



# Tackling Socioeconomic Bias in Machine Learning

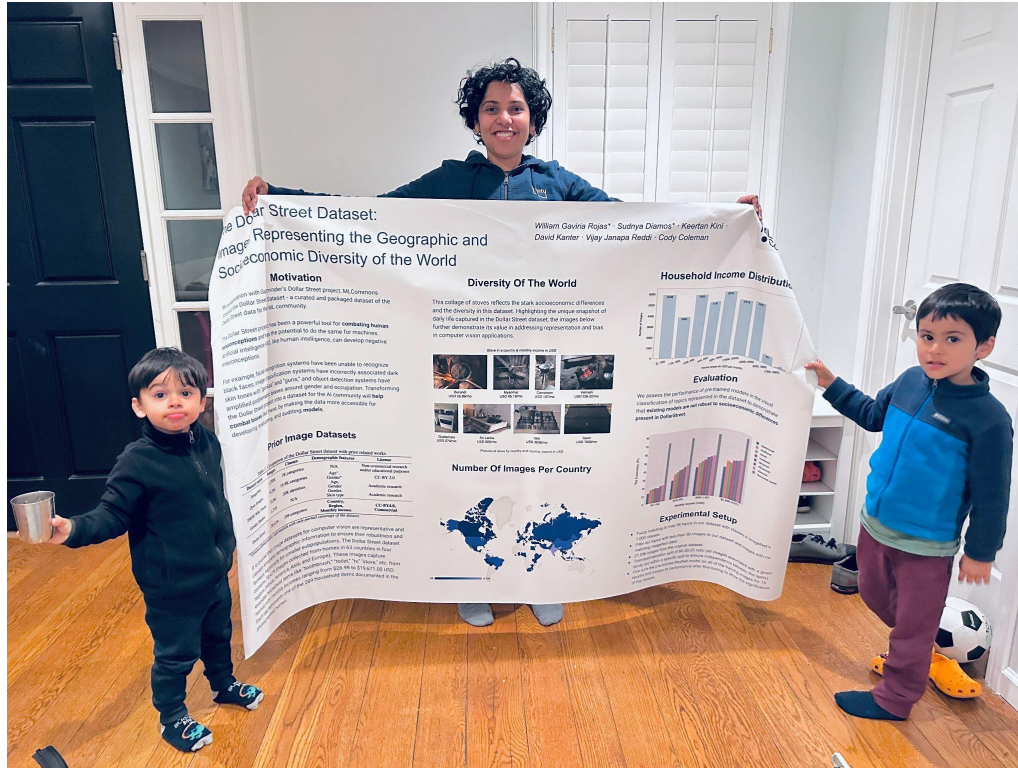
Cody Coleman, PhD  
CEO & Co-Founder | Coactive AI

A collaboration between researchers from



# The Team

This work was an amazing collaboration across, academia, and industry. Thank you!



William A Gaviria Rojas, Sudnya Damos, Keertan Ranjan Kini, David Kanter, Vijay Janapa Reddi, **Cody Coleman**

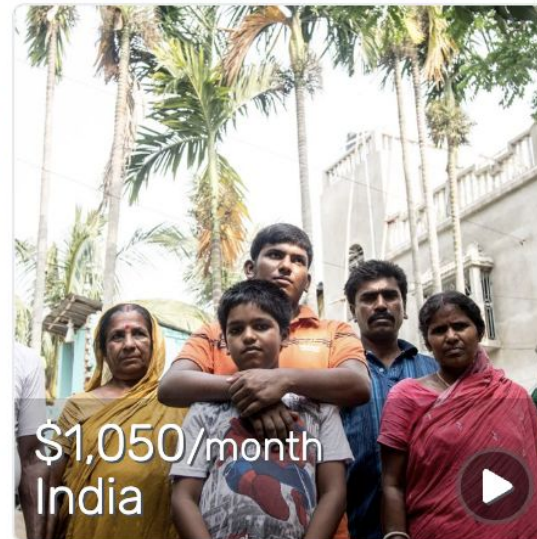
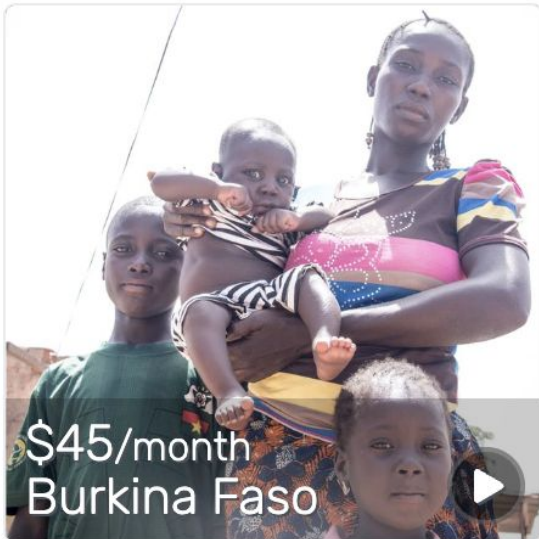
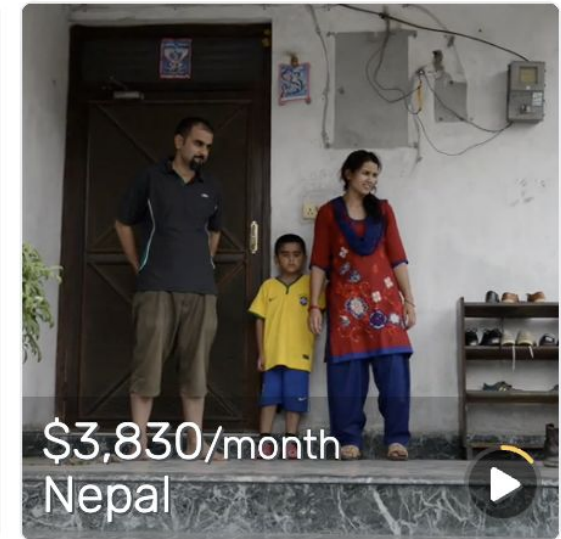
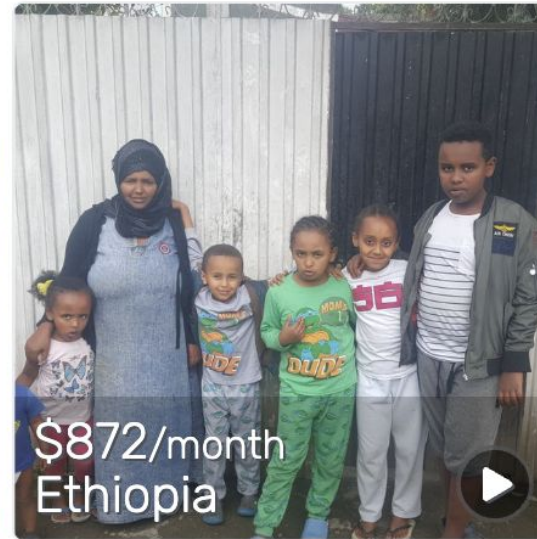
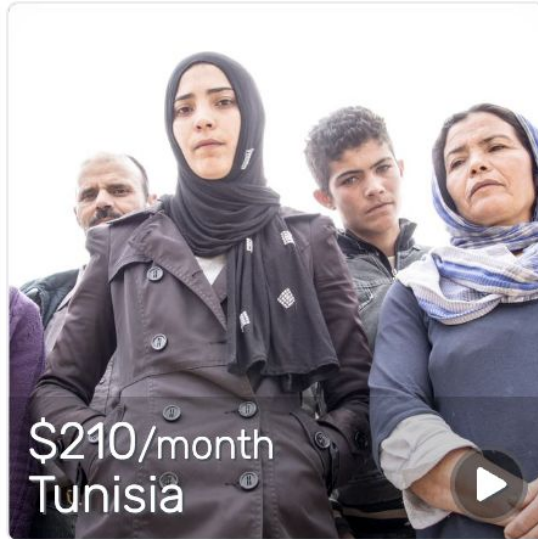


