

---

# Adversarially Learned Representations for Information Obfuscation and Inference

---

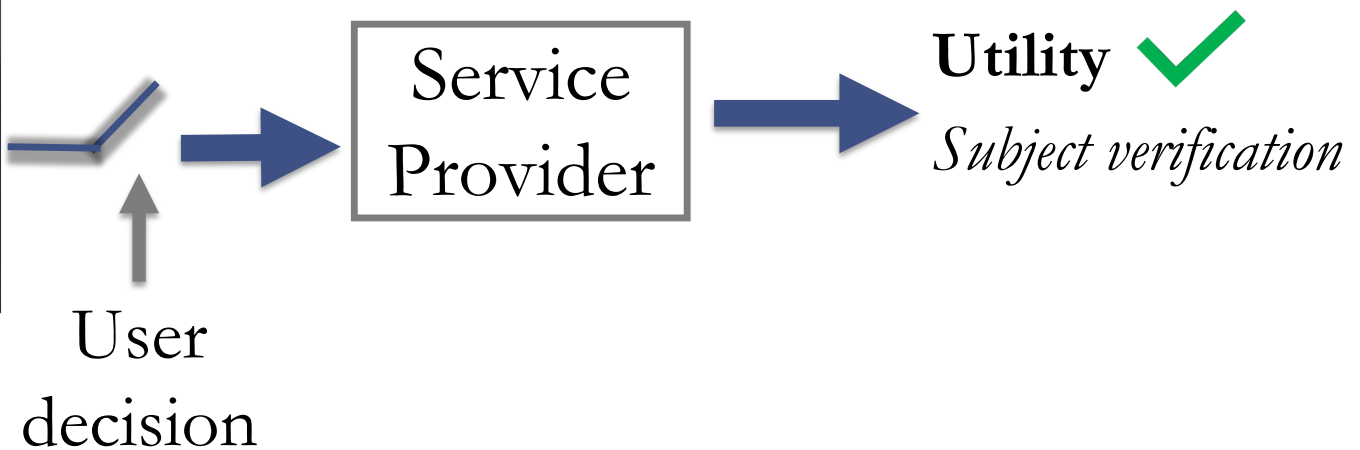
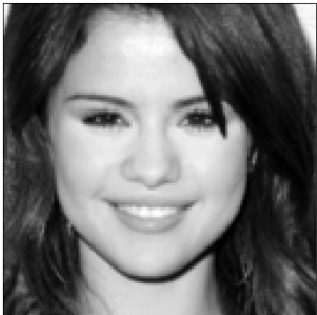
Martin Bertran<sup>1</sup>, Natalia Martinez<sup>1</sup>, Afroditi Papadaki<sup>2</sup>  
Qiang Qiu<sup>1</sup>, Miguel Rodrigues<sup>2</sup>, Galen Reeves<sup>1</sup>, Guillermo Sapiro<sup>1</sup>

1. Duke University
2. University College London

# Motivation

Why do users share their data?

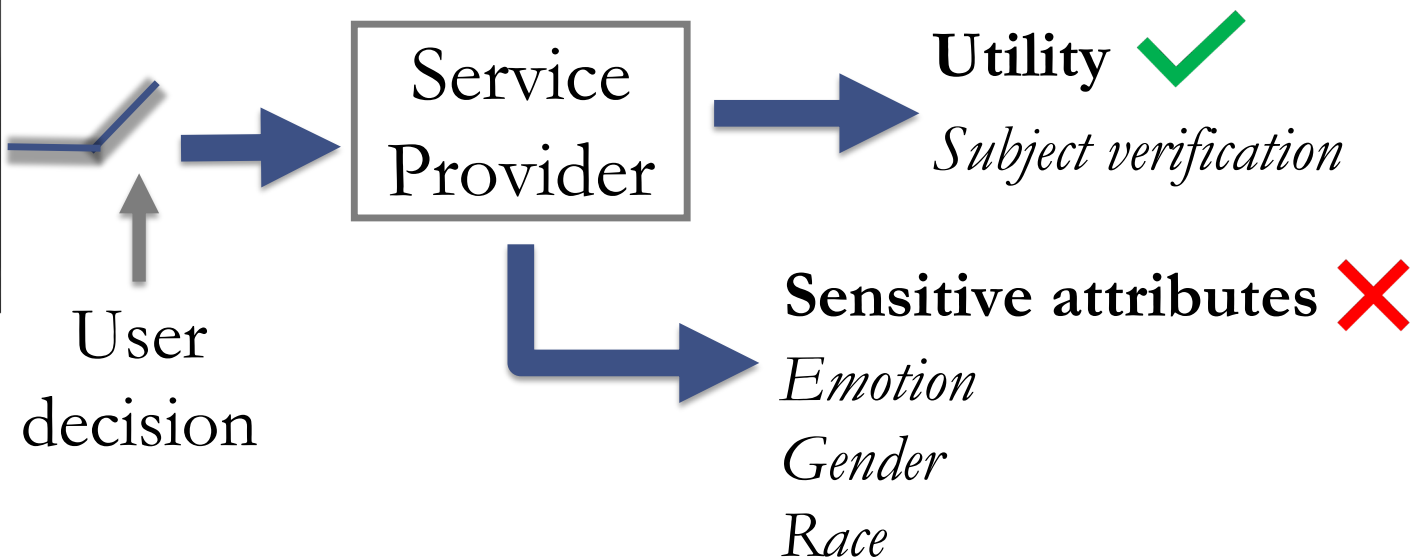
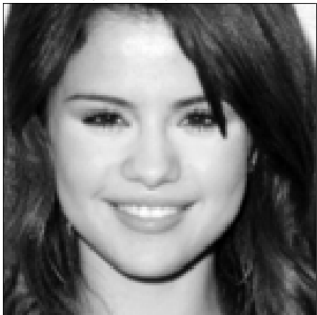
Shared data  
Facial Image



# Motivation

Why do users share their data?

