

# Towards Trustworthy AI with Perspectives on Fairness, Interpretability, Privacy, and Security

2023-05-18

Sungwon Han



Microsoft

**KAIST**

**ibS** 데이터사이언스그룹  
Data Science Group

# Trustworthy AI

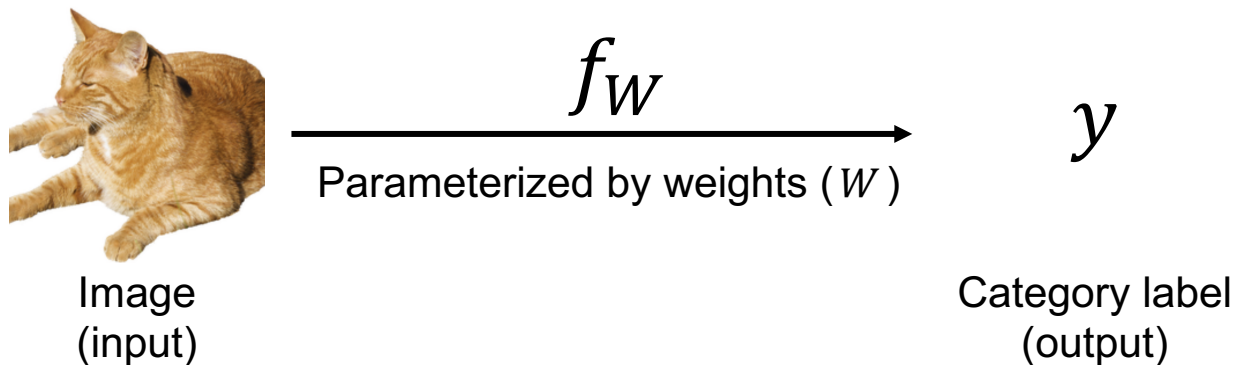
# Trustworthy “Deep Learning”

⇒ Deep learning is a subset of machine learning, which is essentially a neural network with three or more layers.

# I AI with Deep Learning

## How do deep neural networks train and make inference?

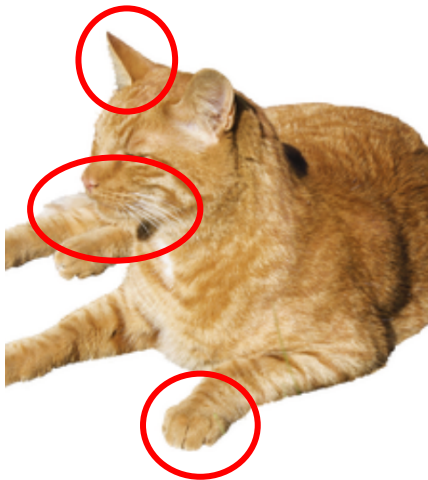
It varies depending on the task... Let's consider a supervised, image classification problem. We aim to learn a mapping function  $f$  that takes an image  $x$  and produces a label  $y$ .



# I AI with Deep Learning

## How is the prediction generated?

There are various features that can be used to distinguish the image!  
Which features would deep learning model select for the prediction?



Ears? Whiskers? Paw?



Texture?

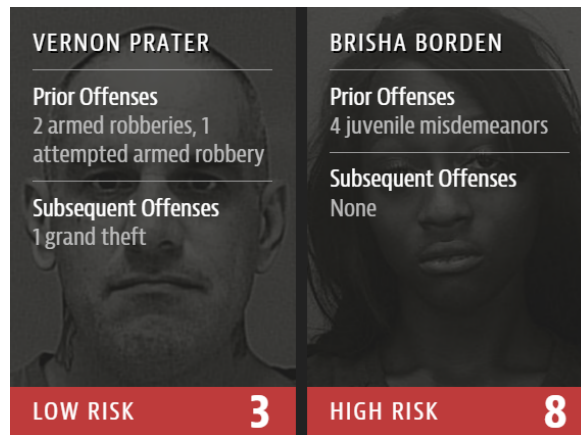
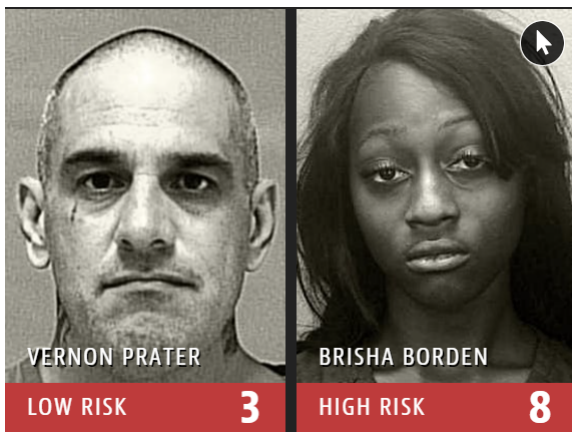


How about this?

# I AI with Deep Learning

## Consider another example..

AI model is trying to predict the recidivism risk from given individuals.  
Which features would deep learning model select for the prediction?



# “Trustworthy” AI

⇒ Describe AI that is lawful, ethically adherent,  
and technically robust.

# I Trustworthy AI

## Increasing public interest on trustworthy AI

- Domains in high-stakes decisions are already using AI models, bringing real risks. How do we ensure safety is built into these systems?
- High-risk domains include those related to safety infra and products, education, employment, justice, immigration, and climate protection (as defined in EU AI Act proposal, Article 6-51).

“

**On artificial intelligence,  
trust is a must, not a nice to have.**

”



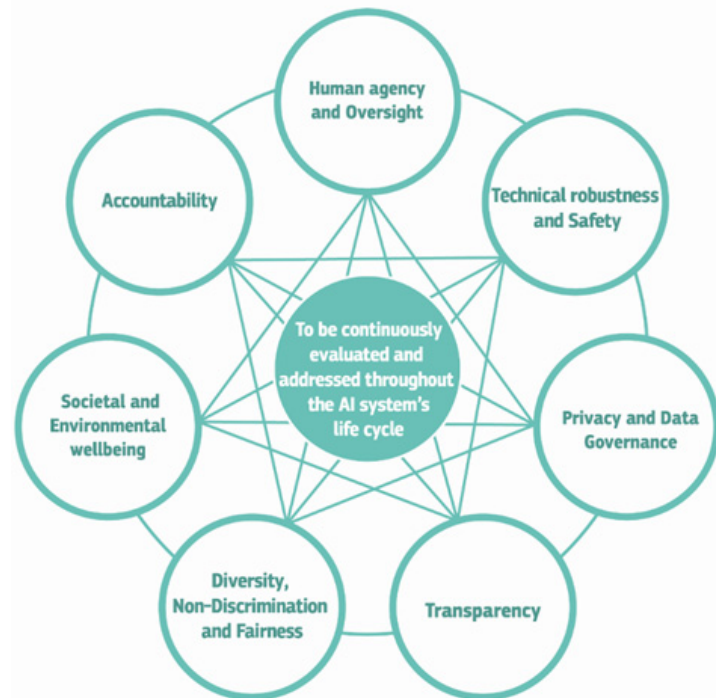


# I Requirements of Trustworthy AI

## Ethics Guidelines for AI from European Commission

Interrelationship of the seven requirements:

- All are of equal importance and should be implemented and evaluated throughout the AI system's lifecycle
- For some applications, they some requirements may be of lesser or greater relevance.



## Ethics Guidelines for AI from European Commission

1. Human agency and oversight
2. Diversity, non-discrimination
3. Transparency
4. Privacy and data governance
5. Technical robustness, security
6. Societal and environmental wellbeing
7. Accountability

## In today's presentation...

1. Human agency and oversight
2. **Diversity, non-discrimination**

➔ Sec. 1) Algorithmic fairness (WWW'23)

3. **Transparency**

➔ Sec. 2) Interpretability (NeurIPS'22)

4. **Privacy and data governance**

➔ Sec. 3) Federated learning (ECCV'22)



**DualFair: Fair Representation Learning at Both Group and Individual Levels via Contrastive Self-supervision (WWW' 2023)**

# I Introduction

## Fairness in machine learning

Why we care about fairness?

