with the minority groups but give a high degree of accuracy on the majority populations [16,37-40]. The errors of facial recognition systems that are trained on the photos of mostly light skinned individuals are a lot higher in processing the photos of darker skinned people, especially women of color. On the same note, the linguistic diversity and cultural views of other settings can be poorly reflected through the use of natural language processing models that are mostly trained when using text representing Western, educated, industrialized, rich, and democratic societies [41-43]. Measurement bias is due to decisions regarding which features to measure, operationalization of abstract features and the selection of proxies when the latter is impossible or even illegal. Making the complicated realities of social constructs measurable and subject to evaluation is inevitable with simplification and possible distortion. As an example, the option to base the treatments on zip codes as their proxy of the socioeconomic status can result in redlining priority, whereas standardized test scores can reflect test-taking ability and access to education instead of innate ability or potential. The poor representation of the heterogeneity of different subgroups in the models lead to the aggregation bias, which is an attempt to approximate different populations, with one model being relevant enough to represent it. A medical diagnostic system that has been trained on the aggregate data might fail to recognize the difference in the manifestation, progression, or responsiveness to medical therapy of diseases among different population groups, resulting in a suboptimal assignment of care to the populations whose features do not match the majority distribution during training. Evaluation bias is a problem in which the benchmark data or evaluation metrics do not provide sufficient evaluation in a variety of populations or scenarios. When validation datasets share the same biases as training data or when the evaluation measures focus on the overall performance as opposed to the even-handed performance across groups, biased models might come out as successful during the development stage but discriminatory when applied to population groups. The mismatches between the contexts of development and the operational contexts or misuse of their models on tasks or population other than that in which they were constructed are what bring about deployment bias. A risk-assessment instrument practiced and supported in one location might be biased when applied in another geographical, cultural, or institutional place using different base rates, population demographics, or decision-making systems.

##### **Understandings of Bias and its different types and manifestations**

The types of mechanism and the effects of an algorithmic bias vary, and so do the effects on different groups [44,45]. Statistical bias is the difference between predicted and actual results of a model that is consistent between demographic groups. Such bias can be measured using fairness metrics that are used to compare the error rate, the false positive rate, the false negative rate, or predictive accuracy between the different groups that are being protected. The social bias involves a wider scope of discriminatory patterns that uphold and enhance society prejudices, stereotypes, and power differences. The

occupational stereotype of some professions being attached by specific genders might be produced by language models and reinforce the stereotypes. Image generation systems have the potential to give out results that help to reinforce racial or gender stereotypes based on the training data used, and can in fact be used to strengthen stereotypes and biases via selective generation. Through the interaction between people and AI systems, interaction bias can be formed with first bias strengthened and increased through a feedback mechanism. Recommendation systems can be the source of filter bubbles because they will recommend content that is similar in nature to past views on the user, depriving them of a varied viewpoint. The search engines can prioritize the results in a way that promotes the mainstream discourses to the exclusion of other opinions.

Allocation bias arises in cases where AI systems make the decisions concerning the allocation of resources, opportunities, or services such that they favor certain groups in a systematic manner. Resume screening software can potentially reject qualified people in underrepresented groups. Even creditworthy individuals in some neighborhoods can be refused credits by credit scoring algorithms as a result of the neighborhood. The systems of healthcare resource allocation can select some categories of patients, leaving the rest of them behind depending on the subjective risk estimates. Quality-of-service bias is a type of inequality in the level of performance between people of different demographics, with the systems showing effective results with people of majority groups and not with the minority groups. There are increased error rates of speech recognition system when speakers have an accent which is not similar to the majority training distribution. The concept of autonomous cars will be not able to identify the dark-skinned people, as there are weaknesses with their sensors and training data. Representational harm prevails when systems liter a stereotype, convey or misunderstand some groups, or encode exhortative connections. The presence of image tagging systems, which are wrong at recognizing the images of individuals belonging to a particular ethnic group or same-sex couples, leads to representational harms. Text generation language models that link some demographic populations with the negative qualities lead to dignitary harms even without affecting consequential decision-making.

### *3.2 Frameworks and Metrics of Fairness.*

Formalizing fairness in machine learning has spurred a wide amount of research building mathematical models of how to operationalize intuitive concepts of equitable treatment. Nonetheless, this writing has also shown the basic tensions and impossibility outcome that make the realization of algorithmic fairness challenging.

Fairness to individuals of an organization versus fairness to a group

Individual fairness holds that individuals are supposed to be treated in the same way and algorithmic systems must generate similar decisions that must be made on individuals who are similar in pertinent aspects [22,30,46-48]. This method coincides with the law of equal protection and non-discrimination against each other with consideration to individual merit. But to apply the individual fairness theory, there is a need to define meaningful metrics of similarity that encompass meaningful similarities and leave out protected attributes, which is a problem that is both technical and normative with respect to which attributes to regards as being meaningful [49-51]. Group fairness aims at statistical parity or equalized outcome of demographic groups within the categories of demographic characteristics of the very attributes like age, gender, race, etc. which are considered as safeguarded attributes. This is because it acknowledges that the structure and past discrimination impacts groups as a unit, and the injustice can be addressed by making right those disadvantages that are systematic, instead of merely treating people equally. Group fairness measures juxtapose statistical measures between groups that are being defended against so that an algorithm does not discriminate against or favor certain groups. The conflict between individual and group fairness contains more profound philosophical conflict concerning the concept of justice and equality. On the one hand, individual fairness focuses on the treatment on merit and formal equality, whereas, on the other hand, group fairness is concerned with substantive equality and corrective justice. No one approach can be considered better than the other, and the necessary framework

will have to be determined by the particular field of application, the rules of law, and the values of stakeholders. The criteria of mathematical fairness are as following.

Demographic parties entail that the ratios of those who got a positive outcome should be equal in the case of the groups being protected. This criterion is to make sure that there is an equal distribution of opportunities or resources among groups irrespective of their features. Nonetheless, demographic parity can come into conflict with accuracy in situations whereby, relevant predictive features are linked with protected characteristics and also it does not consider actual differences in qualification or needs between groups. The equalized odds state that both of these rates received the true positive and false positive rates, on groups, should be equal, such that true positive and false positive rates are also equal irrespective of being in a protected group. The criterion focuses on equal quality of service thus eliminating cases where models are good with the majority groups, but not with the minorities. Equally satisfying odds can however mean compromising the overall performance or even settling with lower performance of all groups. Predictive parity demands that the positive predictive value is the same in every group i.e. those individuals who are subjected to positive predictive values have an equally good chance of actually falling into the positive category irrespective of group membership. This criterion counts on the fact that predictions have the same meaning in groups. Predictive parity may, however, also interfere with other measures of fairness where the base rates vary with groups. The calibration condition is that predicted frequency of results in each group has empirical frequencies that are accurate. A calibrated model makes sure that the probability of an outcome to be assigned is correct to individuals of all groups. Calibration is of interest especially in areas such as medical diagnosis or risk assessment where the probability predictions are utilized in making decisions. Nevertheless, disparate impact is not always avoided by calibration when there are dependency differences in the models that place models with different probability distributions on different groups.

Individual measures of fairness that depend on similarity would demand that the distance between predictions of two individuals be circumscribed by their distance in feature space based on a task-specific measure. This strategy institutionalizes the instinct that like people ought to be treated in similar ways. Nonetheless, a definition of appropriate similarity measures is still a difficult task especially when the relevant similarities require context-specific decisions regarding which aspects an individual ought to incorporate in making decisions. Counterfactual fairness presupposes that individual predictions would not change in case their arguments concerning protecting properties were not the same, everything held constant.

This requirement represents the sense that judgment must not require consideration of the insulated traits. Nonetheless, to impose counterfactual fairness, causal models are needed to think about the possible alternative outcome in the event of attributes being manipulated by the protected factors and possibly fails to deal with discrimination that is proxy-based, or correlated.

#### Impossibility Results and Trade-offs

Such mathematical analysis has also shown that there are some fundamental impossibility results that limit the attempt to achieve a variety of fairness criteria at the same time [52-55]. Unless there is perfect predictability, and equal base rates in the groups that we are interested in, this cannot be done mathematically, even in the face of demographic parity, equalized odds, and predictive parity. These existence of impossibility implications are severe and compel practitioners to make challenging decisions regarding the criteria of fairness to consider depending on the context of application, values of stakeholders, and the law. Another tension which forms a part of fair machine learning is the accuracy-fairness trade-off. Fairness constraints generally have the effect of decreasing the predictive accuracy of overall optimization when compared to unconstrained optimization. The seriousness of this trade-off would depend on the relationship between the attributes that are being protected and legitimate predictive attributes, the fairness rat thus being imposed, and the rates of outcomes within the groups. In fairness constraints, the accuracy costs are very low in one application whereas in other applications, the trade-offs are high. Various stakeholders might differ as to the most suitable criteria of fairness or the most favoured way to be involved in tradeoff between incompatible goals. The developers, users, those affected and society in general might hold different interests and values resulting in difference in



preferences toward the fairness constraints. A participatory approach where stakeholders are involved to create their definition of fairness needs can be used to overcome these conflicts, but it also creates problems of representation, power relations, and the need to reconcile the differing views.

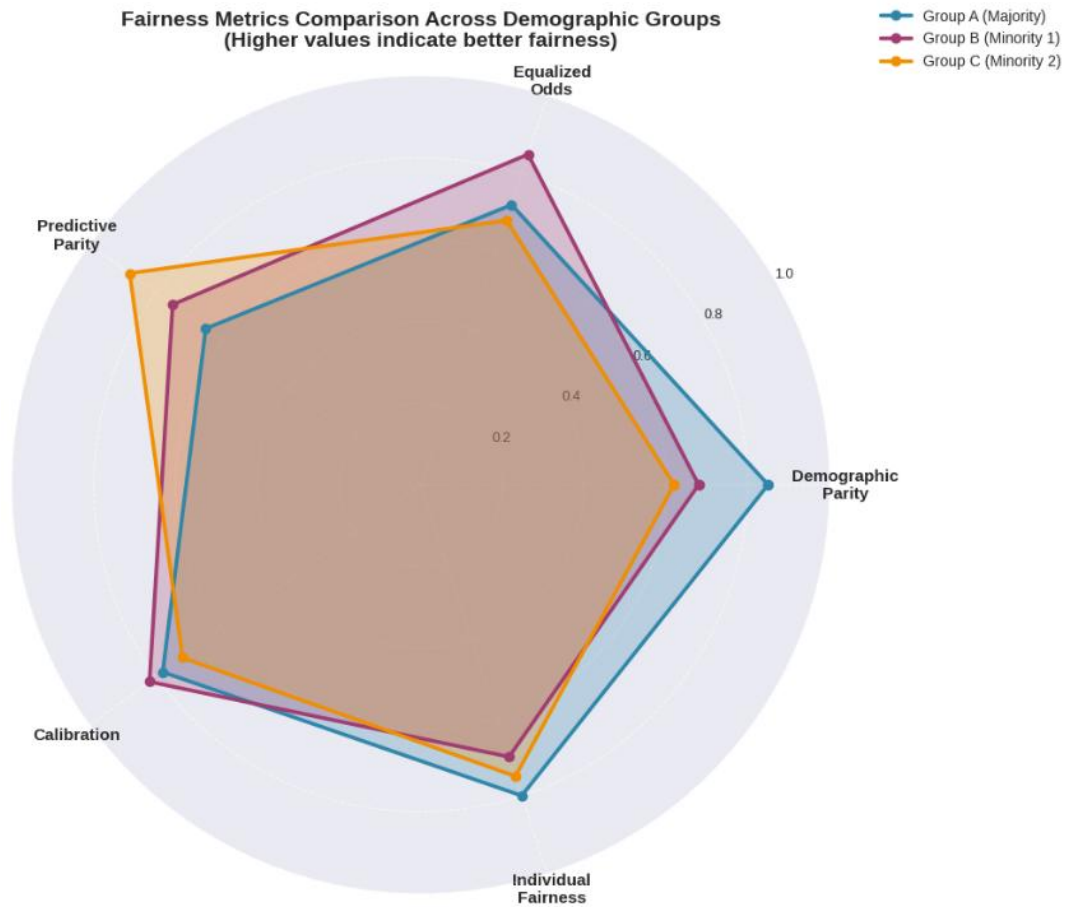


Fig 1: Fairness Metrics Comparison Across Demographic Groups

Fig. 1 compares multiple fairness metrics (Demographic Parity, Equalized Odds, Predictive Parity, and Calibration) across different demographic groups. The radar chart shows how well an AI system performs on different fairness criteria for each group. Values closer to 1.0 indicate better fairness. The plot reveals that achieving all fairness metrics simultaneously is challenging - for instance, Group A shows strong demographic parity (0.85) but weaker predictive parity (0.65), while Group C demonstrates the opposite pattern, illustrating the inherent trade-offs between different fairness definitions.

The values and preferences of the affected communities, legal requirements, certain harm which one type of errors can cause, and the application domain which should be considered are the factors that should guide the selection of fairness criteria. In criminal justice, risk assessment can take preference on reducing racial differences in false positive rates which may result to wrong incarceration whereas medical diagnosis may take preference on calibration where the predicted disease rates should be equally informative regardless of the demographic group. Such context-specific views serve to emphasize the fact that the term of fairness is not a universal principle but, instead, it has to be carefully examined in terms of the specifics of every application and its ethical aspects.

### 3.3 Techniques of Bias Detection and Measurement.

The accumulation and detection of bias in AI systems necessitate advanced frameworks that may reveal unfair tendencies in various points of machine learning and the variety of their unjustified treatment.

#### Data Auditing and Analysis

Data auditing is the process of methodical reviewing of training data to establish possible sources of bias prior to them being coded into models [23,56,57]. The demographic analysis evaluates the coverage of various groups in the datasets, with underrepresentation or excesses of different groups that may cause variations in performance [58-61]. The statistical analysis reveals the correlations between the attributes used as protection and the target results that may be some historical bias or discriminating trends in generating the data. Label auditing addresses attributes of quality and consistency to annotations, especially when subjective tasks are being analyzed, and biases of human annotators can have effect on labeling judgments. The comparative analysis of disagreements between a set of annotators can indicate systematic variability in the perception or categorisation of annotators in a particular way, which may be due to the influence of a given culture, individual prejudice or a lack of clarity in the annotation instructions. By comparing the inter-rater agreement rates by various instances of different kinds, one can find out, whether some specific examples are easier to interpret subjectively. The tasks of feature analysis involve exploring the relationships between input features and those features that are to be protected, what the possible proxy variables might be that would be used as a proxy of the protected features even in the instances of missing these characteristics in the training data. They may be used in correlation and mutual information metrics and causal analysis, which can identify the features that are informative about the attributes that are being protected, and these may be useful to make a single discrimination by using indirect routes. Counterfactual data generation is an approach used to generate synthetic examples, which are examples that are modified in a manner that the rest of the features are kept constant and the relations between demographic characteristics and data distributions as well as model predictions are evaluated. This method may assist in determining the causal association between the attributes under protection and the outcomes against the spurious associations caused due to the presence of confounding factors.

#### Model Auditing Techniques

Post hoc model auditing audits trained models to determine their fairness criteria as well as discriminatory tendencies of predictions [62-64]. Disaggregated evaluation is calculating performance measures in distinct groups of people and show differences in the accuracy, precision, recall or other measurements of quality [1,65,66]. There are systematic dissimilarities in the rate of error among groups that could suggest bias that has to be investigated and corrected.

