

COMP6211: Trustworthy Machine Learning

Fairness

Minhao CHENG

Machine learning ethics

Ad related to latanya sweeney ⓘ

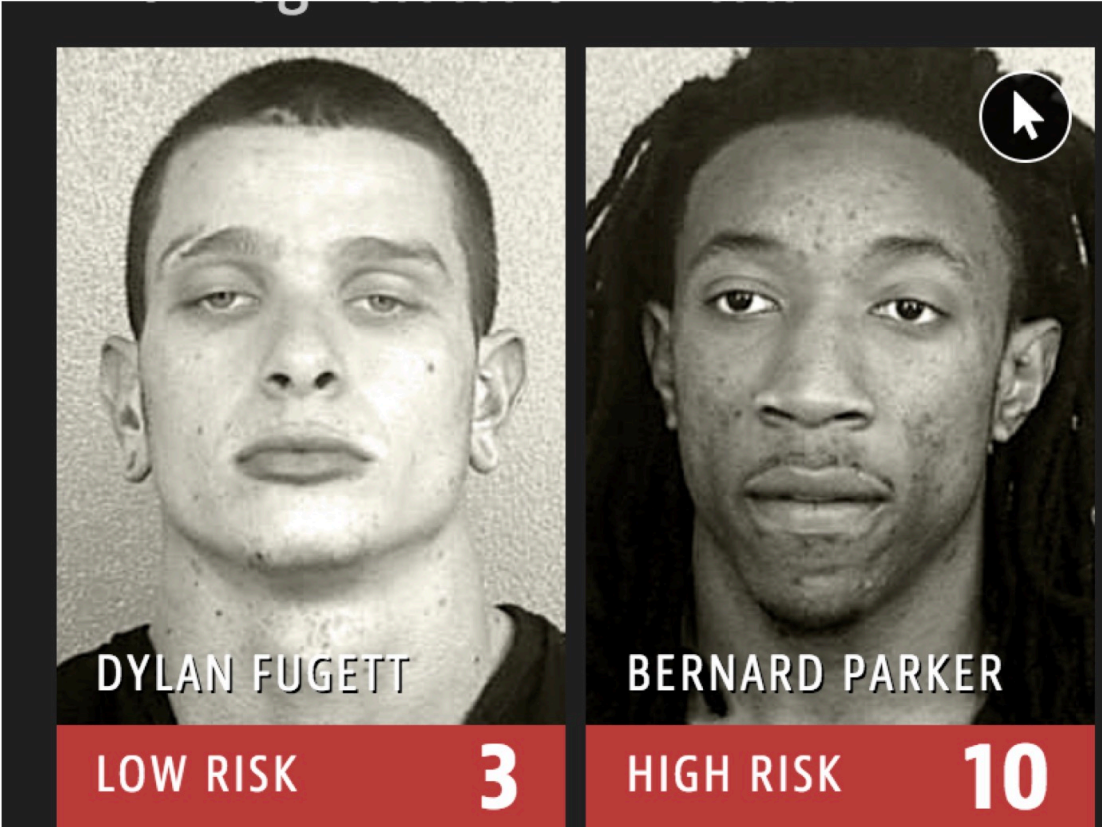
[Latanya Sweeney Truth](#)
www.instantcheckmate.com/
Looking for **Latanya Sweeney**? Check **Latanya Sweeney's** Arrests.

Ads by Google

[Latanya Sweeney, Arrested?](#)
1) Enter Name and State. 2) Access Full Background Checks Instantly.
www.instantcheckmate.com/

[Latanya Sweeney](#)
Public Records Found For: **Latanya Sweeney**. View Now.
www.publicrecords.com/

[La Tanya](#)
Search for La Tanya Look Up Fast Results now!
www.ask.com/La+Tanya



DYLAN FUGETT BERNARD PARKER

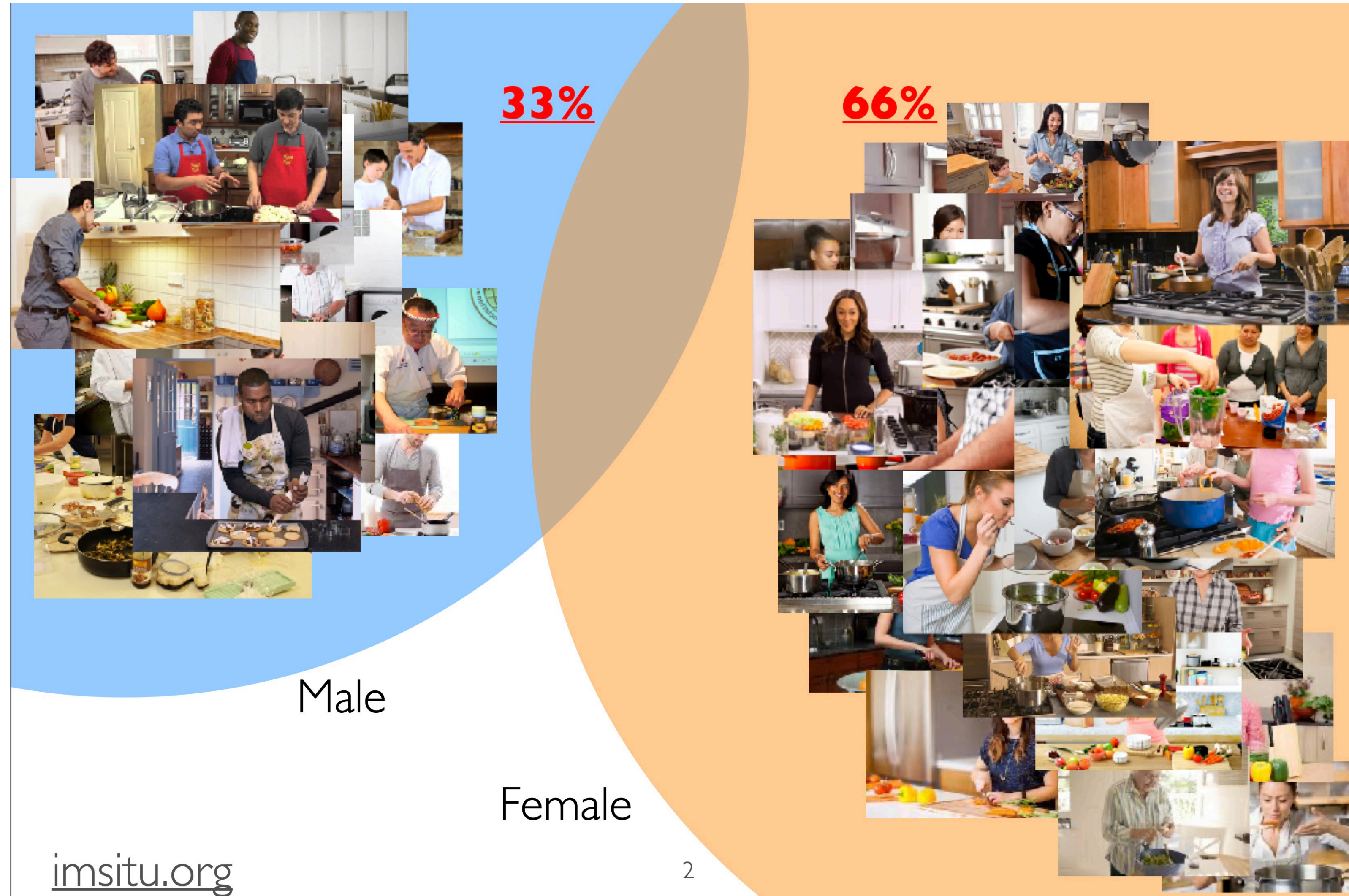
LOW RISK **3** HIGH RISK **10**

Fugett was rated low risk after being arrested with cocaine and marijuana. He was arrested three times on drug charges after that.

Ethical machine learning matters
in **high-stakes** domains



Group bias example: gender bias



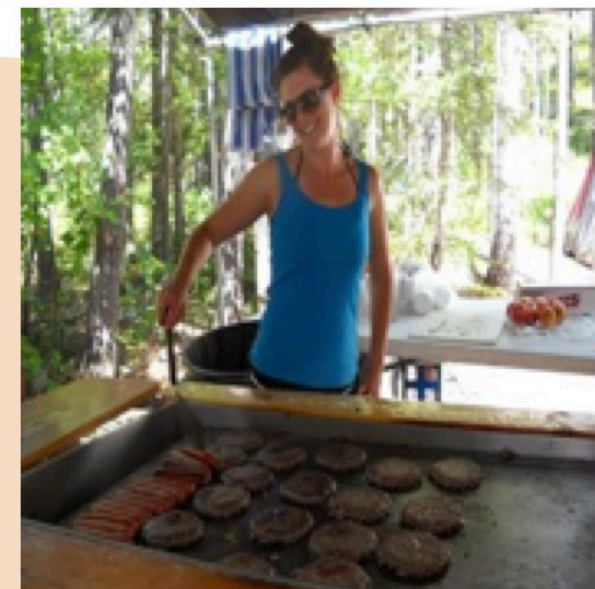
Group bias example: gender bias



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



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



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

Fairness in Machine Learning

- Group fairness
 - Don't discriminate unnecessarily between **protected** groups (race, gender, sexuality, religion, etc.)
- Individual fairness
 - Treat similar individuals similarly

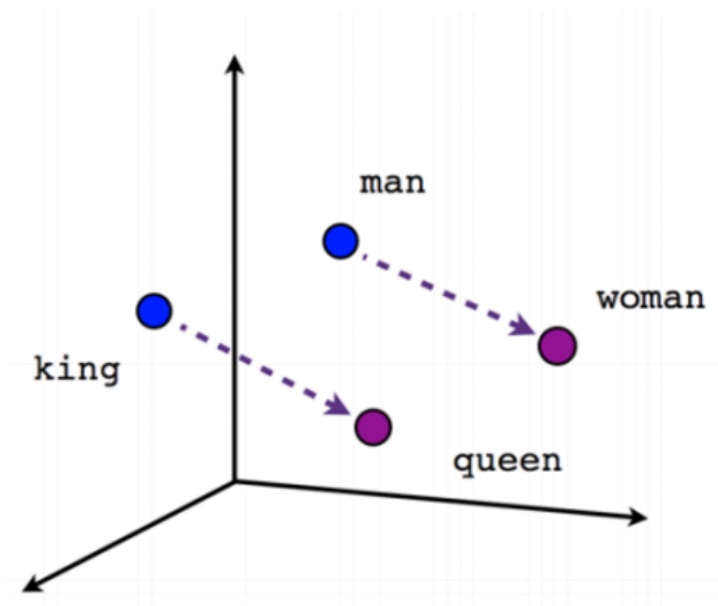
Bias

- Found in language data, learned by humans and ML
- Stereotyped bias: “problematic where such information is derived from aspects of human culture known to lead to harmful behavior”
- Prejudiced actions are taken based on stereotyped bias

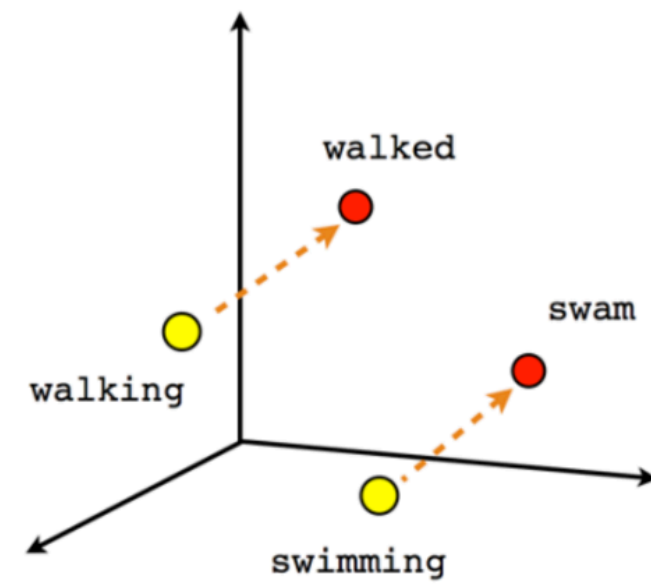
How to measure word embedding bias?