HARVARD  
UNIVERSITY



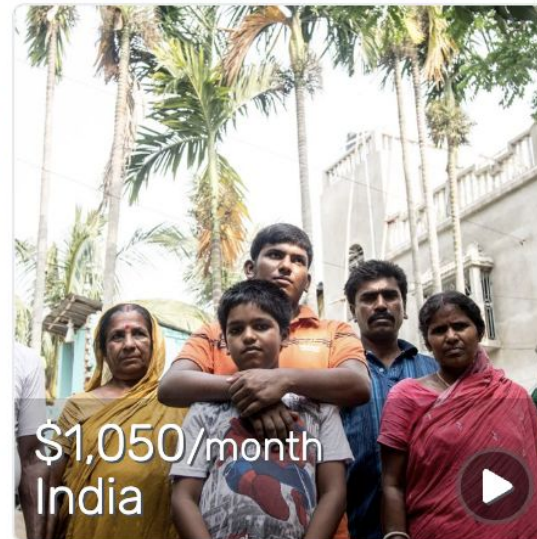
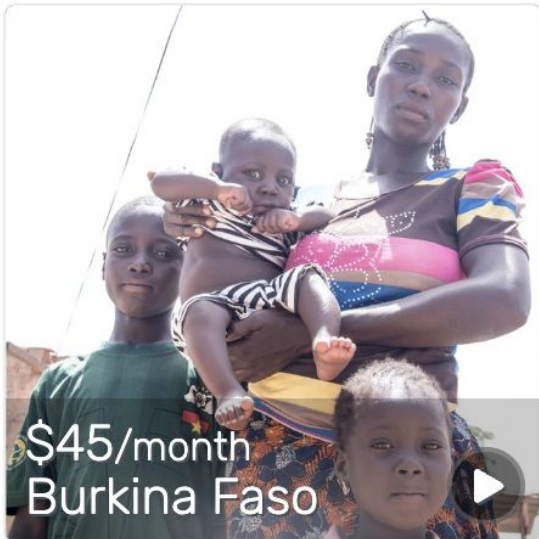
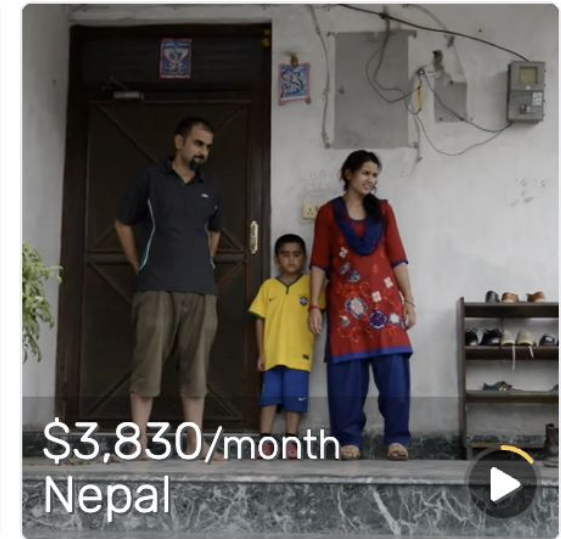
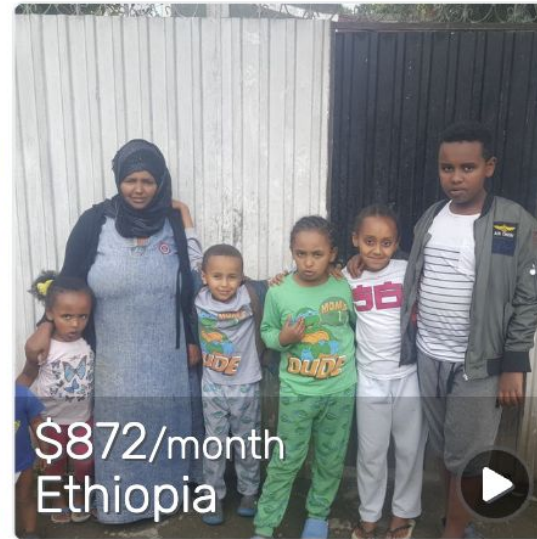
# The Dollar Street Project

Combating human bias through photos as data across incomes and borders



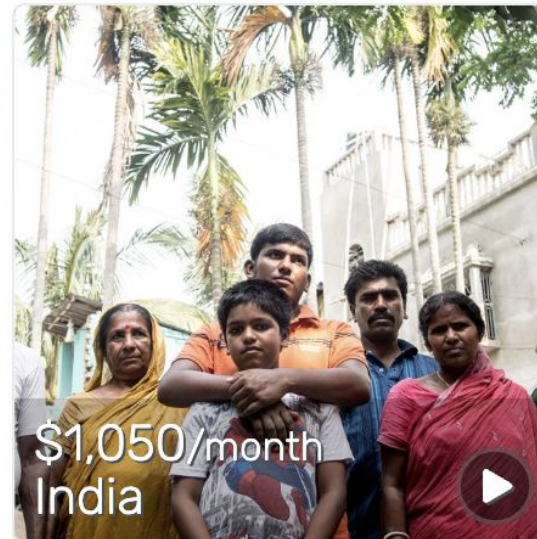
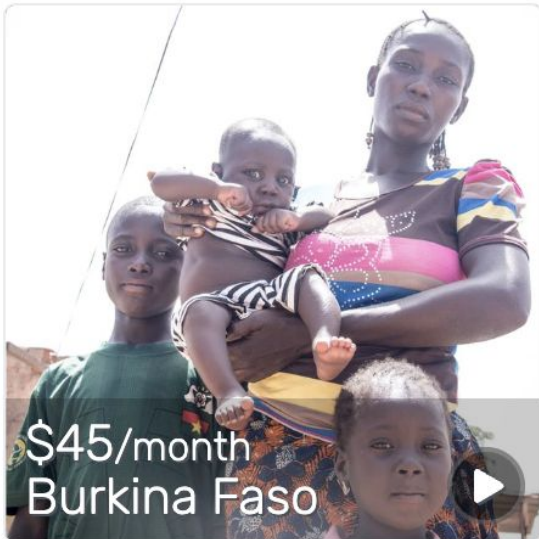
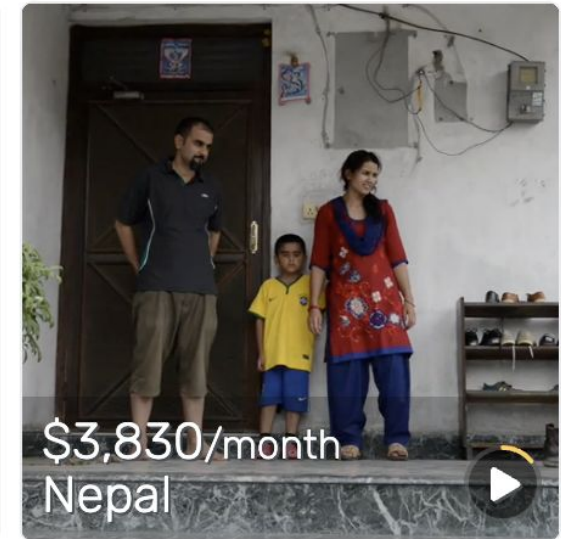
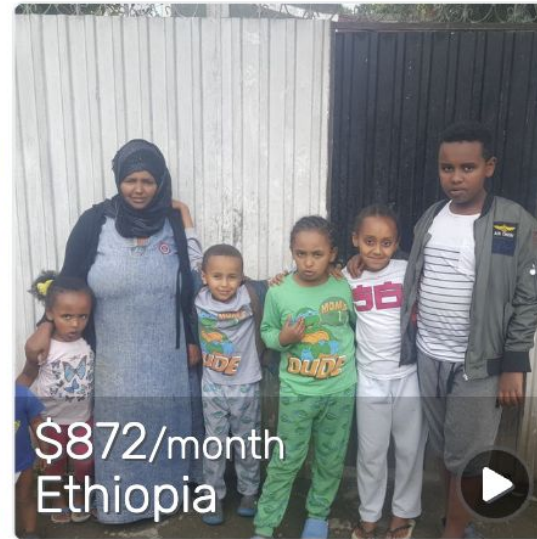
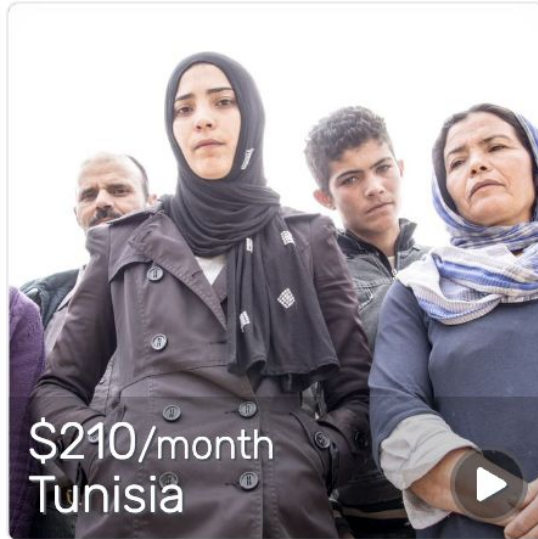
# The Dollar Street Project