Fairness metric calculation is a measure of different mathematical fairness measures to determine whether a particular model is able to meet a certain fairness requirement. Computing demographic parity, equalized odds, predictive parity and calibration measures among the groups that have been protected gives a multifaceted evaluation of equal conditions of fairness. Demonstrations of the comparison of such metrics between groups can serve to determine what fairness criteria is breached and by how much. Adversarial testing is a technique that specifically seeks inputs which can reveal unfair practices or identification of discriminating tendencies in model predictions. There is a possibility of generating edge cases, adversarial examples, or stress-tests approaches to uncover the hidden biases which might not be immediately noticeable in benchmark evaluation. It can be useful to systematically manipulate and control the presence of the attributes of a system by setting other features fixed to understand whether the predictions are dependent on the protected characteristics in an inappropriate manner. Interpretability and explainability methods allow knowing the logic behind model predictions and the characteristics behind discriminating performance. Looking up the analysis of the feature importance shows what input variables contribute to most to the prediction, which may reveal exposure to features that are under protection or proxies. Decision rule extraction tries to generate human interpretable rules approximating model behavior, which is simpler to understand the outcome of discriminatory trends. Causal analysis methods examine causal correlations among features, guarded attributes and results, assisting in differentiating certain forecasting connections with outcomes that are lawful and those that are discriminating. Causal graphs record the structural relationships between variables thus helping to find alteration routes of discrimination which are direct and indirect. Counterfactual reasoning evaluates the effect that the elimination of causal effects of litigated attributes would alter predictions, which demonstrate discrimination.



## Constant Monitoring and Evaluation

Implemented AI systems would need continuous monitoring in order to identify new biases that might appear due to the distribution shift, feedback loops, or the change of the conditions under which it operates. Continuous fairness monitoring provides fairness measures over time, exposing fairness properties developed by a developer to slumping that might need action. On-demand dashboards that visualize inequity indicators at the level of demographic groups allow identifying problematic trends as quickly as possible. The feedback loop analysis looks at the effect of model prediction on future data distribution and how the processes help to amplify the initial biases. In systems where predictions influence the behavior of users or influence their selection, preferential exposure or selective labels or strategic responses may be used to entrench initial biases. Tracking the distributions of data and outcome differences through time does assist in identifying these reinforcing cycles. Drift detection helps to detect data distribution changes or model changes, which can indicate the appearance of bias problems. Concept drift takes place when the association between features and outcomes varies with time, which may have a different impact on varying groups. Covariate shift is the effect that alters feature distributions that could affect a minority group disproportionately provided those models make poor extrapolations to areas of feature space with few train data. The incident response protocols lay out a series of processes to show how bias incidences are to be investigated and dealt with in case they are reported or detected. Processes of systematic investigation establish deep-seated causes of discriminatory results, be it of data problems, model flaws, or context mismatches in their utilization. Remediation plans provide actions of correction such as retraining models with revised information to change the decision-making procedures or introduce human control.

### *3.4 Bias Mitigation Strategies*

To deal with bias in AI systems, the intervention should take place at several phases of the machine learning pipeline, including the stage of data acquisition, as well as model design, deployment, and monitoring [67-69]. Good mitigation plans entail both technical plans and process enhancement with organisational transformations.

#### Pre-processing Techniques

Pre-processing methods alter training data and then model development in order to mitigate bias and enhance the fairness properties of developed models. The resampling methodologies can re-represent unevenly represented groups in training data by either oversampling underrepresented groups, undersampling overrepresented groups, or synthetic data representation. These methods can be used to even out datasets and minimize representation bias, but again, they can also introduce some other types of artifacts or are not able to resolve the underlying data quality problems. Training examples are generated to augment the data in order to add an additional representation of the minority groups or the few cases that the models do not perform well at. The augmentation methods should not alter any of the characters important but only add enough variety to enhance generalization. One should take care not to augment data in the way that it symbolizes authentic diversity instead of perpetuating less diverse stereotypes. Reweighting places a weight of varying significance on training data of various demographic groups to the training process, giving them more weight on the minority groups. This strategy is able to enable models to learn to work across all groups without necessarily varying data distributions. Nonetheless, reweighting must be carefully tuned so as to favor the outcomes of enhancing fairness but not compromising the accuracy.

Direct Discrimination via proxy variables may be lessened by feature transformation and encoding since it is possible to decrease the correlation between attributes that are being removed and other features. The univariate non-response of features (data points) which are proxies of a hidden characteristic can be used to indeed stop model learning patterns of discrimination. Nevertheless, within this approach, there is a likelihood of compromising predictive accuracy in legitimate discarding features and there is a possibility of low efficacy in the event that a set of weak proxies works together to discriminate. Fair data generation makes the use of new training information that is specifically set with the aim of promoting fairness. Specific outreach to underrepresented populations will be effective in enhancing

the proportion of the groups, as well as capture a variety of views and note a wider variety of situations. Data can be better represented through participatory data collection which can involve the affected communities so that the data is more accurate in terms of their experiences and priorities.

### In-processing Methods

In-processing methods are an adjustment of model training algorithms to promote fairness objectives in addition to predictive accuracy. Fairness-obligated learning takes fairness constraints or penalties into the objective function and this aims to promote the model to meet fairness requirements without compromising its good predictive accuracy. Regularization penalizing unfair predictions provides smooth trade-offs between fair and high accuracy which can be adjusted depending on the requirements of the application.

Adversarial debiasing Adversarial networks are trained to predict outcomes and the second branch tries to predict learned representations as the protecting attribute [70,71]. This adversarial system is designed to motivate the models to learn representations that are predictive of the results but do not predict the information on the protected attributes such that fairness is promoted in the process of representation learning [20,72-74]. The adversarial style is robust with complicated models such as deep neural network where it is infeasible to create features seeming and encompassing its context in the formulation of a feature engineering model. Constrained optimization develops model training as optimization problems having fairness constraints made explicit, that has to be achieved. The methods ensure that trained models satisfy desired amounts of fairness but can scale computationally very expensive in any case. The selection of constraints identifies the fairness standards to be implemented and the trade-off among various objectives to be made. Met algorithms methods are those that integrate a collection of models or methods of decisions to imbue them with fairness properties which might be not possessed by individual models. Ensemble mechanisms can work to enhance equity by decreasing cross-demographic variance in prediction or directly integrating those models that are optimized on distinct subpopulations. With methods of meta-learning, it is possible to learn better methods of adapting models to obtain improved properties of fairness to models across tasks or domains. The causal modeling systems entail the use of causal relationships between variables in model structures, which allow principled direct versus indirect discrimination reasoning. Counterfactual causally-motivated criteria of causal fairness can be integrated into training goals and the predictions must meet counterfactual-motivated causal fairness assumptions. These methods need domain knowledge to define causal structures, but can offer greater guarantees of fairness than more statistical methods.

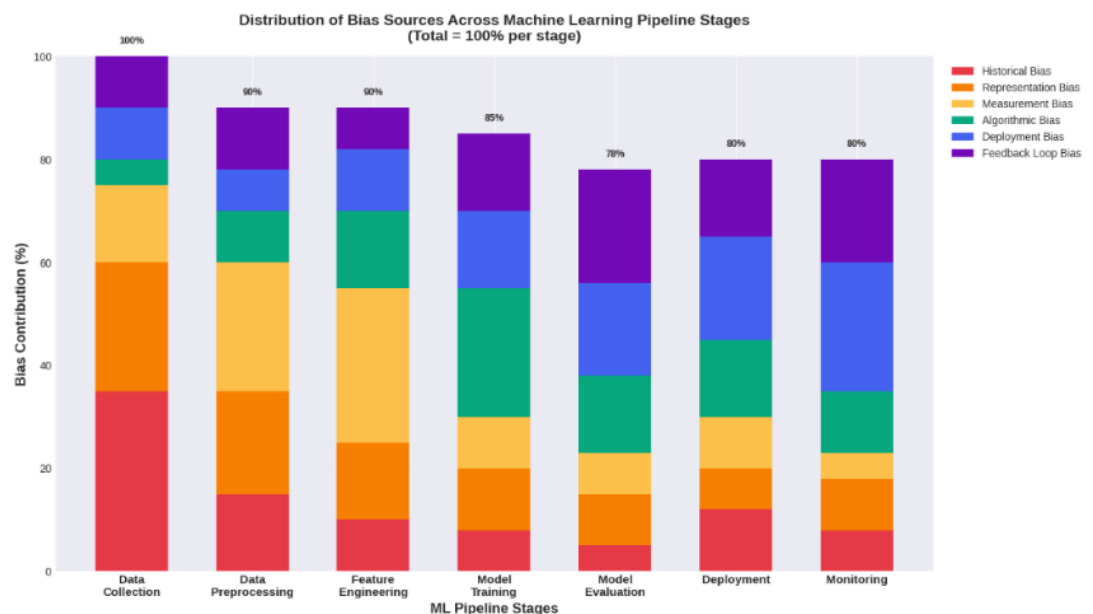


Fig 2: Bias Sources Distribution Across ML Pipeline Stages

Fig. 2 illustrates the relative contribution of different bias sources across various stages of the machine learning pipeline. The data shows that data collection contributes the highest bias (35%), followed by model training (25%) and deployment (20%). Visualization emphasizes that bias mitigation requires interventions throughout the entire pipeline rather than focusing on a single stage. The cumulative view helps practitioners understand where to prioritize their bias mitigation efforts for maximum impact.

#### Post-processing Strategies

Post-processing is applied to alter the outputs of trained models in order to meet fairness constraints without re-training models [75,76]. The concept of threshold optimization can be used to set different decision thresholds between demographic groups to balance the error rates, false positive rate or other measures of fairness. This is computationally efficient and can be applied to any pre-trained model, but has a probability of yielding different quality of service to groups should there be a systematic difference in the calibration properties of the models. Output calibration is the act of increasing or decreasing the prediction probabilities to achieve calibration within each of the demographic groups such that the outcome frequencies in the outcomes as predicted match that which is observed in the data. Fairness properties that exist with regard to meaning and reliability of predictions among groups can be enhanced using calibration techniques. Nonetheless, the process of calibration might not endure other equity issues such as the idea of demographic parity or equalized odds.

The classification of reject option enables the models to avoid making predictions on the hard cases wherein they are not very confident which may decrease the difference in errors [36,77-79]. Selective prediction can make the situation fairer by assigning borderline cases to human decisions which might be more informed with other contexts. Nonetheless, this solution should be designed to make sure that the rates of abstinence are not distributed systematically across the groups in a manner that would cause new fairness concerns. Algorithms based on fair ranking list ranked outputs in a different order to meet fairness requirements of representation of various groups in high position. These methods are especially applicable in the context of the information retrieval, recommendations services, as well as candidate selection when the rankings play a crucial role. Fair ranking should be fair in terms of relevance and representation wherein the highly qualified candidates across all groups are given the right chance to be seen. Ensemble reconciliation is a method that employs predictions obtained by a number of models that have been trained either with varying fairness goals or employed to optimize against various subpopulations. Weighted averaging or stacking methods have ability of generating final predictions to optimize competing fairness requirements, over and above that of individual models. This method offers the capability of trade-offs of fairness through probing ensemble weights according to the priorities of the stakeholders.

### 3.5 Principles and Frameworks of Ethics.

Ethical frameworks have given a normative direction to the responsible development of AI procedures, describing values, principles and considerations that might influence the technical decision making and inform the organizational culture.

#### Core Ethical Principles for AI

Beneficence leads to the need to design AI systems and implement them in ways that are beneficial to people and society, which enhance the wellbeing and prosperity of humans [80-83]. According to this principle, AI has a potential that is positive to deal with significant issues, enhance decision-making, and human abilities. Nevertheless, to achieve positive AI, one should take into account the interests to be prioritized and allocate the benefit to various populations [84-86]. The concept of non-maleficence requires the AI systems not to harm any individual, community or society. This principle covers physical harm, psychological harm, economic harm, the social harm and dignitary harm. When it comes to preventing harm, it is necessary to predict the adverse effects of AI application, take precautions against possible abuses, and attract the means of accountability in the event of any harm. Autonomy acknowledges that people have the right to take informed choices regarding their lives and that they can have a meaningful human agency in situations when AI affects the consequences. Transparency

regarding the moment of AI utilization, the possibility to challenge such decisions made by machines, and the concern of the possible consequences decisions should be under human control, which is supported by this principle. Keeping autonomy is especially relevant when the AI systems start to mediate access to the information, opportunities, and services more and more. Addressing issues of justice would mean equitable allocation of benefits and harms of AI systems whereby technological benefits are received by already advantaged groups and harms and risks are unequally distributed to the disadvantaged groups. This principle covers distributive justice in terms of resource allocation, procedural justice in the terms of decision making and corrective justice in terms of remediation of the harms. Explicability requires the operations of AI systems to be comprehensible to the stakeholders concerned so as to have significant human control and responsibility. In this principle, it is important to provide transparency regarding the abilities and constraints of the system and explain the decisions made about individuals in an interpretable manner, and delineate the responsibilities of the system activities are clear. The various stakeholders might demand various kinds and amounts of the type of explanation that suits them.

#### Value-Sensitive Design

Value-sensitive design methods consider the ethical factors and the value of the stakeholders involved in the design and development process as opposed to the perceived add-on aspect of ethics [6,87,88]. This approach will imply recognizing the stakeholders concerned, evoking their values and issues, and transforming values into technical specifications and achieving design refinement as an iterative process to fit more closely ethical purposes. Participatory design involves the involvement of various stakeholders especially those who have the highest likelihood of being impacted by the AI systems in the formulation of development priorities and design decisions. This strategy acknowledges that the technical professionals might be unaware of some of the most important contextual information and that the technology users must be able to express their opinion about shaping the technology. Effective involvement means the establishment of open avenues to feedback, curbing power tension, and showing the stakeholder feedback as an input to decisions.

Abstract ethical principles may be converted into technical development requirements and constraints through values specification. It consists of converting the high level values such as fairness or privacy into measurable properties, testable criteria and verifiable constraints. Specification of values presents tensions which need to be navigated among the various values, explicit trade-off to be made, as well as, requirements translated to the requirements of a particular application. Ethical impact is a systematic process through which ramifications of AI systems on a number of fronts such as individual rights, social equity, environmental sustainability and democratic values are viewed. Such tests help to determine possible harms, involved groups, protection measures, and monitoring needs. Conducting impact assessment on a regular basis during development and deployment of the product helps to identify and resolve ethical issues in advance.

#### Ethics and the Culture of the Organization as a Professional

The AI professional ethics include ethical duties to act in the best interests of the population, to act without compromise in technical competence, truthfulness in communication with both capabilities and constraints, and ethical standards to promote ethical behaviors despite institutional pressures [6,19,87,88]. Professional codes of conduct offer directions on how to address ethical dilemmas and they offer anticipations of how to practice responsibly. Organizational culture influences the attitude of the ethical consideration value attached, debated and incorporated in the decision making process. Companies that are focused on responsible AI development foster a culture that enhances the disclosure of ethical concerns, implement resources to respond to fairness and safety, recognize ethical actions, and introduce accountability to ethical performances. Organizational culture on AI ethics is affected by leadership commitment, institutional structures as well as the incentive systems.

Ethics training and education enhances the ability to recognize and deal with ethical challenges along the AI development life volumes. Effective training extends beyond the normative principles to offer practical advice, case studies, and decision models of the work of the practitioners. Interdisciplinary training unites the technical and humanistic viewpoint wherein teamwork is encouraged between

disciplines. The use of whistleblower guarantees and ethical escalation processes is also a way of assuring the employees that they can easily speak out concerning unethical practices without any fear of retaliation. Effective reporting channels in ethical matters, well-investigated processes and substantial subsequent actions on ethical matters leads to a commitment by the organization in ethical responsibility. It is necessary to protect the people who raise their voices against the malpractice and ensure that the lapse of morals is avoided to enhance trust. The third element, Transparency and Explainability, involves how information and data are displayed in a manner that is comprehensible and transparent to all stakeholders, including all senior management of the company and the employees.