- Humans:
 - Implicit Association Test
 - Response time differs when humans pair concepts that they find similar compared to concepts that they find different
- Machines:
 - Word embeddings
 - Measure cosine distance between embedding vectors

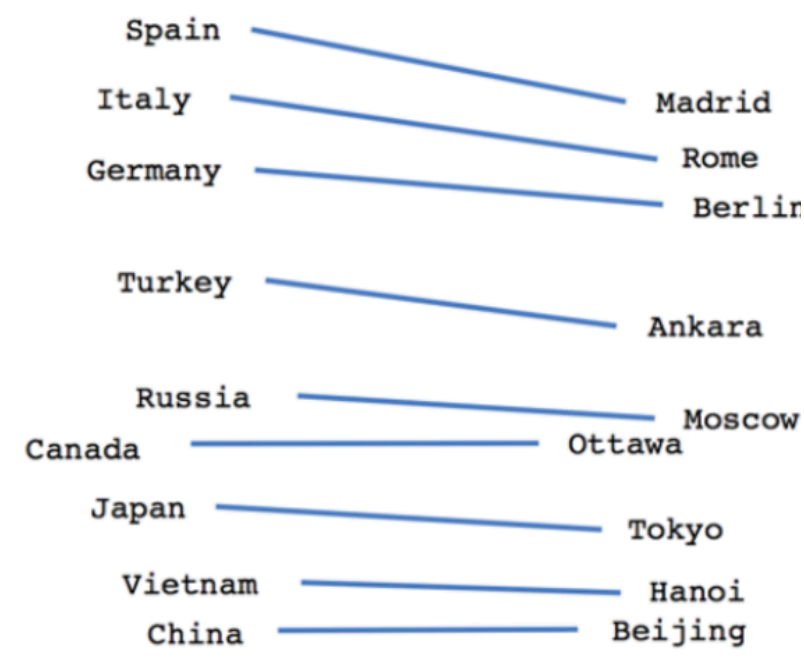
Word embeddings



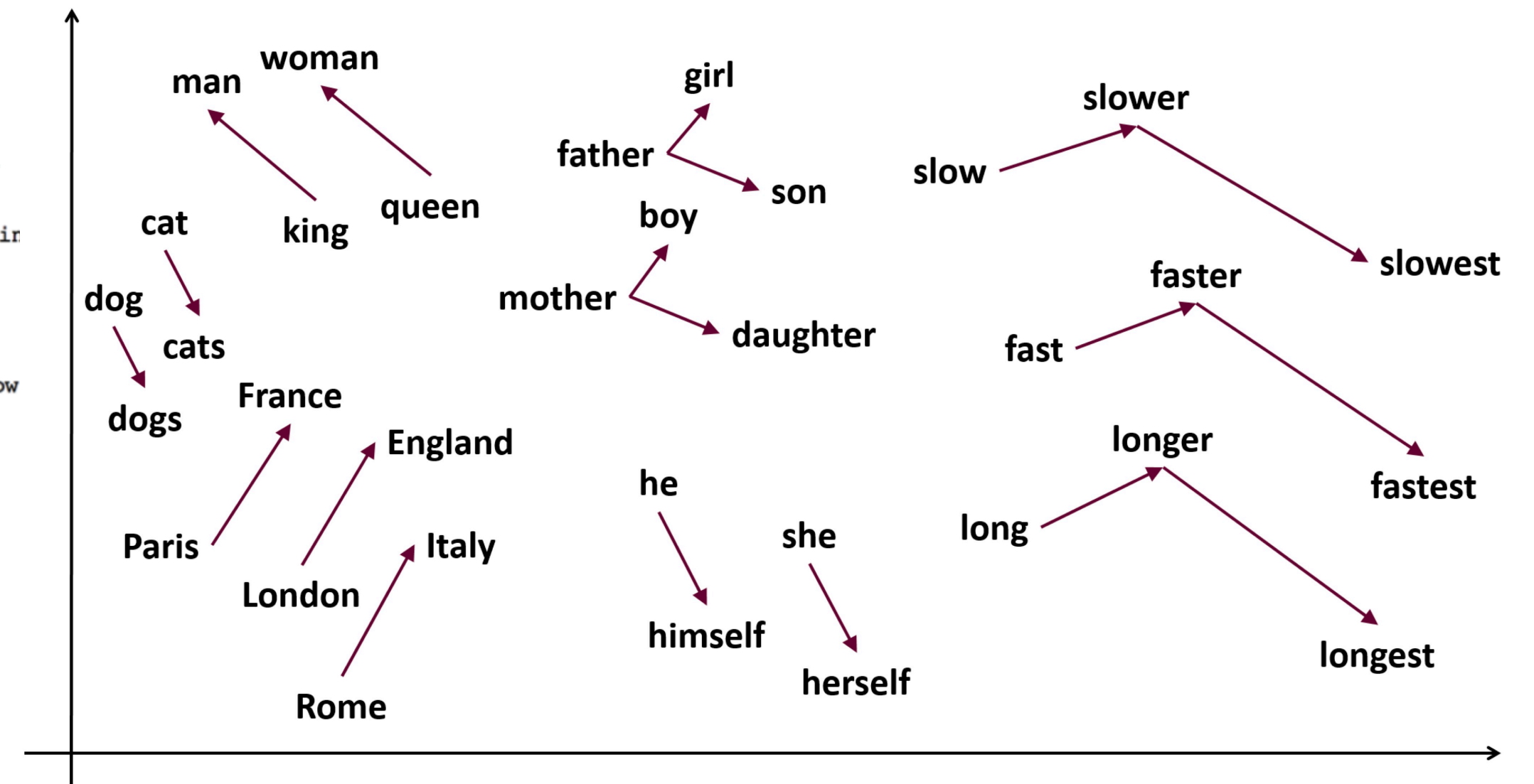
Male-Female



Verb tense



Country-Capital



N: population size

d: effect size

p : p-value

N_T: number of target words

N_A: number of attribute words

Target words	Attrib. words	Original Finding				Our Finding			
		Ref	N	d	p	N _T	N _A	d	p
Flowers vs insects	Pleasant vs unpleasant	(5)	32	1.35	10^{-8}	25 × 2	25 × 2	1.50	10^{-7}
Instruments vs weapons	Pleasant vs unpleasant	(5)	32	1.66	10^{-10}	25 × 2	25 × 2	1.53	10^{-7}
Eur.-American vs Afr.-American names	Pleasant vs unpleasant	(5)	26	1.17	10^{-5}	32 × 2	25 × 2	1.41	10^{-8}
Eur.-American vs Afr.-American names	Pleasant vs unpleasant from (5)	(7)	Not applicable			16 × 2	25 × 2	1.50	10^{-4}
Eur.-American vs Afr.-American names	Pleasant vs unpleasant from (9)	(7)	Not applicable			16 × 2	8 × 2	1.28	10^{-3}
Male vs female names	Career vs family	(9)	39k	0.72	$< 10^{-2}$	8 × 2	8 × 2	1.81	10^{-3}
Math vs arts	Male vs female terms	(9)	28k	0.82	$< 10^{-2}$	8 × 2	8 × 2	1.06	.018
Science vs arts	Male vs female terms	(10)	91	1.47	10^{-24}	8 × 2	8 × 2	1.24	10^{-2}
Mental vs physical disease	Temporary vs permanent	(23)	135	1.01	10^{-3}	6 × 2	7 × 2	1.38	10^{-2}
Young vs old people's names	Pleasant vs unpleasant	(9)	43k	1.42	$< 10^{-2}$	8 × 2	8 × 2	1.21	10^{-2}

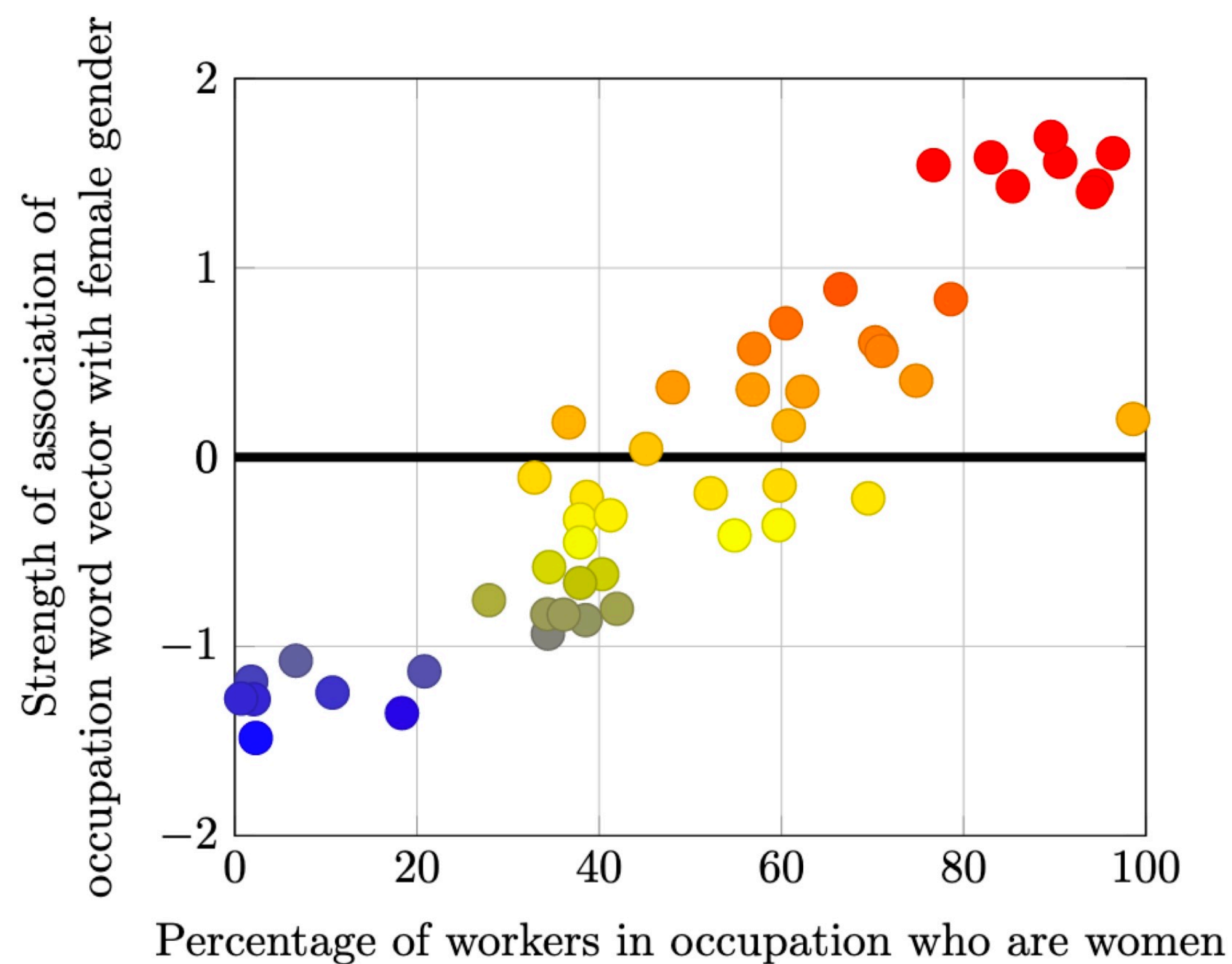


Figure 1: Occupation-gender association. Pearson's correlation coefficient $\rho = 0.90$ with p -value $< 10^{-18}$.

