

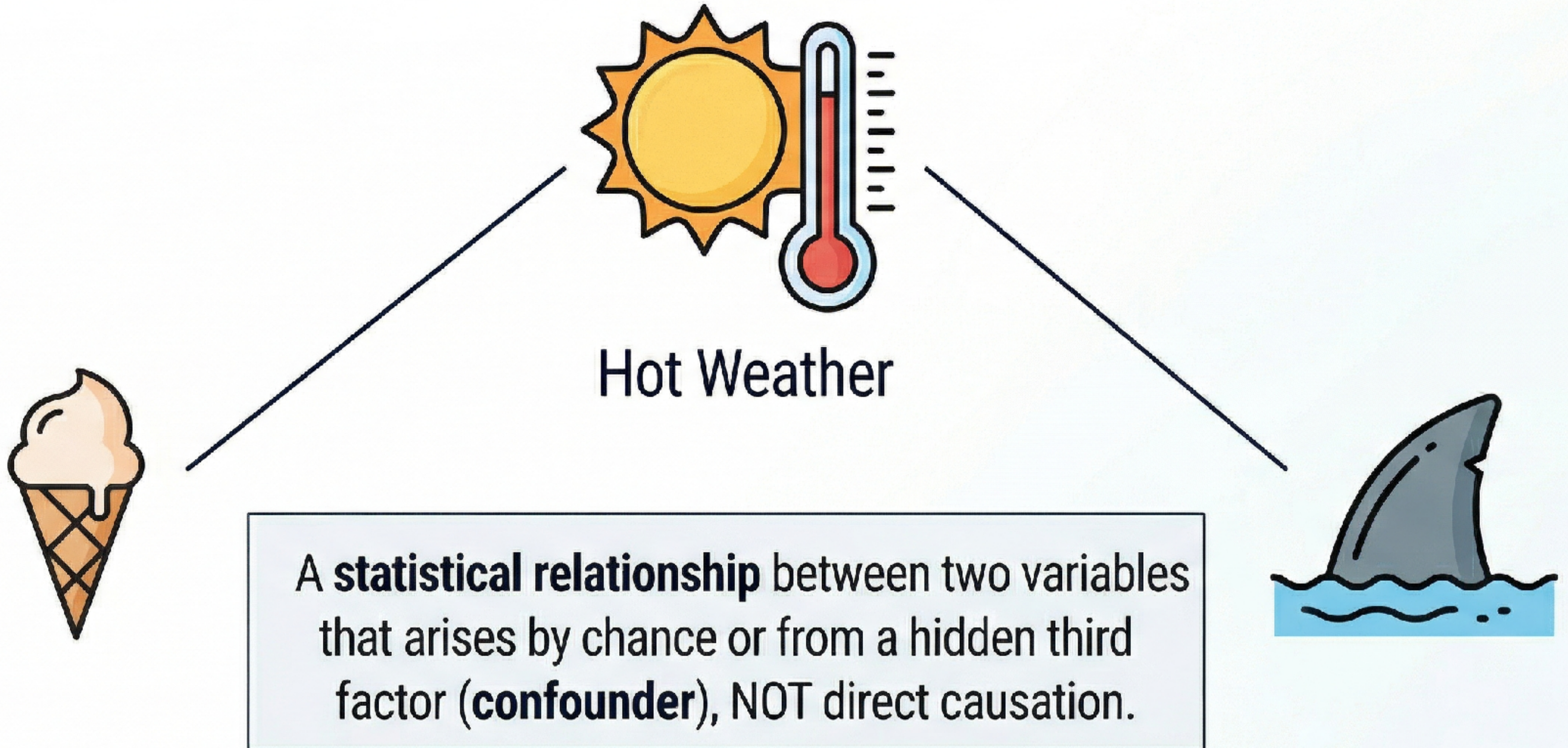
# **CSCI1470**

## **Deep Learning**

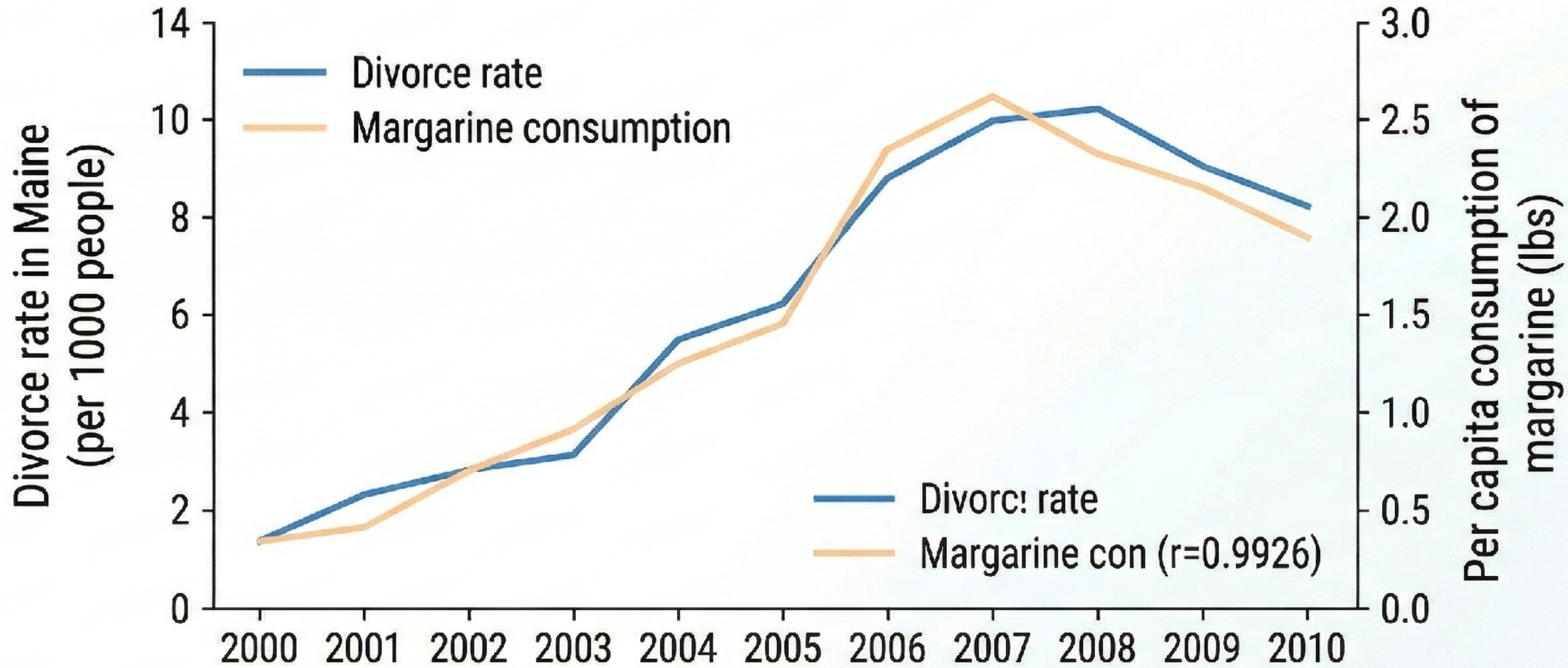
**Randall Balestrieri**

**Recap**

# What is a Spurious Correlation?



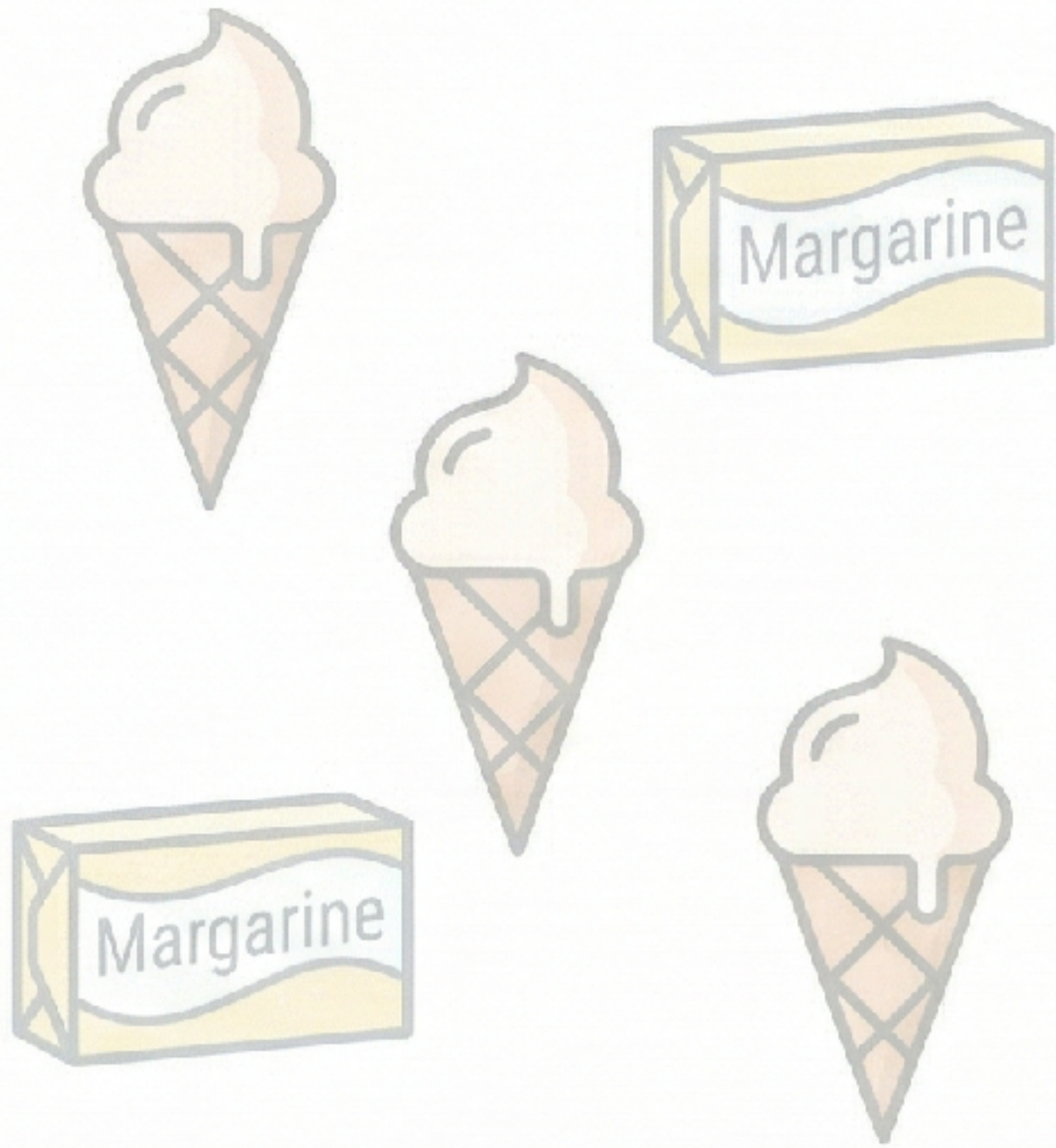
# Famous Spurious Correlation Example



Data source: U.S. Census & USDA (Visualization style inspired by Tyler Vigen).