Shared data  
Facial Image



---

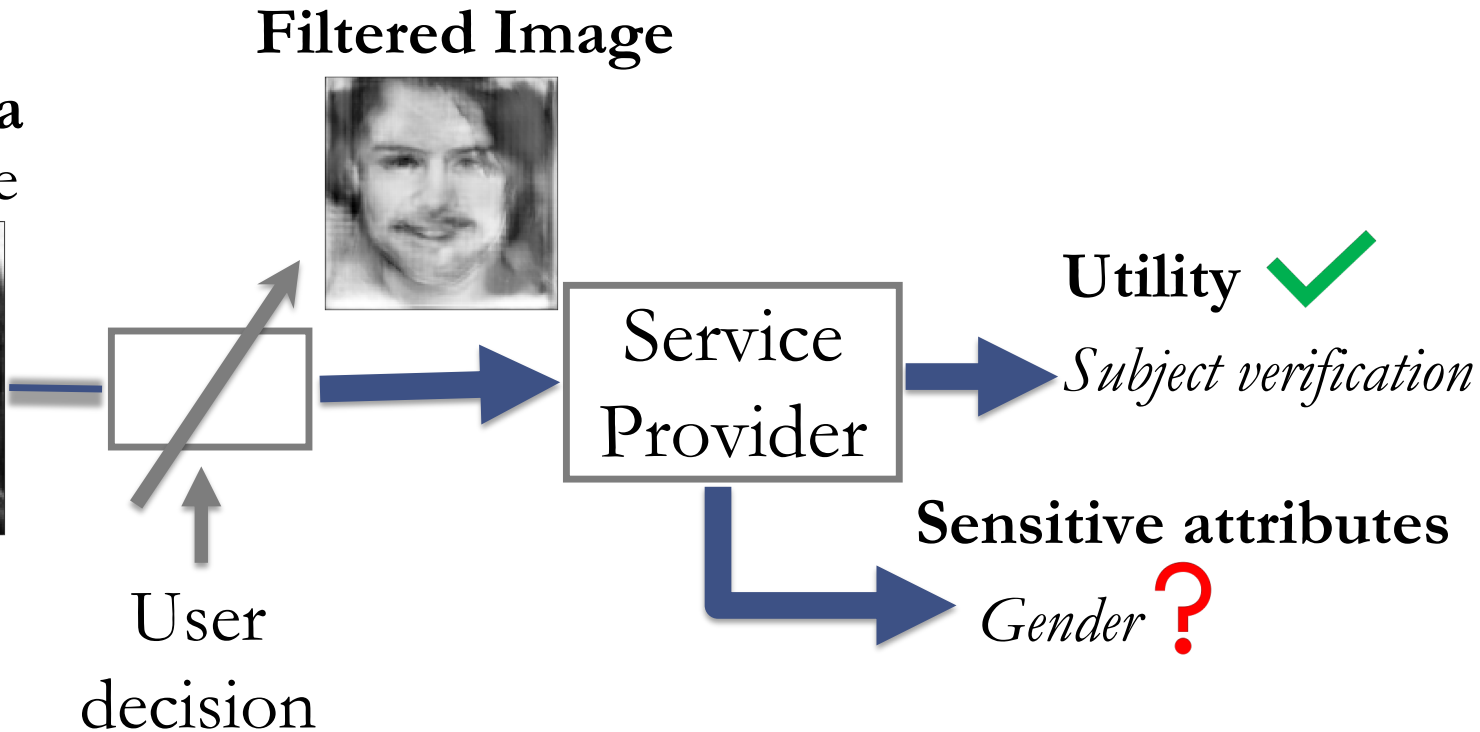
# Motivation

**Can we do better?**

# Motivation

Can we do better?

Shared data  
Facial Image



Learn **space-preserving** representations that **obfuscate** sensitive information while **preserving** utility.

# Motivation

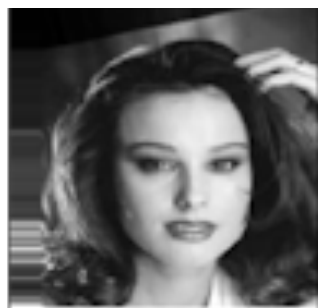
Example: Preserve gender & obfuscate emotion

Original

$P(\text{Male}) = 0.98$   
 $P(\text{Smile}) = 0.78$

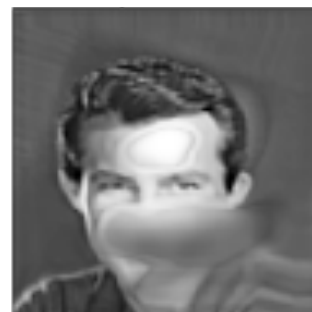


$P(\text{Female}) = 0.99$   
 $P(\text{Serious}) = 0.98$

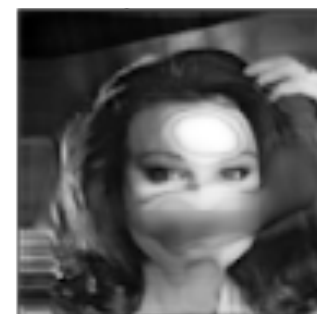


Filtered

$P(\text{Male}) = 0.98$   $\equiv$   
 $P(\text{Smile}) = 0.38$   $\downarrow$



$P(\text{Female}) = 0.99$   $\equiv$   
 $P(\text{Serious}) = 0.31$   $\downarrow$

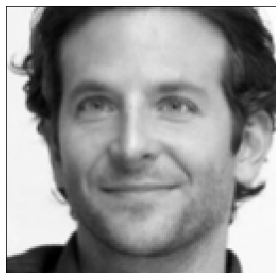


# Motivation

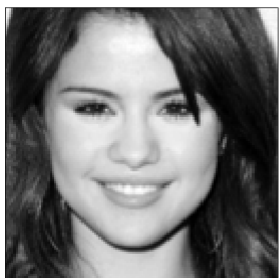
Example: Preserve subject & obfuscate gender

Original

$P(\text{Male}) = 0.99$   
Subject verified ✓



$P(\text{Female}) = 0.99$   
Subject verified ✓



Filtered

$P(\text{Male}) = 0.70$  ↓  
Subject verified ✓



$P(\text{Female}) = 0.54$  ↓  
Subject verified ✓



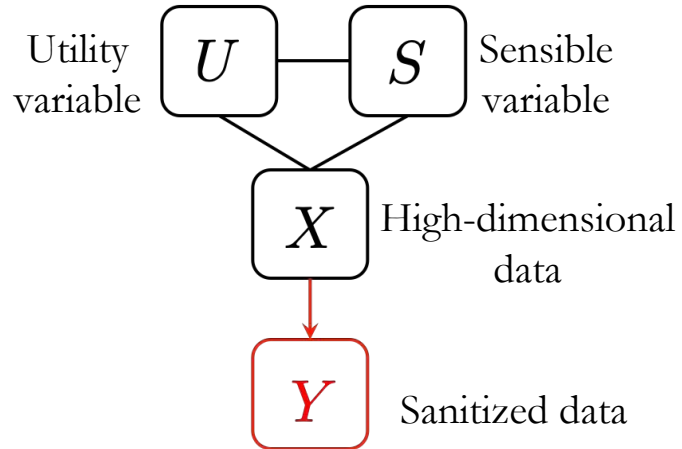
---

# Sample of related work

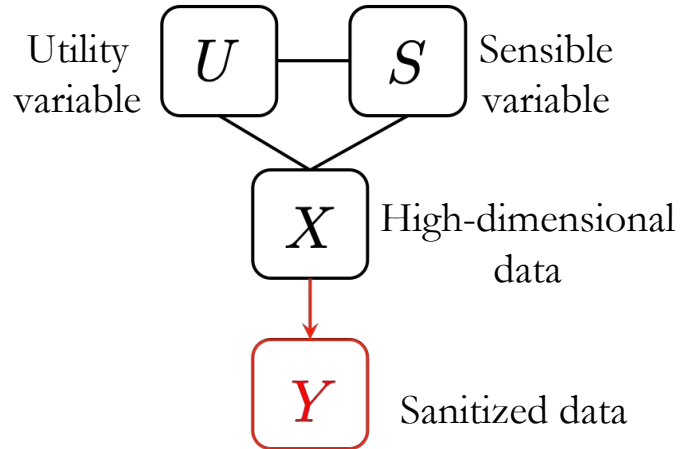
- (2003) Chechik et al. Extracting relevant structures with side information.
- (2016) Basciftci et al. On privacy-utility tradeoffs for constrained data release mechanisms.
- (2018) Madras et al. Learning adversarially fair and transferable representations.
- (2018) Sun et al. A hybrid model for identity obfuscation by face replacement.



