

**“Everyone wants to do the model work, not the data work”**  
**Data Cascades in High-Stakes AI**  
**Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong,**  
**Praveen Paritosh, Lora Aroyo**

**Reviewed by: Rishabh Devgon**

**Critical Review:**

Quality of data dictates the fairness, performance, robustness, safety, and scalability of AI systems [1]. The authors discuss how the quality of data is undervalued and less incentivised when designing these systems. This can create downstream issues, particularly in certain high stakes domains. These issues can have significant humanitarian impacts with the ability to ostracize specific communities and contexts, can lack consistency in data and create problems for collaboration among disciplines. The prevalence of severe data cascades points to a larger systemic problem in terms of the practices, methodologies and incentives in the field of AI. The paper highlights a gap in the understanding of the human impacts of AI models. The paper further reviews related work and expands on data in HCI, politics of data, data quality interventions and machine learning in production. The authors do an incredible job of finding out problems associated with data cascading. The paper continually stresses on how data cascades are often avoidable. The findings of the paper highlight the various insights into data cascades in their overview, a broader landscape and triggers and practices. The discussion takes several sub-themes such as the goodness of fit to the goodness of data, incentives for data excellence, real-world literacy in AI education and better visibility in the AI data lifecycle and data equity in the Global South. The title corresponds to the findings and motivation of the paper and hence is justified.

The research method used by the authors is semi-structured interviews with 53 AI practitioners recruited through a combination of snowball and purposive sampling. The analysis was done through a multi-level coding structure. The authors have clearly mentioned their limitations in the study, outlined their process of anonymity and confidentiality, and I found the research airtight. However, since the investigation is international, European, and other Asian countries could have been included in the sampling. Another way to improve the study could have been to iterate over the results found in previous interview studies and create a positive feedback loop for probes that could provide more in-depth insights. The downside that I see to this analysis is that adopting an international approach may subdue several socio-cultural issues that may pertain to AI practices in these diverse regions. The researchers have acknowledged the fact that the interviewees have experience in inculcating positive habits that are required for high stakes AI usage. A way to make the study more comprehensive could be to have a survey to get a better picture of high stake AI practices in general which would reduce the bias.

Although the focus on data practices is not new, what makes the paper novel is that the paper focuses explicitly on data practices in high stake domains. The research tries to focus on

communities with higher financial constraints and moves towards an international analysis with investigating AI practitioners in India, USA and Eastern and Western Africa.

The paper could further benefit from explorations in HCI to design systems that can facilitate the processes of data collection and be intuitive such that human error is minimal. The authors do mention the importance of evaluating data in terms of social and physical geography but the authors do not particularly talk about the personal biases of the interpreter of data. To make the study stronger, perhaps the authors could have further investigated the relationships between practitioners and other involved stakeholders. Laying stress on communication between these channels could benefit the research. Possibly awareness and further discussion about HCAI among other actors could improve the entire AI ecosystem. Overall, the paper maintains a right balance between the qualitative insights gained, human experiences and computational and resource constraints.

### **References:**

1. Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora Aroyo. 2021. Everyone wants to do the model work, not the data work”: Data Cascades in High-Stakes AI.