# Why Should We Care? The Shift to Fairness

## Harmless Data

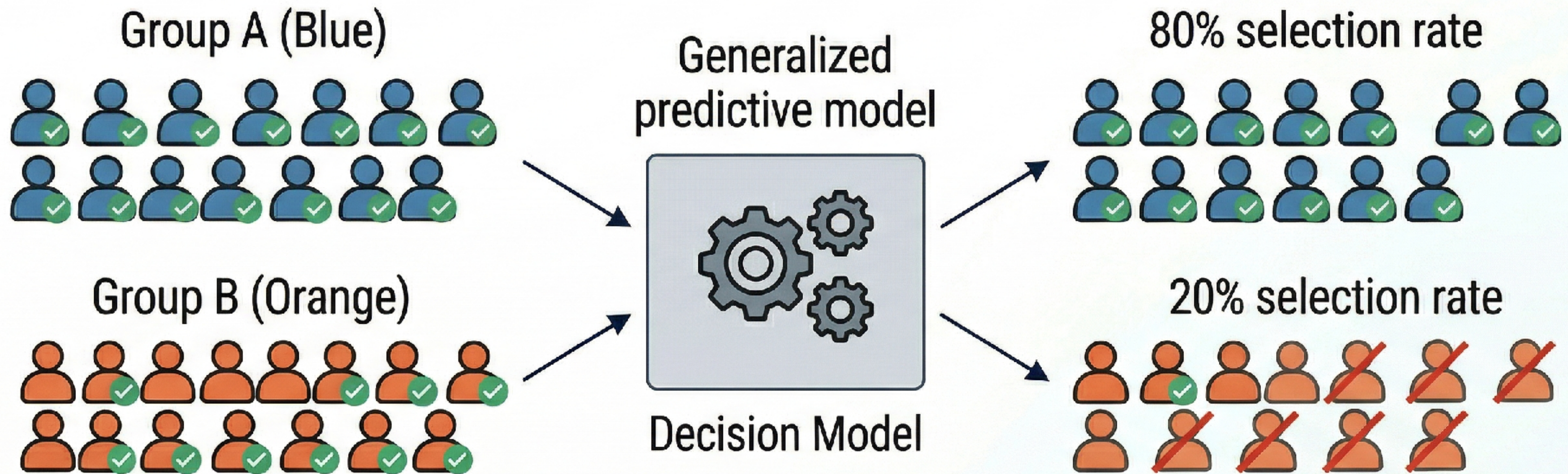


## High-Stakes Decisions



When **data models** use **spurious correlations** (e.g., using a proxy variable relates with a protected group) for decisions about people, it leads to **systemic bias**.