# Problem formulation

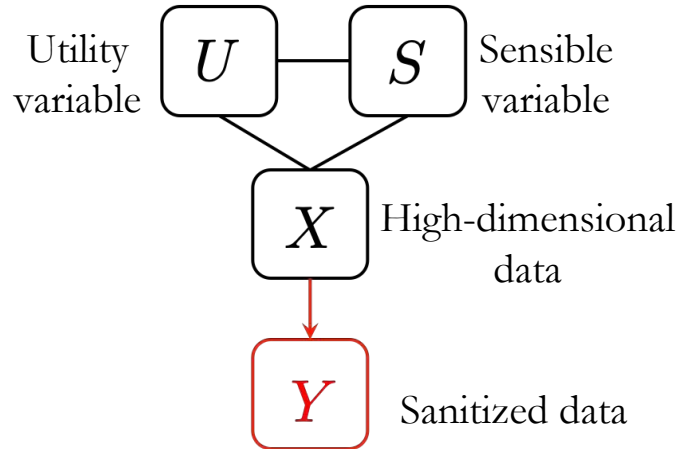


# Problem formulation



$$(U, S) \sim p(U, S)$$

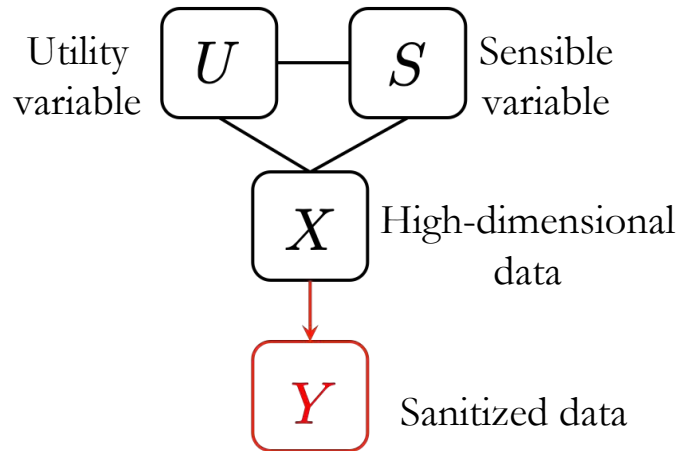
# Problem formulation



$$(U, S) \sim p(U, S)$$

$$X \sim p(X|U, S)$$

# Problem formulation



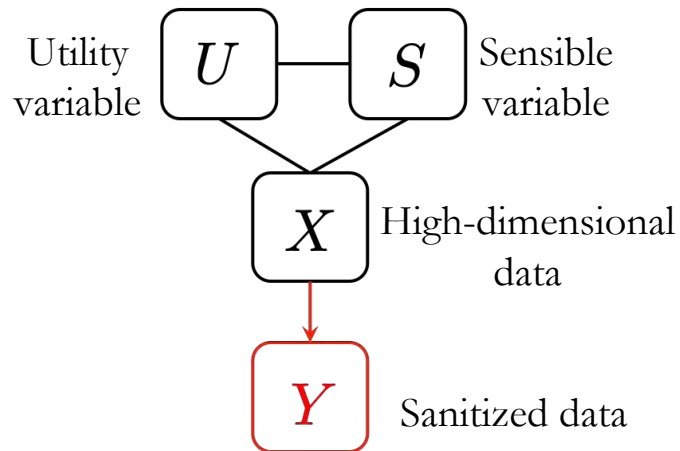
$$(U, S) \sim p(U, S)$$

$$X \sim p(X|U, S)$$

$$Y \sim p(Y|X)$$

***Our objective!***

# Problem formulation



$$(U, S) \sim p(U, S)$$

$$X \sim p(X|U, S)$$

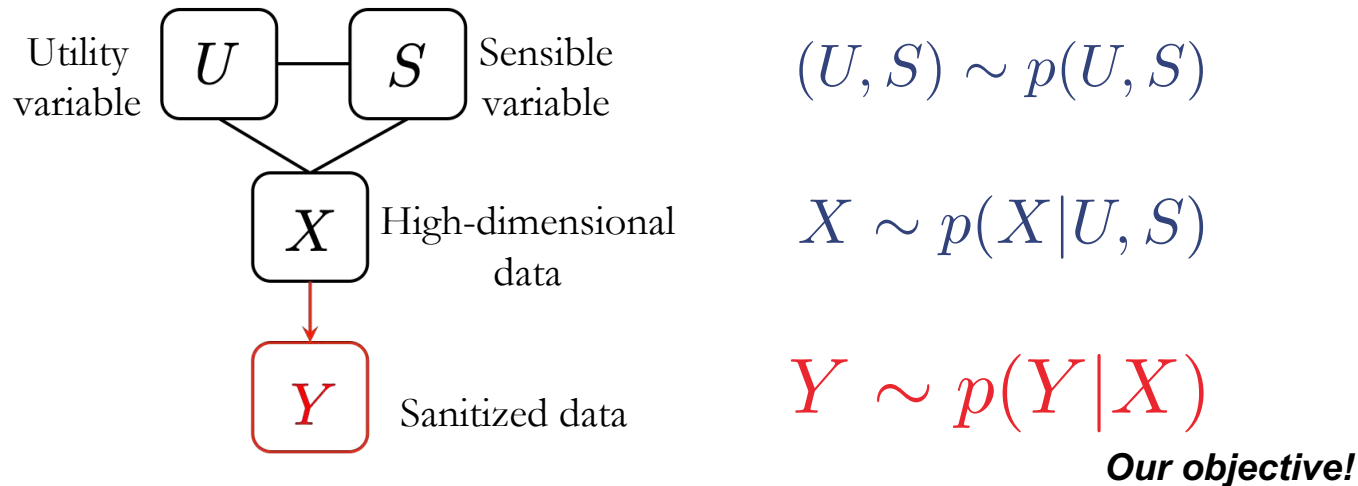
$$Y \sim p(Y|X)$$

***Our objective!***

Want to learn  $Y \sim p(Y|X)$  such that :

- $p(S|Y) \sim p(S)$
- $p(U|Y) \sim p(U|X)$

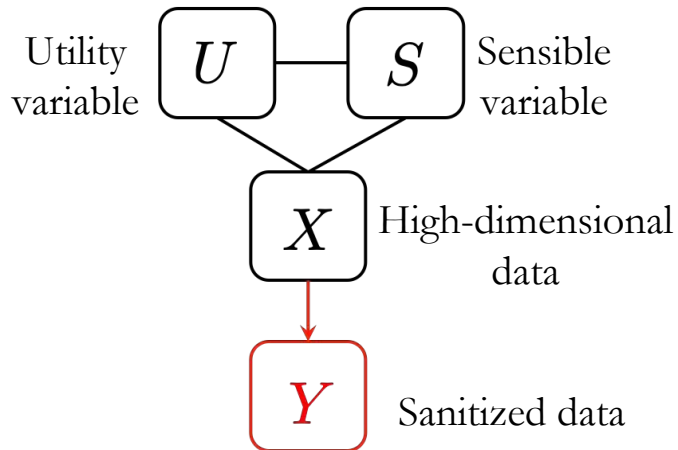
# Problem formulation



Want to learn  $Y \sim p(Y|X)$  such that :

- $p(S|Y) \sim p(S) \longrightarrow \min D_{KL}[p(S|Y)||p(S)]$
- $p(U|Y) \sim p(U|X)$

# Problem formulation



$$(U, S) \sim p(U, S)$$

$$X \sim p(X|U, S)$$

$$Y \sim p(Y|X)$$

***Our objective!***

Want to learn  $Y \sim p(Y|X)$  such that :

