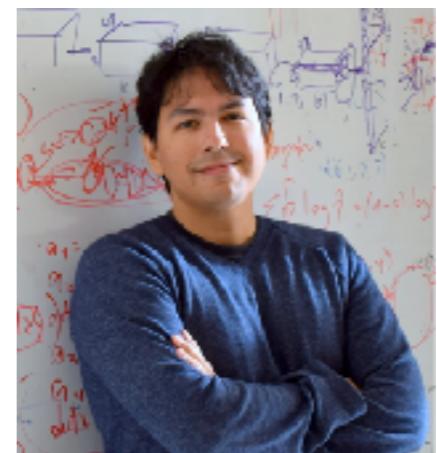


MEN ALSO LIKE SHOPPING

REDUCING GENDER BIAS AMPLIFICATION USING CORPUS-LEVEL CONSTRAINTS

Jieyu Zhao^{1,3}, Tianlu Wang¹, **Mark Yatskar**^{2,4}, Vicente Ordonez¹, Kai-Wei Chang^{1,3}

¹ University of Virginia ² University of Washington ³ UCLA ⁴ Allen Institute for AI



Dataset Gender Bias

Male

Female

33%

66%



Model Bias After Training

16%



84%



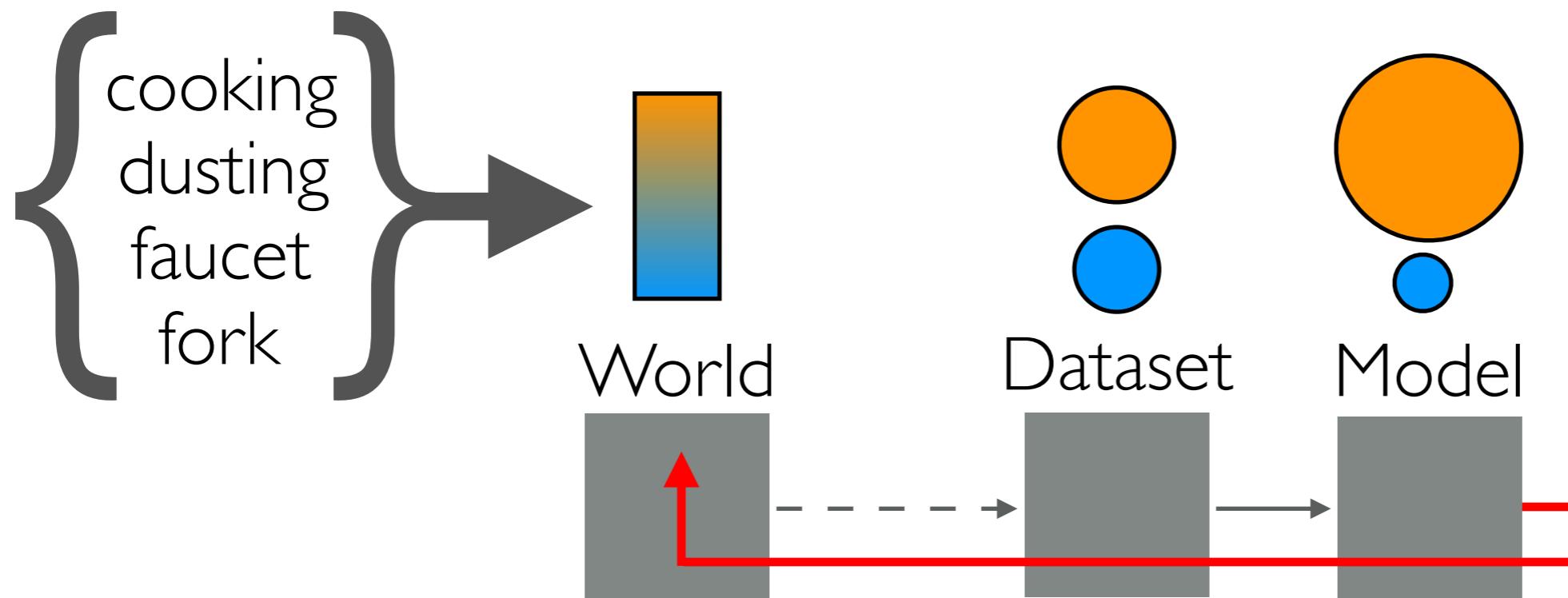
Male

Female

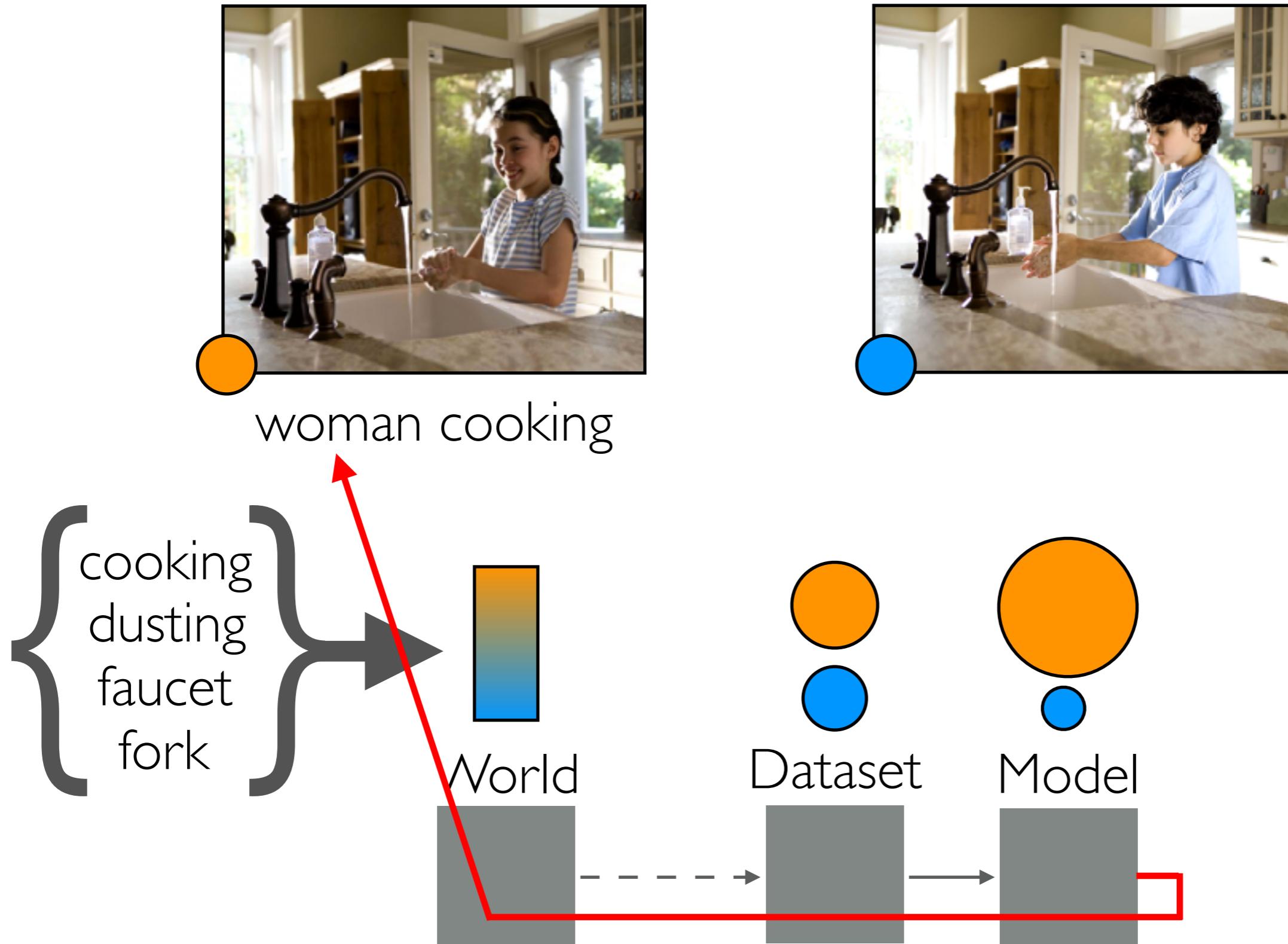
Why does this happen? Good for accuracy



Algorithmic Bias in Grounded Setting



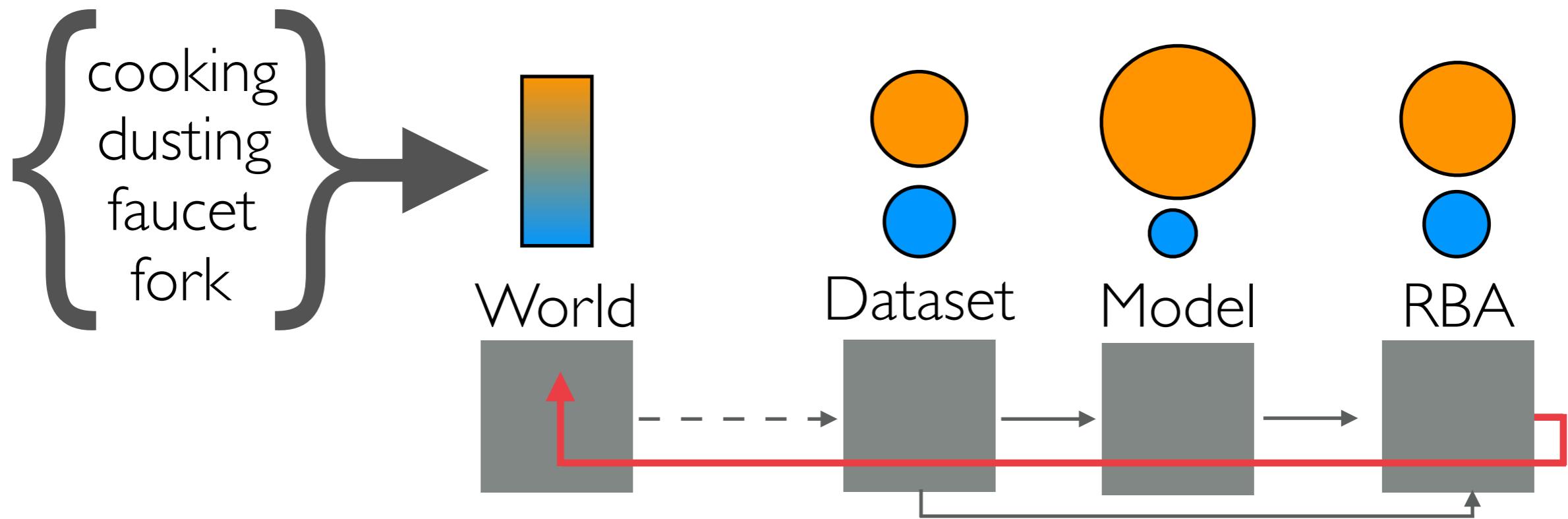
Algorithmic Bias in Grounded Setting



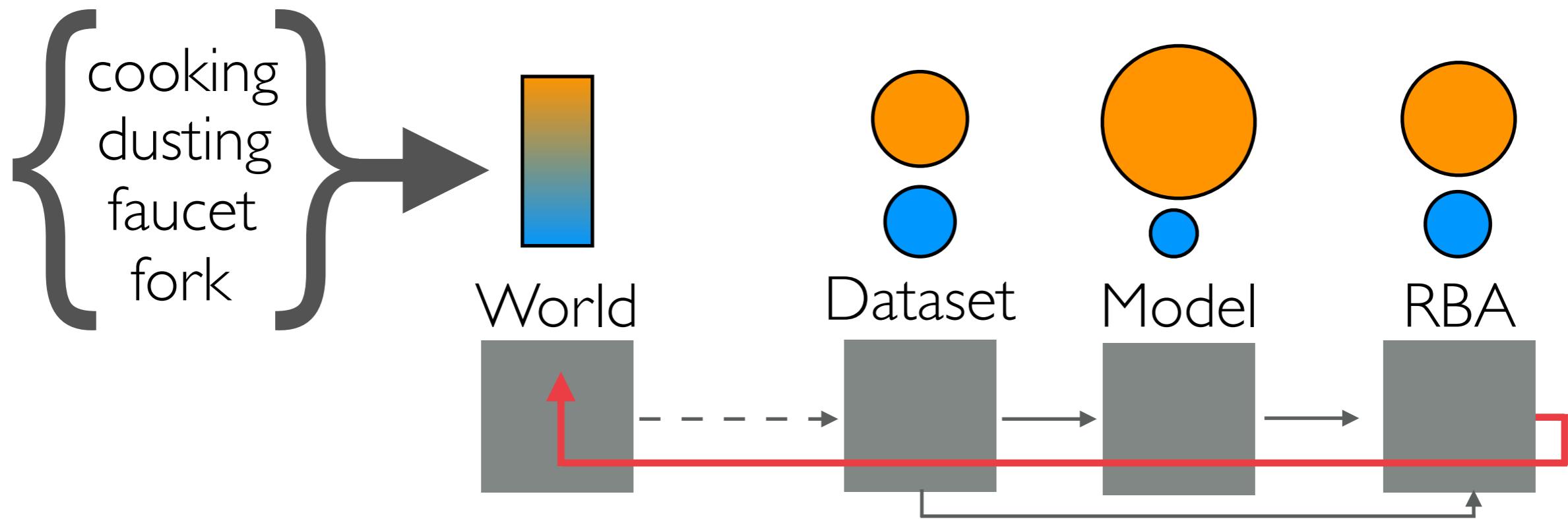
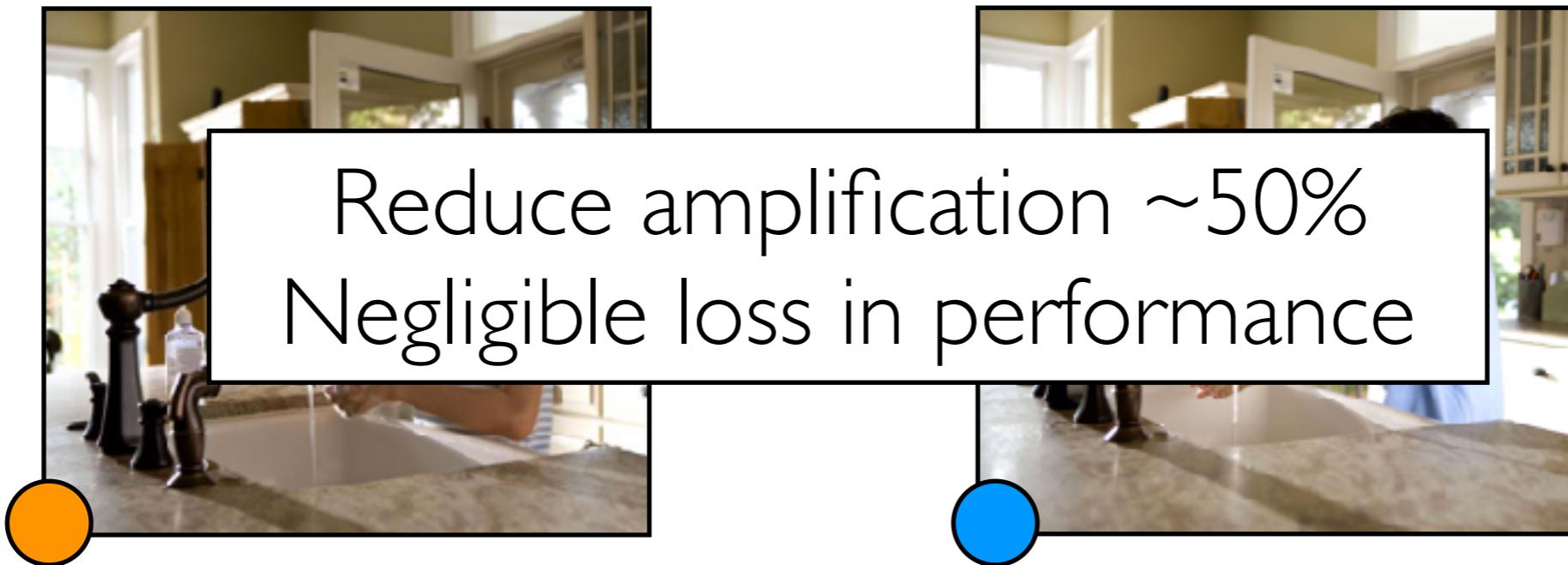
Algorithmic Bias in Grounded Setting



Algorithmic Bias in Grounded Setting



Algorithmic Bias in Grounded Setting



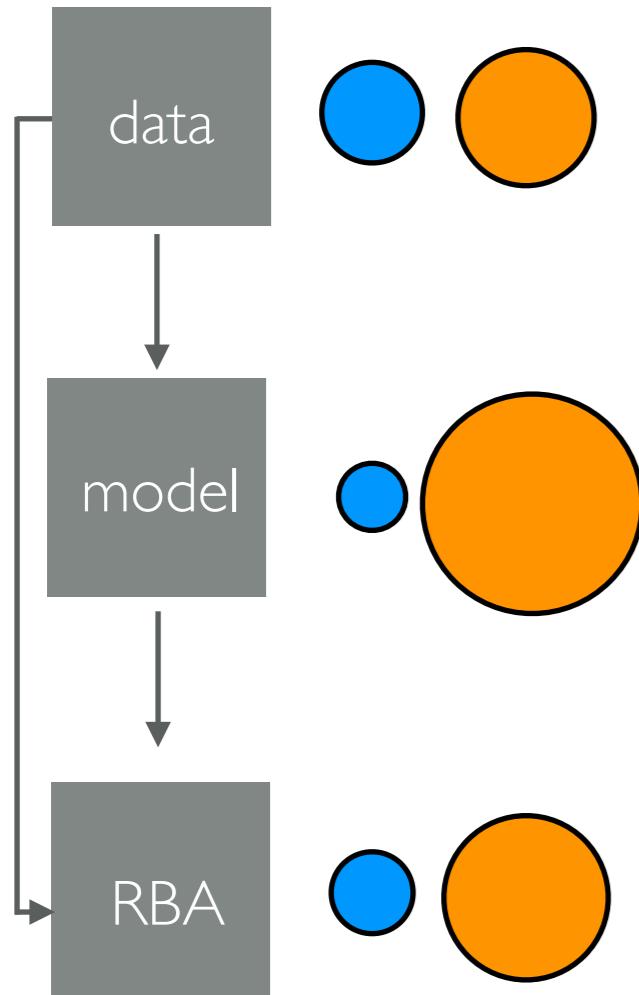
Contributions



imSitu vSRL
(events)



