

Supervised Algorithmic Fairness in Distribution Shifts

Tutorial at the 2024 IEEE International Conference on Big data (IEEE BigData 2024)



Chen Zhao¹
chen_zhao@baylor.edu



Xintao Wu²
xintaowu@uark.edu

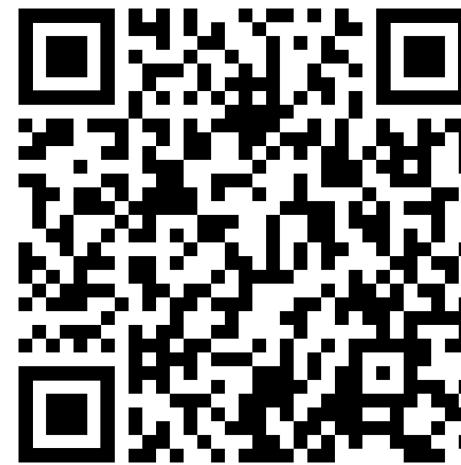
¹Department of Computer Science, Baylor University

²Department of Electrical Engineering and Computer Science, University of Arkansas

Tutorial Website and Survey



**IEEE BigData 2024
Tutorial Website**



**IJCAI 2024
Survey Paper**

Outline

Part I: Introduction to Fairness

Part II: Distribution Shift Undermines Fairness

Part III: Mitigating Unfairness under Distribution Shift (Offline)

Part IV: Mitigating Unfairness under Distribution Shift (Online)

Part V: Open Challenges and Beyond

Why Fairness in Machine Learning is Important?

Insight - Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

October 10, 2018 7:50 PM CDT



The Washington Post
Democracy Dies in Darkness

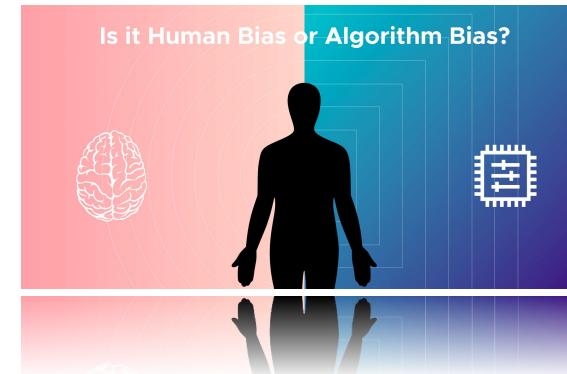
Subscribe for 99¢
every four weeks for the first year

INTERNET CULTURE

Google's algorithm shows prestigious job ads to men, but not to women. Here's why that should worry you.

Medium

TikTok's Addictive and Unethical Algorithm



The New York Times

Art World Takes On A.I. Putting A.I. in Charge A.I. and Hollywood Microsoft-OpenAI Partnership

Facebook Apologizes After A.I. Puts 'Primates' Label on Video of Black Men

Facebook called it "an unacceptable error." The company has struggled with other issues related to race.

BBC News

Home Israel-Gaza war War in Ukraine Climate Video World US & Canada UK Business Tech

Tech

Google apologises for Photos app's racist blunder

1 July 2015

Skyscrapers Airplanes Cars

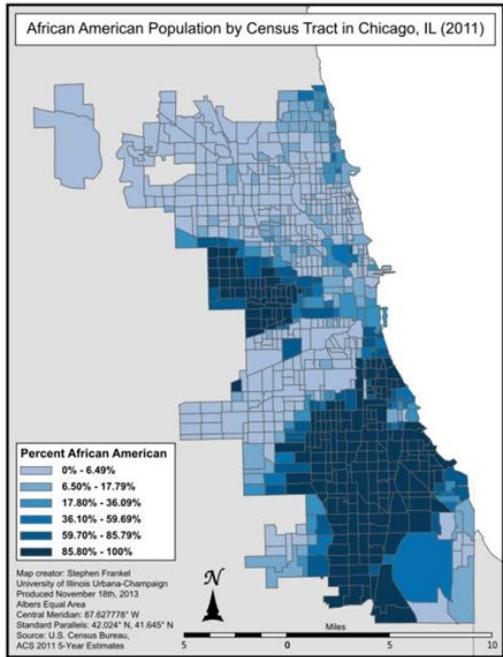
Bikes Gorillas Graduation

diri noir avec banan @jackylaline - Jun 29 Google Photos, y'all... My friend's not a gorilla.

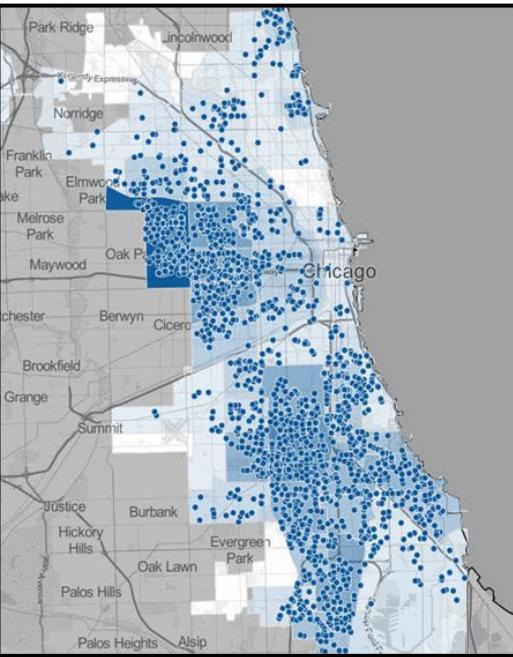
Mr Alcine tweeted Google about the fact its app had misclassified his photo

- <https://www.reuters.com/article/world/insight-amazon-scaps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK0AG/>
- <https://www.nytimes.com/2021/09/03/technology/facebook-ai-race-primates.html>
- <https://www.bbc.com/news/technology-33347866>
- <https://www.washingtonpost.com/news/the-intersect/wp/2015/07/06/googles-algorithm-shows-prestigious-job-ads-to-men-but-not-to-women-heres-why-that-should-worry-you/>
- <https://medium.com/si-410-ethics-and-information-technology/tiktoks-addictive-and-unethical-algorithm-3f44f41ff13c>

Example 1 - Crime Prediction



Black neighborhoods
broken down by community areas



Crime counts prediction

Sensitive Features: Race
Predictions: Crime Counts

NEWS

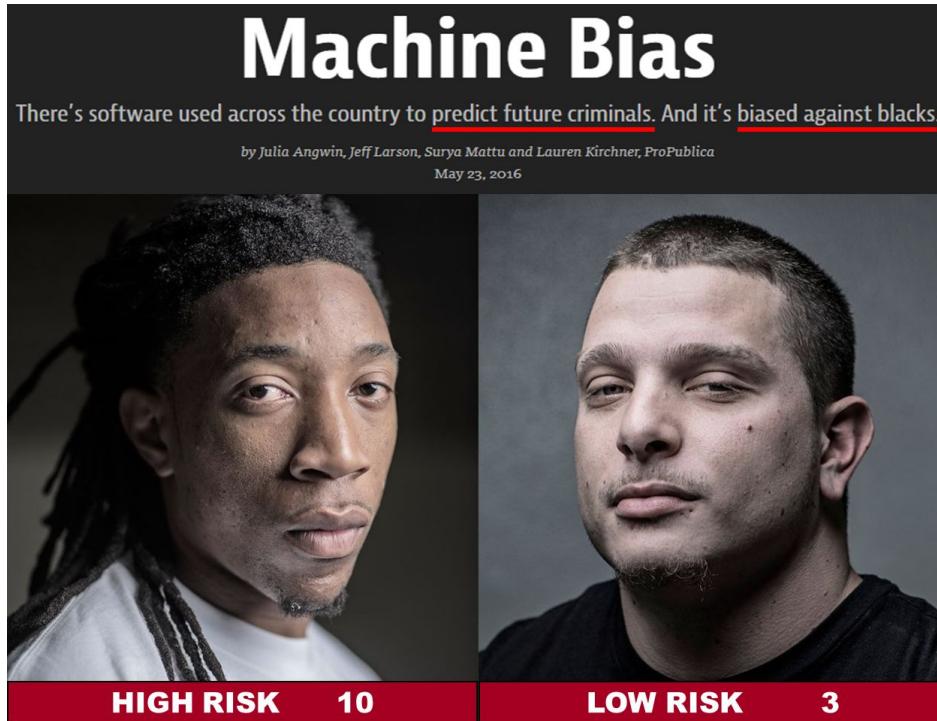
Crime in Chicago: Explore your community



By CHICAGO TRIBUNE
UPDATED: April 1, 2019 at 11:40 AM CST

1. Image source: <https://www.chicagotribune.com/2019/04/01/crime-in-chicago-explore-your-community/>

Example 2 - Recidivism Forecasting



Sensitive Features: Race
Predictions: Risk Level

The black man on the left was predicted as high risk with 1 prior offense and no subsequent offenses.

However, the white man on the right was rated as low risk, even though he was arrested 3 times on drug charges after 1 prior offense.

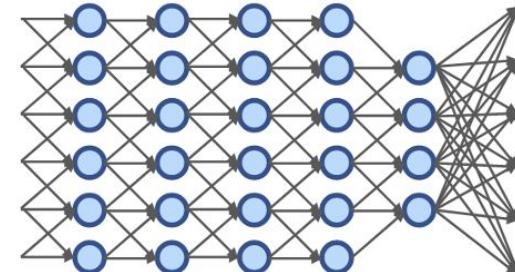
1. Image source: <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

Example 3 - Image Classification

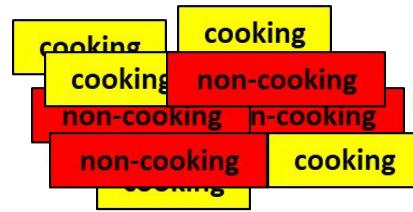
Training images of People



Deep Learning Model $f(x)$



Predictive Outcomes



COOKING	
ROLE	VALUE
AGENT	WOMAN
FOOD	PASTA
HEAT	STOVE
TOOL	SPATULA
PLACE	KITCHEN

COOKING	
ROLE	VALUE
AGENT	WOMAN
FOOD	FRUIT
HEAT	--
TOOL	KNIFE
PLACE	KITCHEN

COOKING	
ROLE	VALUE
AGENT	WOMAN
FOOD	MEAT
HEAT	STOVE
TOOL	SPATULA
PLACE	OUTSIDE

NON-COOKING	
ROLE	VALUE
AGENT	MAN
FOOD	-
HEAT	STOVE
TOOL	SPATULA
PLACE	KITCHEN

NON-COOKING	
ROLE	VALUE
AGENT	MAN
FOOD	-
HEAT	STOVE
TOOL	SPATULA
PLACE	KITCHEN

Sensitive Features: Gender

Predictions: Cooking / non-Cooking

1. Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, Kai-Wei Chang. Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints, EMNLP, 2017

Algorithmic Fairness



X: images

Z: {Male, Female}

Y: {Cooking, non-Cooking}

For example:

Source Domain: All Male images are non-Cooking and Female are Cooking.

Target Domain: How to ensure a classifier that predicts image labels independent of gender?

In fairness-aware machine learning, we mitigate the dependence:

Sensitive Features (Z)
e.g., Male / Female

spuriously correlated to

Predictive Outcomes ($\hat{Y} = f_{\theta}(X)$)
e.g., Cooking / non-Cooking

Algorithmic Fairness:

The spurious dependence between sensitive features to predictive model outcomes is eliminated or restricted to an accepted fairness level.

Definition of Fairness

1 Group Fairness

Group fairness ensures that predictions are equitable across different demographic groups defined by sensitive attributes such as race, gender, or age.

2

Individual Fairness

Individual fairness ensures that similar individuals are treated similarly by the model, regardless of group membership.

3

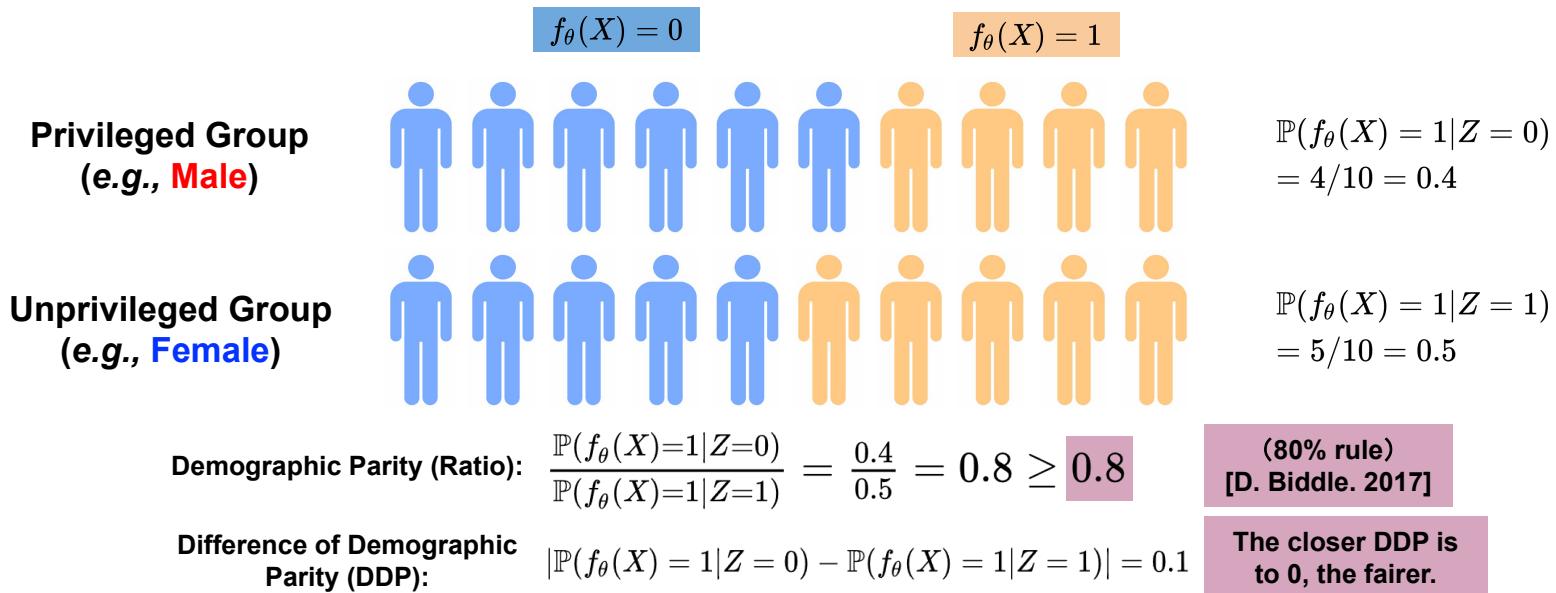
Counterfactual Fairness

Counterfactual fairness evaluates fairness by comparing outcomes across hypothetical "counterfactual" worlds where only the sensitive attribute is altered while keeping everything else the same.

Fairness

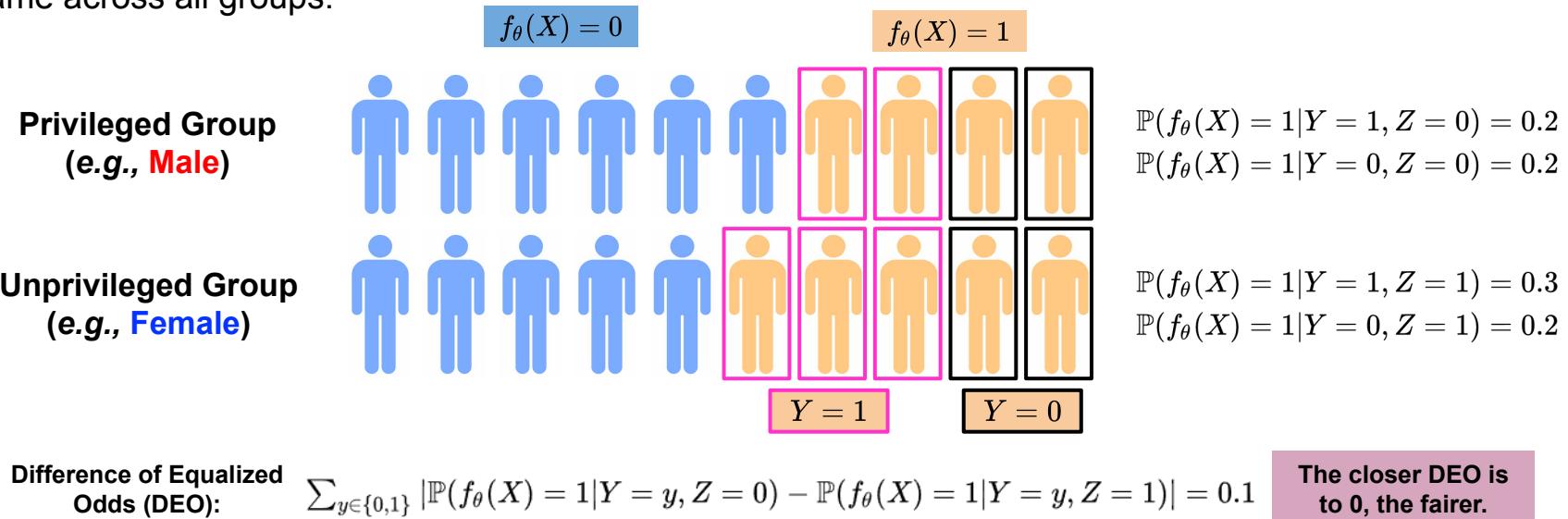
Group Fairness (Demographic Parity)

Demographic Parity (DP) ensures a model's predictions are independent of sensitive attributes. It requires that individuals from different demographic groups have an equal probability of receiving a positive outcome ($\hat{Y} = f_\theta(X) = 1$).



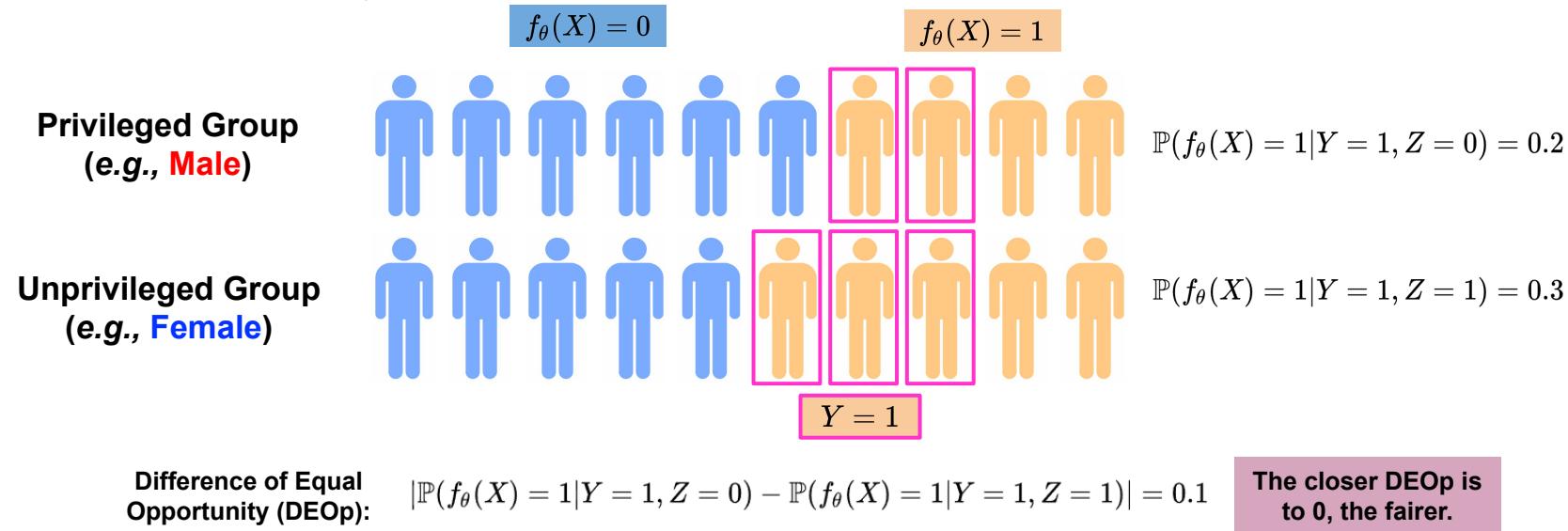
Group Fairness (Equalized Odds)

Equalized odds (EO) ensures a model's predictions are equally accurate across different demographic groups. It focuses on parity of error rates between sensitive groups, for a binary classification task. Specifically, it requires that the model's true positive rate (TPR) and false positive rate (FPR) be the same across all groups.



Group Fairness (Equal Opportunity)

Equal Opportunity (EOp) focuses on ensuring equal True Positive Rates (TPR) across different groups defined by a sensitive attribute. This means that individuals who belong to different sensitive groups but have the same actual positive outcome should have an equal chance of being correctly classified as positive by the model.