### *3.6 Transparency and Explainability*

Accountable AI systems rely on the prerequisite of transparency and explainability, which allow the stakeholders to comprehend, challenge, and oppose automated decisions which can influence their lives.

The interpretable machine learning can be regarded as a type of machine learning that can be understood and interpreted as a model.

#### *Interpretable Machine Learning*

**Interpretable Machine Learning** Interpretable machine learning may be referred to as machine learning whose model can be understood and interpreted. Model-intrinsic interpretability is associated with machine learning models the output of which can be explained by its own clear logic that can be understood [89-91]. This category includes linear models, decision trees, and rule-based systems since their mathematical form can be directly examined regarding how the features are affected in making predictions [6,92-94]. Simplified models with few features are associated with interpretability because the relationships are limited making them easy to understand by the users. Predictive accuracy has a trade-off who inherently interpretable models can only to a certain extent capture things in the data which are complicated to predict. But in most cases, small compromise in accuracy can be rewarded by the benefits of transparency. The decision to use interpretable or opaque models must take into account regulatory requirements and stakeholder needs and implications on the prediction made without them being explained. Generalized additive models are an intermediate between the nonlinear prediction models that are entirely flexible, and linear models, where predictions are given as a sum of univariate functions of individual features. This kind of structure is interpretable as it presents the marginal effect of each feature and nonlinear relationships are allowed. Model behavior can be made transparent using visualization of component functions and thus allowing domain expert validation.

#### *Methods of Post hoc Explanations*

Post-hoc types of explanations have an effect of producing explanations of predictions using models that are complex and opaque, without having to modify model design or training algorithms [95,96]. The importance of features procedures determine how each input feature contributes to predicting, or put differently, what variables had the most significant impact on outcomes with the help of a model [97-99]. Model-agnostic importance estimates Importance Weighted by Permutation Importance. LIME (Local Interpretable Model-agnostic Explanations) model surrogate models assess what will happen when a particular prediction is made by a complex model, and approximates the generally complex behavior of the model in the vicinity of the prediction(s) through interpretable models [100-103]. LIME produces instance-specific explanations by creating synthetic examples close to the example under explanation and the model being approximated by simple models of these instance elements that affected a given prediction. This method applies to any model, although it can give unreliable explanations when local approximations give a bad fit to the model behavior.

SHAP (SHapley Additive exPlanations) scores give feature attribution which is theoretically based, and based on the cooperative game theory, and the credit is distributed among features based on their contributions to the prediction [7,9,104-106]. SHAP values have good properties, such as accuracy, missingness, and consistency, so they are applicable to instance or global explanations. Nevertheless, the calculation of precise SHAP values may be resource consuming with complicated models. The

counterfactual explanations determine the smallest variation of the features in the input that would cause changes in the model prediction to respond to the questions on what other inputs would cause other results. These are explanations that are quite intuitive and practical indicating what tangible alterations that individuals can submit to get better predictions. To come up with valid counterfactuals, it is necessary to ensure that changes, which are proposed to change the situation, are realistic and practicable, and not the ones, which are mathematically possible but practically impossible. The neural network saliency maps and attention visualization methods point to the most influential regions of the input that are important in making a prediction. Saliency maps used in computer vision activities reveal which pixels are the most influential in decisions about the classification. In the case of natural language processing, attention weights would show the words or phrases that the model paid attention to. Such visualizations give intuitive information on model reasoning and might not fully represent all the factors that are of importance. In explanation quality and evaluation, it is important to explain why we should trust the results or declare them one-minute reviews.

#### Quality of Explanation and Evaluation

It is also important that we explain the reasons why we should believe the results or testify them to be one-minute reviews. Proper explanations should be loyal to model behavior but not to give probable yet wrong rationalization. The aspect of faithfulness can be tested by determining whether an explanatory model is able to predict model responses to input perturbations. Explanations that do not reflect the true model thinking in the system can give wrong information to the users concerning the system capability and restrictions. Explanations must be comprehensible to the intended audiences in terms of mental capacities as well as domain knowledge of those audiences. Documentation of technical experts might also consist of a combination of mathematical descriptions and statistical measures whereas documentation with end-users should be presented in a natural language with easy-to-understand visualizations. Customizing explanations to suit the requirements of the audience is increasing their usefulness.

The actionability of explanations is determined by their useful tips of attaining the desired results. The explanations made counterfactually which imply unattainable changes or a feature importance score in absence of context in which features change gives low utility. Actionable explanations relate insights in models with actual action that users can perform.

Contrastive theories answer the question why a given prediction took place instead of an appropriate alternative, which is a natural behavior of humans to seek explanations. Such descriptions can be more educative than merely explaining what supported a prediction because they are more discriminating, and show differences between results. The consistency of the explanation on similar instances creates user confidence and allows one to obtain general patterns in the model behavior. The lack of the similar cases being explained consistently might be a symptom of the inconsistency of the explanations methods or the true complexity of the decision-making process of the model that needs to be investigated.

#### *3.7 Privacy and Data Protection*

Privacy risk overlaps with the fairness and bias issue whereby, privacy-saving strategies might provoke new equity problems and biasness violations whereby the vulnerable populations are usually predominantly impacted.

##### Privacy Menace of Machine Learning

The identification of training data privacy risks is that the sensitive data in training datasets can be deduced or retrieved by trained models [6,107-108]. Membership inference attacks identify whether individual data of particular users has been used in the training sets which can reveal sensitive information about a given characteristic or behaviour of an individual. In model inversion attacks, the attacker is trying to recreate the training data, based on the model parameters, and this is specifically dangerous when it training data contains personal data. Risks of re-identification can occur when what is supposed to be anonymized data can be associated with other sources of information to determine an



individual. Patterns that can be used in re-identification might be acquired by machine learning models trained on anonymized data, especially with the help of auxiliary information. Attributes of the business such as demographics that enable fairness analysis can also enhance re-identification risks, which set conflicting goals between fairness and privacy. The concept of differential privacy offers a mathematical description of privacy data privacy risk quantification and data privacy risk bounding by introducing noise with a precisely measured amount to data or model outputs. Bounding certain data points on the effect of individual data points guaranteed the extent to which pure inference on a single person may be achieved. Nevertheless, the protection of differential privacy has utility costs, their presence imposes noise onto models, and they might disproportionately impact various groups of people.

Table 1: Key Aspects of AI Bias, Fairness Metrics, and Mitigation Techniques

Sr. No	Aspect	Source/Type	Detection Method	Mitigation Technique	Challenge	Opportunity	Future Direction
1	Historical Bias	Training data reflecting past discrimination	Demographic analysis of datasets, outcome distribution comparison	Data reweighting, balanced sampling, synthetic data generation	Accurate historical data perpetuates injustice	Creating fair baselines through deliberate data curation	Developing methods to learn from biased history without encoding bias
2	Representation Bias	Underrepresentation of minority groups	Group-wise performance evaluation, error rate analysis	Oversampling, targeted data collection, data augmentation	Limited minority group data availability	Collaborative data sharing across organizations	Privacy-preserving techniques for expanding representation
3	Measurement Bias	Inappropriate feature selection or proxies	Feature correlation analysis, proxy identification	Feature transformation, removing proxy variables	Difficulty identifying all proxy pathways	Causal modeling to distinguish legitimate from proxy features	Automated proxy detection algorithms
4	Aggregation Bias	Single model failing to represent diverse groups	Subgroup performance analysis, clustering analysis	Separate models per group, mixture of experts	Determining appropriate subgroup definitions	Personalized models adapting to individual needs	Meta-learning approaches for efficient multi-group modeling
5	Evaluation Bias	Biased benchmark datasets	Multi-dataset evaluation, demographic composition analysis	Creating diverse benchmarks, participatory evaluation design	High cost of comprehensive benchmark creation	Crowdsourced diverse dataset development	Automated generation of diverse test cases
6	Deployment Bias	Context mismatch between development and deployment	Continuous monitoring, drift detection	Domain adaptation, model updating, human oversight	Unpredictable operational environments	Robust models generalizing across contexts	Transfer learning with fairness constraints
7	Demographic Parity	Unequal outcome rates across groups	Statistical parity testing	Threshold adjustment, constrained optimization	May sacrifice accuracy or merit-based selection	Ensures proportional representation	Context-aware relaxations of strict parity
8	Equalized Odds	Unequal error rates across groups	TPR and FPR comparison	Post-processing threshold optimization, in-processing constraints	Cannot simultaneously satisfy with other criteria	Equal quality of service across groups	Efficient algorithms for constrained training
9	Predictive Parity	Different PPV across groups	Precision analysis per group	Calibration techniques, score transformation	Conflicts with other fairness definitions	Predictions equally meaningful across groups	Multi-objective optimization frameworks
10	Individual Fairness	Dissimilar treatment of similar individuals	Distance-based similarity analysis	Fairness regularization, Lipschitz constraints	Defining task-appropriate similarity metrics	Respects individual merit	Learning similarity metrics from

11	Counterfactual Fairness	Decisions depend on protected attributes	Causal model analysis	Causal graph-based training, removing discriminatory paths	Requires accurate causal models	Principled approach to discrimination	stakeholder input Causal discovery from observational data
12	Intersectional Bias	Compounded disadvantage across multiple attributes	Multidimensional subgroup analysis	Intersectional fairness constraints, targeted interventions	High-dimensional demographic spaces, limited data	Better captures real-world discrimination	Developing tractable intersectional metrics
13	Allocation Harm	Unfair resource distribution	Resource distribution analysis	Fair division algorithms, randomization	Defining fair allocation with scarce resources	Explicit equity considerations	Participatory approaches to allocation criteria
14	Quality-of-Service Harm	Differential performance quality	Disaggregated accuracy metrics	Balanced training, importance weighting	Accuracy-fairness trade-offs	Minimum performance standards for all groups	Fairness without significant accuracy loss
15	Representational Harm	Stereotyping and demeaning depictions	Content analysis, user feedback	Filtering harmful outputs, diverse training data	Subjective harm definitions	Inclusive and respectful AI systems	Culturally-informed content evaluation
16	Statistical Discrimination	Group stereotypes influencing individual predictions	Feature importance analysis	Individual assessment, reducing reliance on group statistics	Trade-off between accuracy and fairness	Merit-based individual evaluation	Hybrid approaches balancing individual and group considerations
17	Feedback Loops	AI decisions affecting future data	Longitudinal monitoring, simulation modeling	Randomization, exploration strategies	Long time horizons for detection	Breaking cycles of disadvantage	Predictive modeling of feedback dynamics
18	Label Bias	Inconsistent or discriminatory annotations	Inter-annotator agreement analysis	Multiple annotators, consensus methods	Subjectivity in ground truth	Diverse annotator pools improve quality	AI-assisted annotation with bias detection
19	Sample Selection Bias	Non-representative data collection	Comparing sample to population demographics	Stratified sampling, reweighting	Unknown population distributions	Deliberate representative sampling	Active learning for diverse data collection
20	Temporal Bias	Outdated data not reflecting current reality	Monitoring concept drift over time	Regular retraining, online learning	Costs of continuous updates	Models adapting to changing contexts	Efficient incremental learning approaches
21	Geographic Bias	Regional disparities in data and performance	Geographic performance analysis	Location-specific models, geographic reweighting	Cultural and infrastructural variation	Locally-appropriate AI systems	Federated learning across regions
22	Linguistic Bias	Disadvantaging non-standard language varieties	Performance analysis by language/dialect	Multilingual training, dialect-aware models	Limited data for low-resource languages	Inclusive language technology	Multilingual transfer learning
23	Accessibility Bias	Failing to accommodate disabilities	Accessibility testing, assistive technology compatibility	Universal design principles, accessibility features	Diverse disability needs	Technology empowering all users	AI-powered accessibility assistance
24	Age Bias	Differential treatment by age group	Age-stratified evaluation	Age-inclusive training data, age-aware features	Legitimate age-related differences versus discrimination	Appropriate age considerations	Distinguishing correlation from causation
25	Socioeconomic Bias	Disadvantaging lower-income populations	Analysis by socioeconomic indicators	Removing socioeconomic	Socioeconomic factors intertwined	Reducing inequality	Alternative features not

	proxies, needs- based design	with legitimate features	through technology	correlated with wealth
--	---------------------------------	-----------------------------	-----------------------	---------------------------

### Privacy-speaking Machine Learning

Federated learning provides a way to train models on distributed data without sensitive data being centralized. Each participant also trains local models using personal data and does not provide the raw data, which is an advantage of privacy. The process of federated learning enables the cooperation across organizations without violating information sovereignty and privacy policies. Nevertheless, updates of the models can still provide information on local data, and extra protection of privacy is needed. Secure multi-party computation enables several parties to compute functions collectively using their joint data without the input of each party being made public. Based on these cryptographic protocols, privacy-preserving analytics and machine learning are possible with computations being made on encrypted data. Although secure multi-party computation is a privacy-assuring construct, it has high computational costs that could restrict its practical use. Homomorphic encryption allows calculations to be performed on encrypted data without the data being decrypted and model inferences of encrypted inputs to encrypted outputs. This strategy ensures the privacy of data along the computation line. Nonetheless, homomorphic encryption is at presently a highly costly and restrictive procedure in terms of computational resources and the kind of operations that can be effectively done. Synthetic data generation produces artificial data that still has the statistical characteristics of real data but eliminates information on a case-by-case basis. Well planned artificial information may aid in the development of models and analysis and safeguard privacy. Nonetheless, it is difficult to guarantee that synthetic data is sufficient in the representation of significant patterns but gives a sense of true privacy, especially in high-dimensional or structured data.

### Privacy-Fairness Trade-offs

Techniques that protect privacy such as differential privacy can disproportionately hurt the model performance of minority groups which already have low representation in training data. Addition of noise is a more serious problem in smaller groups as compared to larger ones, and may worsen underlying elements of unfairness. Those asymmetric privacy-utility trade-offs should be taken into special consideration when developing privacy-sensitive systems. It is possible that privacy laws and regulations can affect different demographic groups differently based on their data practices, digital literacy and susceptibility to privacy breaches. Privacy frameworks that are based on consent might not work well with those who have little knowledge on the process of data practices or have low bargaining power to stipulate privacy terms. Privacy guarantees that restrict gathering may decrease the depiction of already minority populaces on training information. To balance privacy and fairness, complex strategies need to be considered and in this case these should have a special consideration of the impact of privacy control in the case of various populations and possibly offer divergent levels of privacy to various populations. Other scholars have suggested group privacy definitions that are more protective to vulnerable groups, but such strategies have the difficulty of incurring complex normative dilemmas of justified different treatments.

### 3.8 Accountability and Governance

The company has several accountability and governance policies to be observed to ensure that the protectorate remains effective. An effective governance system and accountability of AI systems must be set to relax responsible deployment and development of AI systems and provide remedies in case of harm.

#### Accountability Frameworks

Responsibility in AI systems entails determining the individual in charge of the system behavior, the system has a way of monitoring and redressing, as well as ensuring that there are incentives that encourage responsible behavior. The attribution of responsibility in multi-actor complex AI supply chains that use data providers in the beginning to model developers at the other end, through deployers,

is complicated. Definite accountability structures are necessary to define tasks throughout the AI lifecycle, and how the various participants can liaise with each other to deliver ethical results.

Human-in-the-loop systems ensure there is a significant degree of human control and management of the consequences of decisions by holding an automated recommendation subject to significant human scrutiny, the approval, or the supervision of the automated decision-making process. The efficiency of the human oversight mechanisms is heavily influenced by the level of their design since a simple audit of the process can turn to rubberstamping the overseers, whereas excessive scrutiny can be dodged or, conversely, can lead to emergence of additional biases. Good human control is achieved based on proper training, reasonable workload and some instructions regarding circumstances when automated recommendation are to be ignored. Audit trails and documentation practice generate records which can be retrospectively used to analyze system behavior and make decisions. Recording of inputs, outputs, and version of the model and the surroundings can be used to investigate accidents, determine trends, and hold oneself responsible. Nevertheless, extensive logging enhances privacy issues and creates huge amounts of data which need complex management.

