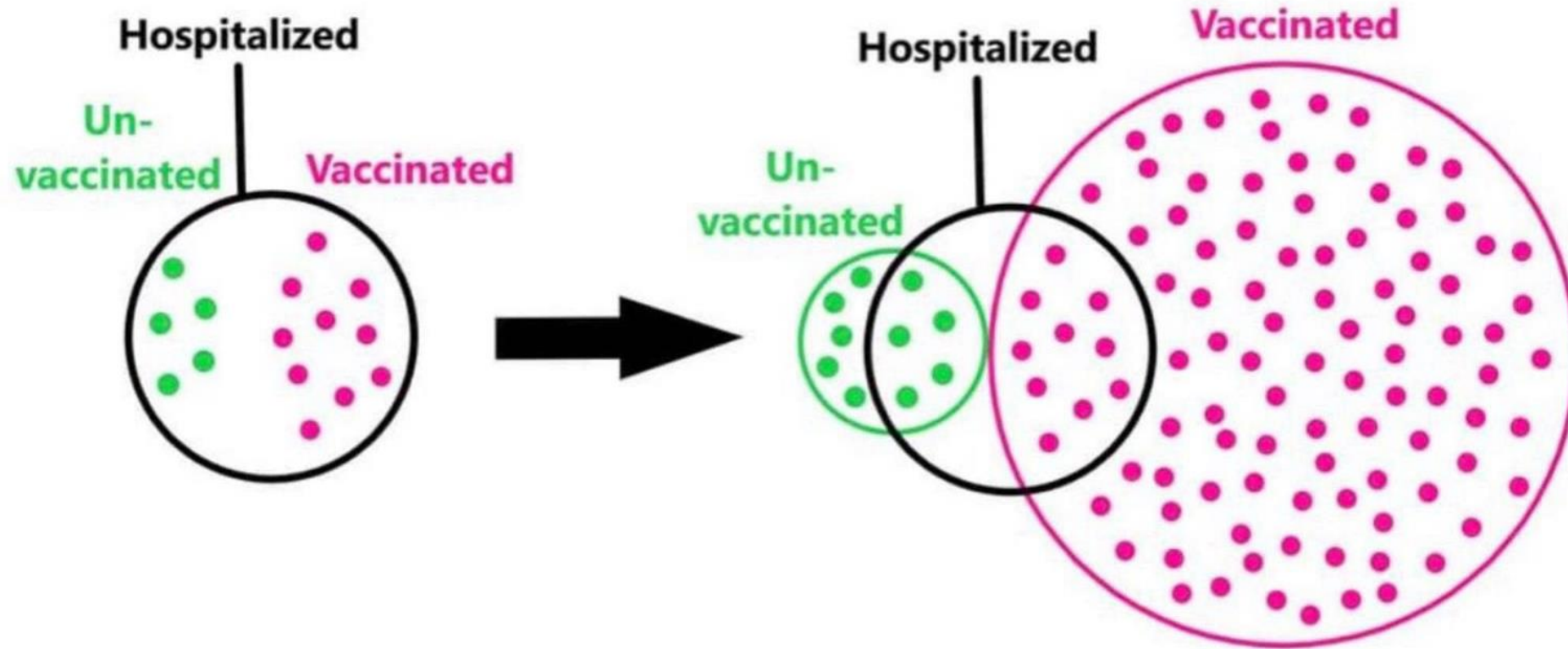


SoK: The Good, The Bad, and The Unbalanced

Measuring Structural Limitations of Deepfake Media Datasets

Seth Layton, Tyler Tucker, Daniel Olszewski, Kevin Warren, Kevin Butler, Patrick Traynor