- $p(S|Y) \sim p(S) \longrightarrow \min D_{KL}[p(S|Y)||p(S)]$
- $p(U|Y) \sim p(U|X) \longrightarrow \min D_{KL}[p(U|X)||p(U|Y)]$

---

# Problem formulation

Want to learn  $Y \sim p(Y|X)$  such that:

- $\min D_{KL}[p(S|Y)||p(S)]$
- $\min D_{KL}[p(U|X)||p(U|Y)]$



# Problem formulation

Want to learn  $Y \sim p(Y|X)$  such that:

- $\min D_{KL}[p(S|Y)||p(S)] \xrightarrow{E_Y[\cdot]} I(S; Y)$
- $\min D_{KL}[p(U|X)||p(U|Y)]$

# Problem formulation

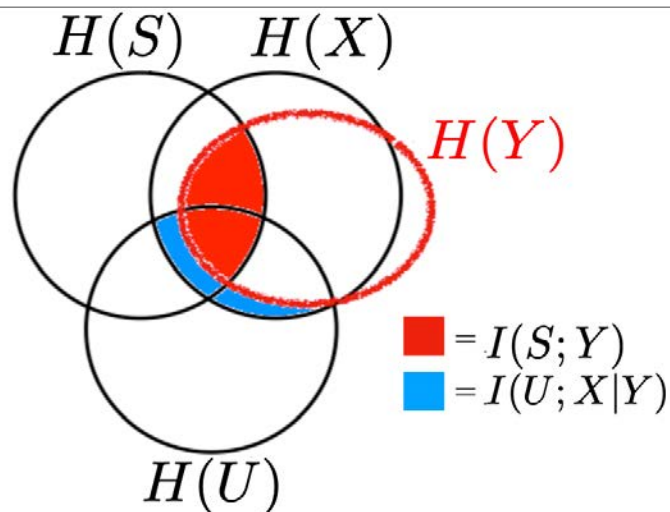
Want to learn  $Y \sim p(Y|X)$  such that:

- $\min D_{KL}[p(S|Y)||p(S)] \xrightarrow{E_Y[\cdot]} I(S; Y)$
- $\min D_{KL}[p(U|X)||p(U|Y)] \xrightarrow{E_{X,Y}[\cdot]} I(U; X|Y)$

# Problem formulation

Want to learn  $Y \sim p(Y|X)$  such that:

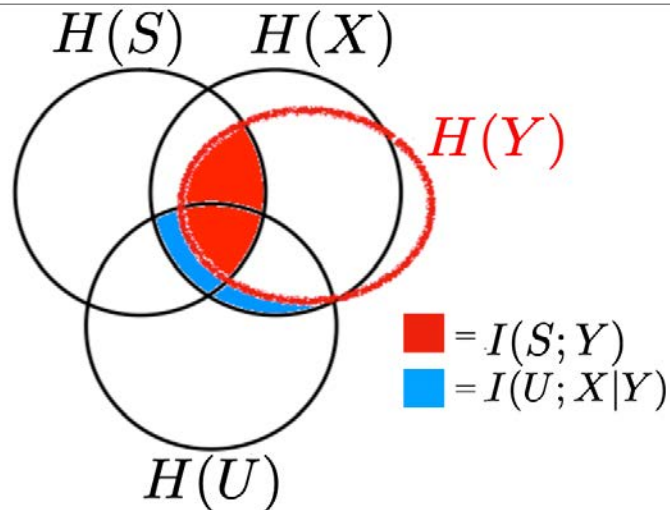
- $\min D_{KL}[p(S|Y)||p(S)] \xrightarrow{E_Y[\cdot]} I(S; Y)$
- $\min D_{KL}[p(U|X)||p(U|Y)] \xrightarrow{E_{X,Y}[\cdot]} I(U; X|Y)$



# Problem formulation

Want to learn  $Y \sim p(Y|X)$  such that:

- $\min D_{KL}[p(S|Y)||p(S)] \xrightarrow{E_Y[\cdot]} I(S; Y)$
- $\min D_{KL}[p(U|X)||p(U|Y)] \xrightarrow{E_{X,Y}[\cdot]} I(U; X|Y)$



**Objective:**

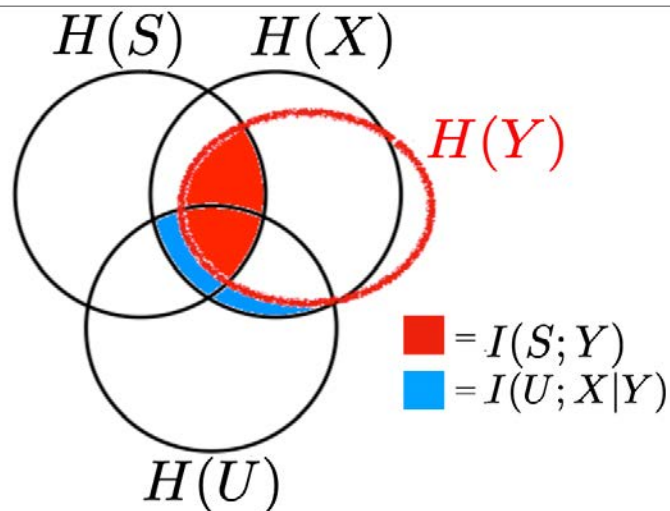
$$\min_{p(Y|X)} I(U; X|Y)$$

$$s.t. \quad I(S; Y) \leq k$$

# Problem formulation

Want to learn  $Y \sim p(Y|X)$  such that:

- $\min D_{KL}[p(S|Y)||p(S)] \xrightarrow{E_Y[\cdot]} I(S; Y)$
- $\min D_{KL}[p(U|X)||p(U|Y)] \xrightarrow{E_{X,Y}[\cdot]} I(U; X|Y)$



**Objective:**

$$\min_{p(Y|X)} I(U; X|Y) \quad \sim \quad \max_{p(Y|X)} I(U; Y)$$

$$s.t. \quad I(S; Y) \leq k$$

---

# Performance bounds

Given the objective  $\min_{p(Y|X)} I(U; X|Y)$  s.t.  $I(S; Y) \leq k$

---

# Performance bounds

Given the objective  $\min_{p(Y|X)} I(U; X|Y)$  s.t.  $I(S; Y) \leq k$

**What are the intrinsic limits on the trade-offs for this problem?**

# Performance bounds

Given the objective  $\min_{p(Y|X)} I(U; X|Y)$  s.t.  $I(S; Y) \leq k$

**What are the intrinsic limits on the trade-offs for this problem?**

**Lemma 1.**

$(U, S) \in \mathcal{U} \times \mathcal{S}$  finite alphabets,  $X \sim p(X|U, S)$ .

Then:

$$\begin{aligned} \min_{p(Y|X)} I(U; X|Y) &\geq \min_{p(Y|U, S)} I(U; X) - I(U; Y) \\ \text{s.t. } I(S; Y) &\leq k & \text{s.t. } I(S; Y) &\leq k \\ & & I(U; Y) &\leq I(U; X) \end{aligned}$$

- With  $|\mathcal{Y}|$  finite we can compute a sequence of upper bounds: Restricted cardinality sequence (RCS).



