

RESEARCH REVIEW 2023

**Carnegie
Mellon
University**
Software
Engineering
Institute

Leveraging Adversarial Machine Learning Techniques to Perform Query-Access Fairness Evaluations

NOVEMBER 2023

Anusha Sinha
Machine Learning Research Scientist

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

©2023

Document Markings

Copyright 2023 Carnegie Mellon University.

This material is based upon work funded and supported by the Department of Defense under Contract No. FA8702-15-D-0002 with Carnegie Mellon University for the operation of the Software Engineering Institute, a federally funded research and development center.

The view, opinions, and/or findings contained in this material are those of the author(s) and should not be construed as an official Government position, policy, or decision, unless designated by other documentation.

NO WARRANTY. THIS CARNEGIE MELLON UNIVERSITY AND SOFTWARE ENGINEERING INSTITUTE MATERIAL IS FURNISHED ON AN "AS-IS" BASIS. CARNEGIE MELLON UNIVERSITY MAKES NO WARRANTIES OF ANY KIND, EITHER EXPRESSED OR IMPLIED, AS TO ANY MATTER INCLUDING, BUT NOT LIMITED TO, WARRANTY OF FITNESS FOR PURPOSE OR MERCHANTABILITY, EXCLUSIVITY, OR RESULTS OBTAINED FROM USE OF THE MATERIAL. CARNEGIE MELLON UNIVERSITY DOES NOT MAKE ANY WARRANTY OF ANY KIND WITH RESPECT TO FREEDOM FROM PATENT, TRADEMARK, OR COPYRIGHT INFRINGEMENT.

[DISTRIBUTION STATEMENT A] This material has been approved for public release and unlimited distribution. Please see Copyright notice for non-US Government use and distribution.

This material may be reproduced in its entirety, without modification, and freely distributed in written or electronic form without requesting formal permission. Permission is required for any other use. Requests for permission should be directed to the Software Engineering Institute at permission@sei.cmu.edu.

Carnegie Mellon® is registered in the U.S. Patent and Trademark Office by Carnegie Mellon University.

DM23-2009

Presentation Outline

- Motivations and project overview
- Understanding bias transfer through stereotypical examples
- Quantifying bias transfer through model finetuning
- Expected impacts and takeaways

Motivations



DoD hiring teams are **overwhelmed by scale** and need **unbiased automation** to manage the world's largest workforce.

Machine learning (ML) models applied to workforce management tasks are increasingly **finetuned from large public models**.

We need to understand **bias transfer** from large public models.

Project Overview

We noticed:

- Bias in ML models contributes to unfair decisions.
- Transfer learning is increasingly being applied.



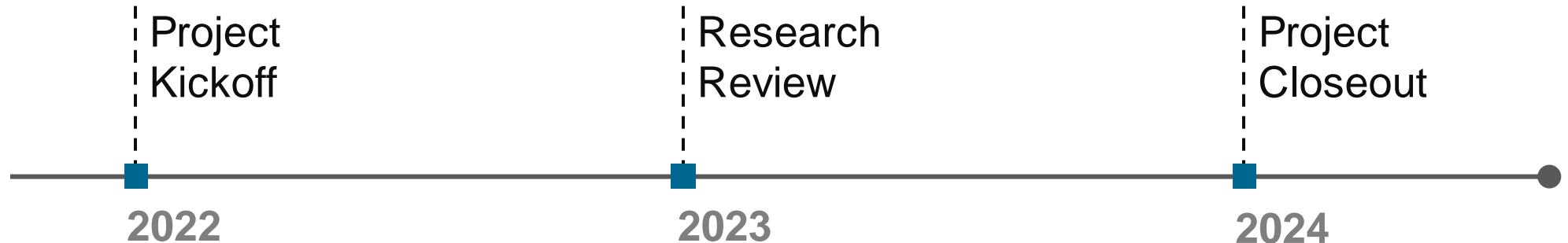
Our goal:

- Understand how bias can transfer from large foundational models to child models finetuned from them.