The University of Florida



Background

Class Distribution and
Bias

Metric Usage Impacts

Base-Rate
Contextualization

Suggestions and
Greater Impacts

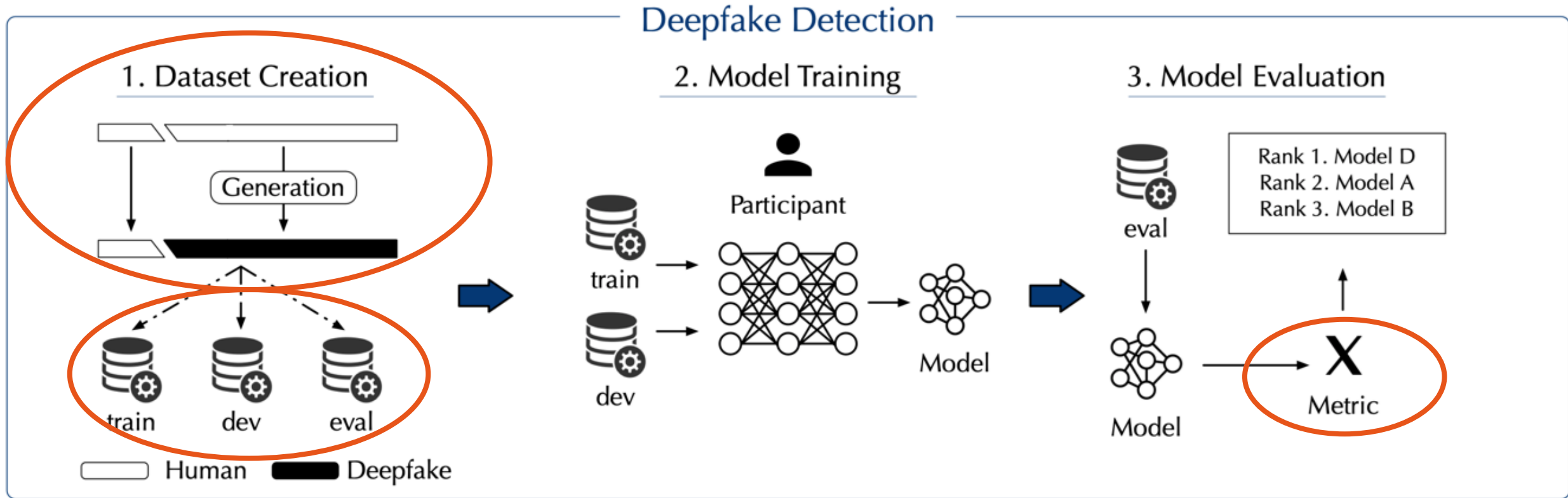