\*Recidivism rate

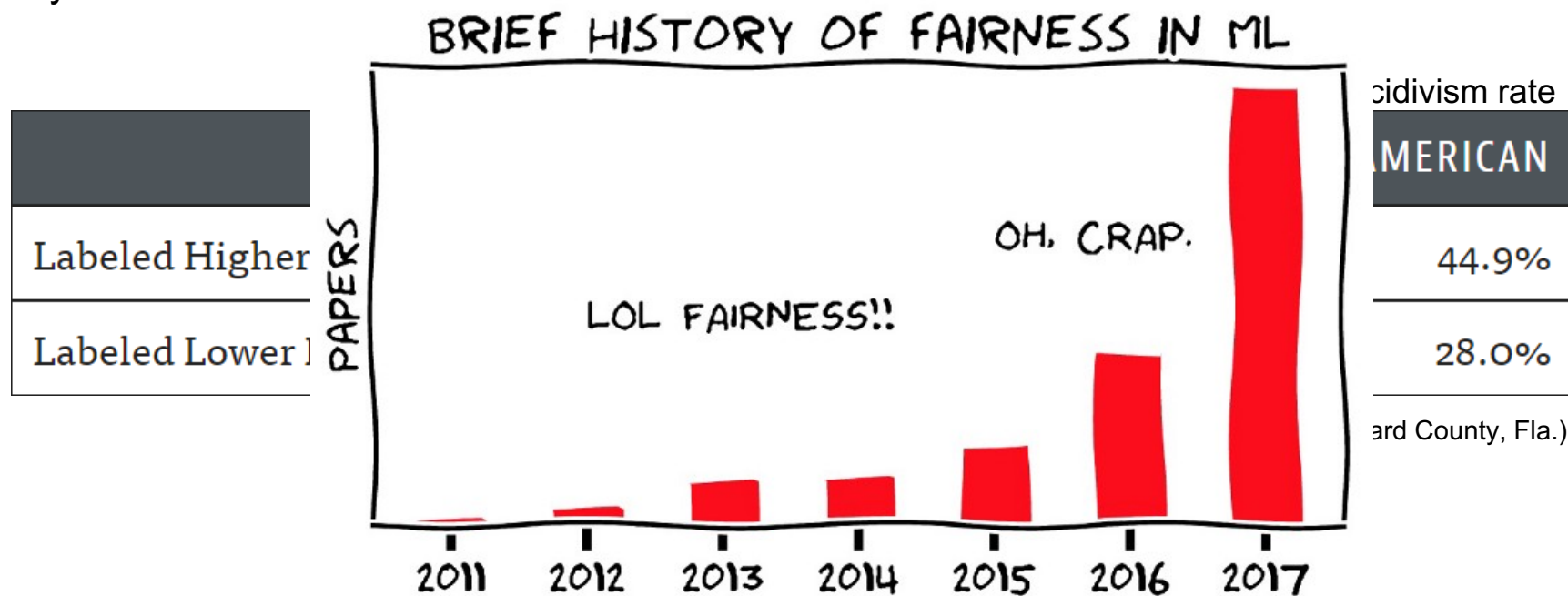
	WHITE	AFRICAN AMERICAN
Labeled Higher Risk, But Didn't Re-Offend	23.5%	44.9%
Labeled Lower Risk, Yet Did Re-Offend	47.7%	28.0%

(Source: ProPublica analysis of data from Broward County, Fla.)

# I Introduction

## Fairness in machine learning

Why we care about fairness?



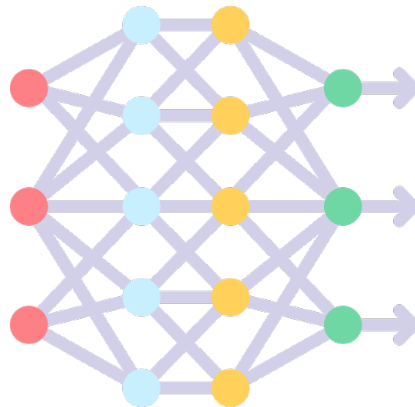
# I Background

## Fair representation learning

Debias sensitive information and generate low-dimensional representation



Private dataset  
with Sensitive information

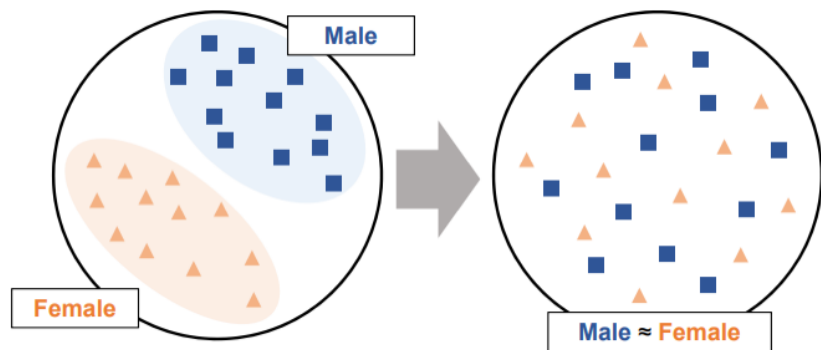


Public dataset  
without Sensitive information

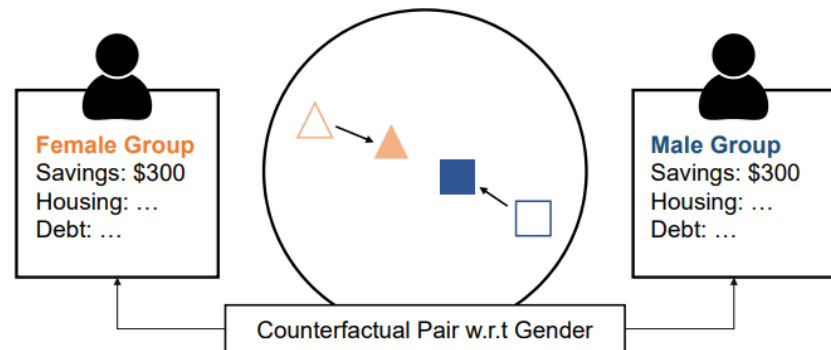
# Motivation: Two fairness criteria

Fairness should be achieved at both group and individual-level

Group fairness vs. Counterfactual fairness in representation learning



**Group fairness:** group-level fairness  
Ex) Groups are indistinguishable in embeddings



**Counterfactual fairness:** individual-level fairness  
Ex) Counterfactual pair from groups should be close



# I Motivation: Contrastive learning

## Contrastive learning for representation learning

Generalized-InfoNCE objective can be decomposed into two terms:

$$L_c(\mathbf{x}, \mathcal{X}_+, \mathcal{X}_-) = -\frac{1}{|\mathcal{X}_+|} \underbrace{\left( \log \sum_{\mathbf{x}' \in \mathcal{X}_+} \exp(\text{sim}(f(\mathbf{x}), f(\mathbf{x}'))/\tau) \right)}_{\text{Alignment loss}} - \underbrace{\log \sum_{\mathbf{x}' \in (\mathcal{X}_+ \cup \mathcal{X}_-)} \exp(\text{sim}(f(\mathbf{x}), f(\mathbf{x}'))/\tau)}_{\text{Distribution loss}}$$

- Alignment loss encourages the embeddings of positive pairs to be placed closer.
- Distribution loss matches all instances' embeddings into the prior with high entropy.

$$L_{\text{gen-c}}(\mathbf{x}, \mathbf{x}_+, \mathcal{X}_-) = -L_{\text{align}}(\mathbf{x}, \mathbf{x}_+) + \text{SWD}(\tilde{\mathcal{Z}}, \mathcal{Z}_{\text{prior}})$$

# I Main idea: Fair contrastive learning

## Contrastive objective for both group & counterfactual fairness

Treat counterfactual pair alike and ensure non-distinguishable embeddings among groups

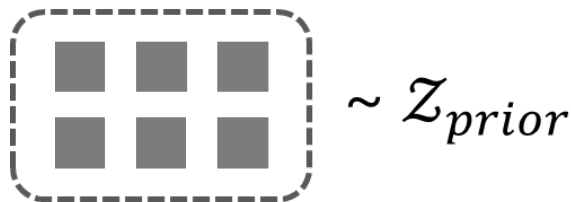
Maximize Agreement



■ : Input sample (Anchor)

▲ : Counterfactual pair of input (Positive)

Distribution Matching



■ : Random sample from the same group with Anchor (Negative)

$$L = -L_{\text{align}}(\blacksquare, \blacktriangle) + \text{SWD}(\blacksquare, \mathcal{Z}_{\text{prior}})$$

# I Main idea: Fair contrastive learning

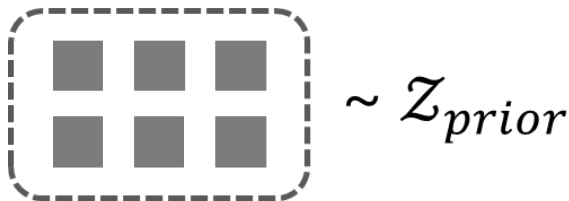
## Contrastive objective for both group & counterfactual fairness