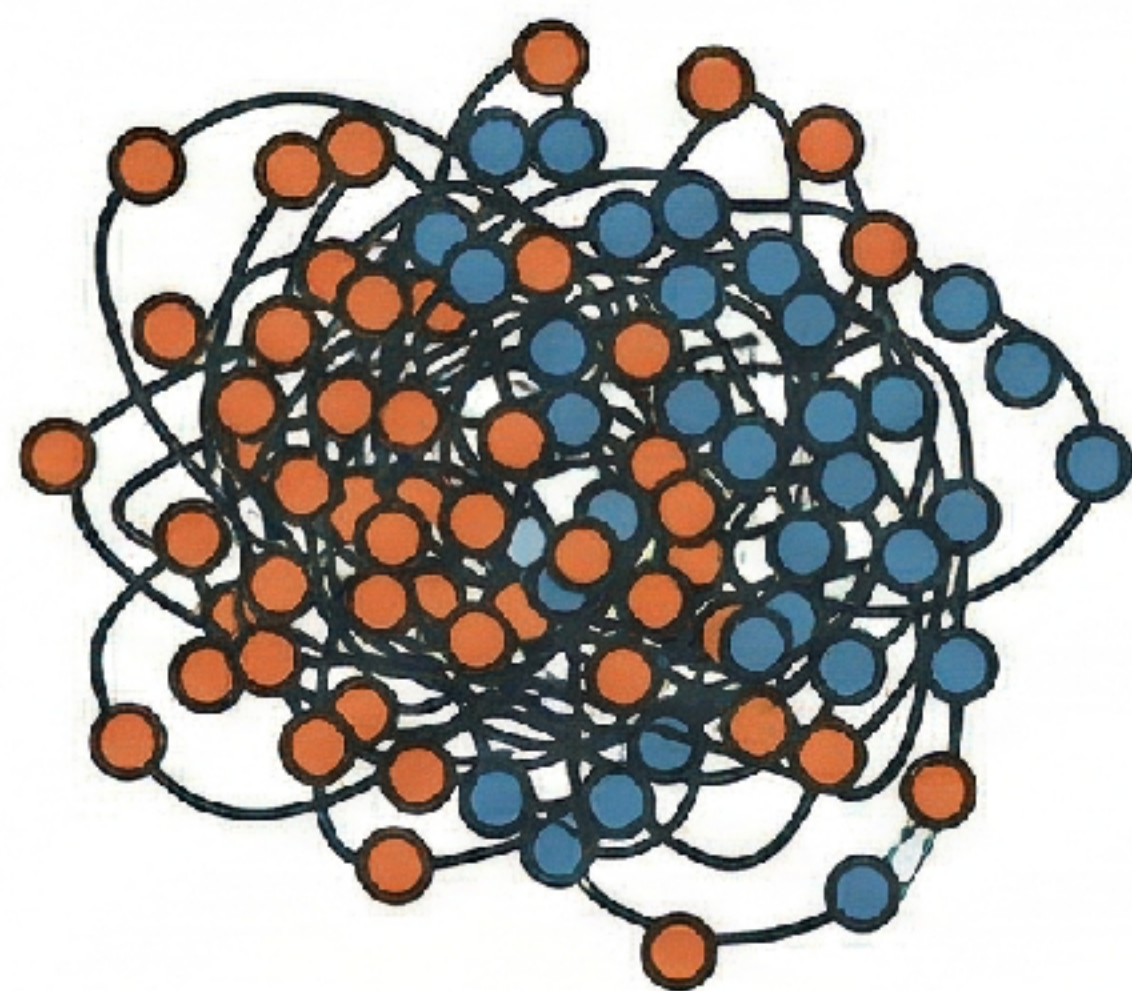
# Defining Group Fairness



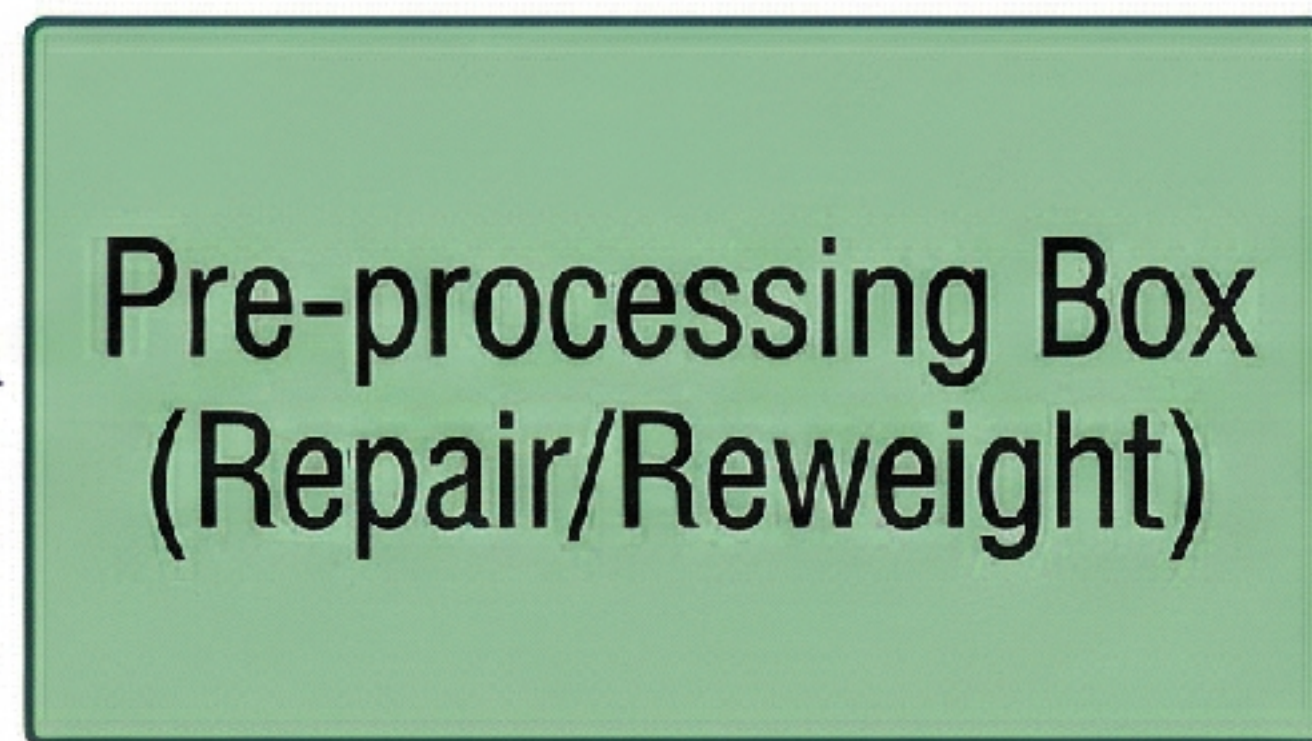
**Goal:** Ensuring that the outcomes of a model are **comparable** across different demographic groups (e.g., race, gender), avoiding disparate impact, regardless of statistical correlation.