Combating **human** bias through photos as data across incomes and borders

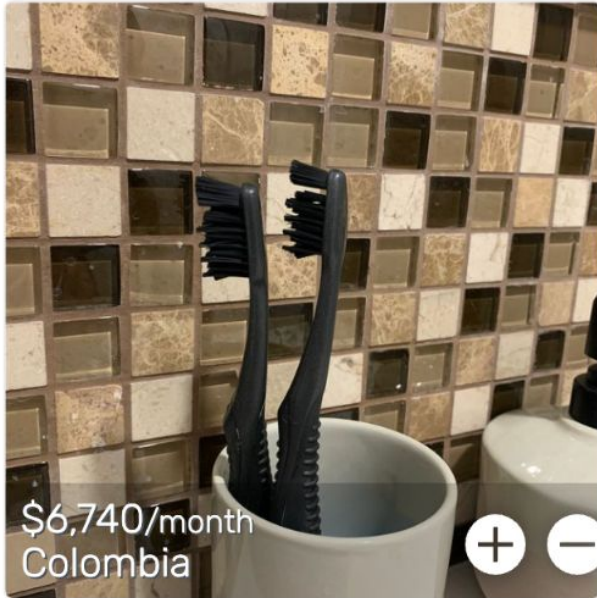
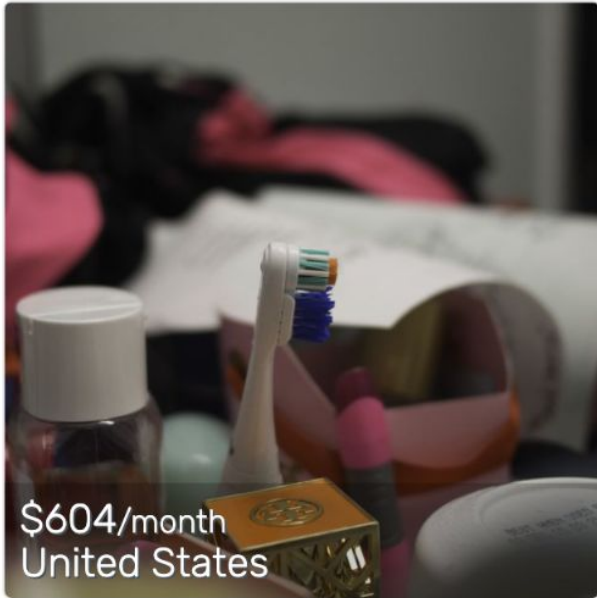
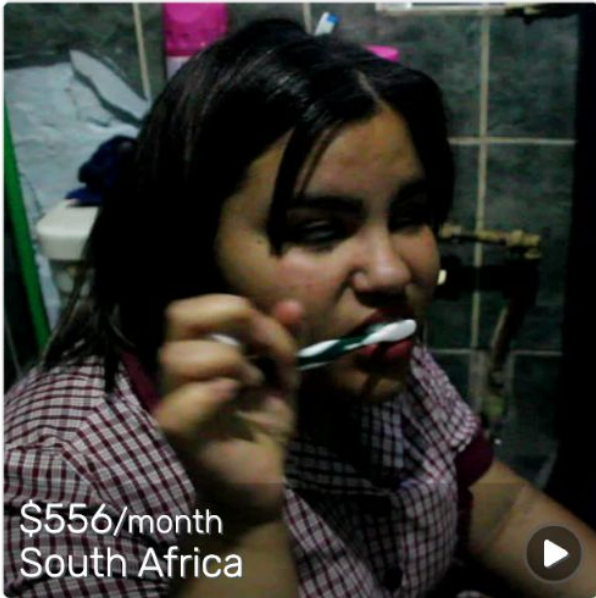
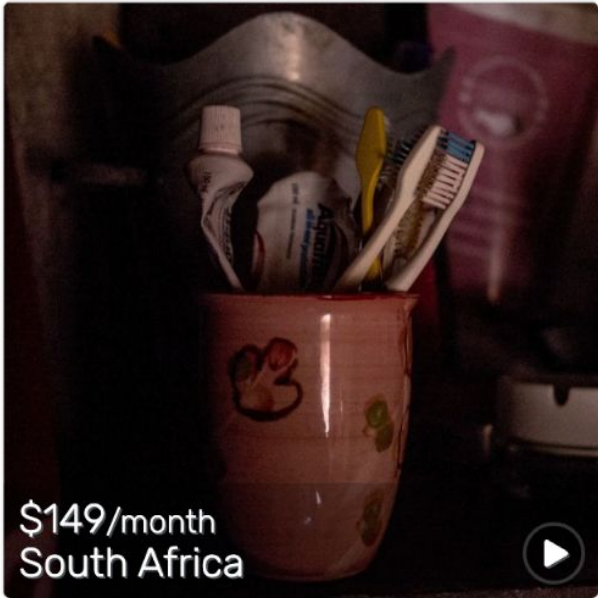
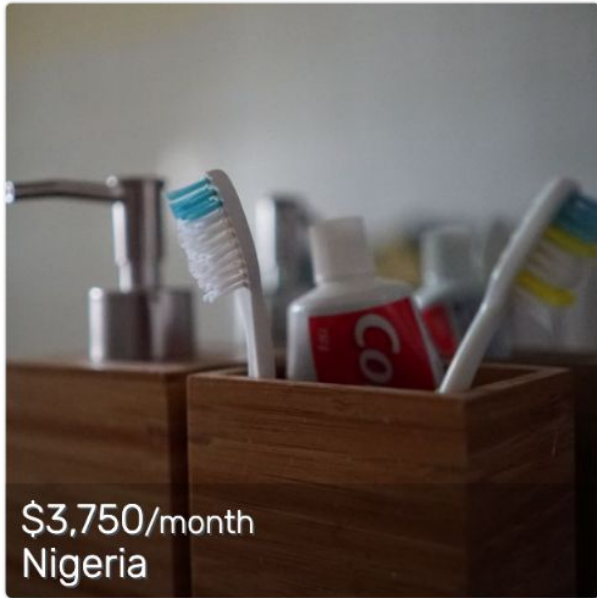
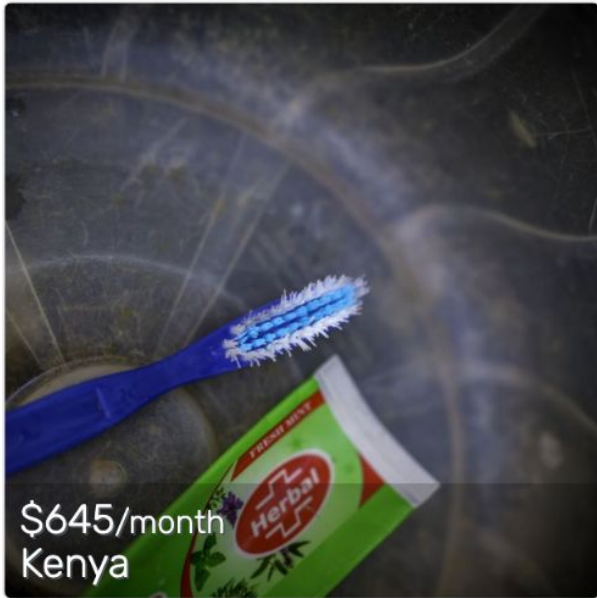
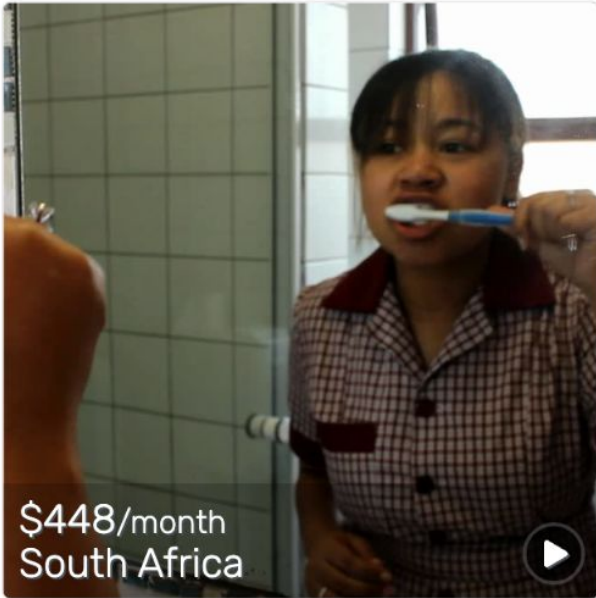
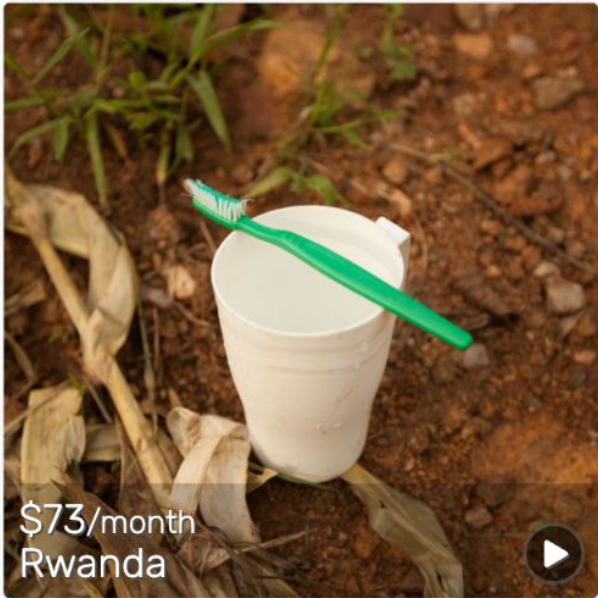