---

# Performance bounds

Given the objective  $\min_{p(Y|X)} I(U; X|Y)$  s.t.  $I(S; Y) \leq k$

**What are the intrinsic limits on the trade-offs for this problem?**

**Lemma 2.** Given  $(X, U, S) \sim p(X, U, S)$

$$I(U; X|Y) \geq -I(S; Y) + I(U; S) - I(U; S|X)$$

# Performance bounds

Given the objective  $\min_{p(Y|X)} I(U; X|Y)$  s.t.  $I(S; Y) \leq k$

What are the intrinsic limits on the trade-offs for this problem?

**Lemma 2.** Given  $(X, U, S) \sim p(X, U, S)$

$$I(U; X|Y) \geq -I(S; Y) + I(U; S) - I(U; S|X)$$

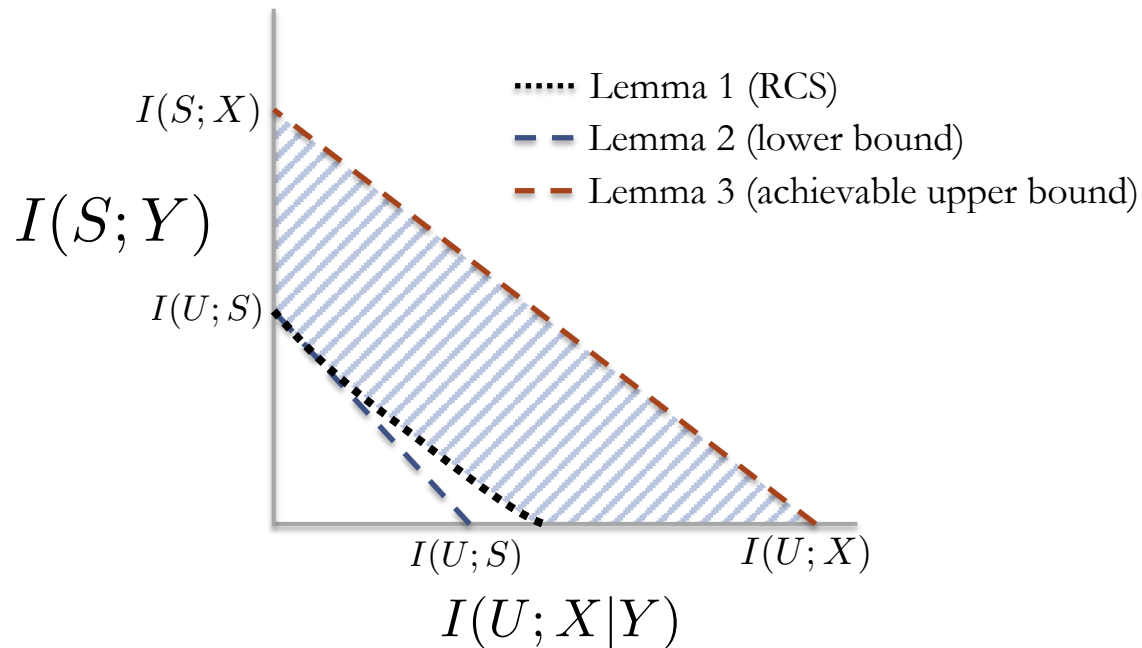
**Lemma 3.** Given  $(X, U, S) \sim p(X, U, S)$ ,  $\forall k \geq 0 \exists p(Y|X)$  such that:

$$I(S; Y) \leq k$$

$$I(U; X|Y) = \max\left(0, 1 - \frac{k}{I(S; X)}\right) I(U; X)$$

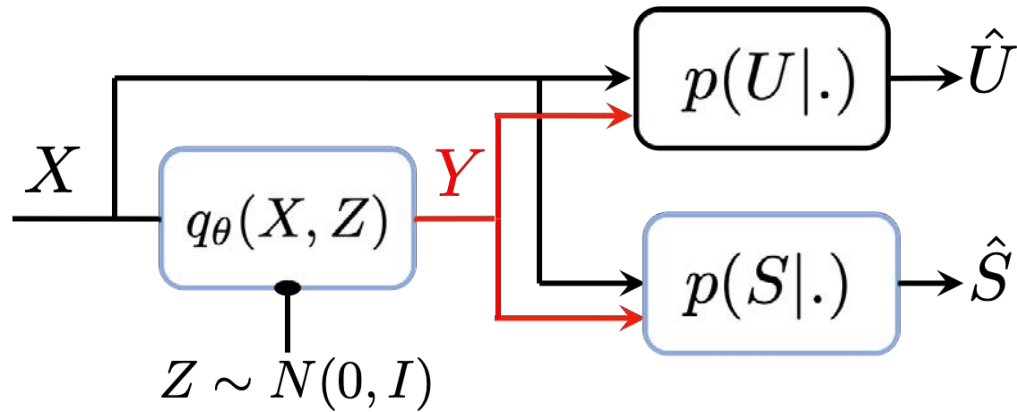
# Performance bounds

Lemmas 1, 2 and 3 can be approximated using contingency tables.