Word Embedding Factual Association Test

Target word

Target attributes

$$s(w, A, B) = \frac{\text{mean}_{a \in A} \cos(\vec{w}, \vec{a}) - \text{mean}_{b \in B} \cos(\vec{w}, \vec{b})}{\text{std-dev}_{x \in A \cup B} \cos(\vec{w}, \vec{x})}$$

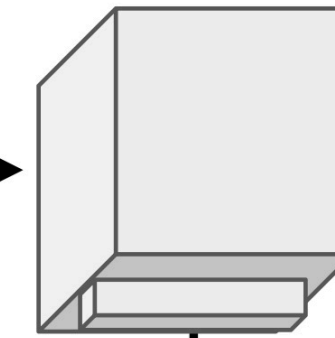
Visual semantic role labeling

imSitu Visual Semantic Role Labeling (vSRL)

[Yatskar et al. 2016]



Convolutional
Neural Network



COOKING (events)

ROLES	NOUNS
AGENT	woman
FOOD	vegetable
CONTAINER	pot
TOOL	spatula



Regression

Conditional Random Field

Identifying data bias

$$b(o, g) = \frac{c(o, g)}{\sum_{g' \in G} c(o, g')},$$

where $c(o, g)$ is the number of occurrences of o and g in a corpus. For example, to analyze how genders of agents and activities are co-related in vSRL, we define the gender bias toward man for each verb $b(\text{verb}, \text{man})$ as:

$$\frac{c(\text{verb}, \text{man})}{c(\text{verb}, \text{man}) + c(\text{verb}, \text{woman})}. \quad (1)$$




If $b(o, g) > 1/\|G\|$, then o is positively correlated with g and may exhibit bias.

Defining dataset bias

Events

Training Gender Ratio ( verb)

Training Set

-  cooking
-  woman
-  man



	COOKING	
	ROLES	NOUNS
	AGENT	woman
	FOOD	stir-fry



	COOKING	
	ROLES	NOUNS
	AGENT	man
	FOOD	noodle

$$\frac{\#(\text{red diamond cooking}, \text{blue circle man})}{\#(\text{red diamond cooking}, \text{blue circle man}) + \#(\text{red diamond cooking}, \text{orange circle woman})} = 1/3$$

Defining dataset bias

Objects

Training Gender Ratio (▲ noun)

Training Set

- ▲ snowboard
- woman
- man



●	MAN	
▲	snowboard	yes
	refrigerator	no
	bowl	no



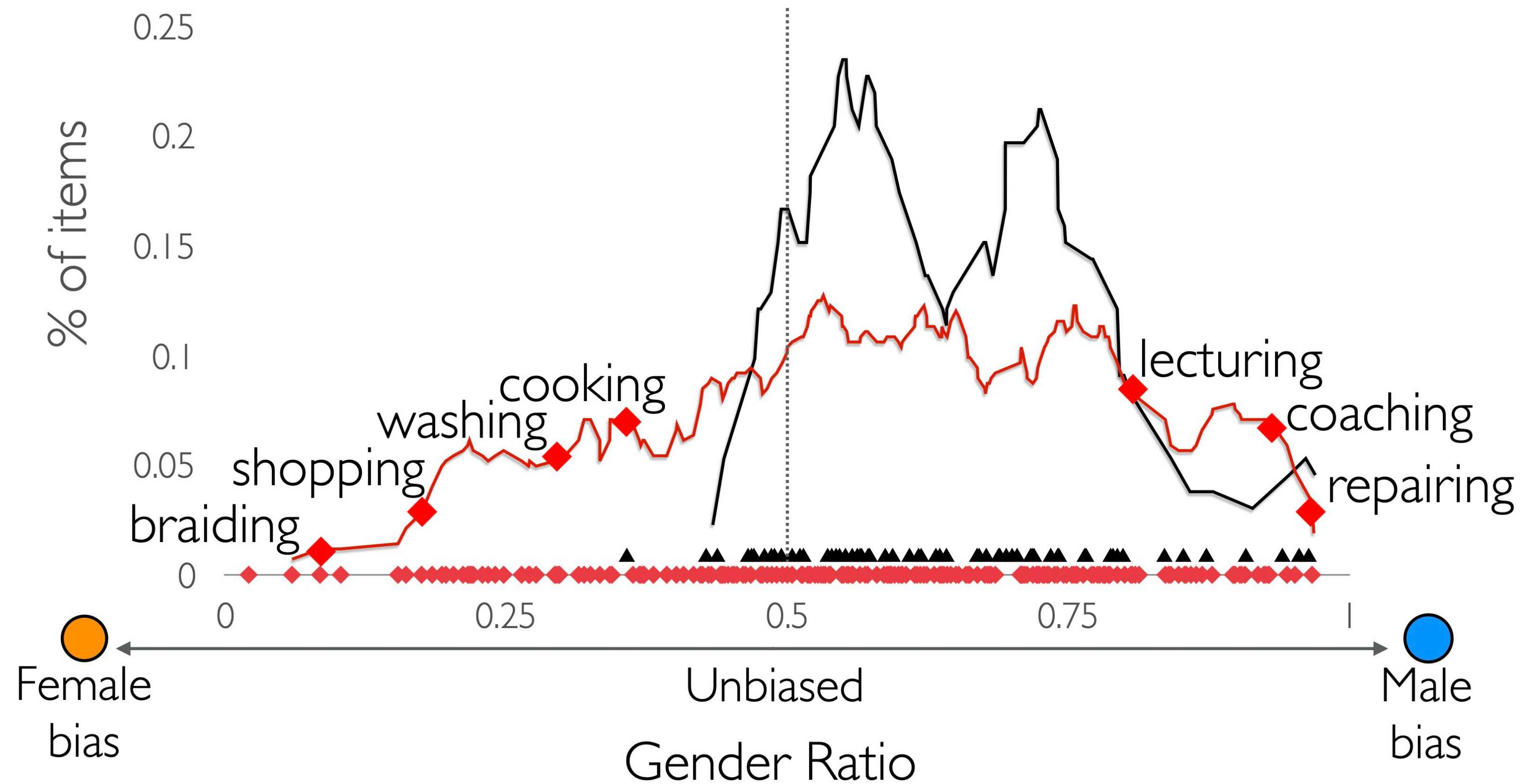
●	WOMAN	
▲	snowboard	yes
	refrigerator	no
	bowl	no

$$\frac{\#(\blacktriangle \text{ snowboard}, \bullet \text{ man})}{\#(\blacktriangle \text{ snowboard}, \bullet \text{ man}) + \#(\blacktriangle \text{ snowboard}, \bullet \text{ woman})} = 2/3$$