# The Dollar Street Project

Combating human bias through photos as data across incomes and borders

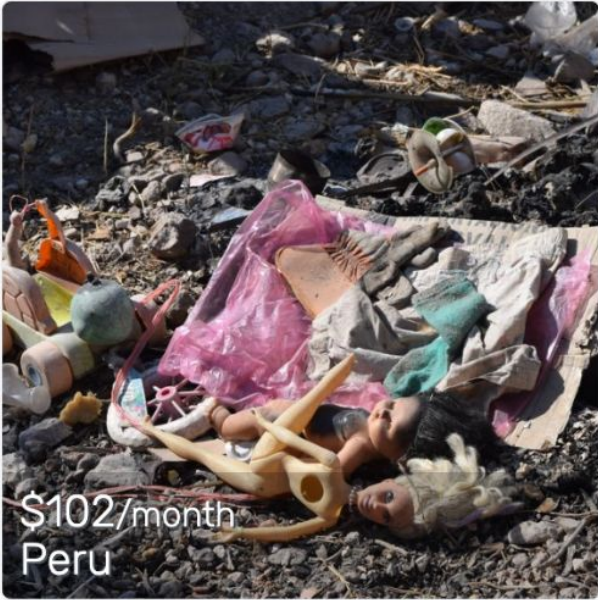












\$102/month  
Peru



\$430/month  
Serbia



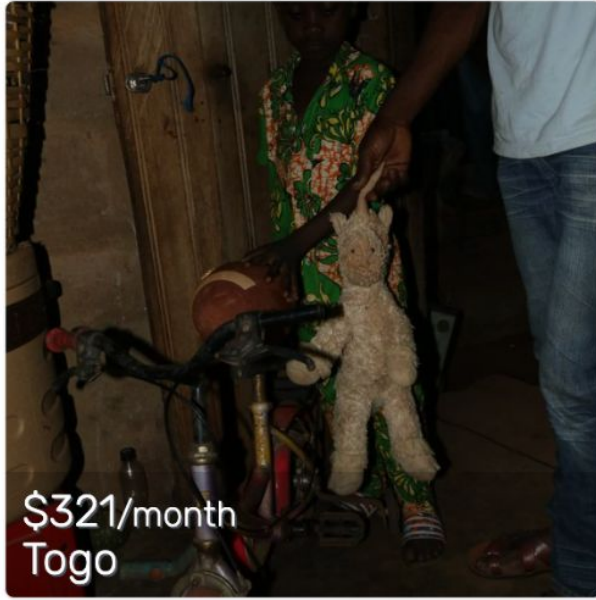
\$1,019/month  
Romania



\$4,553/month  
Spain



\$96/month  
Nepal



\$321/month  
Togo



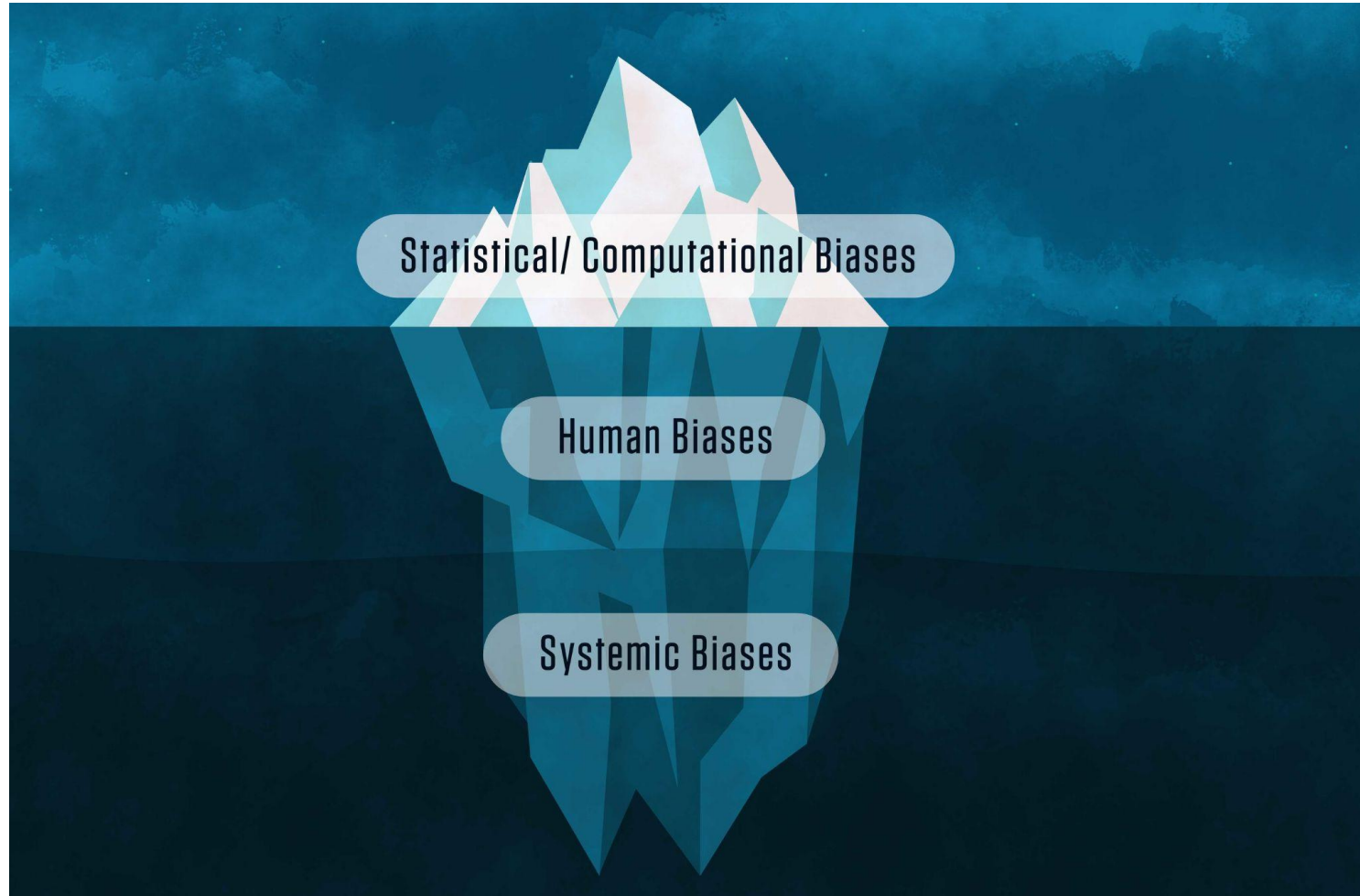
\$719/month  
Ethiopia



\$4,162/month  
Colombia

# Human and algorithmic bias are connected

Human biases can get absorbed in AI models and systems



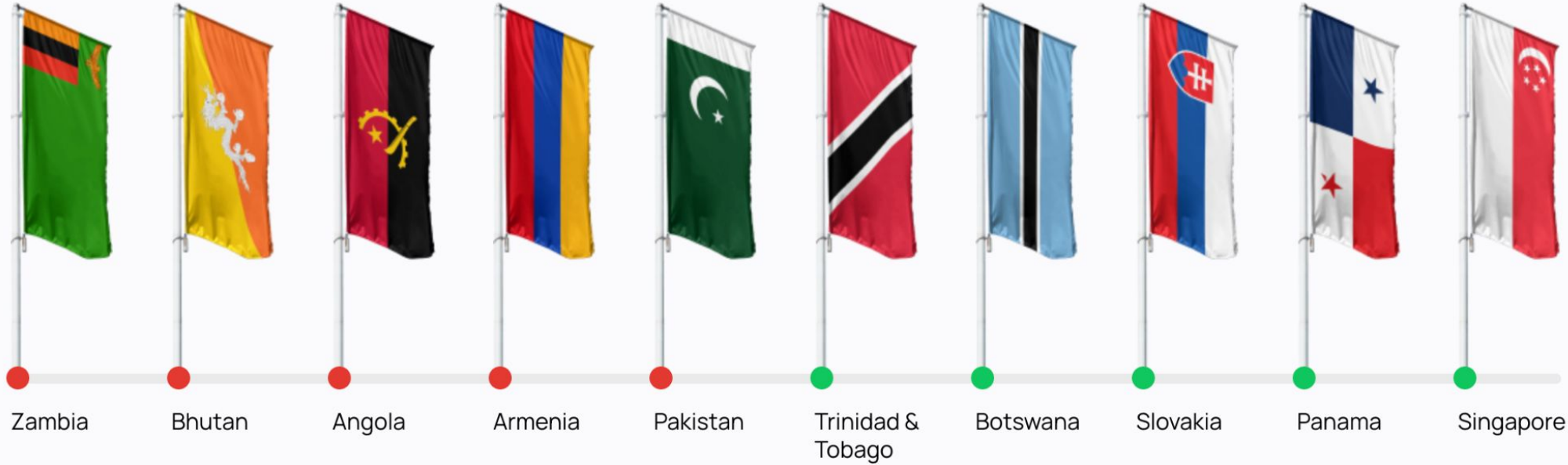
# AI affects all of us

Automation and technology proliferation concerns worldwide