Linear Relaxation of DDP and DEOp

DDP (Difference of Demographic Parity) and **DEOp (Difference of Equal Opportunity)** are fairness metrics that measure disparities in outcomes or opportunities between sensitive groups. When addressing these fairness criteria in optimization problems, exact fairness constraints are often non-convex or difficult to handle directly. **Linear relaxation** refers to approximating these constraints in a linear form to make them tractable in optimization settings.

$$\text{DDP: } \left| \frac{1}{n_0} \sum_{(\mathbf{x}, z, y) \in \mathcal{D}_0} \mathbb{I}_{f_\theta(\mathbf{x}) > 0} - \frac{1}{n_1} \sum_{(\mathbf{x}, z, y) \in \mathcal{D}_1} \mathbb{I}_{f_\theta(\mathbf{x}) > 0} \right| \leq \tau$$
$$\text{DEOp: } \left| \frac{1}{n_0} \sum_{(\mathbf{x}, z, y) \in \mathcal{D}_0} \mathbb{I}_{f_\theta(\mathbf{x}) > 0, Y=1} - \frac{1}{n_1} \sum_{(\mathbf{x}, z, y) \in \mathcal{D}_1} \mathbb{I}_{f_\theta(\mathbf{x}) > 0, Y=1} \right| \leq \tau$$

Re-write
Linear Relaxation

$$\mathbb{E}_{(\mathbf{x}, z, y) \sim \mathcal{D}} \left[\frac{1}{\hat{p}_1(1-\hat{p}_1)} \left(\frac{z+1}{2} - \hat{p}_1 \right) \mathbb{I}_{f_\theta(\mathbf{x}) > 0} \right] \leq \tau$$

$$\left| \frac{1}{n} \sum_{(\mathbf{x}, z, y) \in \mathcal{D}} \frac{1}{\hat{p}_1(1-\hat{p}_1)} \left(\frac{z+1}{2} - \hat{p}_1 \right) f_\theta(\mathbf{x}) \right| \geq \tau$$

For DDP:
 $\hat{p}_1 = \mathbb{P}_{(\mathbf{x}, z, y) \sim \mathcal{D}}(z = 1)$

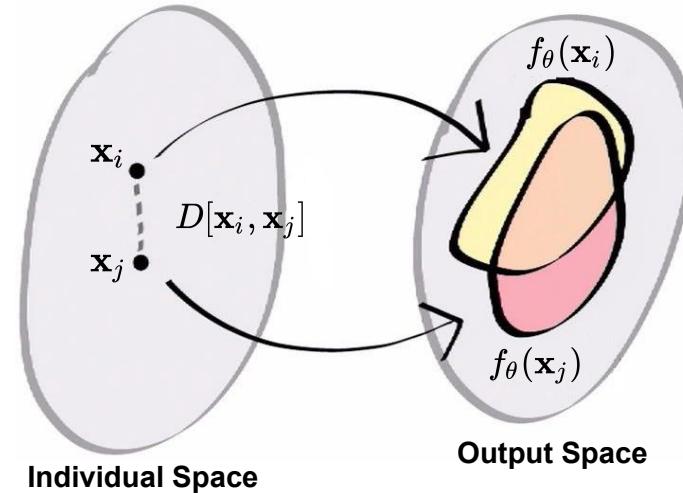
For DEOp:
 $\hat{p}_1 = \mathbb{P}_{(\mathbf{x}, z, y) \sim \mathcal{D}}(y = 1, z = 1)$

1. Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, Krishna P. Gummadi. Fairness Beyond Disparate Treatment & Disparate Impact: Learning Classification without Disparate Mistreatment. WWW, 2017.
2. Michael Lohaus, Michael Perrot, Ulrike Von Luxburg. Too Relaxed to Be Fair. ICML, 2020.

Individual Fairness

Individual fairness emphasizes **treating similar individuals similarly**. It ensures that two individuals who are similar in terms of relevant features receive similar predictions or outcomes from the model.

Search Query	Work Experience	Education Experience	Candidate	Xing Ranking
Brand Strategist	146	57	male	1
Brand Strategist	327	0	female	2
Brand Strategist	502	74	male	3
Brand Strategist	444	56	female	4
Brand Strategist	139	25	male	5
Brand Strategist	110	65	female	6
Brand Strategist	12	73	male	7
Brand Strategist	99	41	male	8
Brand Strategist	42	51	female	9
Brand Strategist	220	102	female	10
...				
Brand Strategist	3	107	female	20
Brand Strategist	123	56	female	30
Brand Strategist	3	3	male	40



$$d[f_\theta(\mathbf{x}_i), f_\theta(\mathbf{x}_j)] \leq \delta \cdot D[\mathbf{x}_i, \mathbf{x}_j], \forall i \neq j$$

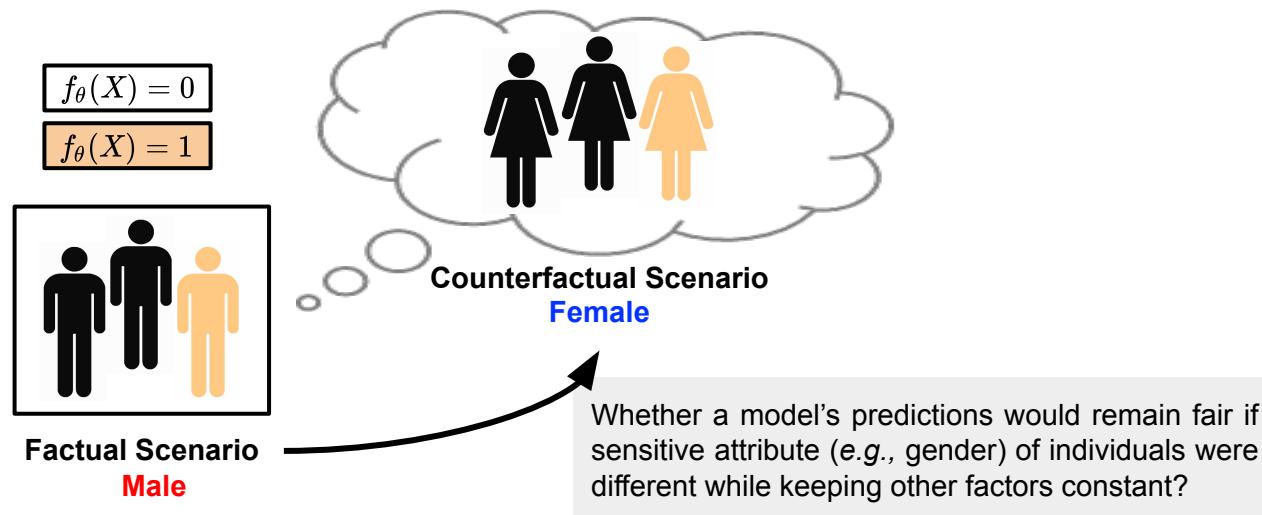
The top-10 results satisfy group fairness with regard to gender. However the outcomes are far from being fair for the individual users: people with very similar qualifications, such as *Work Experience* and *Education Score* ended up on ranks that are far apart (e.g., ranks 5 and 30).

1. Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, Rich Zemel. Fairness Through Awareness. ITCS, 2012.

2. Preethi Lahoti, Krishna P. Gummadi, Gerhard Weikum. iFair: Learning Individually Fair Data Representations for Algorithmic Decision Making. ICDE, 2019.

Counterfactual Fairness

Counterfactual Fairness ensures a model's predictions for an individual are not influenced by changes to a sensitive attribute in a counterfactual scenario. This concept originates from causal inference and is designed to capture how decisions would change (or not) if a person's sensitive attribute were different while keeping all other aspects of their situation the same.



Counterfactual Fairness

The **Total Causal Effect (TCE)** measures the overall influence of sensitive attributes on the model's prediction.

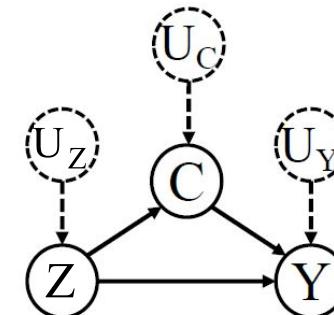
Given context $O = \mathbf{o}$, the **Counterfactual Effect (CE)** measures the influence of sensitive attributes on the model prediction specified on the sub-population by the context only.

These two metrics quantify how much the prediction would change if the sensitive attribute were altered, while allowing all other variables that depend on it to adjust accordingly through the causal pathways.

$$\text{TCE} = |\mathbb{P}(f_{\theta, Z \leftarrow 0}(X) = y) - \mathbb{P}(f_{\theta, Z \leftarrow 1}(X) = y)|, \forall y \quad \text{CE} = |\mathbb{P}(f_{\theta, Z \leftarrow 0}(X) = y | O = \mathbf{o}) - \mathbb{P}(f_{\theta, Z \leftarrow 1}(X) = y | O = \mathbf{o})|, \forall y$$

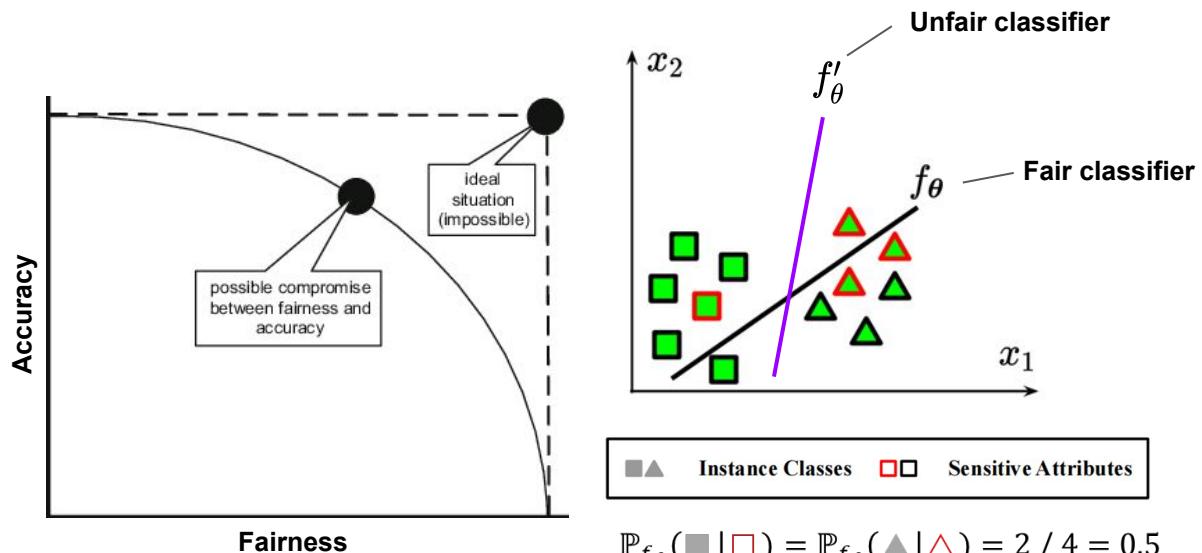
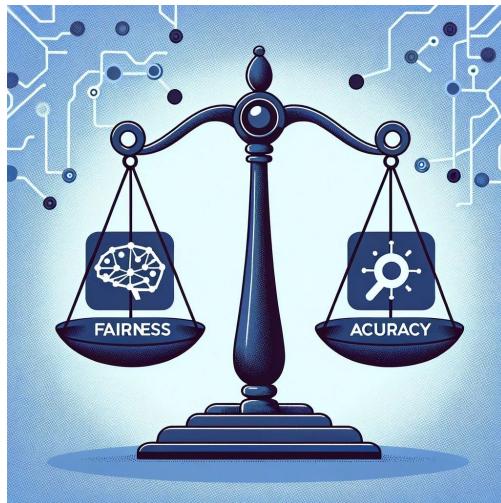
Table 1: Connection between previous fairness notions and PC fairness

Description	Relating to PC fairness
Total effect	$O = \emptyset$ and $\pi = \Pi$
(System) Direct discrimination	$O = \emptyset$ or $\{Z\}$ and $\pi = \pi_d = \{Z \rightarrow \hat{Y}\}$
(System) Indirect discrimination	$O = \emptyset$ or $\{Z\}$ and $\pi = \pi_i \subset \Pi$
Individual direct discrimination	$O = \{Z, X\}$ and $\pi = \pi_d = \{Z \rightarrow \hat{Y}\}$
Group direct discrimination	$O = Q = PA_Y \setminus \{Z\}$ and $\pi = \pi_d = \{Z \rightarrow \hat{Y}\}$
Counterfactual fairness	$O = \{Z, X\}$ and $\pi = \Pi$
Counterfactual error rate	$O = \{Z, Y\}$ and $\pi = \pi_d$ or π_i



1. Judea Pearl. Causality. Cambridge university press, 2009.
2. Ilya Shpitser and Judea Pearl. Complete identification methods for the causal hierarchy. JMLR, 2008.
3. Yongkai Wu, Lu Zhang, Xintao Wu, Hanghang Tong. PC-Fairness: A Unified Framework for Measuring Causality-based Fairness. NeurIPS, 2019.

Fairness-Accuracy Tradeoff



$$\mathbb{P}(f_\theta(X) = 1 | Z = 1) = \mathbb{P}(f_\theta(X) = 1 | Z = 0)$$

1. Image source: https://medium.com/@afiori_78621/the-fairness-accuracy-tradeoff-af71d5d0c38a

2. Wil van der Aalst. Responsible Data Science: Using Event Data in a “People Friendly” Manner. ICEIS, 2017.

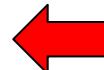
3. Minglai Shao, Dong Li, Chen Zhao, Xintao Wu, Yujie Lin, Qin Tian. Supervised Algorithmic Fairness in Distribution Shifts: A Survey. IJCAI 2024.

Outline

Part I: Introduction to Fairness



Part II: Distribution Shift Undermines Fairness

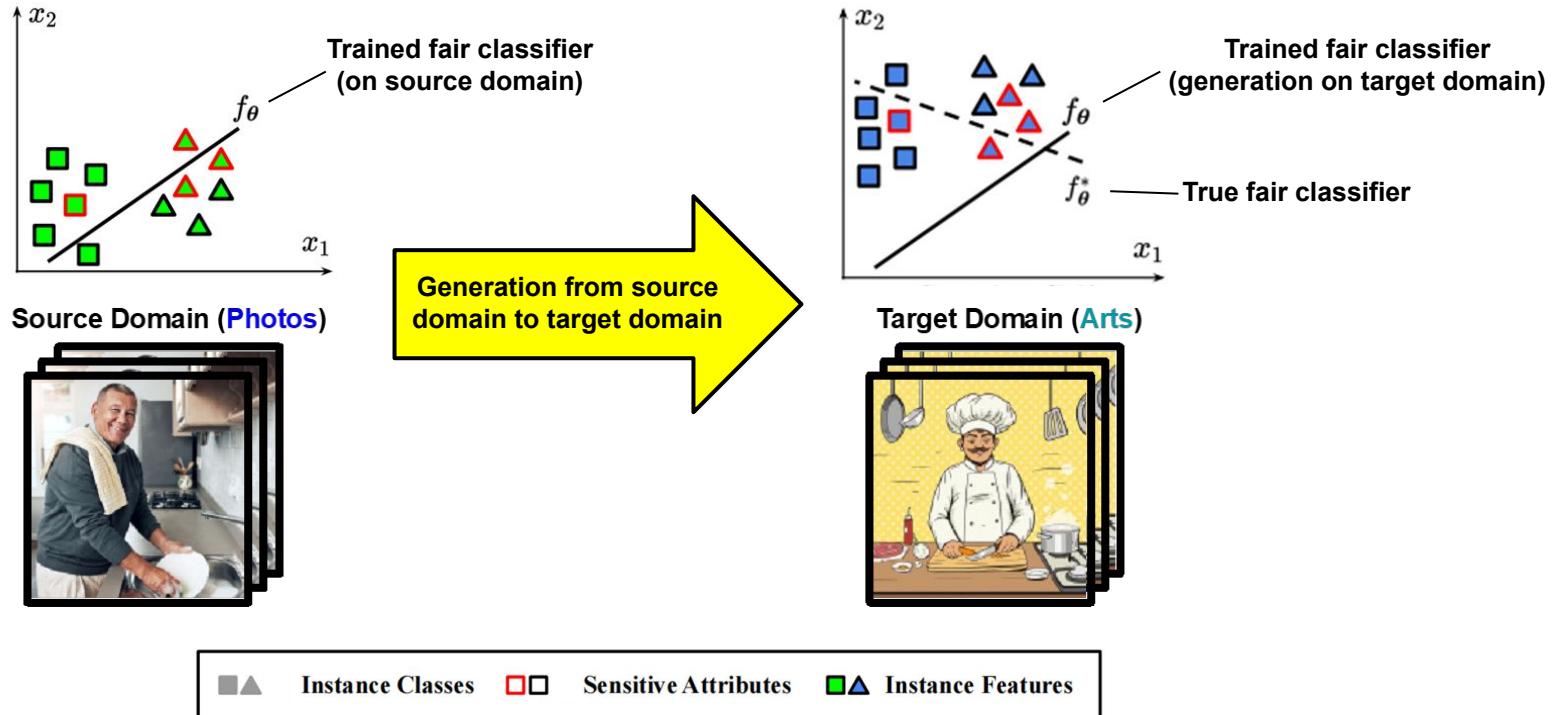


Part III: Mitigating Unfairness under Distribution Shift (Offline)

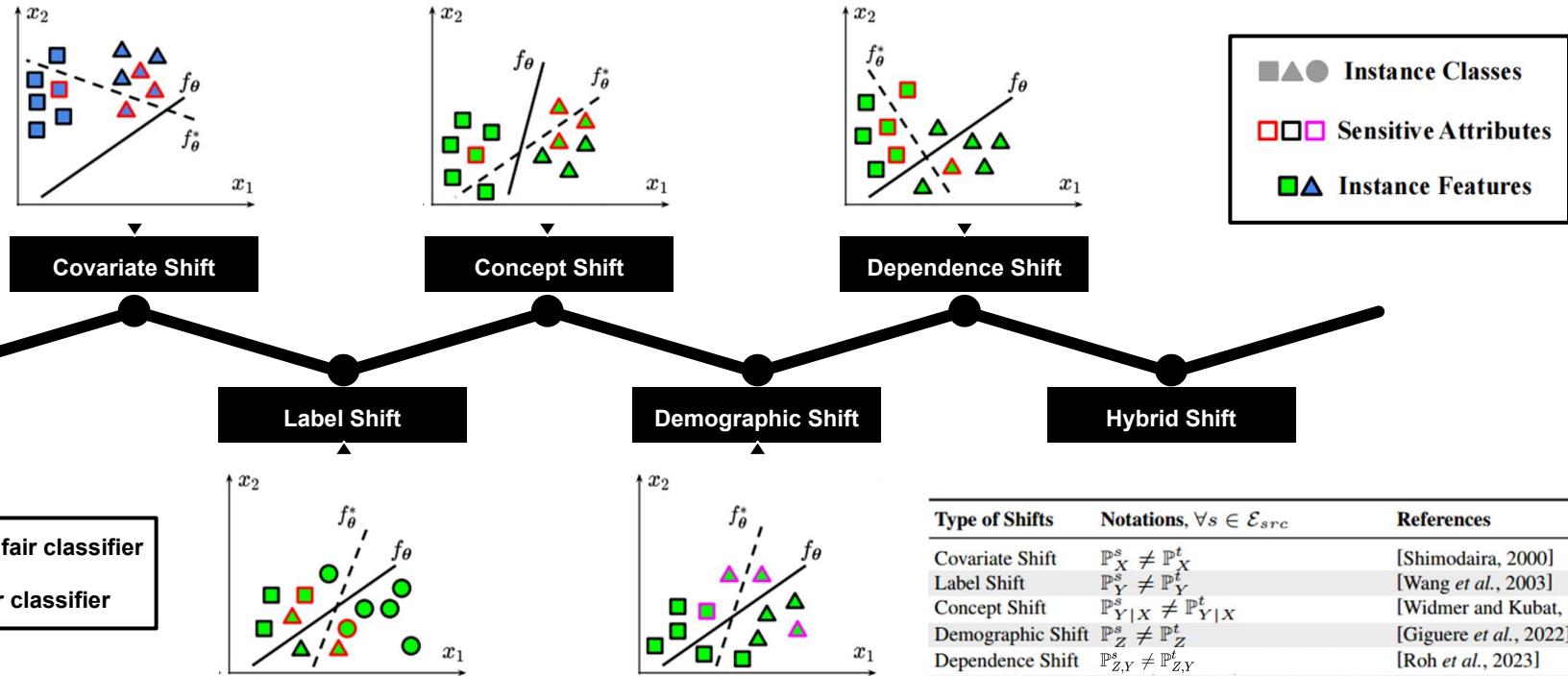
Part IV: Mitigating Unfairness under Distribution Shift (Online)

Part V: Open Challenges and Beyond

Distribution Shift Undermines Fairness



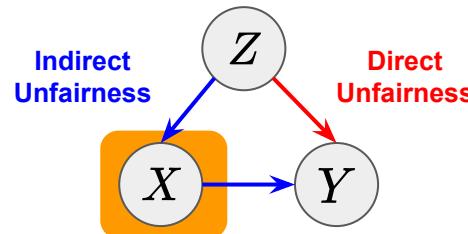
Five Types of Distribution Shifts in Model Fairness



1. Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. JSPI, 2000.
2. Ke Wang, Senqiang Zhou, Chee Ada Fu, and Jeffrey Xu Yu. Mining changes of classification by correspondence tracing. SDM, 2003.
3. Gerhard Widmer and Miroslav Kubat. Learning in the presence of concept drift and hidden contexts. Machine learning, 1996.
4. Stephen Giguere, Blossom Metivier, Yuriy Brun, Bruno Castro da Silva, Philip S Thomas, and Scott Niekum. Fairness guarantees under demographic shift. ICLR, 2022.
5. Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. Improving fair training under correlation shifts. ICML, 2023

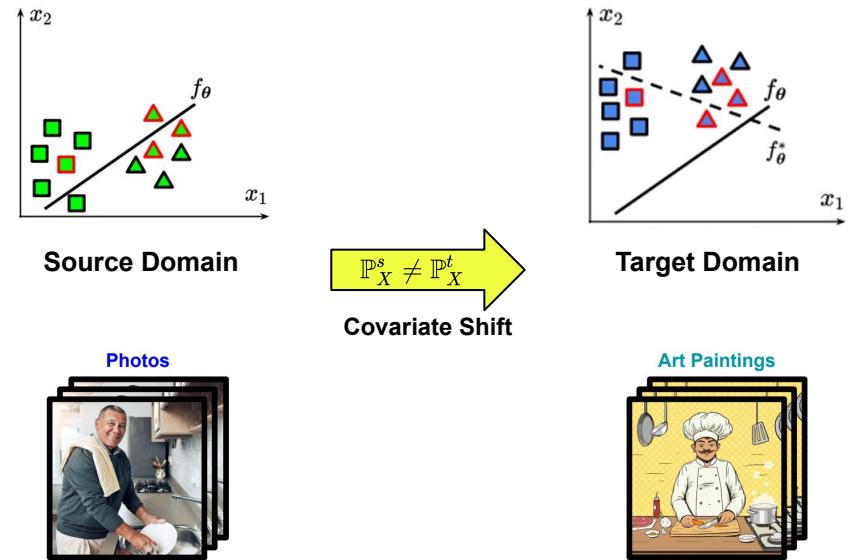
Covariate Shift

When the **Covariate Shift** occurs, the model's assumptions about the relationships between X and Y may no longer hold, leading to biased predictions.