Background

Class Distribution and
Bias

Metric Usage Impacts

Base-Rate
Contextualization

Suggestions and
Greater Impacts



1. Class Distributions

2. Metric Usage

3. Base-Rate
Contextualization

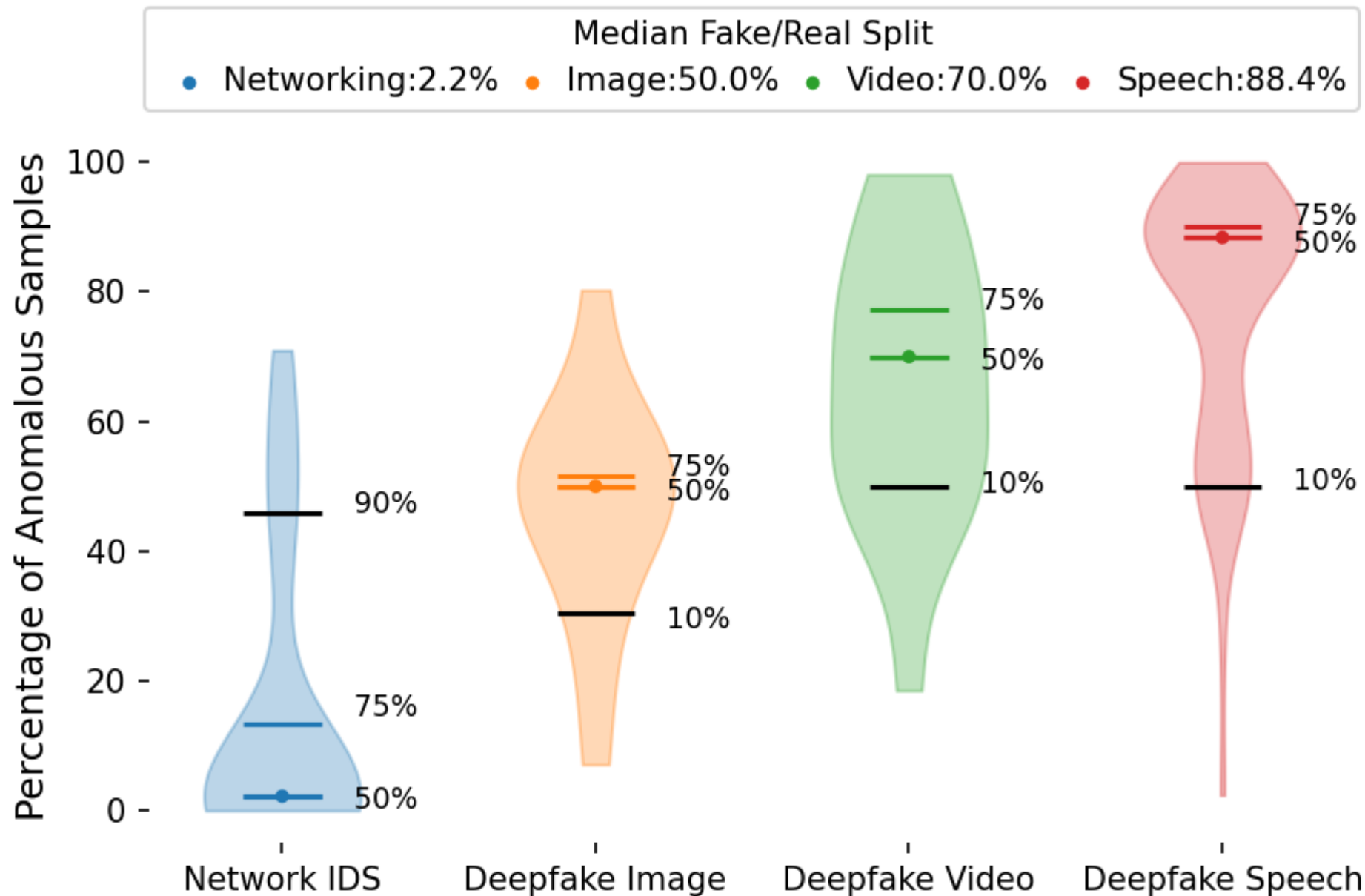
Background

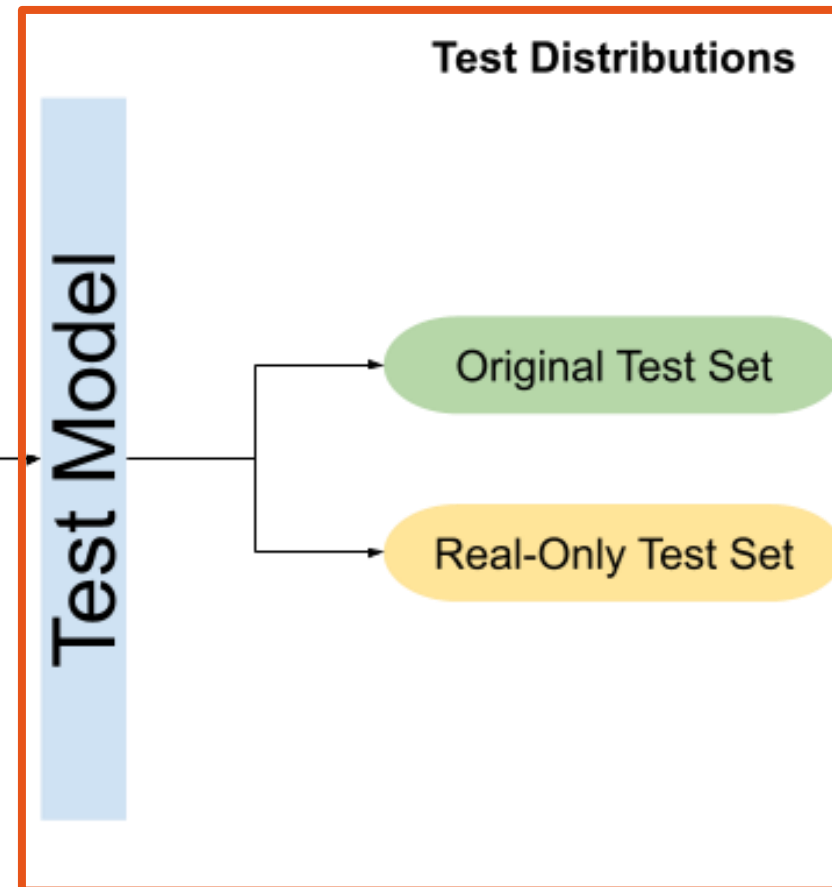
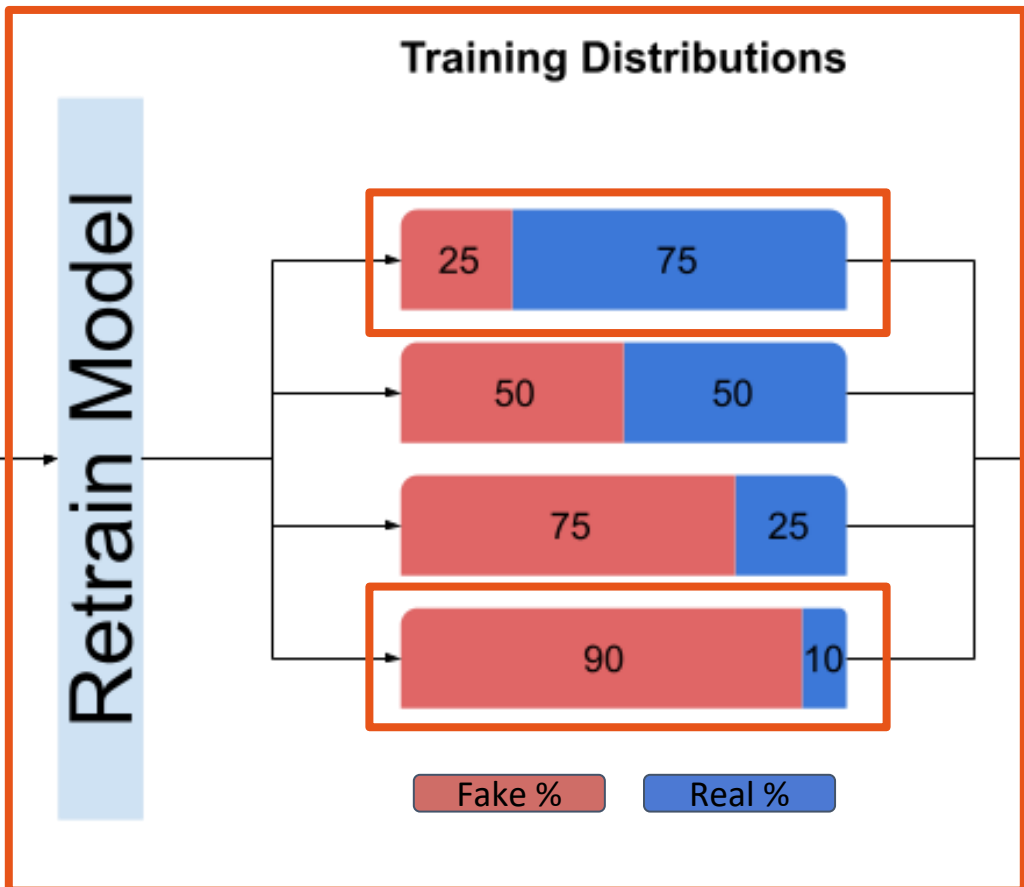
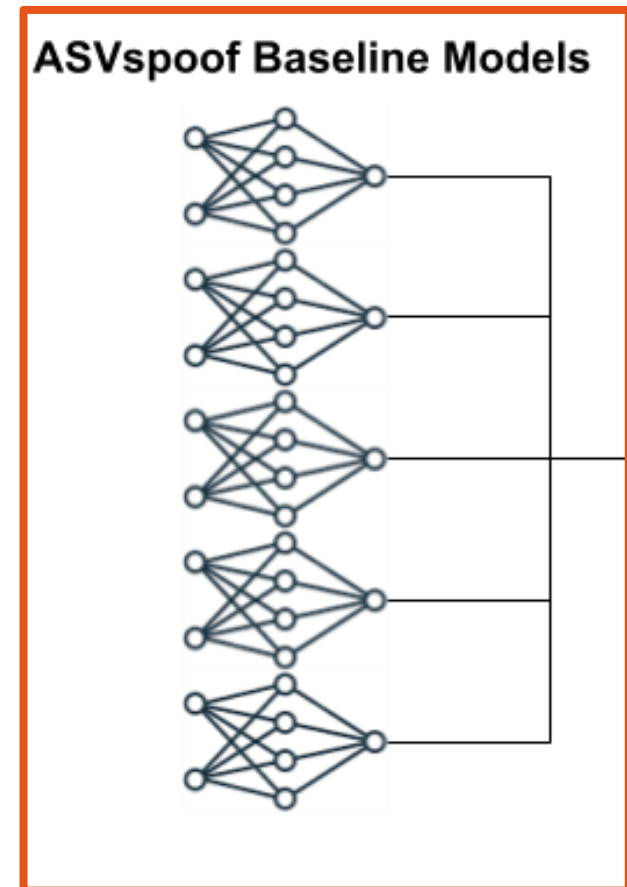
Class Distribution and
Bias