The conditions of impact assessment necessitate the consideration of the possible outcomes prior to the implementation of the AI systems in the sensitive areas. These evaluations would be based on risks to people and the community, the affected populations, evaluation of mitigation action, as well as creating an elucidation strategy. Reassessment on a regular basis adds up to ensuring that new risks are detected and mitigated with changes in the systems and occurrence of new environments.

#### Organizational Governance

The AI governance regimes put in place institutional frameworks of managing AI creation and implementation at the operational levels of organizations. Ethics review boards are developed to analyze the proposed AI projects, evaluate the potential risks, and give advice on resolving ethical issues. These boards are supposed to have varied skills in terms of technical, ethical, legal, and domain points of view, and also the stakeholder groups that are affected must be represented in them. Responsible AI programs articulate the policies, procedures, and resources in the organization that would promote ethical AI development. These programs allow establishing standards of fairness, transparency and accountability, training and tools, and processes of responding to ethical issues. Effective programs must be supported with executive leadership, sufficient allocation of resources and must be included in the normal development processes.

The AI risk management frameworks expand the classic risk management methods to the AI-related issues such as algorithmic bias, opaqueness, emergent behaviors, and dual use potential. The risk assessment helps in identifying the possible harm, the likelihood and the severity of the harms and helps in the prioritization of the mitigation effort. Risk registers follow up on identified risks and mitigation procedures as well as accountable individuals. The stakeholder involvement approaches guarantee that different views are being used in the governance of AI. The external voices in organizational decision-making are through advisory boards, community consultations as well as user research. As meaning engagement, it is not enough to consult stakeholders on the basis of tokenism or be ready to show the impact of stakeholder contributions on decision-making and establish relationships with communities that are to be affected.

#### External Oversight and Regulation

Third-party auditing is the evaluation of the fairness of AI systems, their safety and adherence to the regulations or moral standards conducted by a third party. External auditors help to create objectivity and specialized knowledge, but they need access to systems, data, and documentation which the organizations might be unwilling to provide. It is not an easy task to establish AI auditing standards and certify qualified auditors.

A variety of regulatory measures to govern AI is used in different jurisdictions, and it reflects varying legal traditions and values as well as risk disposition. Certain jurisdictions are interested in sector-focused laws on high-risk applications whereas others are interested in horizontal laws that are

implemented across AI systems. A good regulation should be able to strike a compromise between innovation and protection, offer strong guidelines to the developers, and keep in pace with a fast changing technology.

The activities of certification and standardization come up with technical standards, testing standards, and certification programs of AI systems. Standards offer an understanding over shared words, provide minimum requirements, and bring interoperability. But, inadequate standardization would prevent evaluation and comparison of systems, whereas premature standardization would establish suboptimal methods of doing things.

Legal liability frameworks define the situations and the persons who will be legally liable to the harms caused by the AI systems. The current regimes on liability might not suit very well the nature of AI systems that include their autonomous nature, distributed nature, and emergent nature. The issues concerning the strict liability and negligence, allocation of responsibility through the supply chains, and the proper ways of compensating algorithmic harms are disputed.

### *3.9 Domain-Relevant Applications and Problems.*

#### **Criminal Justice and Policing**

AI solutions used in criminal justice as risk assessment tools to decide on bails, sentencing, and parole offer fundamental concerns of fairness due to the stakes involved and reported disparities in criminal justice results based on racial and socioeconomic lines. Predictive policing systems that can predict crime locations or designate specific people that might receive special scrutiny have a danger to perpetuate discriminatory crime patterns and stability, as more people are monitored, leading to more information about over-policed populations. The bias in criminal justice data of the past is incredibly problematic because the arrests, convictions, and recidivism statistics are clear-cut indicators of both actual criminal activities and bias police and prosecution policies. When predictive models are trained based on this data, injustices are likely to be propagated. This fact makes the application of the arrest records as the outcome, particularly troubling since arrests are an expression of the police actions rather than of individual behavior.

Criminal justice AI must be transparent to promote the right to due process and be able to contest automated criminal justice decisions significantly to keep freedom. There are however a large number of commercial risk assessment tools that are proprietary and, therefore, could not be externally scrutinized. The interpretability of contemporary machine learning model poses further challenges to elucidation of individual risk forms that can be useful in legal procedures. The metric of fairness in criminal justice is a dispute that is symptomatic of larger disputes between the various metrics and antagonistic normative principles. Calibration: When there are varying levels of false positives, then using the same risk scores across racial groups can mean the same thing though it must not be assumed it has the same level of false positivity. Equalized false positive rates guarantee equal chances of mistakenly imprisoning low- threat suggested ones yet might cause differing precision in groups. It has practical trade-offs to freedom of individuals and the security of the public.

#### **Healthcare and Medicine**

Healthcare systems that use AI have an impact on diagnosis, prescription, resource distribution, and participation in clinical trials and malfunction of fairness may lead to severe health risks or even death. Medical research disproportionately adds women, racial minorities, and other people who are historically underrepresented, which results in a lack of balance in data and machine learning systems can exacerbate the situation. The diagnostic systems that are trained on the data of one demographic group can be very poor on the other demographic groups and fail to detect diseases or falsely positive. The healthcare resource allocation algorithmic systems have to maneuver through some hard ethical dilemmas regarding how to provide scarce resources equally to patients with varying profiles and requirements. Models that forecast healthcare expenditure or anticipated future healthcare usage have been racially biased because they have used the spending as a proxy of the health requirements, which

do not account for the disparate provisions of care that lead to decreased spending by black patients in spite of their increased health requirements. The healthcare sector is a field that encompasses privacy issues especially where any medical information is highly sensitive and is regulated by stringent laws. The machine learning approaches implemented to preserve privacy and assist in collaborative research across institutions should not only have intense protections but still remain useful in medical settings. Consent models should provide tradeoffs between personal authority over medical information, and social good of medical research. Clinical decision support systems should be easy to understand by clinicians and be able to question and to override the recommendation in case they suit. Black-box systems rather pure and offering no explanation just recommendations, are not likely to be taken seriously and integrated safely in clinical practice. Nevertheless, too complicated explanations can be ineffective in the situation of time pressure in clinical settings.

#### Employment and Hiring

The AI systems become more and more involved in the process of employment access mediation via automated resume screening, video interview analysis, and candidate ranking. These systems threat to perpetuate past instances of discrimination in employment and form new sources of bias. Computerized screening of resumes that disapproves employment breaks can be discriminatory to women who went on parental leaves and people with disabilities who were on medical leaves. The systems of video interview analysis, which determine personality characteristics, communication skills, or cultural fits through facial expression and speech patterns, beg the question of validity, fairness, and privacy issues. These systems can discriminate against people with disabilities that are speech or facial expression-related, can be discriminating towards accents or other forms of communication typical of a particular demographic, and have poor evidence coverage to validate them. The fact that the AI used in employment is rather obscure leaves the job applicants confused by the reasons why they were denied the job or how they can increase the chances of getting it. Compared to consumer lending where policies would mean that contracts on adverse actions must be sent, employment decisions are characterized by fewer transparency obligations. This information imbalance benefits the employers and leaves the applicants with no avenue. Ethical issues surrounding AIs in hiring can also be concerned with the proper data gathering and drawing of conclusions. Not trying to predict such intangible traits as disability or pregnancy based on other observable traits, it is clear that this practice goes against non-discrimination principles, but possibly such predictions can be made unconsciously in complicated models. There should be a fine line between acceptable assessment and intrusive surveillance.

#### Financial Services

The AI systems used in credit scoring, loan approval, and insurance underwriting have far-reaching implications on the economic opportunity of access to much-needed financial services. Past discrimination of redlining and discriminatory lending practices generate biases of data that can be perpetuated or even increased by machine learning systems. Using alternative sources of data such as activity on social media or internet use begs the question of whether this type of data can be used to make financial decisions or not. Explainability of consumer finance is not only legally mandated in most jurisdictions but of practical value in the facet of allowing consumers to become better credit-worthy. Adverse action notices do have to give reasons as to why credit applications were not approved but it is difficult to make complex machine learning model decisions adapt and be executed. The scores of generic features importances might not be of value to specific applicants. The fairness of insurance pricing should be balanced between actuarial fairness that imposes costs on individuals according to the expected costs and the fairness principles of society, which rejects the idea of discrimination. The predictive variables which may have been correlated with insurable characteristics raise a challenging trade-offs. The ban on the licensed characteristics fails to eliminate discrimination provided that proxy variables are used to discriminate against it indirectly. Due to the temporal dynamics of financial AI, certain issues of fairness exist, since credit score systems determine the financial prospects of people and their future credit progress. The first bads can restrict access to credit and thus people cannot record good credit histories which forms path dependencies and this can further lead to disadvantage. Cycle interruption can be a complicated process that could demand specific intervention in individual models in addition to biases reduction.



## Education

Algorithms in education, such as automated scoring systems and student performance predictors, and personalization learning systems, have long-term effects on educational trajectories. Educational AI bias can solidify available educational inequities in education in terms of socioeconomic status, race, language background, and disability status. The automated essay grading systems can prove disadvantageous to non-native speakers, students of other cultural backgrounds with alternative writing rules, or students having disabilities that will influence their writing. Such systems tend to put more emphasis on features such as vocabulary and structure of sentence which can be inadequate proxy indicators to actual understanding or reasoning capacity. These issues are multiplied by automated scoring of high stakes examinations such as college admissions.

Performance prediction systems that can predict the students at risk of failing the academic system or dropping out of school would allow prompt intervention, but have a danger of making self prophecies. False negatives can deprive flailing students of the resources they require and negative predictions can involve lowered expectations or resources of those who are predicted to perform poorly. These predictions have feedback effects that must be taken into account on the performance of the students. Individualized learning systems which adjust content and speed to the needs of individual students have a potential to improve the needs of different students better but with the risk to sort students into restrictive learning streams. When systems offer inferior content to struggling students in the early ages they can deny such students access to higher level content keeping the achievement gaps increasing. The conditions of relevant content adaptation include the normative and technical judgments.

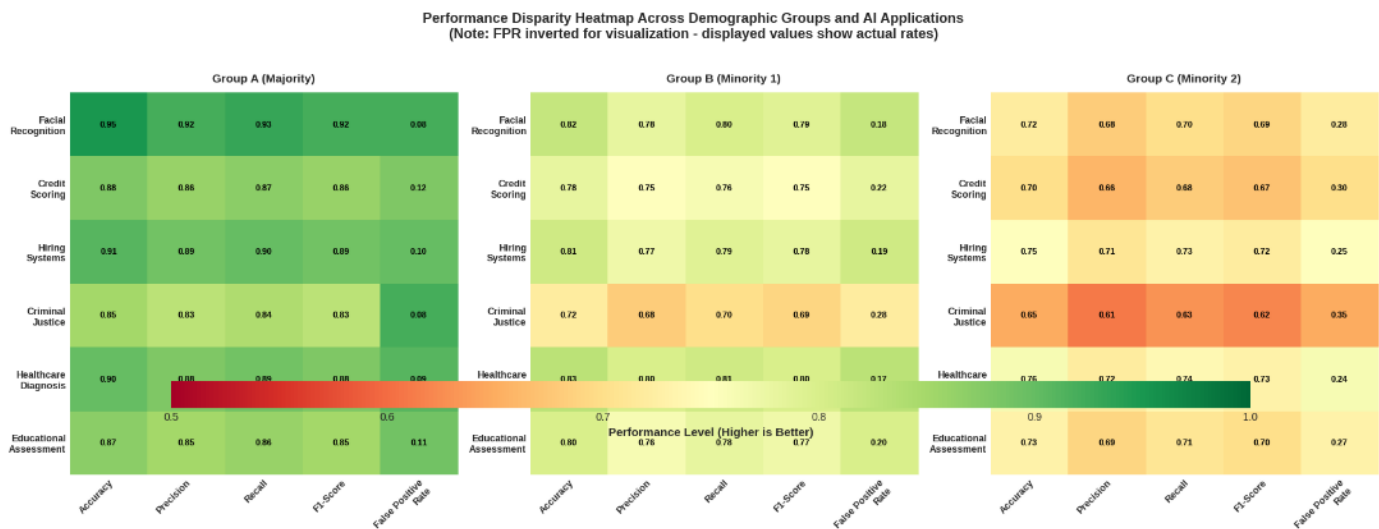


Fig 3: Performance Disparity Heatmap Across Demographics and Applications

Fig. 3 displays performance metrics (Accuracy, Precision, Recall, F1-Score, and False Positive Rate) across different demographic groups and AI application domains. The color intensity represents the metric value - darker red indicates lower performance/higher FPR, while darker green shows better performance. The visualization reveals significant disparities: for instance, facial recognition shows accuracy of 0.95 for Group A but only 0.72 for Group C, while criminal justice systems show dangerously high false positive rates (0.28) for Group B compared to Group A (0.08). This pattern highlights systemic bias requiring immediate attention.

### 3.10 Future and Emerging Technologies and Future Challenges

#### Generative AI and Large Language Models

MLMs that were trained on very large internet text corpus are biased by what they have been trained on, such as stereotypes, toxic language, and false information. The models are able to formulate realistic yet fake data, recreate the copyrighted data and harmful data even with safety in place. The magnitude

and generality of these models imply that they can have a bias influence on large scopes of applications and users. Language models have bias that can be seen as the proximities between demographics and certain traits, suggestive completions of prompts that pertain to particular demographics, and differences in quality of services of different languages or dialects. These prejudices may be implicit and contextual hence not easy to identify in a holistic manner. Any effort to train data that attempts to lessen bias by increasing training data or fine-tuning can present scaling problems, and can also create additional bias or decrease model capacity. Generation AI image or video and audio generators create images, videos, or audio that are concerning in terms of creating damaging stereotypic images, generating non-consent intimate photos, facilitating identity fraud, and propagating false information. The persuasiveness and magnitude of synthetic media resists established strategies in regards to the integrity and confidence of information. Technical responses include watermarking and provenance tracking of synthetic content, which rely on capturing a market share and is difficult to avoid. Large language models have a dual-use character, which makes it difficult to establish any governance: the same tools that allow the utilization of the language models include language translation service, content creation, and code completion can be applied to create spam, propaganda, or malicious code. Keeping access limited so as to avoid misuse restricts good use, whereas free access allows ill use. This conflict has created controversy on the issue of responsible releases and controlled access.

#### Autonomous Systems and Robotics

Autonomous cars have to solve ethical problems concerning the importance of prioritizing various consequences in case accidents cannot be avoided, which brings concerns regarding the safety of whom to prioritize and how to implement ethical standards into systems that make decisions. Autonomous vehicle programming can be interpreted as value judgement concerning the acceptable level of risk, as well as demonstrating a trade-off between the safety of the occupants and that of pedestrians. Disparate demographic safety outcomes might occur because of bias in autonomous vehicle perception systems like the higher failure rates to identify darker-skinned pedestrians. People, in the event that autonomous cars are not as safe as they claim on some groups of the population, their implementation can potentially worsen transportation disparity instead of improving it. Strict testing in different conditions and population are prerequisites that might be difficult due to a lack of variety in testing data sets. Dignity, autonomy and the relationships of care are issues that emerge when dealing with robot systems in service jobs like the care robots assisting the elderly or the disabled. Implementation of care robots should take into consideration the question of human interaction being rightfully augmented against improperly replicated, individuals having meaningful choice concerning robotic care, and the reproduction of problematic stereotypes or expectations by robot.

Robotics in military and security services become particularly acute as soon as autonomous weapons systems, surveillance robots, and decisions on the use of lethal forces are possible. Delegating life-and-death algorithms to autonomous systems provokes the core concepts of human agency and responsibility. The effort globally to put a system of governance of autonomous weapons faces challenging issues of technological capability, check, and accountability.