Treat counterfactual pair alike and ensure non-distinguishable embeddings among groups

Maximize Agreement



Distribution Matching



■ : Input sample (Anchor)

▲ : Counterfactual pair of input (Positive)

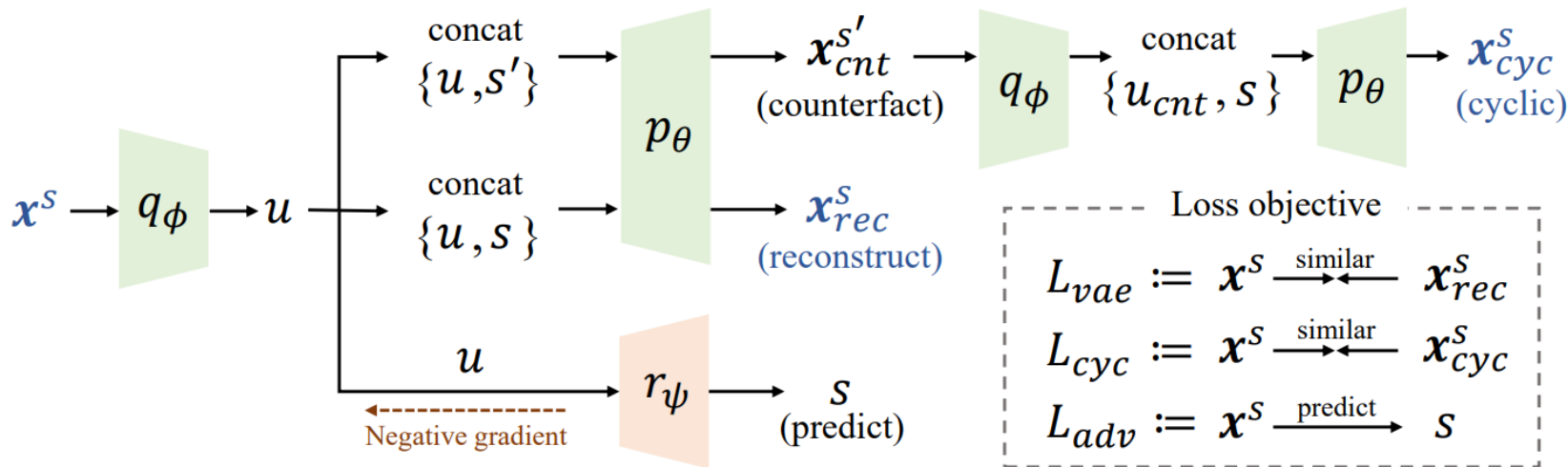
■ : Random sample from the same group with Anchor (Negative)

$$L = -L_{\text{align}}(\blacksquare, \blacktriangle) + \text{SWD}(\blacksquare, \mathcal{Z}_{\text{prior}})$$

# I Main idea: C-VAE

## Generation of counterfactual samples with C-VAE

Variational autoencoder with adversarial training for counterfactual sample generator



# I Main idea: Fair contrastive learning

## Contrastive objective for both group & counterfactual fairness

Treat counterfactual pair alike and ensure non-distinguishable embeddings among groups

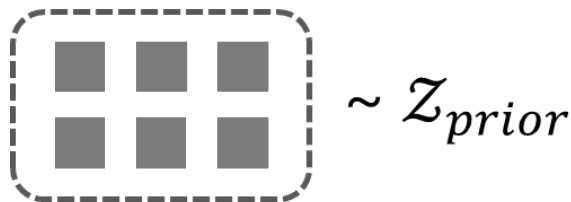
Maximize Agreement



■ : Input sample (Anchor)

▲ : Counterfactual pair of input (Positive)

Distribution Matching



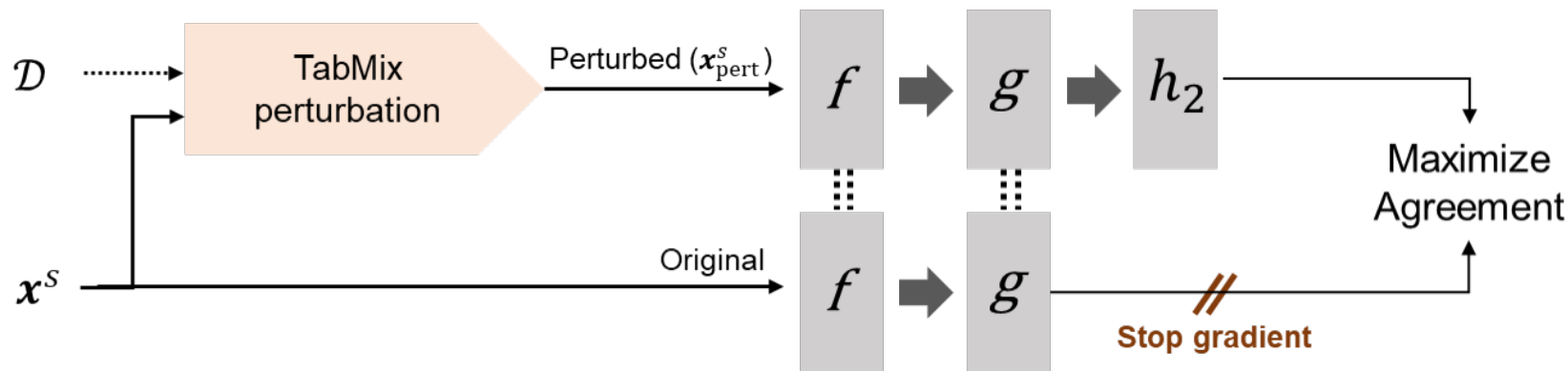
■ : Random sample from the same group with Anchor (Negative)

$$L = -L_{\text{align}}(\blacksquare, \blacktriangle) + \text{SWD}(\blacksquare, \mathcal{Z}_{\text{prior}})$$

# ■ Main idea: Self-knowledge distillation

## Self-knowledge distillation to maintain representation quality

Reduce the discrepancy between original & perturbations to learn data semantics



$$L_{\text{total}} = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x}^s \in \mathcal{D}} (L_{\text{fair-cl}}(\mathbf{x}^s) + L_{\text{self-kd}}(\mathbf{x}^s))$$

# I Experiments

## Datasets

Six fairness-required datasets with various kinds of downstream tasks for evaluation

Dataset	# samples	# attr.	Sensitive attr.	Split	Task
Adult	48,842	14	gender, race	2:1	classification
Credit	1,000	20	gender	4:1	classification
COMPAS	6,172	7	gender, race	4:1	classification
LSAC	22,407	12	gender, race	4:1	classification
Students	649	33	gender	4:1	regression
Communities	1,994	128	race	4:1	regression

# I Experiments

## Performance evaluation

Performance comparison summaries among fairness-aware baselines and DualFair. Averaged rank for each evaluation metric across six datasets is reported.

Method	AUC/RMSE	$\Delta DP$	$\Delta EO$	$\Delta CP$	Total
VFAE	3.2	5.8	5.0	4.2	4.6
LAFTTR	<b>2.7</b>	6.2	5.5	2.8	4.3
MIFR	6.0	3.0	3.8	4.6	4.4
L-MIFR	5.3	3.2	3.3	3.8	3.9
C-InfoNCE	3.2	4.5	3.3	5.2	4.1
WeaC-InfoNCE	4.7	2.8	4.3	4.8	4.2
<b>DualFair</b>	3.0	<b>2.5</b>	<b>2.5</b>	<b>2.2</b>	<b>2.6</b>

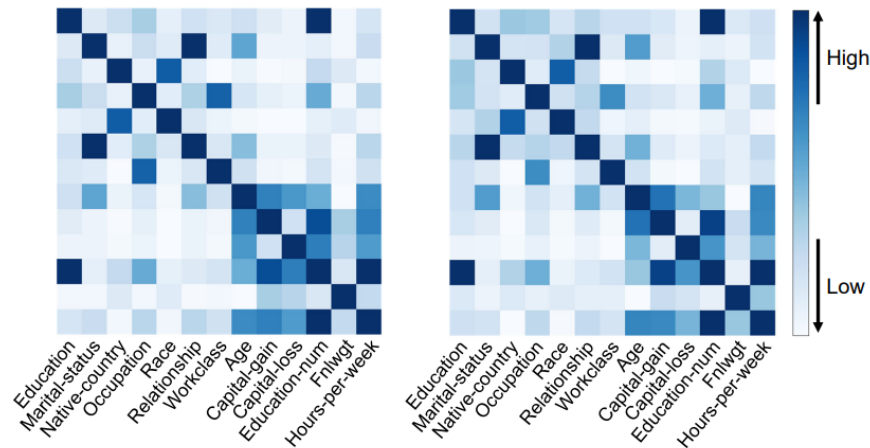


# Experiments

## Quality of counterfactual samples

The original relationship between features is well-maintained in the counterfactual samples.