■▲● Instance Classes
□□□ Sensitive Attributes
■▲ Instance Features

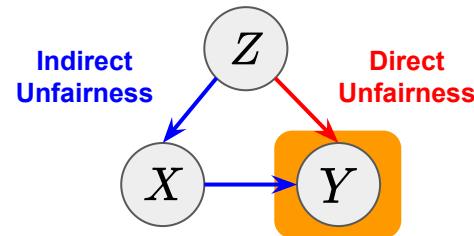
f_θ Trained fair classifier
 f_θ^* True fair classifier



1. Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. JSPI, 2000.
2. Minglai Shao, Dong Li, Chen Zhao, Xintao Wu, Yujie Lin, Qin Tian. Supervised Algorithmic Fairness in Distribution Shifts: A Survey. IJCAI 2024.

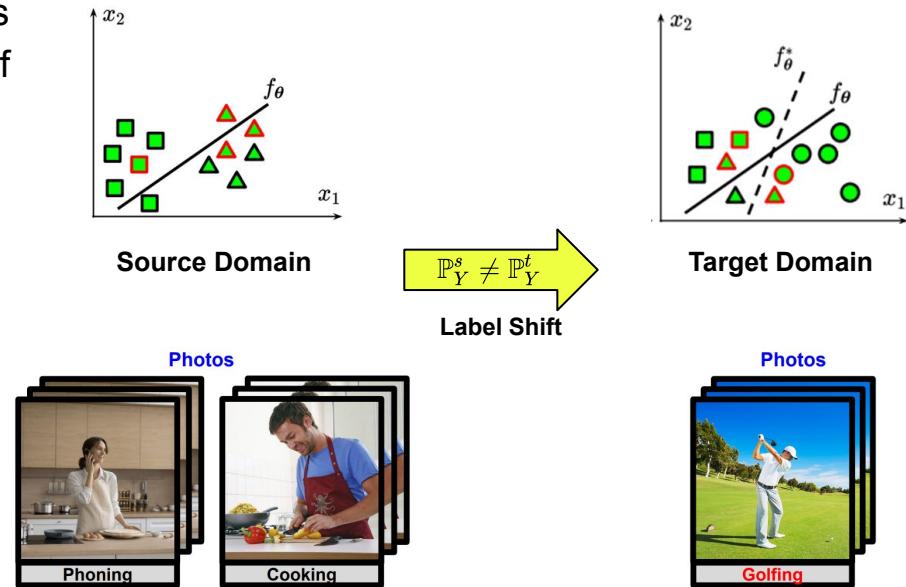
Label Shift

The challenge of fairness learning under **Label Shift** primarily involves mitigating predictive bias in studies related to outlier detection or out of distribution (OOD) detection.



	Instance Classes
	Sensitive Attributes
	Instance Features

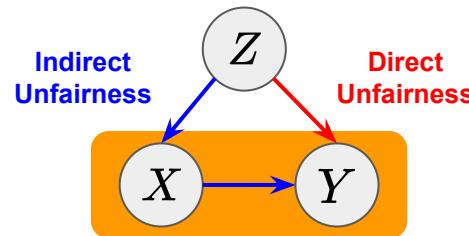
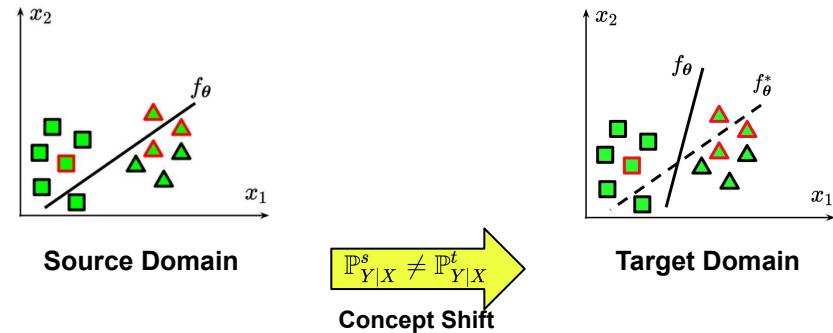
f_θ Trained fair classifier
 f_θ^* True fair classifier



1. Ke Wang, Senqiang Zhou, Chee Ada Fu, and Jeffrey Xu Yu. Mining changes of classification by correspondence tracing. SDM, 2003.
2. Minglai Shao, Dong Li, Chen Zhao, Xintao Wu, Yujie Lin, Qin Tian. Supervised Algorithmic Fairness in Distribution Shifts: A Survey. IJCAI 2024.

Concept Shift

Concept Shift occurs when the underlying concept defining the labels evolves, leading to differences in how the input features correspond to the output labels. It specifically refers to changes in $P(Y|X)$. When this relationship shifts, models trained on historical data may not only perform poorly but also produce predictions that exacerbate biases or unfair outcomes.



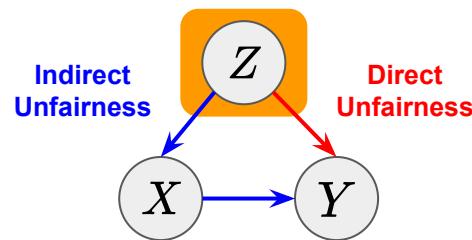
■▲●	Instance Classes
□□□	Sensitive Attributes
■▲△	Instance Features

f_θ Trained fair classifier
 f_θ^* True fair classifier

1. Gerhard Widmer and Miroslav Kubat. Learning in the presence of concept drift and hidden contexts. Machine learning, 1996.
2. Minglai Shao, Dong Li, Chen Zhao, Xintao Wu, Yujie Lin, Qin Tian. Supervised Algorithmic Fairness in Distribution Shifts: A Survey. IJCAI 2024.

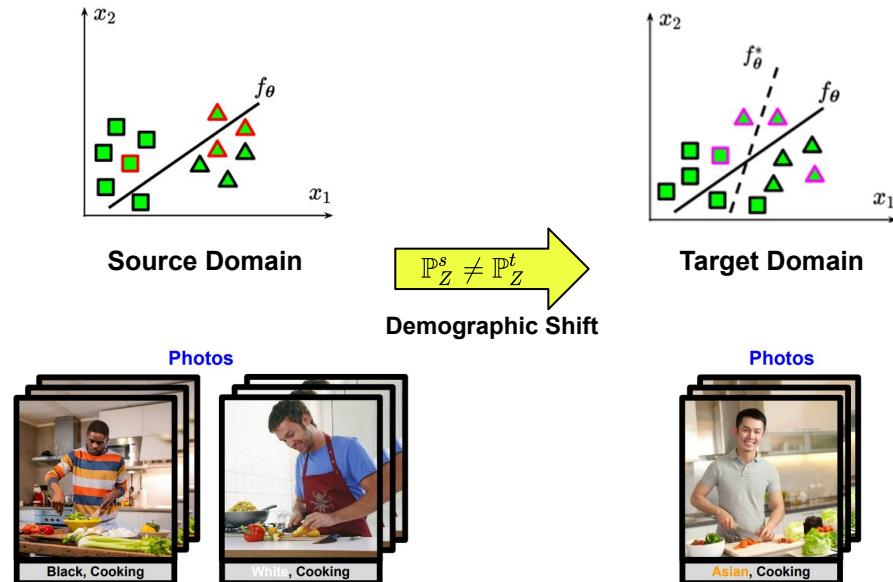
Demographic Shift

A **Demographic Shift** occurs when a specific sensitive subgroup of the population becomes more or less probable in the target domain.



■▲●	Instance Classes
□□□	Sensitive Attributes
■▲	Instance Features

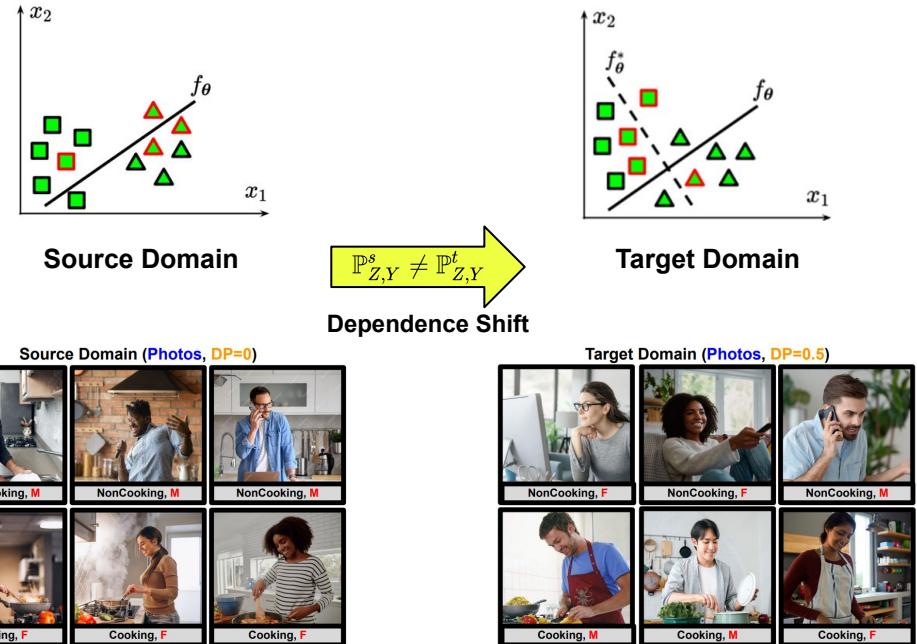
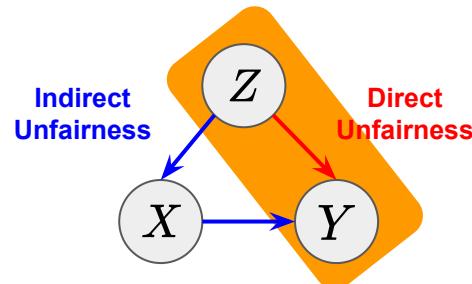
f_θ Trained fair classifier
 f_θ^* True fair classifier



1. Stephen Giguere, Blossom Metevier, Yuriy Brun, Bruno Castro da Silva, Philip S Thomas, and Scott Niekum. Fairness guarantees under demographic shift. ICLR, 2022.
2. Minglai Shao, Dong Li, Chen Zhao, Xintao Wu, Yujie Lin, Qin Tian. Supervised Algorithmic Fairness in Distribution Shifts: A Survey. IJCAI 2024.

Dependence (Correlation) Shift

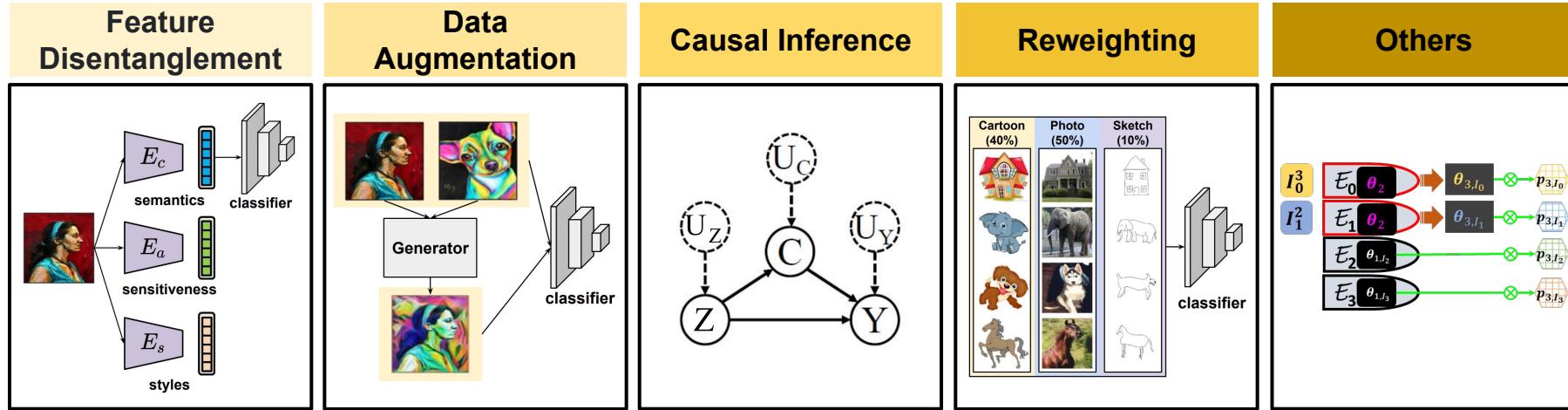
Dependence Shift can explicitly capture the alteration in *the correlation between Y and Z* of the data across the source and target domains.



1. Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. Improving fair training under correlation shifts. ICML, 2023

2. Minglai Shao, Dong Li, Chen Zhao, Xintao Wu, Yujie Lin, Qin Tian. Supervised Algorithmic Fairness in Distribution Shifts: A Survey. IJCAI 2024.

Approaches for Supervised Algorithmic Fairness in Distribution Shifts



Outline

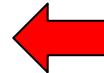
Part I: Introduction to Fairness



Part II: Distribution Shift Undermines Fairness



Part III: Mitigating Unfairness under Distribution Shift (Offline)



Part IV: Mitigating Unfairness under Distribution Shift (Online)

Part V: Open Challenges and Beyond

Break

- 15 Minutes

Outline

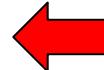
Part I: Introduction to Fairness



Part II: Distribution Shift Undermines Fairness



Part III: Mitigating Unfairness under Distribution Shift (Offline)

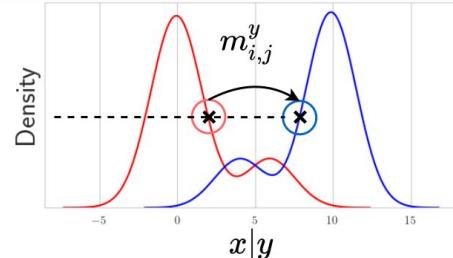


Part IV: Mitigating Unfairness under Distribution Shift (Online)

Part V: Open Challenges and Beyond

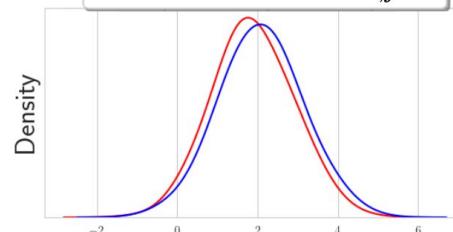
Fairness and Accuracy Transfer by Density Matching (FATDM)

— source domain i — source domain j



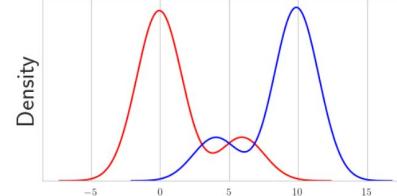
Stage 1: Finding $m_{i,j}^y : P_{D_i^S}^{x|y} = P_{D_j^S}^{x|y} \forall y \in \mathcal{Y}$

— $z|x$ — $z|m_{i,j}^y(x)$



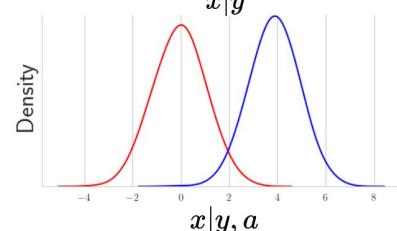
Stage 2: Enforcing $P^{(z|x)} = P^{z|m_{i,j}^y(x)} \forall x \in \mathcal{X}, y \in \mathcal{Y}$

— source domain i — source domain j



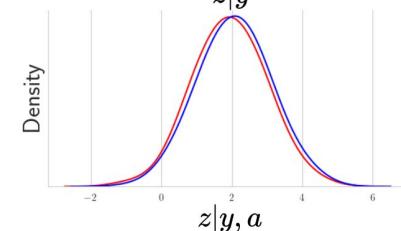
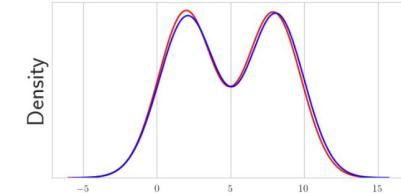
$$P_{D_i^S}^{z|y} = P_{D_j^S}^{z|y}$$

$$\forall y \in \mathcal{Y}$$



$$P_{D_i^S}^{z|y,a} = P_{D_j^S}^{z|y,a}$$

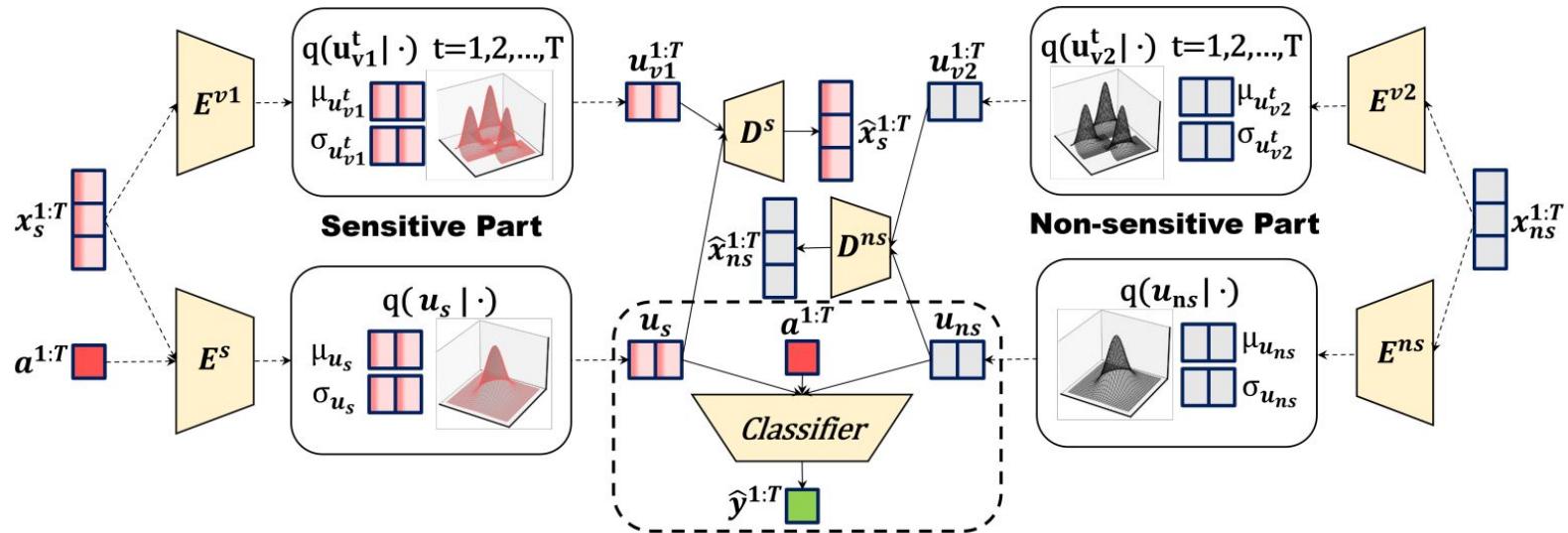
$$\forall y \in \mathcal{Y}, a \in \mathcal{A}$$