#### Component models and Transfer Learning

Broadly trained foundation models and fine-tuned or transfer learning based on broad data focus power into the companies with sufficient resources to train large models and establish dependencies on downstream applications. Biases within foundation models are replicated to all downstream applications increasing their effects. The current existence of foundation model development within few organizations prompts the question of the values and priorities, of whom these powerful systems take, to be developed.

Their generality and the fact that it is impossible to test all types of applications and prompts of foundation models complicate the assessment of their fairness and safety. This is because red-teaming and adversarial testing can determine some undesirable usage patterns, although the space of potential usage is too large to evaluate comprehensively. The post-implementation of harmful capabilities or biases is what requires fast response systems. Such a practice of providing to foundation models a fine-tuning or prompt to particular applications creates uncertainty on the question of liability concerning

downstream harms. It is controversial when the use of the general-purpose models leads to destructive outputs that are caused by either wrong or malicious application of the models. The policies of terms of service and use are trying to divide the responsibility, yet their implementation under the law and their sufficiency are not entirely clear. Large models with emergent capabilities, where capacities for greater scale are introduced that small versions do not have, make safety and fairness testing difficult. Only smaller models can be tested and it may not show problems that arise in larger ones and the large models are constituted by their computational costs which makes it prohibitive to test safety through iteration. This scaling dynamic introduces a problem in terms of predicting and avoiding harms before deployment.

#### Multimodal and Dual-domain AI

The interplay of modalities in AI systems generating and processing more than one aspect of data (e.g., text to image generation or visual question answering) creates a novel set of issues in fairness at the confluence of disparate data types. Prejudice can be in the form of associations between images and words, like the creation of images of typical stereotypical appearance to the textual cues regarding specific work or activity. Cross-domain transfer AI models i.e. when systems that have been trained in one domain are applied to another domain can create or enhance biases when the domains are different or when the transfer process is not good enough to accommodate the differences in domains. A model developed based on the data of a certain culture will not work or turn out to be biased when implemented in another culture with other norms, values, or practices. To pose multimodal models of fairness, the datasets must be diverse with multiple dimensions at the same time, i.e., there should be images of people of different demographics occupying various environments doing different things. The development of these datasets is energy consuming and has privacy issues especially when dealing with images of recognizable entities. The growing complexity of multimodal synthesis allows one to produce very realistic fake content, representing people in falsified situations, which is of great concern in the context of consent, damage to reputation, and the evidence validity of the visual image. The synthetic multimodal content is hard to detect technically and can be potentially undergenerative.

### *3.11 Regulatory Environment and Policy Interventions.*

#### Current Regulations and Law Systems

In most jurisdictions, the decision made by anti-discrimination laws is against an action on the basis of a protected characteristic that may be race, gender, age, disability, and religion. These rules are applicable to AI systems employed in the employment, housing, credit, and other areas, according to which the civil rights legislation applies. Nevertheless, when implementing current anti-discrimination paradigms to algorithmic systems, there are interpretive problems of indirect discrimination, disparate impact, as well as the applicability of statistical evidence. Such rules as the European Union General Data Protection Regulation provide transparency of data processing, purpose limitation, data minimization, and the rights of individuals such as automated decisions explanation. The requirements influence the development of AI in the form of limiting the data collection and processing, requiring documentation, and necessitating being able to explain. But the real application of such provisions as the right to the explanation is disputed.

Regulations in sectors such as the financial industry, healthcare and telecommunication sector contain fairness and non-discrimination standards to which AI systems can be applied to the relevant sector. The credit denials that financial regulations may impose need an adverse action notice, action against healthcare regulations safeguard patient information and the informed consent, and communications regulations are concerned with discriminatory practices. These diverse needs are a complication to AI systems that work in different industries. The liability of products and consumer protection regulations provide responsibility in terms of injuries that defective product causes or deceitful practices. The application of these frameworks to the AI systems leads to doubt of what is considered a defect in an AI system, how to decide the causation when an algorithmic decision can contribute to harms with human decisions and whether current liability regimes will be sufficient to encourage the production of safe AI systems.

### Emerging AI-Specific Regulations

The AI Act suggested by the European Union is the set of regulations that fully classify AI systems based on risk and has requirements commensurate to the risk. Certain systems require heightened risk evaluation necessities in the areas of work, education, legal frameworks, and vital infrastructure that can comprise cargo of selection, information quality, openness, human control, and auditing. Some of the prohibited applications are social scoring and real-time biometric identification in the public places unless there are any special law enforcement reasons. The innovation/protection balance that the risk-based approach of the AI Act should provide is to issue more regulation on those applications with the highest risk of harm and less regulation on those with less risk of harm. Nevertheless, there are still doubts regarding the feasibility of the risk categorization implementation, the sufficiency of transparency proposals, and enforcement tools in the regulator states of the EU. Alternative jurisdictions are building their own AI governance systems based on askew regulatory ideologies and priorities. Others focus on self-regulation in the industry and setting of rules based on principle and others prefer compulsory conditions and prior approval procedures. The variety of regulatory methods poses difficulties to the development of international AI and can result in regulatory fragmentation or even race to the bottom.

The accountability mechanisms applied to AI systems would consist of impact assessment, auditing, transparency reporting and other mechanisms of algorithmic accountability as suggested by the algorithmic accountability acts proposed in different jurisdictions. Such proposals are usually aimed at high-stakes applications or large organizations, trying to scale oversight and non-imposing too large a burden on small organizations or useful applications. Determining right levels to which requirements must be implemented and setting audit standards are still problems.

### International Cooperation and Standards

Development of AI systems technologic standards in areas such as bias testing, transparency documentation, risk management, etc. is being developed by international standards-setting organizations. Standards offer universal patterns through which comparisons can be made across systems across jurisdictions as well as compatibility with several regulatory frameworks by utilizing harmonized standards. Nonetheless, the processes of standard-setting can be ineffective at both the economic global and global level and too slow to stay abreast of technology. Such multilateral efforts as the OECD AI Principles, UNESCO Recommendation on AI Ethics, and other multi-stakeholder forums are attempting to create common ethical standards and align national governance practices. These activities aim at avoiding regulatory fragmentation besides setting minimum standards of responsible AI. Non-binding recommendations, however, may not have enforcing policies and it is hard to arrive at an agreement between different countries whose values and interests differ. The global effort towards AI regulation faces the challenges of geopolitical stress, inconsistency between regulatory ideologies, and the rivalry in AI innovation by nations that want to gain an edge. Giving the right balance between cooperating on safety, ethical and competing technologies to lead in technology necessitates diplomacy and trust that it might be hard to maintain. International controls and checks on the AI development practices are subject to technical and political hurdles.

The data flows across the borders necessary to train AI systems intersect with the needs of the data localization and the issues of sovereignty which poses a possible conflict between the interests of development of AI and the laws of data protection. The regulatory interoperability instruments, e.g., the adequacy decisions or mutual recognition agreements, can ensure the positive data-sharing exercises but still considering the local needs, which is difficult to negotiate.

### *3.12 Best Practices and Implementation Strategies*

To convert ethical principles and regulation demands into practice, definite plan of action that should be followed in the lifecycle of AI development is needed.

### Responsible AI Development Lifecycle

The needs analysis at the project start and stakeholders consultation must identify possible issues of fairness, beneficiaries of the project, the regulations applicable, and measures of success involving both ethical aspects and technical goals. Early consultation with various stakeholders such as potential users and communities that are impacted, may help unearth issues and priorities that technical teams may not take note of. Curation and development of datasets are advised to pay attention to the aspects of representation, quality, and bias in training data by using a planned sampling methodology, repeating annotation sessions with different annotators, evaluating label quality and recording the data properties and constraints. The documentation of databases should encompass details of collection procedures, demographic distribution, biases which are familiar, and the right uses. Improved fairness objectives together with predictive precision ought to be included in the development of models in ways that provide fairness-sensitive learning algorithms, regularization to address fairness constraints, and assessment of multiple fairness measures. The model used ought to be selected based on fairness-accuracy trade-offs and with the selection it ought to be based on the needs and input of the stakeholders. Model behavior in a variety of scenarios, for different demographic groups, and edge cases should be tested and validated using disaggregated testing as well as fairness metrics calculation, adversarial testing and stress testing. Robustness to distribution shift testing should also be done, and an assessment to realistic operation data should be done instead of test to curated benchmarks only.

Preparation of deployment must encompass the recording of model coverage and constraints, system development of monitoring infrastructure, human oversight process and the development of incident response schemes. Applications to the user must include the relevant transparency concerning AI participation and decision challenge. Model performance and fairness measures in the long term should be monitored in the post-deployment, identify instances of distribution drift or bias, interpret user feedback and complaints, and review the real-world effects. Frequent review must identify that deployed systems remain fair, and that new circumstances may require that systems be altered. The capabilities and resources of the organization will be described. To develop organizational capabilities of responsible AI, investments in people, processes, and tools are needed. Both orientation and avoidance of ethical issues can be realized through technical team training or via the commitment of experts with ethical skills such as AI ethics. Legal knowledge will guarantee adherence to the laws and determination of the liability risks. Domain knowledge is a basis of the technical growth in the contextual knowledge of application domains.

Inclusiveness by a team that involves people with varied demographics, discipline, and views is more likely to be in a better position to spot any form of biasness and address various stakeholder needs. There is need to have diversity not only in terms of demographic aspects but also in terms of professional overages, including those who are technical and those who are social scientists, domain specialists, ethicists, and affected members of the community. Responsible AIs Tools and infrastructure Responsible AI tools and infrastructure include fairness testing libraries, bias detection systems, explanation generation systems and documentation frameworks, and monitoring systems. The advantages of investing in these resources consist in decreasing the barriers one will face on executing ethical practices and making a consistent application in projects.

Ethical processes and procedures that entrench ethics in the development work processes are in a way a guarantee that ethical issues are not handled on ad-hoc basis. The checks of fairness can be present during the code review, the release can be approved with the need of the ethics review, and the procedures of the incident response can ensure that the discovered harms are properly addressed.

#### Stakeholder Engagement and Participation

Significant stakeholder engagement is what lies beyond the consultation aspects of token participation and entails the inclusion of varied views in decision-making. Discussion in the early days of identifying the problem, requirements collection guarantees that the priorities of the stakeholders influence the course of the projects. Constant participation during growth allows projecting of the development through feedback. The participatory design approaches engage stakeholders in the process of solution development in co-design workshops, user research, and community-based participatory research. These methods acknowledge communities with affected ones as people experienced and knowing their

needs are the experts who are needed to find the right solutions. To resolve the issue of power imbalances in stakeholder engagement, conscious effort can be made to give voice to the marginalized groups, remunerate community members who utilize their time and expertise, and show that feedback-related responses are heeded. Development of open participation channels that should not be technical can help increase involvement. The existence of divergent perspectives of the stakeholders necessitates open procedures to the synthesis of input, making of trade-offs, and recording of justifications of the decisions. The preferences of all the stakeholders may not be met at the same time but the reasonability of why input was taken into account will create trust even in cases where certain demands are not fulfilled.

#### Education and Training

Responsible AI education should be taught across computer science and engineering spectrums as opposed to being focused in a time-out course about ethics. Fairness metrics in machine learning classes as well as privacy in security courses and human-computer interaction in systems design are to be covered in core technical courses. Learning based on case studies prepares the abstract ethical principles with real-life situations that practitioners might find themselves and create judgment and decision-making skills. It is necessary to study failures of AI in the real world, talk about challenging trade-offs, and learn how to think ethically to prepare practitioners to experience complex situations. Providing interdisciplinary education that combines the humanistic and the technical viewpoint helps to avoid antagonism between the domains of work. Having computer scientists study with ethicists, social scientists, and domain experts through joint programs, courses taught collaboratively, and group projects allow students in these disciplines to be exposed to these experts and imparts technical literacy to them. Continuing education and professional development holds important that the practicing AI developers are kept updated on new ethical concerns, new practices, and new regulatory concerns that emerge. There are accessible mechanisms of continued learning provided through industry conferences, workshops, certifications and online courses.

#### *3.13 Future Research Directions*

##### Technical Research Needs

The problem of coming up with fairness methods in complex AI systems such as reinforcement learning, generative models and multi-agent systems is still a topic of discussion. Most of the fairness studies have been conducted on the paradigm of supervised classification, whereas other learning paradigms have different problems. It is necessary to take into account the long term and sequential decision making when implementing reinforcement learning fairness. Generative model fairness has two aspects on the presentation of training data and the diversity in generated output.

The fact that intersectional fairness considers a variety of demographic dimensions at once as opposed to evaluating bias on a limited set of axes is also a technical and theoretical challenge. People who pertain to more than one marginalized identity can experience compounded disadvantages, which is not effectively reflected by a consideration of each one of the characteristics. To create fairness measures and lessening strategies to address intersectional fairness, there is a need to handle high-dimensional demographic spaces and dearth of data. Fairness strategies that rely on causality to draw the line between valid and invalid causal channels between features and effects offer more principled reduction of bias. Nonetheless, the problem of identifying the causal structures, estimating the causal consequences based on observational data, and creating causal fairness requirements into working algorithms are all areas of active research. The causal fairness models can be tested using robustness to causal misspecification and confounding but more research is needed on it. The long-term effects and feedback mechanisms of implemented AI systems should be the subject of further research since most of the prevailing studies consider the fixed fairness properties. To foresee accumulation damages, it is crucial to understand the impact of AI systems on upcoming data distributions, individual behavior, and societal structures to prevent and lower them. Dynamics may be made out of simulation studies and longitudinal empirical research. Reasonableness in resource constrained systems where resource constraints make complex bias mitigation methods infeasible is worth consideration, because it turns out that much AI deployment



is in the systems where access to massive computer resources is unavailable. Enhancing access to responsible AI practices by developing efficient fairness-aware algorithms and assessing fairness-accuracy trade-offs in the presence of computational constraints can be achieved.