# De-biasing Techniques: Pre-processing

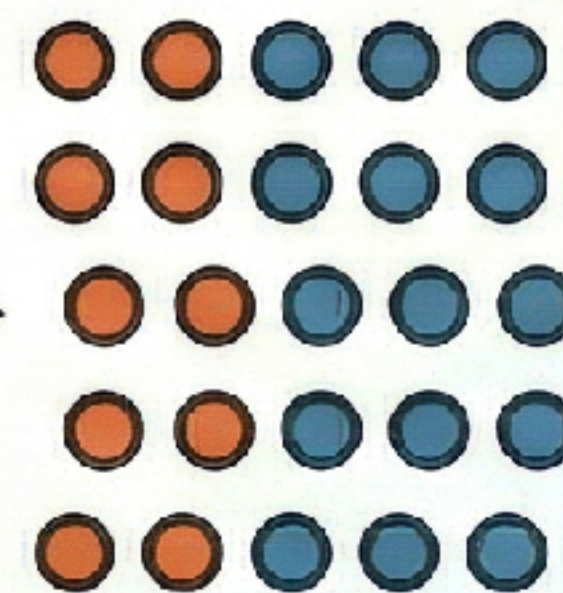
Raw (Biased) Training Data



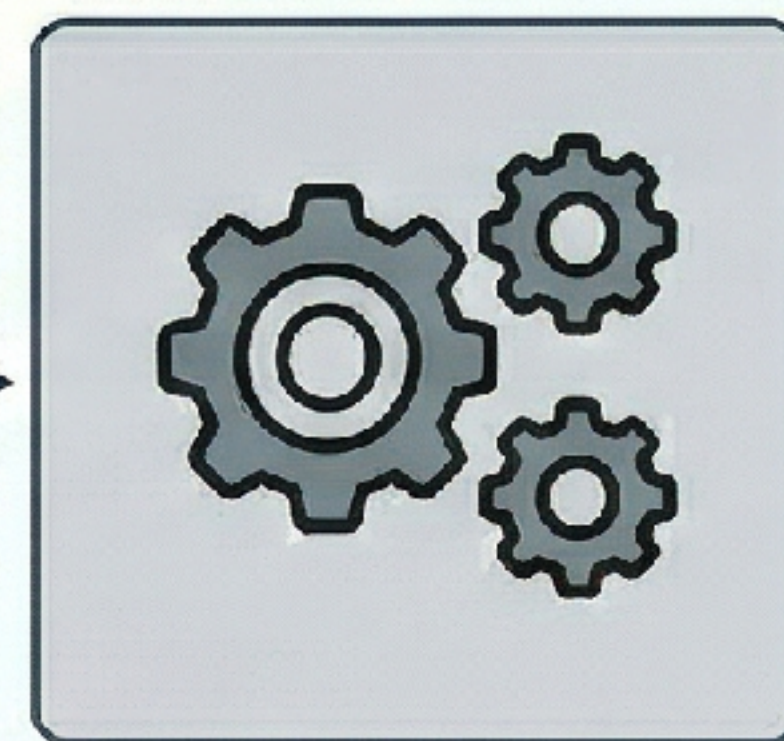
Raw (Biased)  
Training Data



Debiased  
(Repaired) Data

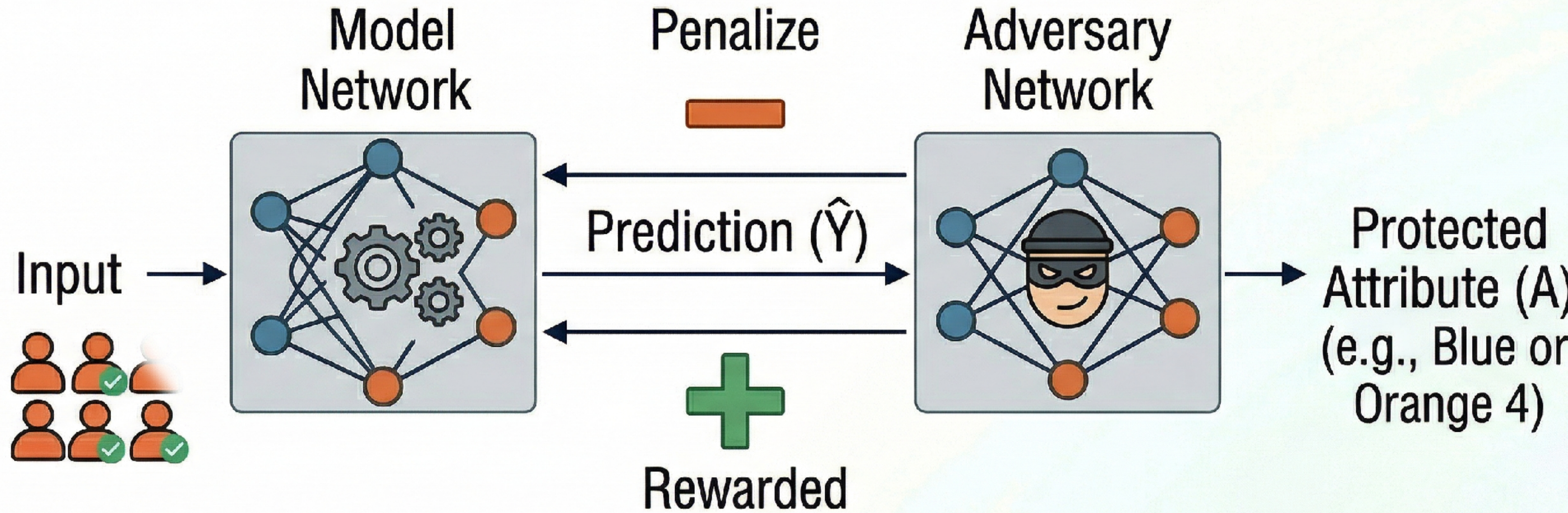


Debiased  
(Repaired) Data



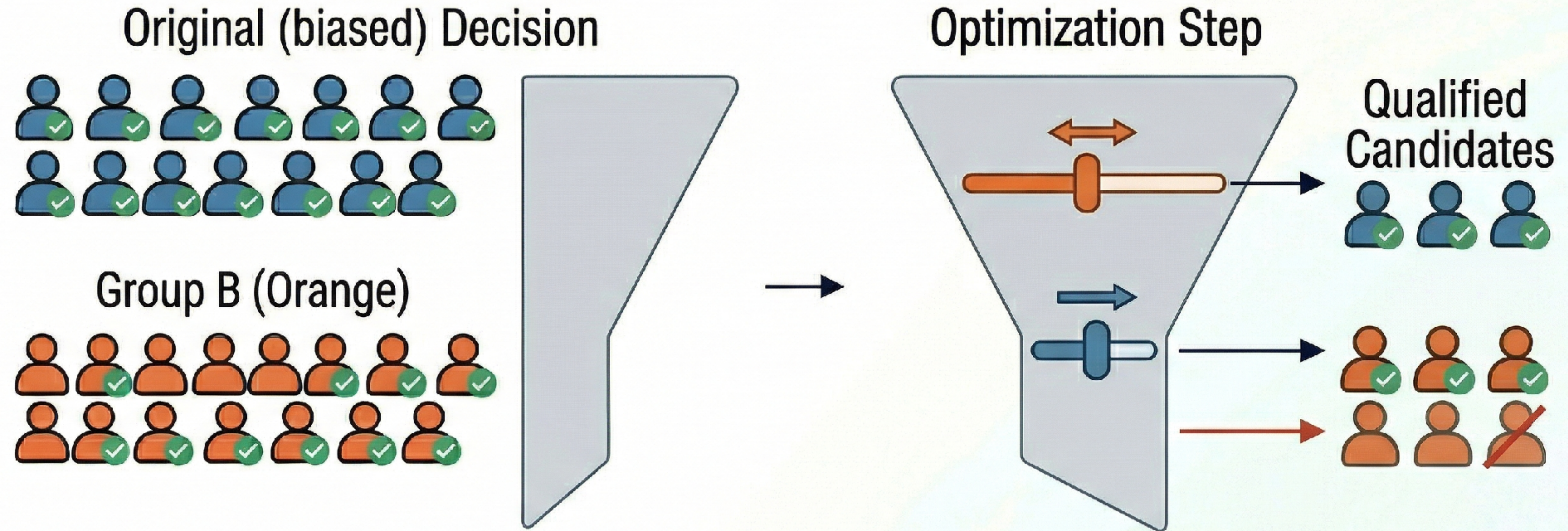
**Focus:** Fix the *data* distribution *before* training. Useful when data collection causes the spurious correlation (Image 1 confounder).