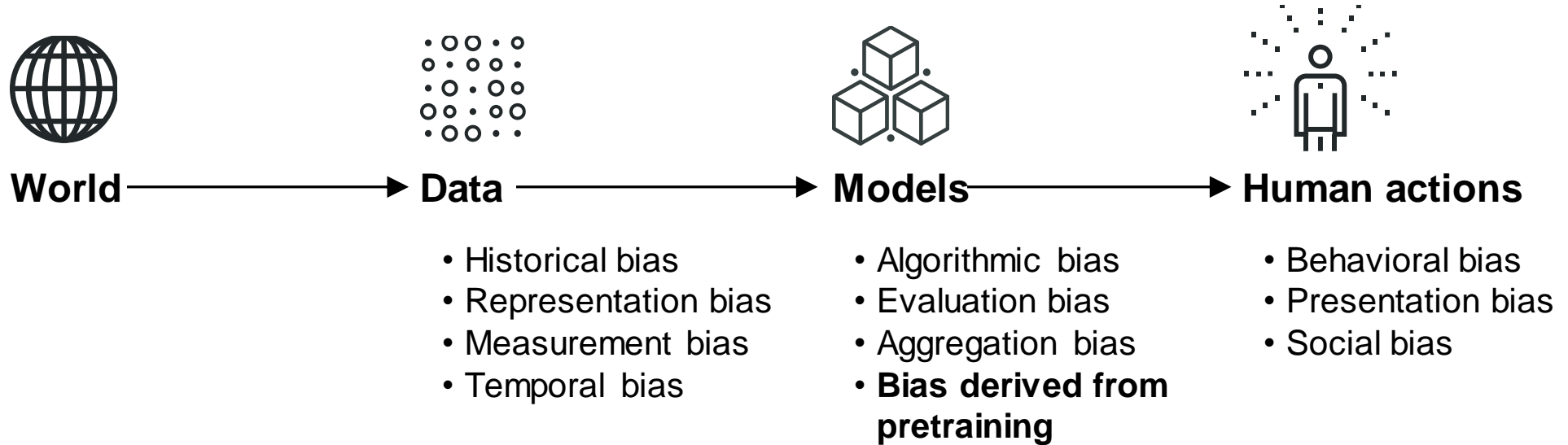


Intended impact:

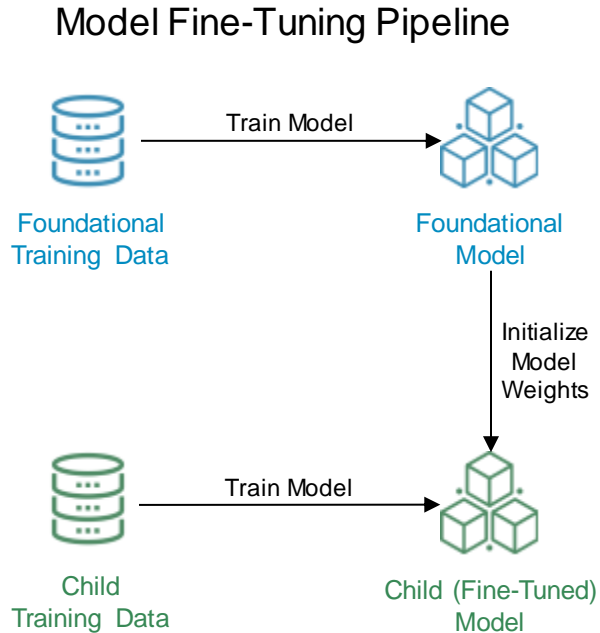
- Inform best practices for bias mitigation in model finetuning pipelines.



Sources of Bias in ML Decision-Making Pipelines



Overview of Experiments



- We wanted to understand how properties of foundational models transfer to child models.
 - Bias 1: strong stimuli
 - Bias 2: performance disparity
- Can we use our methods to quantify bias transfer to inform best practices for model finetuning?