Table 2: Ethical Frameworks, Governance Approaches, and Application Domains

Sr. No	Application Domain	Ethical Issue	Current Approach	Regulation/Policy	Implementation Challenge	Impact	Future Direction
1	Criminal Justice	Risk assessment bias	Algorithmic risk tools for bail, parole, sentencing	Anti-discrimination laws, due process requirements	Historical data reflects discriminatory practices	Potential for perpetuating racial disparities	Participatory development with affected communities
2	Healthcare	Diagnostic accuracy disparities	AI-assisted diagnosis and treatment recommendation	HIPAA, medical device regulations, clinical validation	Underrepresentation in medical training data	Health outcome disparities across demographics	Diverse clinical trials and validation studies
3	Employment	Hiring discrimination	Automated resume screening, video interview analysis	Equal employment opportunity laws, disability accommodation	Validity of personality inferences from video	Barriers to economic opportunity	Transparency in hiring algorithms, audit rights
4	Financial Services	Credit and lending bias	Automated underwriting, alternative credit scoring	Fair lending laws, adverse action notices	Use of proxy variables for protected attributes	Economic exclusion of disadvantaged groups	Explainable credit decisions, alternative data validation
5	Education	Performance prediction bias	Early warning systems, automated grading	Educational equity requirements, disability protections	Self-fulfilling prophecies from predictions	Educational tracking and opportunity gaps	Student-centered design, teacher-in-loop systems
6	Housing	Rental and mortgage discrimination	Property valuation, tenant screening	Fair housing laws	Historical redlining patterns in data	Residential segregation perpetuation	Geographic fairness constraints
7	Social Media	Content moderation bias	Automated content filtering, recommendation algorithms	Platform liability, content regulations	Cultural variation in content norms	Disproportionate silencing of marginalized voices	Culturally-sensitive moderation, appeal mechanisms
8	Advertising	Discriminatory ad targeting	Behavioral targeting, lookalike audiences	Anti-discrimination in advertising	Correlation between interests and demographics	Disparate access to opportunities	Fairness constraints on audience selection
9	Insurance	Unfair risk assessment	Predictive modeling for premiums and coverage	Insurance anti-discrimination laws	Tension between actuarial fairness and social fairness	Unaffordable coverage for high-risk groups	Risk pooling approaches, subsidization
10	Facial Recognition	Demographic performance gaps	Identification and verification systems	Biometric privacy laws, law enforcement restrictions	Higher error rates for women and darker-skinned individuals	Wrongful arrests, surveillance disparities	Improved datasets, accuracy thresholds by demographic
11	Voice Assistants	Accent and dialect bias	Speech recognition and natural language understanding	Accessibility requirements	Training data dominated by standard dialects	Digital divide for linguistic minorities	Multilingual and multidialectal training
12	Autonomous Vehicles	Safety disparities	Pedestrian detection, collision avoidance	Vehicle safety standards, liability frameworks	Sensor and algorithm performance variation	Differential accident risk across demographics	Diverse testing scenarios and populations
13	Smart Cities	Surveillance and privacy	Predictive policing, traffic management	Privacy regulations, public sector accountability	Over-policing of minority neighborhoods	Erosion of civil liberties	Privacy-by-design, community oversight
14	Mental Health	Diagnostic bias	Symptom screening, treatment	Mental health parity, informed consent	Cultural variation in symptom expression	Misdiagnosis and	Culturally-validated

			recommendations			inappropriate treatment	assessment tools
15	Child Welfare	Risk assessment in family services	Predictive models for intervention decisions	Child protection standards, family rights	Disproportionate intervention in marginalized families	Family separation disparities	Human-centered decision support
16	Immigration	Visa and asylum decision support	Document verification, risk assessment	Immigration law, due process	Language barriers, cultural context gaps	Arbitrary denials with life-altering consequences	Multilingual systems, cultural competence
17	Content Creation	Generative AI bias	Text, image, and video generation	Copyright, right of publicity, deepfake regulations	Reproducing stereotypes, generating harmful content	Spread of misinformation and stereotypes	Safety filtering, watermarking, provenance tracking
18	Search and Retrieval	Ranking bias	Search engine results, information retrieval	Platform transparency requirements	Reinforcing dominant narratives	Information access disparities	Diversity-aware ranking, personalization controls
19	Public Benefits	Eligibility determination	Automated screening for social services	Public benefits regulations, due process	Error-prone exclusion of eligible individuals	Denial of essential services	Error analysis, human review of denials
20	Energy and Utilities	Resource allocation bias	Smart grid management, pricing	Utility regulation, affordability requirements	Differential service quality by neighborhood	Energy poverty in disadvantaged areas	Equitable infrastructure investment
21	Agriculture	Precision agriculture disparities	Crop monitoring, yield prediction	Agricultural support policies	Technology access gaps for small farmers	Concentration in industrial agriculture	Affordable technology for smallholders
22	Legal Services	Case outcome prediction	Legal research, case assessment	Attorney competence requirements	Replicating judicial biases	Unequal justice system outcomes	Decision support not replacement
23	Emergency Response	Resource allocation	Dispatch optimization, resource prioritization	Emergency service standards	Response time disparities by neighborhood	Life-threatening delays in underserved areas	Equity constraints in optimization
24	Environmental Justice	Pollution monitoring gaps	Environmental sensor placement, exposure assessment	Environmental protection regulations	Monitoring gaps in marginalized communities	Invisible pollution exposures	Community-based monitoring, environmental equity
25	Democratic Participation	Voter influence	Microtargeting, persuasion modeling	Election regulations, political advertising rules	Manipulation of democratic processes	Erosion of informed deliberation	Transparency in political advertising, platform accountability

### Conceptual and Theoretical Development

The idea of creating context-specific frameworks of fairness that can be extended to various spheres, cultures, and uses is still a valuable concept in terms of theoretical significance. The definitions of universal fairness might not be suitable at all in the context where other values are upheld or circumstances where trade-offs need to be resolved in different ways. The models that direct contextualization of equity and uphold fundamental positions towards non-discrimination would help the right AI development. The research on the alignment of one value to another, including the way its specification and maximization should be done to achieve the benefits of AI systems, is essential to the work of beneficial AI. Existing reward specification and objective functions give crude approximations of value of what human beings rightly value. Also formulation of ways to elicit value, aggregate value over multi-stakeholders, and value learning would enhance congruence between human interests and AI systems. Philosophical and inquiry into the concept of fairness, the correlation between various fairness criteria, and the norms underpinning various definitions may help to explain conceptual premises and make practical decisions. The challenging trade-offs between competing fairness goals

and efficiency and other values in the society can be ethicalised to provide information required to make tough decisions regarding system design. Theoretical and institutional challenge is having democratic structures of AI governance that will allow meaningful public input to the resulting effects of technology choices. Any system of effective and legitimate public conduct over AI needs to meet the complexity and scale associated with technology, and power disparities amongst creators and populaces.

Empirical evidence on the fairness properties monitored by deployed AI systems, user experiences, and social impacts over time would help to obtain the needed evidence that is currently absent. The majority of the fairness research is based on benchmark dataset and simulated deployment, without answering whether it works in the real world or not. Such studies may be implemented through partnerships between scholars and entities implementing AI without breaching commercial and privacy limitations. Research with comparative effectiveness, assessing various bias mitigation strategies in other domains, datasets, and various criteria of fairness would inform practitioners about which approach to use. Although at the present, it can hardly be compared systematically, it is hard to forecast which methods will be effective in particular applications. Joint evaluation schemes and benchmarks development might help in accumulation of knowledge.

Proper fairness objectives might be informed by user studies on how various stakeholders feel and appreciate various criterion in fairness. Technical optimization of fairness measure might not be in accordance with the priorities of affected people or the subjective measures of fairness. The participatory research processes would be able to base fairness research on user requirements. The research on the effectiveness of regulations that would investigate the impact of a variety of governance options on the AI development practices and performance would guide the policy design. The differences between jurisdictions in their regulatory choices, discussions with practitioners regarding the effects of regulation on their practice and reports of incidences of fairness may help explain the effects of regulation. AI ethics and fairness can be studied using cross-cultural research methods where the idea of fairness, privacy, and the use of technology what is expected in each given culture is exposed. The AI systems deployed around the world have to be sensitive to a variety of values, although most of the research is based on the Western and educated visions. The co-operation of researchers in different cultures across the globe can enhance literacy.

#### **4. Conclusions**

This overview of the literature has analyzed the complex circumstances of morality, discrimination and justice in the artificial intelligence and machine learning systems, and found that they are complex, but they have to be confronted because AI is already making consequential decisions that can influence the lives of humans and the social order of society. This analysis indicates that to enable ethical, unbiased, and fair AI systems, collaborative work is necessary to address the issues through the technical innovations, theoretical formulation, policy intervention, organizational change, and relevant stakeholder engagement. The origins and causes of algorithmic bias can be varied, and manifest across all points of the machine learning system, including the biases in the past that are systematized in the training data, measurement decisions and model selection, through deployment conditions and feedback. These mechanisms of bias are important mechanisms of understanding how to implement effective mitigation strategies and according to the analysis, no technical solution exists that can accommodate all the issues dealing with fairness. Various sources of bias should be dealt with through different means and mitigation techniques should rely on the situations of their implementation and the uniqueness of the fairness goals. The mathematical description of fairness has produced valuable information as well as revealed some fundamental limitations. The fairness metric proliferation reflects real pluralism concerning the concept of fairness in various situations but impossibility results prove that several intuitive fairness concepts cannot be met in any situation other than trivial ones. Such mathematical constraints lead to hard decisions concerning the priority of fairness goals, decisions which concern not only technical factors but also normative decisions concerning reasonableness of justice, values and priorities of people. It is necessary that the field should not pursue a generalized definition of fairness but rather provide the framework to enforce the consideration of the specificity of fairness through the prism of affected stakeholders.

The existing bias mitigation approaches touch the entire machine learning process, including pre-processing strategies that mitigate data quality using in-processing strategies that train models with fairness constraints to post-processing strategies to modify the outputs of models. Although these methods indicate technical proof of enhancing fairness along particular dimensions, the literature shows a lack of empirical support of their usefulness in the real world, insights on fairness-accuracy compromise in applications, as well as helpfulness on practitioners facing complex technical and ethical decisions. The creation of more efficient, stronger and practicable bias mitigation strategies is a dynamic research field that has significant potentials to be innovative. Ethical standards of AI development have spread, and many organizations, governments, as well as multi-stakeholder projects, have advanced their principles and guidelines of responsible AI. Although there is seemingly either unanimous agreement over such high-level principles as beneficence, justice, transparency, and accountability, it is difficult to convert abstractions into technical specifications and organizational behaviors. The disjuncture between aspirational values and operation reality symbolizes certain technical restrictions as well as institutional constraints, such as inappropriate incentives, lack of resources, and ineffective ways to give valuable accountability in cases of harms. The means to fill this gap has to do not only with the improved tools and methodologies but also with organizational cultures that truly emphasize the importance of ethics percolated with the measures of performance. The legal framework of AI ethics and fairness is a swiftly advancing and irregular concept under development. Varied jurisdictions are following divergent paths that represent different legal traditions and policy priorities, which presents complications in compliance to systems implemented in a global fashion and may provide a way to deregulate. The rate of AI advancement surpasses sites of regulation such that once a new regulation is introduced, there is still a gap in governance. Regulation should be sustainable by exercising a variety of goals; benefiting individuals and society by shielding people and society against algorithmic evils, allowing meaningful innovation, guiding developers properly, and keeping up with technological change. To strike this balance, continuous consultations among technologists, policy makers, legal advisers, ethicists and those communities who are affected is necessary.

And domain-specific analyses indicate that the problem of fairness has different manifestations in applications, and the industrial agreement in criminal justice, healthcare, employment, finance, education, and others has both unique technical and ethical concerns. The context is critical: fit fairness standards, permissible accuracy/ fairness trade-offs, and effective stakeholder interaction systems are sensitive to context, to the field of operation, legal provisions, possible harm, and target groups. This context-sensitivity claims that generic solutions should be avoided and claims that domain-specific best practices and governance models should be created. New emerging AI technologies such as large language models, generative AI, autonomous system, and multimodal models raise new ethical issues and change the old ones. The magnitude, ability and generality of those systems enhance both the potential usefulness and the dangers of these systems. The foundation models have a tendency of concentrating the power in organisations that have the capacity to train them and it forms a dependence of the downstream use and therefore there is question as to how it should be governed and accessed. Generative AI features of producing convincing part of accessory material dispute the nature of information and presents new representational evils. The above developments highlight the importance of on-going research, active risk evaluation and responsive governance frameworks.

It takes a long-term effort on all fronts to advance AI towards becoming more ethical, fair and responsible. Technical research should also improve techniques to detect and reduce bias and build more advanced frameworks of fairness to consider intersectionality and dynamics, develop tools to make responsible AI practices more open. Future theoretical proposals are strongly needed to further establish ethical frameworks that are suitable in an AI setting, to define how various concepts of fairness are in relationship to one another, and to give advice on how to balance between conflicting trade-offs. Facts should be tested, and the effectiveness of the offered solutions should be measured in practice in the empirical research, one should comprehend the impact of AI systems on people and society in the long run, and one should base the technical development on the facts about the one that actually works. On top of technical and research innovations, to have responsible AI, organizational, professional, and societal adjustments are needed. Companies creating and implementing AI have to establish internal capabilities of finding and resolving ethical concerns, develop a governance system that holds people

accountable, and develop a culture that truly cares about ethics and performance. Professional society needs to formulate and implement responsible practice standards, train and educate practitioners on AI ethics, and assist practitioners through issues in ethical decision-making. The society needs to make an informed discussion on how AI should be developed and implemented, come up with proper regulatory frameworks, and make sure that the development of AI would disfavor specific commercial or institutional interests as opposed to the overall societal interests. The theme of stakeholder involvement comes out as a critical theme in this review. Victims of AI systems have essential information regarding their requirements, principles, and dangers that he/she might not know since it is costly and technical specialists. It will be critical to conduct meaningful interaction with various stakeholders, especially marginalized groups, which are immediately at a disadvantage when it comes to algorithmic harms, to create AI systems that will indeed represent human interests. Nonetheless, the participation must be effective, which must be based on considering power asymmetries, providing the involvement processes to be available, responsive to the input, and resourceful to be engaged in the participation in the long term. As a challenge and as an ethical duty, constructing these participatory processes is a challenge.

The ethical aspects of AI on a global scale should be given more of the attention that it has not received. The AI systems created in the Western environment, which is based on specific cultural value and assumptions, are implemented all over the world, which can either impose inappropriate standards or simply not take local values into consideration. The field of AI ethics needs to be more globalized, allowing the inclusion of the visions of other cultures, law practices, and economic backgrounds. The discussion of AI ethics cross-culturally has a potential to enhance the knowledge and determine both general and local correctness. In the future, a number of areas of concern can be identified to enhance responsible AI. To start with, it is necessary to work on the practical methods of participatory design and value-sensitive development that should be executed on a massive scale in various settings. Second, establishing appropriate tools of continuous supervision and audit of the deployed AI systems to reveal arising harms and guarantee further observance of the requirements of fairness. Third, defining more effective accountability structures that define accountability throughout the convoluted intelligence AI supply chains and provide effective redresses in case of damages. Fourth, institutional capacity to promote AI ethics by educating, training, and professionalizing and structuring the institution. Fifth, ensuring interdisciplinary cooperation between technical skills and knowledge on the one hand and on the other hand, ethics, law, social sciences, and concerning communities. Sixth the creation of governance structures that are not only innovative and protective but also open to technological fast churning.

The problems of ethical, unbiased, and fair AI are also big and entrenched, and they may be addressed through technical complexity, conceptual pluralism, institutional inertia and underlying conflict of competing values and ideas. The impossibility of such fairness challenges can arise because of mathematical impossibility and due to inherent trade-offs, such problems cannot be addressed purely on the technical level. Nevertheless, the complexity of these issues must motivate more effort as opposed to an underling attitude. Although absolute fairness might be impossible, it is realistic to make great progress with regard to the current practices with the help of a long-term and considerate effort. The stakes are high. Artificial intelligences are becoming more influential in defining the opportunities of lives, services distribution, access to information and amongst social networks in a manner that ultimately can ease or widen the prevailing inequalities. The use of AI technology as an agent of democratization and empowerment or as a tool to concentrate power and keep the injustice going lies in the decisions of the technology developers, organizations, policymakers, and societies. Ethics, prejudice, and equitability in AI are not only technical issues but also constantly humanistic and need wisdom, humility, and a sense of human dignity and social justice. The review is added to the current work on the development of responsible AI since it attempts to synthesize the existing knowledge base, determines gaps and limitations in the current solutions, and formulates priority directions in future endeavors. The journey to the reliable AI systems that would be aligned with human values is lengthy and disordered, and it would involve further learning and adjustment as well as cross-disciplinary and cross-cultural collaboration. The technical innovation is not all that success requires, but also ethical clarity, institutional responsibility, administrative insight and substantial involvement of the individuals

whose lives AI systems touch. The task is immense, yet also tremendous is the chance to influence transformative technology in a manner that will support human blossom and social justice.

### **Author Contributions**

RR: Methodology, software, resources, visualization, writing original draft, writing review and editing, and supervision. NLR: Conceptualization, software, resources, visualization, writing original draft, writing review and editing. AC: Analysis, data collection, methodology, software, resources, visualization, writing review and editing, and supervision. JR: Conceptualization, study design, analysis, data collection.

### **Conflict of interest**

The authors declare no conflicts of interest.