# In-processing: Adversarial Debiasing



The Model is trained to *hide* the protected attribute from the Adversary.

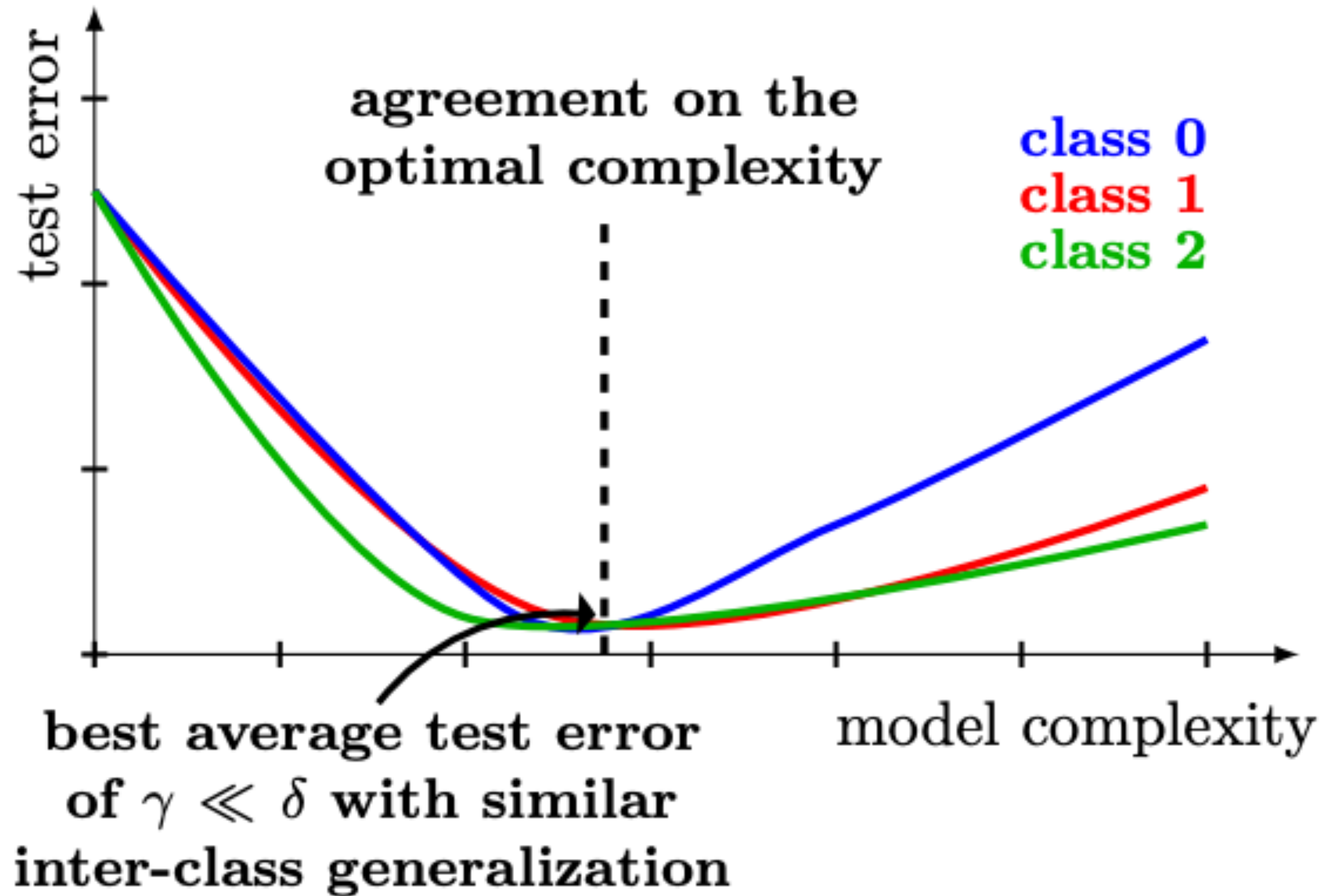
# Post-processing: Threshold Optimization (Equalized Odds)