COCO MLC
(objects)



High dataset gender bias
38% (objects) 47% (events) exhibit strong bias

Models amplify existing gender bias
~70% objects and events have bias amplification

Reducing bias amplification
~50% reduction in amplification
Insignificant loss in performance

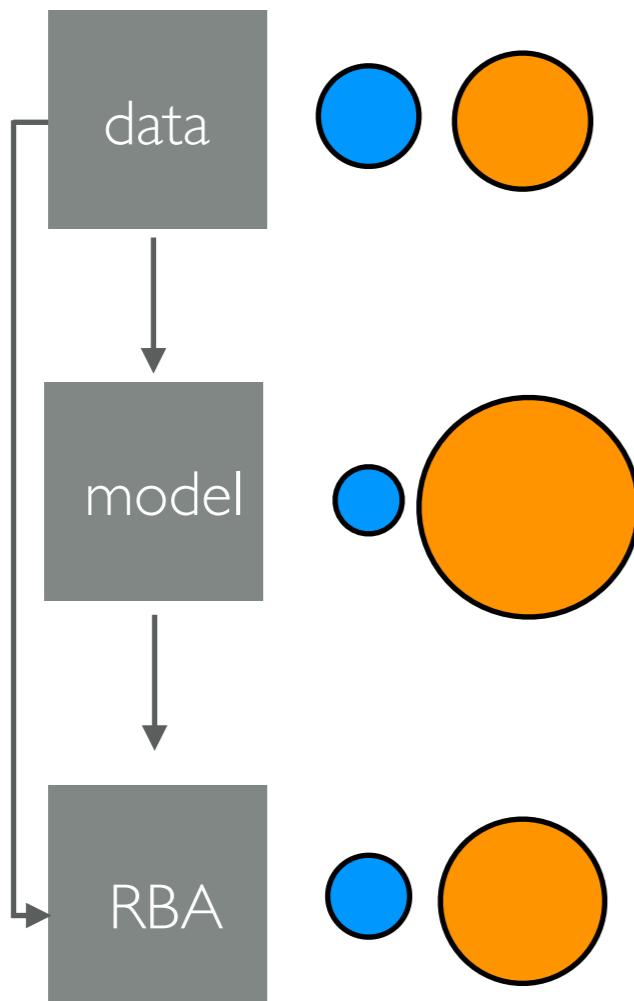
Outline



imSitu vSRL
(events)

COCO MLC
(objects)

I. Background

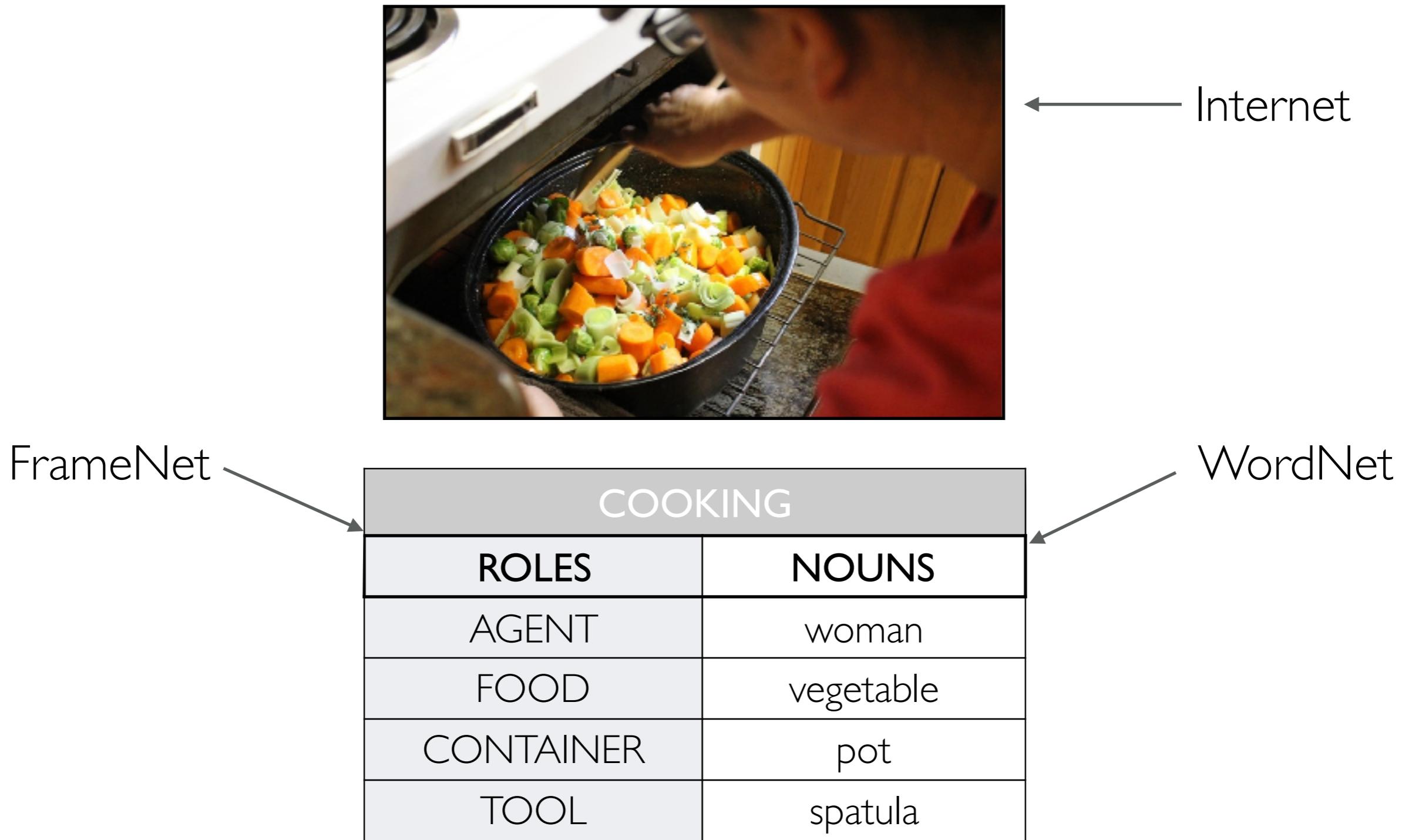


2. Dataset Bias

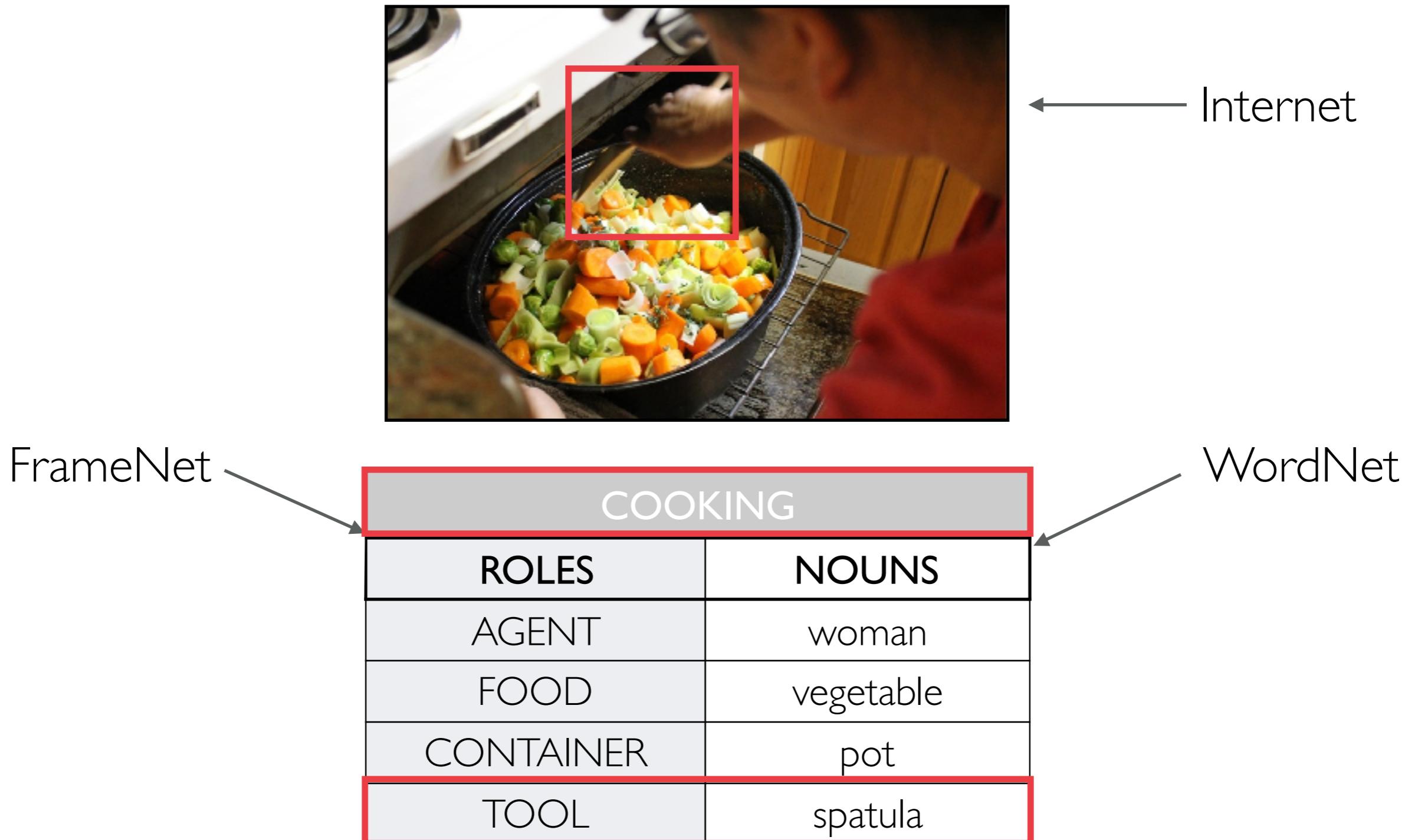
3. Bias Amplification

4. Reducing Bias Amplification

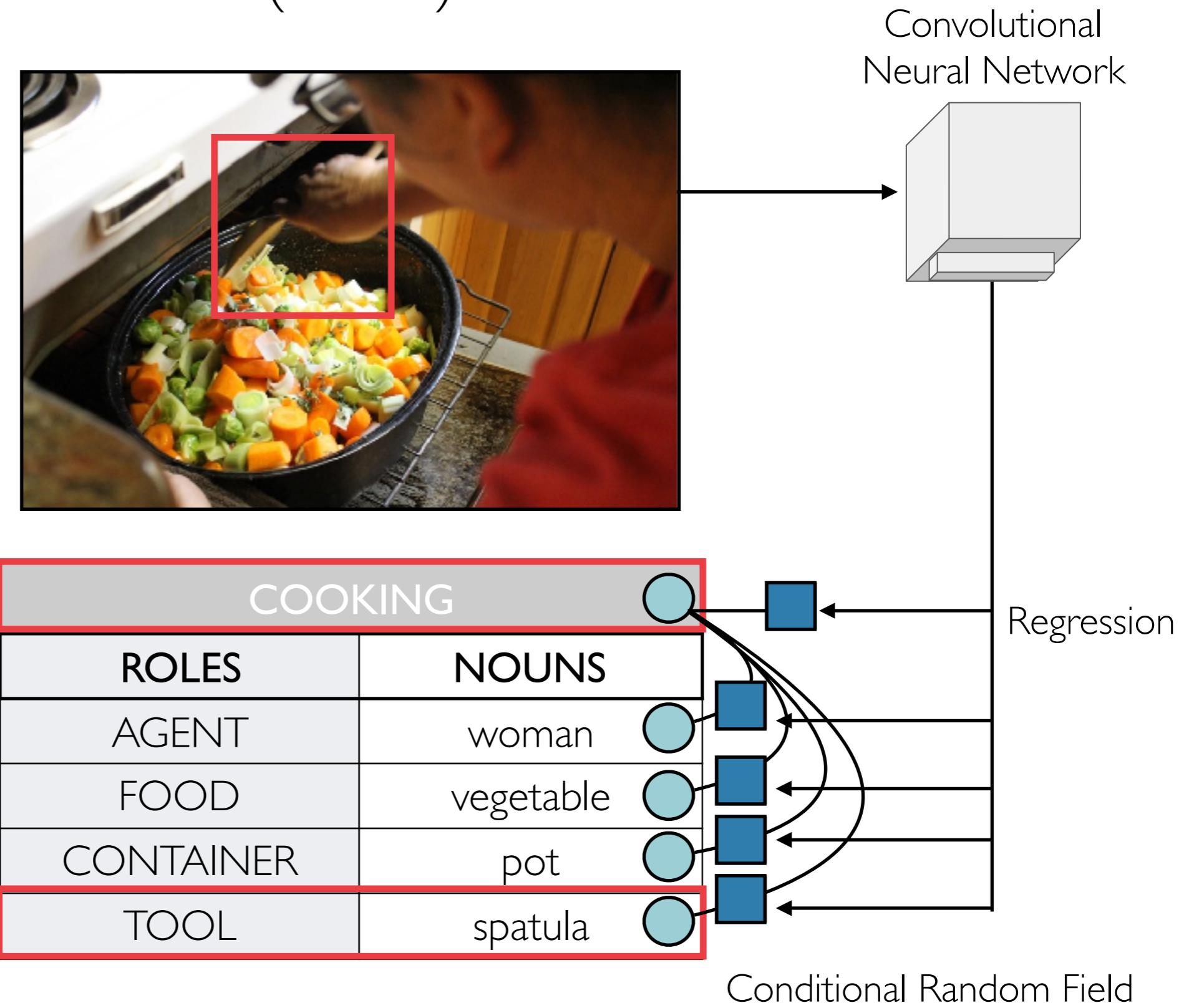
imSitu Visual Semantic Role Labeling (vSRL) (events)



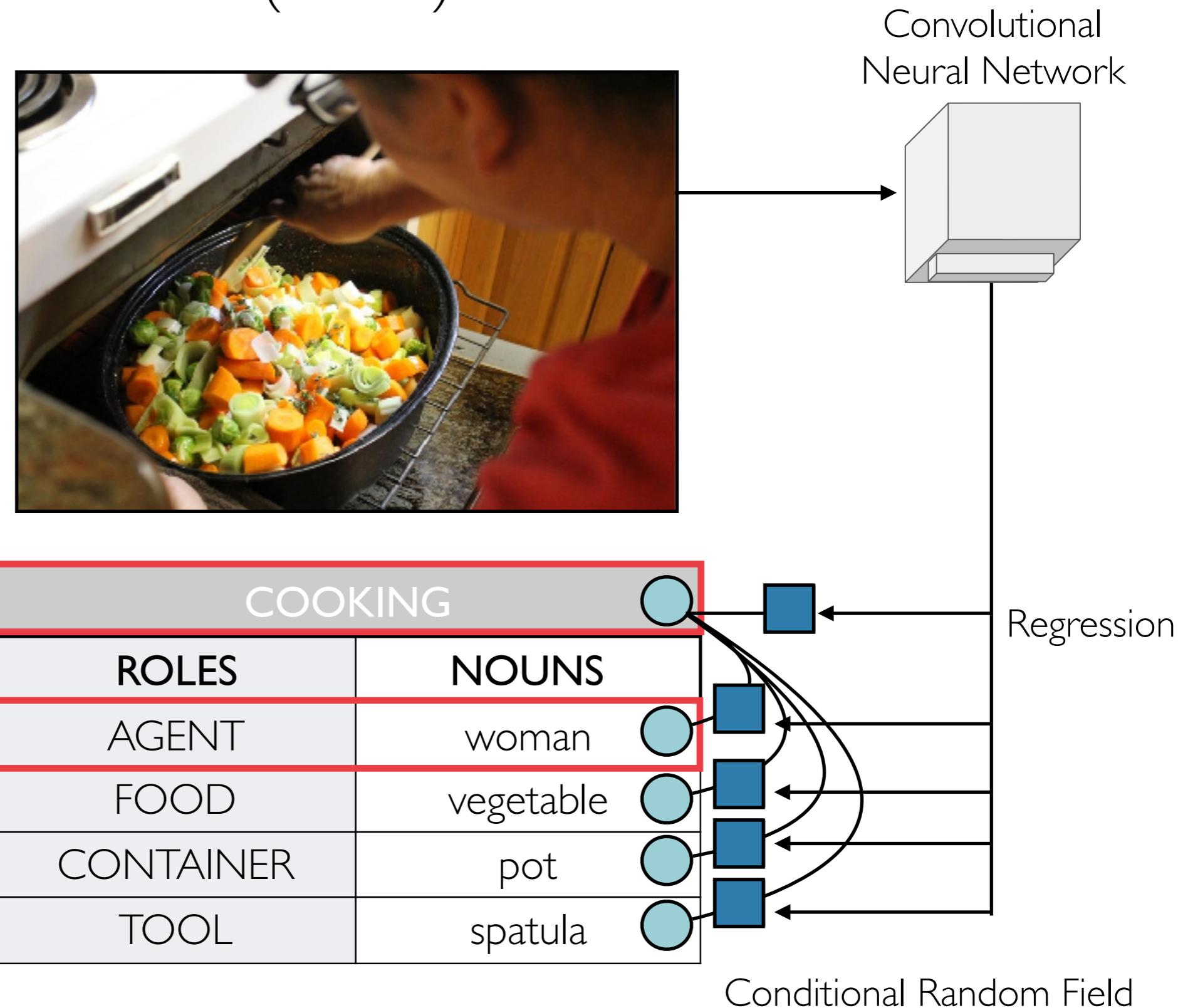
imSitu Visual Semantic Role Labeling (vSRL) (events)