## TOP 5 NATIONS MOST / LEAST AFFECTED BY AI AUTOMATION

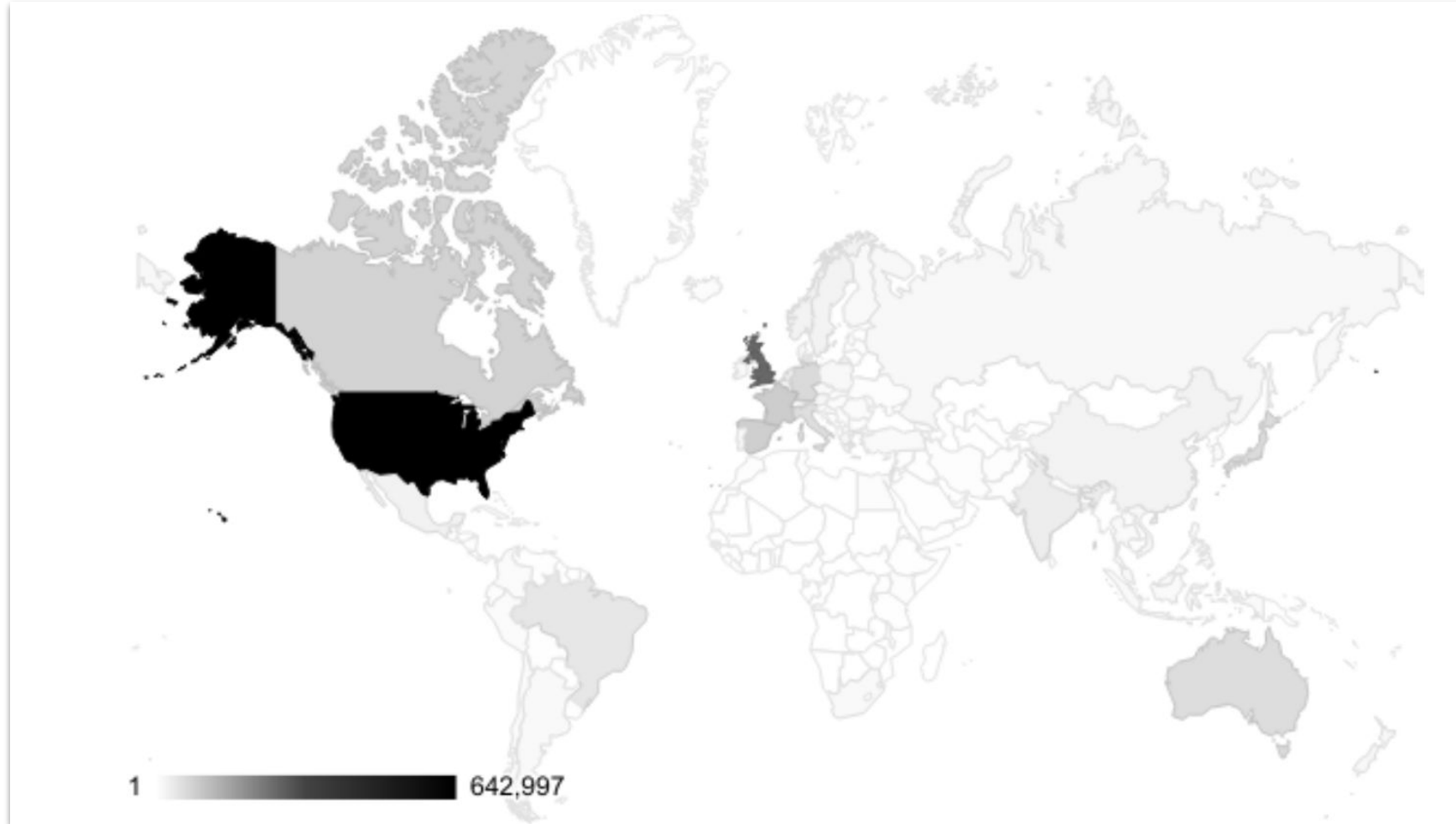
MOST →

← LEAST



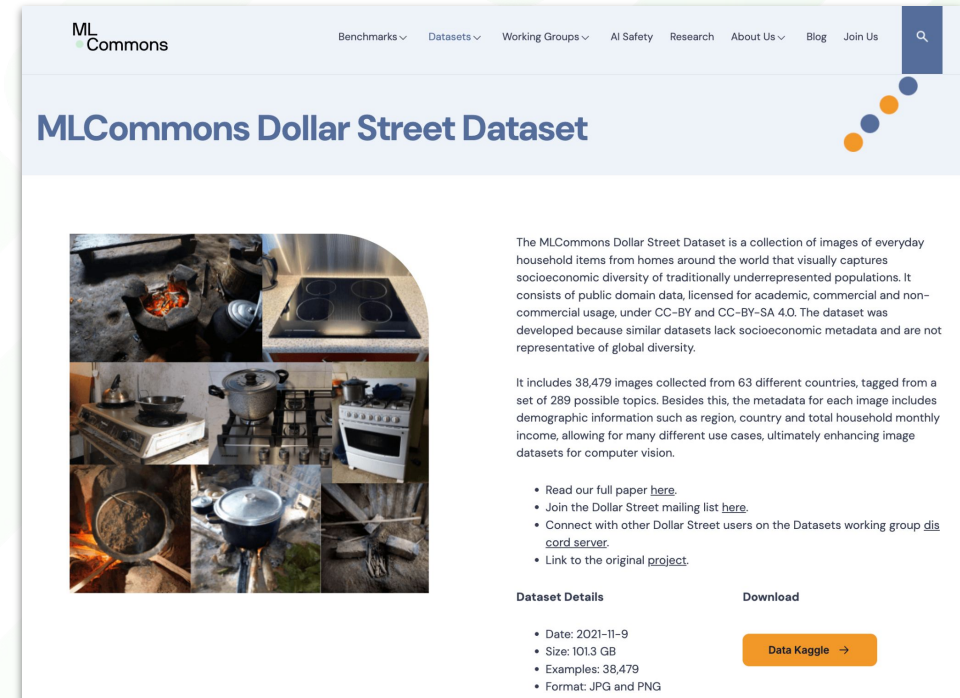
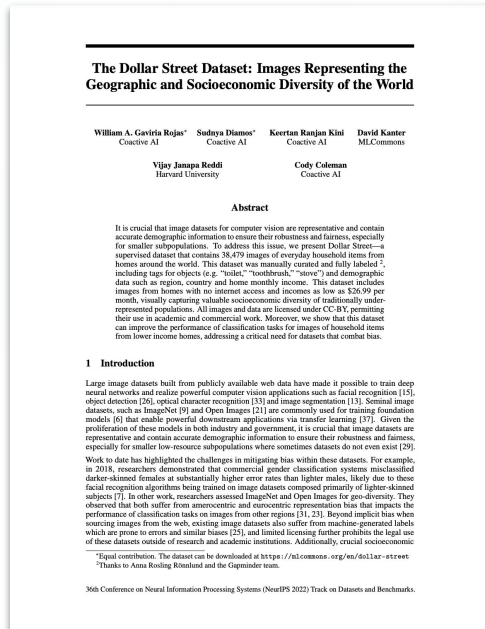
# Yet AI training data is not often representative

Training data sourced from the internet: 33% from US, 60% from six countries



# The Dollar Street dataset for AI training

Combating AI training bias through photos as data across incomes and borders



Published in top AI academic conference  
(NeurIPS 2022)

And made freely available to the world through  
MLCommons, with a CC-BY-4.0 license