Gender dataset bias

◆ imSitu Verb

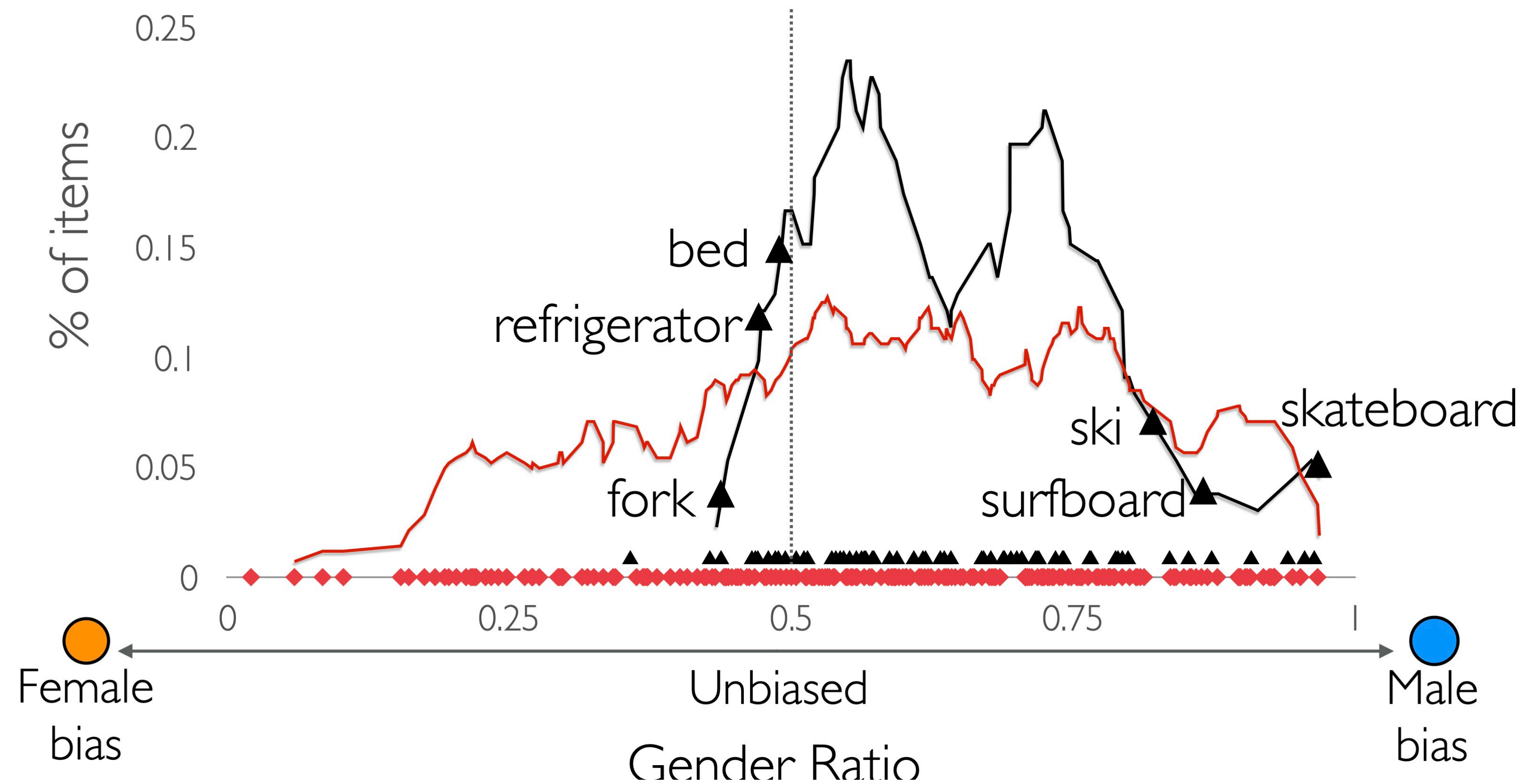
▲ COCO Noun



Gender dataset bias

◆ imSitu Verb

▲ COCO Noun



Gender dataset bias

◆ imSitu Verb

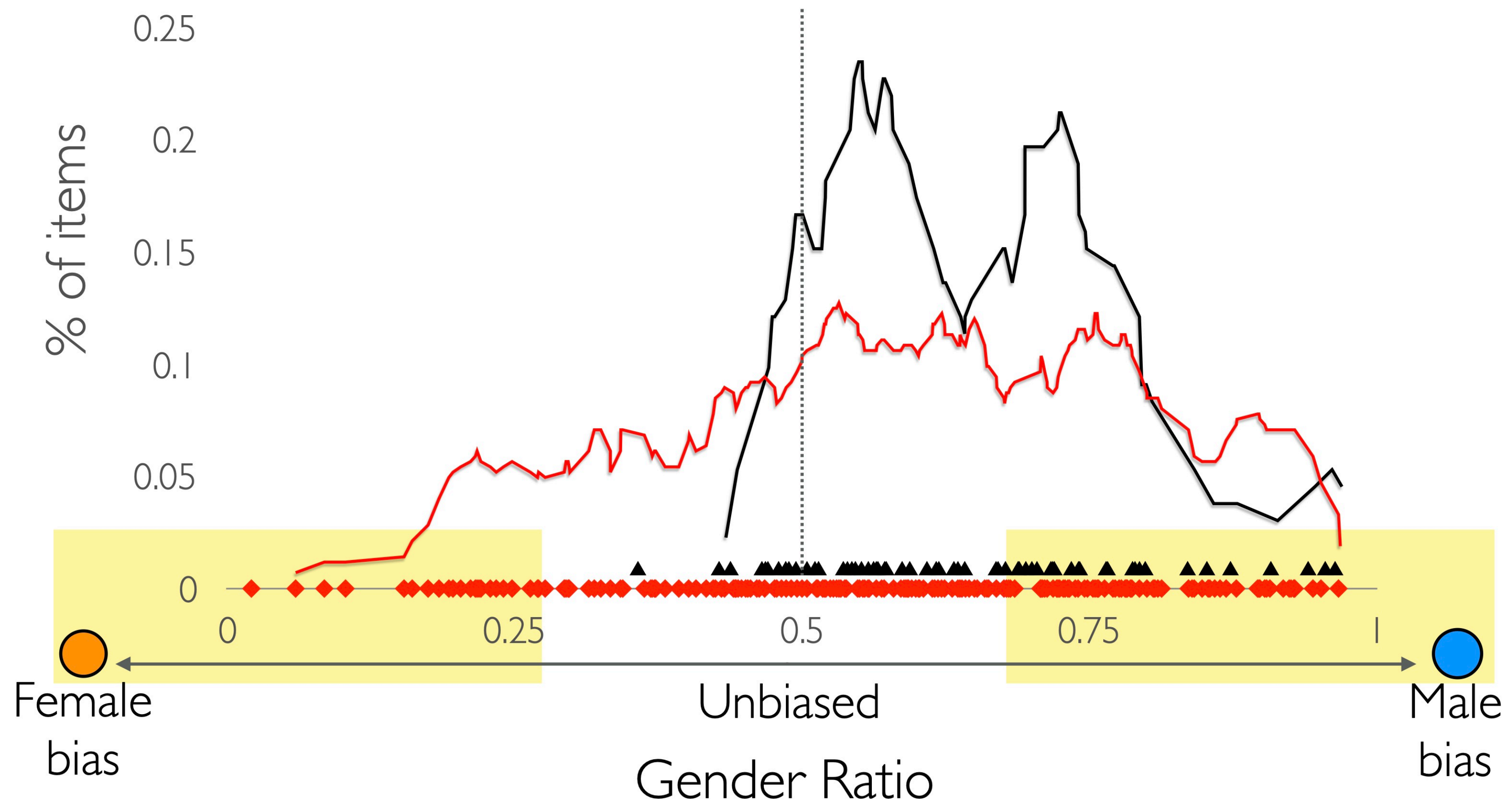
64.6% bias

46.9% strong bias (>2:1)

▲ COCO Noun

86.6% bias

37.9% strong bias (>2:1)



Evaluating bias amplification

$$\frac{1}{|O|} \sum_g \sum_{o \in \{o \in O \mid b^*(o, g) > 1/\|G\|\}} \tilde{b}(o, g) - b^*(o, g).$$

-

- $\tilde{b}(o, g)$: bias score on unlabeled evaluation set of images that has been annotated by a predictor
- $b^*(o, g)$: bias score on training set

Evaluating bias amplification

Predicted Gender Ratio (◆ verb)

Development Set

- ◆ cooking
- woman
- man

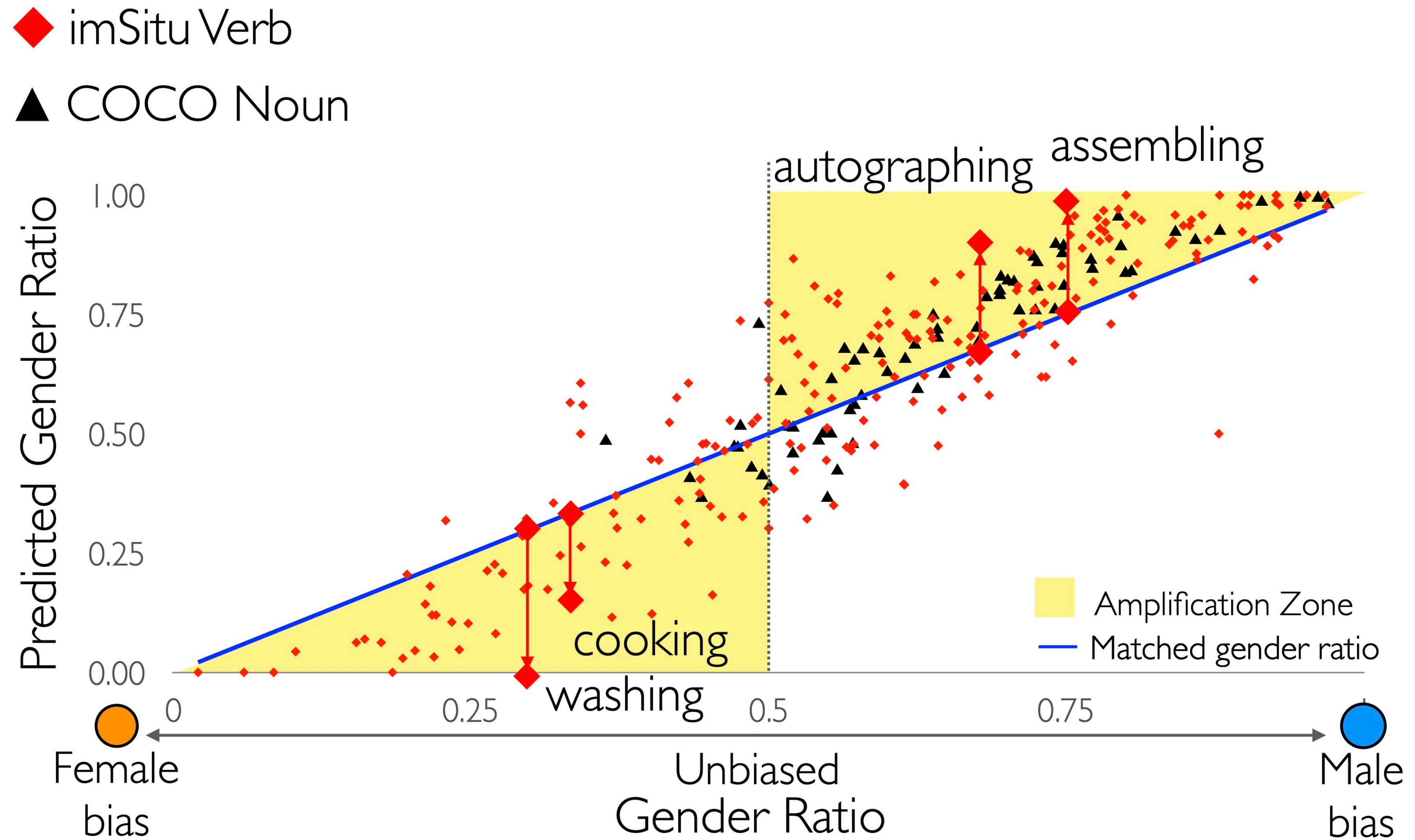


COOKING	
ROLES	NOUNS
● AGENT	woman
FOOD	stir-fry

COOKING	
ROLES	NOUNS
● AGENT	man
FOOD	noodle

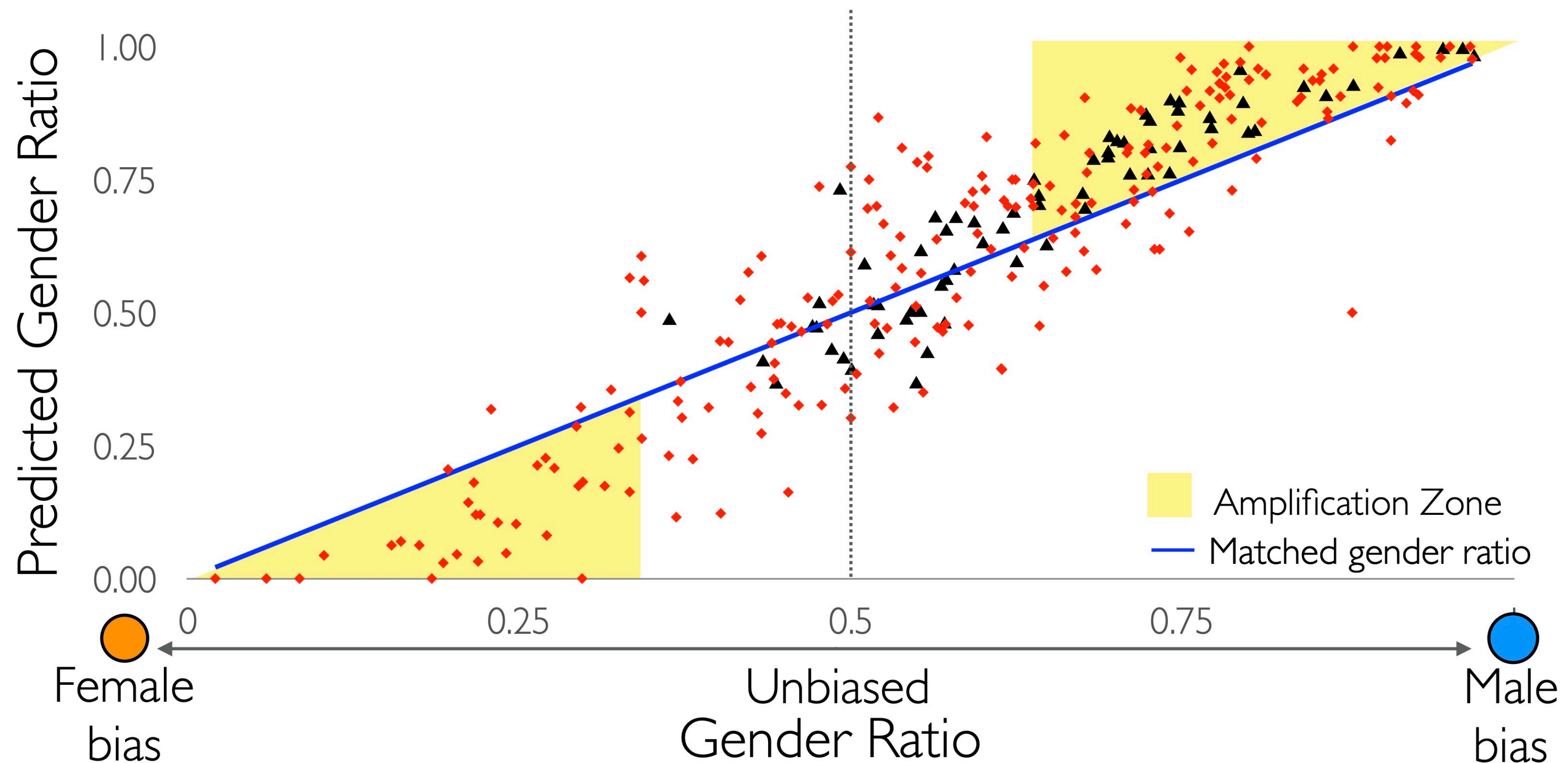
$$\frac{\#(\text{◆ cooking}, \text{● man})}{\#(\text{◆ cooking}, \text{● man}) + \#(\text{◆ cooking}, \text{● woman})} = 1/6$$