Metric Usage Impacts

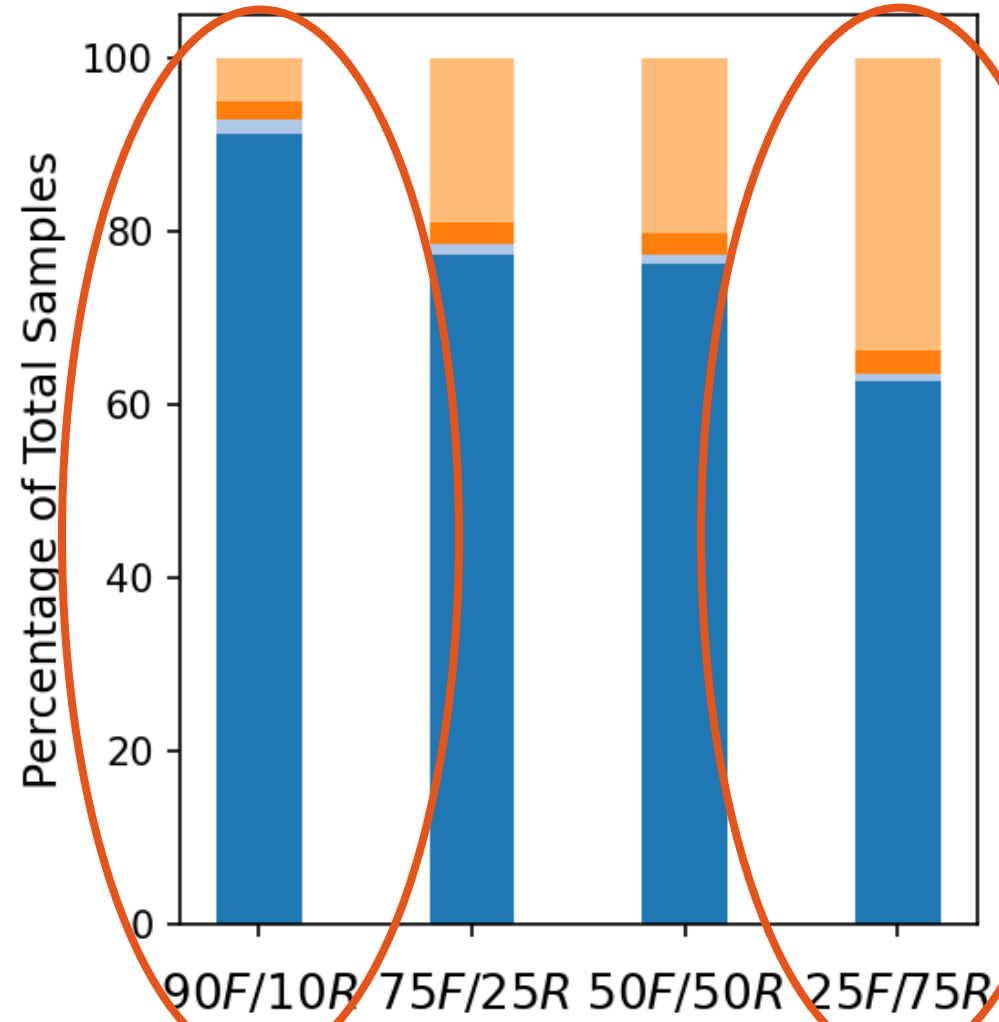
Base-Rate
Contextualization

Suggestions and
Greater Impacts

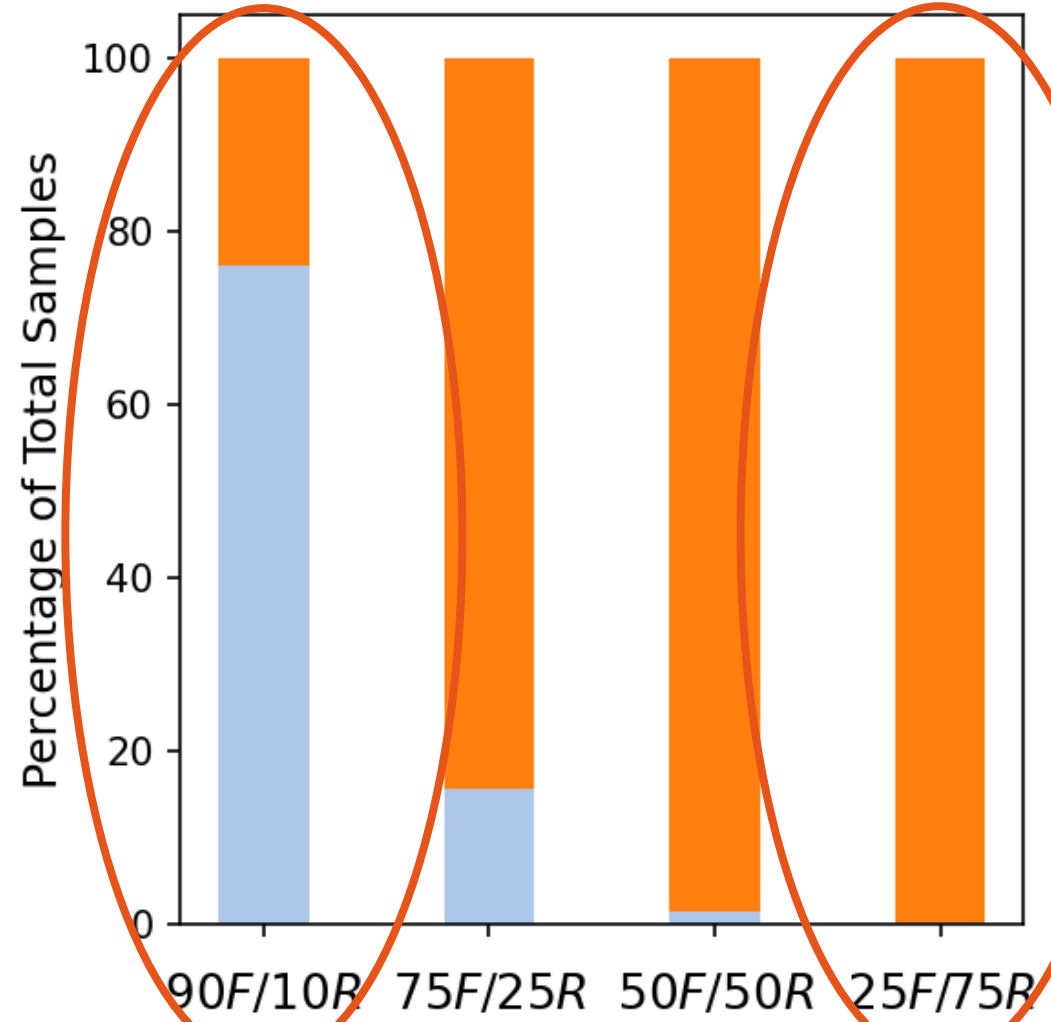




Original Test Set



Real-Only Test Set



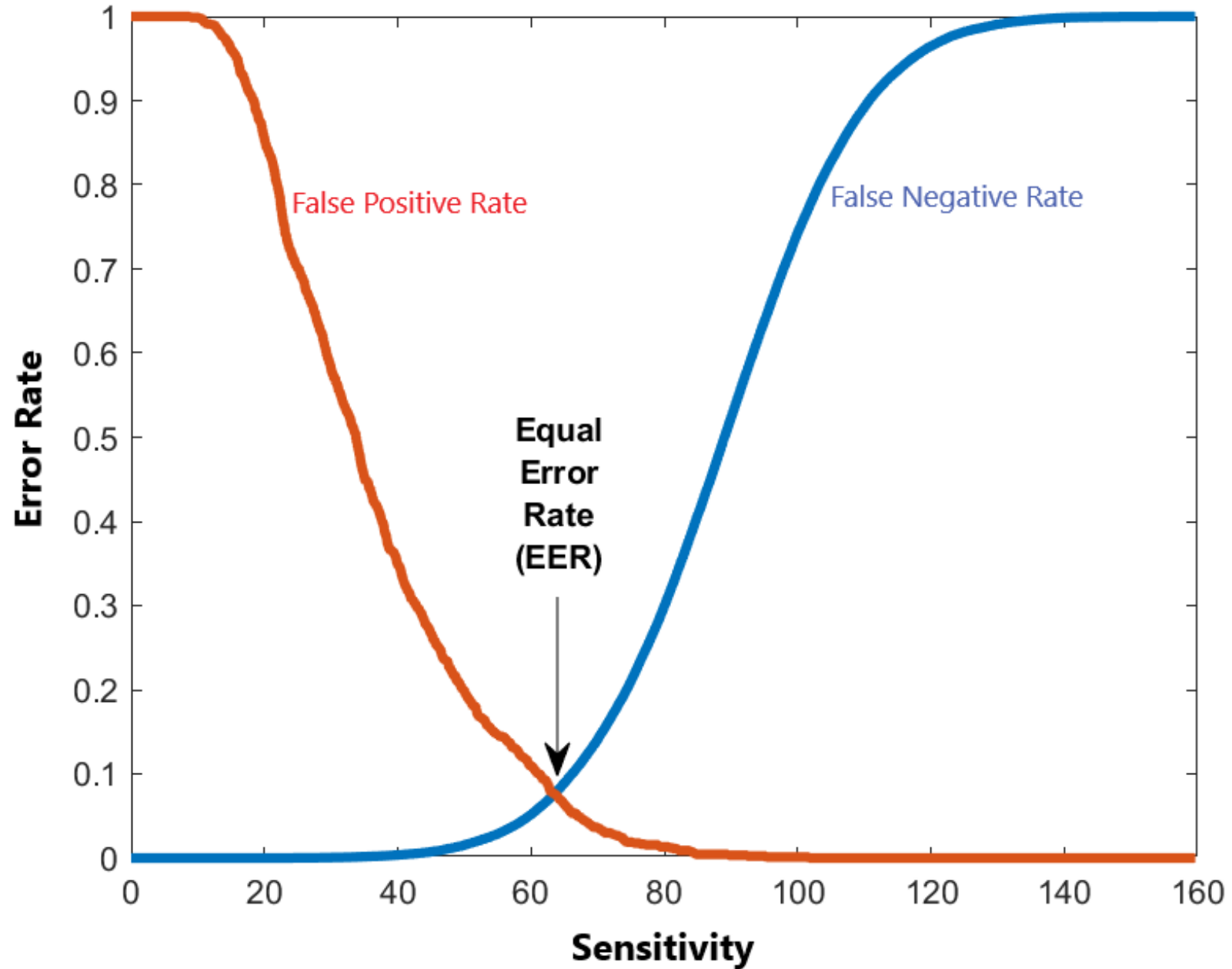
Background

Class Distribution and
Bias

Metric Usage Impacts

Base-Rate
Contextualization

Suggestions and
Greater Impacts



Current Deepfake Detectors claim EERs $< 1\%$

Does that mean this is a “solved” space?