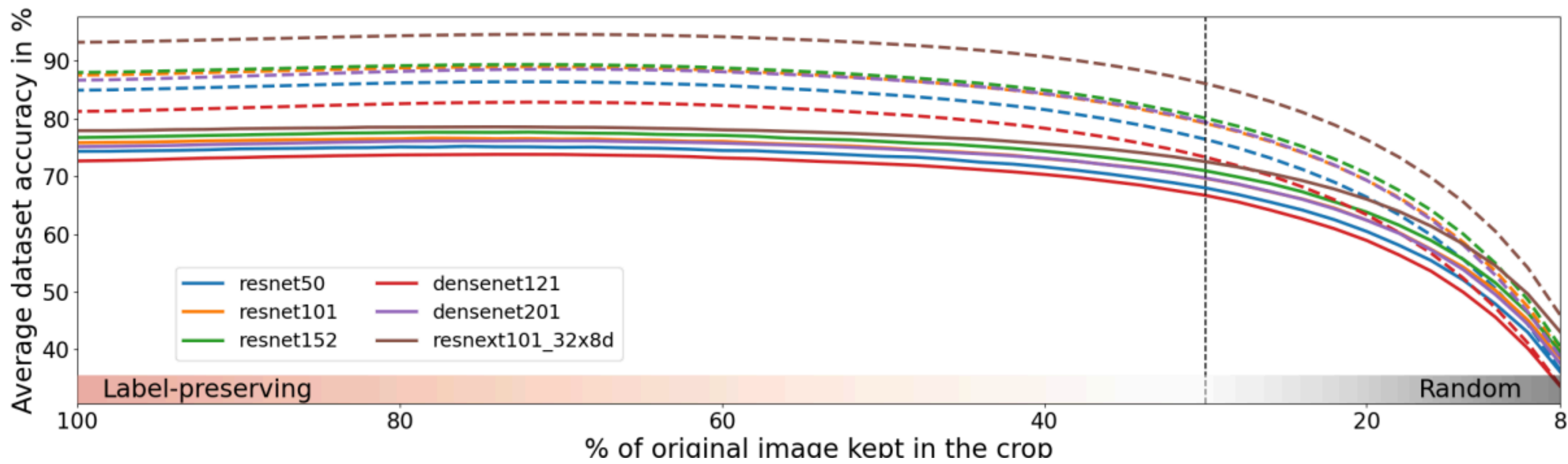


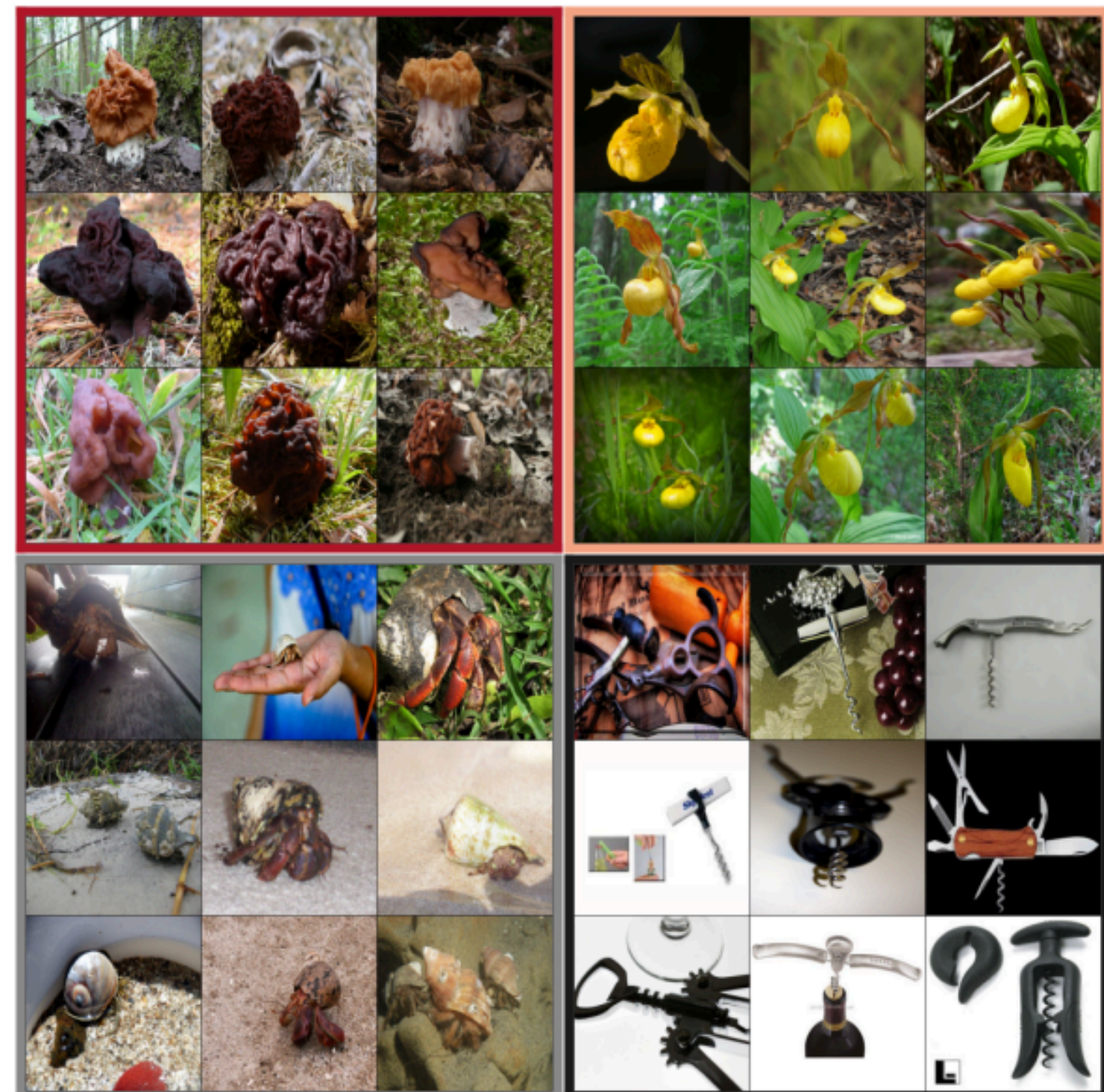
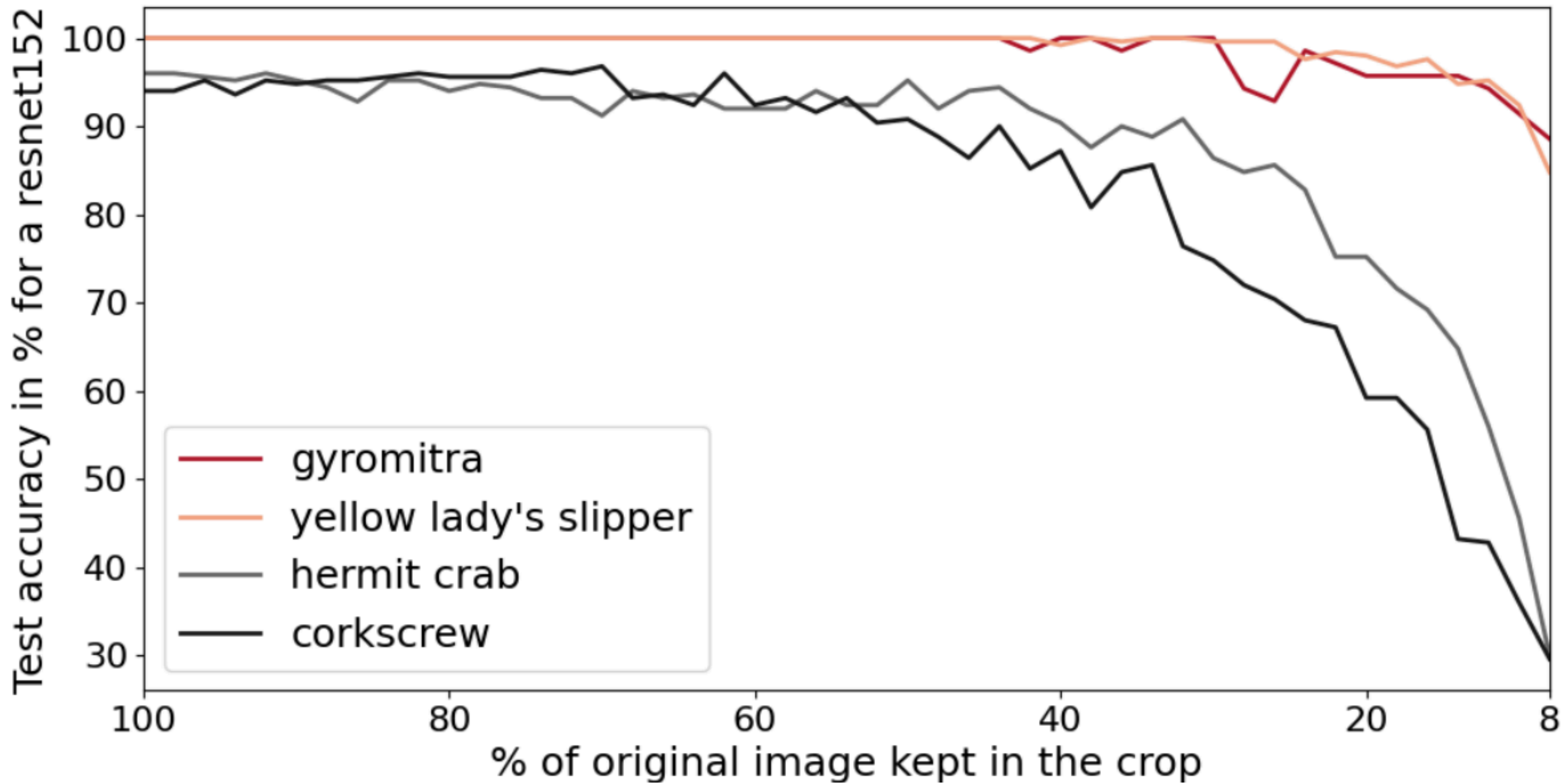
Adjust thresholds *after* deployment to meet fairness targets (Image 10 metrics)