Model bias amplification



Model bias amplification

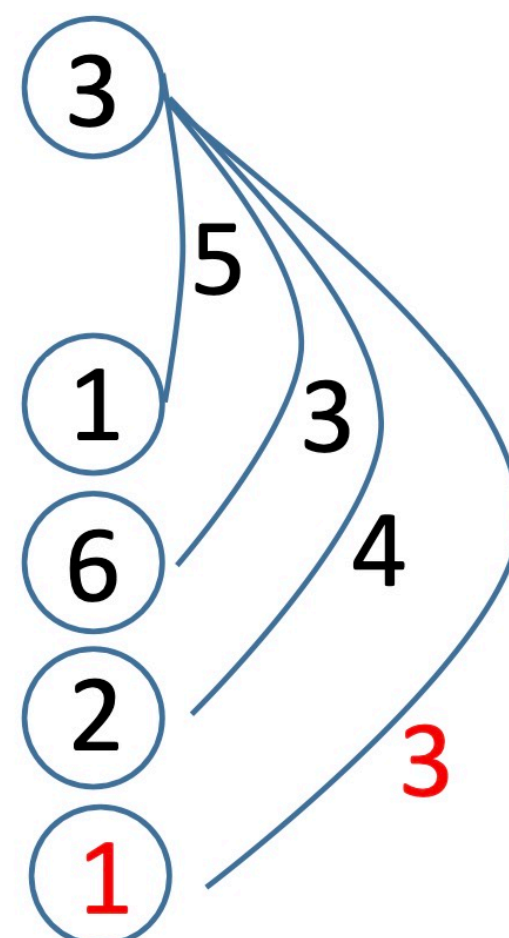
- ◆ imSitu Verb 69% bias↑ .05 |bias↑| > 2:1 initial bias : .07 |bias↑|
- ▲ COCO Noun 73% bias↑ .04 |bias↑| > 2:1 initial bias : .08 |bias↑|



Decomposition of scoring function

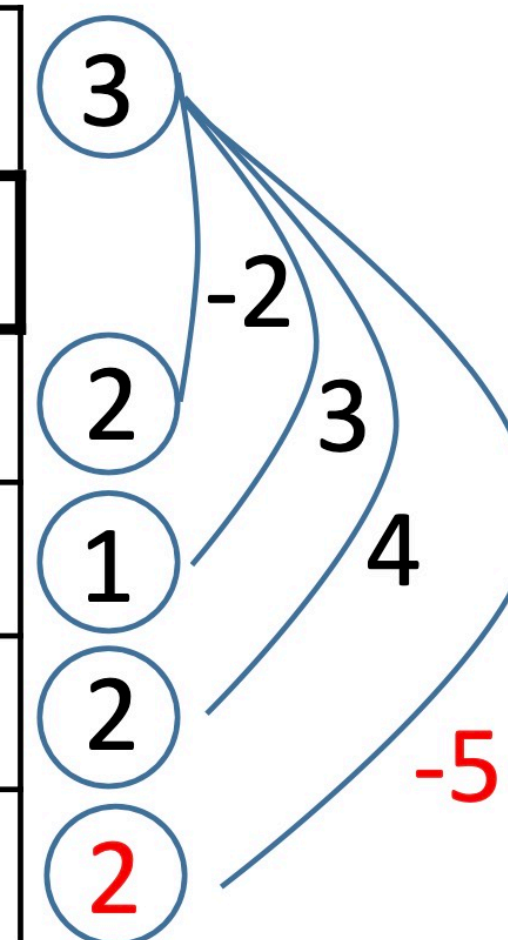


COOKING	
ROLES	NOUNS
AGENT	woman
FOOD	vegetable
CONTAINER	pot
TOOL	spatula



...

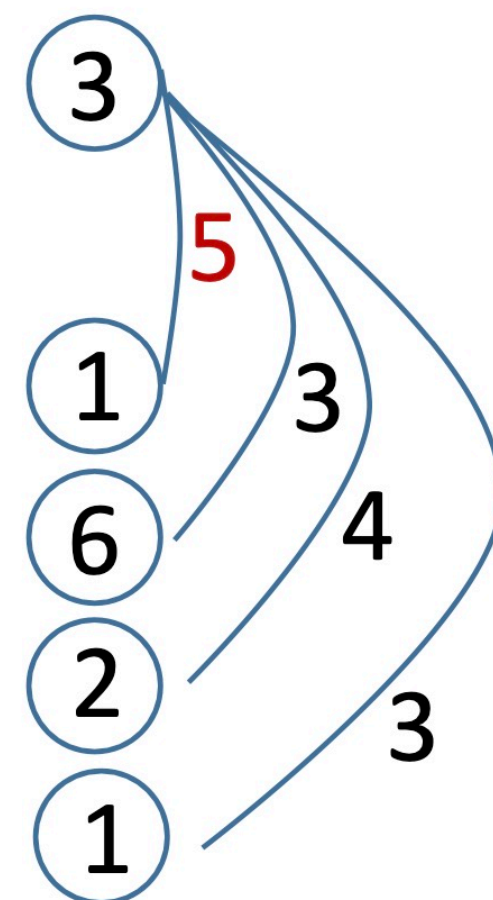
COOKING	
ROLES	NOUNS
AGENT	man
FOOD	meat
CONTAINER	pot
TOOL	screwdriver



Decomposition of scoring function

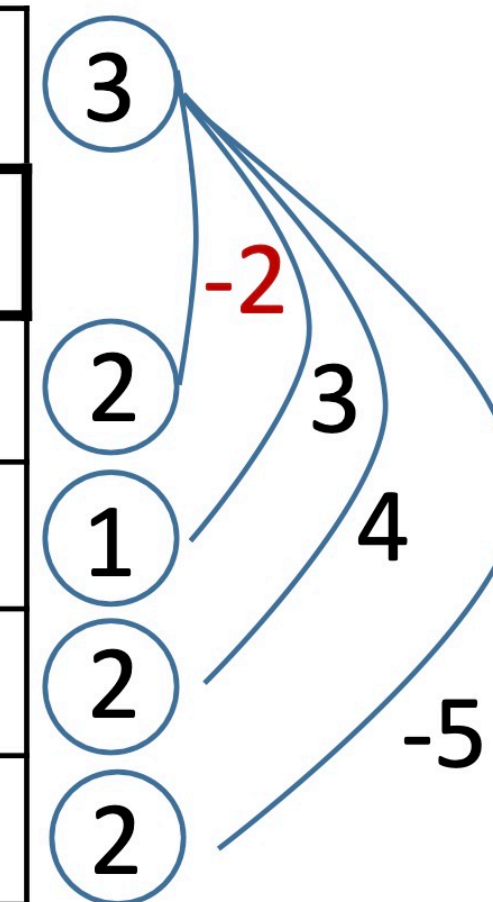


COOKING	
ROLES	NOUNS
AGENT	woman
FOOD	vegetable
CONTAINER	pot
TOOL	spatula



...

COOKING	
ROLES	NOUNS
AGENT	man
FOOD	meat
CONTAINER	pot
TOOL	screwdriver



Decomposition of scoring function

Intuition of Calibration

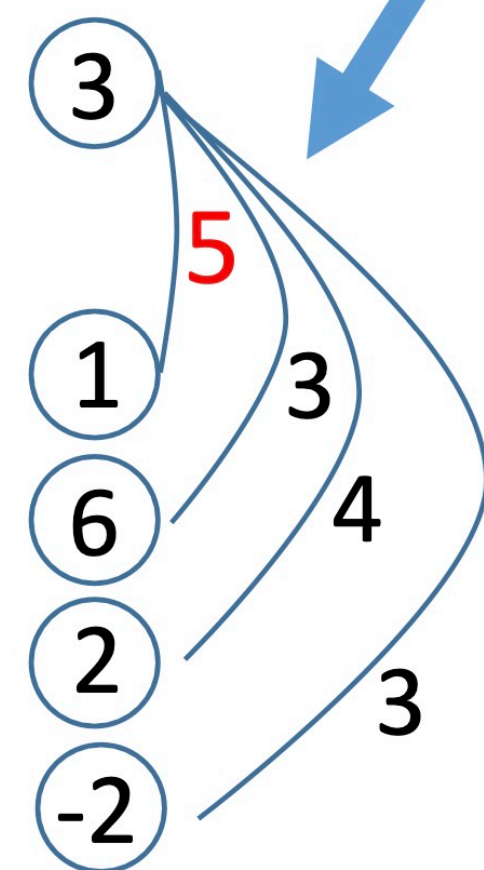
$$\lambda_1, \lambda_2 > 0$$



$$5 \rightarrow 5 - \lambda_1$$

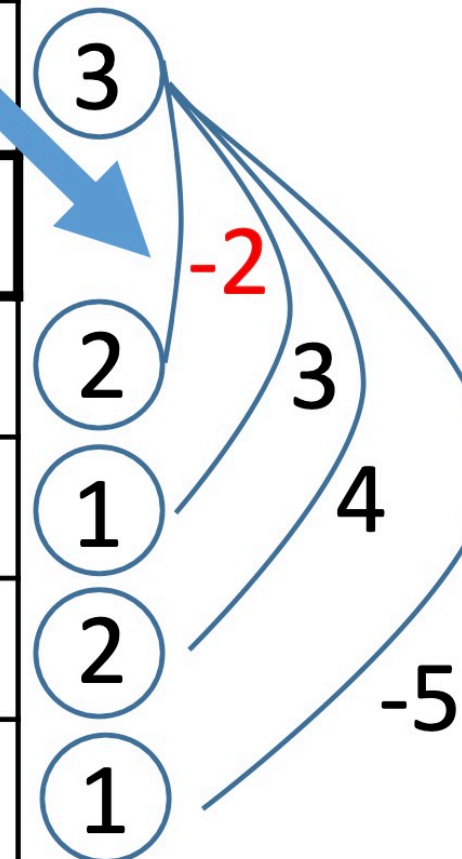
$$-2 \rightarrow -2 + \lambda_2$$

COOKING	
ROLES	NOUNS
AGENT	woman
FOOD	vegetable
CONTAINER	pot
TOOL	spatula



...