Training set	Adult		Credit		Compas		LSAC	
	AUC	F1	AUC	F1	AUC	F1	AUC	F1
Original	0.91	0.78	0.76	0.68	0.74	0.68	0.85	0.65
Counterfactual	0.89	0.76	0.75	0.62	0.73	0.67	0.84	0.63



(a) Original UCI Adult

(b) Counterfactual synthetic data

# I Conclusion

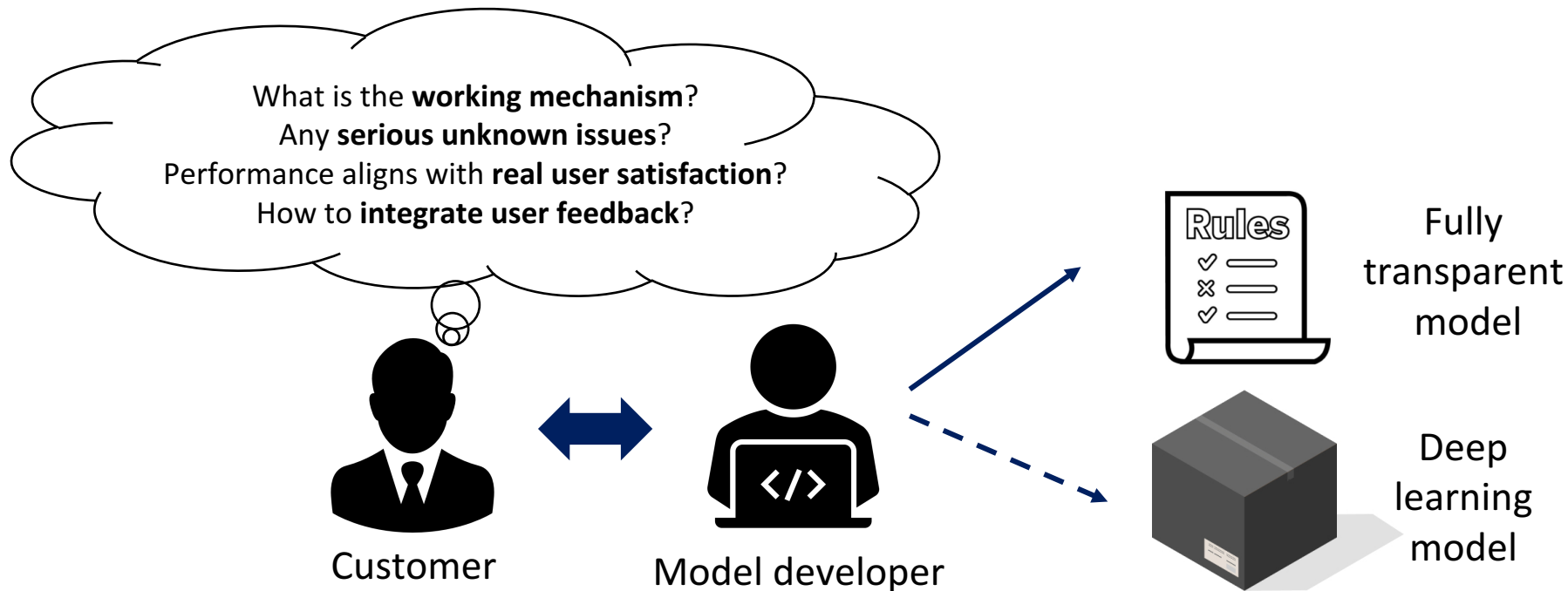
- We propose a self-supervised learning framework, **DualFair**, that simultaneously **debias sensitive attributes at both group and individual levels**.
- We introduce the C-VAE model to generate counterfactual samples and propose **fairness-aware contrastive loss** to meet the two fairness criteria jointly.
- We design the **self-knowledge distillation loss** to maintain representation quality by minimizing the embedding discrepancy between original and perturbed instances.
- Experiments confirm that DualFair generates a fair embedding with high representation quality. We further show **a synergistic effect of the two fairness criteria**.



# Self-explaining deep models with logic rule reasoning (NeurIPS' 2022)

# I Background

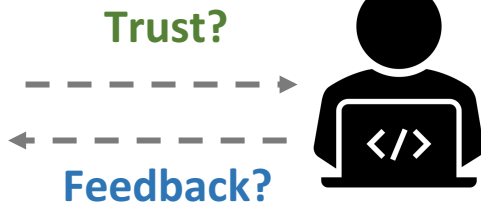
## Trust Issues with Deep learning models



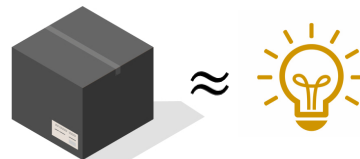
# I Motivation

## Limitation of Post-Hoc Explanations

Post-hoc explanations



**Can we trust the explanations?**



- Always an approximation [1]
- “General uneasiness” of practitioners [2]

**How to integrate user feedback?**

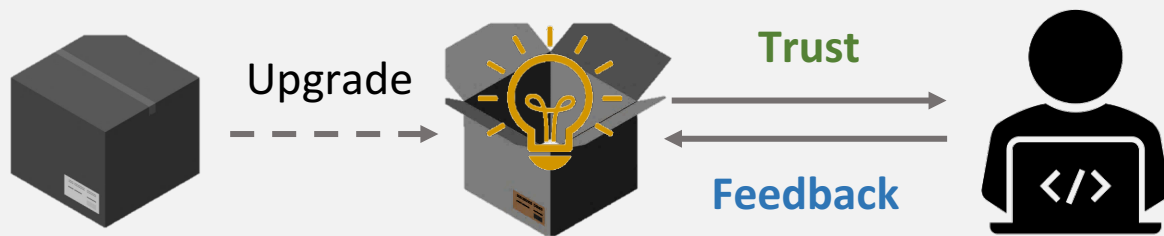
- No systematic method for direct control
- Requires model retraining
- No guarantee for satisfying user demands

[1] “Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead”, *Nature Machine Intelligence*, 2019

[2] “Human Factors in Model Interpretability: Industry Practices, Challenges, and Needs”, *ACM HCI 2020*

# I Main Approach

## SELOR: Self-Explaining with **L**OGic rule **R**easoning

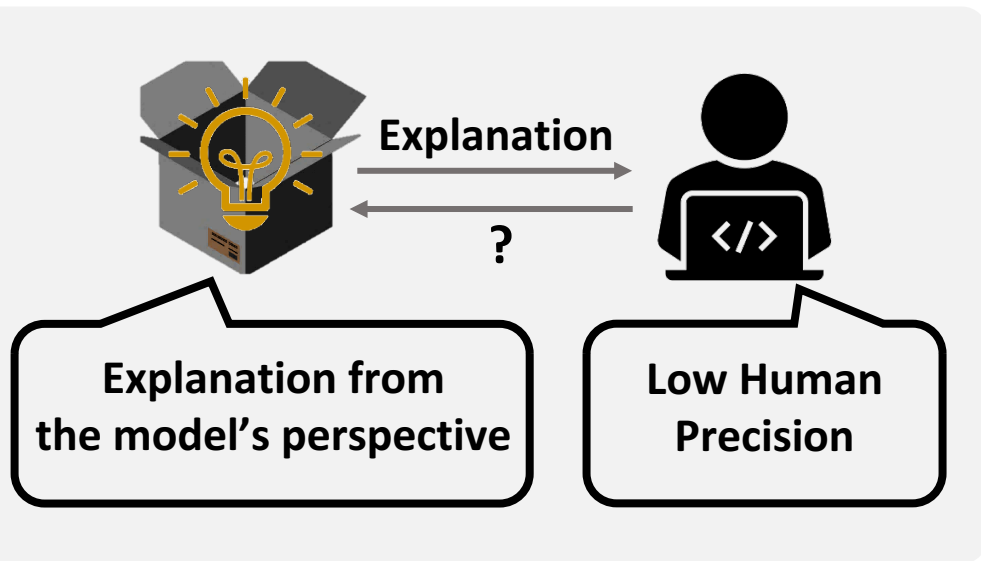


Lays the **foundation** for close collaboration

- **Trust**: explanations **faithful** to the model
- **Feedback**: explanations as **handle** for control