# The Dollar Street captures unique geographic and socioeconomic information

38k

Images

289

Everyday household  
items

404

Homes

63

Countries

- Data was manually collected and verified
- Includes homes with no internet access
- Includes households with incomes as low as \$26.99 per month
- All data is licensed under CC-BY 4.0

# The Dollar Street dataset captures visual diversity



Burundi  
\$26.99 / month



Myanmar  
\$45.18 / month



Cameroon  
\$137.00 / month



Vietnam  
\$236.22 / month



Guatemala  
\$270.00 / month



Sri Lanka  
\$909.00 / month



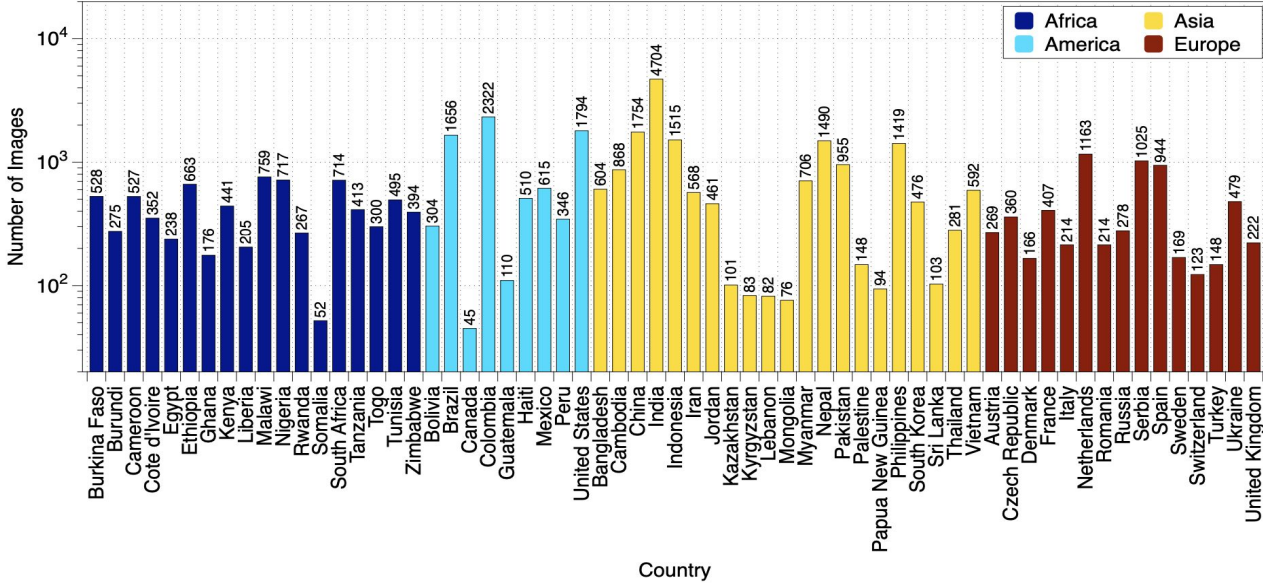
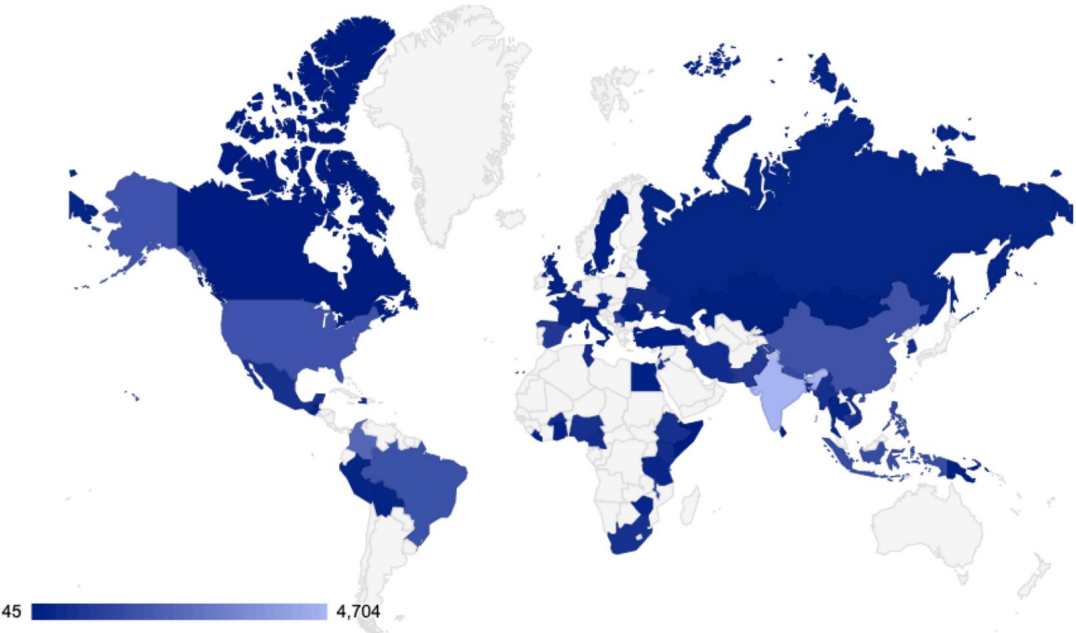
Italy  
\$3,636.00 / month