Types of Distribution Shifts and Methods

Covariate Shift ✓	Label Shift	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement	Augmentation ✓	Causal	Reweighting	Others

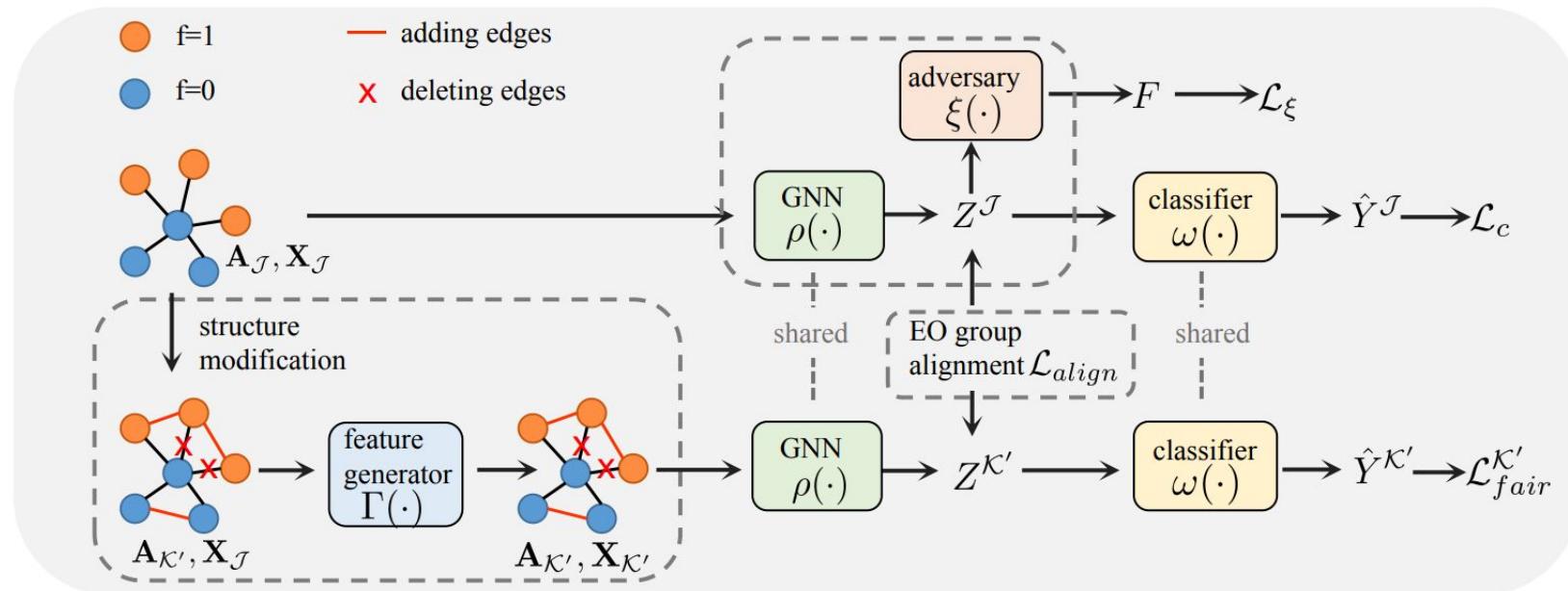
Disentanglement for Counterfactual Fairness-aware Domain Generalization (DCF DG)



Types of Distribution Shifts and Methods

Covariate Shift ✓	Label Shift	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement ✓	Augmentation	Causal ✓	Reweighting	Others ✓

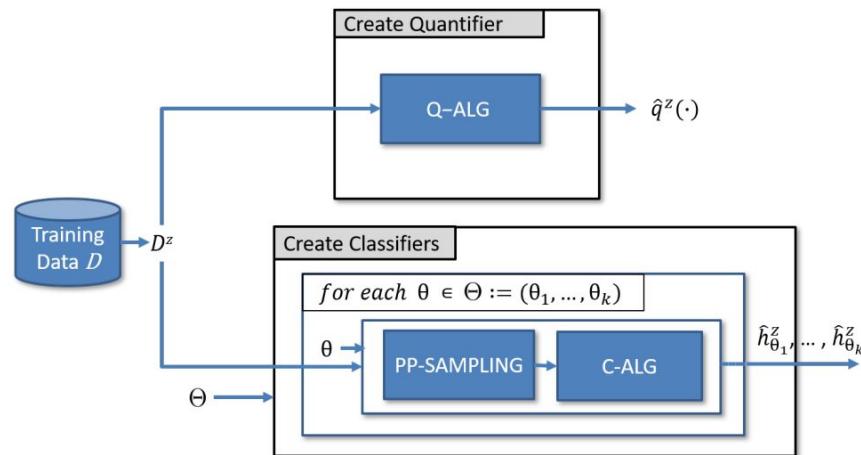
Graph Fairness Learning under Distribution Shifts (FatraGNN)



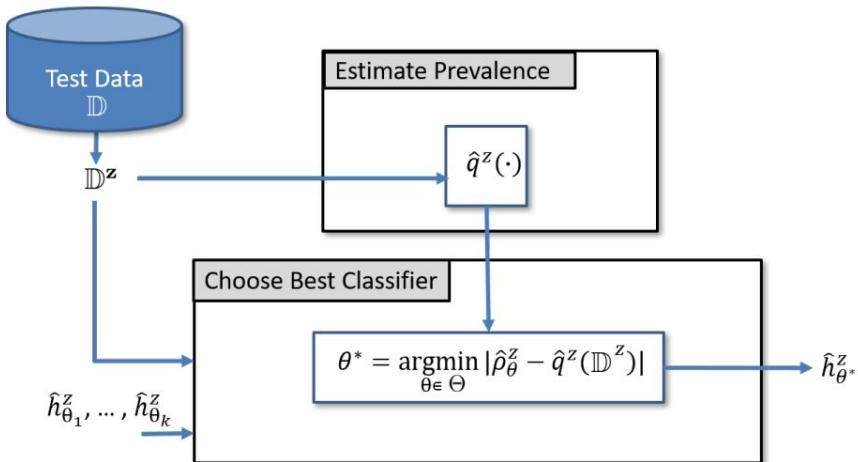
Types of Distribution Shifts and Methods

Covariate Shift ✓	Label Shift	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement	Augmentation ✓	Causal	Reweighting	Others

Combinatorial Algorithm for Proportional Equality (CAPE)



Training Phase

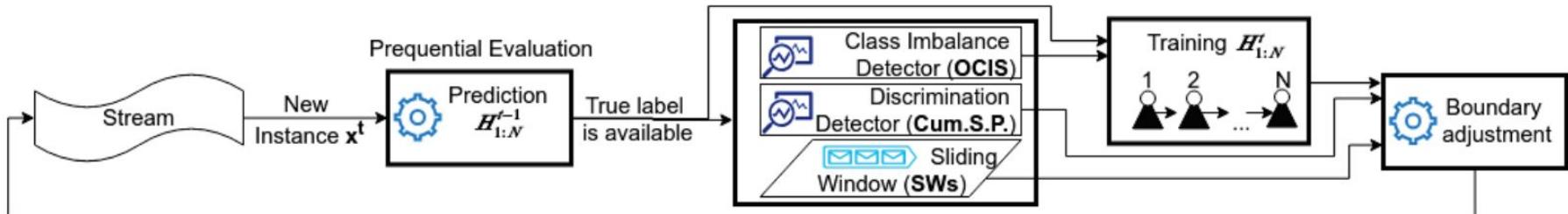


Testing Phase

Types of Distribution Shifts and Methods

Covariate Shift	Label Shift ✓	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement	Augmentation	Causal	Reweighting	Others ✓

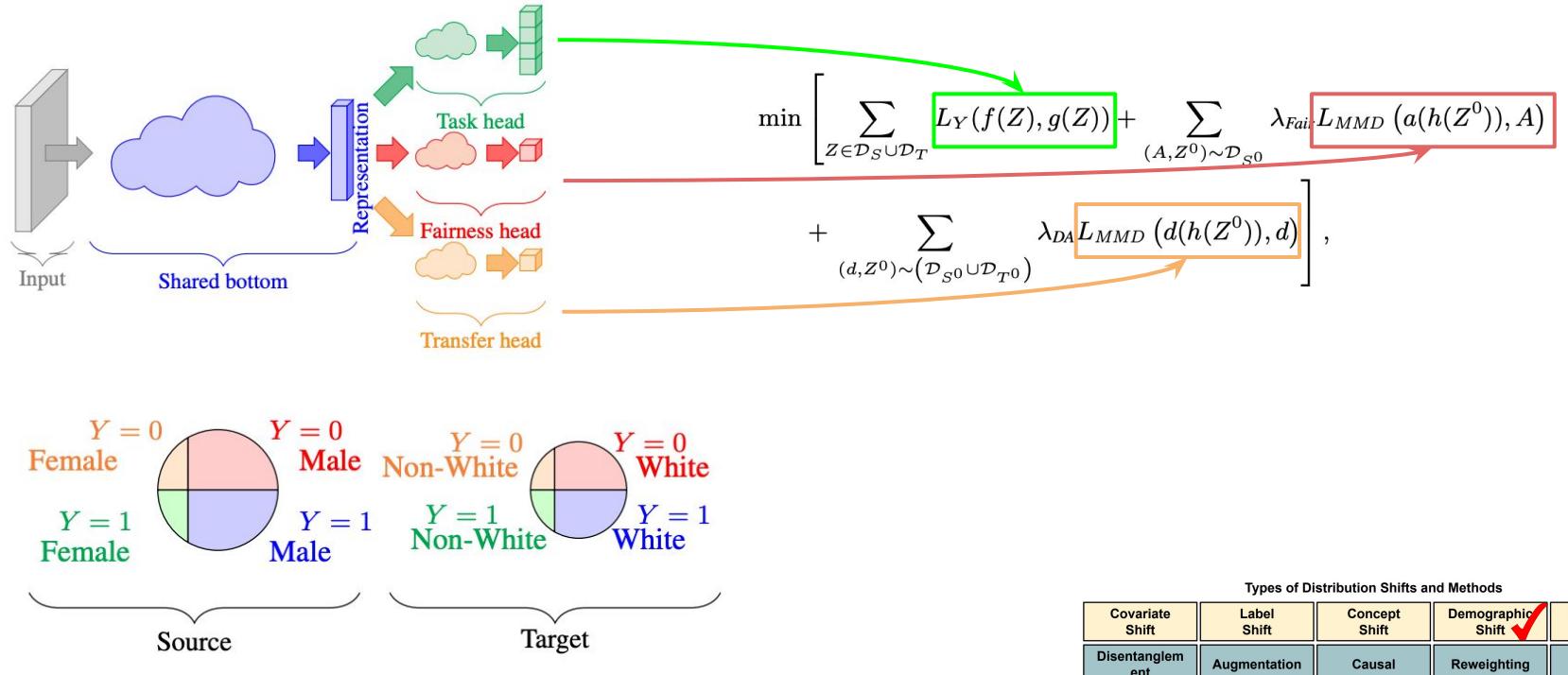
FAirness and class imBalance-aware BOOsting (FABBOO)



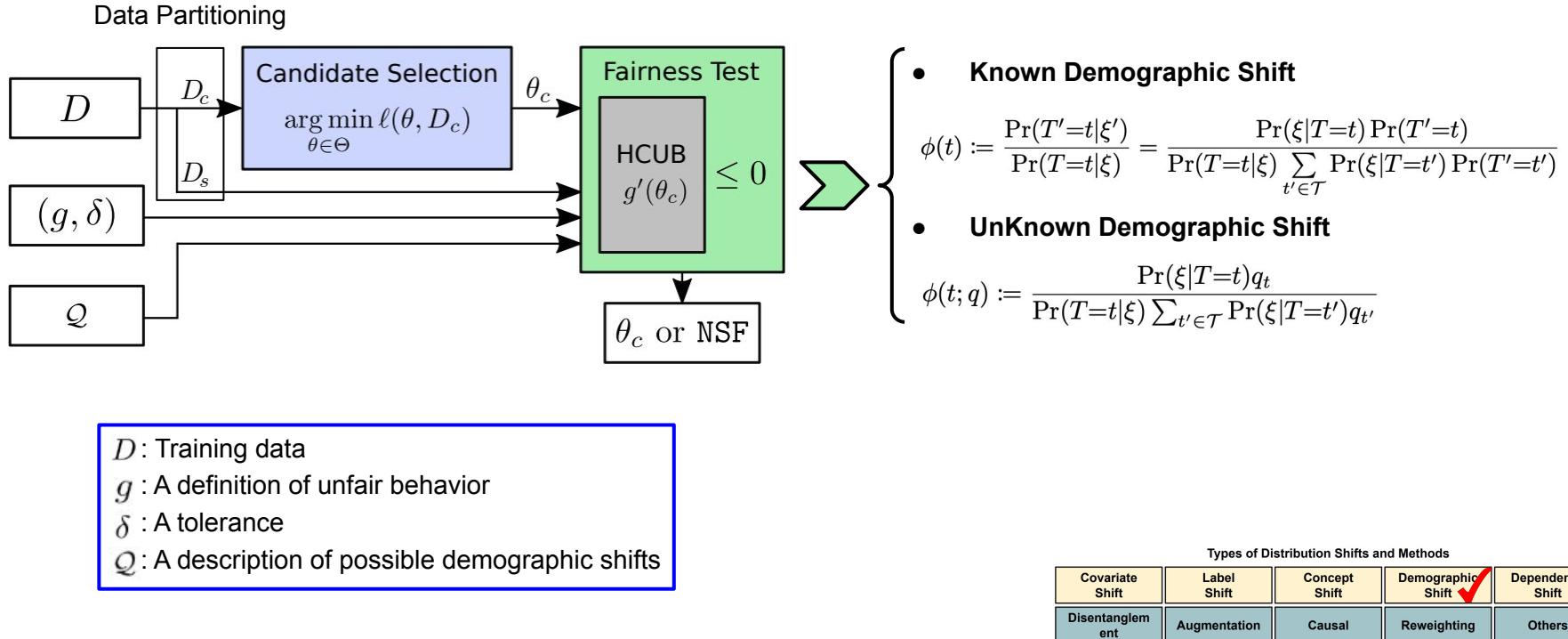
Types of Distribution Shifts and Methods

Covariate Shift	Label Shift	Concept Shift ✓	Demographic Shift	Dependence Shift
Disentanglement	Augmentation	Causal	Reweighting	Others ✓

Transfer of Machine Learning Fairness across Domains



Fairness Guarantees under Demographic Shift (Shifty)



Environment Inference for Invariant Learning (EIIL)



(a) **Inferred environment 1**
*(mostly) landbirds on land, and
waterbirds on water*



(b) **Inferred environment 2**
*(mostly) landbirds on water,
and waterbirds on land*

EIIL involves the following sequential approach:

1. Input *reference model* $\tilde{\Phi}$;
2. Fix $\Phi \leftarrow \tilde{\Phi}$ and optimize the EI objective to infer environments: $\mathbf{q}^* = \arg \max_{\mathbf{q}} C^{EI}(\tilde{\Phi}, \mathbf{q})$;
3. Fix $\tilde{\mathbf{q}} \leftarrow \mathbf{q}^*$ and optimize the IL objective to yield the new model: $\Phi^* = \arg \min_{\Phi} C^{IL}(\Phi, \tilde{\mathbf{q}})$

$$C^{EI}(\Phi, \mathbf{q}) = \|\nabla_{\bar{w}} \tilde{R}^e(\bar{w} \circ \Phi, \mathbf{q})\|,$$

$$\tilde{R}^e(\Phi, \mathbf{q}) = \frac{1}{\sum_{i'} \mathbf{q}_{i'}(e)} \sum_i \mathbf{q}_i(e) \ell(\Phi(x_i), y_i)$$

In practice, $C^{IL} \in \{C^{IRM}, C^{GroupDRO}\}$

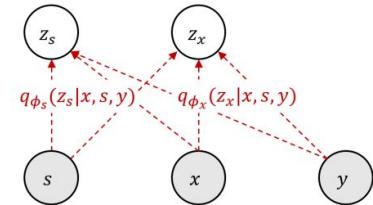
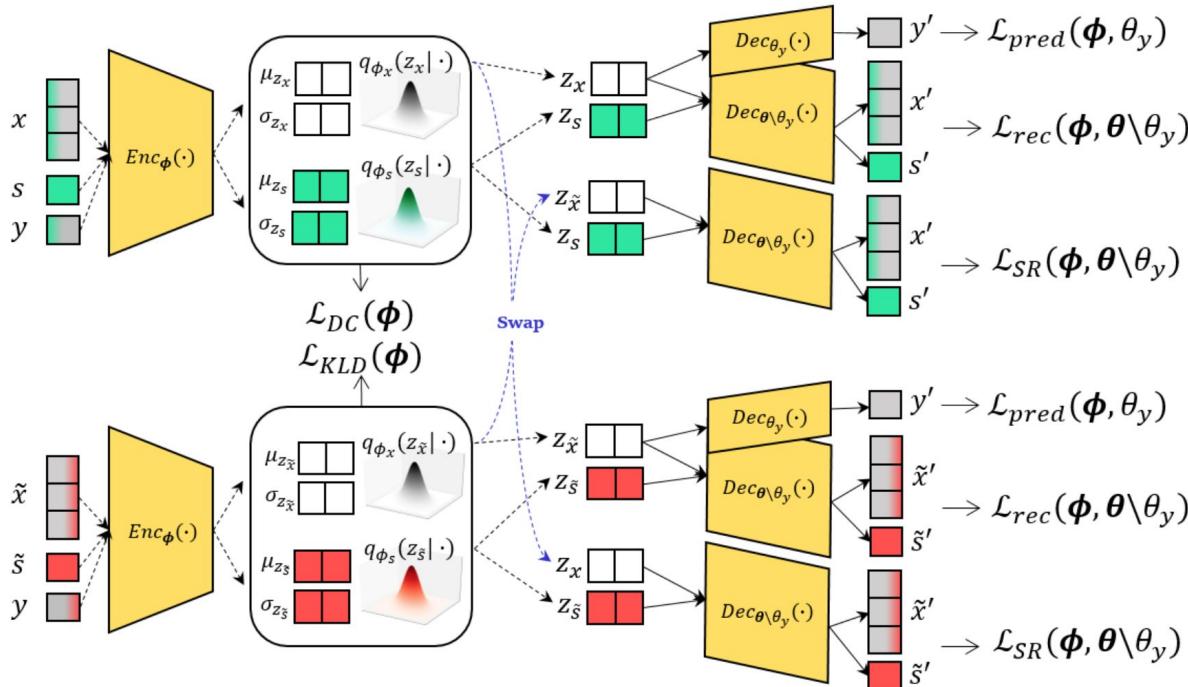
$$C^{IRM}(\Phi) = \sum_{e \in \mathcal{E}^{obs}} R^e(\Phi) + \lambda \|\nabla_{\bar{w}} R^e(\bar{w} \circ \Phi)\|$$

$$C^{GroupDRO}(\Phi) = \max_g \mathbb{E}_{g(e)}[R^e(\Phi)]$$

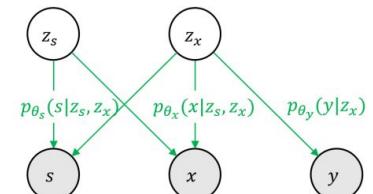
Types of Distribution Shifts and Methods

Covariate Shift	Label Shift	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement	Augmentation	Causal	Reweighting	Others

FAir Representation via distributional CONtrastive Variational AutoEncoder (FarconVAE)