# I Main Approach

## SELOR: Self-Explaining with LOGic rule Reasoning



### Human Precision:

Whether the explanation naturally leads to the prediction according to human perception

#### Low Human Precision:

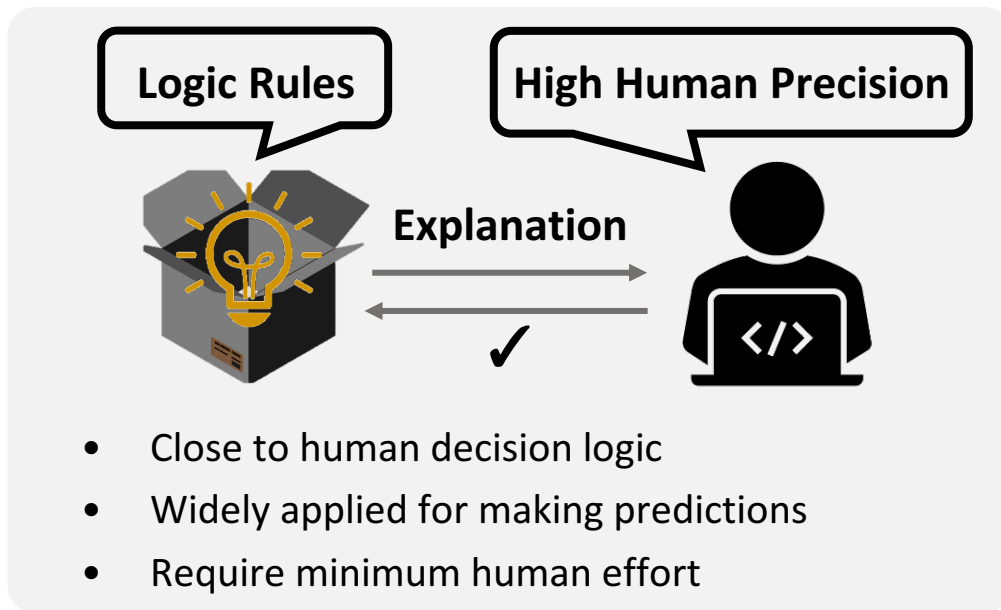
*is, an => positive sentiment*

#### High Human Precision:

*Awesome => positive sentiment*

# I Main Approach

## SELOR: Self-Explaining with LOGic rule Reasoning



*awesome AND tasty*

*Antecedent*

(condition to apply)



*positive sentiment*

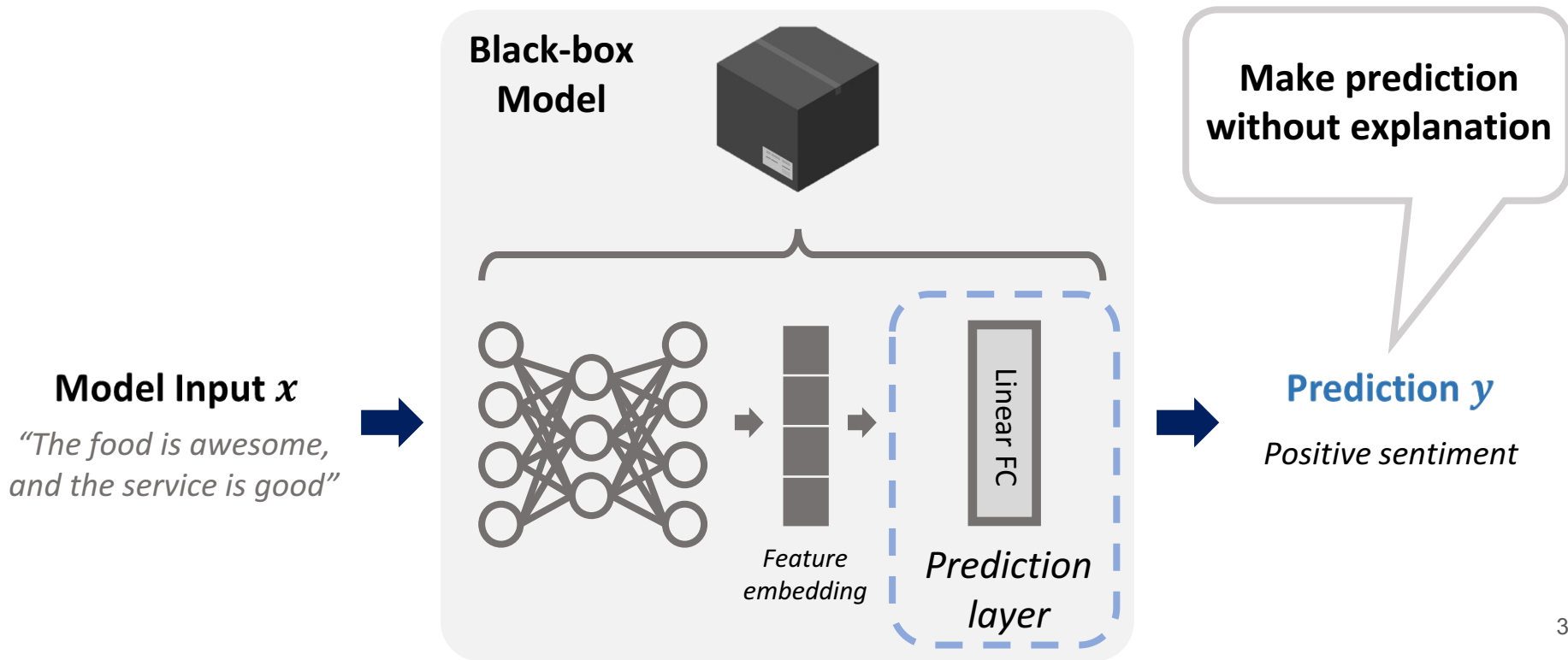
*Consequent*

(prediction result)



