# Beyond Engagement: Engineering Algorithmic Resilience Against Societal Collapse

## I. Diagnosis of Algorithmic Collapse: Homogenization and Extremism

The emergence of profound political polarization, the rise of echo chambers, and the facilitation of extreme discourse within digital ecosystems signals a fundamental structural failure in current machine learning (ML) paradigms. This systemic failure can be understood as an engineered fragility resulting from the misalignment between optimization objectives and long-term societal stability.

### A. The Technical Parallel: Mode Collapse in Generative Models and Recommender Systems

The societal collapse into homogenous digital environments is the macrocosm of a well-documented technical failure state in deep learning known as *mode collapse*. Mode collapse occurs primarily in Generative Adversarial Networks (GANs) and other generative models. When models are not properly regularized, they tend to produce only a limited variety of outputs—or "modes"—failing to capture the full diversity and richness of the real data distribution. This phenomenon risks perpetuating and reinforcing existing stereotypes and unfairness derived from biases in the training data.

In a social context, the recommendation algorithms, functioning as generative systems for user experience, have collapsed into modes defined by high engagement metrics, sacrificing the breadth of human thought and perspective. This convergence towards predictable, high-interaction content confirms that optimizing solely for extrinsic, narrow signals causes the ML system to inherently trend toward simplified, limited distributions, resulting in a homogenous and fragile informational landscape.

### B. The Social Manifestation: Optimization Myopia and the Extremism Trap

The core driver of this collapse is **optimization myopia**, a focus on short-term, quantifiable metrics at the expense of unintended societal consequences. Algorithmic recommender systems (RS) optimize overwhelmingly for metrics such as clicks, likes, and time spent (predicted engagement), a strategy intrinsically linked to advertising-based business models. This aggressive optimization strategy exploits fundamental human psychological biases. Humans are predisposed to prioritizing morally and emotionally charged information for social learning, a mechanism evolved for group cooperation and survival. However, in modern, diverse, and complex digital communities, this predisposition is exploited by algorithms that select information primarily to boost engagement metrics. The causal chain is clear: optimizing

for engagement incentivizes the promotion of sensationalist and divisive content. This reduction in cognitive diversity, driven by optimization for virality, correlates directly with increased polarization and the spread of misinformation. If this divisiveness escalates, the threat to democratic stability becomes existential.

### C. Nuanced Causality: The Socio-Technical Problem of Polarization

A critical technical understanding is required to move beyond the simplistic "filter bubble" hypothesis, which posits that algorithms automatically isolate passive users. Empirical research suggests that the evidence supporting the simple filter bubble claim is mixed. Studies in various countries indicate that algorithmic selection often leads to slightly *more diverse* news consumption for the general public, the opposite of the filter bubble thesis.

However, this complexity does not exonerate the algorithms. The issue is better framed as a **reinforcement hypothesis**: the algorithm is not necessarily the *primary cause* of isolation; rather, it acts as a hyper-efficient amplifier of pre-existing homophily and user *selective exposure*. The research indicates that self-selection by a small minority of highly partisan individuals is the primary driver of echo chambers. For this minority, the algorithm's engagement objective maximizes their radicalization potential by feeding their bias.

This understanding implies that interventions cannot simply revert to chronological or non-personalized feeds, which merely offer a false choice. Instead, interventions must focus on disrupting these user-driven feedback loops. Furthermore, users who become critically aware of how these systems shape their experience often develop **algorithmic cynicism**, a state of resignation rooted in perceived lack of agency. This cynicism leads to digital disengagement rather than constructive, corrective behavior. Therefore, future algorithmic design must transition from acting as an amplifier of existing bias to acting as a *counter-agent* designed to foster resilience and constructive friction, demanding solutions that restore transparency and user control.

# II. Paradigm Shift 1: Exploration-Based Learning (EBL) for Viewpoint Diversity

Exploration-Based Learning (EBL) offers a foundational technical mechanism to counter the optimization-driven exploitation loop by systematically injecting novelty and uncertainty into the recommender system's decision process.