COOKING	
ROLES	NOUNS
AGENT	man
FOOD	meat
CONTAINER	pot
TOOL	screwdriver



Reducing bias amplification

Integer Linear Program

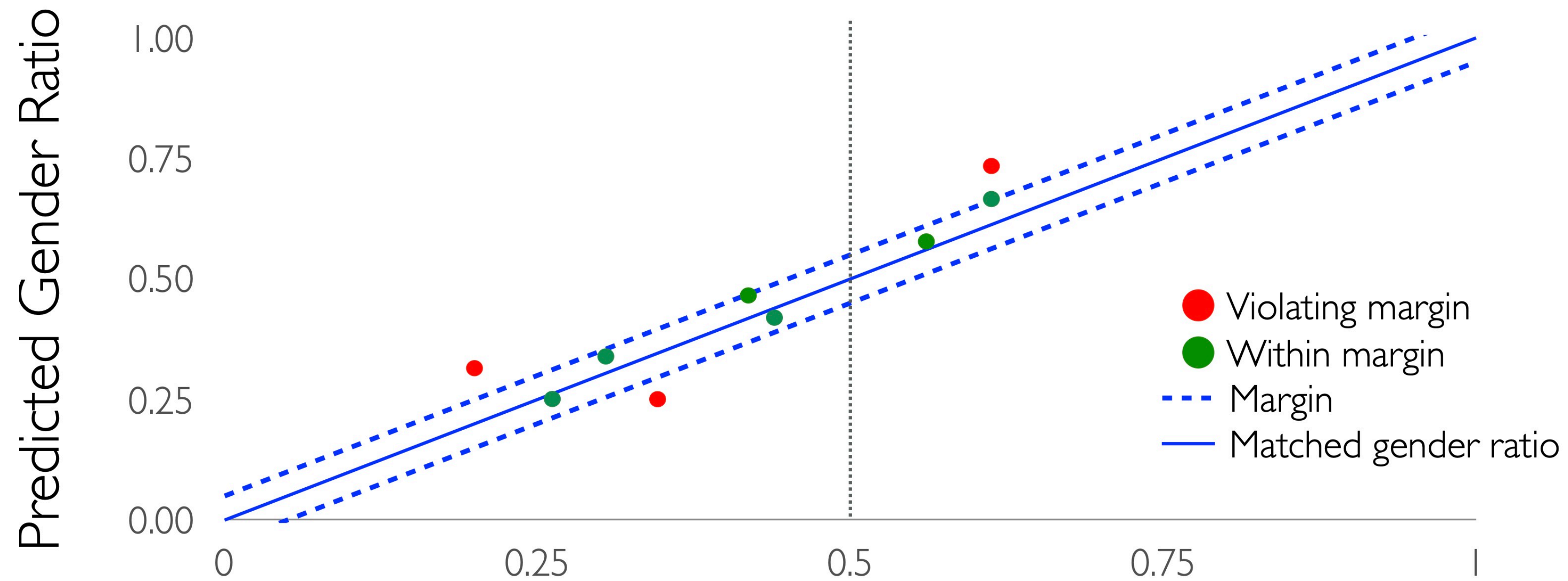
$$\sum_i \max_{y_i} s(y_i, \text{image})$$

\forall points

$$\left| \text{Training Ratio} - \text{Predicted Ratio} \right| \leq \text{margin}$$

$f(y_1 \dots y_n)$

\leq margin



Reducing bias amplification

Integer Linear Program

$$\sum_i \max_{y_i} s(y_i, \text{image})$$

\forall points

$$\left| \text{Training Ratio} - \text{Predicted Ratio} \right|$$

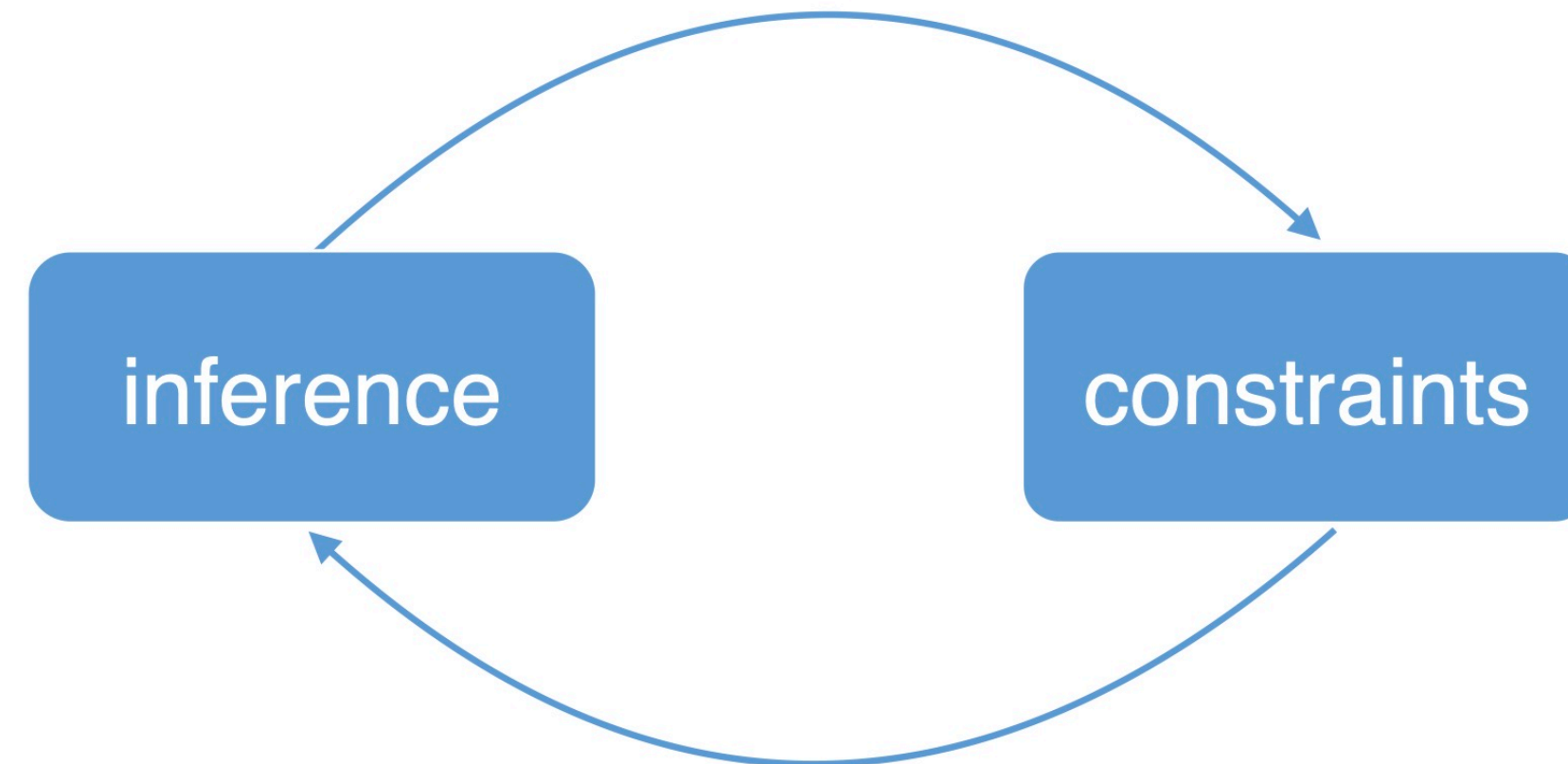
$f(y_1 \dots y_n)$

\leq margin

Lagrangian Relaxation

inference

constraints



Reducing bias amplification

$$\max_{y_i} \sum_i s(y_i, \text{image})$$
$$\left| \text{Training Ratio} - \frac{\text{Predicted Ratio}}{f(y_1 \dots y_n)} \right| \leq \text{margin}$$

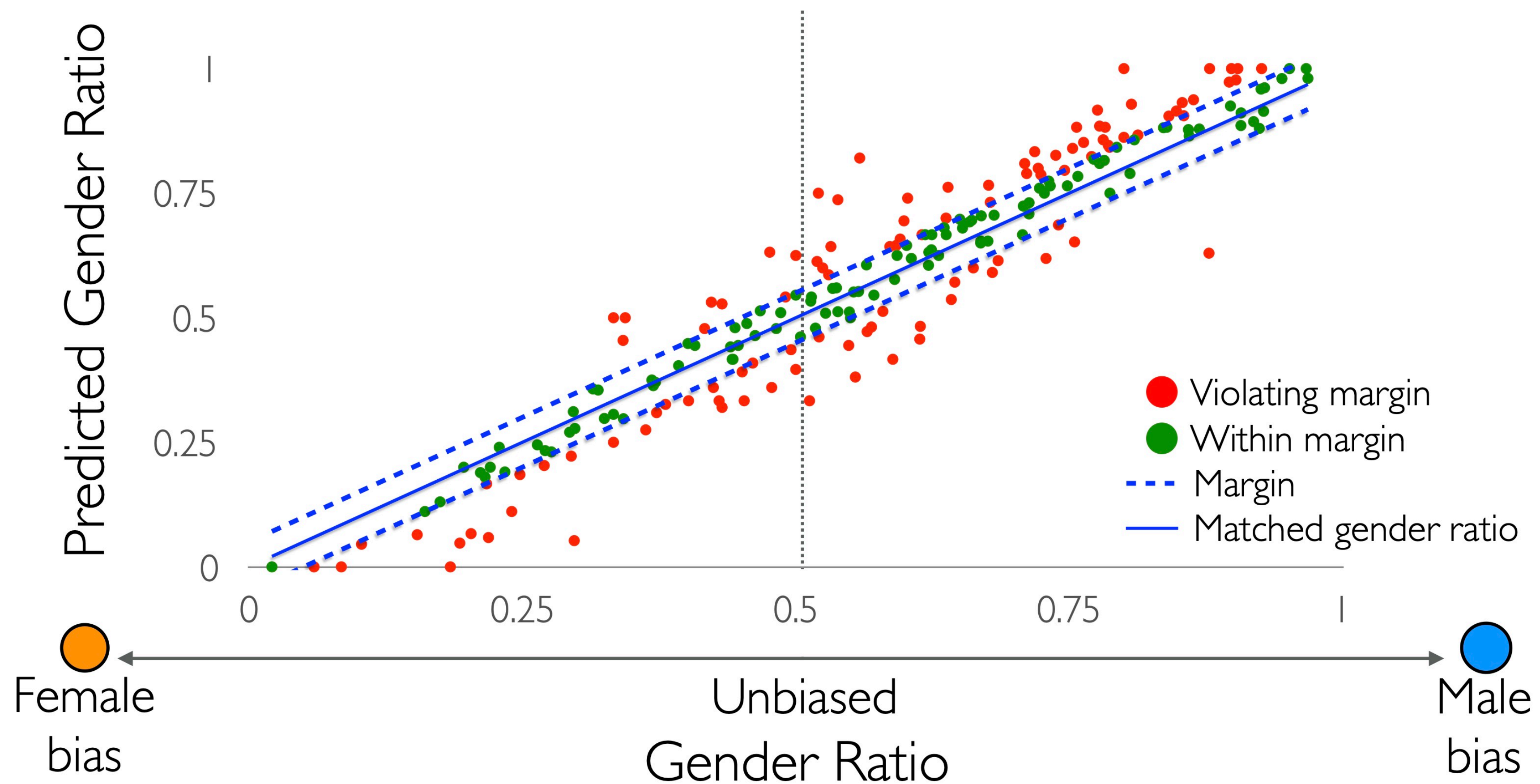
Lagrangian Relaxation

$$\max_{\{y^i\} \in \{Y^i\}} \sum_i f_{\theta}(y^i, i), \quad \text{s.t.} \quad A \sum_i y^i - b \leq 0$$

Lagrangian : $\sum_i f_{\theta}(y^i) - \sum_{j=1}^l \lambda_j (A_j \sum_i y^i - b_j) \quad \lambda_j \geq 0$