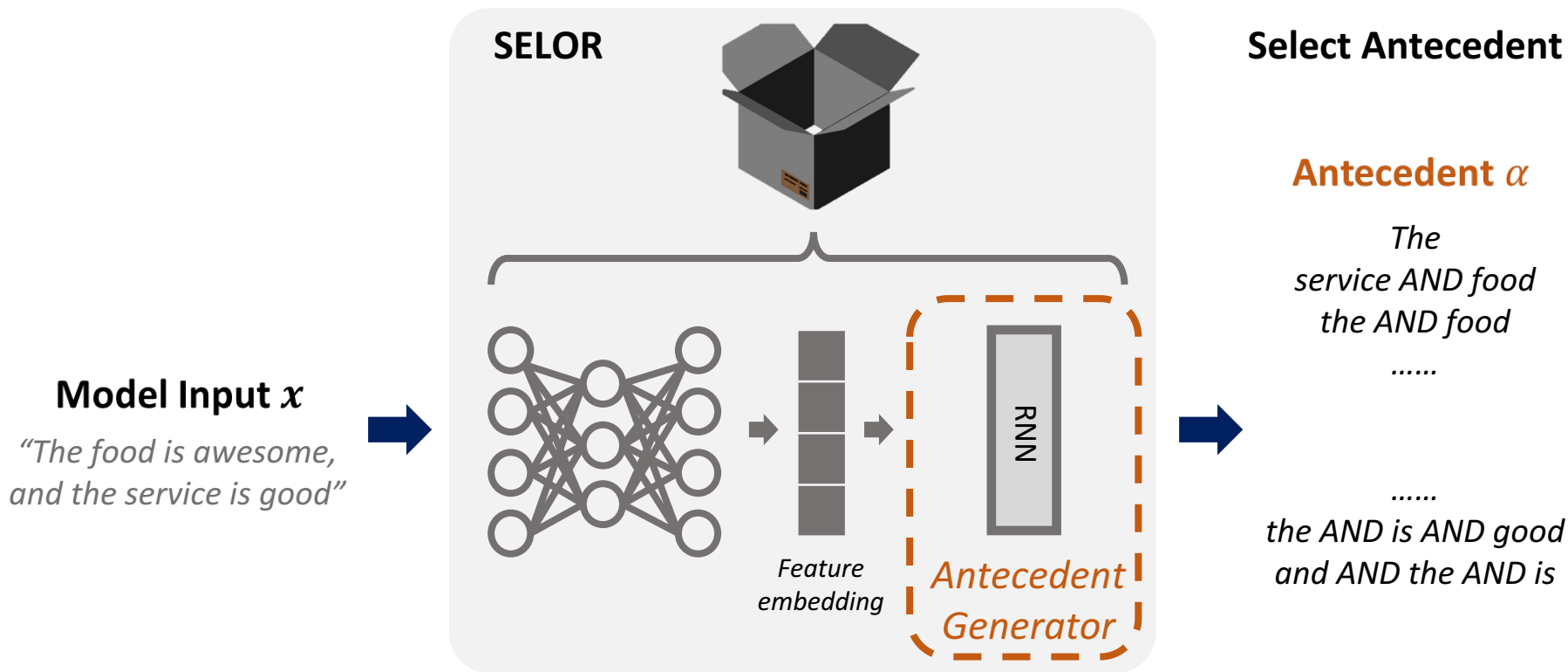
# I Main Model Framework

## SELOR Framework: Black-box Model



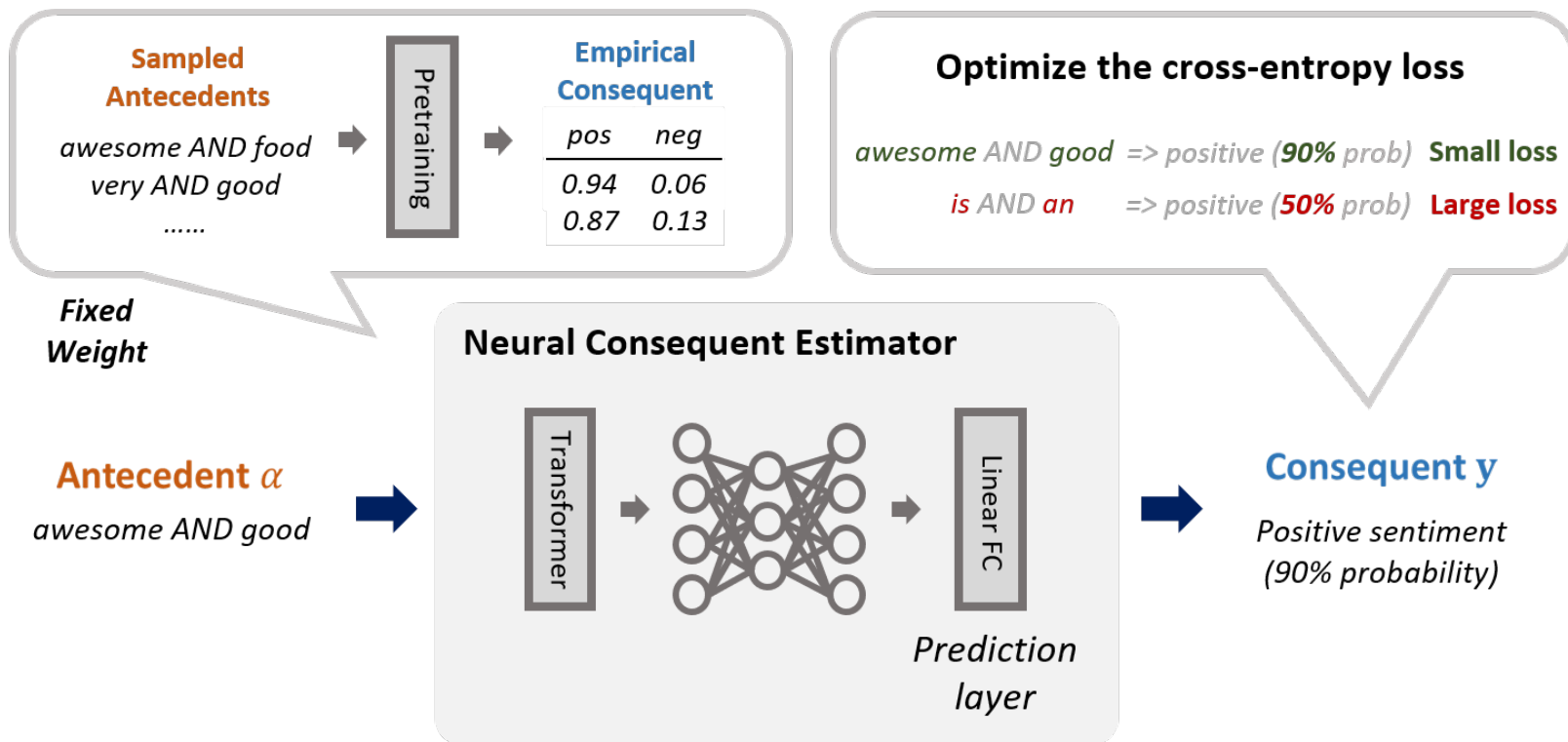
# I Main Model Framework

## SELOR Framework: Antecedent Generator



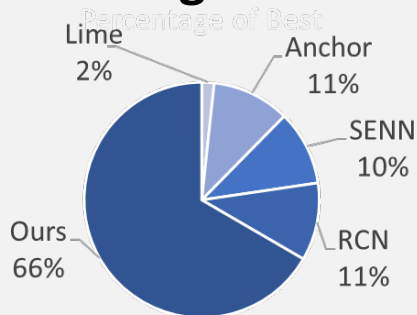
# I Main Model Framework

## SELOR Framework: Consequent Estimator



# I Results

## High Human Precision



User study  
Percentage of best  
**+500%**  
(Adult dataset)

## Good Prediction Performance



**SELOR**

≈



**Black-box**

## Training Cost

- **Efficient, differentiable** training
- **Slightly slower** than black-box model



**SELOR**



**BERT**

### Training Time



# I Conclusion

- We propose **SELOR**, which incorporates self-explanatory capabilities into a deep model to **provide high human precision by explaining logic rules**.
- SELOR does not require predefined rule sets and can be learned in a differentiable way.
- Extensive tests involving human evaluation show that our method achieves **high prediction performance and human precision in explanation**.

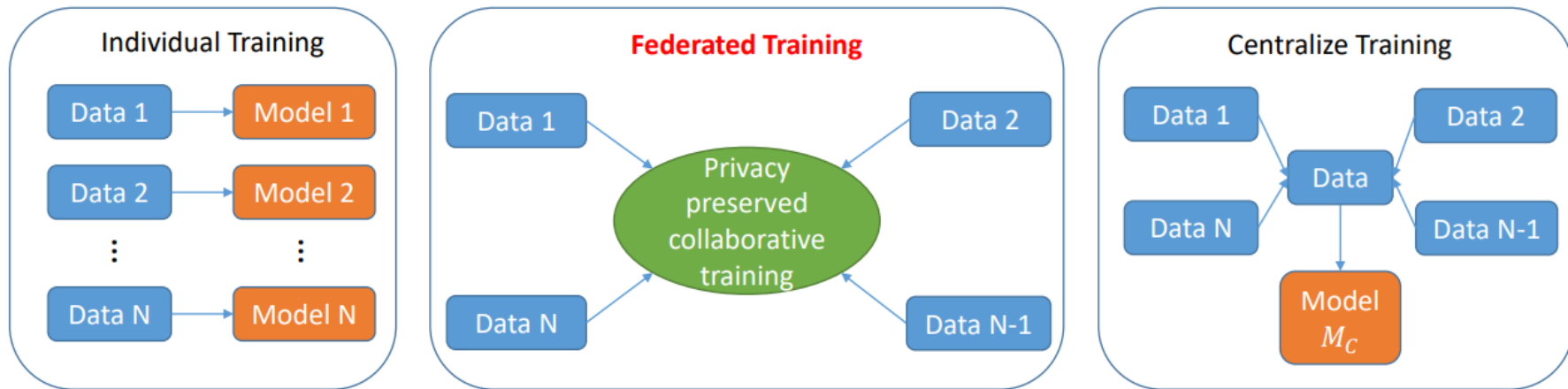


## **FedX: Unsupervised Federated Learning with Cross Knowledge Distillation (ECCV' 2022)**

# Background

## Federated learning

- Concerns on user data privacy and confidentiality.
- Inability to build an ML model due to inadequate data or training cost on ML implementation of the computational cost involved for training an ML model.