Denmark  
\$5,343.00 / month



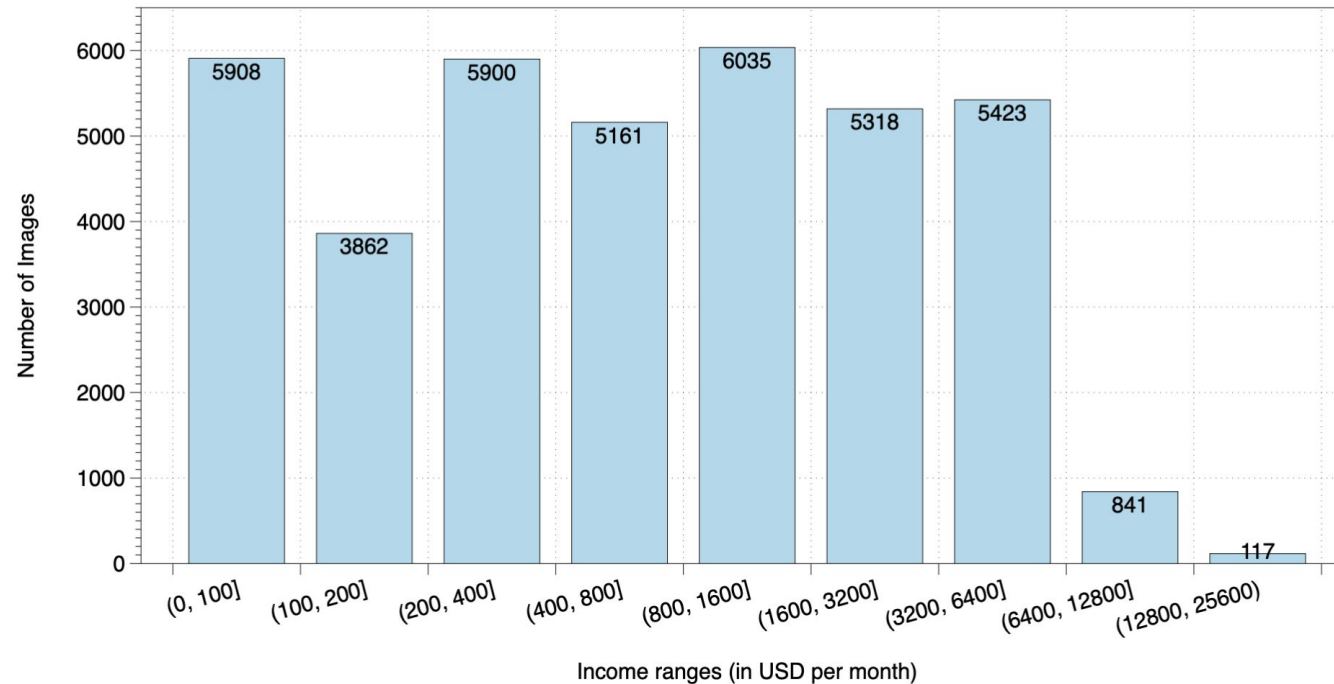
# The Dollar Street dataset captures geographical representation



Number of images per country



# The Dollar Street dataset captures socioeconomic representation



## Number of images by monthly income

The Dollar Street dataset contains images from homes with a wide range of monthly incomes (in USD).

At a glance: Min = \$26.99, Median = \$685.00, Max = \$19,671.00

# The results: equity in AI performance

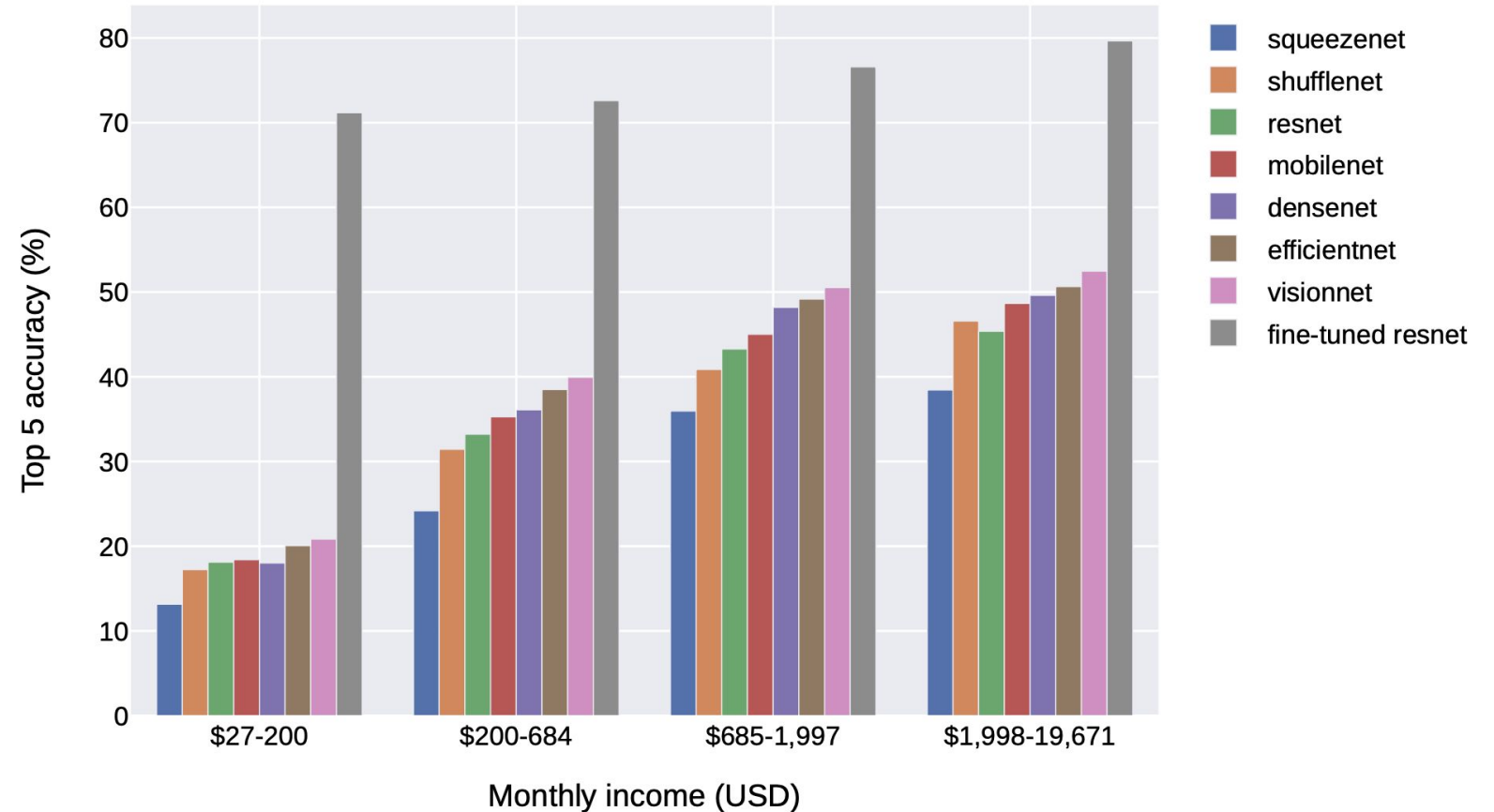
Visual classification accuracy improves by +30-50% across various income levels

## Key results:

Pre-trained models show **monotonically decreasing classification accuracy** with **decreasing income**.

Fine-tuning using the Dollar Street dataset → substantially performance improvement

The large improvement for the lowest-income quartile highlights the importance of the Dollar Street dataset.



# The results: inspiring others to tackle AI bias

## The **dollar street dataset**: Images representing the geographic and socioeconomic diversity of the world

[WAG Rojas](#), [S Damos](#), [KR Kini](#), [D Kanter](#)... - ... -sixth Conference on ..., 2022 - openreview.net

... **Dollar Street**—a supervised **dataset** that contains 38,479 images of everyday household items from homes around the world. This **dataset** ... This **dataset** includes images from homes with ...

☆ Save 📄 Cite **Cited by 54** 📄 Related articles All 2 versions 🔗

**50+ citations across research teams in academia, industry and startups**

## Facet: Fairness in computer vision evaluation benchmark

[L Gustafson](#), [C Rolland](#), [N Ravi](#)... - Proceedings of the ..., 2023 - openaccess.thecvf.com

Computer vision models have known performance disparities across attributes such as gender and skin tone. This means during tasks such as classification and detection, model ...