Presentation Outline

- Motivations and project overview
- Understanding bias transfer through stereotypical examples
- Quantifying bias transfer through model finetuning
- Expected impacts and takeaways

Initial Experiments Using Stereotypical Examples



Stereotypical Example for ImageNet
Class: Indian Cobra

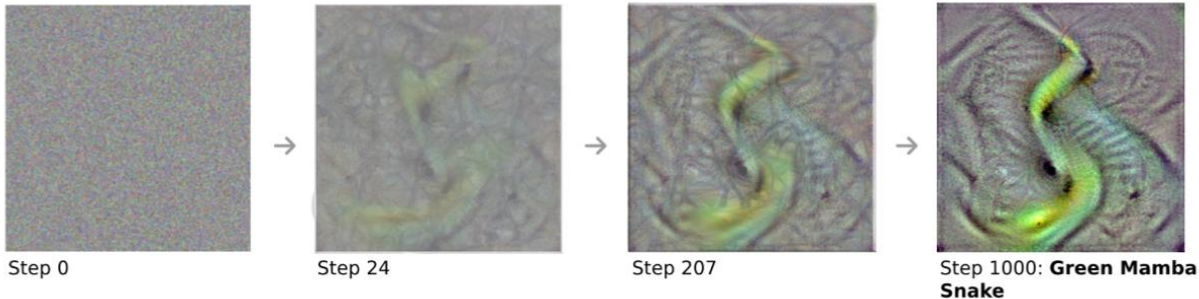
Initial research question: How do stereotypical examples of classes as defined by foundational models persist in child models?

Process:

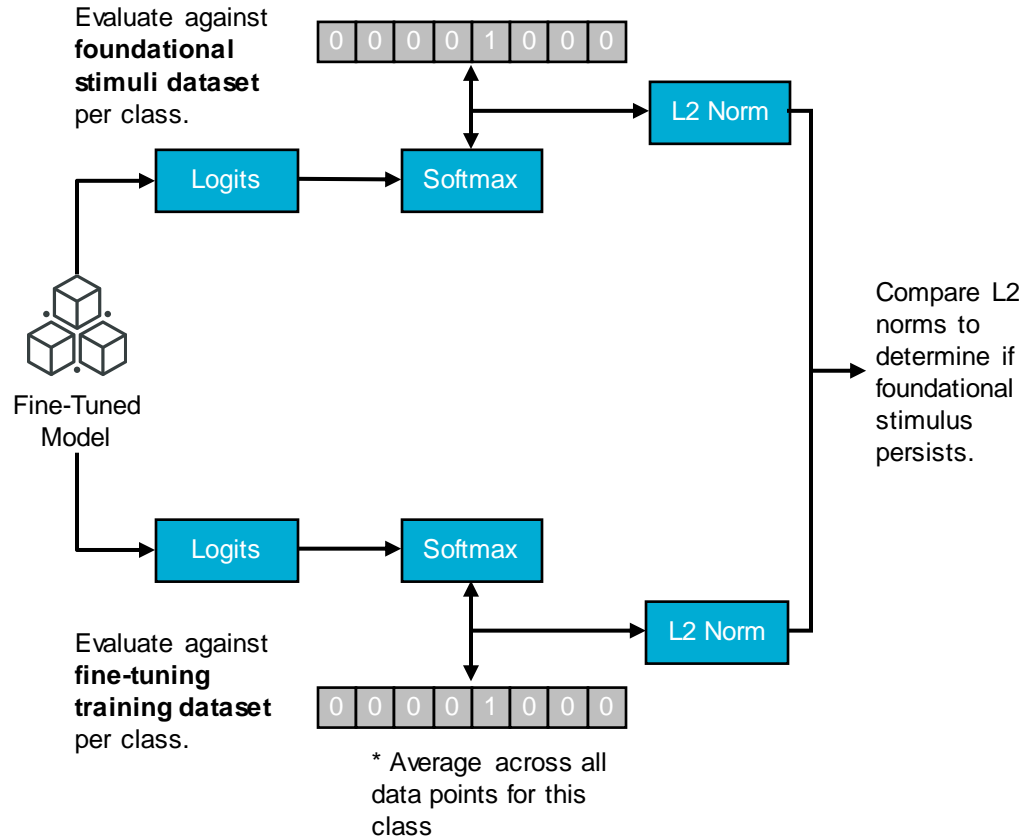
1. Produce strong stimuli against foundational models.
2. Evaluate against spectra of finetuned child models.
3. Determine if finetuning methods result in bias transfer.