### **References**

- [1] Choudhary OP, Infant SS, Chopra H, Manuta N. Exploring the potential and limitations of artificial intelligence in animal anatomy. *Annals of Anatomy-Anatomischer Anzeiger*. 2025 Feb 1;258:152366. <https://doi.org/10.1016/j.aanat.2024.152366>
- [2] Ooi KB, Tan GW, Al-Emran M, Al-Sharafi MA, Capatina A, Chakraborty A, Dwivedi YK, Huang TL, Kar AK, Lee VH, Loh XM. The potential of generative artificial intelligence across disciplines: Perspectives and future directions. *Journal of Computer Information Systems*. 2025 Jan 2;65(1):76-107. <https://doi.org/10.1080/08874417.2023.2261010>
- [3] Raisch S, Fomina K. Combining human and artificial intelligence: Hybrid problem-solving in organizations. *Academy of Management Review*. 2025 Apr;50(2):441-64. <https://doi.org/10.5465/amr.2021.0421>
- [4] Korteling JE, van de Boer-Visschedijk GC, Blankendaal RA, Boonekamp RC, Eikelboom AR. Human-versus artificial intelligence. *Frontiers in artificial intelligence*. 2021 Mar 25;4:622364. <https://doi.org/10.3389/frai.2021.622364>
- [5] Berente N, Gu B, Recker J, Santhanam R. Managing artificial intelligence. *MIS quarterly*. 2021 Sep 1;45(3):1433-50. <https://doi.org/10.25300/MISQ/2021/16274>
- [6] Nie J, Jiang J, Li Y, Wang H, Ercisli S, Lv L. Data and domain knowledge dual-driven artificial intelligence: Survey, applications, and challenges. *Expert Systems*. 2025 Jan;42(1):e13425. <https://doi.org/10.1111/exsy.13425>
- [7] Nenni ME, De Felice F, De Luca C, Forcina A. How artificial intelligence will transform project management in the age of digitization: a systematic literature review. *Management Review Quarterly*. 2025 Jun;75(2):1669-716. <https://doi.org/10.1007/s11301-024-00418-z>
- [8] Qin C, Zhang L, Cheng Y, Zha R, Shen D, Zhang Q, Chen X, Sun Y, Zhu C, Zhu H, Xiong H. A comprehensive survey of artificial intelligence techniques for talent analytics. *Proceedings of the IEEE*. 2025 Jun 6. <https://doi.org/10.1109/JPROC.2025.3572744>
- [9] Kaack LH, Donti PL, Strubell E, Kamiya G, Creutzig F, Rolnick D. Aligning artificial intelligence with climate change mitigation. *Nature Climate Change*. 2022 Jun;12(6):518-27. <https://doi.org/10.1038/s41558-022-01377-7>
- [10] Ali S, Abuhmed T, El-Sappagh S, Muhammad K, Alonso-Moral JM, Confalonieri R, Guidotti R, Del Ser J, Díaz-Rodríguez N, Herrera F. Explainable Artificial Intelligence (XAI): What we know and what is left to attain Trustworthy Artificial Intelligence. *Information fusion*. 2023 Nov 1;99:101805. <https://doi.org/10.1016/j.inffus.2023.101805>
- [11] Kelly S, Kaye SA, Oviedo-Trespalacios O. What factors contribute to the acceptance of artificial intelligence? A systematic review. *Telematics and informatics*. 2023 Feb 1;77:101925. <https://doi.org/10.1016/j.tele.2022.101925>
- [12] Bhattamisra SK, Banerjee P, Gupta P, Mayuren J, Patra S, Candasamy M. Artificial intelligence in pharmaceutical and healthcare research. *Big Data and Cognitive Computing*. 2023 Jan 11;7(1):10. <https://doi.org/10.3390/bdcc7010010>
- [13] Pallathadka H, Ramirez-Asis EH, Loli-Poma TP, Kaliyaperumal K, Ventayen RJ, Naved M. Applications of artificial intelligence in business management, e-commerce and finance. *Materials Today: Proceedings*. 2023 Jan 1;80:2610-3. <https://doi.org/10.1016/j.matpr.2021.06.419>
- [14] McDonald N, Johri A, Ali A, Collier AH. Generative artificial intelligence in higher education: Evidence from an analysis of institutional policies and guidelines. *Computers in Human Behavior: Artificial Humans*. 2025 Mar 1;3:100121. <https://doi.org/10.1016/j.chbah.2025.100121>

- [15] Emon MM, Khan T. The mediating role of attitude towards the technology in shaping artificial intelligence usage among professionals. *Telematics and Informatics Reports*. 2025 Mar 1;17:100188. <https://doi.org/10.1016/j.teler.2025.100188>
- [16] Shao X, Cai B, Zou Z, Shao H, Yang C, Liu Y. Artificial intelligence enhanced fault prediction with industrial incomplete information. *Mechanical Systems and Signal Processing*. 2025 Feb 1;224:112063. <https://doi.org/10.1016/j.ymssp.2024.112063>
- [17] Rožanec JM, Novalija I, Zajec P, Kenda K, Tavakoli Ghinani H, Suh S, Bian S, Veliou E, Papamartzivanos D, Giannetsos T, Menesidou SA. Human-centric artificial intelligence architecture for industry 5.0 applications. *International journal of production research*. 2023 Oct 18;61(20):6847-72. <https://doi.org/10.1080/00207543.2022.2138611>
- [18] Chander B, John C, Warriar L, Gopalakrishnan K. Toward trustworthy artificial intelligence (TAI) in the context of explainability and robustness. *ACM Computing Surveys*. 2025 Feb 10;57(6):1-49. <https://doi.org/10.1145/3675392>
- [19] Horani OM, Al-Adwan AS, Yaseen H, Hmoud H, Al-Rahmi WM, Alkhalifah A. The critical determinants impacting artificial intelligence adoption at the organizational level. *Information Development*. 2025 Sep;41(3):1055-79. <https://doi.org/10.1177/02666669231166889>
- [20] Varriale V, Cammarano A, Michelino F, Caputo M. Critical analysis of the impact of artificial intelligence integration with cutting-edge technologies for production systems. *Journal of Intelligent Manufacturing*. 2025 Jan;36(1):61-93. <https://doi.org/10.1007/s10845-023-02244-8>
- [21] Fu C, Chen Q. The future of pharmaceuticals: Artificial intelligence in drug discovery and development. *Journal of Pharmaceutical Analysis*. 2025 Feb 26:101248. <https://doi.org/10.1016/j.jpha.2025.101248>
- [22] Vasishta P, Dhingra N, Vasishta S. Application of artificial intelligence in libraries: a bibliometric analysis and visualisation of research activities. *Library Hi Tech*. 2025 May 19;43(2/3):693-710. <https://doi.org/10.1108/LHT-12-2023-0589>
- [23] Samala AD, Rawas S, Wang T, Reed JM, Kim J, Howard NJ, Ertz M. Unveiling the landscape of generative artificial intelligence in education: a comprehensive taxonomy of applications, challenges, and future prospects. *Education and Information Technologies*. 2025 Feb;30(3):3239-78. <https://doi.org/10.1007/s10639-024-12936-0>
- [24] Vora LK, Gholap AD, Jetha K, Thakur RR, Solanki HK, Chavda VP. Artificial intelligence in pharmaceutical technology and drug delivery design. *Pharmaceutics*. 2023 Jul 10;15(7):1916. <https://doi.org/10.3390/pharmaceutics15071916>
- [25] Loureiro SM, Guerreiro J, Tussyadiah I. Artificial intelligence in business: State of the art and future research agenda. *Journal of business research*. 2021 May 1;129:911-26. <https://doi.org/10.1016/j.jbusres.2020.11.001>
- [26] Najjar R. Redefining radiology: a review of artificial intelligence integration in medical imaging. *Diagnostics*. 2023 Aug 25;13(17):2760. <https://doi.org/10.3390/diagnostics13172760>
- [27] Su J, Ng DT, Chu SK. Artificial intelligence (AI) literacy in early childhood education: The challenges and opportunities. *Computers and Education: Artificial Intelligence*. 2023 Jan 1;4:100124. <https://doi.org/10.1016/j.caeai.2023.100124>
- [28] Soori M, Arezoo B, Dastres R. Artificial intelligence, machine learning and deep learning in advanced robotics, a review. *Cognitive Robotics*. 2023 Jan 1;3:54-70. <https://doi.org/10.1016/j.cogr.2023.04.001>
- [29] Bhuyan SS, Sateesh V, Mukul N, Galvankar A, Mahmood A, Nauman M, Rai A, Bordoloi K, Basu U, Samuel J. Generative artificial intelligence use in healthcare: opportunities for clinical excellence and administrative efficiency. *Journal of medical systems*. 2025 Jan 16;49(1):10. <https://doi.org/10.1007/s10916-024-02136-1>
- [30] Cooper G. Examining science education in ChatGPT: An exploratory study of generative artificial intelligence. *Journal of science education and technology*. 2023 Jun;32(3):444-52. <https://doi.org/10.1007/s10956-023-10039-y>
- [31] Aldoseri A, Al-Khalifa KN, Hamouda AM. Re-thinking data strategy and integration for artificial intelligence: concepts, opportunities, and challenges. *Applied Sciences*. 2023 Jan;13(12):7082. <https://doi.org/10.3390/app13127082>
- [32] George B, Wooden O. Managing the strategic transformation of higher education through artificial intelligence. *Administrative Sciences*. 2023 Aug 29;13(9):196. <https://doi.org/10.3390/admsci13090196>
- [33] Wang B, Rau PL, Yuan T. Measuring user competence in using artificial intelligence: validity and reliability of artificial intelligence literacy scale. *Behaviour & information technology*. 2023 Jul 4;42(9):1324-37. <https://doi.org/10.1080/0144929X.2022.2072768>
- [34] Morandini S, Fraboni F, De Angelis M, Puzzo G, Giusino D, Pietrantonio L. The impact of artificial intelligence on workers' skills: Upskilling and reskilling in organisations. *Informing Science*. 2023;26:39-68. <https://doi.org/10.28945/5078>
- [35] Topaz M, Peltonen LM, Michalowski M, Stiglic G, Ronquillo C, Pruinelli L, Song J, O'connor S, Miyagawa S, Fukahori H. The ChatGPT effect: nursing education and generative artificial intelligence. *Journal of Nursing Education*. 2025 Jun 1;64(6):e40-3. <https://doi.org/10.3928/01484834-20240126-01>
- [36] Noy S, Zhang W. Experimental evidence on the productivity effects of generative artificial intelligence. *Science*. 2023 Jul 14;381(6654):187-92. <https://doi.org/10.1126/science.adh2586>



- [37] Nguyen A, Ngo HN, Hong Y, Dang B, Nguyen BP. Ethical principles for artificial intelligence in education. *Education and information technologies*. 2023 Apr;28(4):4221-41. <https://doi.org/10.1007/s10639-022-11316-w>
- [38] Malatji M, Tolah A. Artificial intelligence (AI) cybersecurity dimensions: a comprehensive framework for understanding adversarial and offensive AI. *AI and Ethics*. 2025 Apr;5(2):883-910. <https://doi.org/10.1007/s43681-024-00427-4>
- [39] Shata A, Hartley K. Artificial intelligence and communication technologies in academia: faculty perceptions and the adoption of generative AI. *International Journal of Educational Technology in Higher Education*. 2025 Mar 14;22(1):14. <https://doi.org/10.1186/s41239-025-00511-7>
- [40] Yuan L, Liu X. The effect of artificial intelligence tools on EFL learners' engagement, enjoyment, and motivation. *Computers in Human Behavior*. 2025 Jan 1;162:108474. <https://doi.org/10.1016/j.chb.2024.108474>
- [41] Malhotra G, Ramalingam M. Perceived anthropomorphism and purchase intention using artificial intelligence technology: examining the moderated effect of trust. *Journal of Enterprise Information Management*. 2025 Feb 25;38(2):401-23. <https://doi.org/10.1108/JEIM-09-2022-0316>
- [42] Abbasi BN, Wu Y, Luo Z. Exploring the impact of artificial intelligence on curriculum development in global higher education institutions. *Education and Information Technologies*. 2025 Jan;30(1):547-81. <https://doi.org/10.1007/s10639-024-13113-z>
- [43] Baig MI, Yadegaridehkordi E. Factors influencing academic staff satisfaction and continuous usage of generative artificial intelligence (GenAI) in higher education. *International Journal of Educational Technology in Higher Education*. 2025 Feb 3;22(1):5. <https://doi.org/10.1186/s41239-025-00506-4>
- [44] Bewersdorff A, Hartmann C, Hornberger M, Seßler K, Bannert M, Kasneci E, Kasneci G, Zhai X, Nerdel C. Taking the next step with generative artificial intelligence: The transformative role of multimodal large language models in science education. *Learning and Individual Differences*. 2025 Feb 1;118:102601. <https://doi.org/10.1016/j.lindif.2024.102601>
- [45] Fang B, Yu J, Chen Z, Osman AI, Farghali M, Ihara I, Hamza EH, Rooney DW, Yap PS. Artificial intelligence for waste management in smart cities: a review. *Environmental Chemistry Letters*. 2023 Aug;21(4):1959-89. <https://doi.org/10.1007/s10311-023-01604-3>
- [46] Kaur R, Gabrijelčič D, Klobučar T. Artificial intelligence for cybersecurity: Literature review and future research directions. *Information Fusion*. 2023 Sep 1;97:101804. <https://doi.org/10.1016/j.inffus.2023.101804>
- [47] Kamalov F, Santandreu Calonge D, Gurrib I. New era of artificial intelligence in education: Towards a sustainable multifaceted revolution. *Sustainability*. 2023 Aug 16;15(16):12451. <https://doi.org/10.3390/su151612451>
- [48] Naik N, Hameed BM, Shetty DK, Swain D, Shah M, Paul R, Aggarwal K, Ibrahim S, Patil V, Smriti K, Shetty S. Legal and ethical consideration in artificial intelligence in healthcare: who takes responsibility?. *Frontiers in surgery*. 2022 Mar 14;9:862322. <https://doi.org/10.3389/fsurg.2022.862322>
- [49] Le Berre C, Sandborn WJ, Aridhi S, Devignes MD, Fournier L, Smail-Tabbone M, Danese S, Peyrin-Biroulet L. Application of artificial intelligence to gastroenterology and hepatology. *Gastroenterology*. 2020 Jan 1;158(1):76-94. <https://doi.org/10.1053/j.gastro.2019.08.058>
- [50] Rawashdeh A. The consequences of artificial intelligence: an investigation into the impact of AI on job displacement in accounting. *Journal of Science and Technology Policy Management*. 2025 Mar 6;16(3):506-35. <https://doi.org/10.1108/JSTPM-02-2023-0030>
- [51] Derakhshan A, Teo T, Khazaie S. Investigating the usefulness of artificial intelligence-driven robots in developing empathy for English for medical purposes communication: The role-play of Asian and African students. *Computers in Human Behavior*. 2025 Jan 1;162:108416. <https://doi.org/10.1016/j.chb.2024.108416>
- [52] Tamsah A, Alhasan K, Altamimi I, Jamal A, Al-Eyadhy A, Malki KH, Tamsah MH. DeepSeek in healthcare: revealing opportunities and steering challenges of a new open-source artificial intelligence frontier. *Cureus*. 2025 Feb 18;17(2). <https://doi.org/10.7759/cureus.79221>
- [53] Wang Y, Wang L, Siau KL. Human-centered interaction in virtual worlds: A new era of generative artificial intelligence and metaverse. *International Journal of Human-Computer Interaction*. 2025 Jan 17;41(2):1459-501. <https://doi.org/10.1080/10447318.2024.2316376>
- [54] Huynh-The T, Pham QV, Pham XQ, Nguyen TT, Han Z, Kim DS. Artificial intelligence for the metaverse: A survey. *Engineering Applications of Artificial Intelligence*. 2023 Jan 1;117:105581. <https://doi.org/10.1016/j.engappai.2022.105581>
- [55] Briganti G, Le Moine O. Artificial intelligence in medicine: today and tomorrow. *Frontiers in medicine*. 2020 Feb 5;7:509744. <https://doi.org/10.3389/fmed.2020.00027>
- [56] Swiecki Z, Khosravi H, Chen G, Martinez-Maldonado R, Lodge JM, Milligan S, Selwyn N, Gašević D. Assessment in the age of artificial intelligence. *Computers and Education: Artificial Intelligence*. 2022 Jan 1;3:100075. <https://doi.org/10.1016/j.caeai.2022.100075>