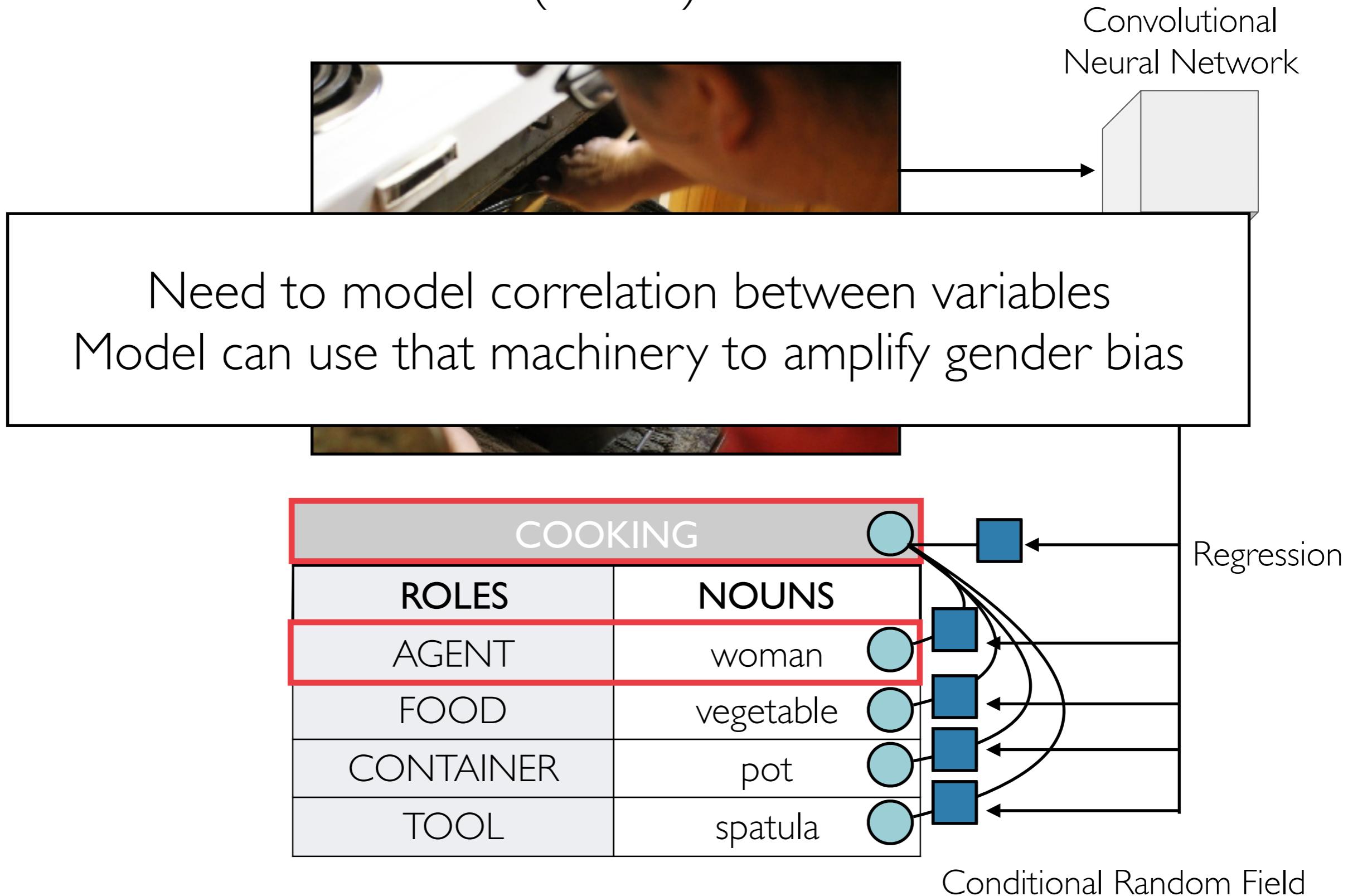
imSitu Visual Semantic Role Labeling (vSRL) (events)



imSitu Visual Semantic Role Labeling (vSRL) (events)



imSitu Visual Semantic Role Labeling (vSRL) (events)



COCO Multi-Label Classification (MLC) (objects)



← Internet

a woman is smiling in a kitchen near a pizza on a stove

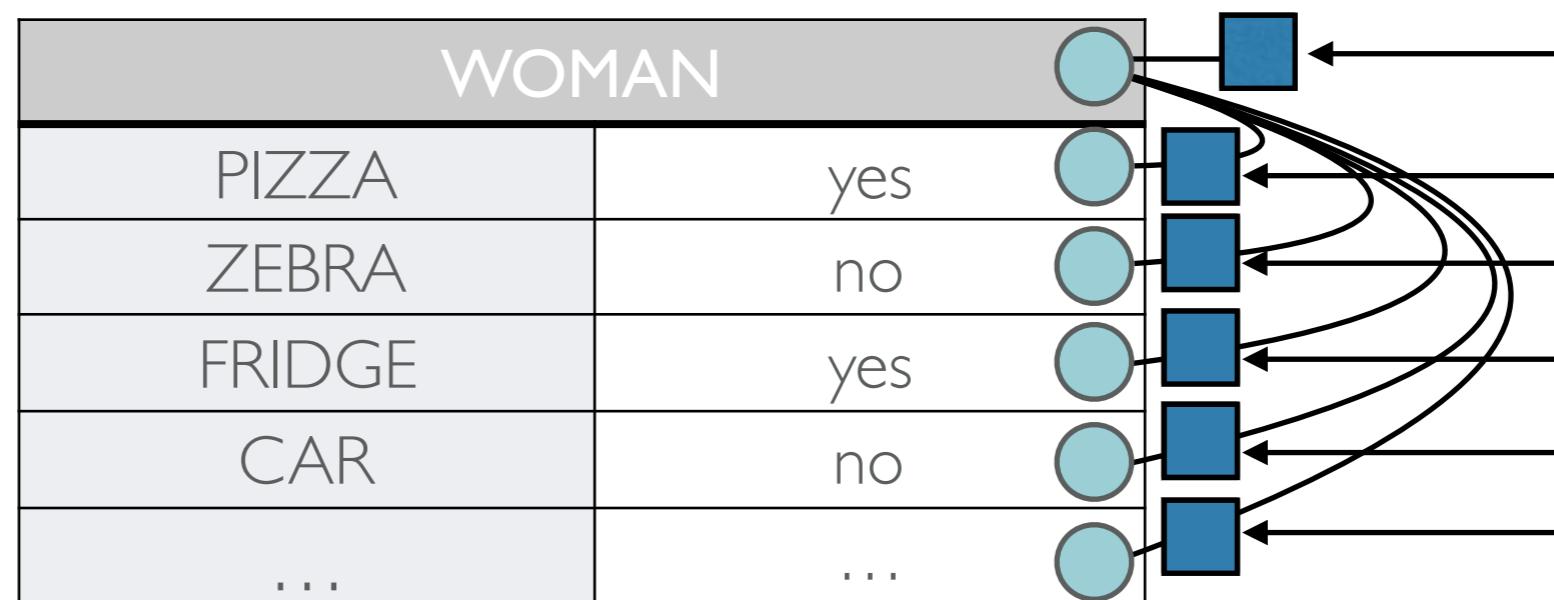
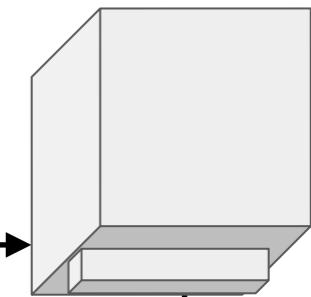
COCO
Objects →

WOMAN	
PIZZA	yes
ZEBRA	no
FRIDGE	yes
CAR	no
...	...

← Caption Inferred
Label

COCO Multi-Label Classification (MLC) (objects)

Convolutional
Neural Network



Conditional Random Field

Related Work

- Implicit Bias
 - image search (Kay et al., 2015)
 - search advertising (Sweeny, 2013)
 - online news (Ross and Carter, 2011)
 - credit score (Hardt et al., 2016)
 - word vector (Bolukbasi et al., 2016)
- Classifier class imbalance
 - Barocas and Selbst, 2014; Dwork et al., 2012;
 - Feldman et al., 2015; Zliobaite, 2015

Outline

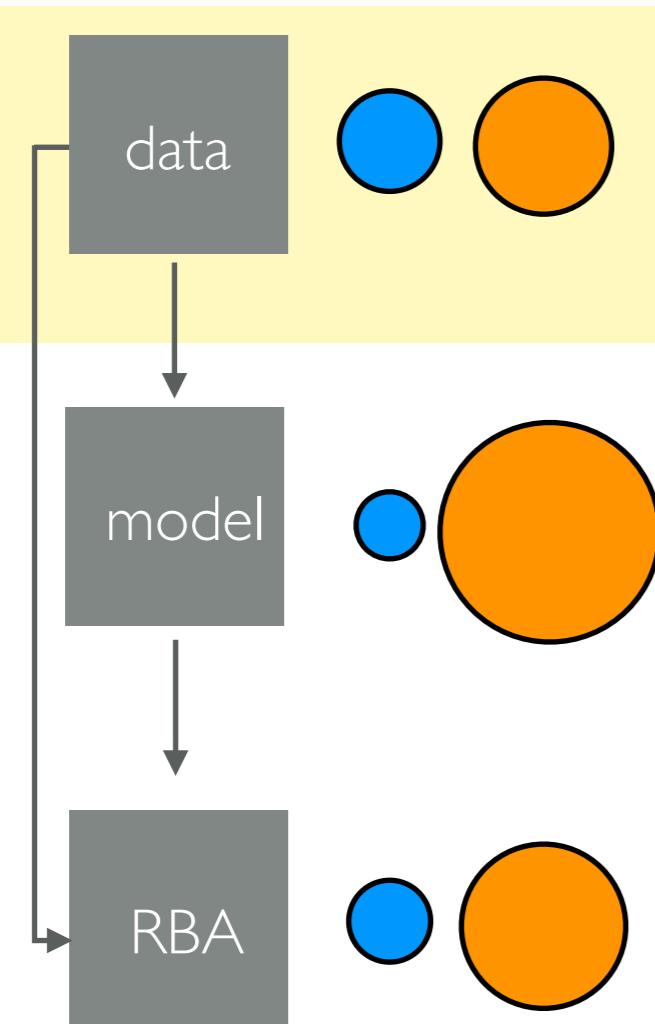


imSitu vSRL
(events)

COCO MLC
(objects)

I. Background

2. Dataset Bias



3. Model Bias Amplification

4. Reducing Bias Amplification

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

$$\frac{\#(\text{◆ cooking}, \text{● man})}{\#(\text{◆ cooking}, \text{● man}) + \#(\text{◆ cooking}, \text{○ woman})} = 1/3$$

Defining Dataset Bias (objects)

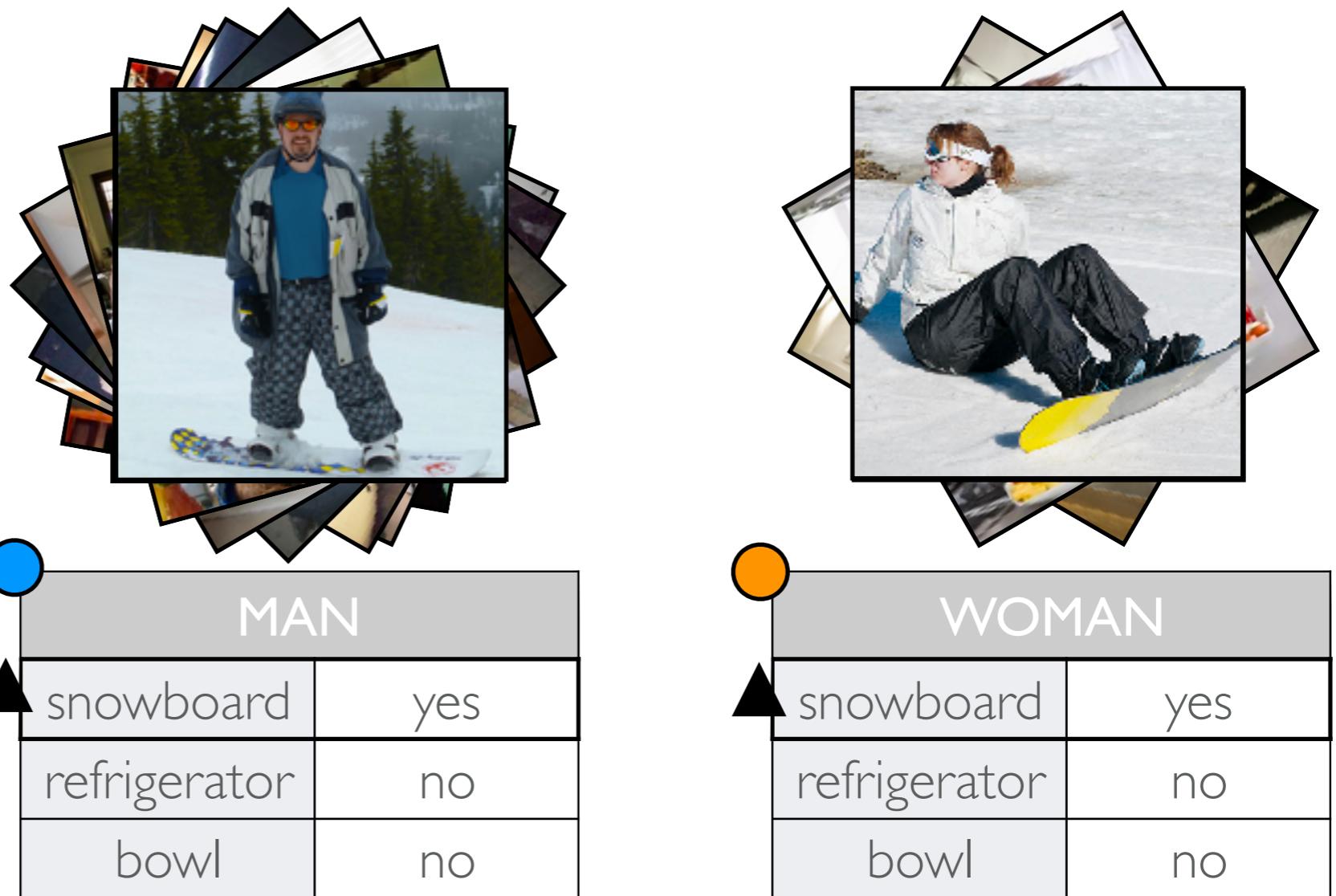
Training Gender Ratio (\blacktriangle noun)