(a) Recognition model

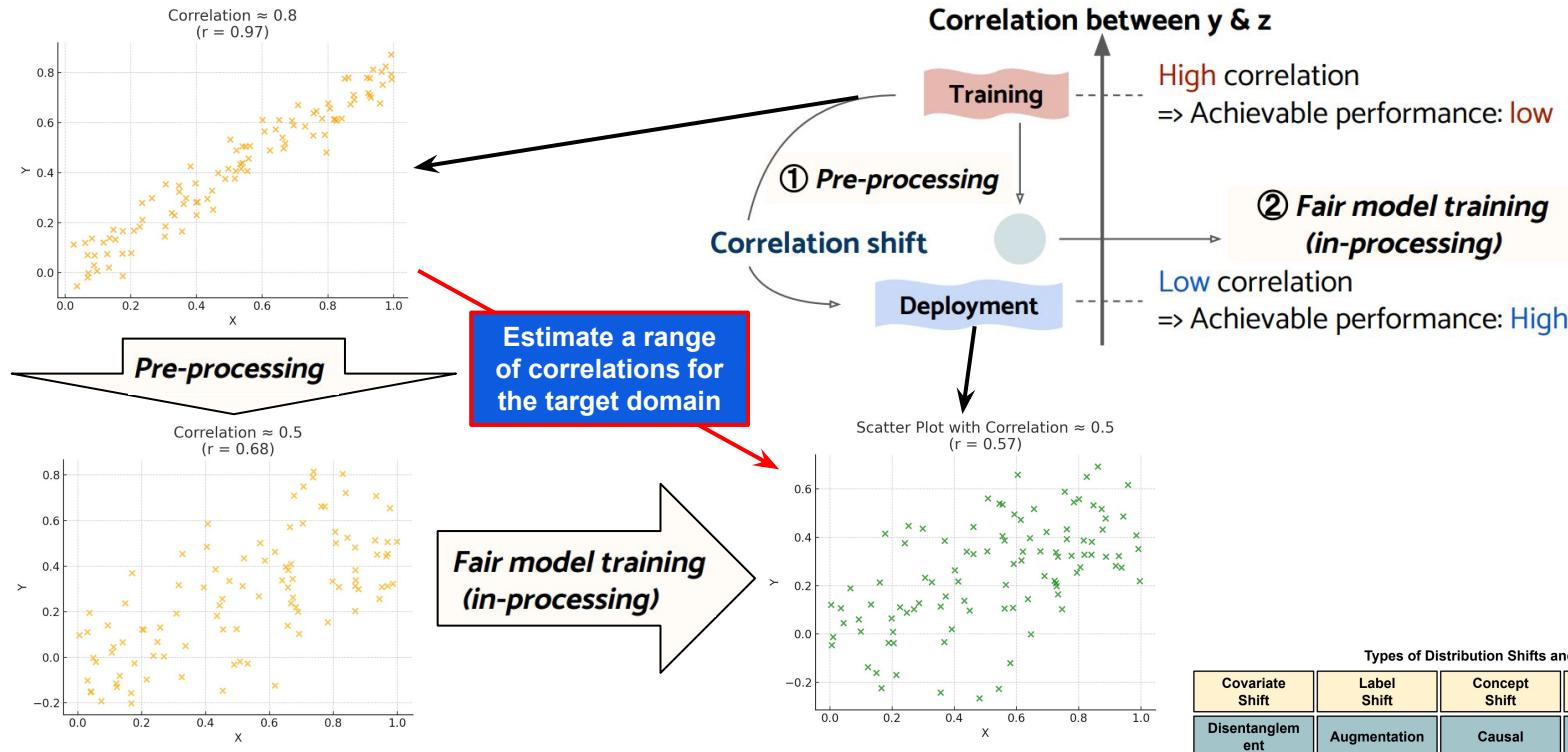


(b) Generative model

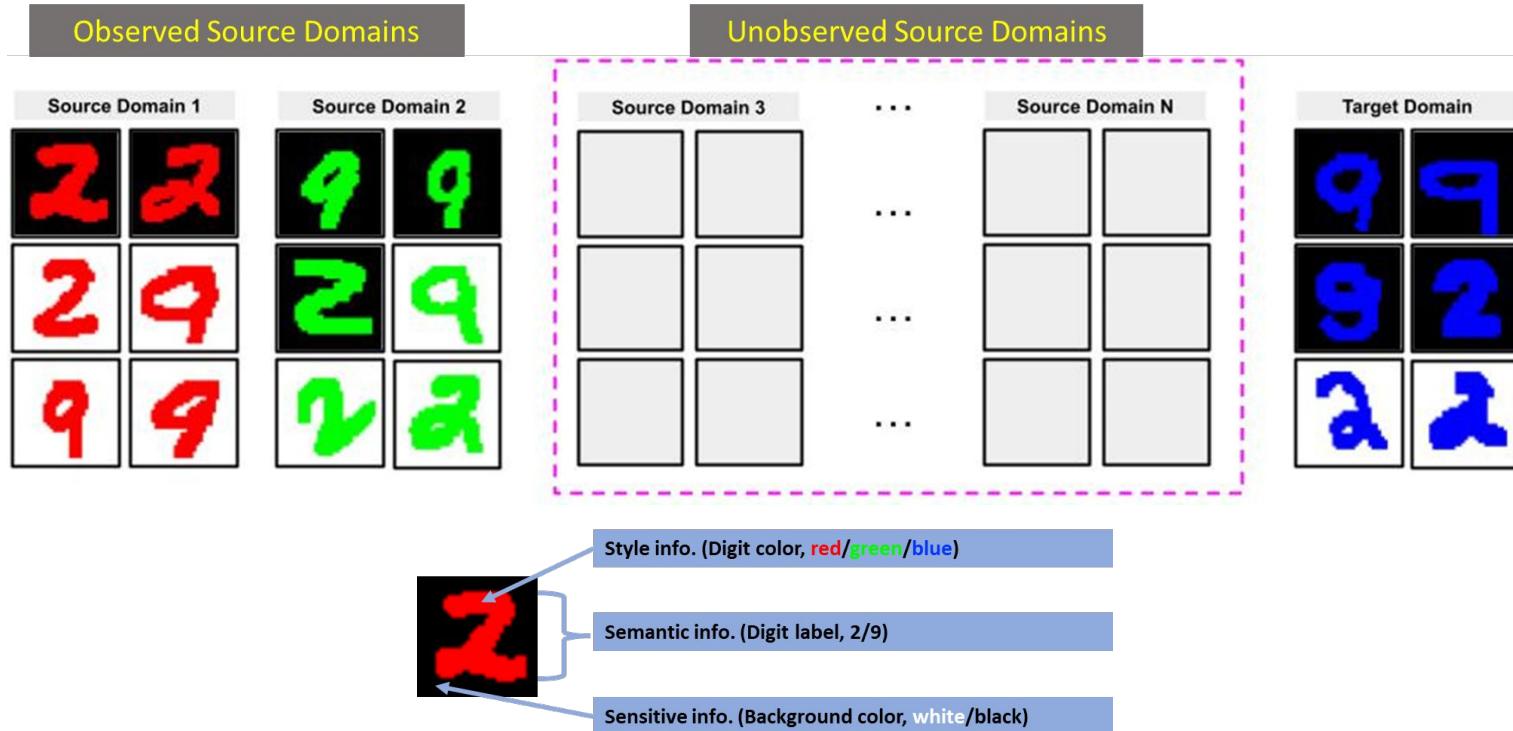
Types of Distribution Shifts and Methods				
Covariate Shift	Label Shift	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement ✓	Augmentation	Causal	Reweighting	Others

- Changdae Oh, Heeji Won, Junhyuk So, Taero Kim, Yewon Kim, Hosik Choi, and Kyungwoo Song. Learning Fair Representation via Distributional Contrastive Disentanglement. KDD, 2022.

Improving Fair Training under Correlation Shifts

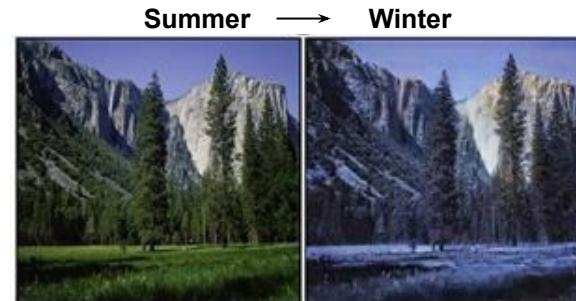
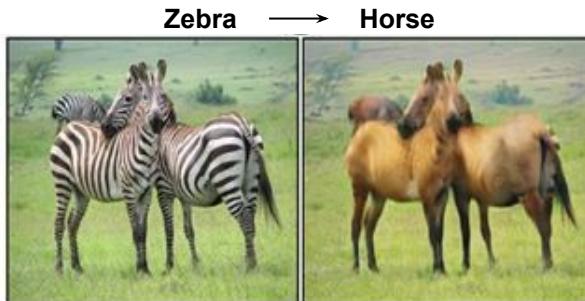


The More You See; The More You Know...

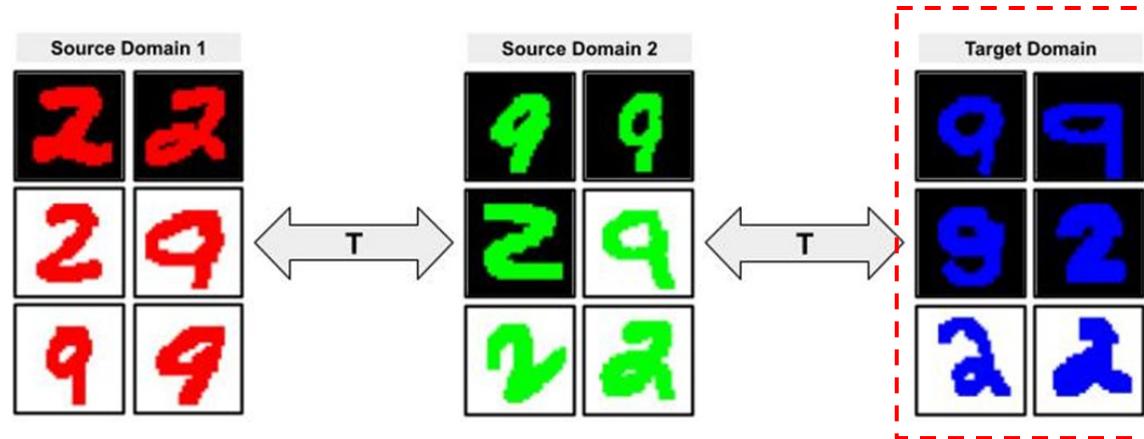


- Chen Zhao, Kai Jiang, Xintao Wu, Haoliang Wang, Latifur Khan, Christian Grant, Feng Chen. Algorithmic Fairness Generalization under Covariate and Dependence Shifts Simultaneously. KDD 2024

Transformation Models



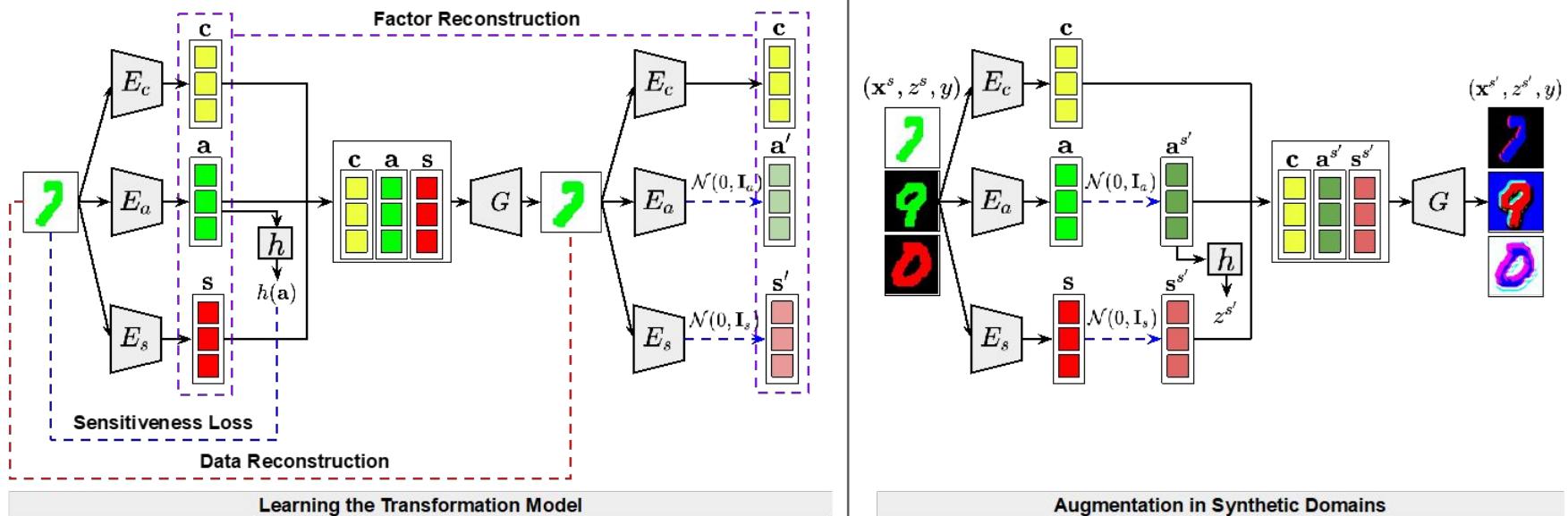
Transformation Models



There exists a transformation model T that:

1. Distribution shifts can characterize generalization tasks across domains.
2. Generate augmented data in unobserved domains by perturbing existing samples with various variations.

Fair disEntangled DOmain geneRALization (FEDORA)

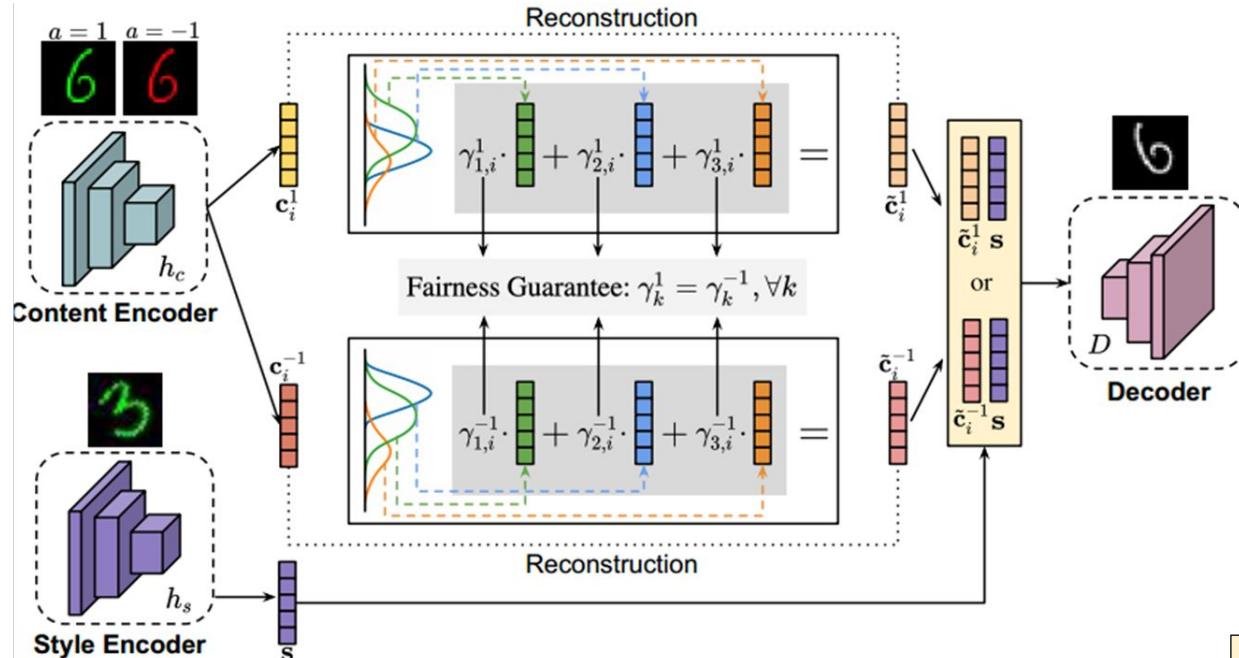


Types of Distribution Shifts and Methods

Covariate Shift ✓	Label Shift	Concept Shift	Demographic Shift	Dependence Shift ✓
Disentanglement ✓	Augmentation ✓	Causal	Reweighting	Others ✓

- Chen Zhao, Kai Jiang, Xintao Wu, Haoliang Wang, Latifur Khan, Christian Grant, Feng Chen. Algorithmic Fairness Generalization under Covariate and Dependence Shifts Simultaneously. KDD 2024

Fairness-aware LeArning Invariant Representations (FLAIR)



- Minimize the information disclosure related to a specific sensitive subgroup
- Maximize the preservation of significant information within non-sensitive representations

Types of Distribution Shifts and Methods				
Covariate Shift ✓	Label Shift	Concept Shift	Demographic Shift	Dependence Shift ✓
Disentanglement ✓	Augmentation	Causal	Reweighting	Others ✓

Fairness-aware Outlier Detection (FairOD)

PROBLEM 1 (FAIRNESS-AWARE OUTLIER DETECTION). Given samples X and protected variable values \mathcal{PV} , estimate outlier scores S and assign outlier labels O , to achieve

$$(i) P(Y = 1|O = 1) > P(Y = 1),$$

[Detection effectiveness]

$$(ii) P(O | X, PV = v) = P(O | X), \forall v \in \{a, b\},$$

[Treatment parity]

$$(iii) P(O = 1|PV = a) = P(O = 1|PV = b),$$

[Statistical parity]

$$(iv) \pi_{PV=v}^{\text{BASE}} = \pi_{PV=v}, \forall v \in \{a, b\}, \text{ where BASE is a fairness-agnostic detector.}$$

[Group fidelity proxy]

$$\mathcal{L} = \alpha \underbrace{\mathcal{L}_{\text{BASE}}}_{\text{Reconstruction}} + (1 - \alpha) \underbrace{\mathcal{L}_{\text{SP}}}_{\text{Statistical Parity}} + \gamma \underbrace{\mathcal{L}_{\text{GF}}}_{\text{Group Fidelity}}$$

$$\mathcal{L}_{\text{BASE}} = \sum_{i=1}^N \|X_i - G(X_i)\|_2^2$$

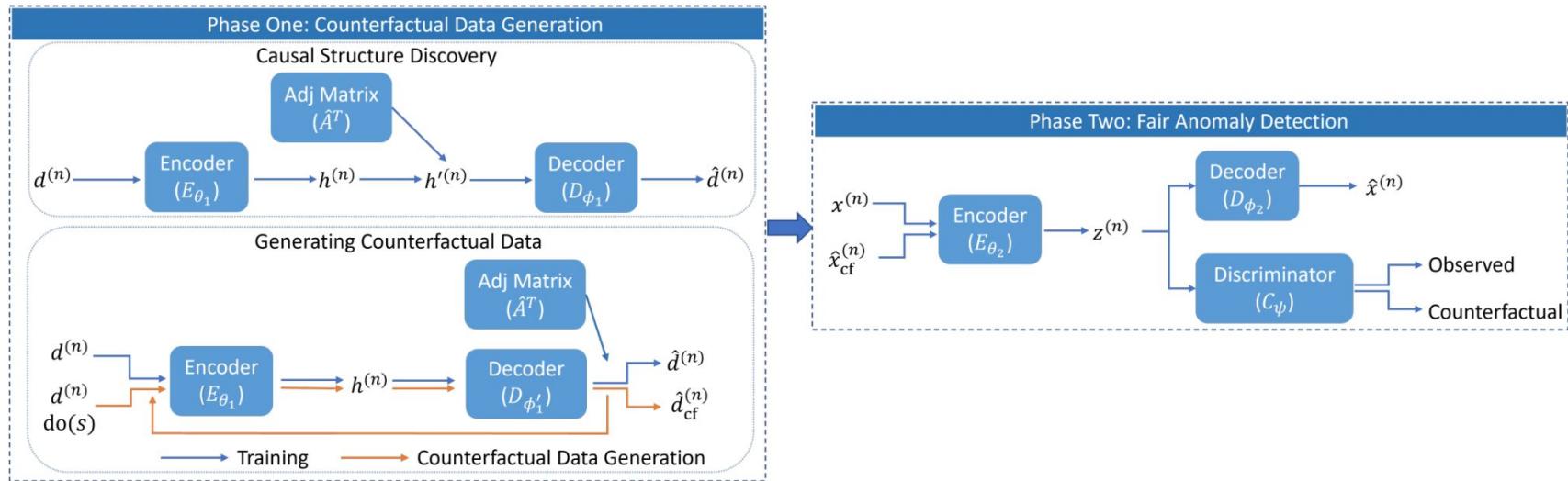
$$\mathcal{L}_{\text{SP}} = \left| \frac{(\sum_{i=1}^N s(X_i) - \mu_s) (\sum_{i=1}^N PV_i - \mu_{PV})}{\sigma_s \sigma_{PV}} \right|$$

$$\mathcal{L}_{\text{GF}} = \sum_{v \in \{a, b\}} \left(1 - \sum_{X_i \in \mathcal{X}_{PV=v}} \frac{2^{s^{\text{BASE}}(X_i)} - 1}{\text{DNM}} \right)$$