## A. Intrinsic Motivation and Curiosity-Driven Exploration (CDE)

In traditional Reinforcement Learning (RL), the sparsity of external rewards presents a major hurdle for effective learning. EBL addresses this by introducing **intrinsic rewards**—internal signals derived from factors such as novelty, state visitation frequency, or prediction error—to motivate the agent to actively explore novel states, even in the absence of external incentives. This concept is formalized through **Curiosity-Driven Exploration (CDE)**. CDE frameworks typically derive curiosity signals from both the actor and the critic components of the RL architecture. The intuition is that agents, similar to children in early cognitive development, display uncertainty and exploratory behavior when confronted with novel situations, propelled by an intrinsic drive to learn and explore.

In the application to social systems, the "novel state" is defined as exposure to information that is divergent from the user's historical profile or counter-attitudinal to their established political cluster. The intrinsic reward mechanism translates directly to rewarding the algorithm for

presenting content that is **highly uncertain in terms of immediate user acceptance but high in informational divergence**. Crucially, the bonus weight ($\omega_t$) is often dynamically annealed to facilitate aggressive exploration in early stages before shifting focus towards exploitation as the policy converges.

## B. Novelty Search and Diversity Maximization

To effectively counter the homogeneity of the echo chamber, EBL leverages principles from novelty and diversity maximization research. Novelty search, defined as identifying prior art that discloses similar or identical inventions, is translated in RS to finding items highly dissimilar to items the user has already interacted with. Maximizing novelty and diversity is critical for mitigating filter bubbles and increasing long-term user engagement and satisfaction. Diversity, specifically intra-list diversity (items from different sources or viewpoints), enhances user satisfaction.

The traditional challenge lies in the apparent dilemma between maximizing recommendation accuracy (exploitation) and maximizing diversity (exploration). Research has shown that this trade-off can be successfully resolved using **hybrid algorithms** that combine accuracy-focused and diversity-focused components, achieving simultaneous gains in both metrics without relying on semantic or context-specific information.

Diversity can also be imposed through **technical regularization** during training. For instance, diversity-promoting priors can be imposed over the multiple choices obtained from inference, ensuring that the model provides representations with more complementary information, rather than merely slight variations of existing choices (e.g., preventing outputs that are simply "one-pixel shifted versions" of others). Furthermore, spectral analysis of model weights reveals that increasing data diversity alters the weight distribution similarly to other regularization techniques, confirming data diversity as a robust form of regularization.

## C. Robustness through Multi-Agent Reinforcement Learning and Mutual Information Regularization

The social environment, especially when addressing polarization, must be modeled as a complex, non-stationary system, making Multi-Agent Reinforcement Learning (MARL) the appropriate framework for analysis. In cooperative MARL settings, ensuring robustness against adversarial or unpredictable actions—such as coordinated campaigns or the polarizing behavior of highly partisan actors—is crucial for real-world deployment. The number of potential threat scenarios grows exponentially with the number of agents.

A promising technical defense against this complexity is the framing of robust MARL as a control-as-inference problem, leveraging **Mutual Information Regularization as Robust Regularization ($\text{MIR}^3$)**. $\text{MIR}^3$ maximizes a lower bound on robustness during routine training, implicitly optimizing worst-case robustness across all threat scenarios. This regularization functions as an information bottleneck, preventing the agents (i.e., the components of the recommender system) from over-reacting to others' actions and aligning their policies with robust action priors. This mechanism effectively implements an architectural "caution," ensuring the system does not brittly over-optimize for the highest, often most extreme, engagement signal, thereby defending against the extremism trap and ensuring improved generalization across varied scenarios.