Training Set

\blacktriangle snowboard

\bullet woman

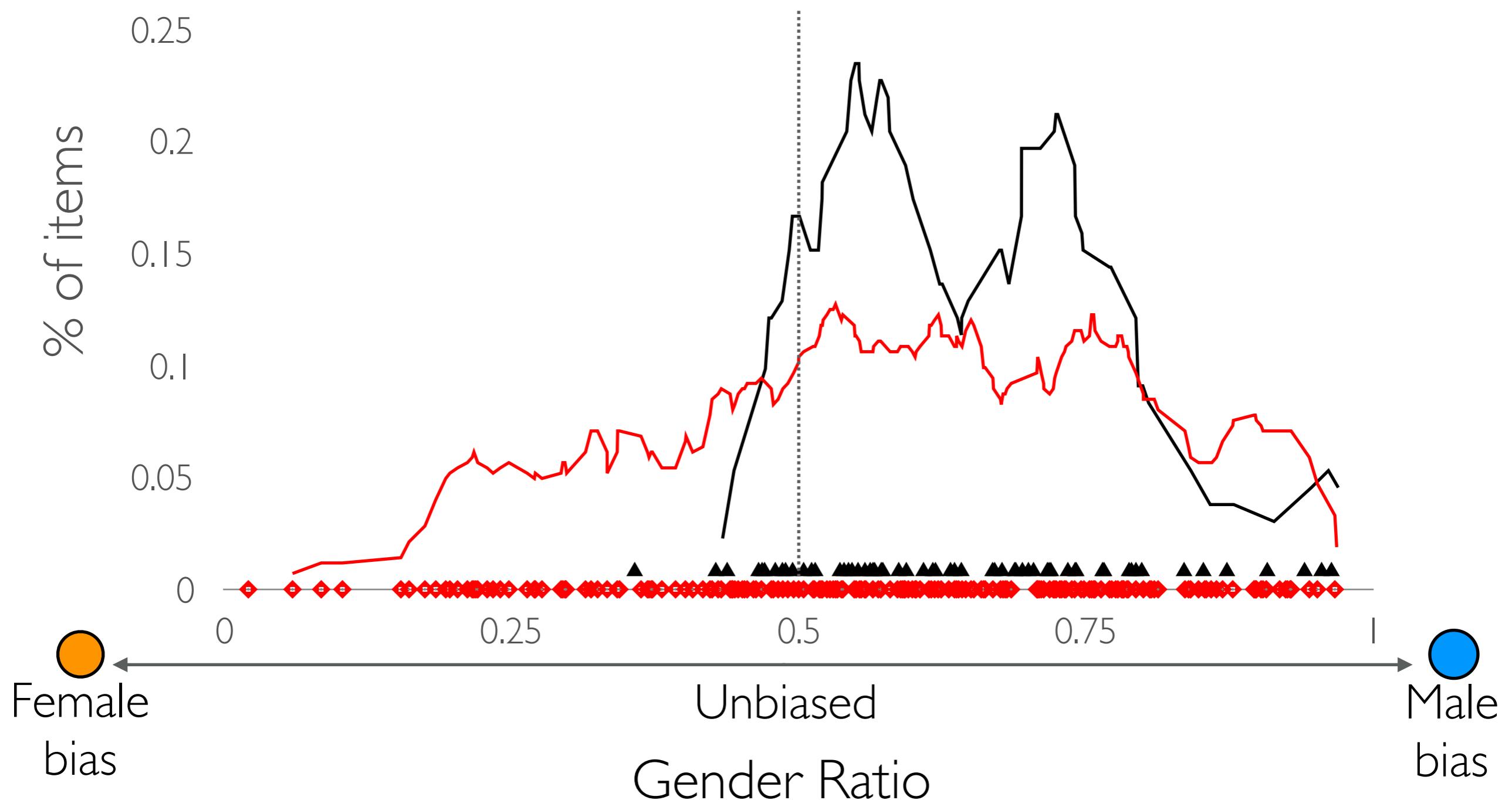
\circ man



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

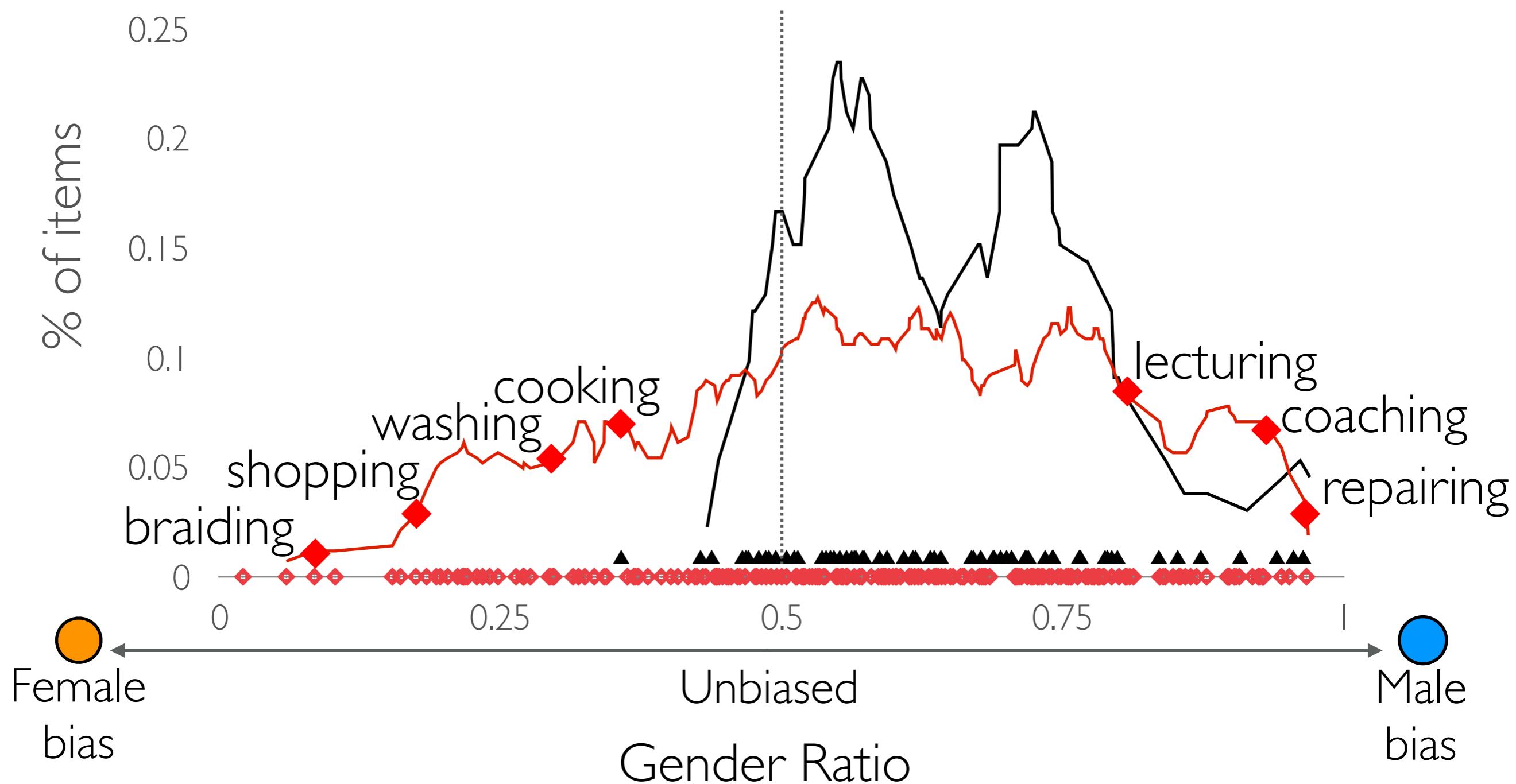
Gender Dataset Bias

- ◆ imSitu Verb
- ▲ COCO Noun



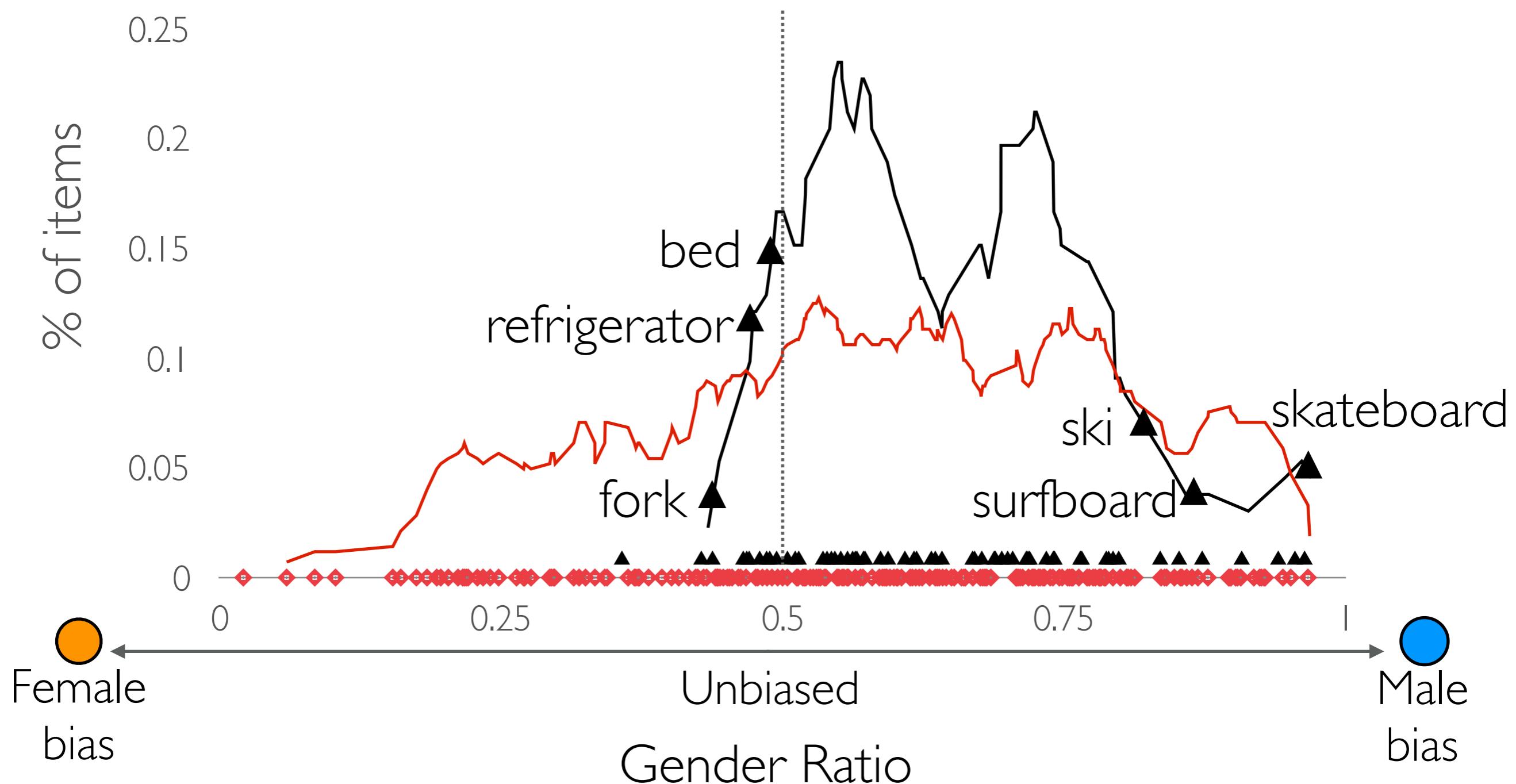
Gender Dataset Bias

- ◆ imSitu Verb
- ▲ COCO Noun



Gender Dataset Bias

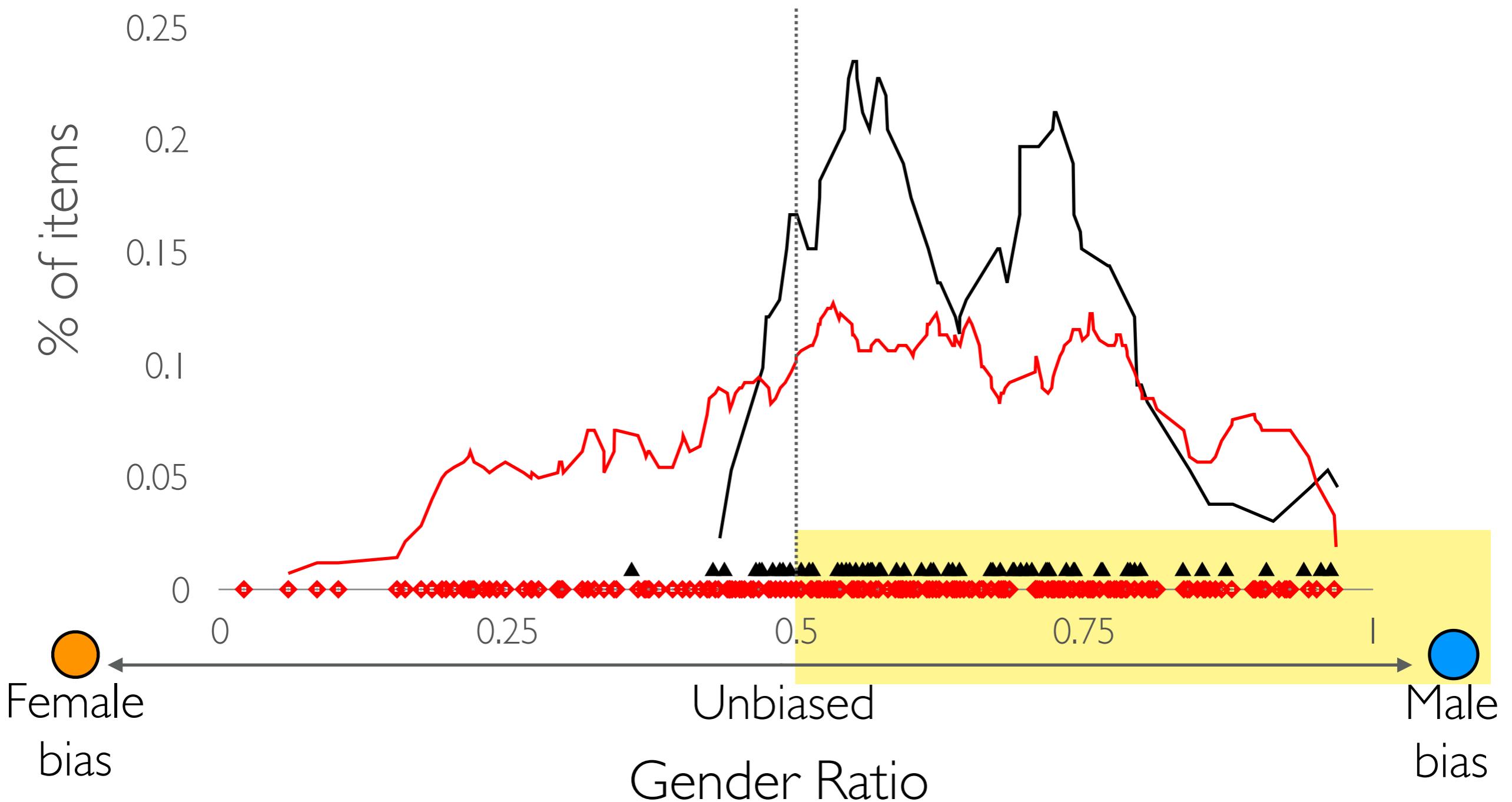
- ◆ imSitu Verb
- ▲ COCO Noun



Gender Dataset Bias

- ◆ imSitu Verb
- ▲ COCO Noun

64.6% bias
86.6% bias



Gender Dataset Bias

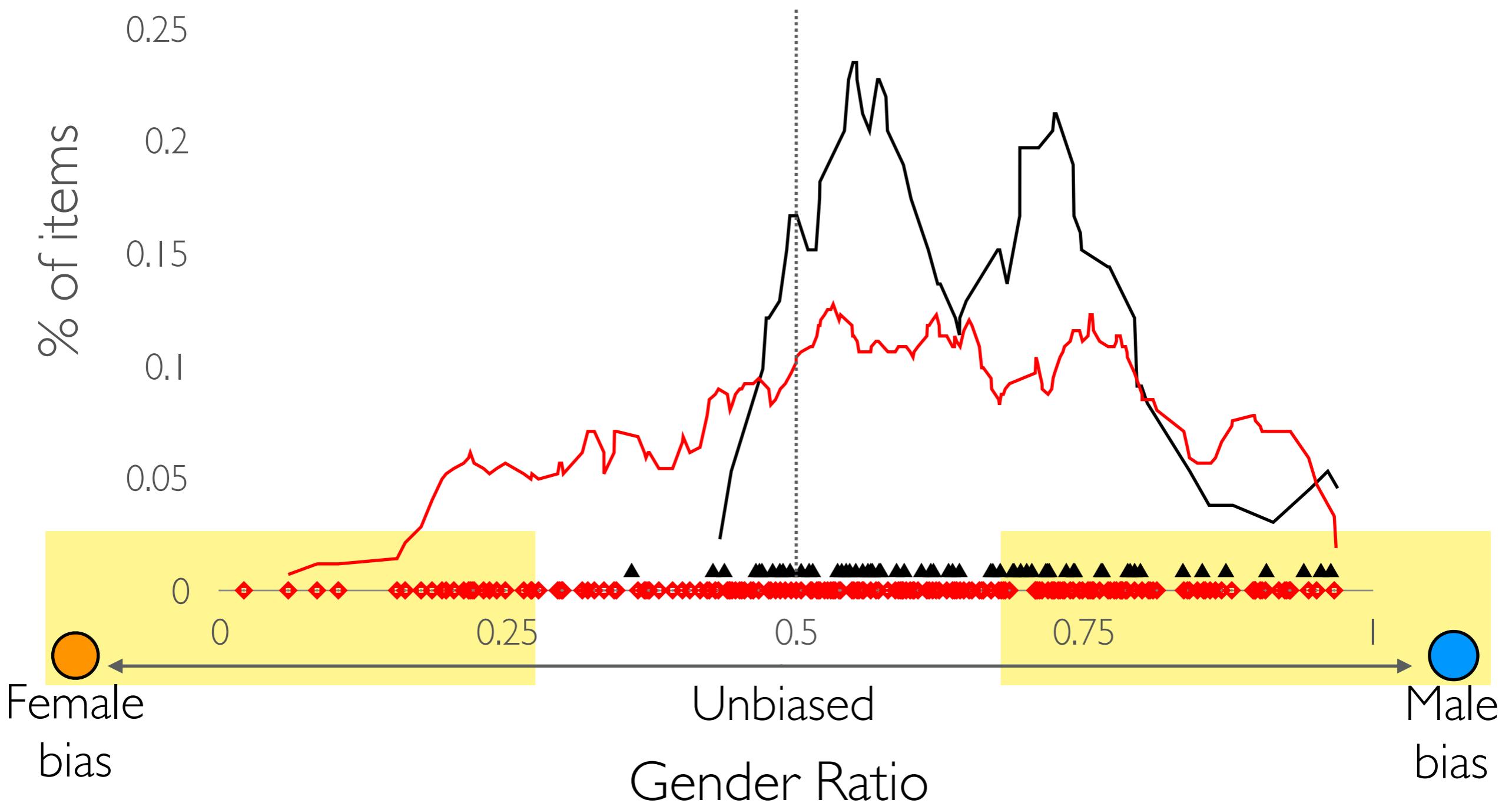
- ◆ imSitu Verb
- ▲ COCO Noun

64.6% bias

86.6% bias

46.9% strong bias ($>2:1$)

37.9% strong bias ($>2:1$)



Outline



imSitu vSRL
(events)

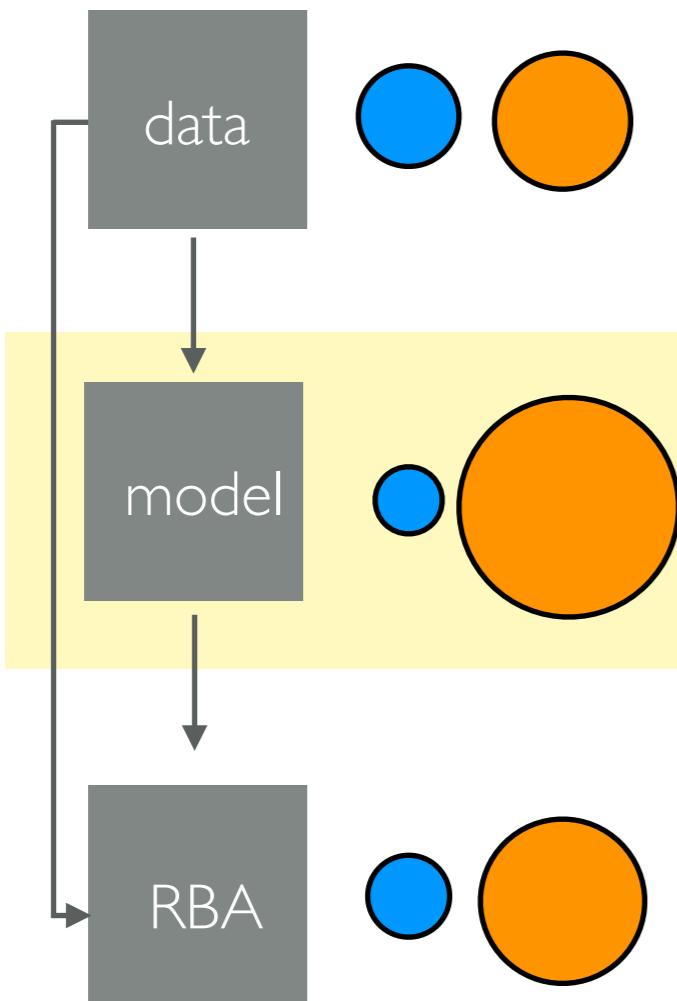
COCO MLC
(objects)