Producing Stereotypical Examples

- We produced strong foundational model stimuli using gradient-based optimization.
- We iteratively applied gradient updates to minimize a cost function between model outputs and the strongest output for a target class (one-hot vector).



Transfer of Stereotypical Examples

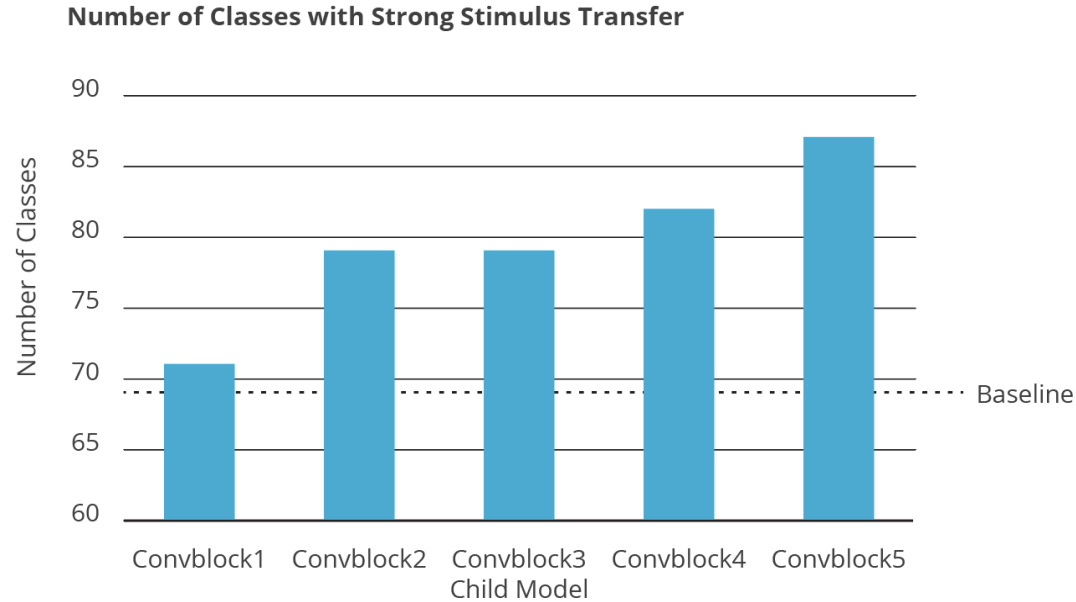


We first determine child model confidence on strong stimuli produced against foundational model.

We then use the centroid of the training data for the target class as a baseline for stimulus performance.

Discussion of Results and Implications for Bias Transfer

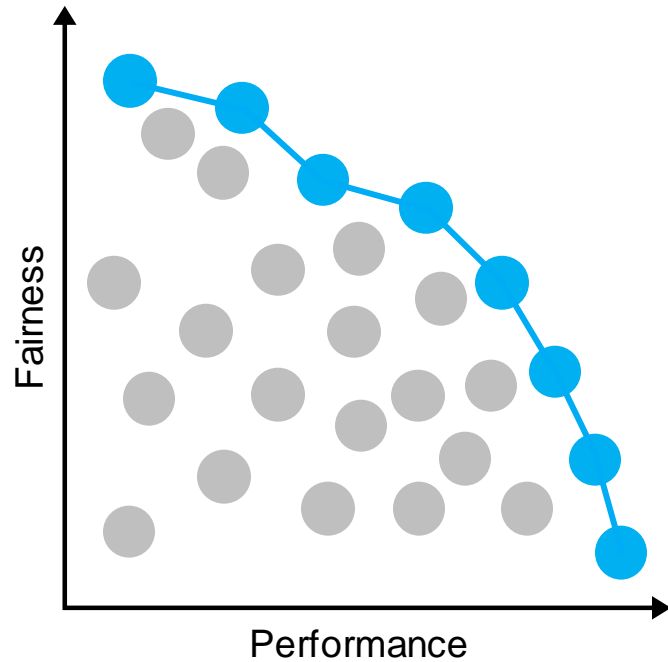
Main trend: A greater number of strong stimuli outperform training centroids when fewer layers of the model are retrained.