Types of Distribution Shifts and Methods

Covariate Shift ✓	Label Shift ✓	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement	Augmentation	Causal	Reweighting	Others ✓

Counterfactually Fair Anomaly Detection (CFAD)



Types of Distribution Shifts and Methods

Covariate Shift ✓	Label Shift ✓	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement ✓	Augmentation	Causal ✓	Reweighting	Others

Outline

Part I: Introduction to Fairness



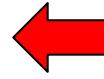
Part II: Distribution Shift Undermines Fairness



Part III: Mitigating Unfairness under Distribution Shift (Offline)

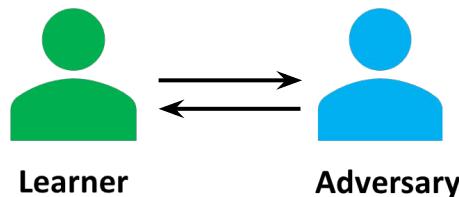
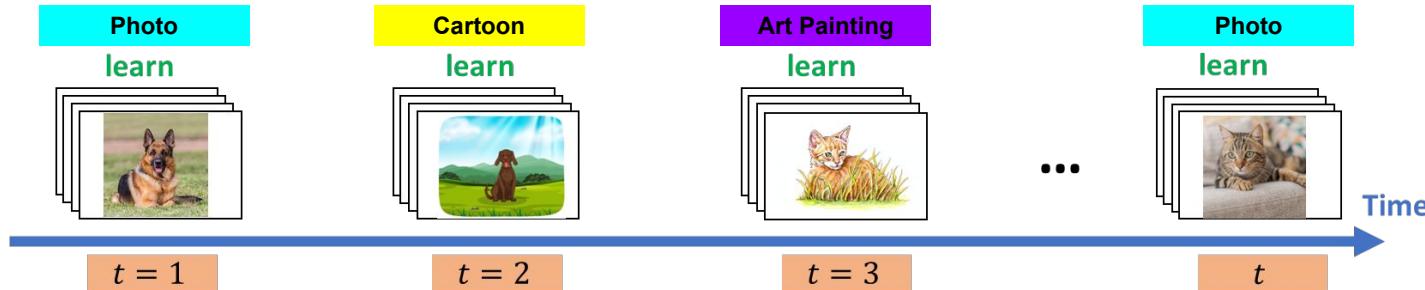


Part IV: Mitigating Unfairness under Distribution Shift (Online)



Part V: Open Challenges and Beyond

Fair Online Learning under Distribution Shifts



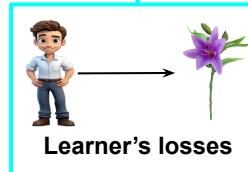
- Step 1: The learner **SELECTS** θ_t in the fair domain Θ
- Step 2: The adversary reveals a loss function $f_t(\cdot)$ and fairness function $g(\cdot)$
- Step 3: The learner incurs an instantaneous loss $f_t(\theta_t, \mathcal{D}_t)$ and fairness notion $g_i(\theta_t, \mathcal{D}_t)$
- Step 4: Advance to round $t + 1$

Regrets

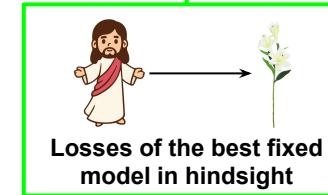


Static Regret:

$$R_S = \sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta)$$



$$\min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta)$$



- Static regret works for data coming from fixed distribution (with non-stationary task distributions).
- The goal is to design algorithms such that this regret grows with T as slowly as possible.
- In particular, an algorithm whose regret grows sub-linearly in T is non-trivially learning and adapting.

Regrets

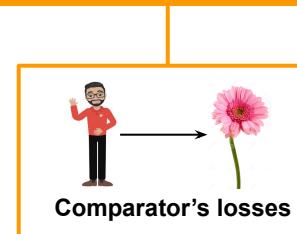
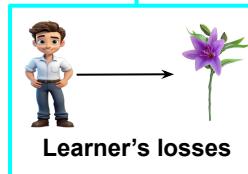


Static Regret:

$$R_S = \sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta)$$

Dynamic Regret:

$$\begin{aligned} R_D &= \sum_{t=1}^T \ell_t(\theta_t) - \sum_{t=1}^T \ell_t(u_t) \\ &= \boxed{\sum_{t=1}^T \ell_t(\theta_t)} - \boxed{\sum_{t=1}^T \min_{\theta \in \Theta} \ell_t(\theta)} \end{aligned}$$



- Dynamic regret works for data coming from changing environments.

Regrets



Static Regret:

$$R_S = \sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta)$$

Dynamic Regret:

$$R_D = \sum_{t=1}^T \ell_t(\theta_t) - \sum_{t=1}^T \ell_t(u_t)$$

Adaptive Regret:

$$R_A = \max_{[s,q] \subseteq [T]} \left(\sum_{t=s}^q \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=s}^q \ell_t(\theta) \right)$$

static regrets in short intervals

- Both dynamic and adaptive regrets work for data coming from changing environments.
- Dynamic regret addresses changing environments from a **global** perspective, while adaptive regret takes a **local** perspective by focusing on comparators within short intervals.

1. Martin Zinkevich. Online Convex Programming and Generalized Infinitesimal Gradient Ascent. ICML, 2003.

2. Elad Hazan, C. Seshadhri. Adaptive Algorithms for Online Decision Problems. Electronic Colloquium on Computational Complexity, 2007.

Regrets



Static Regret:

$$R_S = \sum_{t=1}^T \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta)$$

Dynamic Regret:

$$R_D = \sum_{t=1}^T \ell_t(\theta_t) - \sum_{t=1}^T \ell_t(u_t)$$

Adaptive Regret:

$$R_A = \max_{[s,q] \subseteq [T]} \left(\sum_{t=s}^q \ell_t(\theta_t) - \min_{\theta \in \Theta} \sum_{t=s}^q \ell_t(\theta) \right)$$

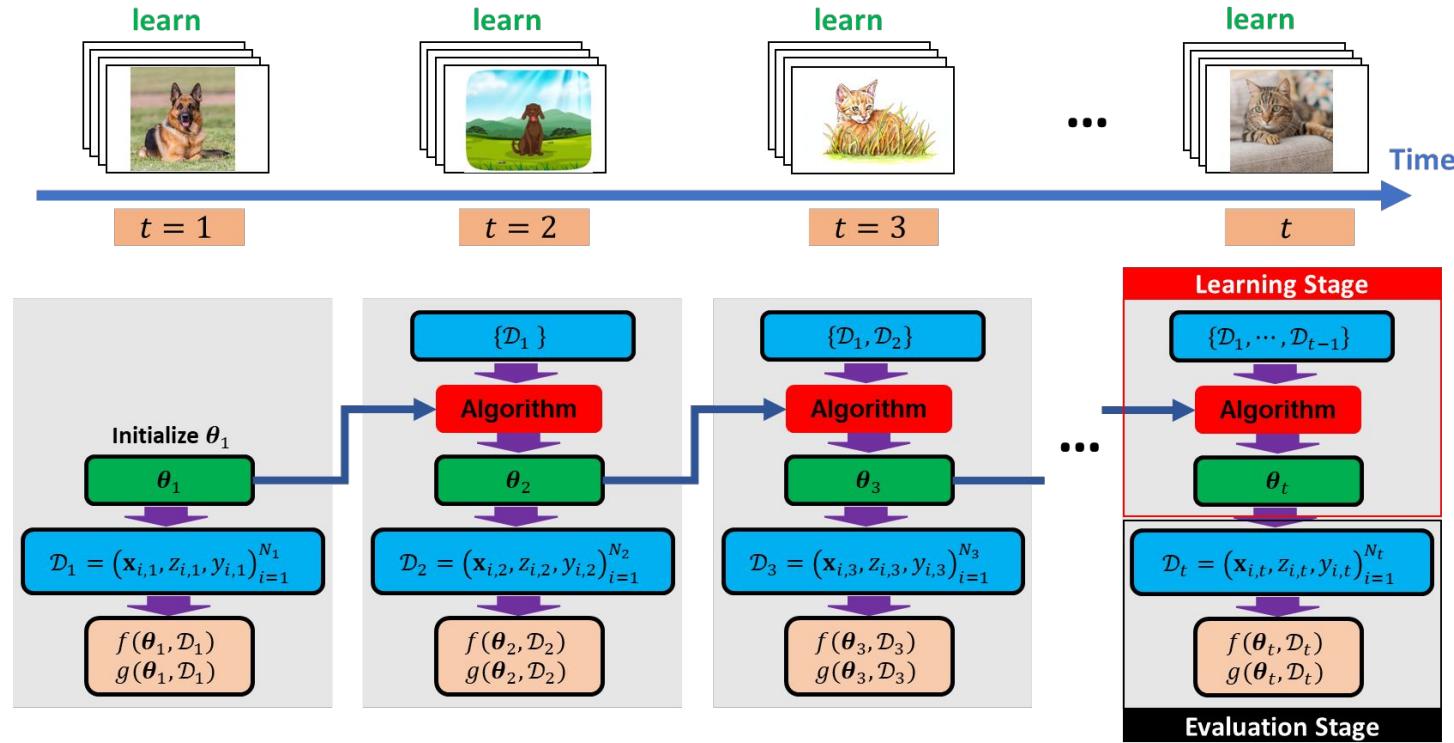
Long-Term Fair Constraints:

$$\sum_{t=1}^T \left\| [g_i(\theta_t)]_+ \right\| \leq \epsilon, i \in [m]$$

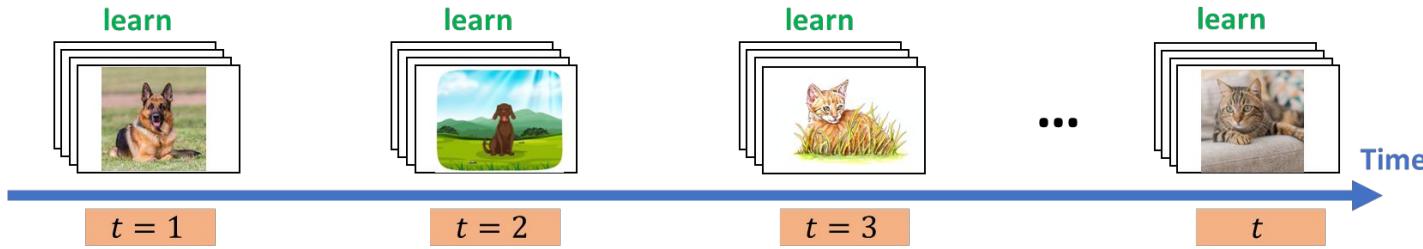
add to these regrets

1. Martin Zinkevich. Online Convex Programming and Generalized Infinitesimal Gradient Ascent. ICML, 2003.
2. Elad Hazan, C. Seshadhri. Adaptive Algorithms for Online Decision Problems. Electronic Colloquium on Computational Complexity, 2007.
3. Mehrdad Mahdavi, Rong Jin, and Tianbao Yang. Trading Regret for Efficiency: Online Convex Optimization with Long Term Constraints. JMLR, 2012.

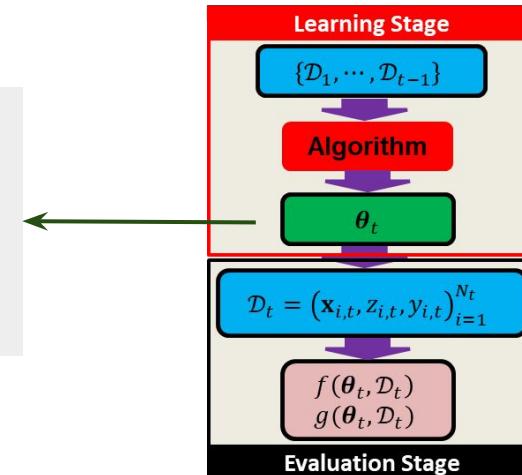
Fair Online Learning under Distribution Shifts



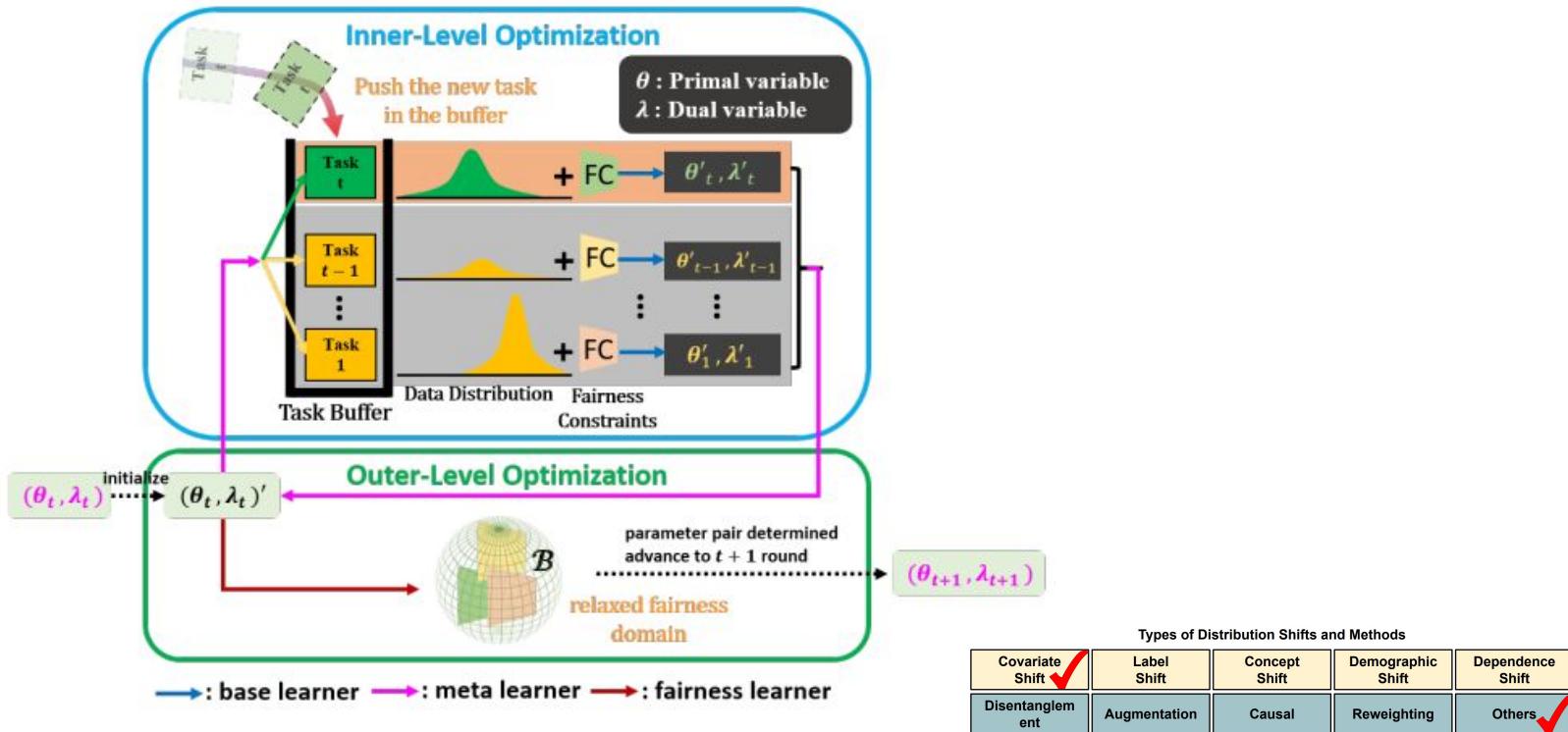
Fair Online Learning under Distribution Shifts



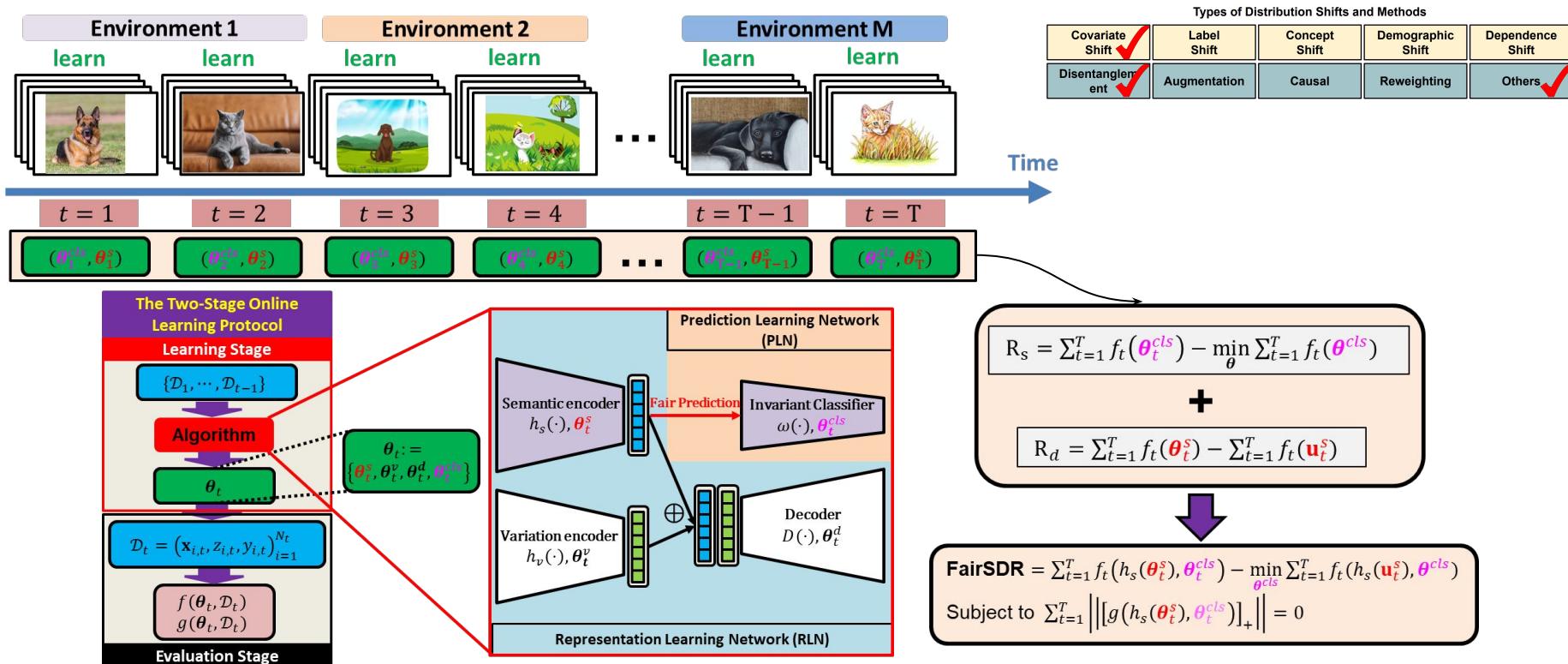
The goal is to develop an effective fairness-aware online algorithm for changing environments, ensuring that the learned parameter sequence $\{\theta_1, \dots, \theta_t\}$ minimizes both **the loss regret** and **the violation of long-term fairness constraints**.