The combined strength of EBL approaches is summarized below:

Table 1: Comparison of EBL Metrics and Their Anti-Polarization Function

| EBL Metric | Technical Definition | Societal Benefit |
|---|---|---|
| Novelty (Recommendation) | Dissimilarity to user's historical interaction profile. | Breaks user habit loops; exposes niche/out-of-party content, mitigating filter bubbles. |
| Intrinsic Reward (CDE) | Prediction error or state visitation count bonus in novel information clusters. | Encourages proactive exploration of less-traveled, diverse information spaces, driving curiosity. |
| Robustness Regularization ($\text{MIR}^3$) | Maximizing lower bound on robustness against adversarial agent policies. | Promotes systemic caution, prevents over-reaction to divisive/extreme content, ensuring non-brittleness. |

# III. Paradigm Shift 2: Creation-Based Learning (CBL) for Contested Discourse

Creation-Based Learning (CBL) moves beyond passively recommending existing content (EBL) to actively leveraging generative models to create necessary informational resources, enabling a proactive intervention in discourse dynamics.

## A. Generative Models as Proactive Intervention Tools

Generative AI models learn the patterns and distributions of training data and apply those understandings to produce novel content—text, images, audio, or code. This capability, central to computational creativity , allows the algorithm to fill informational voids that EBL identifies but cannot satisfy through existing content retrieval. Where EBL excels at identifying the *gap* in the information landscape, CBL provides the capacity to *create* the resource needed to form bridging narratives or introduce specific, necessary viewpoints.

The success of this approach is highly contingent on structural regulation. If deployed without architectural safeguards, generative AI is susceptible to the same mode collapse that plagues recommender systems, risking the perpetuation of bias rather than mitigating it.

## B. Diversity Regulation in Generative Outputs

To prevent the CBL engine from simply generating novel content that conforms to existing bias modes, **diversity regulation** is mandatory. Generative Adversarial Networks (GANs), for example, often suffer from mode collapse where an input is always mapped to a single output despite the existence of multi-modal conditional distributions.

The technical solution is to propose explicit **diversity-sensitive regularization** on the generator. This regularization is simple, general, and can be integrated into most conditional GAN objectives. Its purpose is to explicitly regularize the generator to produce diverse outputs depending on the latent codes, thereby controlling the balance between output quality and diversity. Applied to large language models (LLMs) used for content generation, this regularization would prevent the system from consistently generating narratives that align solely

with the statistical majority or the most extreme, high-engagement modes of the training data.

## C. Conceptualizing Counter-Attitudinal Content Synthesis (C-ACS)

A core testable intervention leveraging CBL is the systematic generation of civil, viewpoint-diverse content, termed **Counter-Attitudinal Content Synthesis (C-ACS)**. This synthesized content would be actively introduced into polarized social networks in a manner optimized to shift discourse dynamics.

The fundamental goal of C-ACS is **conflict transformation**, which seeks not to suppress or eliminate conflict, but to make it more constructive. The generated narratives should bridge existing divides, offer civil counter-attitudinal arguments, and provide missing context.

However, the deployment of GenAI tools in social contexts is complex. Studies show that while some AI tools may increase user engagement and volume of content, they can simultaneously decrease the perceived quality and authenticity of discussion, potentially introducing negative spill-over effects. Therefore, C-ACS deployment must adhere to strict ethical and technical principles:

1.  **Transparency and Disclosure:** There must be a transparent disclosure mechanism, ensuring end-users are aware that they are interacting with AI-generated content (e.g., deepfake or LLM labeling).
2.  **Civility Metric Alignment:** Generated content must be rigorously aligned with civility metrics to avoid sensationalizing outrage. The generative process must be constrained by an objective function that penalizes content likely to trigger uncivil responses.
3.  **Context-Sensitivity:** The generation process must be context-sensitive, accounting for both the topic and the user's intent to maximize the pedagogical or depolarizing value of the synthesized content.

The synergy between EBL and CBL is crucial: EBL, driven by intrinsic rewards, searches for informational gaps, and the C-ACS generative engine provides the necessary high-quality, civil, missing viewpoints that the EBL system is rewarded for exploring.