I. Background

2. Dataset Bias

3. Bias Amplification

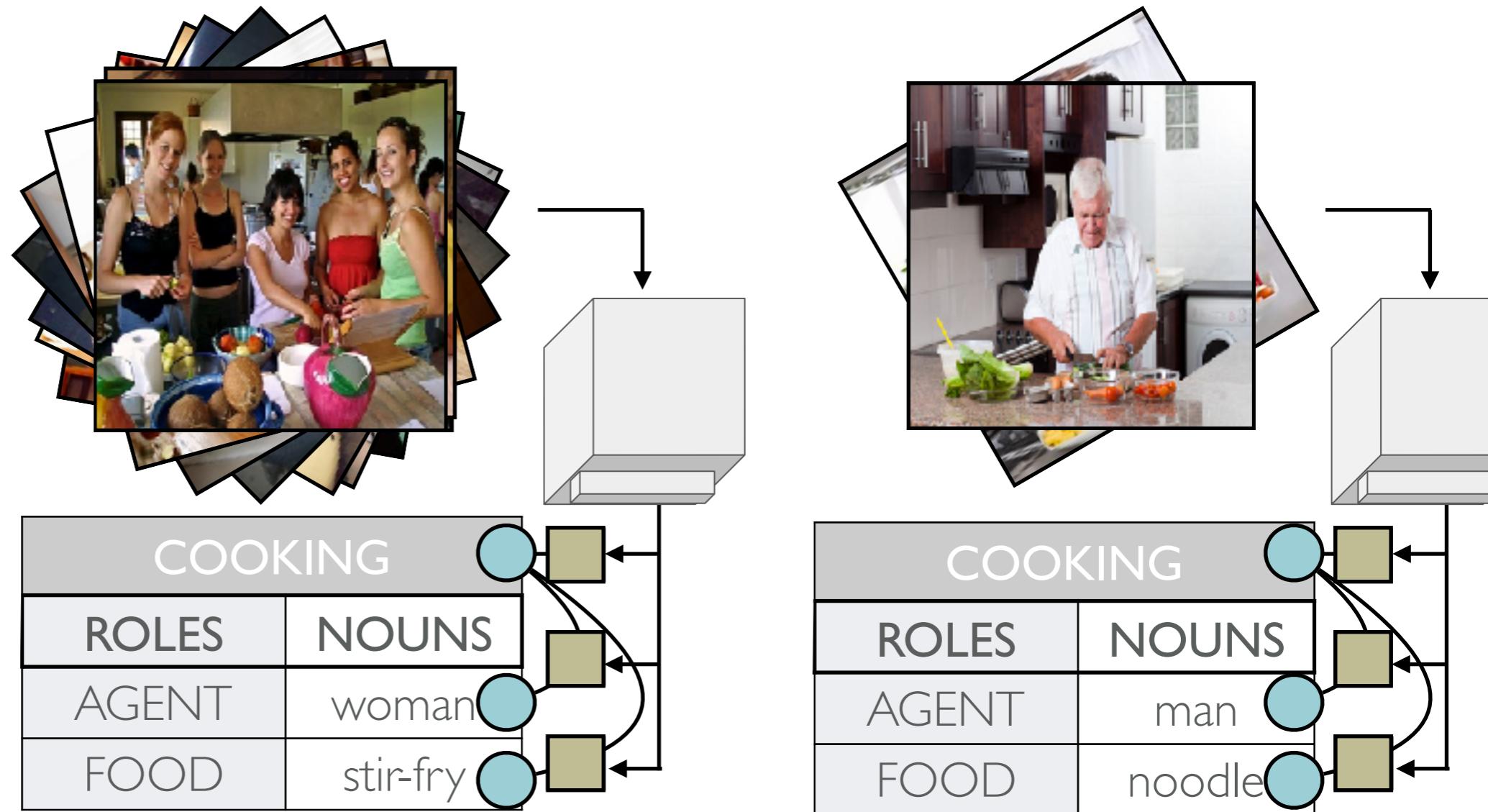
4. Reducing Bias Amplification



Defining Bias Amplification (events)

Predicted Gender Ratio (◆ verb)

Development Set



What does the model predict on unseen data?

Defining Bias Amplification (events)

Predicted Gender Ratio (\diamond verb)

Development Set

- \diamond cooking
- \circ woman
- \bullet man



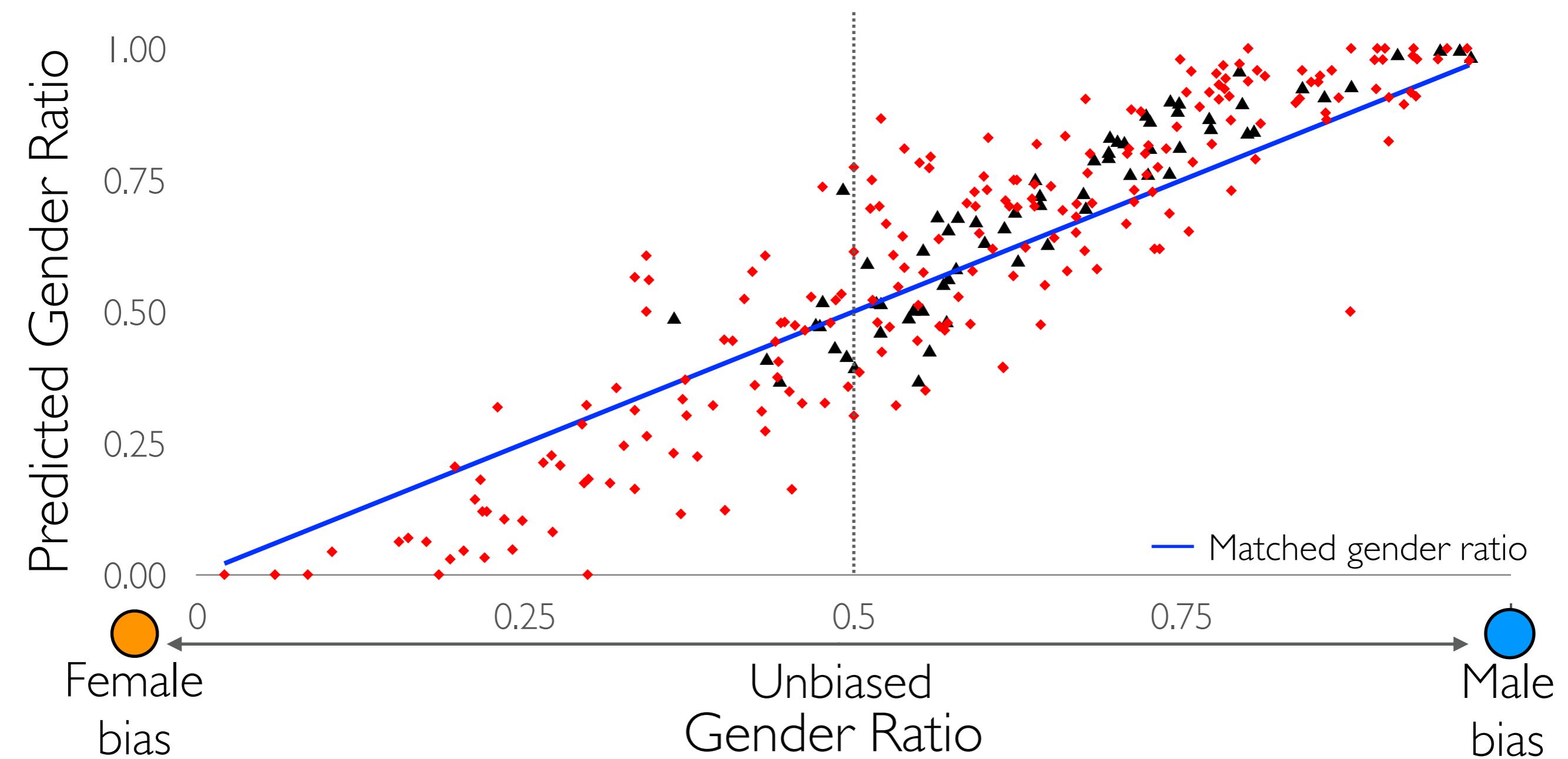
COOKING	
ROLES	NOUNS
AGENT	woman
FOOD	stir-fry

COOKING	
ROLES	NOUNS
AGENT	man
FOOD	noodle

$$\frac{\#(\diamond \text{ cooking}, \bullet \text{ man})}{\#(\diamond \text{ cooking}, \bullet \text{ man}) + \#(\diamond \text{ cooking}, \circ \text{ woman})} = 1/6$$