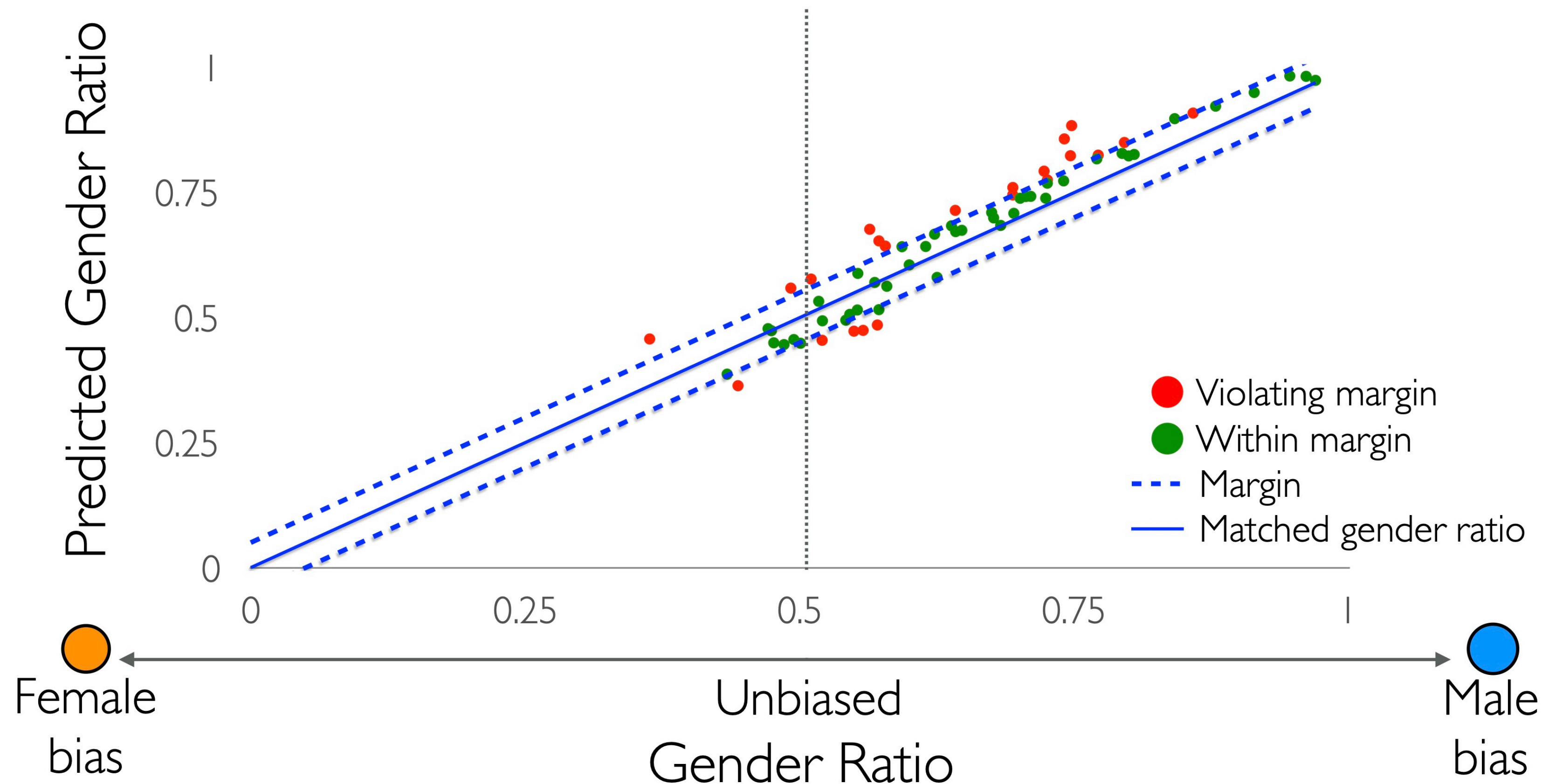
Reducing bias amplification

imSitu Verb	Violation: 72.6%	.050 bias↑	24.07 acc.
w/ RBA	Violation: 50.5%	.024 bias↑	23.97 acc.



Reducing bias amplification

COCO Noun	Violation: 60.6%	.032 bias↑	45.27	mAP
w/ RBA	Violation: 36.4%	.022 bias↑	45.19	mAP



Credit Application

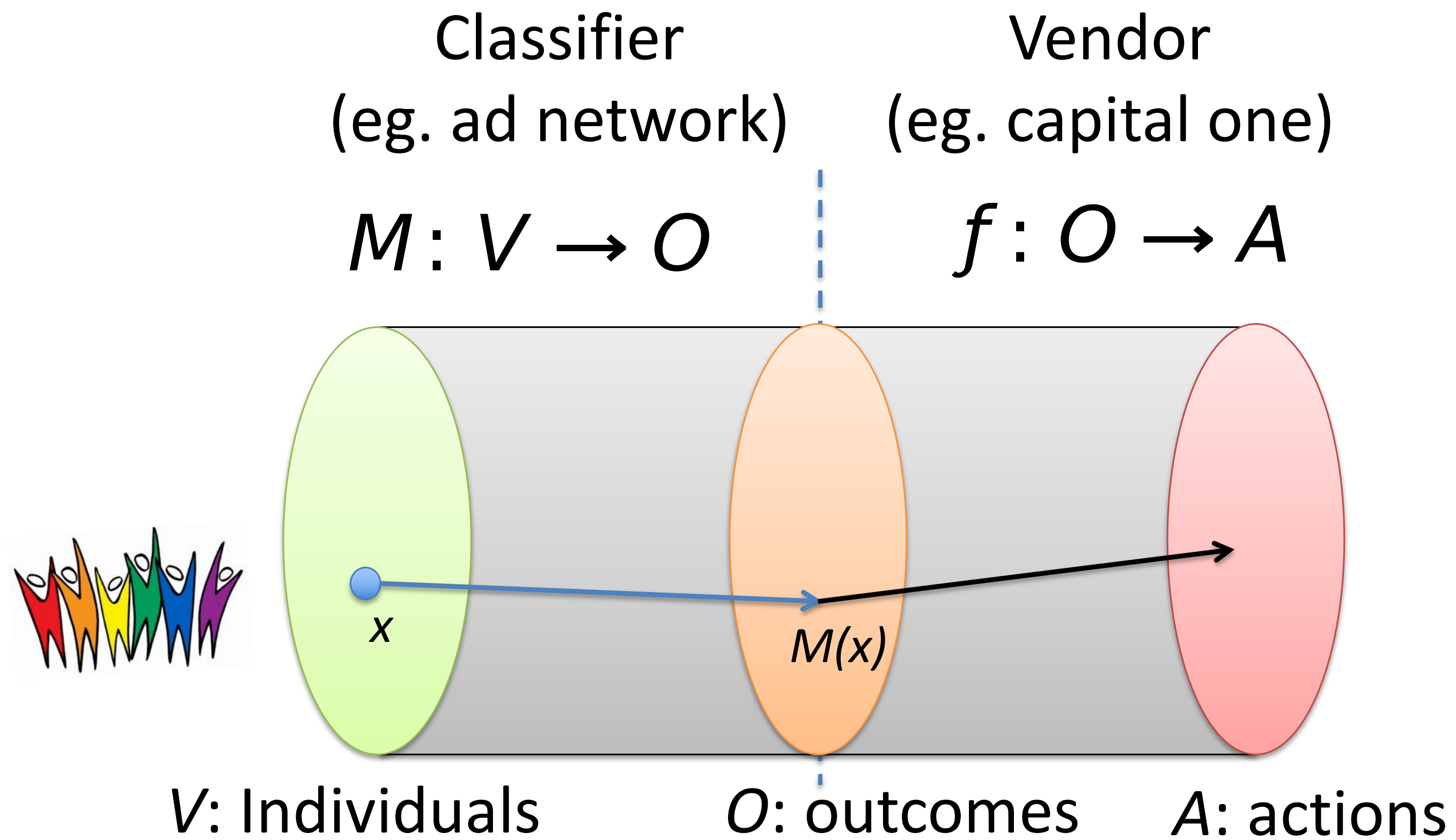


User visits `capitalone.com`

Capital One uses tracking information provided by the tracking network [x+1] to personalize offers

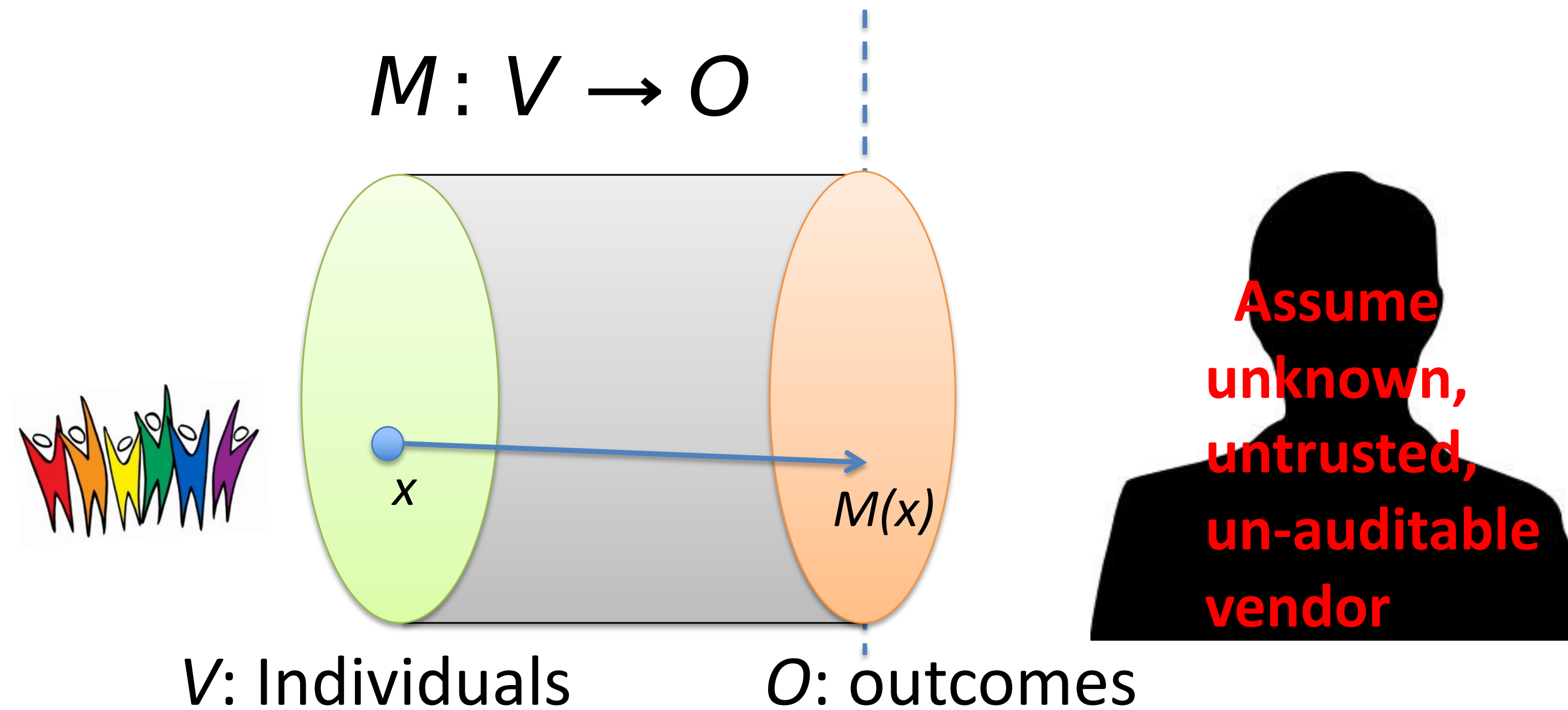
Concern: Steering minorities into higher rates (illegal)

WSJ 2010



Goal:

Achieve Fairness in the classification step



Through blindness

- Ignore all irrelevant/protected attributes
 - You don't need to see an attribute to be able to predict it with high accuracy
 - E.g.: User visits artofmanliness.com ... 90% chance of being male