☆ Save 📄 Cite Cited by 17 Related articles All 6 versions 🔗

## A whac-a-mole dilemma: Shortcuts come in multiples where mitigating one amplifies others

[Z Li](#), [I Evtimov](#), [A Gordo](#), [C Hazirbas](#)... - Proceedings of the ..., 2023 - openaccess.thecvf.com

Abstract Machine learning models have been found to learn shortcuts---unintended decision rules that are unable to generalize---undermining models' reliability. Previous works address ...

☆ Save 📄 Cite Cited by 40 Related articles All 9 versions 🔗

## Survey of social bias in vision-language models

[N Lee](#), [Y Bang](#), [H Lovenia](#), [S Cahyawijaya](#)... - arXiv preprint arXiv ..., 2023 - arxiv.org

In recent years, the rapid advancement of machine learning (ML) models, particularly transformer-based pre-trained models, has revolutionized Natural Language Processing ...

☆ Save 📄 Cite Cited by 5 Related articles All 2 versions 🔗

## The casual conversations v2 dataset

[B Porgali](#), [V Albiero](#), [J Ryda](#)... - Proceedings of the ..., 2023 - openaccess.thecvf.com

This paper introduces a new large consent-driven dataset aimed at assisting in the evaluation of algorithmic bias and robustness of computer vision and audio speech models ...

☆ Save 📄 Cite Cited by 19 Related articles All 6 versions 🔗

## Representation in AI evaluations

[AS Bergman](#), [LA Hendricks](#), [M Rauh](#), [B Wu](#)... - Proceedings of the ..., 2023 - dl.acm.org

Calls for representation in artificial intelligence (AI) and machine learning (ML) are widespread, with "representation" or "representativeness" generally understood to be both ...

☆ Save 📄 Cite Cited by 15 Related articles

## Overwriting pretrained bias with finetuning data

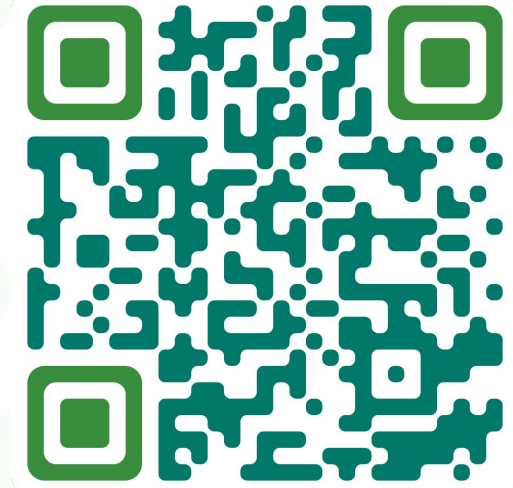
[A Wang](#), [O Russakovsky](#) - Proceedings of the IEEE/CVF ..., 2023 - openaccess.thecvf.com

Transfer learning is beneficial by allowing the expressive features of models pretrained on large-scale datasets to be finetuned for the target task of smaller, more domain-specific ...

☆ Save 📄 Cite Cited by 17 Related articles All 7 versions 🔗

# How you can support this work

- **Download** the Dollar Street dataset from MLcommons at: [mlcommons.org/en/dollar-street](https://mlcommons.org/en/dollar-street)
- **Contact** us at [dollarstreet@mlcommons.org](mailto:dollarstreet@mlcommons.org)
- **Learn** more about the Dollar Street project by visiting: [gapminder.org/dollar-street](https://gapminder.org/dollar-street)
- **Spread the word!**



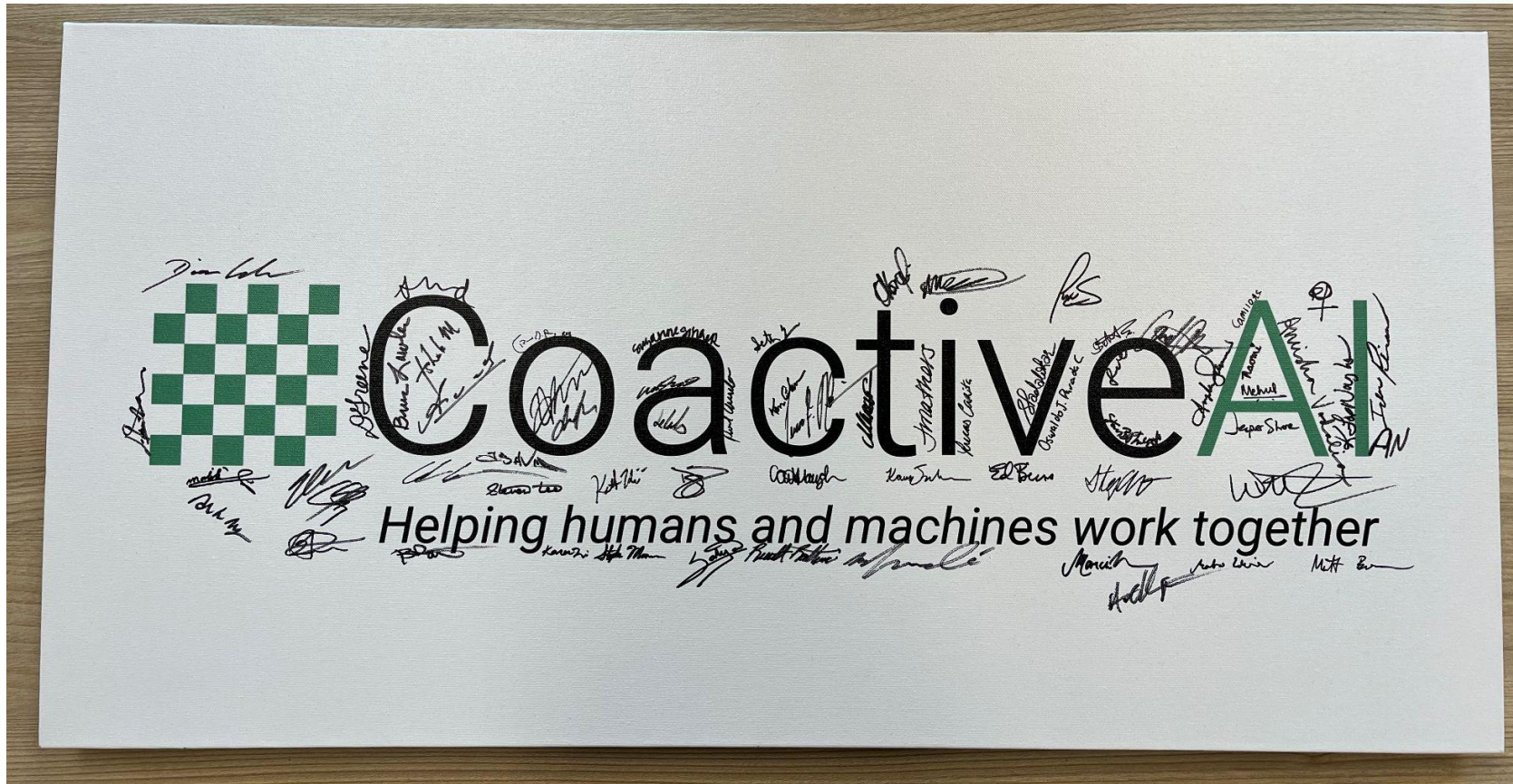
**Download the Dollar Street dataset**



**Visit the Dollar Street project**

# We're hiring!

Build multimodal AI systems that help humans and machines work together.



<https://www.coactive.ai/careers>

# How you can support this work

- **Download** the Dollar Street dataset from MLcommons at: [mlcommons.org/en/dollar-street](https://mlcommons.org/en/dollar-street)
- **Contact** us at [dollarstreet@mlcommons.org](mailto:dollarstreet@mlcommons.org)
- **Learn** more about the Dollar Street project by visiting: [gapminder.org/dollar-street](https://gapminder.org/dollar-street)
- **Spread the word!**



**Download the Dollar Street dataset**



**Visit the Dollar Street project**