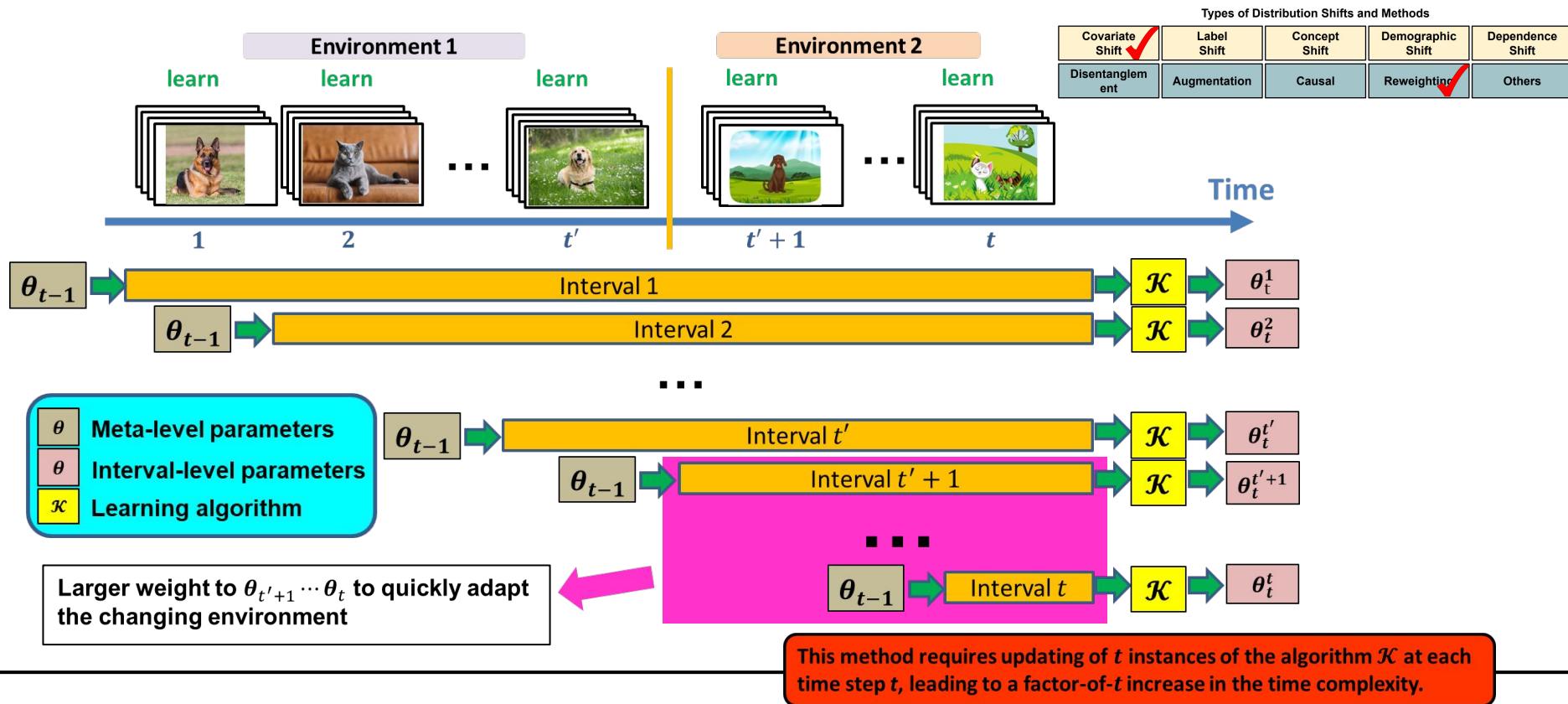
Follow the Fair Meta Leader (FFML)



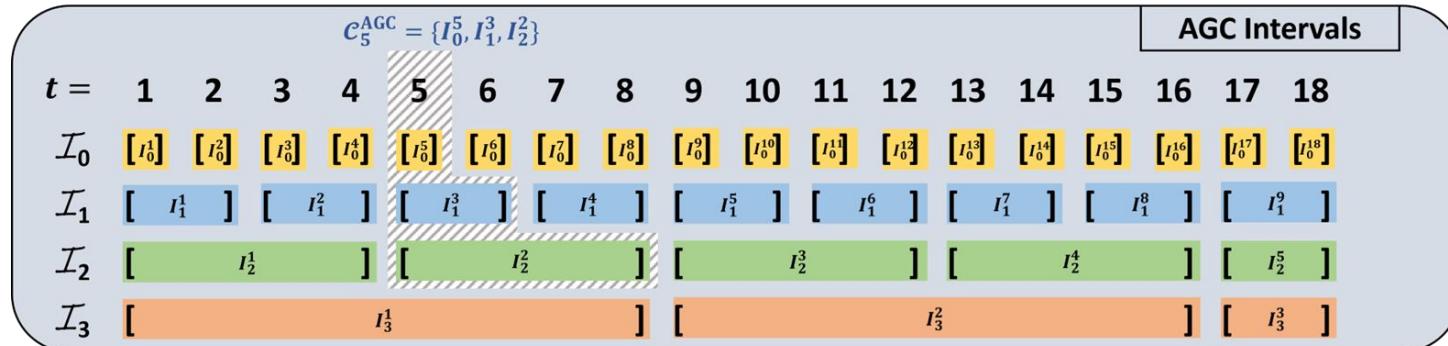
Fair Disentanglement Online Learning for Changing Environments (FairDolce)



Setting of Dynamic Intervals (DI)



Setting of Adaptive Geometric Covering (AGC) Intervals



- AGC divides intervals into N sets of intervals (i.e., $\{\mathcal{I}_k\}_{k=0}^{N-1}$), where $N = \lceil \log_2 T \rceil$.
- Each \mathcal{I}_k contains equal length intervals, where $\mathcal{I}_k = \{I_k^i | [(i-1) \cdot 2^k + 1, \min\{T, i \cdot 2^k\}] : i \in \mathbb{N}\}$.
- At each time t , a target set \mathcal{C}_t is selected, where \mathcal{C}_t includes all intervals starting with t . It is mathematically defined as $\mathcal{C}_t = \{I | I \in \mathcal{I}, t \in I, (t-1) \notin I\}$.
- AGC requires T needs to be known and fixed in advance.**

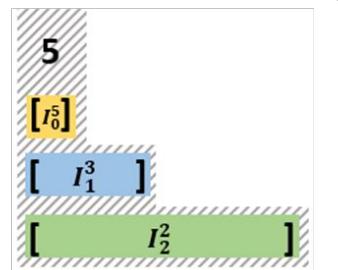
Types of Distribution Shifts and Methods					
Covariate Shift ✓	Label Shift	Concept Shift	Demographic Shift	Dependence Shift	
Disentanglement	Augmentation	Causal	Reweighting ✓	Others	

Setting of Adaptive Geometric Covering (AGC) Intervals

	AGC Intervals																	
$t =$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
\mathcal{I}_0	$[I_0^1]$	$[I_0^2]$	$[I_0^3]$	$[I_0^4]$	$[I_0^5]$	$[I_0^6]$	$[I_0^7]$	$[I_0^8]$	$[I_0^9]$	$[I_0^{10}]$	$[I_0^{11}]$	$[I_0^{12}]$	$[I_0^{13}]$	$[I_0^{14}]$	$[I_0^{15}]$	$[I_0^{16}]$	$[I_0^{17}]$	$[I_0^{18}]$
\mathcal{I}_1	$[I_1^1]$	$[I_1^2]$	$[I_1^3]$	$[I_1^4]$	$[I_1^5]$	$[I_1^6]$	$[I_1^7]$	$[I_1^8]$	$[I_1^9]$	$[I_1^{10}]$	$[I_1^{11}]$	$[I_1^{12}]$	$[I_1^{13}]$	$[I_1^{14}]$	$[I_1^{15}]$	$[I_1^{16}]$	$[I_1^{17}]$	$[I_1^{18}]$
\mathcal{I}_2	$[I_2^1]$	$[I_2^2]$	$[I_2^3]$	$[I_2^4]$	$[I_2^5]$	$[I_2^6]$	$[I_2^7]$	$[I_2^8]$	$[I_2^9]$	$[I_2^{10}]$	$[I_2^{11}]$	$[I_2^{12}]$	$[I_2^{13}]$	$[I_2^{14}]$	$[I_2^{15}]$	$[I_2^{16}]$	$[I_2^{17}]$	$[I_2^{18}]$
\mathcal{I}_3	$[I_3^1]$	$[I_3^2]$	$[I_3^3]$	$[I_3^4]$	$[I_3^5]$	$[I_3^6]$	$[I_3^7]$	$[I_3^8]$	$[I_3^9]$	$[I_3^{10}]$	$[I_3^{11}]$	$[I_3^{12}]$	$[I_3^{13}]$	$[I_3^{14}]$	$[I_3^{15}]$	$[I_3^{16}]$	$[I_3^{17}]$	$[I_3^{18}]$

For example, when $t = 5$

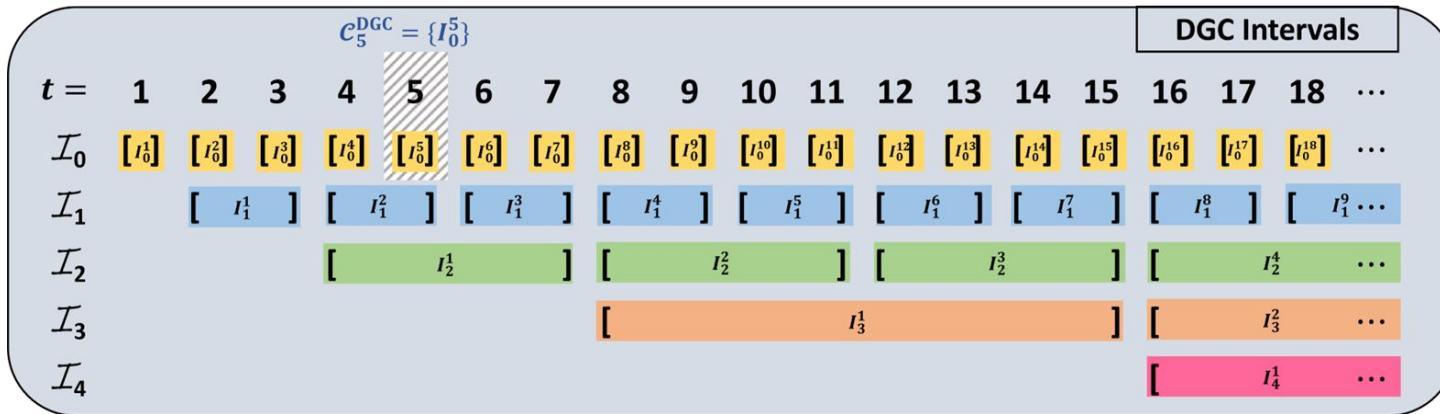
$$\mathcal{C}_5 = \{[I_0^5], [I_1^3], [I_2^2]\}$$



Types of Distribution Shifts and Methods

Covariate Shift ✓	Label Shift	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement	Augmentation	Causal	Reweighting ✓	Others

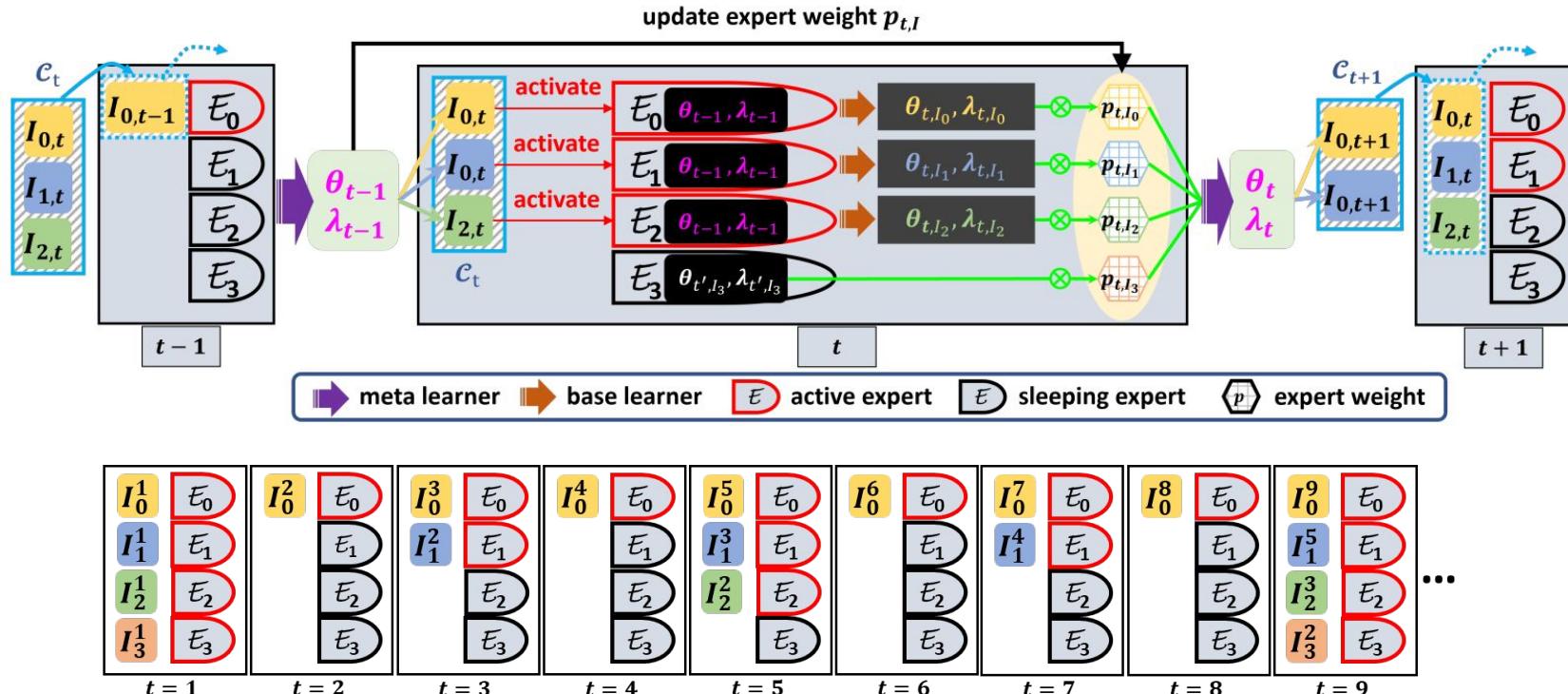
Setting of Dynamic Geometric Covering (DGC) Intervals



- DGC **dynamically** divides intervals into $\lceil \log_2 t \rceil - 1$ sets of intervals.
- Each \mathcal{I}_k contains equal length intervals, where $\mathcal{I}_k = \{I_k^i | [i \cdot 2^k, (i + 1) \cdot 2^k - 1] : i \in \mathbb{N}\}$.
- At each time t , a target set \mathcal{C}_t is selected, where \mathcal{C}_t includes all intervals starting with t . It mathematically defined as $\mathcal{C}_t = \{I | I \in \mathcal{I}, t \in I, (t - 1) \notin I\}$.

Types of Distribution Shifts and Methods				
Covariate Shift ✓	Label Shift	Concept Shift	Demographic Shift	Dependence Shift
Disentanglement	Augmentation	Causal	Reweighting ✓	Others

Learning Experts for AGC and DGC

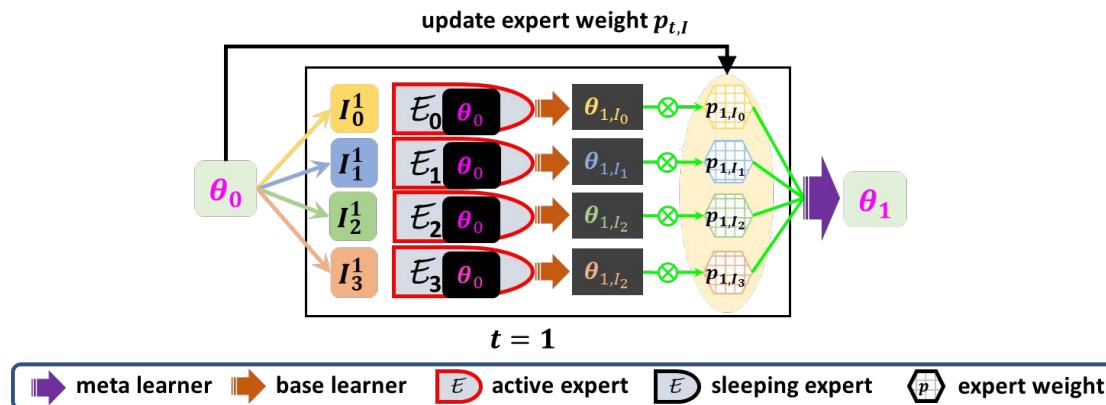


- Chen Zhao, Feng Mi, Xintao Wu, Kai Jiang, Latifur Khan, Feng Chen. Adaptive Fairness-Aware Online Meta-Learning for Changing Environments. KDD, 2022.

Learning with Experts (When t=1)

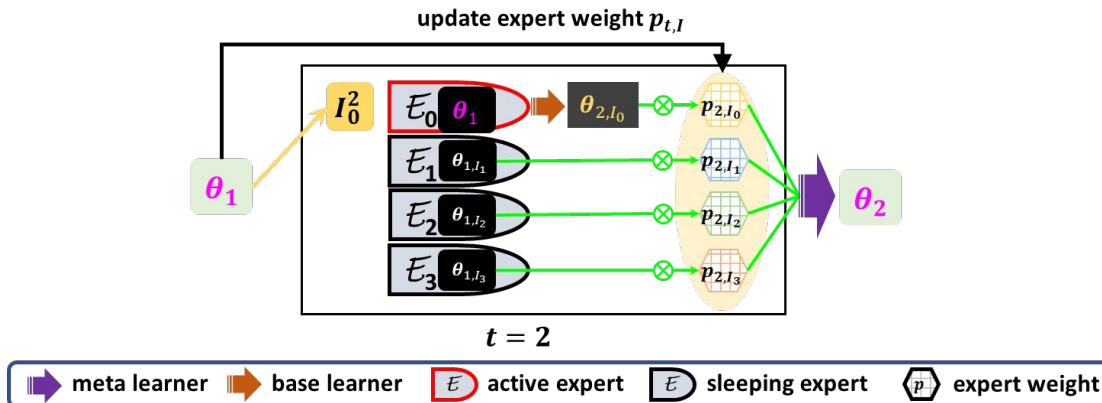
$$\mathcal{C}_1 = \{I_0^1, I_1^1, I_2^1, I_3^1\}$$

t	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
\mathcal{I}_0	$[I_0^1]$	$[I_0^2]$	$[I_0^3]$	$[I_0^4]$	$[I_0^5]$	$[I_0^6]$	$[I_0^7]$	$[I_0^8]$	$[I_0^9]$	$[I_0^{10}]$	$[I_0^{11}]$	$[I_0^{12}]$	$[I_0^{13}]$	$[I_0^{14}]$	$[I_0^{15}]$	$[I_0^{16}]$	$[I_0^{17}]$	$[I_0^{18}]$
\mathcal{I}_1	$[I_1^1]$	I_1^2	I_1^3	I_1^4	I_1^5	I_1^6	I_1^7	I_1^8	I_1^9	I_1^{10}	I_1^{11}	I_1^{12}	I_1^{13}	I_1^{14}	I_1^{15}	I_1^{16}	I_1^{17}	I_1^{18}
\mathcal{I}_2	I_2^1	I_2^2	I_2^3	I_2^4	I_2^5	I_2^6	I_2^7	I_2^8	I_2^9	I_2^{10}	I_2^{11}	I_2^{12}	I_2^{13}	I_2^{14}	I_2^{15}	I_2^{16}	I_2^{17}	I_2^{18}
\mathcal{I}_3	I_3^1	I_3^2	I_3^3	I_3^4	I_3^5	I_3^6	I_3^7	I_3^8	I_3^9	I_3^{10}	I_3^{11}	I_3^{12}	I_3^{13}	I_3^{14}	I_3^{15}	I_3^{16}	I_3^{17}	I_3^{18}



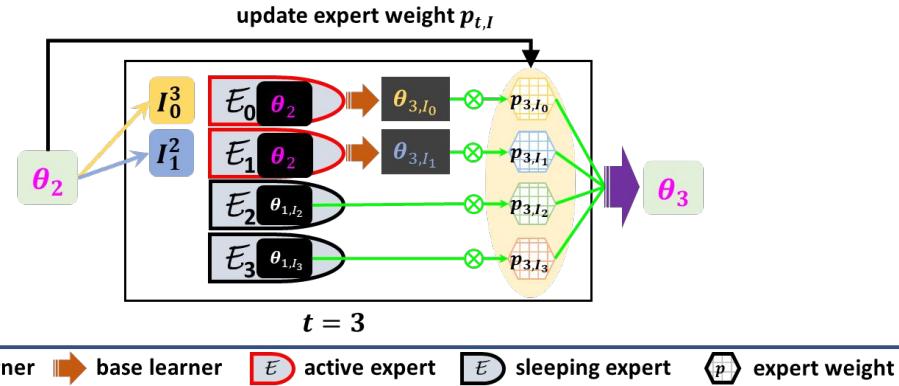
Learning with Experts (When t=2)

t	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
\mathcal{I}_0	$[I_0^1]$	$[I_0^2]$	$[I_0^3]$	$[I_0^4]$	$[I_0^5]$	$[I_0^6]$	$[I_0^7]$	$[I_0^8]$	$[I_0^9]$	$[I_0^{10}]$	$[I_0^{11}]$	$[I_0^{12}]$	$[I_0^{13}]$	$[I_0^{14}]$	$[I_0^{15}]$	$[I_0^{16}]$	$[I_0^{17}]$	$[I_0^{18}]$
\mathcal{I}_1	$[I_1^1]$	$[I_1^2]$	$[I_1^3]$	$[I_1^4]$	$[I_1^5]$	$[I_1^6]$	$[I_1^7]$	$[I_1^8]$	$[I_1^9]$	$[I_1^{10}]$	$[I_1^{11}]$	$[I_1^{12}]$	$[I_1^{13}]$	$[I_1^{14}]$	$[I_1^{15}]$	$[I_1^{16}]$	$[I_1^{17}]$	$[I_1^{18}]$
\mathcal{I}_2	I_2^1	I_2^2	I_2^3	I_2^4	I_2^5													
\mathcal{I}_3	I_3^1		I_3^2		I_3^3		I_3^4		I_3^5		I_3^6		I_3^7		I_3^8		I_3^9	



Learning with Experts (When t=3)