**What other sources of  
unfairness?**

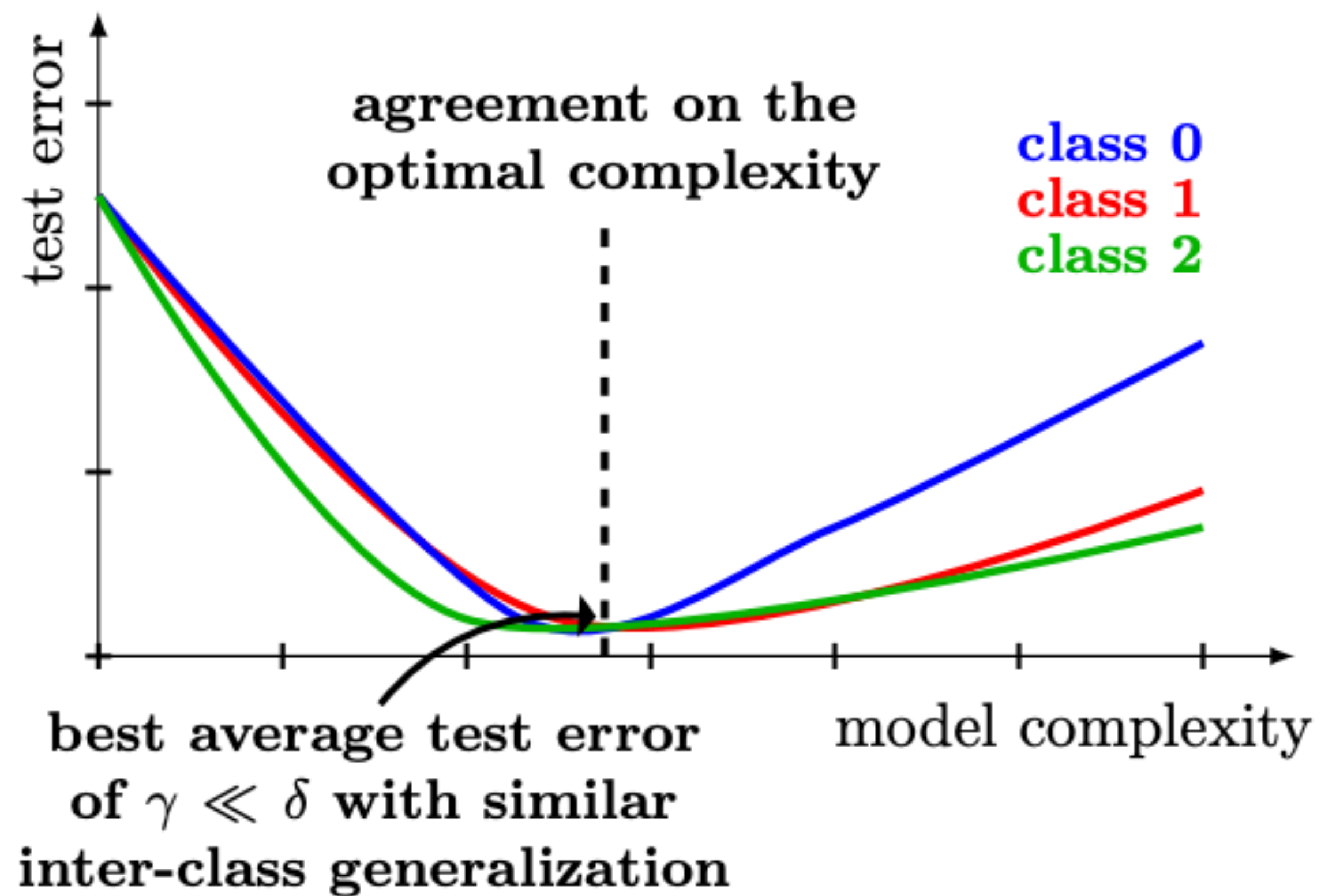
## Ideal Complexity Measure



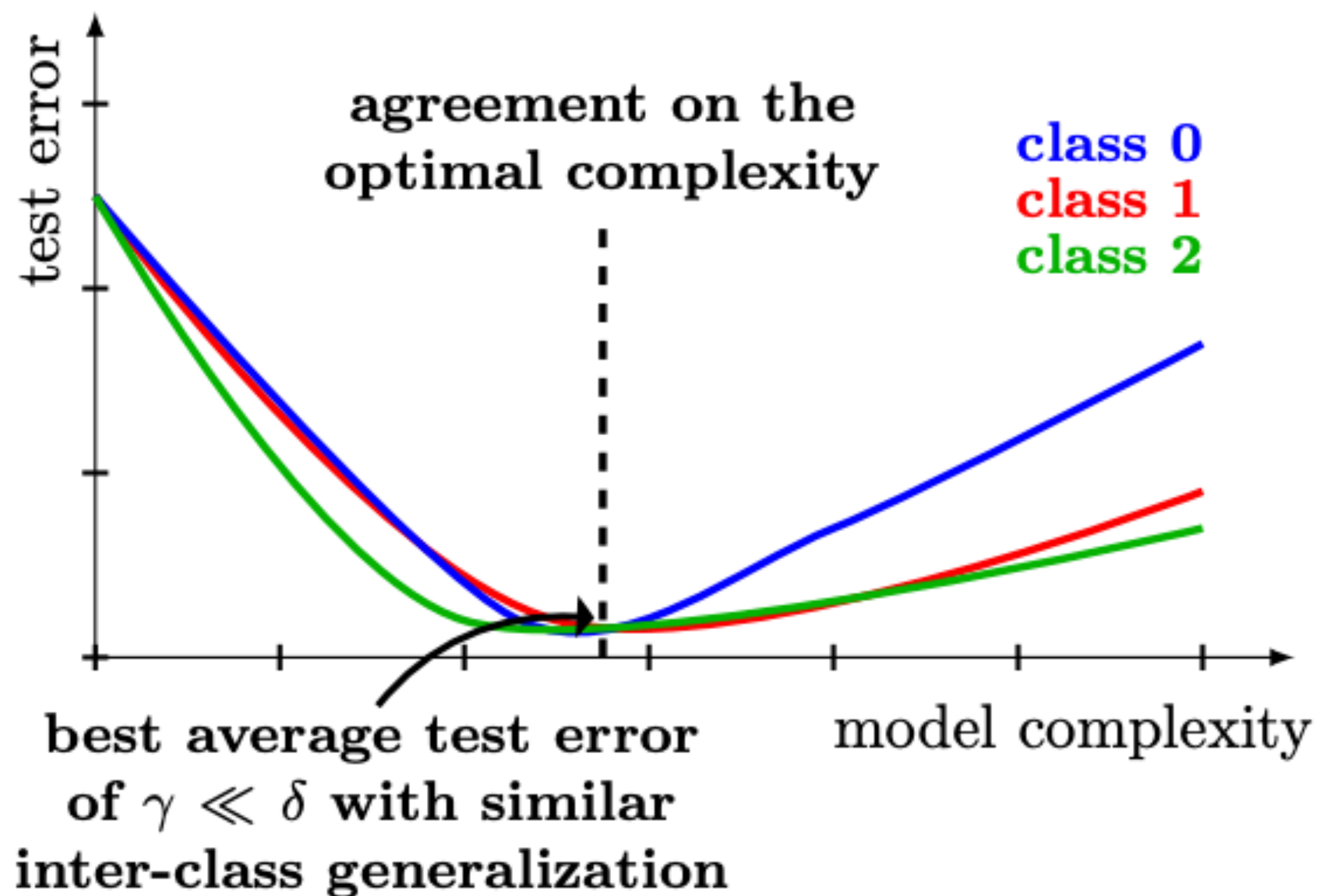




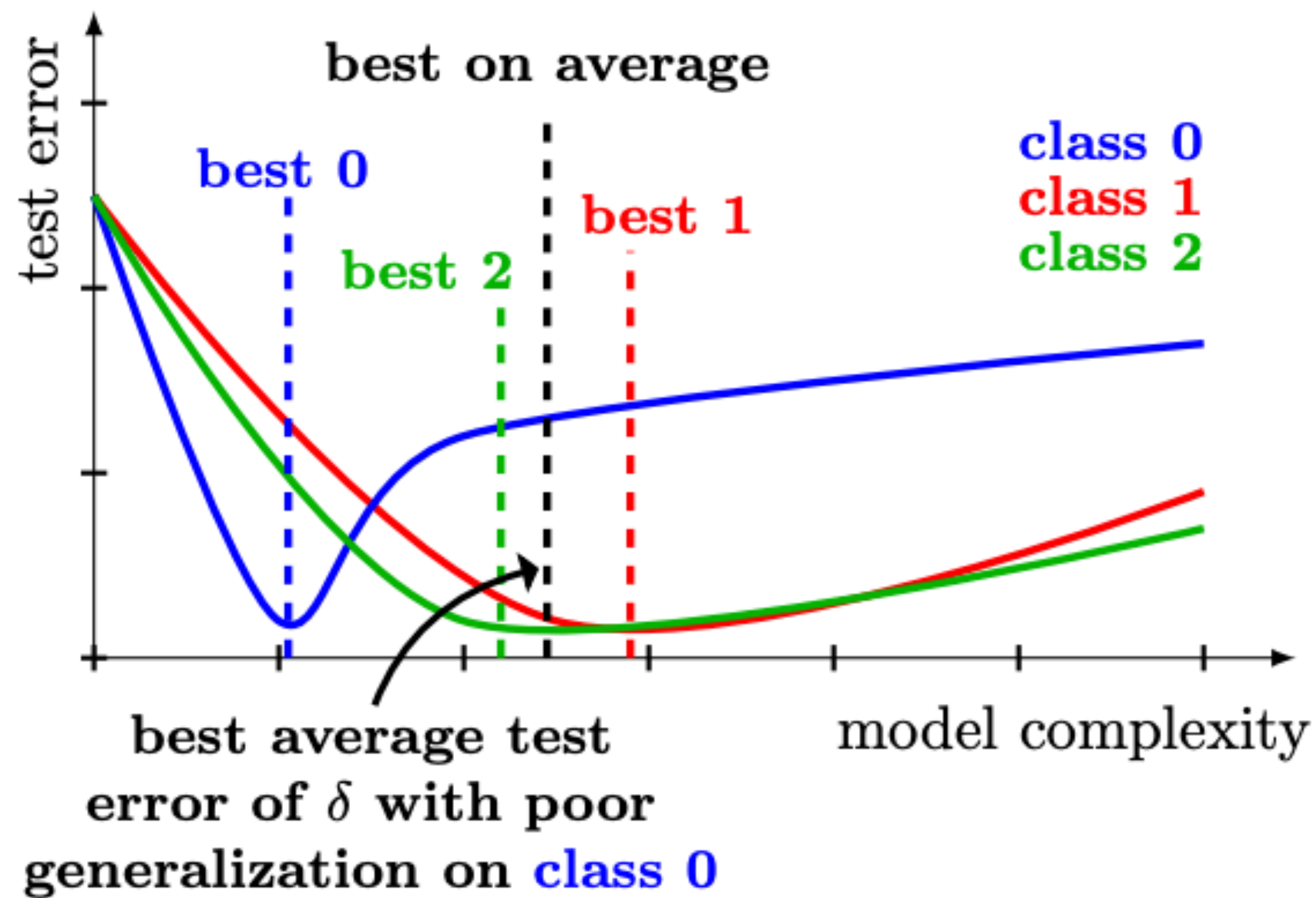