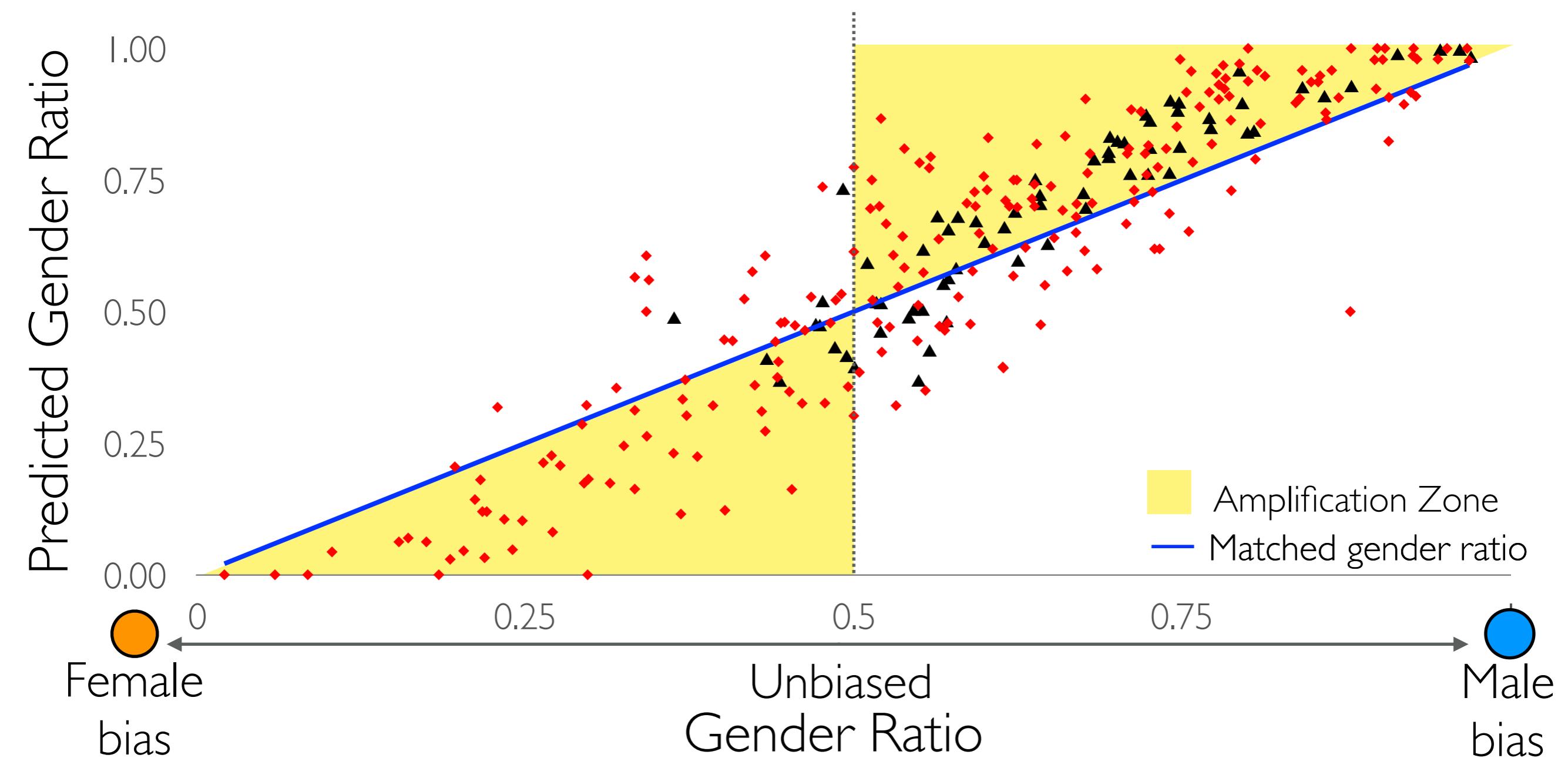
Model Bias Amplification

◆ imSitu Verb
▲ COCO Noun

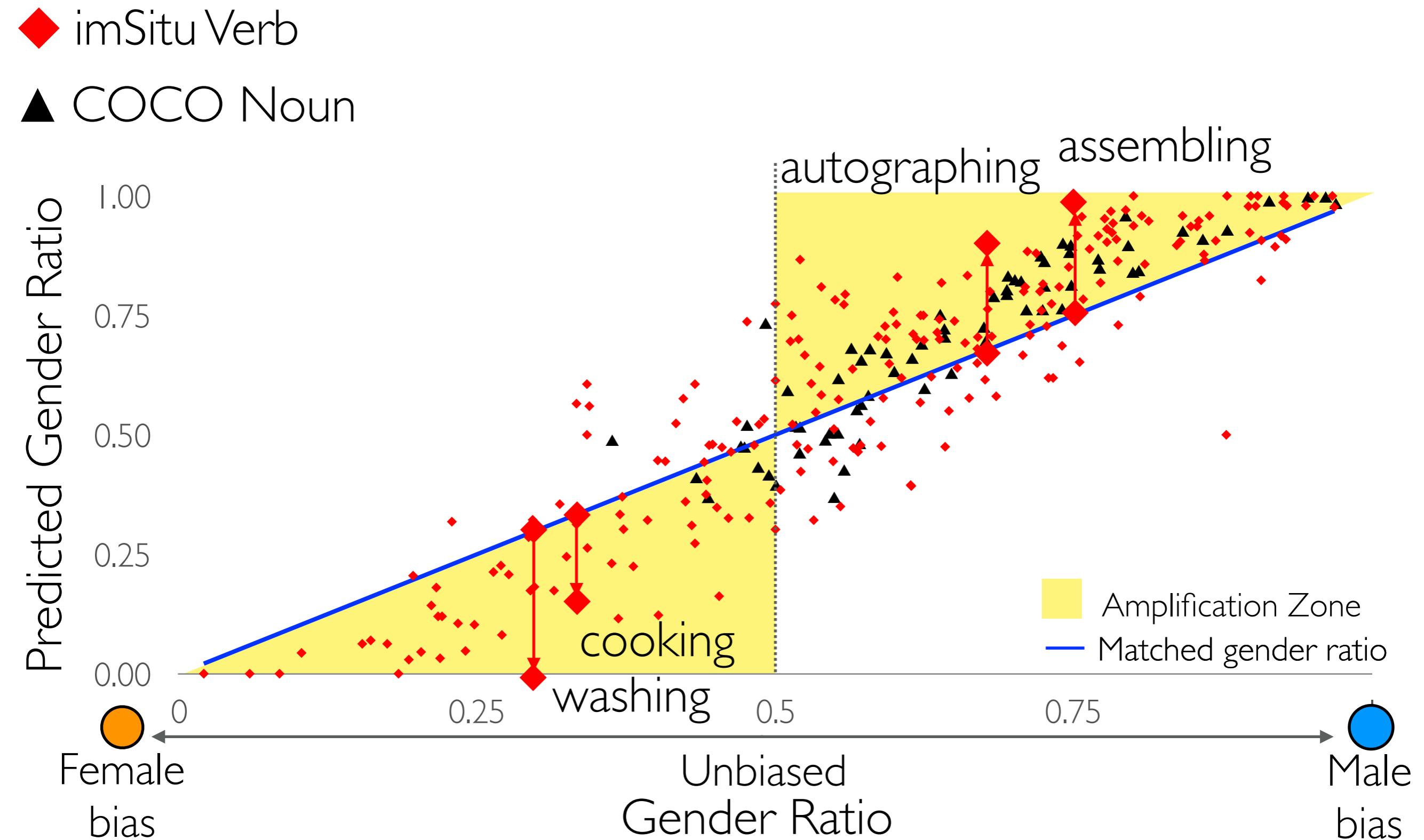


Model Bias Amplification

- ◆ imSitu Verb
- ▲ COCO Noun



Model Bias Amplification

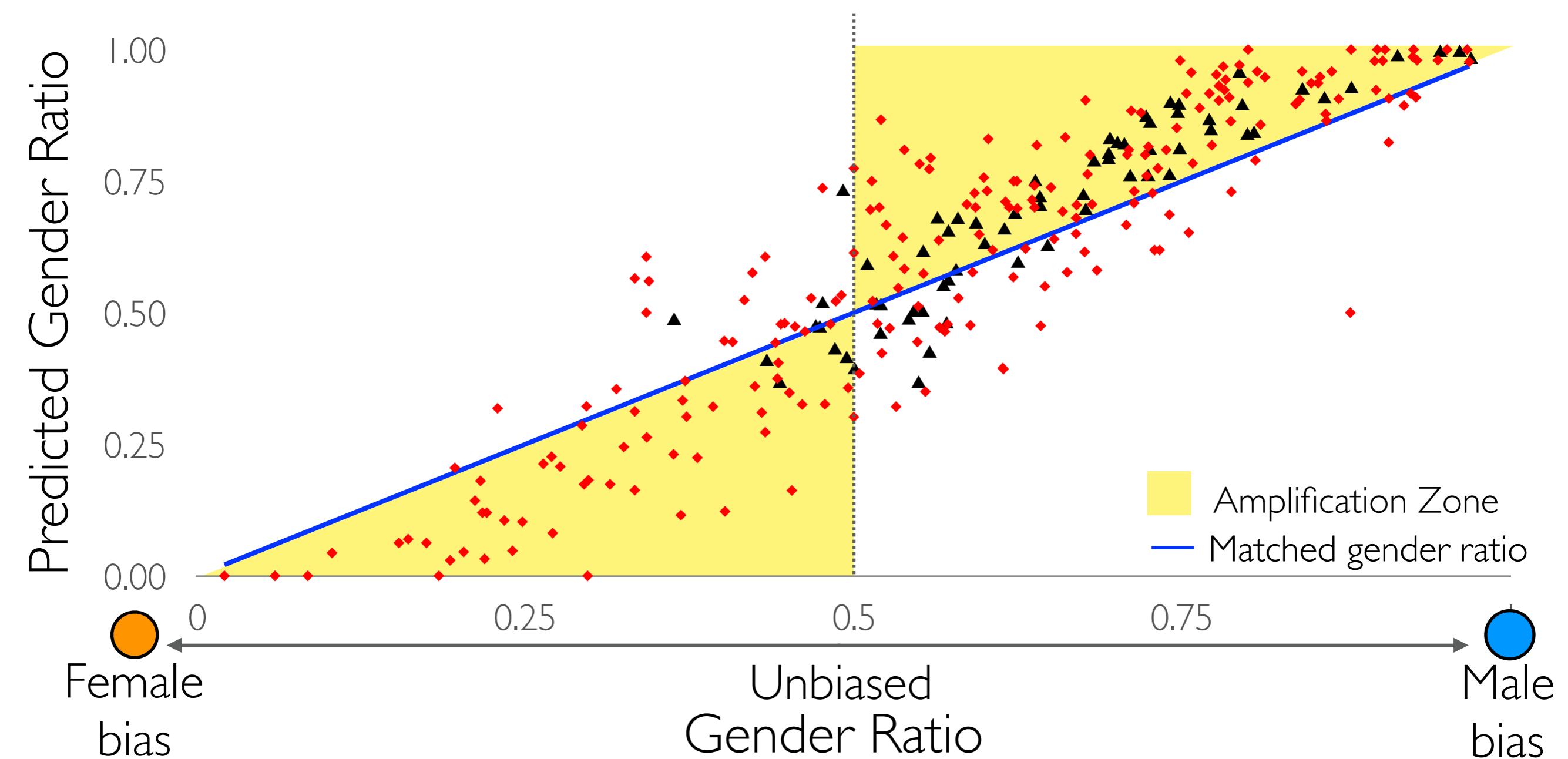


Model Bias Amplification

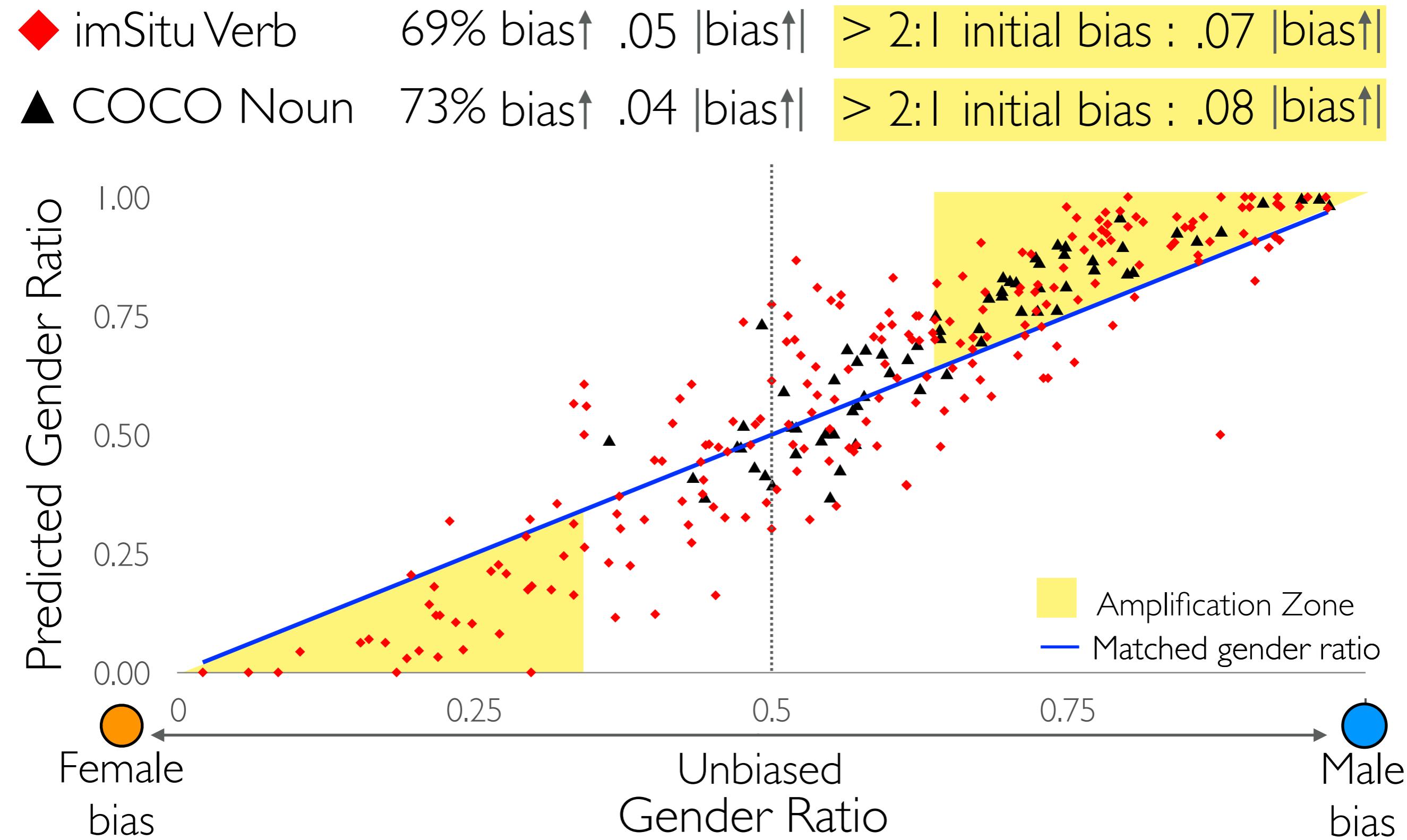
◆ imSitu Verb
▲ COCO Noun

69% bias↑ .05 |bias↑

73% bias↑ .04 |bias↑

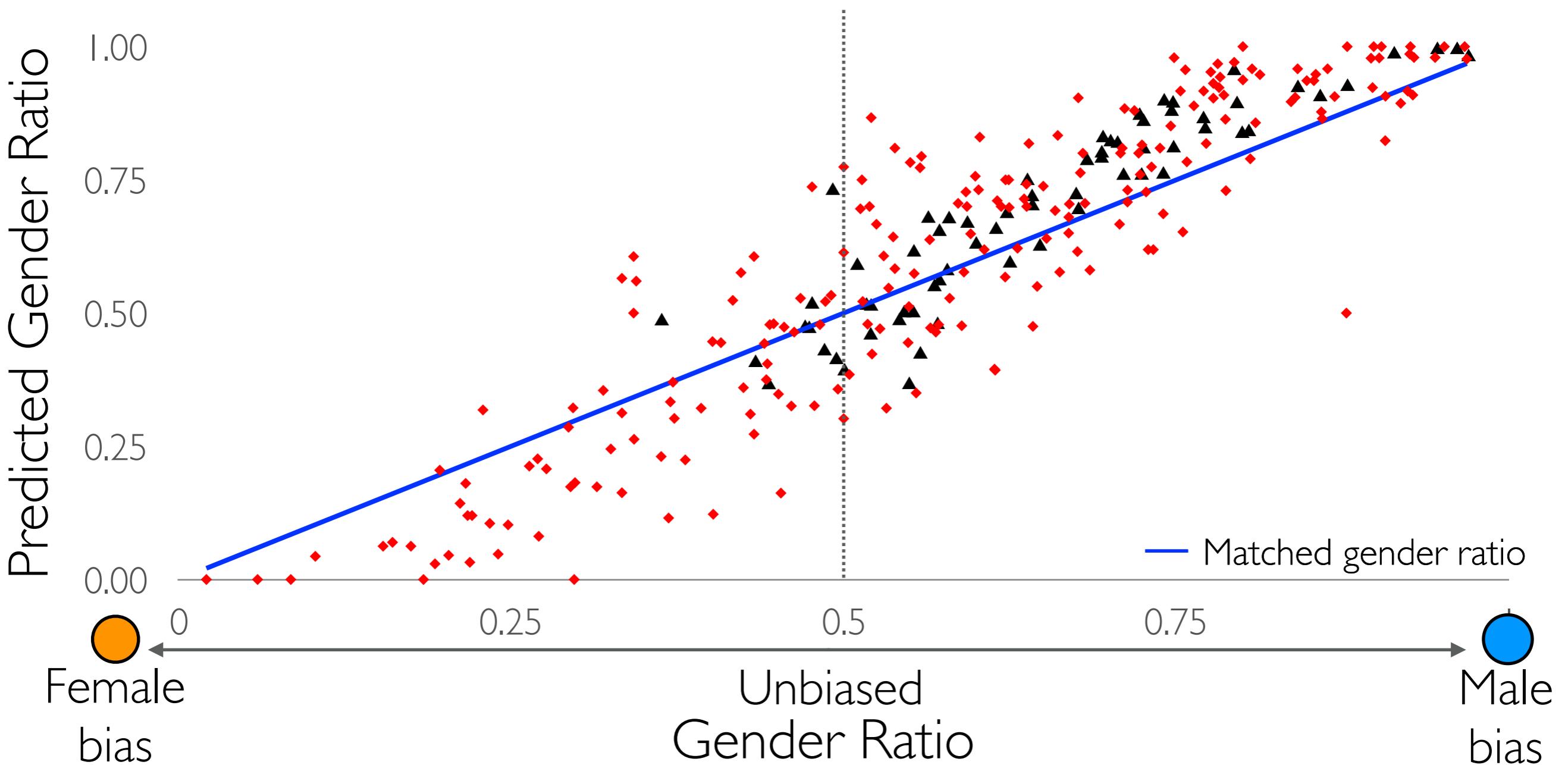


Model Bias Amplification



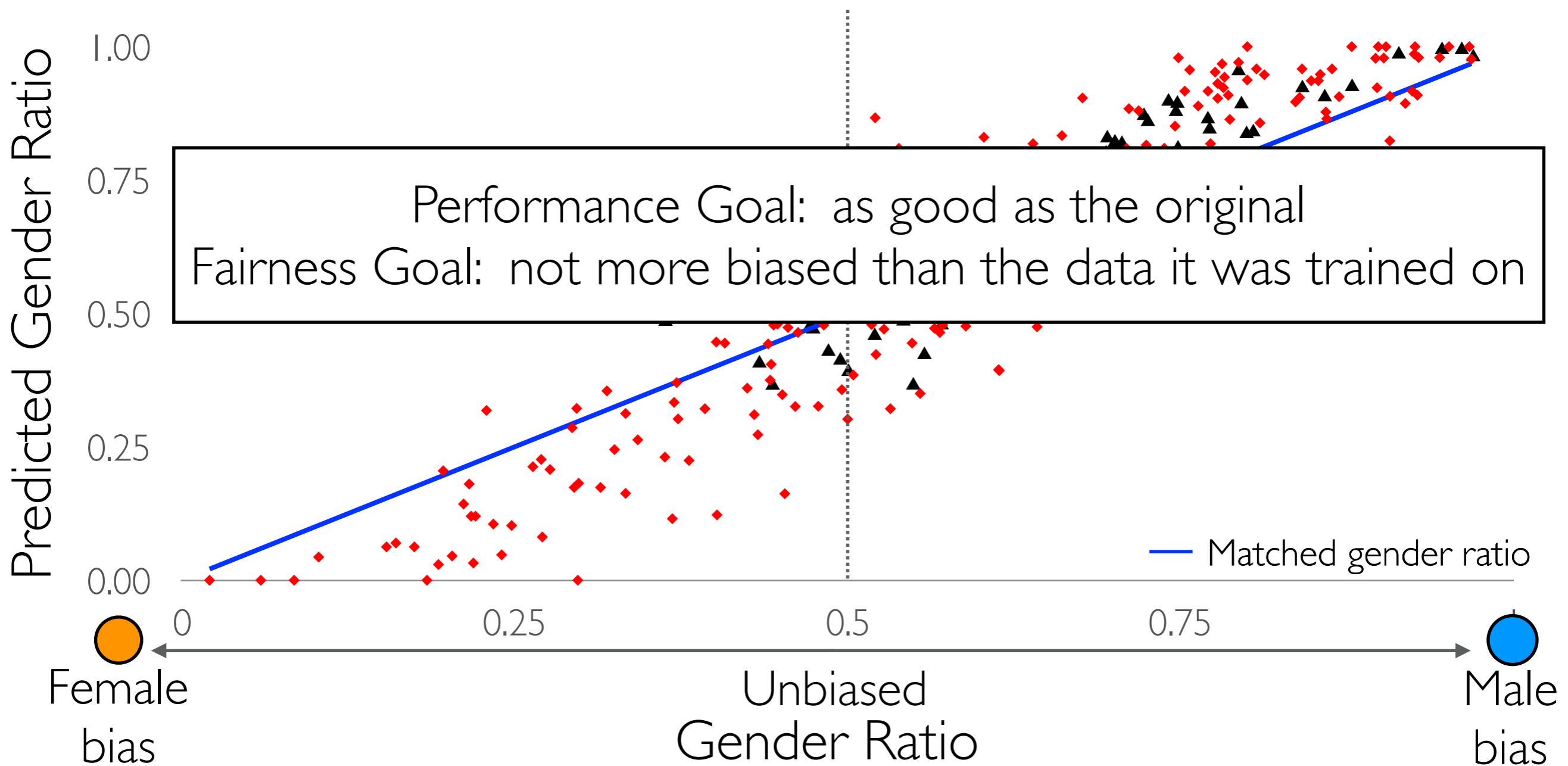
Summary

Can we remove gender bias amplification and still maintain performance?



Summary

Can we remove gender bias amplification and still maintain performance?



Outline



imSitu vSRL
(events)

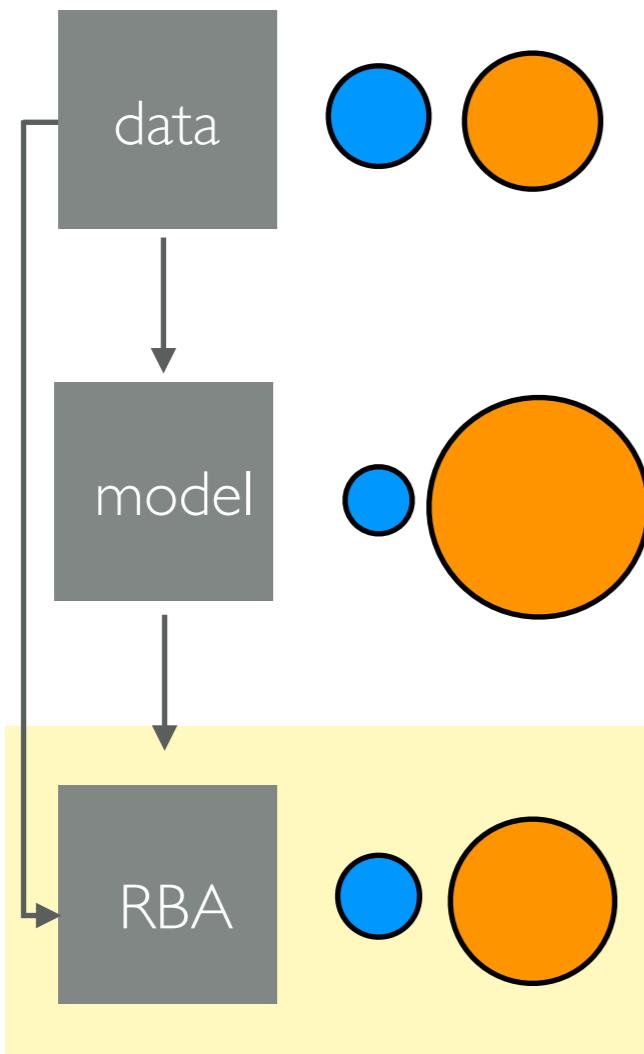
COCO MLC
(objects)

I. Background

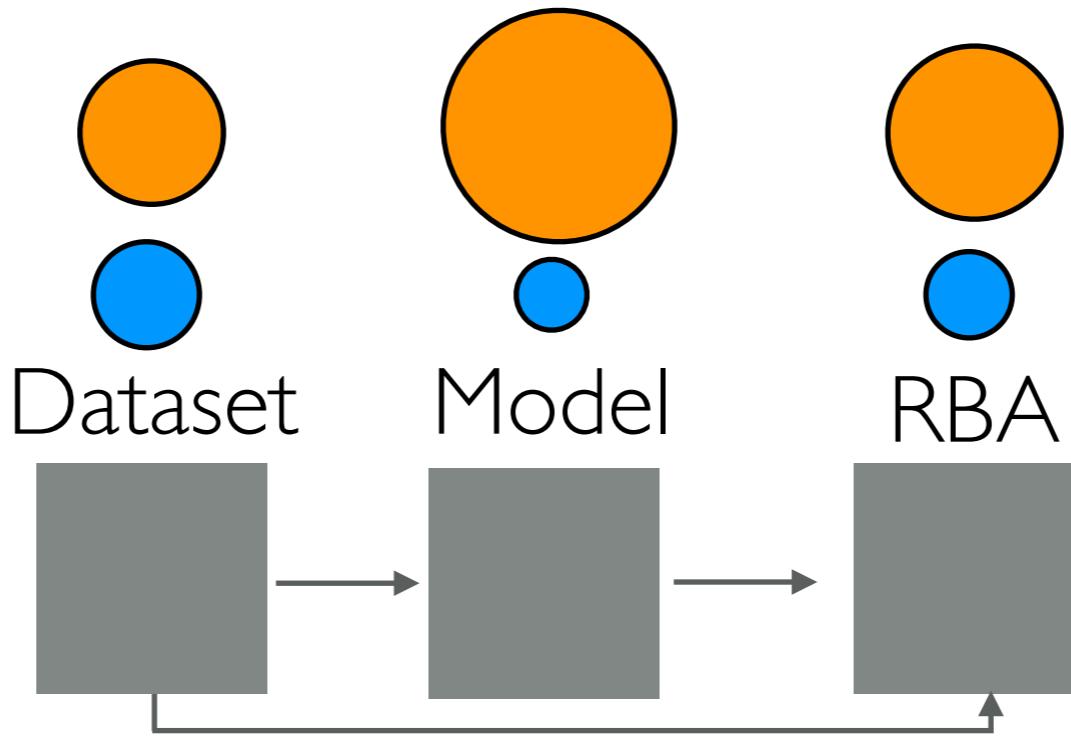
2. Dataset Bias

3. Bias Amplification

4. Reducing Bias Amplification



Reducing Bias Amplification (RBA)



- Corpus level constraints on model output (ILP)
 - ★ Doesn't require model retraining
- Reuse model inference through Lagrangian relaxation
 - ★ Can be applied to any structured model