	$\mathcal{C}_3 = \{I_0^3, I_1^2\}$																	
t	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
\mathcal{I}_0	$[I_0^1]$	$[I_0^2]$	$[I_0^3]$	$[I_0^4]$	$[I_0^5]$	$[I_0^6]$	$[I_0^7]$	$[I_0^8]$	$[I_0^9]$	$[I_0^{10}]$	$[I_0^{11}]$	$[I_0^{12}]$	$[I_0^{13}]$	$[I_0^{14}]$	$[I_0^{15}]$	$[I_0^{16}]$	$[I_0^{17}]$	$[I_0^{18}]$
\mathcal{I}_1	$[I_1^1]$	$[I_1^2]$		$[I_1^3]$	$[I_1^4]$	$[I_1^5]$	$[I_1^6]$	$[I_1^7]$	$[I_1^8]$	$[I_1^9]$		$[I_1^{12}]$	$[I_1^{13}]$	$[I_1^{14}]$	$[I_1^{15}]$	$[I_1^{16}]$	$[I_1^{17}]$	$[I_1^{18}]$
\mathcal{I}_2	$[I_2^1]$			$[I_2^2]$			$[I_2^3]$			$[I_2^4]$		$[I_2^5]$						
\mathcal{I}_3				I_3^1				I_3^2			I_3^2		I_3^3					



Outline

Part I: Introduction to Fairness



Part II: Distribution Shift Undermines Fairness



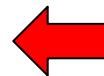
Part III: Mitigating Unfairness under Distribution Shift (Offline)



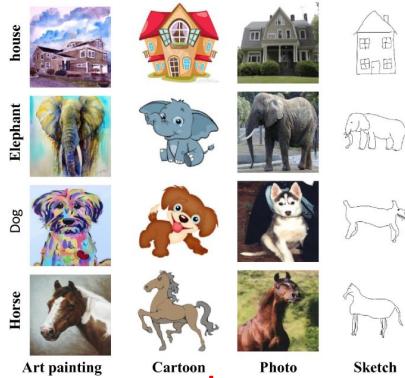
Part IV: Mitigating Unfairness under Distribution Shift (Online)



Part V: Open Challenges and Beyond



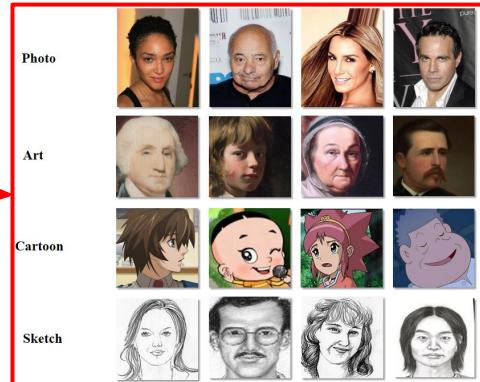
Benchmark Dataset for Fairness Learning under Distribution Shifts



PACS Dataset
• 4 Domains
• 7 Classes
• 9991 Images



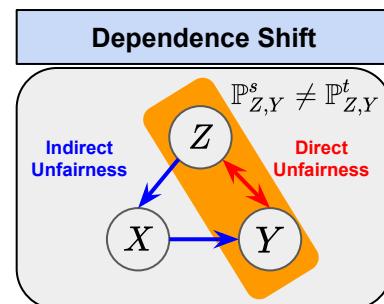
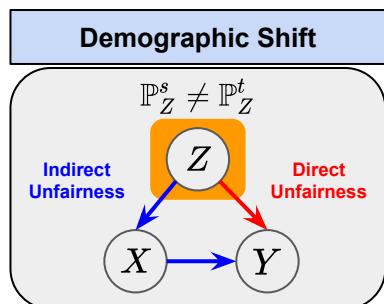
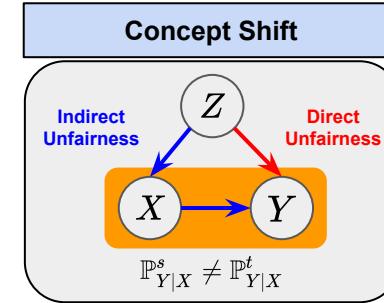
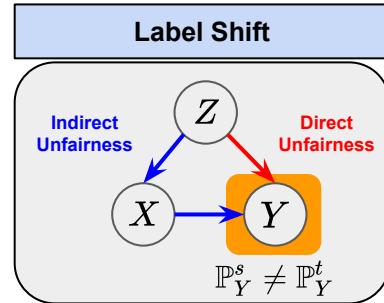
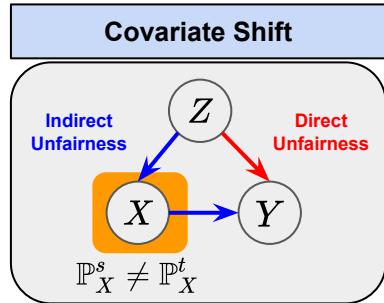
CelebA Dataset
• 40 Annotations
• 202K Images
• Fairness



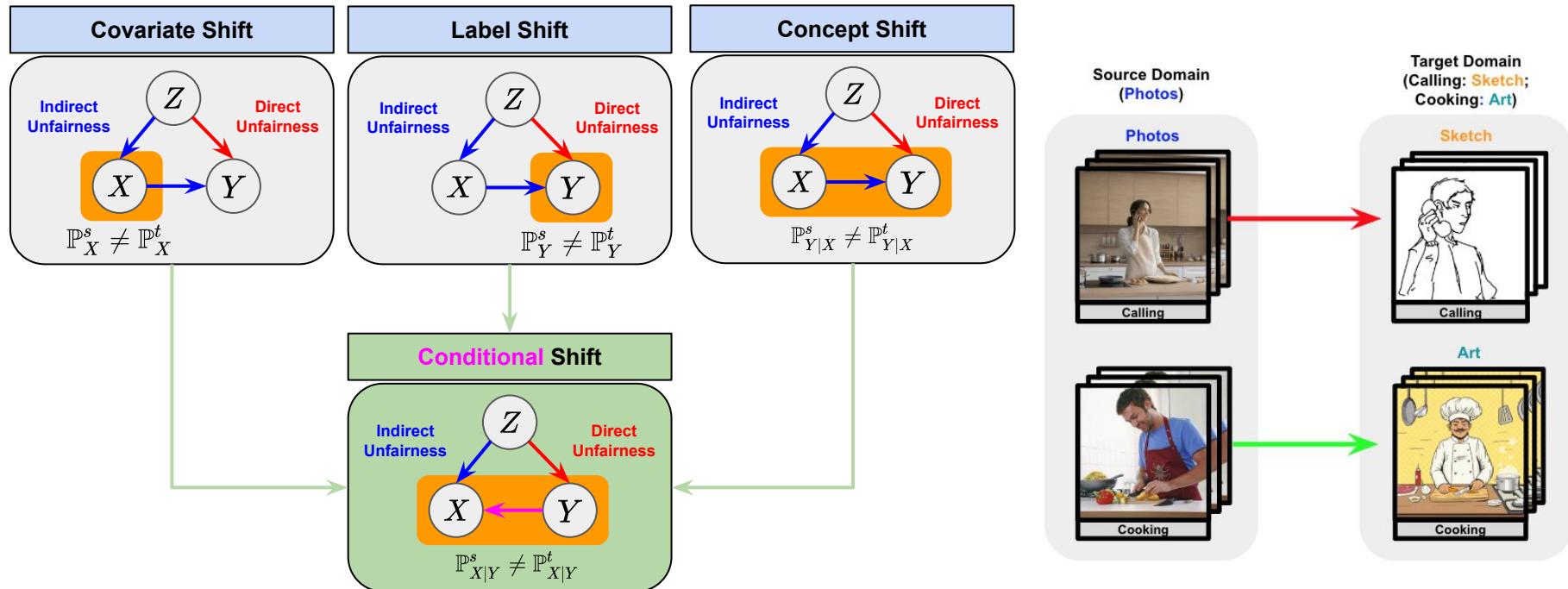
Open Challenge 1:
• Fairness-related Attribute Annotations
• 100K - 200K Images
• Multiple Domains
• Various Types of Shifts

1. Da Li, Yongxin Yang, Yi-Zhe Song, Timothy M. Hospedales. Deeper, Broader and Artier Domain Generalization. ICCV, 2017.
2. Liu, Ziwei and Luo, Ping and Wang, Xiaogang and Tang, Xiaoou. Deep Learning Face Attributes in the Wild. ICCV, 2015.

Algorithmic Fairness under Five Types of Distribution Shifts

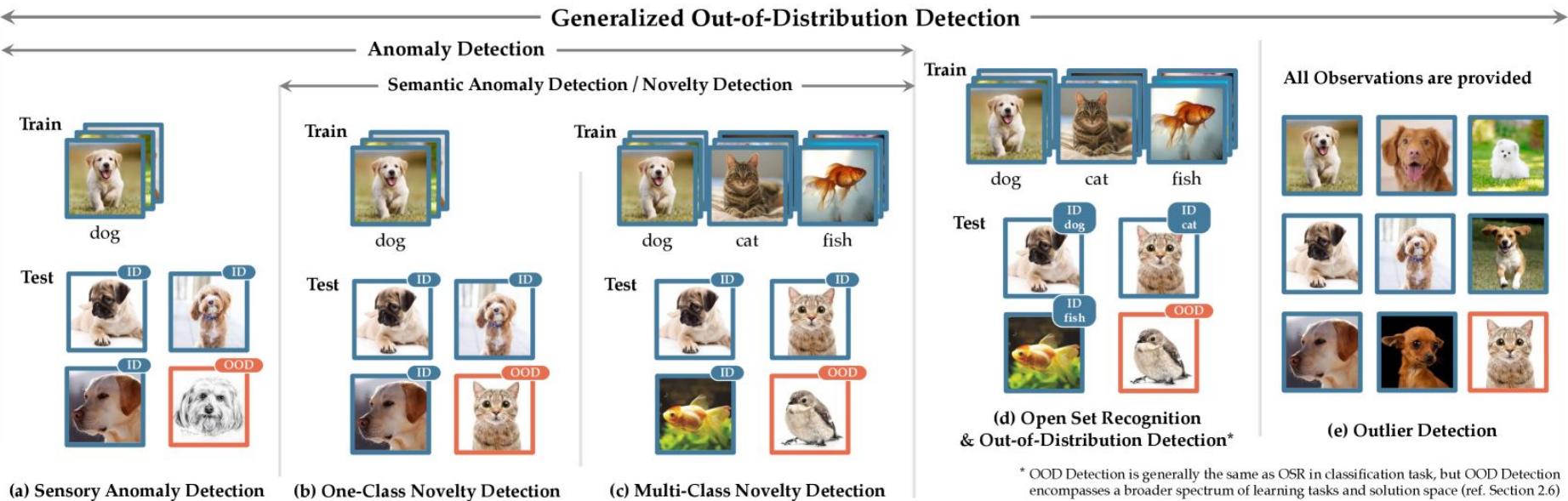


Algorithmic Fairness under Conditional Shift



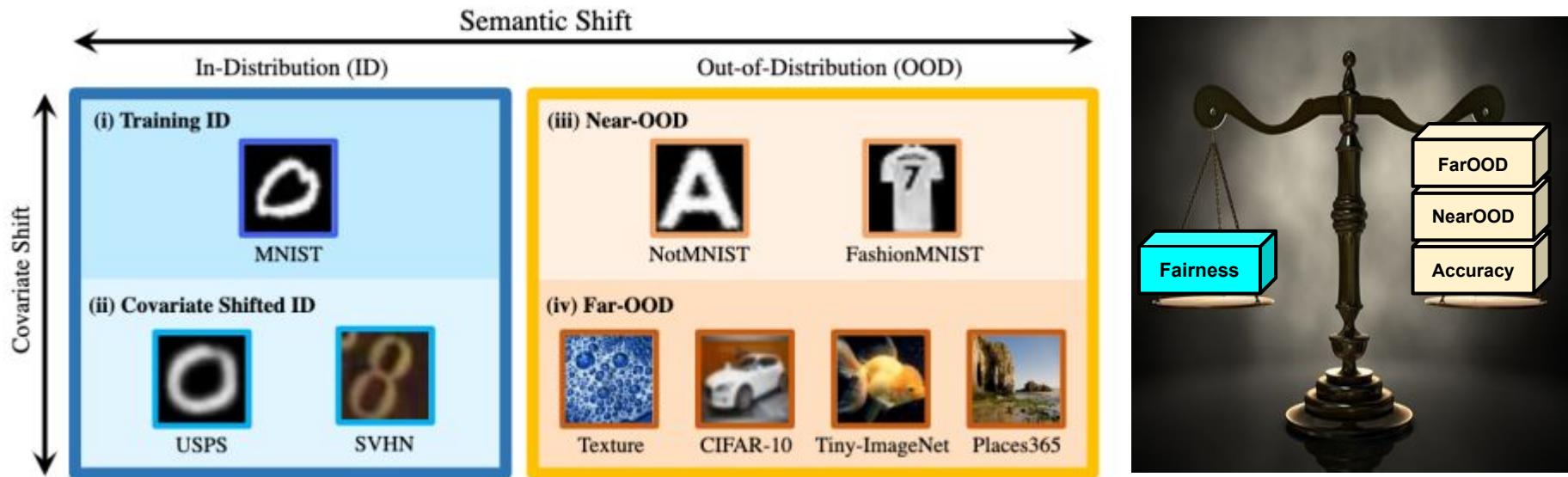
1. Xiaofeng Liu, Bo Hu, Linghao Jin, Xu Han, Fangxu Xing, Jinsong Ouyang, Jun Lu, Georges EL Fakhri, Jonghye Woo. Domain Generalization under Conditional and Label Shifts via Variational Bayesian Inference. IJCAI, 2021.

Algorithmic Fairness for Out-of-Distribution (OOD) Detection

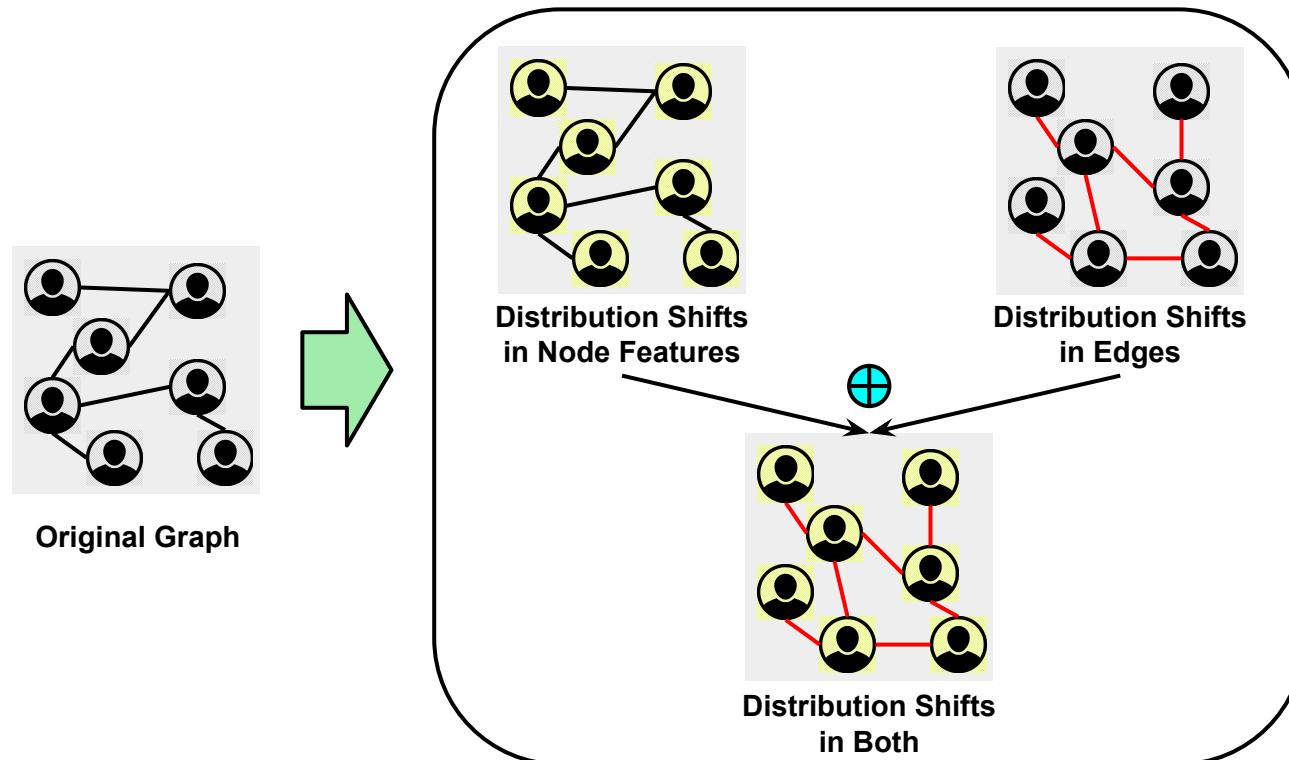


* OOD Detection is generally the same as OSR in classification task, but OOD Detection encompasses a broader spectrum of learning tasks and solution space (ref. Section 2.6)

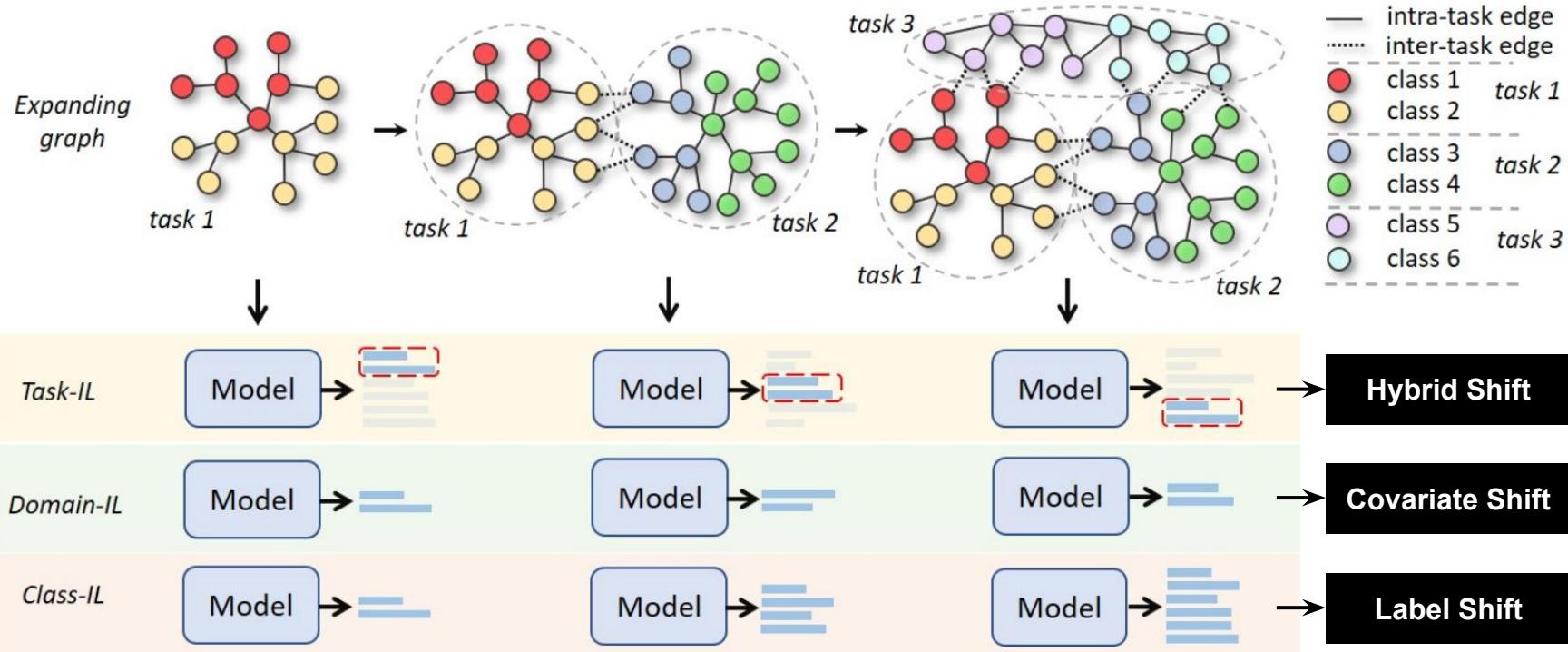
Algorithmic Fairness for Out-of-Distribution (OOD) Detection



Algorithmic Fairness Domain Generalization on Graphs



Algorithmic Fairness on Incremental Learning on Graphs



Supervised Algorithmic Fairness in Distribution Shifts

Tutorial at the 2024 IEEE International Conference on Big data (IEEE BigData 2024)



Chen Zhao¹
chen_zhao@baylor.edu



Xintao Wu²
xintaowu@uark.edu

¹Department of Computer Science, Baylor University

²Department of Electrical Engineering and Computer Science, University of Arkansas