## Ideal Complexity Measure

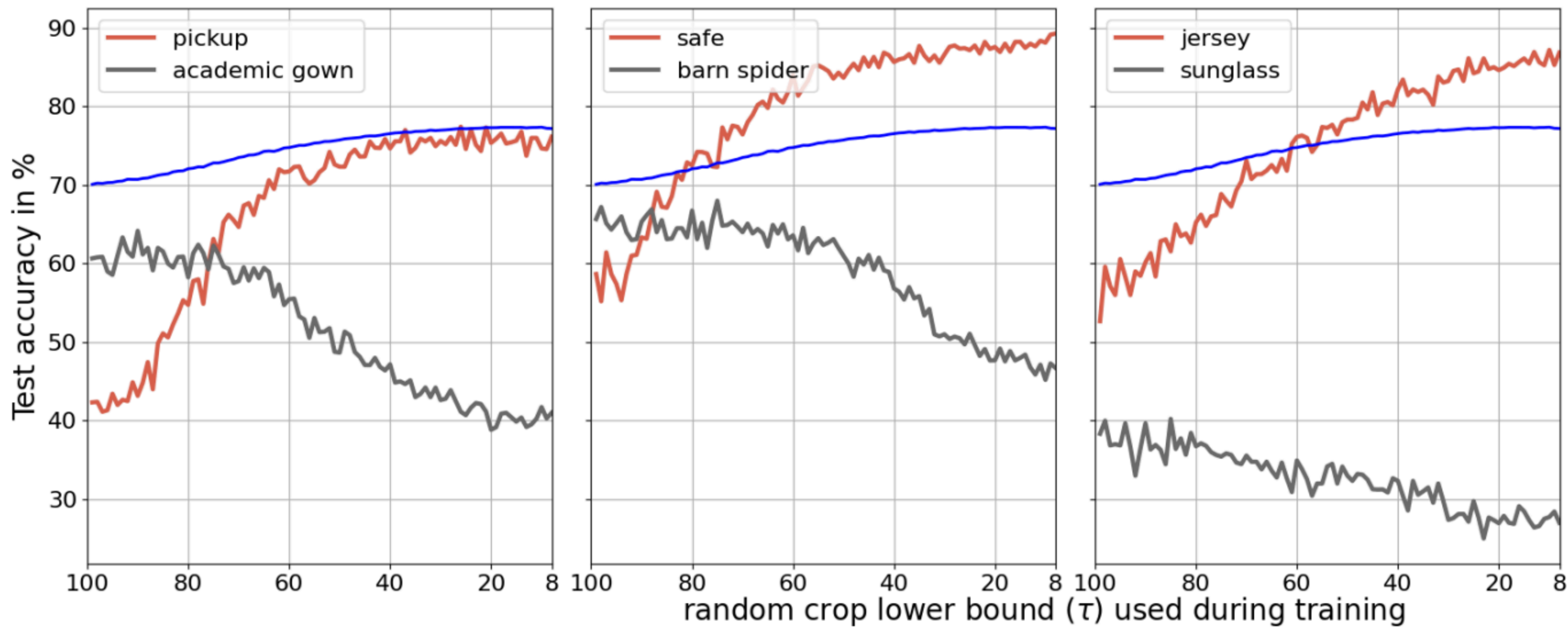


### Ideal Complexity Measure



### Current Deep Learning





**Thank you!**  
**See you Wednesday!**