	EER	TPR	FPR
LFCC-GMM	25.5%	44.9%	8.80%
LFCC-LCNN	22.9%	94.7%	41.7%
Absolute Difference	2.6%	49.8%	32.9%

EER Inherently Obscures Results

Background

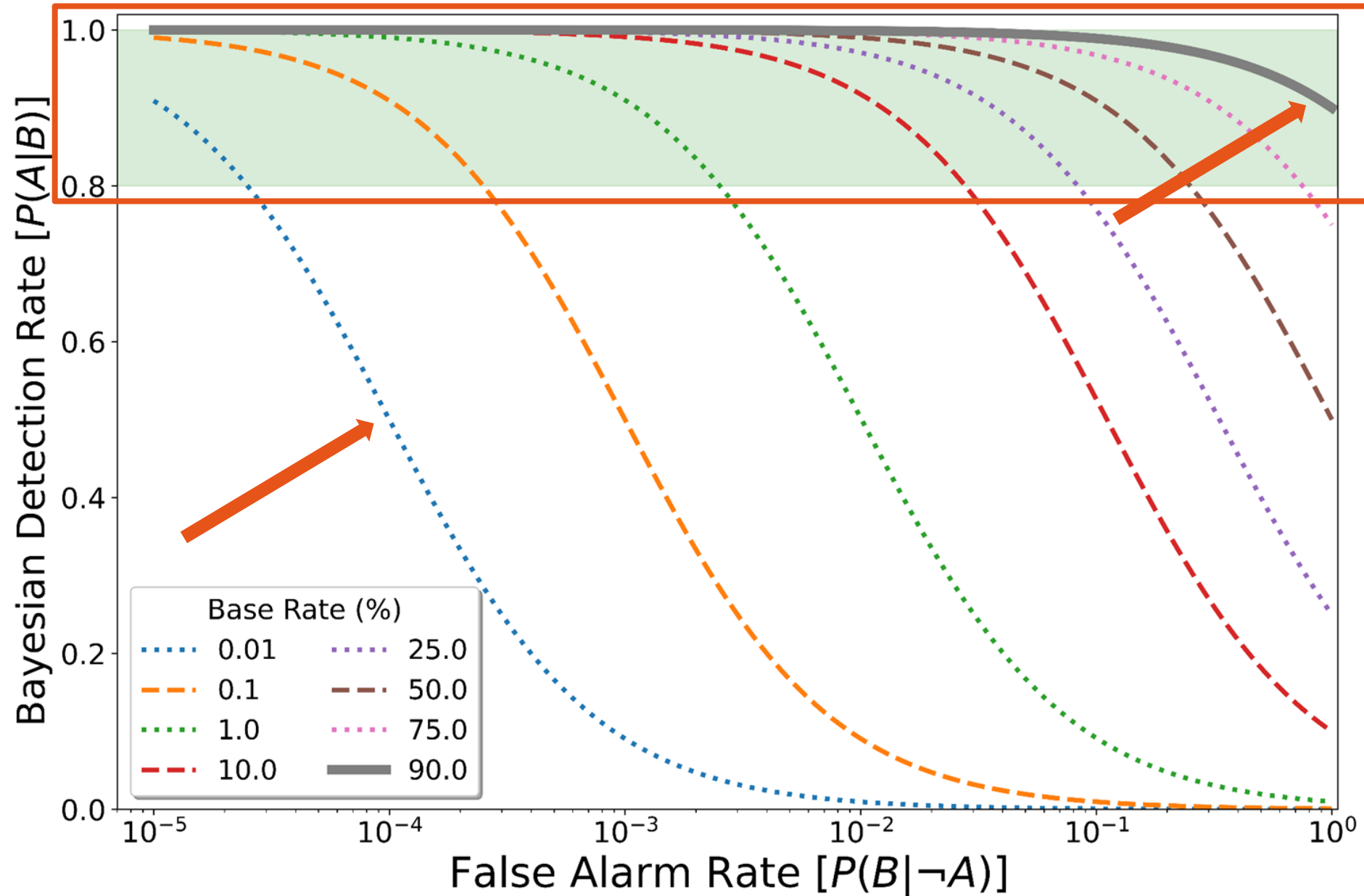
Class Distribution and
Bias

Metric Usage Impacts

Base-Rate
Contextualization

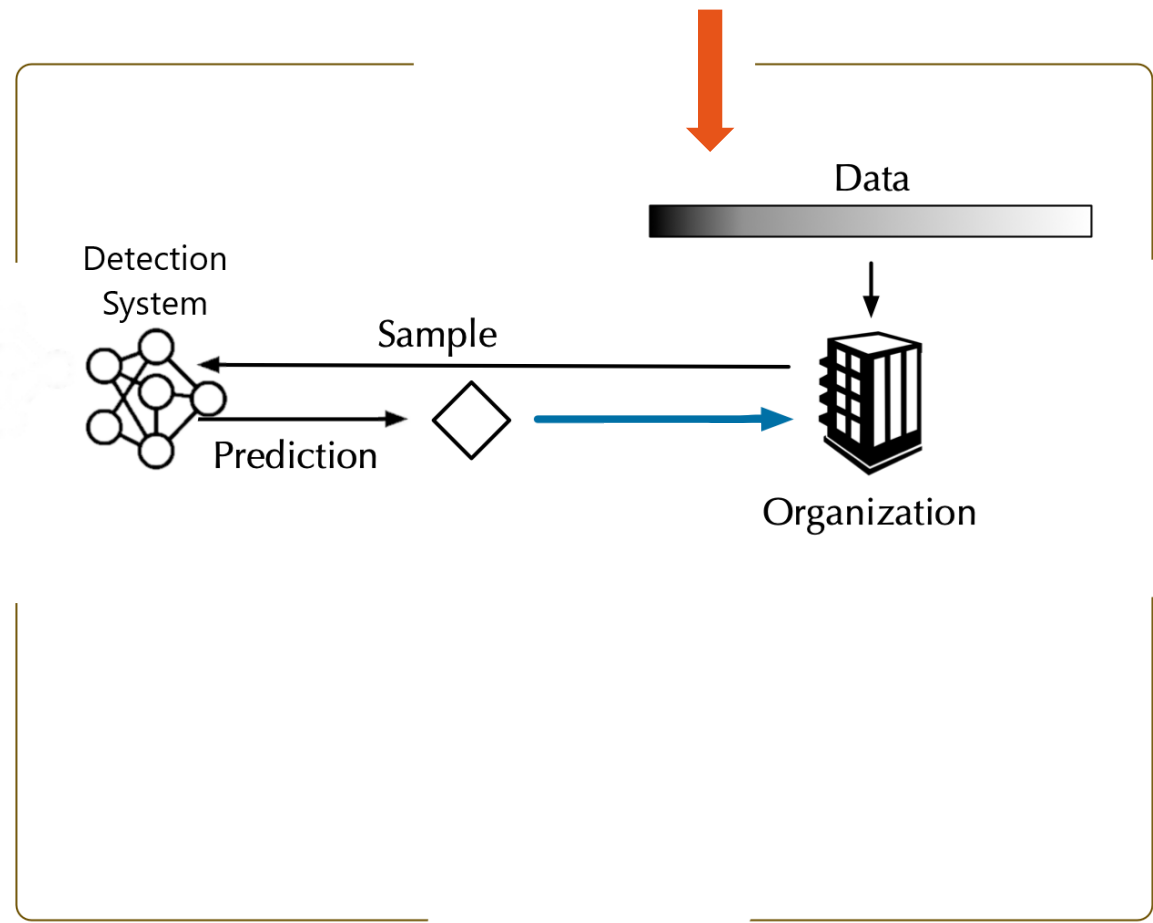
Suggestions and
Greater Impacts

Characterizing Model Efficacy



Call Center Scenario

1/1074 Incoming Calls are Deepfakes



~4,400 Calls Monthly

~4 Deepfake Calls Monthly

	EER	True Positive	False Positive
M_{LG}	25.5%	2	387
M_{SW}	4.14%	4	1182
M_{SW}	4.14%		