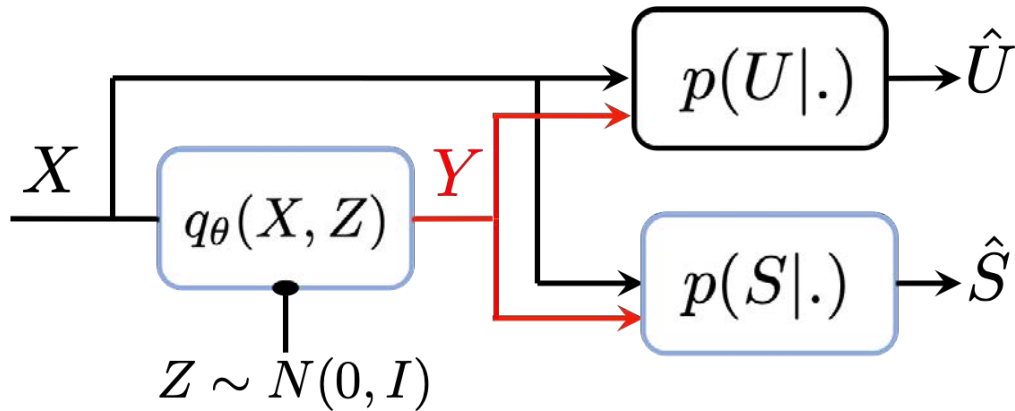


\* Sketch under the assumption that  $I(U; S|X) = 0$

# Proposed framework



# Proposed framework



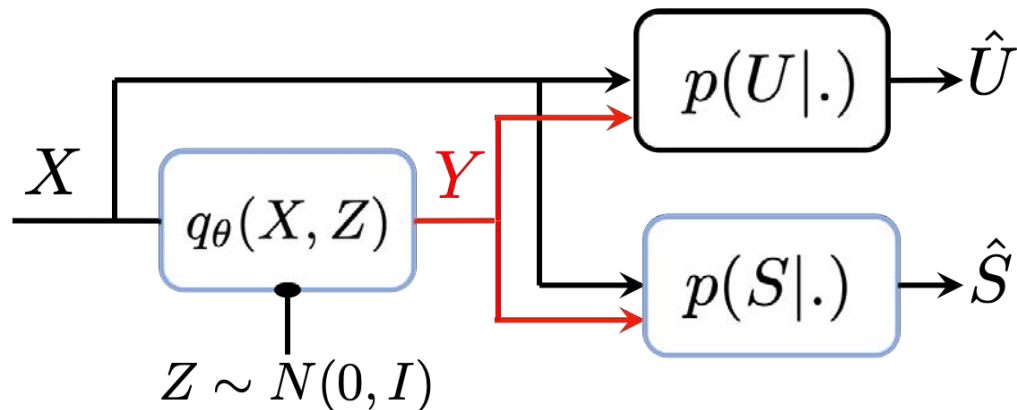
**Objective:**

$$\min I(U; X|Y)$$

$$p(Y|X) \sim q_\theta(X, Z)$$

$$s.t. : I(S; Y) \leq k$$

# Proposed framework



**Objective:**

$$\min I(U; X|Y)$$

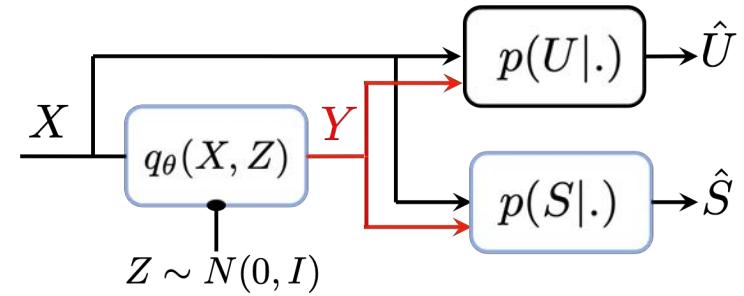
$$p(Y|X) \sim q_\theta(X, Z)$$

$$s.t. : I(S; Y) \leq k$$

Optimization objective:

$$\min [I(U; X|Y) + \lambda \max\{I(S; Y) - k, 0\}^2]$$
$$p(Y|X) \sim q_\theta(X, Z)$$

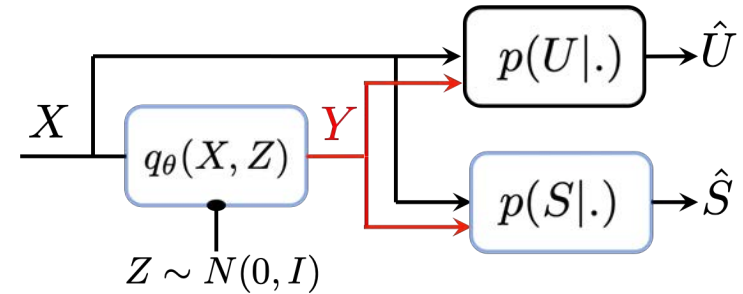
# Implementation



Optimization objective:

$$\min_{q_\theta(X, Z)} [I(U; X|Y) + \lambda \max\{I(S; Y) - k, 0\}^2]$$

# Implementation



Optimization objective:

$$\min_{q_\theta(X, Z)} [I(U; X|Y) + \lambda \max\{I(S; Y) - k, 0\}^2]$$

Learning the stochastic mapping  $Y = q_\theta(X, Z)$  :

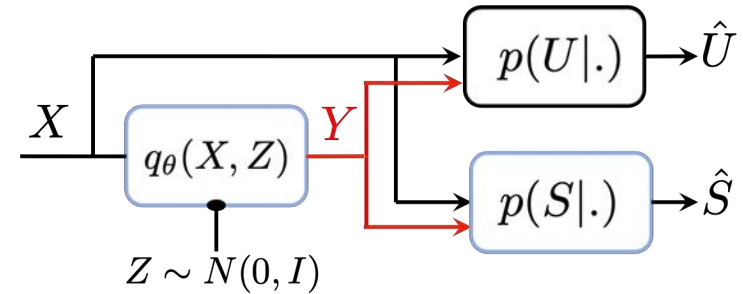
$$p(U|X) \sim p_\phi(U|X) \longrightarrow \hat{\phi} = \operatorname{argmin}_\phi E_{X,U} [-\log(p_\phi(U|X))]$$

$$p(U|Y) \sim p_\psi(U|Y) \longrightarrow \hat{\psi} = \operatorname{argmin}_\psi E_{X,U,Z} [-\log(p_\psi(U|q_\theta(X, Z)))]$$

$$p(S|Y) \sim p_\eta(S|Y) \longrightarrow \hat{\eta} = \operatorname{argmin}_\eta E_{X,S,Z} [-\log(p_\eta(S|q_\theta(X, Z)))]$$



# Implementation



Optimization objective:

$$\min_{q_\theta(X, Z)} [I(U; X|Y) + \lambda \max\{I(S; Y) - k, 0\}^2]$$

Learning the stochastic mapping  $Y = q_\theta(X, Z)$  :

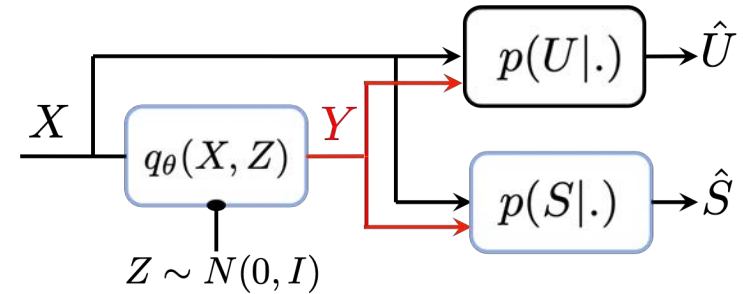
$$p(U|X) \sim p_\phi(U|X) \longrightarrow \hat{\phi} = \operatorname{argmin}_\phi E_{X,U} [-\log(p_\phi(U|X))]$$

$$p(U|Y) \sim p_\psi(U|Y) \longrightarrow \hat{\psi} = \operatorname{argmin}_\psi E_{X,U,Z} [-\log(p_\psi(U|q_{\hat{\theta}}(X, Z)))]$$

$$p(S|Y) \sim p_\eta(S|Y) \longrightarrow \hat{\eta} = \operatorname{argmin}_\eta E_{X,S,Z} [-\log(p_\eta(S|q_{\hat{\theta}}(X, Z)))]$$

$$\hat{\theta} = \operatorname{argmin}_\theta E_{X,Z} [D_{KL}[p_{\hat{\phi}}(U|X) || p_{\hat{\psi}}(U|q_\theta(X, Z))]] \\ + \lambda \max(E_{X,Z} [D_{KL}[p_{\hat{\eta}}(S|q_\theta(X, Z)) || P(S)] - k, 0)^2$$

# Implementation



Optimization objective:

$$\min_{q_\theta(X, Z)} [I(U; X|Y) + \lambda \max\{I(S; Y) - k, 0\}^2]$$

Learning the stochastic mapping  $Y = q_\theta(X, Z)$  :

$$p(U|X) \sim p_\phi(U|X) \longrightarrow \hat{\phi} = \operatorname{argmin}_\phi E_{X,U} [-\log(p_\phi(U|X))]$$

$$p(U|Y) \sim p_\psi(U|Y) \longrightarrow \hat{\psi} = \operatorname{argmin}_\psi E_{X,U,Z} [-\log(p_\psi(U|q_{\hat{\theta}}(X,Z)))]$$

$$p(S|Y) \sim p_\eta(S|Y) \longrightarrow \hat{\eta} = \operatorname{argmin}_\eta E_{X,S,Z} [-\log(p_\eta(S|q_{\hat{\theta}}(X,Z)))]$$

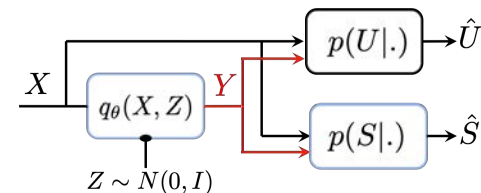
**Xception  
Networks**

$$\hat{\theta} = \operatorname{argmin}_\theta E_{X,Z} [D_{KL}[p_{\hat{\phi}}(U|X) || p_{\hat{\psi}}(U|q_\theta(X,Z))]]$$