## Federated learning: problem statement

- Assume that data are distributed over every different party.
- Parties aim to train a single model  $F$  that can apply for various downstream tasks with the help of central server.
- Then, the global objective function to solve is as follows:

$$\arg \min_{\phi} \mathcal{L}(\phi) = \sum_{m=1}^M \frac{|\mathcal{D}^m|}{|\mathcal{D}|} \mathcal{L}_m(\phi),$$

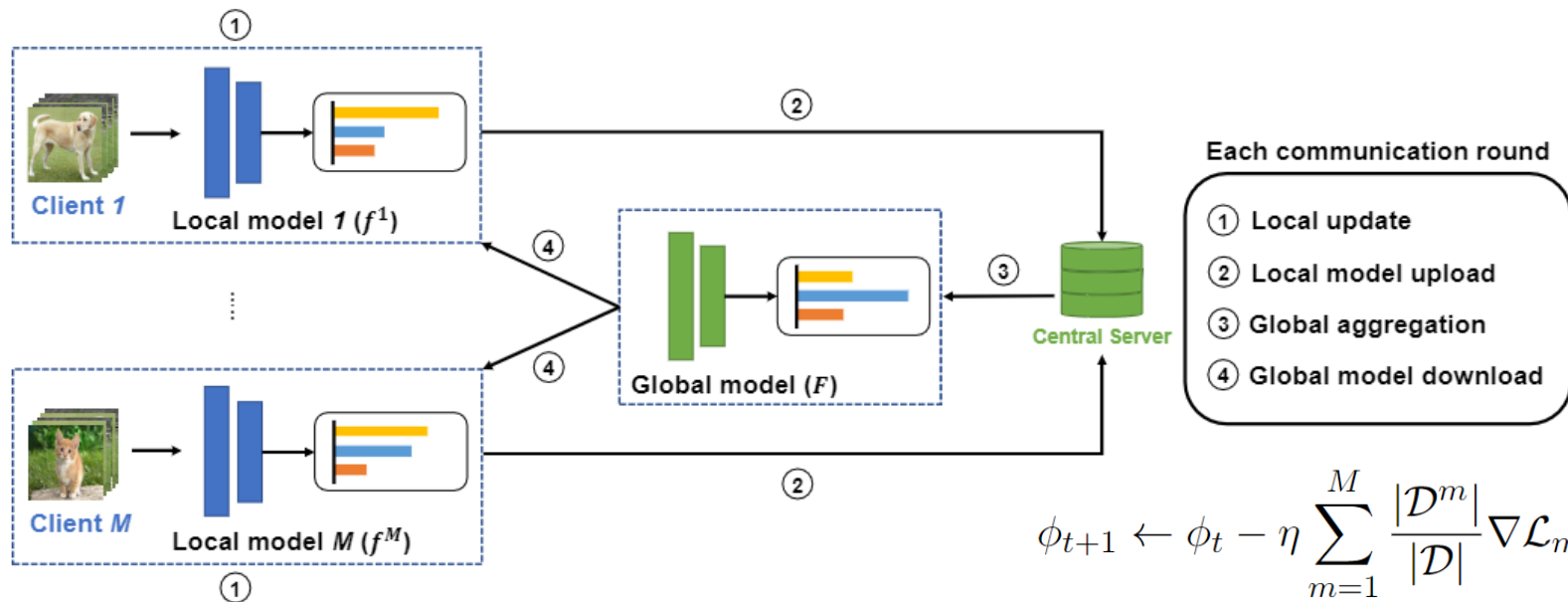
where  $\mathcal{L}_m(\phi) = \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}^m} [l_m(\mathbf{x}, y; \phi)]$ .



# Background

## Basic architecture for federated learning: FedAvg

- Four processes run in each communication round



## Extension to unsupervised federated learning

- We can run unsupervised representation learning on this federated framework.
- Then, the global objective function to solve is as follows:

$$\arg \min_{\phi} \mathcal{L}(\phi) = \sum_{m=1}^M \frac{|\mathcal{D}^m|}{|\mathcal{D}|} \mathcal{L}_m(\phi),$$

$$\text{where } \mathcal{L}_m(\phi) = \mathbb{E}_{\mathbf{x} \in \mathcal{D}^m} [l_m(\mathbf{x}; \phi)].$$

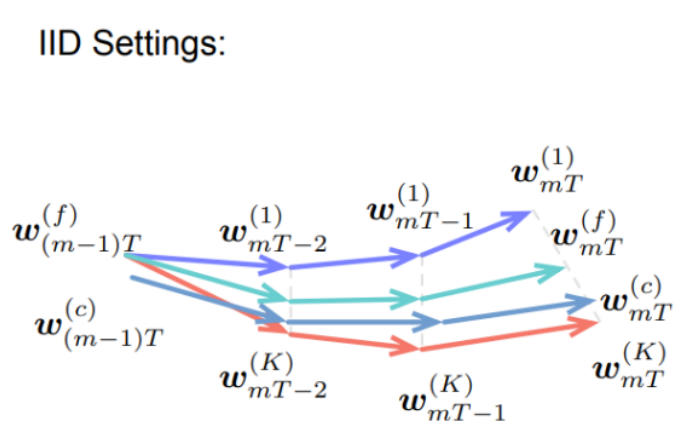
- Loss function could be excerpted from
  - InfoNCE-based model (e.g., SimCLR, MoCo, NNCLR)
  - Asymmetric siamese-based model (e.g., BYOL, SimSiam)
  - ...

# Motivation

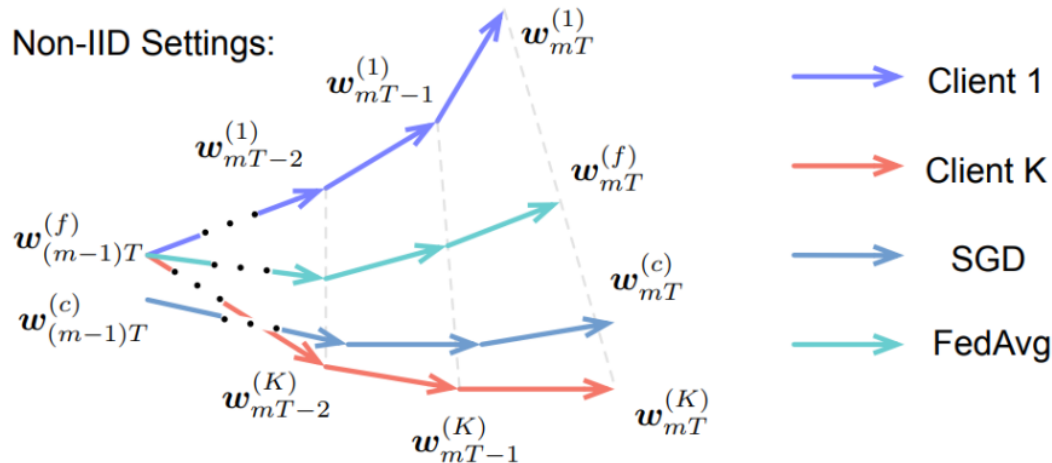
## Challenges in federated learning

- Non-IID distribution of local-data leads to the biased results

IID Settings:



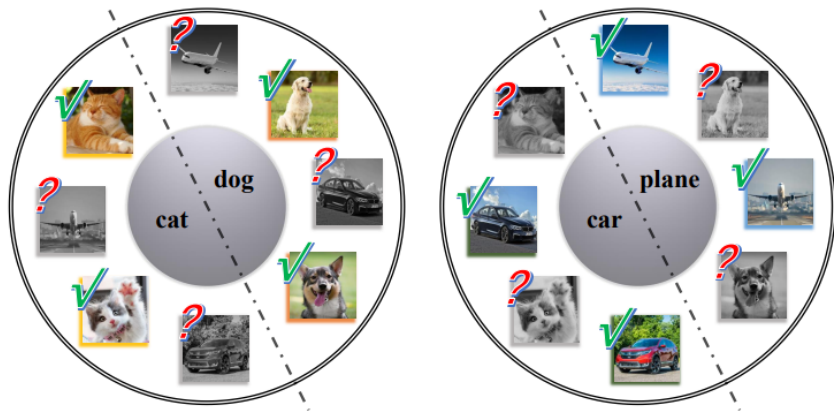
Non-IID Settings:



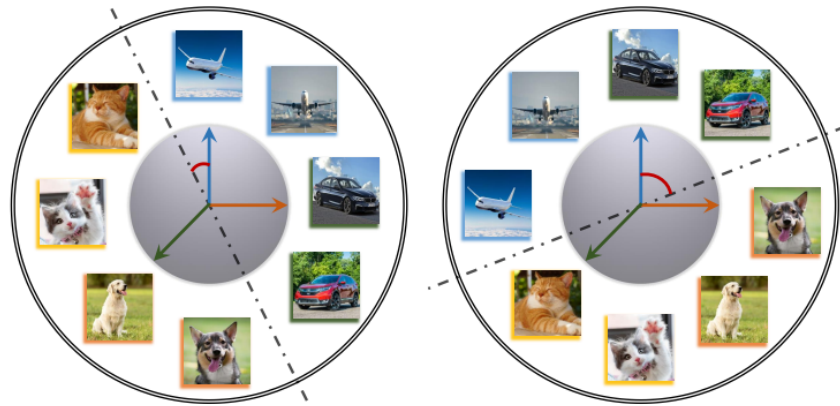
# Motivation

## Extra challenges in “unsupervised” federated learning

- Non-IID distribution without ground-truth labels amplifies the embedding divergence



(a) Inconsistency of representation spaces.