# IV. Redefining Algorithmic Objectives and Feedback Loops using Causal Inference

To achieve durable societal alignment, algorithms must shift their basis of decision-making from correlation (what users click on) to causation (what interventions reliably lead to improved social outcomes).

## A. Causal Machine Learning (CML) and Breaking the Loop

Recommender systems exist in a dynamic, iterative environment characterized by feedback loops: an item is recommended, the user reacts, and that reaction influences future recommendations. If unchecked, these cyclic processes can lead to narrowed content exposure and the formation of echo chambers.

Relying solely on correlation without considering the underlying causal mechanisms leads to a range of practical issues, including problems related to fairness, explainability, robustness, bias, and echo chamber formation. **Causal Machine Learning (CML)** methods are necessary to evaluate the heterogeneous effects and practical impact of social interventions, allowing decision-makers to determine best practices for mitigation.

A key technical advancement is the **Dynamic Causal Collaborative Filtering (\partialCCF)** model. \partialCCF explicitly designs a causal graph that incorporates these feedback loops (which are typically assumed to be Directed Acyclic Graphs, but in RS are cyclic) to accurately model the dynamic data generation process. \partialCCF mitigates echo chambers by estimating users' post-intervention preference on items based on **back-door adjustment** and mitigating the echo chamber effect through **counterfactual reasoning**. This approach allows the system to determine the true causal effect of an exposure, effectively breaking the correlation-causation fallacy that mistakes social media exposure for the primary cause of radicalization.

## B. Anti-Polarization Loss Functions (APLF) for Continuous Alignment

Given that polarization is implicated in the erosion of democracy and the progression to violence, the polarization properties of large algorithmic systems are a matter of concern for peace and security.

A radical technical intervention involves engineering polarization measures as explicit **algorithmic objectives**. Polarization metrics, such as the widely used "feeling thermometer" (a measure of affective polarization), must be continuously monitored and integrated into the loss function. This strategy allows polarization survey feedback to be used as a direct **reward signal** for Reinforcement Learning-based recommendation algorithms. By optimizing against this feedback, the system is designed to learn what interventions actually depolarize and can continuously adapt to complex social changes.

Crucially, the **Anti-Polarization Loss Function (APLF)** must incorporate *civility metrics* alongside diversity. Empirical studies of online conflict suggest that diversity-based interventions alone can sometimes worsen polarization. Using civility metrics in conjunction with diversity is a potentially more effective approach, as it de-prioritizes uncivil presentations of outgroup opinions. Therefore, the APLF must be treated as a necessary **constraint** on engagement maximization, actively penalizing the promotion of uncivil conflict while simultaneously rewarding content that facilitates civil cross-cutting exposure.

## C. Counterfactual Fairness and Systemic Bias Mitigation

Algorithms trained on historical data often embed and magnify societal biases related to race, gender, or other demographic groups. This lack of fairness can lead to negative outcomes in high-stakes decisions across financial, healthcare, and legal systems.

**Counterfactual Fairness (CF)** offers a rigorous causal methodology to address this. CF ensures that the outcome of an algorithm-led decision remains the same for an individual in the actual world and a 'counterfactual world' where their sensitive demographic attributes are altered. Implementing CF allows model researchers to explicitly identify embedded biases, for example, by demonstrating that a lending model suggests a change in gender is required for loan approval.

When applied to content regulation and social dynamics, CF ensures that the promotion or suppression of dissent is based strictly on the content and its adherence to civility/polarization metrics, rather than the sensitive demographic identity of the user or community. This is vital for mitigating unintentional algorithmic censorship or downranking that could silence marginalized communities or legitimate dissent.

# V. Architectural and Governance Resilience

Algorithmic collapse cannot be treated merely as a modeling error; it is a systemic failure rooted in centralized architectural control and misaligned incentives. True resilience requires shifting power dynamics and enforcing continuous oversight.