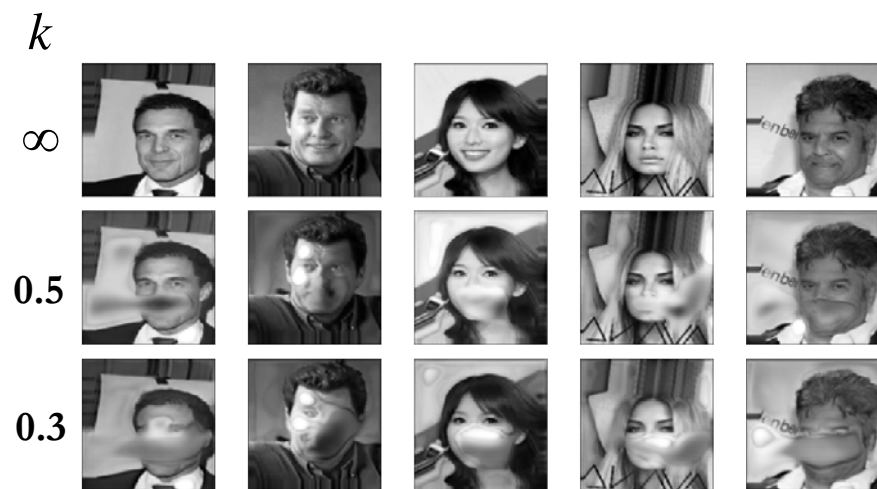
$$+ \lambda \max(E_{X,Z} [D_{KL}[p_{\hat{\eta}}(S|q_\theta(X,Z)) || P(S)] - k, 0)^2$$

**U-NET + noise**

# Experiments

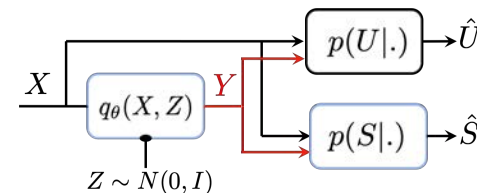


## Emotion obfuscation vs gender detection

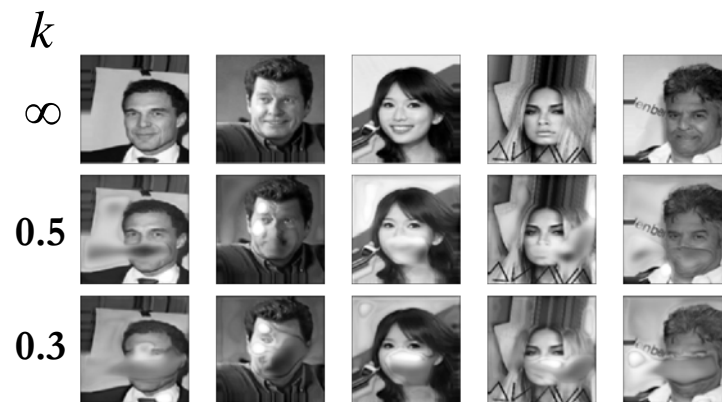
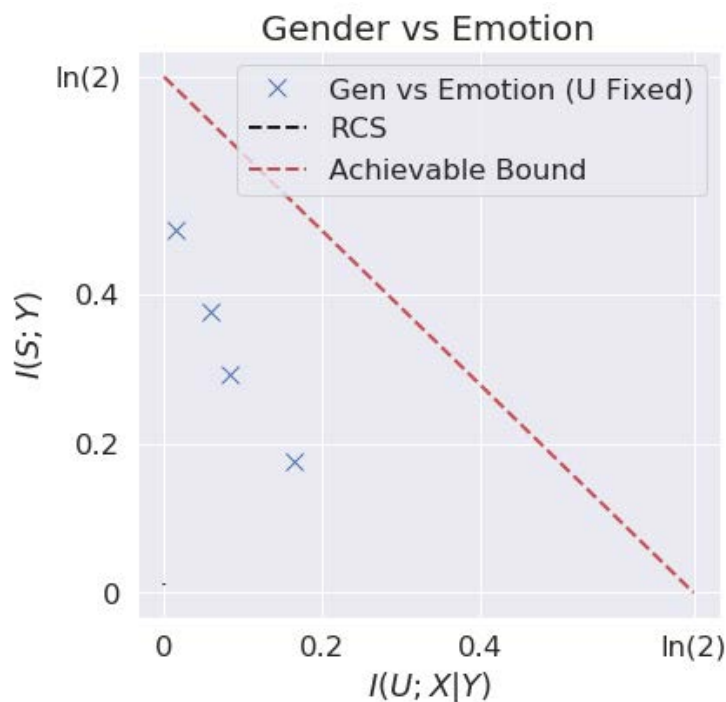


SENSITIVITY TOLERANCE	FIXED EMOTION ACCURACY	ADVERSARIAL EMOTION ACCURACY	FIXED GENDER ACCURACY
$k$			
$\infty$	<b>91.8%</b>	<b>91.8%</b>	<b>94.9%</b>
0.5	68.4%	91.4%	89.3%
0.4	58.6%	85.8%	88.0%
0.3	56.8%	81.5%	86.7%
<b>0.2</b>	<b>51.9%</b>	<b>74.3%</b>	<b>83.9%</b>
GUESSING	51.9%	51.9%	60.7%

# Experiments

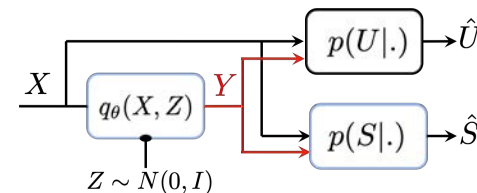


## Emotion obfuscation vs gender detection

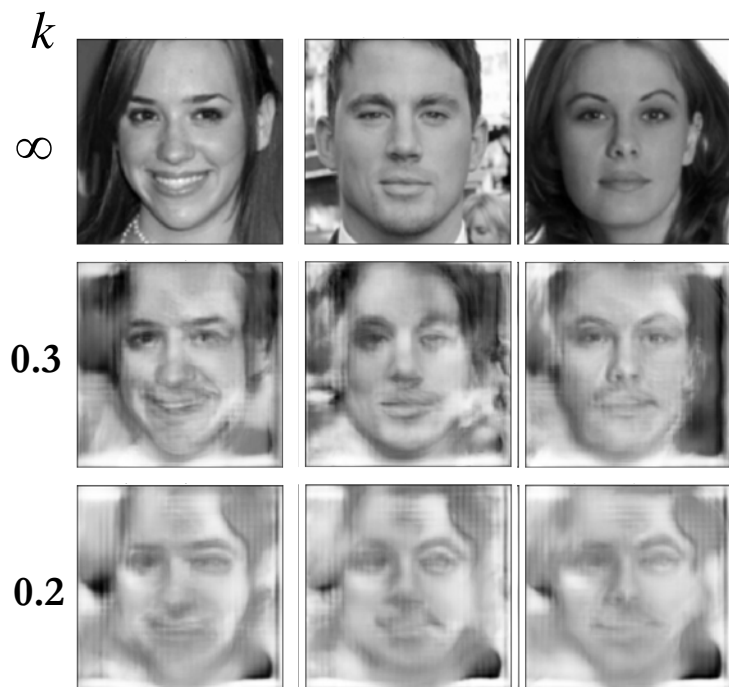


SENSITIVITY TOLERANCE	FIXED EMOTION	ADVERSARIAL EMOTION	FIXED GENDER
$k$	ACCURACY	ACCURACY	ACCURACY
$\infty$	<b>91.8%</b>	<b>91.8%</b>	<b>94.9%</b>
0.5	68.4%	91.4%	89.3%
0.4	58.6%	85.8%	88.0%
0.3	56.8%	81.5%	86.7%
<b>0.2</b>	<b>51.9%</b>	<b>74.3%</b>	<b>83.9%</b>
GUESSING	51.9%	51.9%	60.7%