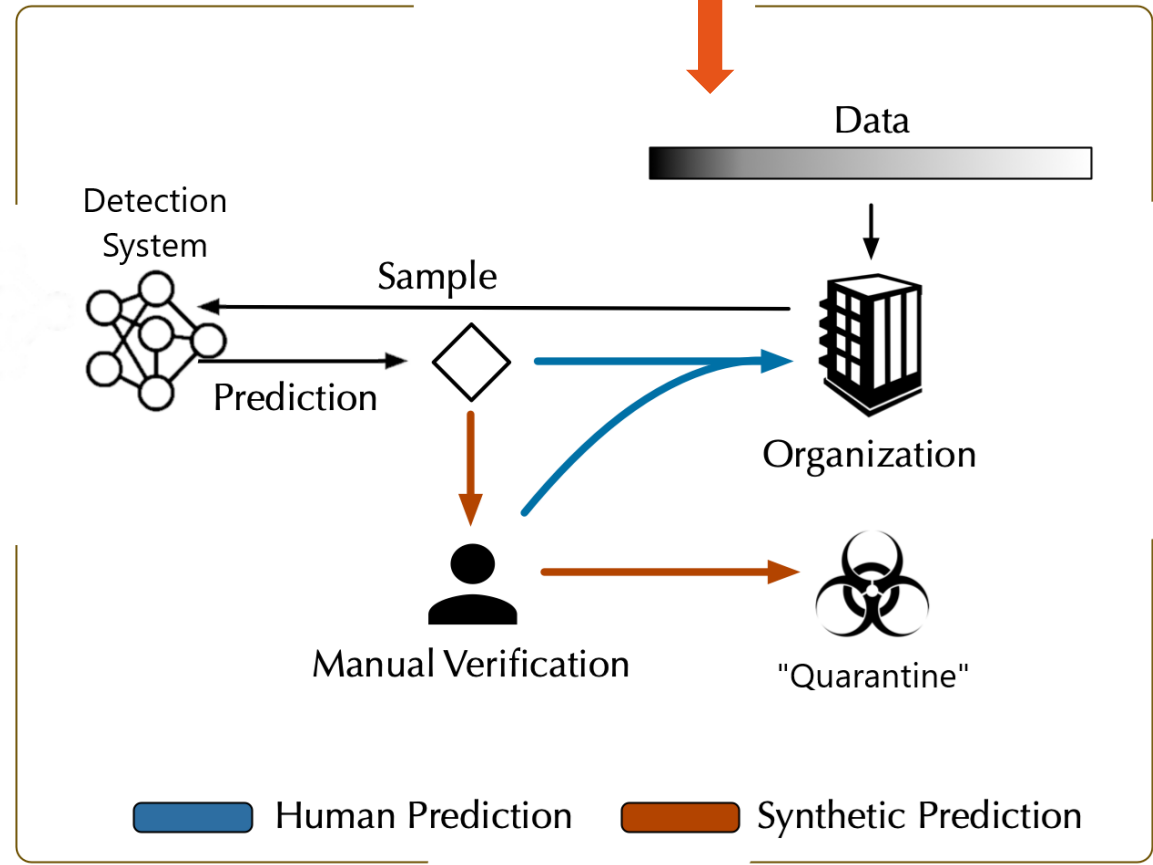
M_{LG} = 1 DF per 200 Alarms

M_{SW} = 1 DF per 333 Alarms

1/11 Incoming Calls are Deepfakes

~4,400 Calls Monthly

~400 Deepfake Calls Monthly



	EER	True Positive	False Positive
M_{LG}	25.5%	180	352
M_{SW}	4.14%	400	1076

$M_{LG} = \sim 1 \text{ DF per } 2 \text{ Alarms}$

$M_{SW} = \sim 2 \text{ DF per } 5 \text{ Alarms}$

Background

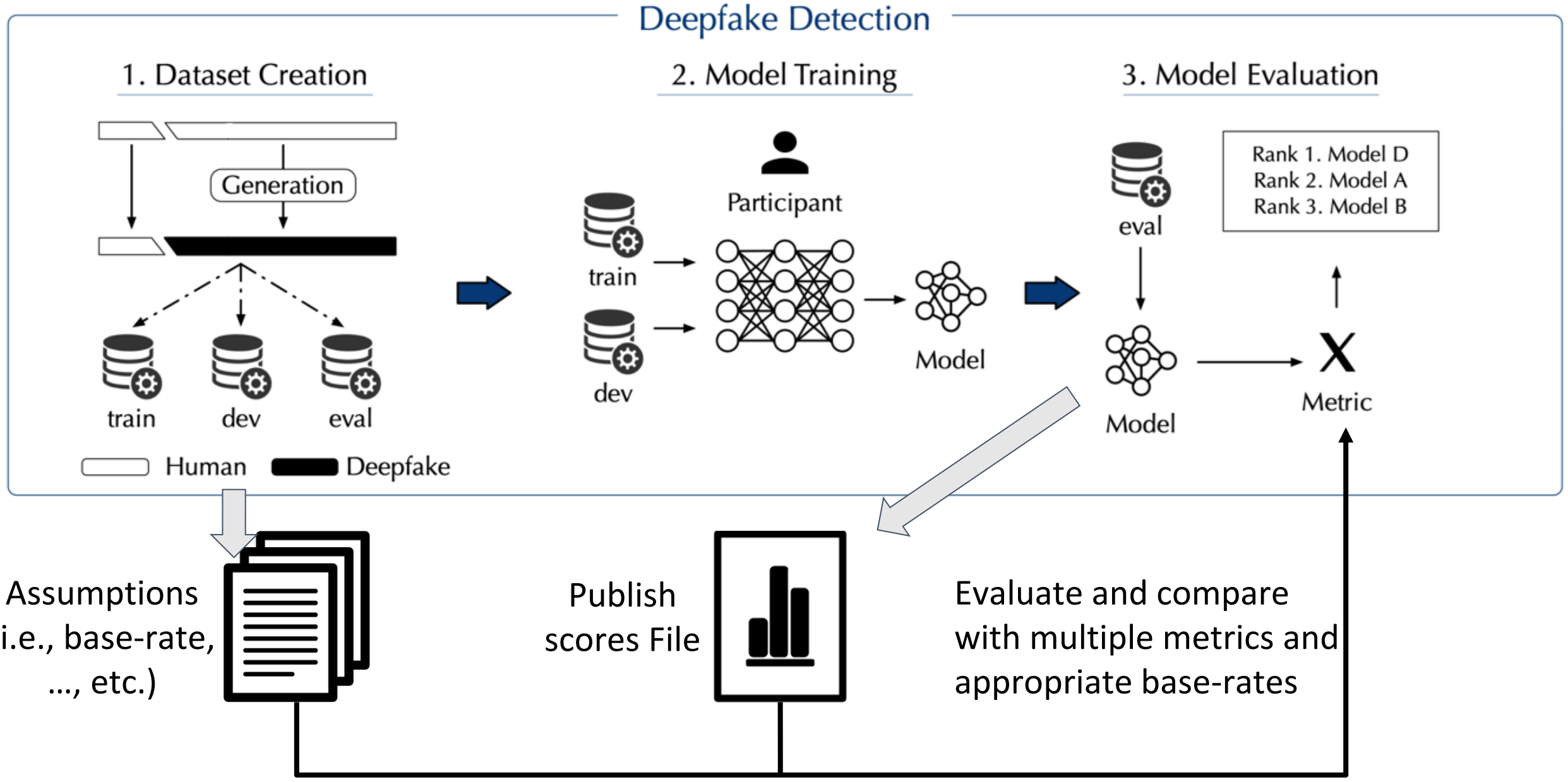
Class Distribution and
Bias

Metric Usage Impacts

Base-Rate
Contextualization

Suggestions and
Greater Impacts

Improving Deepfake Detection



1. This applies to more than deepfakes
2. Honestly characterize model performance
3. Facilitate more research
4. Facilitate more meaningful research

1. Currently used metrics obfuscate results
2. The class distribution in datasets can impose bias on results
3. Contextualization of a dataset is important
4. Current standards of reproducibility and comparability are lacking



sethlayton@ufl.edu