## A. Decentralization and Federated Learning (FL)

Centralized architectures, where data is consolidated on a single server, are a fundamental enabler of systemic collapse, concentrating power and making algorithms vulnerable to broad bias.

The architectural solution is the shift toward **Decentralized Machine Learning**, exemplified by **Federated Learning (FL)**. FL is a distributed training process that allows multiple client devices to collaborate on training AI models while keeping sensitive raw data localized. FL fundamentally changes the privacy landscape by eliminating the need for raw data transmission, thereby reducing exposure risks and regulatory compliance challenges.

This decentralization goes beyond privacy; it addresses bias by promoting justice, accountability, and transparency. By actively discovering and correcting biases across diverse, local datasets, decentralized AI networks can develop more inclusive and reliable models. Moreover, FL serves as a structural defense against algorithmic authoritarianism. By preventing the mass centralization of raw data, FL makes it significantly more challenging for platforms or regimes to repurpose algorithms for mass surveillance, predictive policing, or the systemic silencing of dissent, thereby protecting activists who rely on digital networks.

## B. Systemic Risk Governance and Continuous Auditing

The magnitude of algorithmic collapse, potentially leading to social violence and democratic erosion, mandates that algorithmic governance treat it as a **systemic risk**. Regulatory frameworks, such as the EU AI Act, classify AI based on risk, imposing stringent obligations on systems that present a systemic risk (General Purpose AI, GPAI, models).

For providers of GPAI models with systemic risk, the regulation mandates rigorous compliance, including: conducting model evaluations, tracking and reporting serious incidents, and performing and documenting **adversarial testing** to identify and mitigate systemic risk. This regulatory requirement provides the necessary impetus for developing robust MARL solutions (like those utilizing \text{MIR}^3) capable of anticipating worst-case adversarial behavior in complex social simulations.

Achieving continuous societal alignment requires institutionalized auditing frameworks. Frameworks like SMACTR mandate an end-to-end internal auditing process that emphasizes detailed documentation of algorithm design and training data, standardized impact assessment protocols, and systematic technical testing procedures throughout the AI development lifecycle. This ensures that ethical considerations and risk mitigation are embedded proactively, rather than addressed reactively after deployment failure.

## C. Designing for User Agency and Veto Power

User experience data consistently shows that people are fatigued by and unhappy with the overrepresentation of extreme or controversial content. This fatigue, coupled with a sense of

powerlessness, leads to algorithmic cynicism. Restoring user agency is essential for mitigating this cynicism and improving user experience.

While simple chronological feeds are proposed as an alternative, they represent a false choice by eliminating the benefits of personalization. The technical design must offer sophisticated controls. One mechanism is providing users with explicit **veto power** or an "advanced search option." If a service user is dissatisfied with an algorithm's output, they can refer the decision to a more detailed assessment level. This detailed assessment would trigger a more complex, personalized scrutiny—one governed by the newly defined APLF and CML models—allowing the user to explicitly override the engagement maximization drive and prioritize exploration or fine-grained diversity. This institutionalizes continuous user feedback and alignment within the core decision process.

# VI. Conclusion and Proposed Testable Hypotheses

The phenomenon of algorithmic collapse into polarization and extremism is an engineering problem caused by misaligned incentives rooted in extrinsic, short-term reward maximization. The solution requires a fundamental architectural transformation, shifting the optimization function toward intrinsic rewards (novelty) and causal objectives (depolarization and civility). This framework necessitates the integrated deployment of Exploration-Based Learning (EBL), Creation-Based Learning (CBL), and Causal Machine Learning (CML) within robust, decentralized architectures. By treating societal alignment as a technical requirement, we force ML systems to optimize for democratic resilience rather than simply maximizing commercial engagement.