(b) Misalignment of representations.

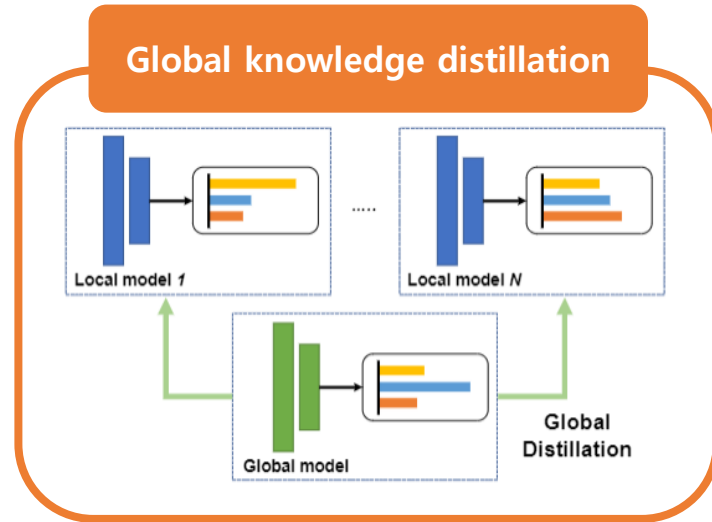
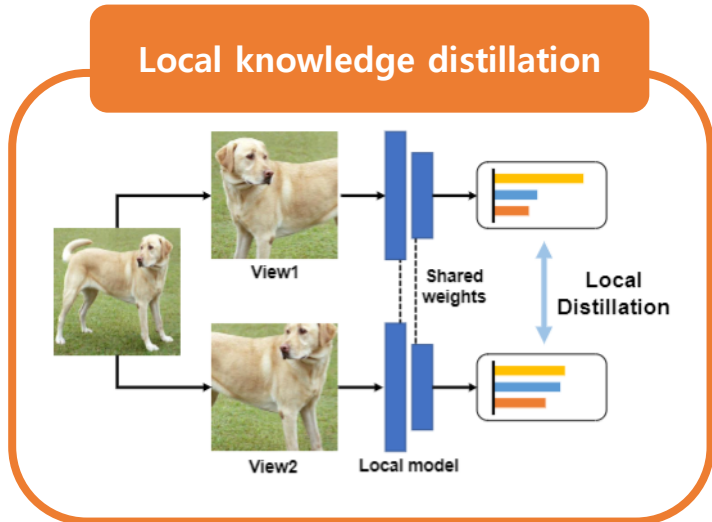
# ■ Main Idea: knowledge distillation

Knowledge distillation to convey global knowledge to local client!

- Knowledge distillation refers to the idea of model compression by teaching a smaller network, step by step, exactly what to do using a bigger already trained network.
- **Global model** has knowledge from entire data distribution (usually has higher performance), which can be regarded as a **teacher model**.
- **Local model** has knowledge only from locally biased data and needs entire data's information, which can be regarded as a **student model**.
- Distill knowledge from global model to regularize the local model's training in an unsupervised fashion!

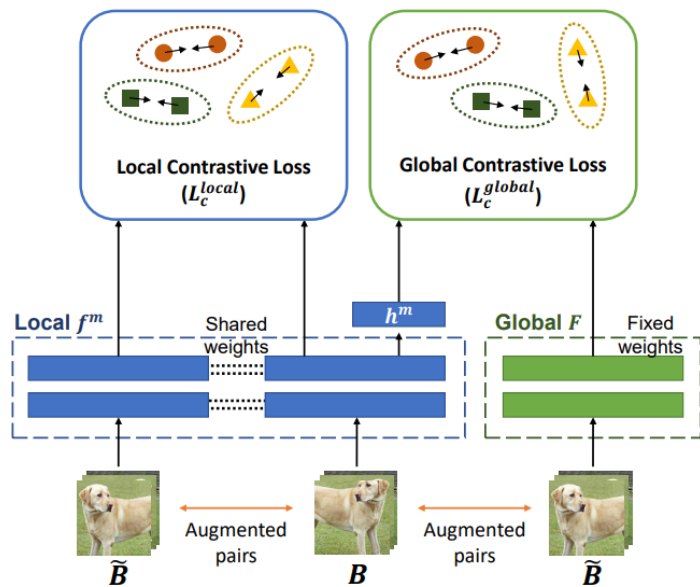
# ■ Main Idea: Two-sided knowledge distillation

Learns representation from local data and refines the central server's knowledge via **two-sided knowledge distillation**

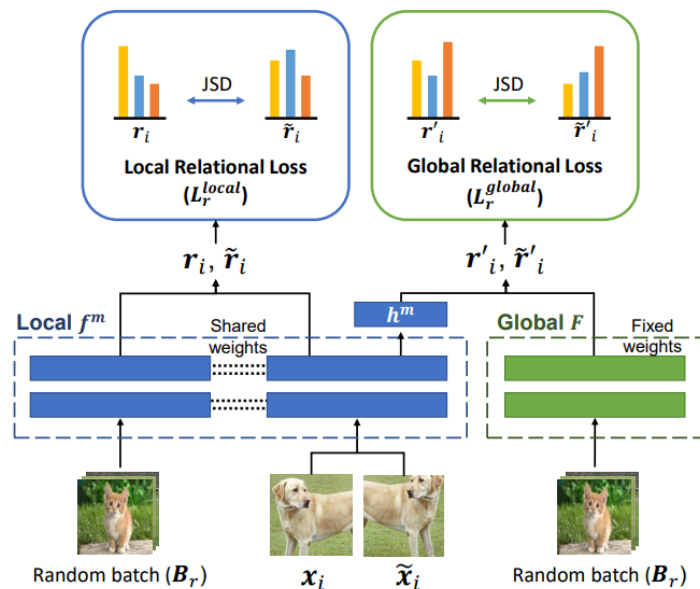


# Model overview

## Two-sided contrastive loss and relational loss



(a) Two-sided contrastive loss



(b) Two-sided relational loss

## Performance Evaluation

Table 1: Performance improvement with FedX on classification accuracy over three datasets. Both the final round accuracy and the best accuracy show that our model brings substantial improvement for all baseline algorithms.

Method	CIFAR-10		SVHN		F-MNIST	
	Last	Best	Last	Best	Last	Best
FedSimCLR	51.31	52.88	75.19	76.50	77.66	79.44
+ FedX	<b>56.88</b>	<b>57.95</b>	<b>77.19</b>	<b>77.70</b>	<b>81.98</b>	<b>82.47</b>
FedMoCo	56.74	57.82	70.69	70.99	82.31	83.58
+ FedX	<b>58.23</b>	<b>59.43</b>	<b>73.57</b>	<b>73.92</b>	<b>83.62</b>	<b>84.65</b>
FedBYOL	52.24	53.14	65.95	67.32	81.45	82.37
+ FedX	<b>56.49</b>	<b>57.79</b>	<b>68.94</b>	<b>69.05</b>	<b>83.18</b>	<b>84.30</b>
FedProtoCL	51.33	52.12	49.85	50.19	81.76	<b>83.57</b>
+ FedX	<b>55.36</b>	<b>56.76</b>	<b>69.31</b>	<b>69.75</b>	<b>82.74</b>	83.34
FedU	50.79	50.79	66.02	66.22	80.59	82.03
+ FedX	<b>56.15</b>	<b>57.26</b>	<b>68.13</b>	<b>68.39</b>	<b>83.73</b>	<b>84.12</b>



# I Conclusion

- We propose FedX, a new advance in unsupervised federated learning that learns data representations via a **unique two-sided knowledge distillation** at local and global levels.
- FedX can be applied to extant algorithms to **enhance performance by 1.58–5.52pp** in top-1 accuracy and further **enhance training speed**.
- **FedX preserves privacy** between clients and does not share data directly. It is also **does not require complex communication** for sending data features.



# Takeaways

# I Conclusion

- This presentation introduces four research branches for trustworthy AI:
  - (1) Fairness
  - (2) Interpretability
  - (3) Data privacy
  - (4) Security
- There still has a large room of improvements in jointly achieving multiple human-centered properties for trustworthy AI:
  - (1) Fairness vs. Performance tradeoff
  - (2) Interpretability vs. Performance tradeoff
  - (3) Transparency vs. Security tradeoff
  - (4) Privacy vs. Security tradeoff

....

# Thank you

**Sungwon Han** (lion4152@gmail.com)