- [57] Khalid N, Qayyum A, Bilal M, Al-Fuqaha A, Qadir J. Privacy-preserving artificial intelligence in healthcare: Techniques and applications. *Computers in Biology and Medicine*. 2023 May 1;158:106848. <https://doi.org/10.1016/j.combiomed.2023.106848>
- [58] Resnik DB, Hosseini M. The ethics of using artificial intelligence in scientific research: new guidance needed for a new tool. *AI and Ethics*. 2025 Apr;5(2):1499-521. <https://doi.org/10.1007/s43681-024-00493-8>
- [59] Kraemer MU, Tsui JL, Chang SY, Lytras S, Khurana MP, Vanderslott S, Bajaj S, Scheidwasser N, Curran-Sebastian JL, Semenova E, Zhang M. Artificial intelligence for modelling infectious disease epidemics. *Nature*. 2025 Feb 20;638(8051):623-35. <https://doi.org/10.1038/s41586-024-08564-w>
- [60] Demaidi MN. Artificial intelligence national strategy in a developing country. *Ai & Society*. 2025 Feb;40(2):423-35. <https://doi.org/10.1007/s00146-023-01779-x>
- [61] Neumann O, Guirguis K, Steiner R. Exploring artificial intelligence adoption in public organizations: a comparative case study. *Public Management Review*. 2024 Jan 2;26(1):114-41. <https://doi.org/10.1080/14719037.2022.2048685>
- [62] Vishwakarma LP, Singh RK, Mishra R, Kumari A. Application of artificial intelligence for resilient and sustainable healthcare system: Systematic literature review and future research directions. *International Journal of Production Research*. 2025 Jan 17;63(2):822-44. <https://doi.org/10.1080/00207543.2023.2188101>
- [63] Manikandan S, Kaviya RS, Shreeharan DH, Subbaiya R, Vickram S, Karmegam N, Kim W, Govarthan M. Artificial intelligence-driven sustainability: Enhancing carbon capture for sustainable development goals-A review. *Sustainable Development*. 2025 Apr;33(2):2004-29. <https://doi.org/10.1002/sd.3222>
- [64] Choi J, Woo S, Ferrell A. Artificial intelligence assisted telehealth for nursing: A scoping review. *Journal of Telemedicine and Telecare*. 2025 Jan;31(1):140-9. <https://doi.org/10.1177/1357633X231167613>
- [65] Buyuktepe O, Catal C, Kar G, Bouzembrak Y, Marvin H, Gavai A. Food fraud detection using explainable artificial intelligence. *Expert Systems*. 2025 Jan;42(1):e13387. <https://doi.org/10.1111/exsy.13387>
- [66] Pereira V, Hadjielias E, Christofi M, Vrontis D. A systematic literature review on the impact of artificial intelligence on workplace outcomes: A multi-process perspective. *Human Resource Management Review*. 2023 Mar 1;33(1):100857. <https://doi.org/10.1016/j.hrmr.2021.100857>
- [67] Cows J, Tsamados A, Taddeo M, Floridi L. The AI gambit: leveraging artificial intelligence to combat climate change-opportunities, challenges, and recommendations. *Ai & Society*. 2023 Feb;38(1):283-307. <https://doi.org/10.1007/s00146-021-01294-x>
- [68] Dogan ME, Goru Dogan T, Bozkurt A. The use of artificial intelligence (AI) in online learning and distance education processes: A systematic review of empirical studies. *Applied sciences*. 2023 Feb 27;13(5):3056. <https://doi.org/10.3390/app13053056>
- [69] Mariani MM, Machado I, Magrelli V, Dwivedi YK. Artificial intelligence in innovation research: A systematic review, conceptual framework, and future research directions. *Technovation*. 2023 Apr 1;122:102623. <https://doi.org/10.1016/j.technovation.2022.102623>
- [70] Wach K, Duong CD, Ejdy J, Kazlauskaitė R, Korzynski P, Mazurek G, Paliszkievicz J, Ziemba E. The dark side of generative artificial intelligence: A critical analysis of controversies and risks of ChatGPT. *Entrepreneurial Business and Economics Review*. 2023 Jun 30;11(2):7-30. <https://doi.org/10.15678/EBER.2023.110201>
- [71] Huang C, Zhang Z, Mao B, Yao X. An overview of artificial intelligence ethics. *IEEE Transactions on Artificial Intelligence*. 2022 Jul 28;4(4):799-819. <https://doi.org/10.1109/TAI.2022.3194503>
- [72] Son JB, Ružić NK, Philpott A. Artificial intelligence technologies and applications for language learning and teaching. *Journal of China Computer-Assisted Language Learning*. 2025 May 23;5(1):94-112. <https://doi.org/10.1515/jccall-2023-0015>
- [73] Ahmed SR, Baghdadi R, Bernadskiy M, Bowman N, Braid R, Carr J, Chen C, Ciccarella P, Cole M, Cooke J, Desai K. Universal photonic artificial intelligence acceleration. *Nature*. 2025 Apr 10;640(8058):368-74. <https://doi.org/10.1038/s41586-025-08854-x>
- [74] Camps-Valls G, Fernández-Torres MÁ, Cohrs KH, Höhl A, Castelletti A, Pacal A, Robin C, Martinuzzi F, Papoutsis I, Prapas I, Pérez-Aracil J. Artificial intelligence for modeling and understanding extreme weather and climate events. *Nature Communications*. 2025 Feb 24;16(1):1919. <https://doi.org/10.1038/s41467-025-56573-8>
- [75] Kamila MK, Jasrotia SS. Ethical issues in the development of artificial intelligence: recognizing the risks. *International Journal of Ethics and Systems*. 2025 Jan 30;41(1):45-63. <https://doi.org/10.1108/IJOES-05-2023-0107>
- [76] Wang Q, Li Y, Li R. Integrating artificial intelligence in energy transition: a comprehensive review. *Energy Strategy Reviews*. 2025 Jan 1;57:101600. <https://doi.org/10.1016/j.esr.2024.101600>
- [77] Zhai X, Chu X, Chai CS, Jong MS, Istenic A, Spector M, Liu JB, Yuan J, Li Y. A Review of Artificial Intelligence (AI) in Education from 2010 to 2020. *Complexity*. 2021;2021(1):8812542. <https://doi.org/10.1155/2021/8812542>

- [78] Zhang C, Lu Y. Study on artificial intelligence: The state of the art and future prospects. *Journal of Industrial Information Integration*. 2021 Sep 1;23:100224. <https://doi.org/10.1016/j.jii.2021.100224>
- [79] Goralski MA, Tan TK. Artificial intelligence and sustainable development. *The International Journal of Management Education*. 2020 Mar 1;18(1):100330. <https://doi.org/10.1016/j.ijme.2019.100330>
- [80] Radanliev P. Artificial intelligence: reflecting on the past and looking towards the next paradigm shift. *Journal of Experimental & Theoretical Artificial Intelligence*. 2025 Oct 3;37(7):1045-62. <https://doi.org/10.1080/0952813X.2024.2323042>
- [81] Zhang X, Wang L, Helwig J, Luo Y, Fu C, Xie Y, Liu M, Lin Y, Xu Z, Yan K, Adams K. Artificial intelligence for science in quantum, atomistic, and continuum systems. *Foundations and Trends® in Machine Learning*. 2025 Jul 20;18(4):385-912. <https://doi.org/10.1561/22000000115>
- [82] McIntosh TR, Susnjak T, Arachchilage N, Liu T, Xu D, Watters P, Halgamuge MN. Inadequacies of large language model benchmarks in the era of generative artificial intelligence. *IEEE Transactions on Artificial Intelligence*. 2025 May 13. <https://doi.org/10.1109/TAI.2025.3569516>
- [83] Hasanzadeh F, Josephson CB, Waters G, Adedinsewo D, Azizi Z, White JA. Bias recognition and mitigation strategies in artificial intelligence healthcare applications. *NPJ Digital Medicine*. 2025 Mar 11;8(1):154. <https://doi.org/10.1038/s41746-025-01503-7>
- [84] Alqahtani N, Wafula Z. Artificial intelligence integration: Pedagogical strategies and policies at leading universities. *Innovative Higher Education*. 2025 Apr;50(2):665-84. <https://doi.org/10.1007/s10755-024-09749-x>
- [85] Xu Y, Liu X, Cao X, Huang C, Liu E, Qian S, Liu X, Wu Y, Dong F, Qiu CW, Qiu J. Artificial intelligence: A powerful paradigm for scientific research. *The Innovation*. 2021 Nov 28;2(4). <https://doi.org/10.1016/j.xinn.2021.100179>
- [86] Galante N, Cotroneo R, Furci D, Lodetti G, Casali MB. Applications of artificial intelligence in forensic sciences: Current potential benefits, limitations and perspectives. *International journal of legal medicine*. 2023 Mar;137(2):445-58. <https://doi.org/10.1007/s00414-022-02928-5>
- [87] Sheikh H, Prins C, Schrijvers E. Artificial intelligence: definition and background. In *Mission AI: The new system technology* 2023 Jan 31 (pp. 15-41). Cham: Springer International Publishing. [https://doi.org/10.1007/978-3-031-21448-6\\_2](https://doi.org/10.1007/978-3-031-21448-6_2)
- [88] Stone P, Brooks R, Brynjolfsson E, Calo R, Etzioni O, Hager G, Hirschberg J, Kalyanakrishnan S, Kamar E, Kraus S, Leyton-Brown K. Artificial intelligence and life in 2030: the one hundred year study on artificial intelligence. *arXiv preprint arXiv:2211.06318*. 2022 Oct 31.
- [89] Taeihagh A. Governance of artificial intelligence. *Policy and society*. 2021 Jun;40(2):137-57. <https://doi.org/10.1080/14494035.2021.1928377>
- [90] Kulkarni S, Seneviratne N, Baig MS, Khan AH. Artificial intelligence in medicine: where are we now?. *Academic radiology*. 2020 Jan 1;27(1):62-70. <https://doi.org/10.1016/j.acra.2019.10.001>
- [91] Giuggioli G, Pellegrini MM. Artificial intelligence as an enabler for entrepreneurs: a systematic literature review and an agenda for future research. *International Journal of Entrepreneurial Behavior & Research*. 2023 May 4;29(4):816-37. <https://doi.org/10.1108/IJEBR-05-2021-0426>
- [92] Kshetri N, Dwivedi YK, Davenport TH, Panteli N. Generative artificial intelligence in marketing: Applications, opportunities, challenges, and research agenda. *International Journal of Information Management*. 2024 Apr 1;75:102716. <https://doi.org/10.1016/j.ijinfomgt.2023.102716>
- [93] Ahmad SF, Han H, Alam MM, Rehmat M, Irshad M, Arraño-Muñoz M, Ariza-Montes A. Impact of artificial intelligence on human loss in decision making, laziness and safety in education. *Humanities and Social Sciences Communications*. 2023 Jun 9;10(1):1-4. <https://doi.org/10.1057/s41599-023-01842-4>
- [94] Bidyalakshmi T, Jyoti B, Mansuri SM, Srivastava A, Mohapatra D, Kalnar YB, Narsaiah K, Indore N. Application of artificial intelligence in food processing: Current status and future prospects. *Food Engineering Reviews*. 2025 Mar;17(1):27-54. <https://doi.org/10.1007/s12393-024-09386-2>
- [95] Law R, Ye H, Lei SS. Ethical artificial intelligence (AI): principles and practices. *International Journal of Contemporary Hospitality Management*. 2025 Jan 2;37(1):279-95. <https://doi.org/10.1108/IJCHM-04-2024-0482>
- [96] Cukurova M. The interplay of learning, analytics and artificial intelligence in education: A vision for hybrid intelligence. *British Journal of Educational Technology*. 2025 Mar;56(2):469-88. <https://doi.org/10.1111/bjet.13514>
- [97] Agha RA, Mathew G, Rashid R, Kerwan A, Al-Jabir A, Sohrabi C, Franchi T, Nicola M, Agha M. Transparency in the reporting of artificial intelligence-the TITAN guideline. *Premier Journal of Science*. 2025;10:100082. <https://doi.org/10.70389/PJS.100082>

- [98] Hanna MG, Pantanowitz L, Dash R, Harrison JH, Deebajah M, Pantanowitz J, Rashidi HH. Future of artificial intelligence (AI)-machine learning (ML) trends in pathology and medicine. *Modern Pathology*. 2025 Jan 4;100705. <https://doi.org/10.1016/j.modpat.2025.100705>
- [99] Chen E, Prakash S, Janapa Reddi V, Kim D, Rajpurkar P. A framework for integrating artificial intelligence for clinical care with continuous therapeutic monitoring. *Nature Biomedical Engineering*. 2025 Apr;9(4):445-54. <https://doi.org/10.1038/s41551-023-01115-0>
- [100] Liu SY. Artificial intelligence (AI) in agriculture. *IT professional*. 2020 May 21;22(3):14-5. <https://doi.org/10.1109/MITP.2020.2986121>
- [101] Shimizu H, Nakayama KI. Artificial intelligence in oncology. *Cancer science*. 2020 May;111(5):1452-60. <https://doi.org/10.1111/cas.14377>
- [102] Schwendicke FA, Samek W, Krois J. Artificial intelligence in dentistry: chances and challenges. *Journal of dental research*. 2020 Jul;99(7):769-74. <https://doi.org/10.1177/0022034520915714>
- [103] Verganti R, Vendraminelli L, Iansiti M. Innovation and design in the age of artificial intelligence. *Journal of product innovation management*. 2020 May;37(3):212-27. <https://doi.org/10.1111/jpim.12523>
- [104] Novelli C, Taddeo M, Floridi L. Accountability in artificial intelligence: What it is and how it works. *Ai & Society*. 2024 Aug;39(4):1871-82. <https://doi.org/10.1007/s00146-023-01635-y>
- [105] Rashidi HH, Pantanowitz J, Hanna MG, Tafti AP, Sanghani P, Buchinsky A, Fennell B, Deebajah M, Wheeler S, Pearce T, Abukhiran I. Introduction to artificial intelligence and machine learning in pathology and medicine: generative and nongenerative artificial intelligence basics. *Modern Pathology*. 2025 Apr 1;38(4):100688. <https://doi.org/10.1016/j.modpat.2024.100688>
- [106] Waisberg E, Ong J, Kamran SA, Masalkhi M, Paladugu P, Zaman N, Lee AG, Tavakkoli A. Generative artificial intelligence in ophthalmology. *Survey of ophthalmology*. 2025 Jan 1;70(1):1-1. <https://doi.org/10.1016/j.survophthal.2024.04.009>
- [107] Wang H, Fu T, Du Y, Gao W, Huang K, Liu Z, Chandak P, Liu S, Van Katwyk P, Deac A, Anandkumar A. Scientific discovery in the age of artificial intelligence. *Nature*. 2023 Aug 3;620(7972):47-60. <https://doi.org/10.1038/s41586-023-06221-2>
- [108] Banh L, Strobel G. Generative artificial intelligence. *Electronic Markets*. 2023 Dec;33(1):63. <https://doi.org/10.1007/s12525-023-00680-1>