## A. Synthesis of the Resilience Framework

The architecture for algorithmic resilience is defined by four interdependent pillars:
1. **Objective Function Redefinition:** The incorporation of CML models (like \partialCCF) and Anti-Polarization Loss Functions (APLF) to model the true, causal impact of recommendations and continuously guide the system toward constructive conflict.
2. **Exploration Engine (EBL):** Utilizing CDE and Novelty Search, regulated by $\text{MIR}^3$, to ensure the system is continuously rewarded for proactively seeking and promoting diverse information and identifying informational gaps.
3. **Intervention Engine (CBL):** Deploying C-ACS models, regulated by diversity-sensitive priors, to generate high-quality, civil, bridging narratives necessary to fill the gaps identified by EBL.
4. **Architectural Backbone:** Implementing Decentralized Architectures (FL) and ensuring strict adherence to Counterfactual Fairness and systemic adversarial testing mandates to structurally prevent bias amplification and guard against authoritarian capture.

## B. Concrete Hypotheses for Scientific Testing

The following hypotheses offer concrete, scientifically testable avenues for validating the efficacy of this new algorithmic paradigm. These can be rigorously tested using agent-based models on synthetic and empirical network topologies (e.g., modeling Facebook or Twitter dynamics ) to assess efficacy and potential unintended consequences, adhering to the principle that all scientific proposals must be testable.

Table 2: Testable Hypotheses for Next-Generation Algorithmic Systems

| Hypothesis ID | Mechanism Tested | Testable Outcome Metric | Expected Result (If Verified) |
|---|---|---|---|
| H1 (EBL) | Implementation of Curiosity-Driven Exploration (CDE) in Recommender Systems. | Increase in user consumption entropy, measured by Novelty Score and Intra-List Similarity of consumed items. | Reduced formation rate of polarized clusters and increased dissolution rate of clusters on synthetic social networks. |
| H2 (CBL) | Deployment of Civil Counter-Attitudinal Content Synthesis (C-ACS) via GenAI, constrained by APLF civility metrics. | Change in user "feeling thermometer" metric (affective polarization) toward out-group members and increase in constructive conflict (civility metric adherence). | Algorithmic fostering of constructive conflict, demonstrating C-ACS avoids negative spill-over effects observed in general AI intervention. |
| H3 (Causal) | Integration of Dynamic Causal Collaborative Filtering ($\partial$CCF) using counterfactual reasoning into ranking models. | Reduction in observed echo chamber persistence across ideologically mixed user cohorts, measured by network homophily/triadic closure rates. | Mitigation of echo chambers by modeling and interrupting the causal feedback loop without sacrificing baseline recommendation performance. |

## C. Final Conclusion: The Path to Algorithmic Resilience

The current crisis of algorithmic collapse demands an immediate shift in focus from engineering efficiency to engineering resilience. By adopting the EBL/CBL/CML framework, which mandates diversity, causality, and robustness as primary design objectives, it is technologically feasible to build socio-technical systems that actively cultivate cognitive friction and civil discourse, transforming algorithms from vectors of societal fragmentation into pillars of democratic stability.

**Works cited**

1. What is AI model collapse? Definition, examples, and challenges - Infobip, https://www.infobip.com/glossary/model-collapse 2. What is a Generative Model? - DataCamp, https://www.datacamp.com/blog/what-is-a-generative-model 3. Optimizing myopia prediction in children and adolescents using machine learning: a multi-factorial risk assessment model - PubMed, https://pubmed.ncbi.nlm.nih.gov/41322210/ 4. Better Feeds: Algorithms That Put People First - Knight-Georgetown Institute, https://kgi.georgetown.edu/research-and-commentary/better-feeds/ 5. Can Algorithmic Recommendation Systems Be Good For Democracy? (Yes! & Chronological Feeds May Be Bad) | The Belfer Center for Science and International Affairs, https://www.belfercenter.org/publication/can-algorithmic-recommendation-systems-be-good-democracy-yes-chronological-feeds-may-be 6. Social media algorithms exploit how we learn from our peers - Northwestern Now,