Reducing Bias Amplification (RBA)

Integer Linear Program

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

← base model
CRF Inference

Reducing Bias Amplification (RBA)

Integer Linear Program

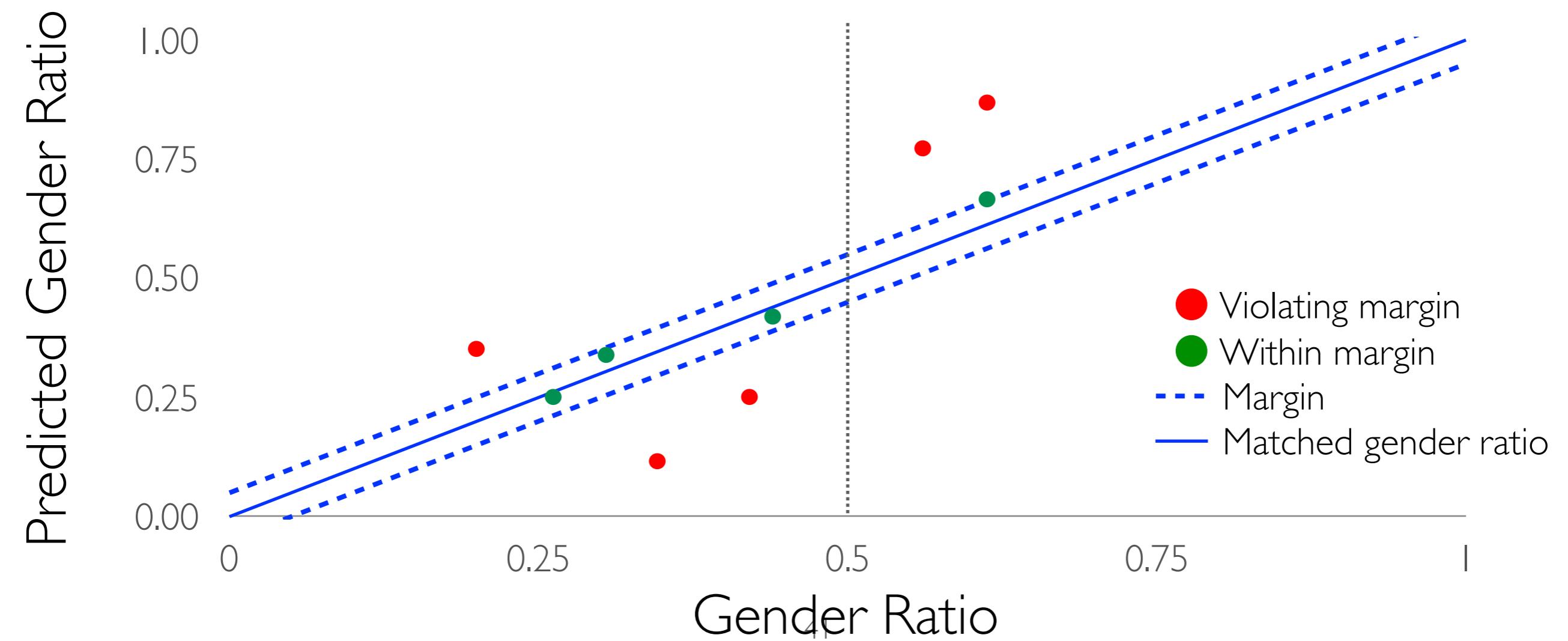
$$\sum_i \max_{y_i} s(y_i, \text{image})$$

\forall points

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

$$f(y_1 \dots y_n)$$

\leq margin



Reducing Bias Amplification (RBA)

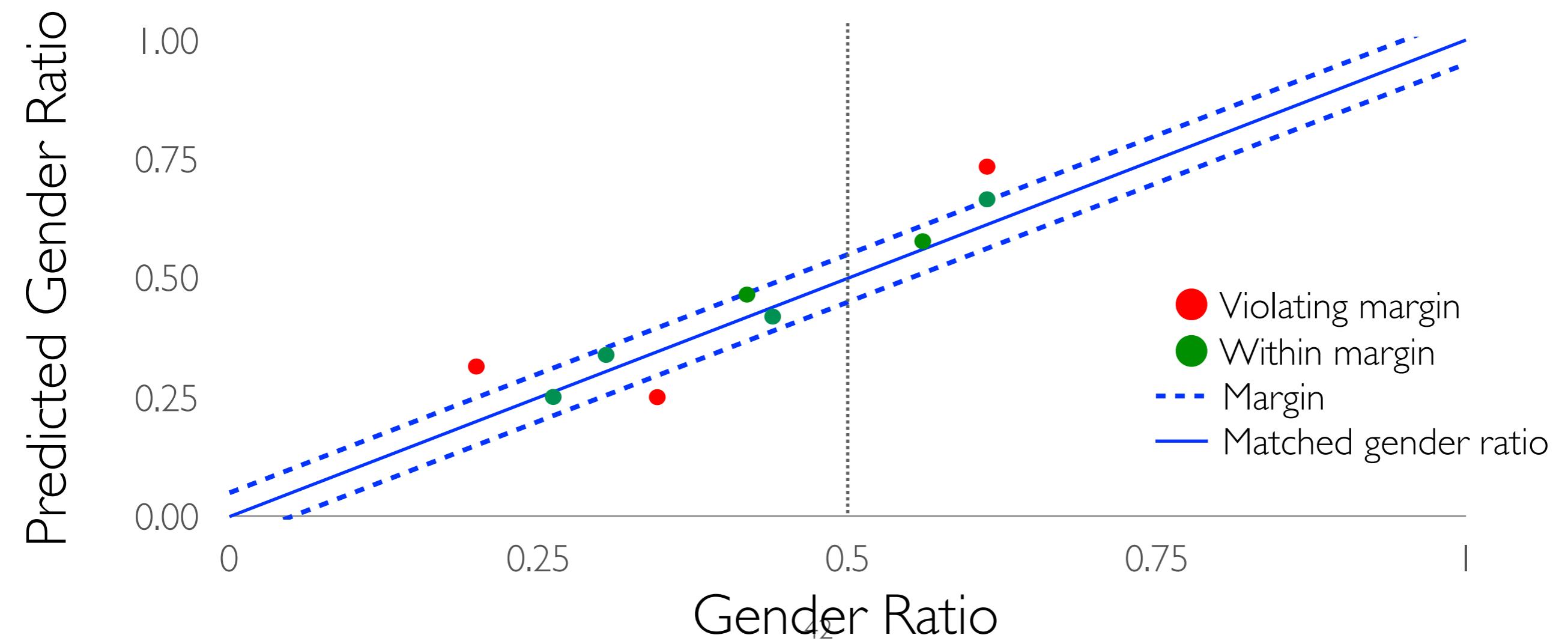
Integer Linear Program

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

\forall points

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

$$f(y_1 \dots y_n)$$



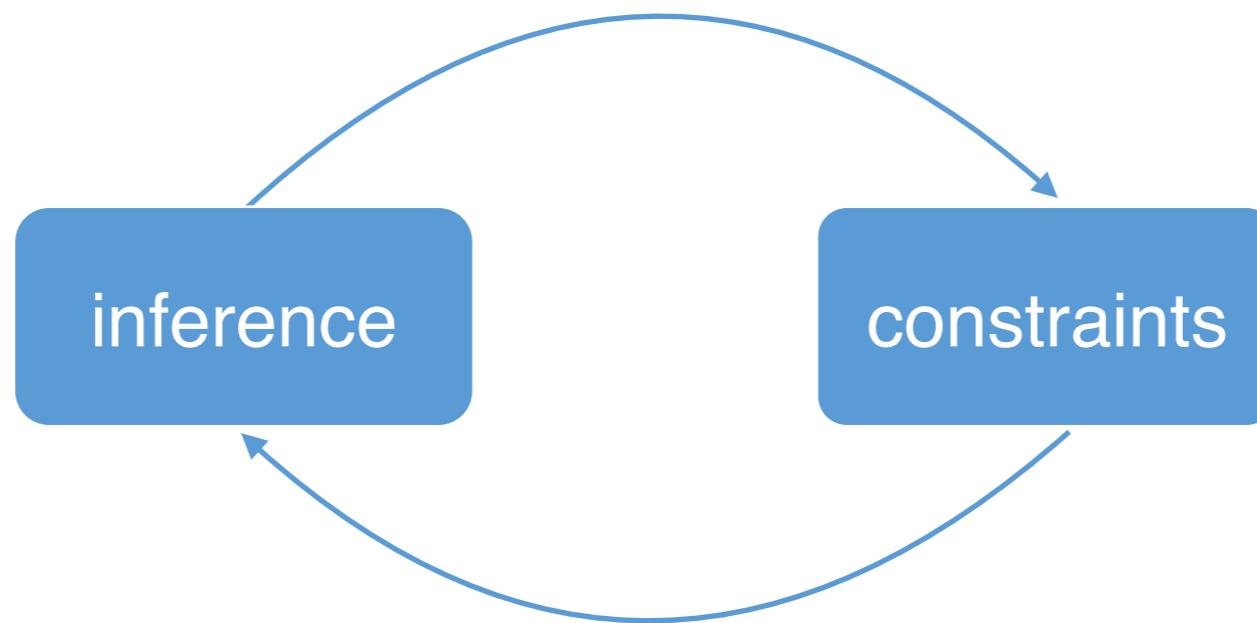
Reducing Bias Amplification (RBA)

Integer Linear Program

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

$$\forall \text{ points } \quad \left| \frac{\text{Training Ratio}}{\text{Predicted Ratio}} - f(y_1 \dots y_n) \right| \leq \text{margin}$$

Lagrangian Relaxation



Lagrangian Relaxation



COOKING	
ROLES	NOUNS
AGENT	woman
FOOD	pancake

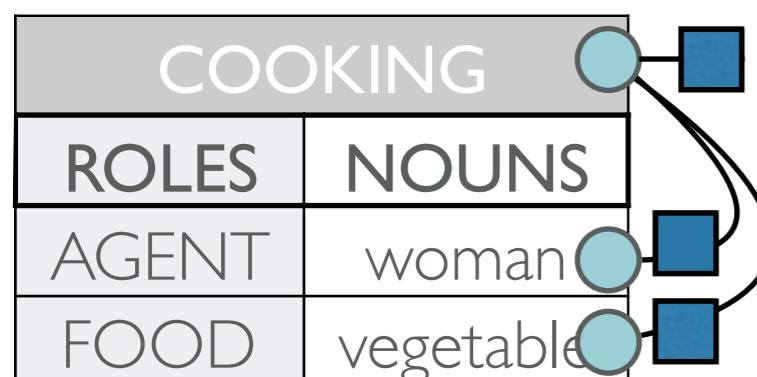
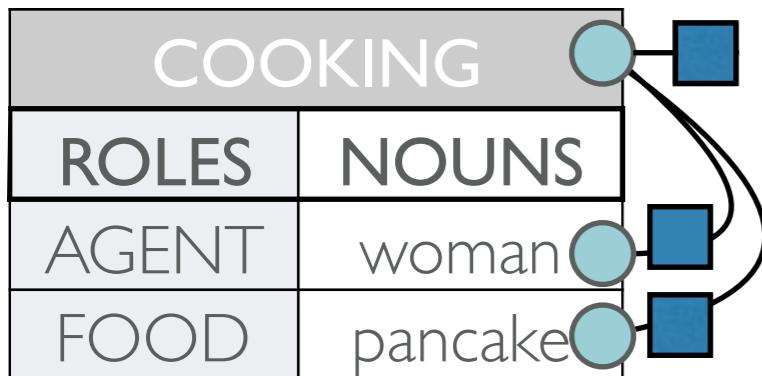


COOKING	
ROLES	NOUNS
AGENT	woman
FOOD	vegetable

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

$$|\text{Training Ratio} - \text{Predicted Ratio}| \leq \text{margin}$$
$$(1/2)$$

Lagrangian Relaxation



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

$$|\text{Training Ratio} - \text{Predicted Ratio}| \leq \text{margin}$$
$$(1/2)$$

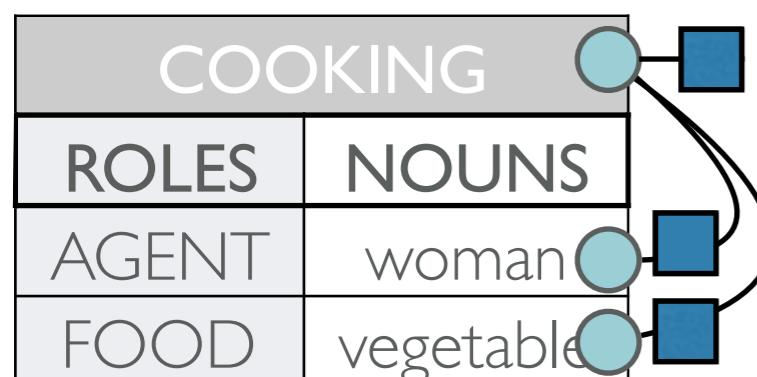
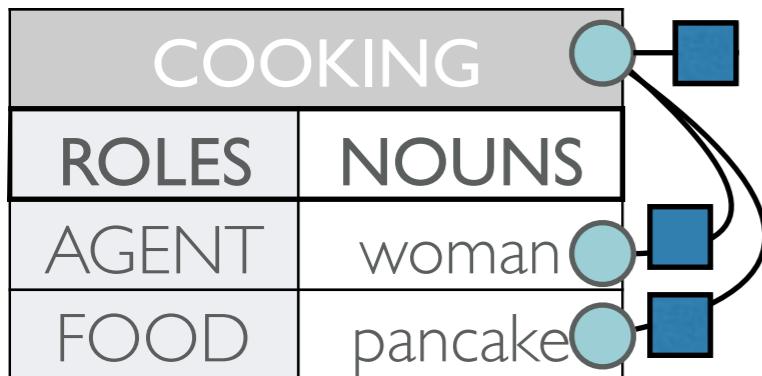
- Lagrange Multiplier (λ) Per Constraint

inference

update λ

update
potentials

Lagrangian Relaxation



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

$$|\text{Training Ratio} - \text{Predicted Ratio}| \leq \text{margin}$$
$$(1/2)$$

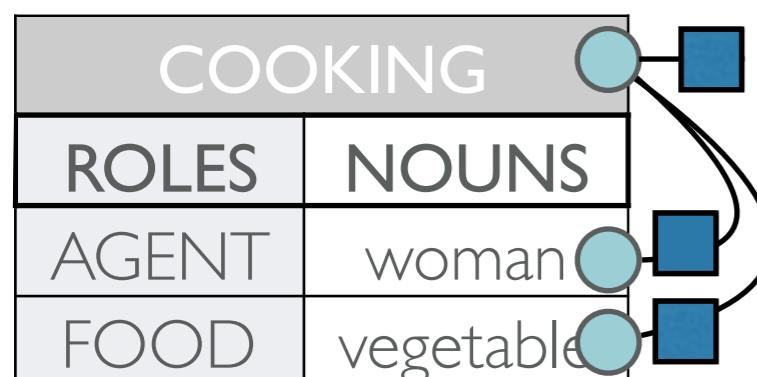
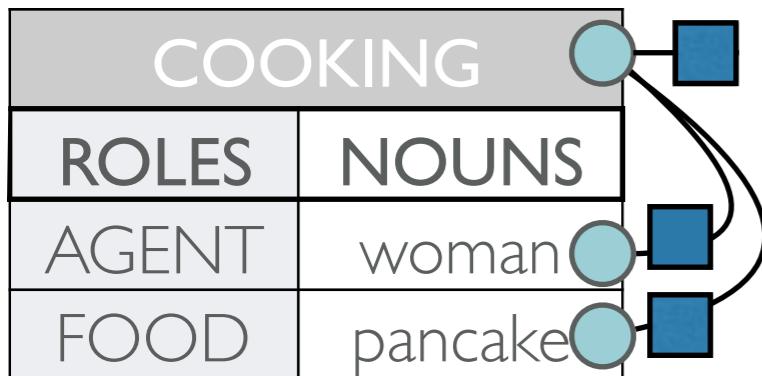
- Lagrange Multiplier (λ) Per Constraint

inference

update λ

update
potentials

Lagrangian Relaxation



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

$$|\text{Training Ratio} - \text{Predicted Ratio}| \leq \text{margin}$$
$$(1/2)$$

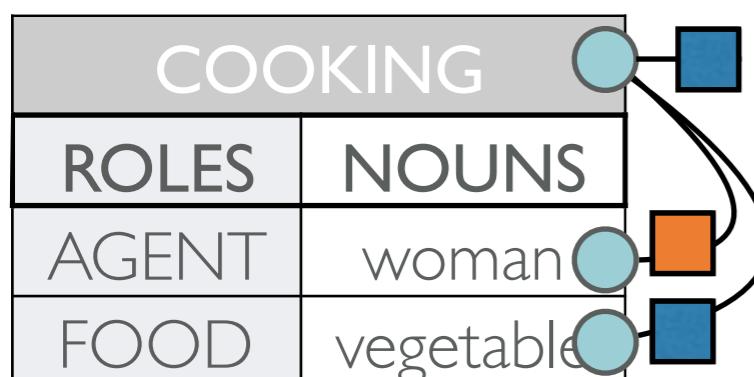
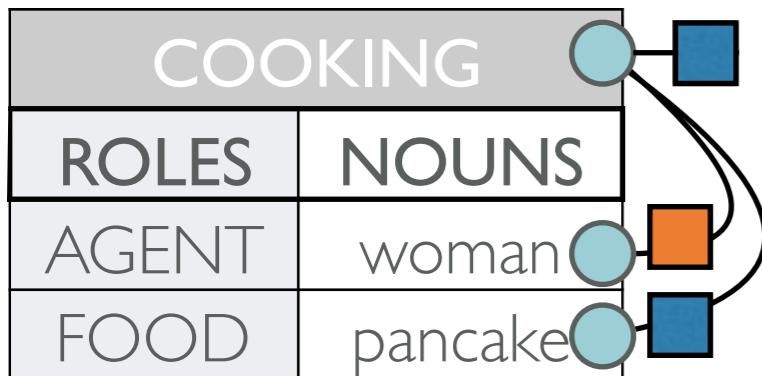
- Lagrange Multiplier (λ) Per Constraint

inference

update λ

update potentials

Lagrangian Relaxation



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

$$|\text{Training Ratio} - \text{Predicted Ratio}| \leq \text{margin}$$
$$(1/2)$$