# Experiments

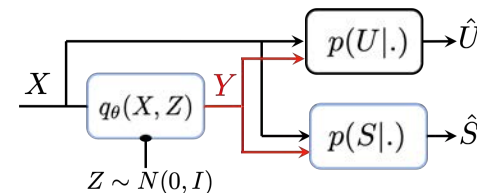


## Gender obfuscation vs subject verification

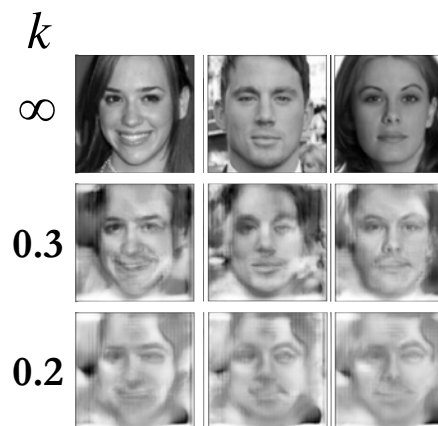
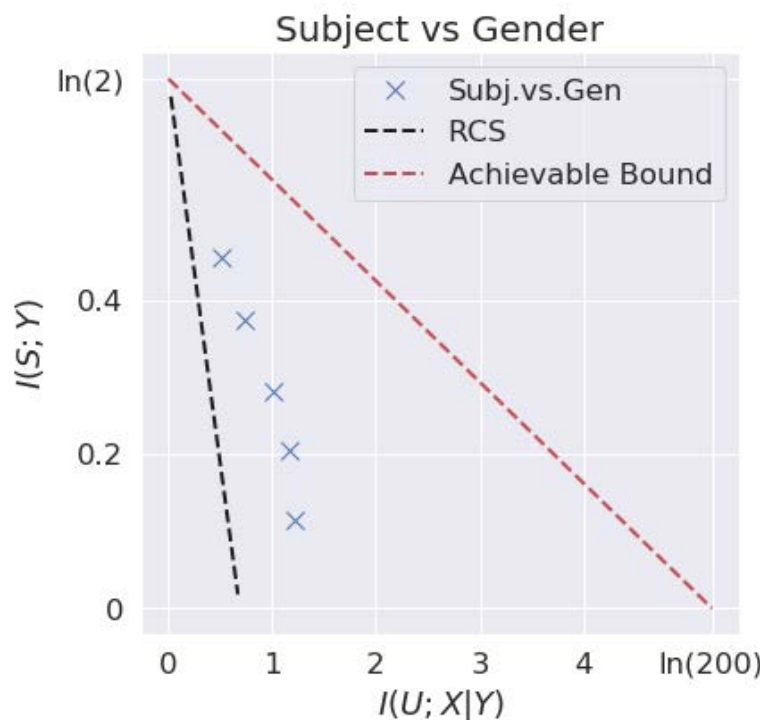


SENSITIVITY TOLERANCE	FIXED GENDER ACCURACY	ADVERSARIAL GENDER ACCURACY	FIXED SUBJECT TOP-5 ACCURACY	RETRAINED SUBJECT TOP-5 ACCURACY
$k$				
$\infty$	<b>98.6%</b>	<b>98.6%</b>	<b>98.8%</b>	<b>98.8%</b>
0.5	59.5%	90.2%	93.5%	96.8%
0.4	60.3%	85.3%	88.1%	94.9%
0.3	54.0%	79.4%	81.4%	92.8%
<b>0.2</b>	<b>56.1%</b>	<b>74.6%</b>	<b>81.6%</b>	<b>91.0%</b>
0.1	51.6%	67.1%	74.5%	89.6%
GUESSING	54.8%	54.8%	2.5%	2.5%

# Experiments

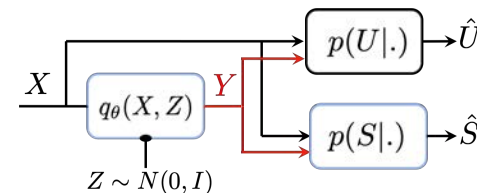


## Gender obfuscation vs subject verification



SENSITIVITY TOLERANCE	FIXED GENDER	ADVERSARIAL GENDER	FIXED SUBJECT	RETRAINED SUBJECT
$k$	ACCURACY	ACCURACY	TOP-5 ACCURACY	TOP-5 ACCURACY
$\infty$	<b>98.6%</b>	<b>98.6%</b>	<b>98.8%</b>	<b>98.8%</b>
0.5	59.5%	90.2%	93.5%	96.8%
0.4	60.3%	85.3%	88.1%	94.9%
0.3	54.0%	79.4%	81.4%	92.8%
<b>0.2</b>	<b>56.1%</b>	<b>74.6%</b>	<b>81.6%</b>	<b>91.0%</b>
0.1	51.6%	67.1%	74.5%	89.6%
GUESSING	54.8%	54.8%	2.5%	2.5%

# Experiments



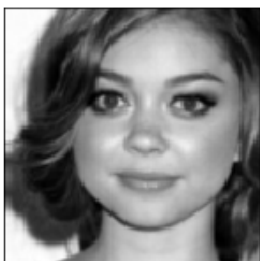
## Subject within Subject

Consenting User

Nonconsenting User

$k$  Subject verified ✓

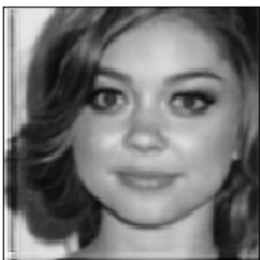
Subject verified ✓



$\infty$

Subject verified ✓

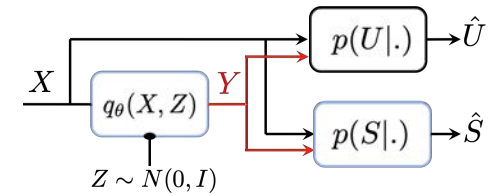
Subject verified ✗



0.5

SENSITIVITY TOLERANCE	CONSENTING USERS	NONCONSENTING USERS
	TOP-5 ACCURACY	TOP-5 ACCURACY
$k$		
$\infty$	<b>98.7%</b>	<b>97.9%</b>
3	98.3%	9.38%
1	97.8%	6.25%
<b>0.5</b>	<b>97.6%</b>	<b>4.69%</b>
GUESSING	2.5%	2.5%

# Concluding remarks



- Learned representations that *preserve utility* and *obfuscate sensitive information*.
- Transformations are *space-preserving*. Can reuse existing pipelines.
- Derived easy-to-compute bounds.
- Experimental results show representations compare favorably against derived bounds.

## Limitations:

- Expectation-based approach.
- Reliance on adversary as a proxy for information.



---

**Thanks!**

Please visit us at poster #81