https://news.northwestern.edu/stories/2023/08/social-media-algorithms-exploit-how-humans-learn-from-their-peers 7. Algorithmic influence and media legitimacy: a systematic review of social media's impact on news production - Frontiers, https://www.frontiersin.org/journals/communication/articles/10.3389/fcomm.2025.1667471/full 8. Filter bubble - Wikipedia, https://en.wikipedia.org/wiki/Filter_bubble 9. The Role of Bubbles and Echo Chambers in Polarization - IEA-USP, https://www.iea.usp.br/midiateca/video/videos-2024/the-role-of-bubbles-and-echo-chambers-in-polarization 10. The Role of Bubbles and Echo Chambers in Polarization - YouTube, https://www.youtube.com/watch?v=GBE_n1_HdTQ 11. Echo chambers, filter bubbles, and polarisation: a literature review ..., https://reutersinstitute.politics.ox.ac.uk/echo-chambers-filter-bubbles-and-polarisation-literature-review 12. Through the Newsfeed Glass: Rethinking Filter Bubbles and Echo Chambers - PMC - PubMed Central, https://pmc.ncbi.nlm.nih.gov/articles/PMC8923337/ 13. The Role of the Internet and Social Media on Radicalization - Office of Justice Programs, https://www.ojp.gov/pdffiles1/nij/305797.pdf 14. When knowing more means doing less: Algorithmic knowledge and digital (dis)engagement among young adults, https://misinforeview.hks.harvard.edu/article/when-knowing-more-means-doing-less-algorithmic-knowledge-and-digital-disengagement-among-young-adults/ 15. The impact of intrinsic rewards on exploration in Reinforcement Learning - arXiv, https://arxiv.org/html/2501.11533v1 16. CDE: Curiosity-Driven Exploration for Efficient Reinforcement Learning in Large Language Models - arXiv, https://arxiv.org/html/2509.09675v1 17. Leveraging AI for Enhanced Patentability and Novelty Search Techniques - Patentskart, https://patentskart.com/leveraging-ai-for-enhanced-patentability-and-novelty-search-techniques/ 18. Not your average RecSys metrics Part 2: Novelty | Shaped Blog, https://www.shaped.ai/blog/not-your-average-recsys-metrics-part-2-novelty 19. How do recommender systems handle diversity and novelty? - Milvus, https://milvus.io/ai-quick-reference/how-do-recommender-systems-handle-diversity-and-novelty 20. Recommender system - Wikipedia, https://en.wikipedia.org/wiki/Recommender_system 21. Solving the apparent diversity-accuracy dilemma of recommender systems - PMC - NIH, https://pmc.ncbi.nlm.nih.gov/articles/PMC2842039/ 22. Diversity in Machine Learning - arXiv, https://arxiv.org/pdf/1807.01477 23. Data Diversity as Implicit Regularization: How Does Diversity Shape the Weight Space of Deep Neural Networks? - arXiv, https://arxiv.org/html/2410.14602v2 24. Multi-Agent Reinforcement Learning: A Review of Challenges and Applications - MDPI, https://www.mdpi.com/2076-3417/11/11/4948 25. Robust Multi-Agent Reinforcement Learning by Mutual Information Regularization, https://ieeexplore.ieee.org/document/11074764/ 26. Robust Multi-Agent Reinforcement Learning with Model Uncertainty, https://proceedings.neurips.cc/paper/2020/hash/774412967f19ea61d448977ad9749078-Abstract.html 27. [2508.20818] cMALC-D: Contextual Multi-Agent LLM-Guided Curriculum Learning with Diversity-Based Context Blending - arXiv, https://arxiv.org/abs/2508.20818 28. What is a Generative Model? | IBM, https://www.ibm.com/think/topics/generative-model 29. Generative artificial intelligence - Wikipedia, https://en.wikipedia.org/wiki/Generative_artificial_intelligence 30. (PDF) Artificial Intelligence and Creativity - ResearchGate, https://www.researchgate.net/publication/389520762_Artificial_Intelligence_and_Creativity 31. [1901.09024] Diversity-Sensitive Conditional Generative Adversarial Networks - arXiv, https://arxiv.org/abs/1901.09024 32. View of Designing recommender systems to depolarize - First Monday, https://firstmonday.org/ojs/index.php/fm/article/view/12604/10636 33. Designing Recommender Systems to Depolarize 16 - ArXiv, https://arxiv.org/abs/2107.04953 34.

