

Date	Topic	Title	Link
Jan 16	Fairness	Inherent Trade-Offs in the Fair Determination of Risk Scores	<a href="https://arxiv.org/abs/1609.05807">https://arxiv.org/abs/1609.05807</a>
Jan 16	Fairness	The Possibility of Fairness: Revisiting the Impossibility Theorem in Practice	<a href="https://arxiv.org/abs/2302.06347">https://arxiv.org/abs/2302.06347</a>
Jan 23	Ethics	On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?	<a href="https://dl.acm.org/doi/10.1145/3442188.3445922">https://dl.acm.org/doi/10.1145/3442188.3445922</a>
Jan 23	Ethics	The Values Encoded in Machine Learning Research	<a href="https://arxiv.org/abs/2106.15590">https://arxiv.org/abs/2106.15590</a>
Jan 23	Fairness	Stereotyping Norwegian Salmon: An Inventory of Pitfalls in Fairness Benchmark Datasets:	<a href="https://www.microsoft.com/en-us/research/uploads/prod/2021/06/The_Salmon_paper.pdf">https://www.microsoft.com/en-us/research/uploads/prod/2021/06/The_Salmon_paper.pdf</a>
Jan 30	Fairness	Image Cropping on Twitter: Fairness Metrics, their Limitations, and the Importance of Representation, Design, and Agency	<a href="https://arxiv.org/abs/2105.08667">https://arxiv.org/abs/2105.08667</a>
Jan 30	Accountability	Closing the AI Accountability Gap: Defining an End-to-End Framework for Internal Algorithmic Auditing	<a href="https://arxiv.org/abs/2001.00973">https://arxiv.org/abs/2001.00973</a>
Jan 30	Fairness	Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup Fairness	<a href="https://arxiv.org/abs/1711.05144">https://arxiv.org/abs/1711.05144</a>
Feb 6	Fairness	Feedback Loop and Bias Amplification in Recommender Systems	<a href="https://dl.acm.org/doi/pdf/10.1145/3340531.3412152">https://dl.acm.org/doi/pdf/10.1145/3340531.3412152</a>

Feb 6	Fairness	Two-Sided Fairness for Repeated Matchings in Two-Sided Markets: A Case Study of a Ride-Hailing Platform	<a href="https://dl.acm.org/doi/10.1145/3292500.3330793">https://dl.acm.org/doi/10.1145/3292500.3330793</a>
Feb 6	Fairness	Fairness is not static: deeper understanding of long term fairness via simulation studies	<a href="https://dl.acm.org/doi/abs/10.1145/3351095.3372878">https://dl.acm.org/doi/abs/10.1145/3351095.3372878</a>
Feb 20	Privacy	<u>Membership Inference Attacks against Machine Learning Models</u>	<a href="https://arxiv.org/abs/1610.05820">https://arxiv.org/abs/1610.05820</a>
Feb 20	Privacy	The Privacy Onion Effect: Memorization is Relative	<a href="https://arxiv.org/abs/2206.10469">https://arxiv.org/abs/2206.10469</a>
Feb 20	Privacy	Machine Unlearning	<a href="https://arxiv.org/abs/1912.03817">https://arxiv.org/abs/1912.03817</a>
March 19	Robustness	On Evaluating Adversarial Robustness	<a href="https://arxiv.org/abs/1902.06705">https://arxiv.org/abs/1902.06705</a>
March 19	Robustness	Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback	<a href="https://arxiv.org/abs/2307.15217">https://arxiv.org/abs/2307.15217</a>
March 19	Safety	Jailbroken: How Does LLM Safety Training Fail?	<a href="https://arxiv.org/abs/2307.02483">https://arxiv.org/abs/2307.02483</a>
March 26	Explainability	on the privacy risks of model explanations	<a href="https://arxiv.org/abs/1907.00164">https://arxiv.org/abs/1907.00164</a>
March 26	Explainability	Counterfactual Explanations and Algorithmic Recourses for Machine Learning: A Review	<a href="https://arxiv.org/abs/2010.10596">https://arxiv.org/abs/2010.10596</a>
March 26	Transparency/ Safety	Alignment of Language Agents	<a href="https://arxiv.org/abs/2103.14659">https://arxiv.org/abs/2103.14659</a>