Individual Fairness

Treat *similar* individuals *similarly*



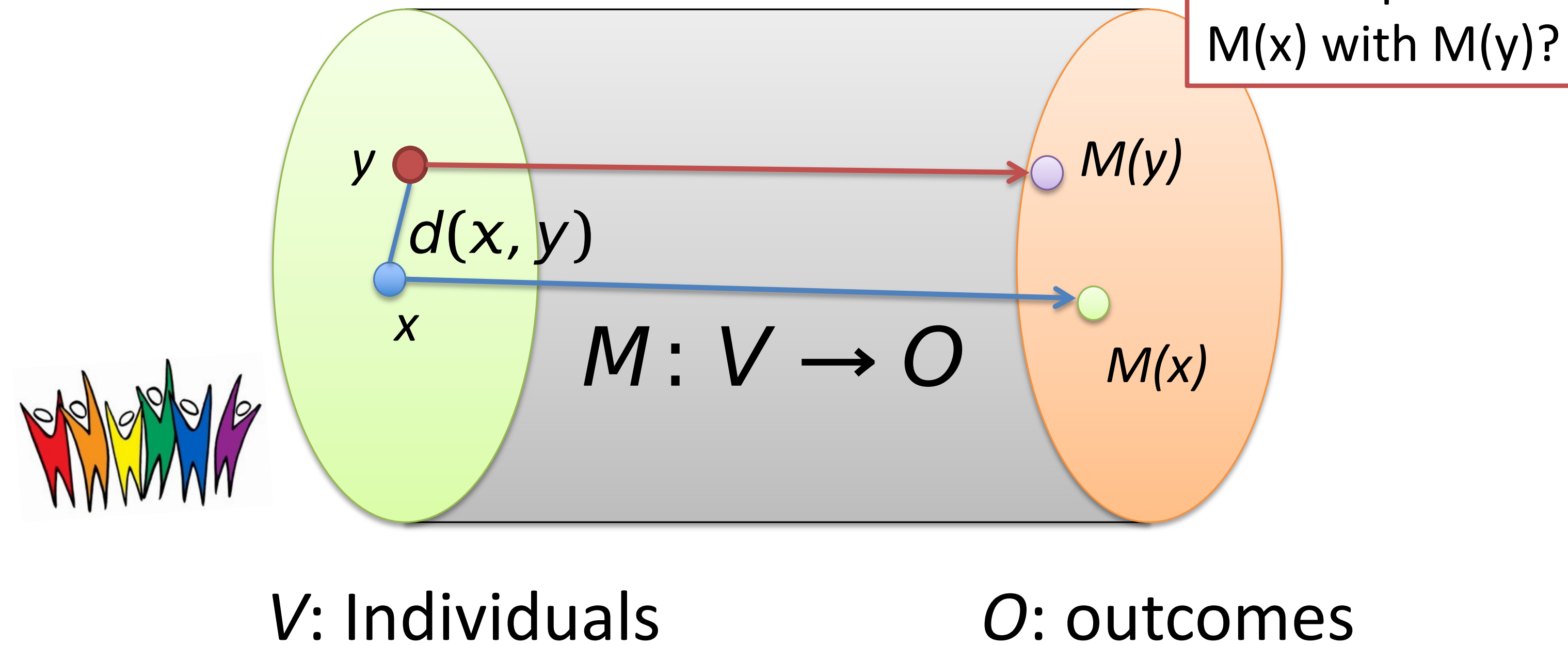
Similar for the purpose of
the classification task



Similar distribution
over outcomes

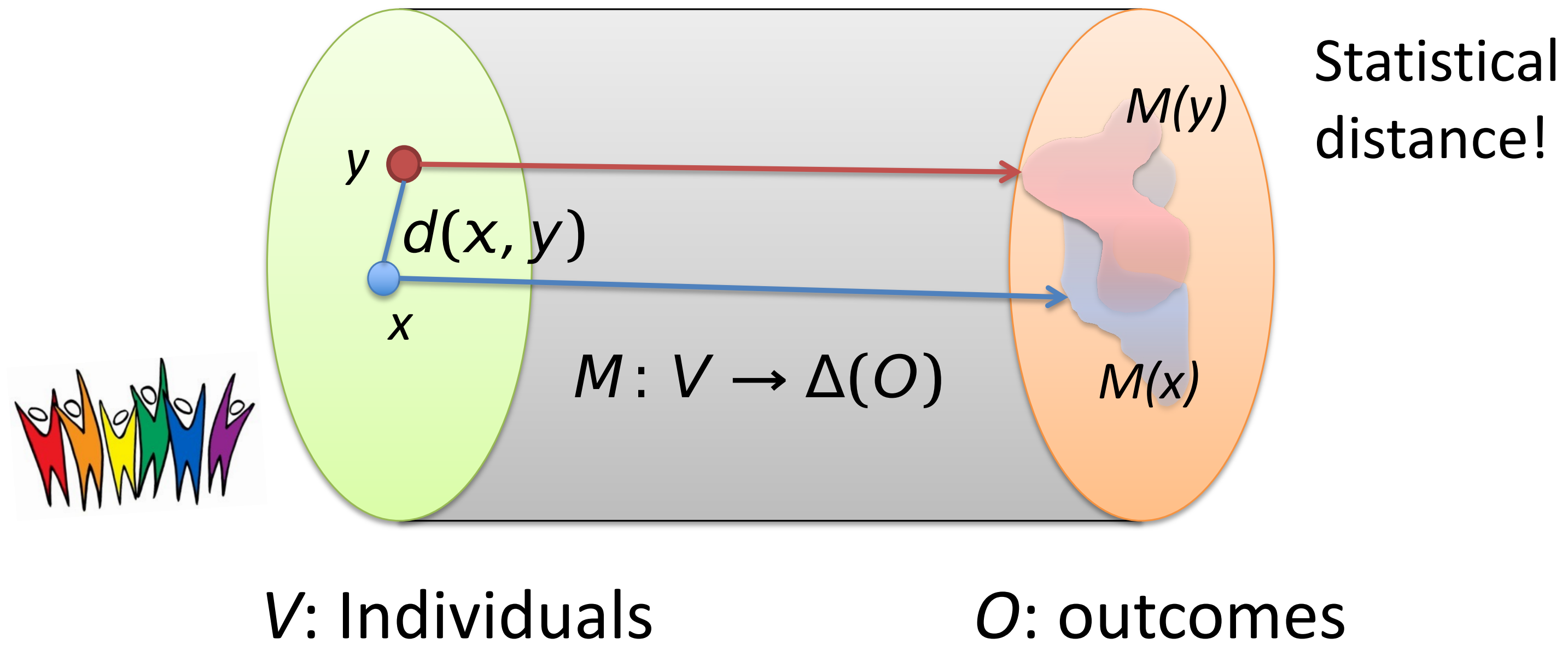
How to formalize this?

Think of V as space
with metric $d(x,y)$
similar = small $d(x,y)$



Distributional outcomes

How can we compare $M(x)$ with $M(y)$?

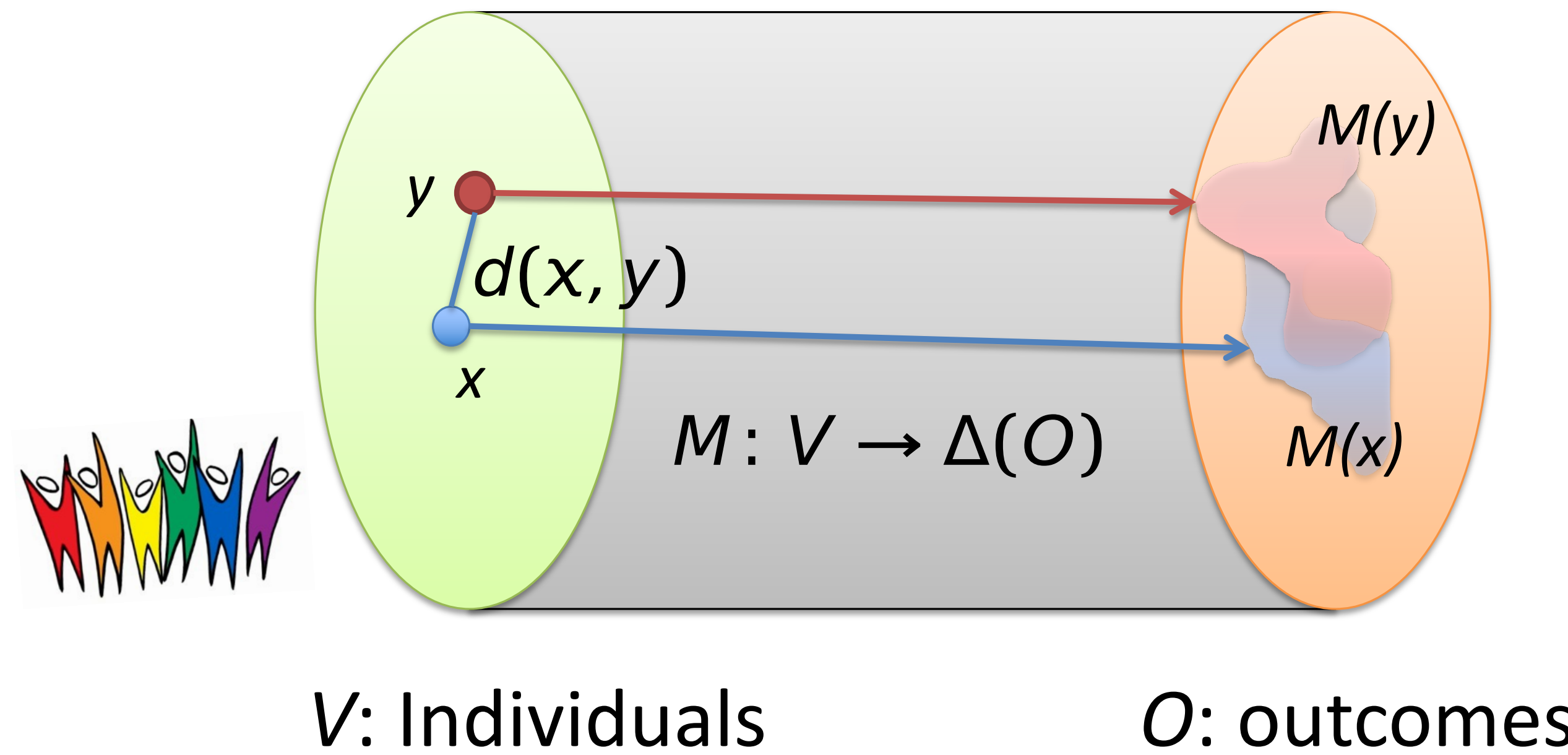


Metric $d: V \times V \rightarrow \mathbb{R}$

Lipschitz condition $\|M(x) - M(y)\| \leq d(x, y)$

This talk: Statistical distance

in $[0,1]$



Statistical Distance

P, Q denote probability measures on a finite domain A . The *statistical distance* between P and Q is denoted by

$$D_{\text{tv}}(P, Q) = \frac{1}{2} \sum_{a \in A} |P(a) - Q(a)|.$$

Example: Mid D

$$A = \{0, 1\}$$

$$P(0) = P(1) = \frac{1}{2}$$

$$Q(0) = \frac{3}{4}, Q(1) = \frac{1}{4}$$

$$D(P, Q) = \frac{1}{4}$$

Utility Maximization

Vendor can specify **arbitrary utility function**

$$U: V \times O \rightarrow \mathbb{R}$$

$U(v,o)$ = Vendor's utility of giving individual v
the outcome o

Maximize vendor's expected utility subject to Lipschitz condition

$$\max_{M(x)} \mathbb{E}_{x \sim V} \mathbb{E}_{o \sim M(x)} U(x, o)$$

s.t. M is d -Lipschitz

$$\|M(x) - M(y)\| \leq d(x, y)$$