[2506.14295] The Impact of Generative AI on Social Media: An Experimental Study - arXiv, https://arxiv.org/abs/2506.14295 35. High-level summary of the AI Act | EU Artificial Intelligence Act, https://artificialintelligenceact.eu/high-level-summary/ 36. The Repressive Power of Artificial Intelligence - Freedom House, https://freedomhouse.org/report/freedom-net/2023/repressive-power-artificial-intelligence 37. Dynamic Causal Collaborative Filtering - arXiv, https://arxiv.org/abs/2208.11094 38. Causal Inference for Recommendation: Foundations, Methods and Applications - arXiv, https://arxiv.org/pdf/2301.04016 39. (PDF) Causal Machine Learning in Social Impact Assessment - ResearchGate, https://www.researchgate.net/publication/374760761_Causal_Machine_Learning_in_Social_Impact_Assessment 40. Causal Inference and Machine Learning for Social Science, https://www.nationalacademies.org/event/docs/D35747FBFF9F5EF3273D787093CFA3BA0C0A10E5E51B?noSaveAs=1 41. Counterfactual fairness - The Alan Turing Institute, https://www.turing.ac.uk/research/research-projects/counterfactual-fairness 42. Ethical algorithm design should guide technology regulation - Brookings Institution, https://www.brookings.edu/articles/ethical-algorithm-design-should-guide-technology-regulation/ 43. A Comprehensive Review on Understanding the Decentralized and Collaborative Approach in Machine Learning - arXiv, https://arxiv.org/html/2503.09833v1 44. 15 Counterfactual Explanations – Interpretable Machine Learning - Christoph Molnar, https://christophm.github.io/interpretable-ml-book/counterfactual.html 45. The Pros and Cons of Social Media Algorithms - Bipartisan Policy Center, https://bipartisanpolicy.org/wp-content/uploads/2023/10/BPC_Tech-Algorithm-Tradeoffs_R01.pdf 46. Federated Learning: Training AI Models Without Centralizing Data - Coderio, https://www.coderio.com/innovation/federated-learning-training-ai-models/ 47. AI Bias Mitigation: Addressing and Reducing Bias in Decentralized AI Models | by DcentAI, https://medium.com/coinmonks/ai-bias-mitigation-addressing-and-reducing-bias-in-decentralized-ai-models-706b86fa893e 48. Using AI as a weapon of repression and its impact on human rights - European Parliament, https://www.europarl.europa.eu/RegData/etudes/IDAN/2024/754450/EXPO_IDA(2024)754450_EN.pdf 49. Silencing Dissent: How the Tech Backlash Threatens Transitional Justice, https://www.globalgovernance.eu/publications/silencing-dissent-how-the-tech-backlash-threatens-transitional-justice 50. Key Issue 3: Risk-Based Approach - EU AI Act, https://www.euaiact.com/key-issue/3 51. Closing the AI accountability gap: defining an end-to-end framework for internal algorithmic auditing - Digital Government Hub, https://digitalgovernmenthub.org/library/closing-the-ai-accountability-gap-defining-an-end-to-end-framework-for-internal-algorithmic-auditing/ 52. Essential Framework For Algorithmic Transparency Audits - Troy Lendman, https://troylendman.com/essential-framework-for-algorithmic-transparency-audits/ 53. Algorithms, Artificial Intelligence, and the Law - Judicature - Duke University, https://judicature.duke.edu/articles/algorithms-artificial-intelligence-and-the-law/ 54. Guided rewiring of social networks reduces polarization and accelerates collective action, https://arxiv.org/html/2309.12141v2