Presentation Outline

- Motivations and project overview
- Understanding bias transfer through stereotypical examples
- Quantifying bias transfer through model finetuning
- Expected impacts and takeaways

Broad Research Questions



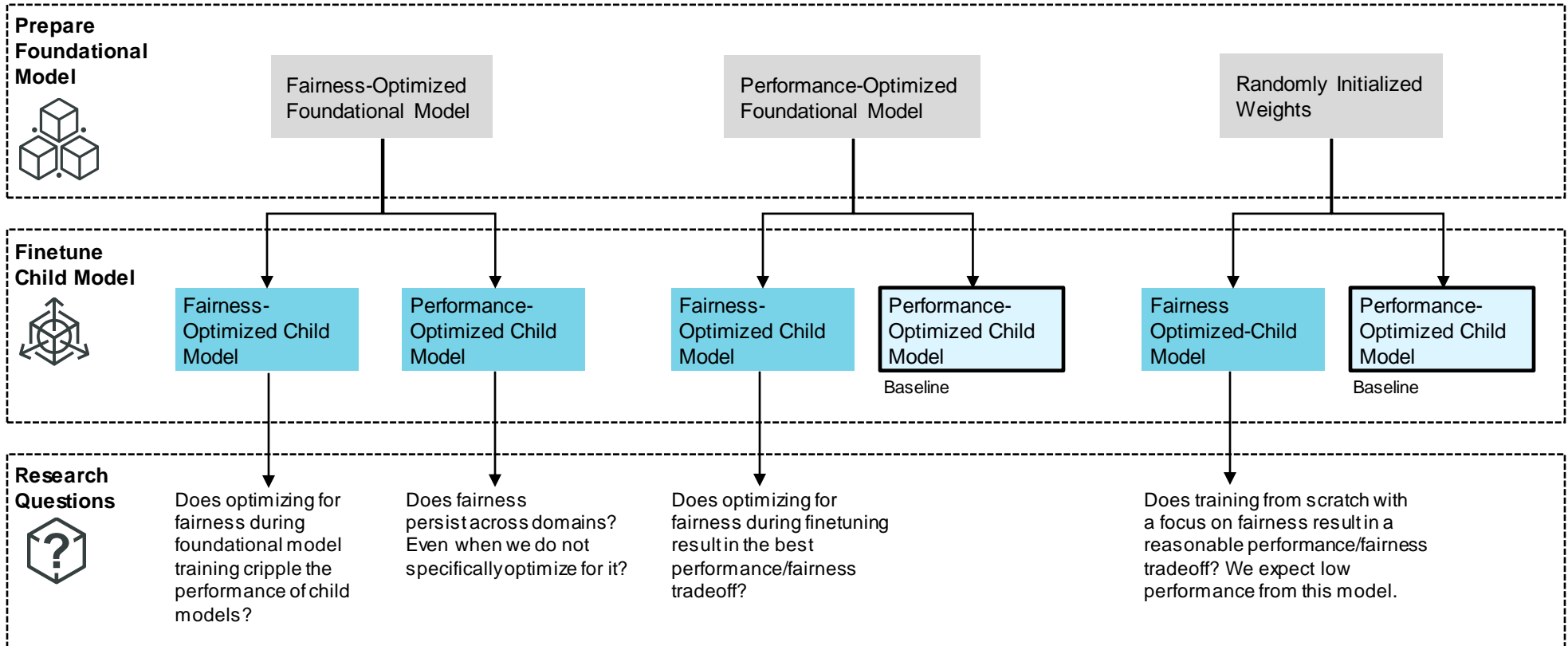
We want to quantify:

The transfer of data-derived biases from foundational to child models

Context:

The tradeoff between bias mitigation and model performance

Experimental Design



Using a Fairness-Aware Image Classification Dataset



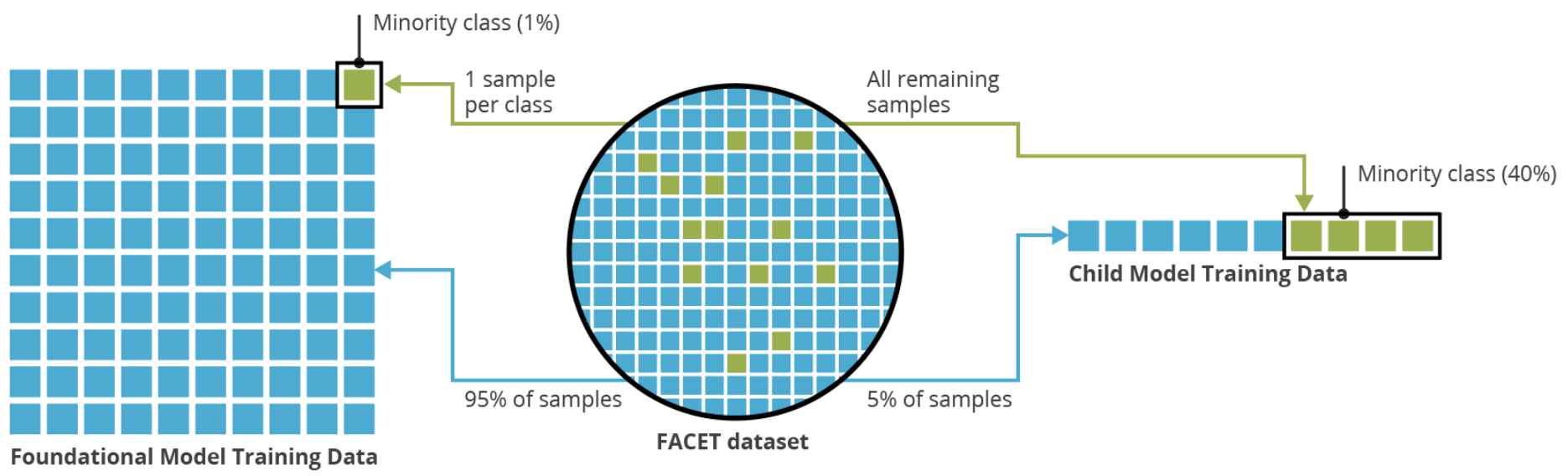
Samples from Class 4 (Basketball Player)

We use a custom split of the Fairness in Computer Vision Evaluation Benchmark (FACET) dataset produced by Meta for our experiments

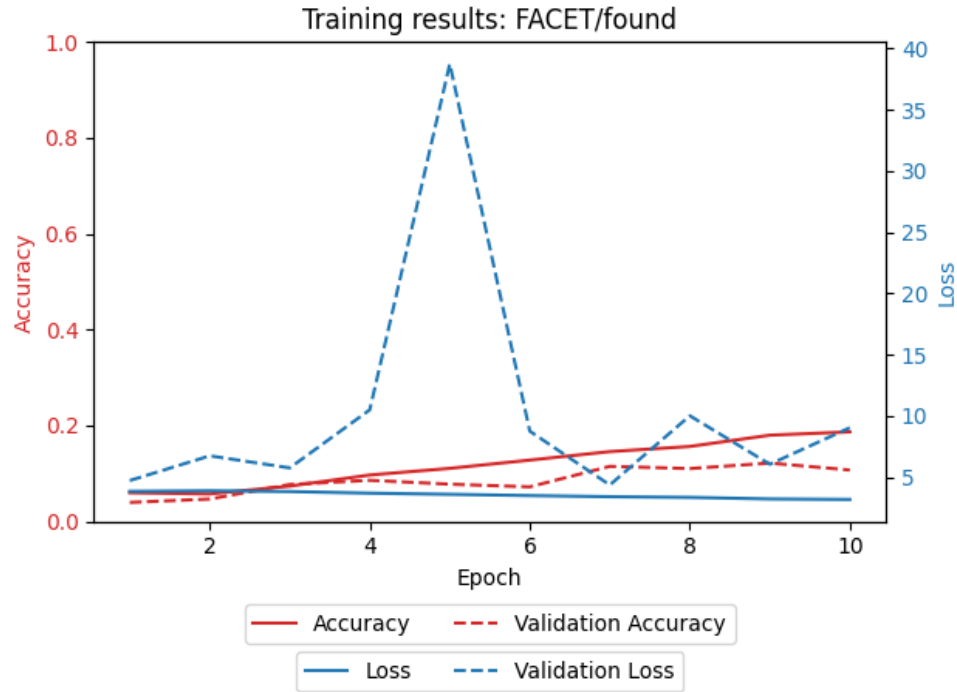
- 50,000 images of people
- 52 roles, such as astronaut, judge, basketball player, etc.
- Accompanied by secondary labels describing pictured subjects
 - skin tone
 - gender presentation
 - age

Dataset Preparation

We produced foundational and child training sets by splitting the FACET dataset based on protected attributes (samples highlighted below).



Training Foundational Models



We are training models on FACET data using cross-entropy loss and focal loss.

Challenges:

- Training on FACET data from scratch is unstable.

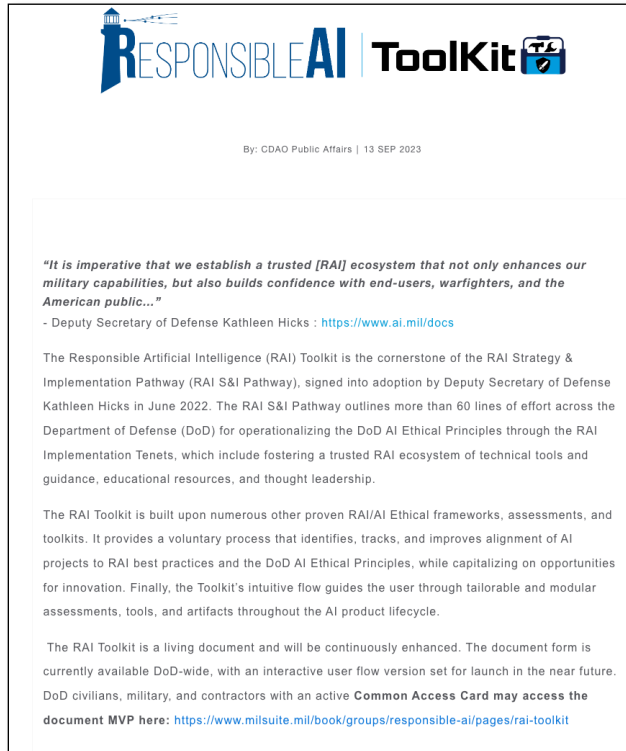
Currently:

- Working on hyperparameter optimization for both foundational models.
- *Fallback option*: Initialize models with ImageNet weights, but this may muddy results.

Presentation Outline

- Motivations and project overview
- Understanding bias transfer through stereotypical examples
- Quantifying bias transfer through model finetuning
- Expected impacts and takeaways

Expected Impacts and Takeaways



- Bias transfer between foundational and child models exists and can be quantified.
- *Hypothesis*: Bias mitigation strategies may be most effective when applied at the final stages of model training pipelines
- Implications for Chief Digital and Artificial Intelligence Office (CDAO) Responsible AI Toolkit

Contact



Anusha Sinha (PI)
Machine Learning
Research Scientist
SEI AI
AI-E Engineering
Center



**Nathan VanHoudnos
(advisor)**
Senior Machine
Learning Research
Scientist and Lab Lead
SEI AI AI-E Adversarial
ML Lab



Hayden Moore
Associate Software
Developer
SEI AI AI-E Engineering
Center



Swati Rallapalli
Senior Machine Learning
Research Scientist
SEI AI
AI-E Engineering Center



Hoda Heidari
Assistant Professor in
the School of Computer
Science at CMU



Steven Wu
Assistant Professor in
the School of Computer
Science at CMU

Telephone: +1 412.268.5800

Email: info@sei.cmu.edu