- Lagrange Multiplier (λ) Per Constraint

inference

update λ

update
potentials

Lagrangian Relaxation



COOKING	
ROLES	NOUNS
AGENT	woman
FOOD	pancake



COOKING	
ROLES	NOUNS
AGENT	man
FOOD	vegetable

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

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

$$(1/2)$$

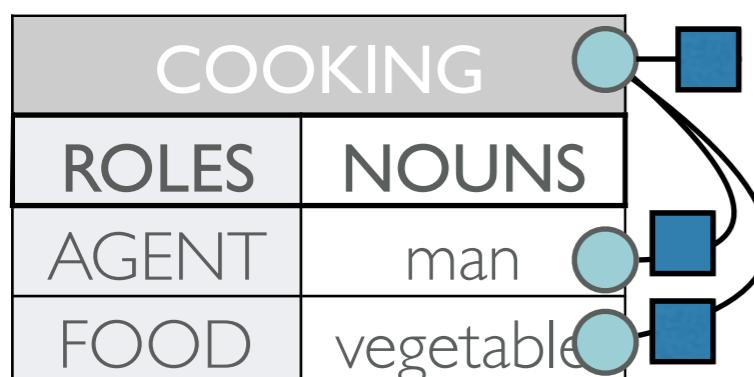
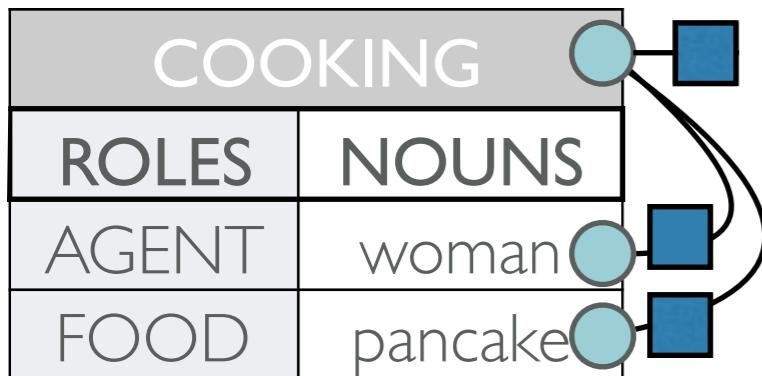
- Lagrange Multiplier (λ) Per Constraint

inference

update λ

update
potentials

Lagrangian Relaxation



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

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

$$(1/2)$$

- Lagrange Multiplier (λ) Per Constraint

inference

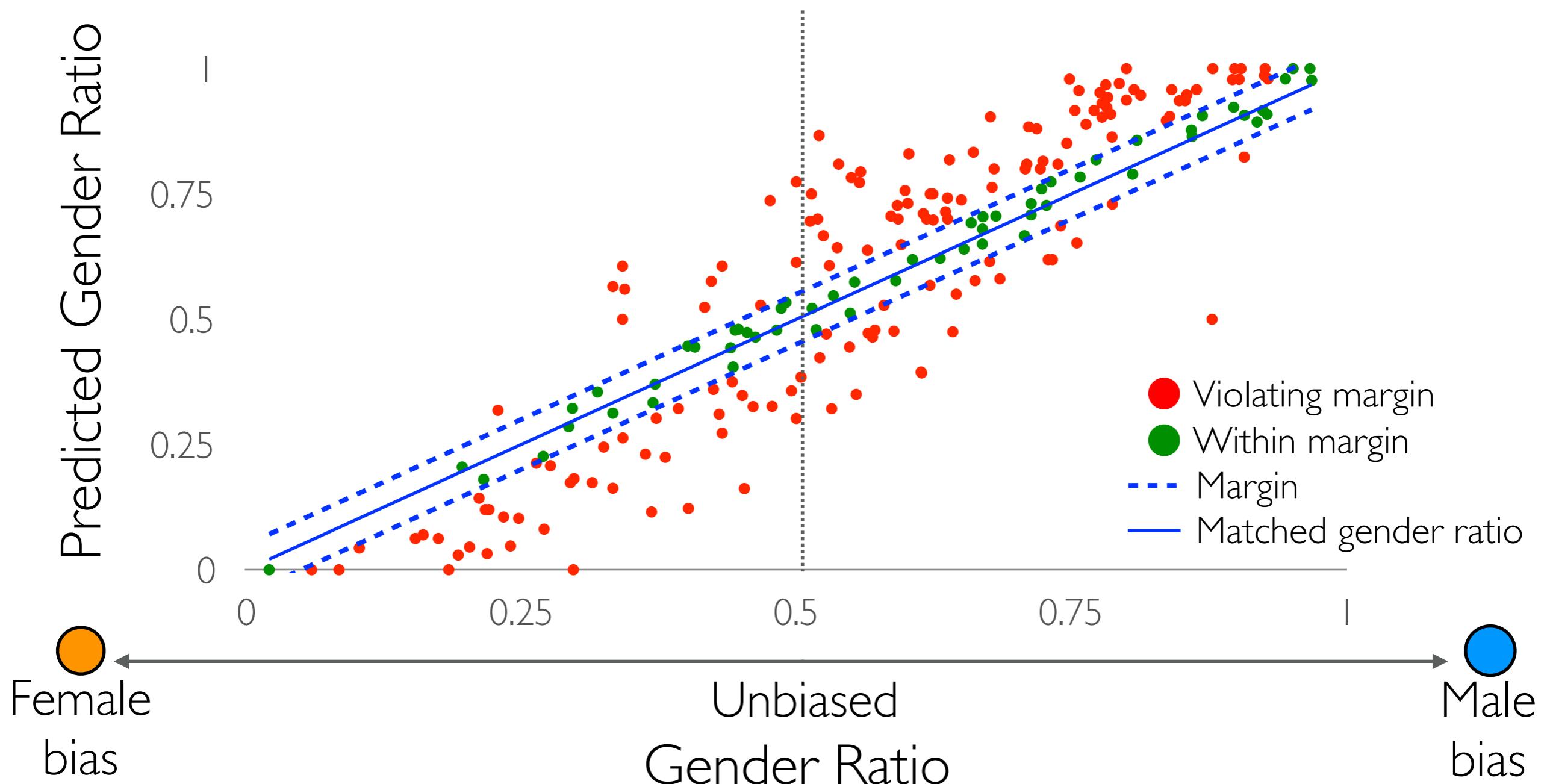
update λ

update potentials



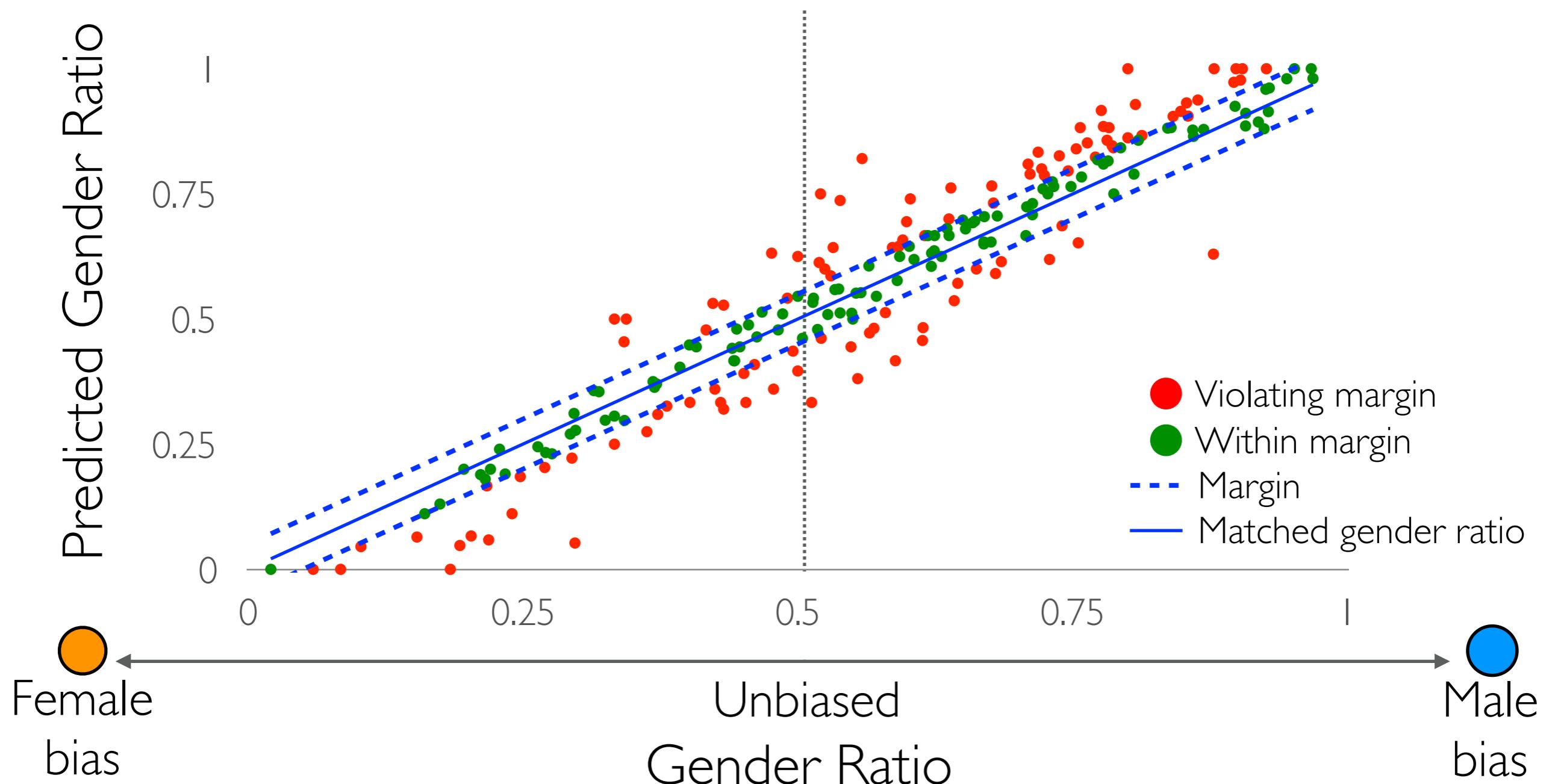
Gender Bias De-amplification in imSitu

imSitu Verb Violation: 72.6% .050 |bias↑| 24.07 acc.



Gender Bias De-amplification in imSitu

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



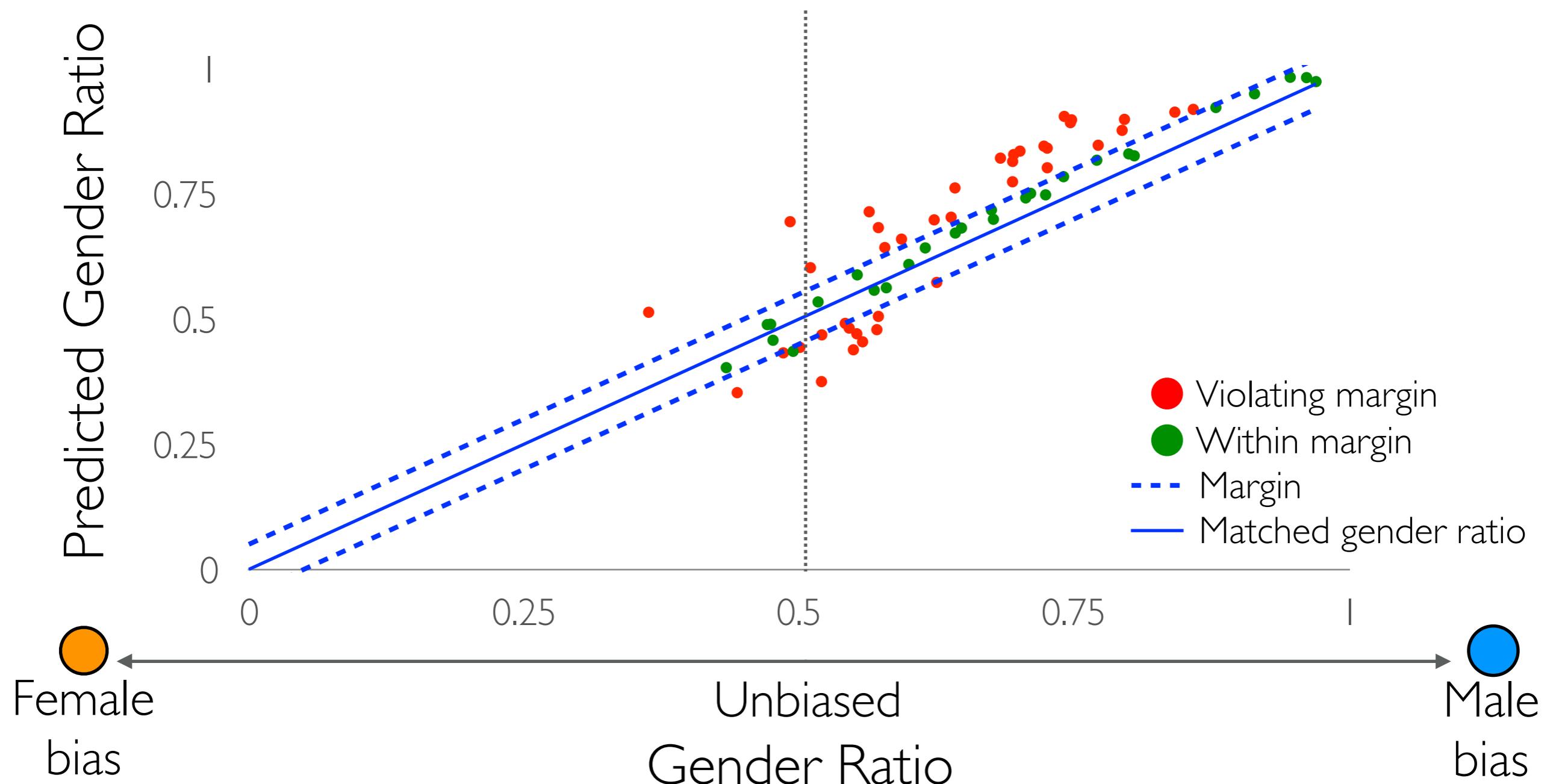
Gender Bias De-amplification in COCO

COCO Noun

Violation: 60.6%

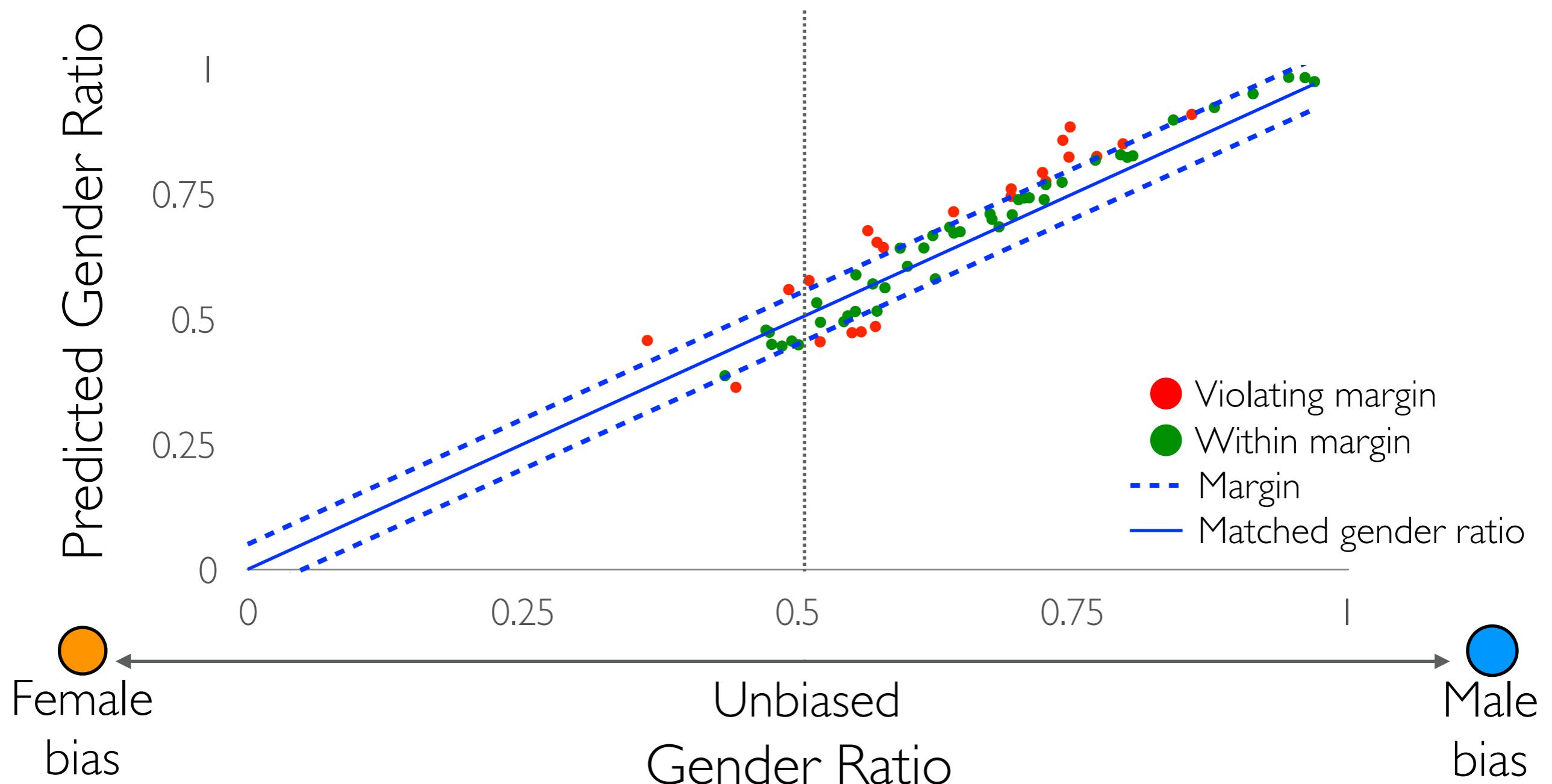
.032 |bias↑|

45.27 mAP



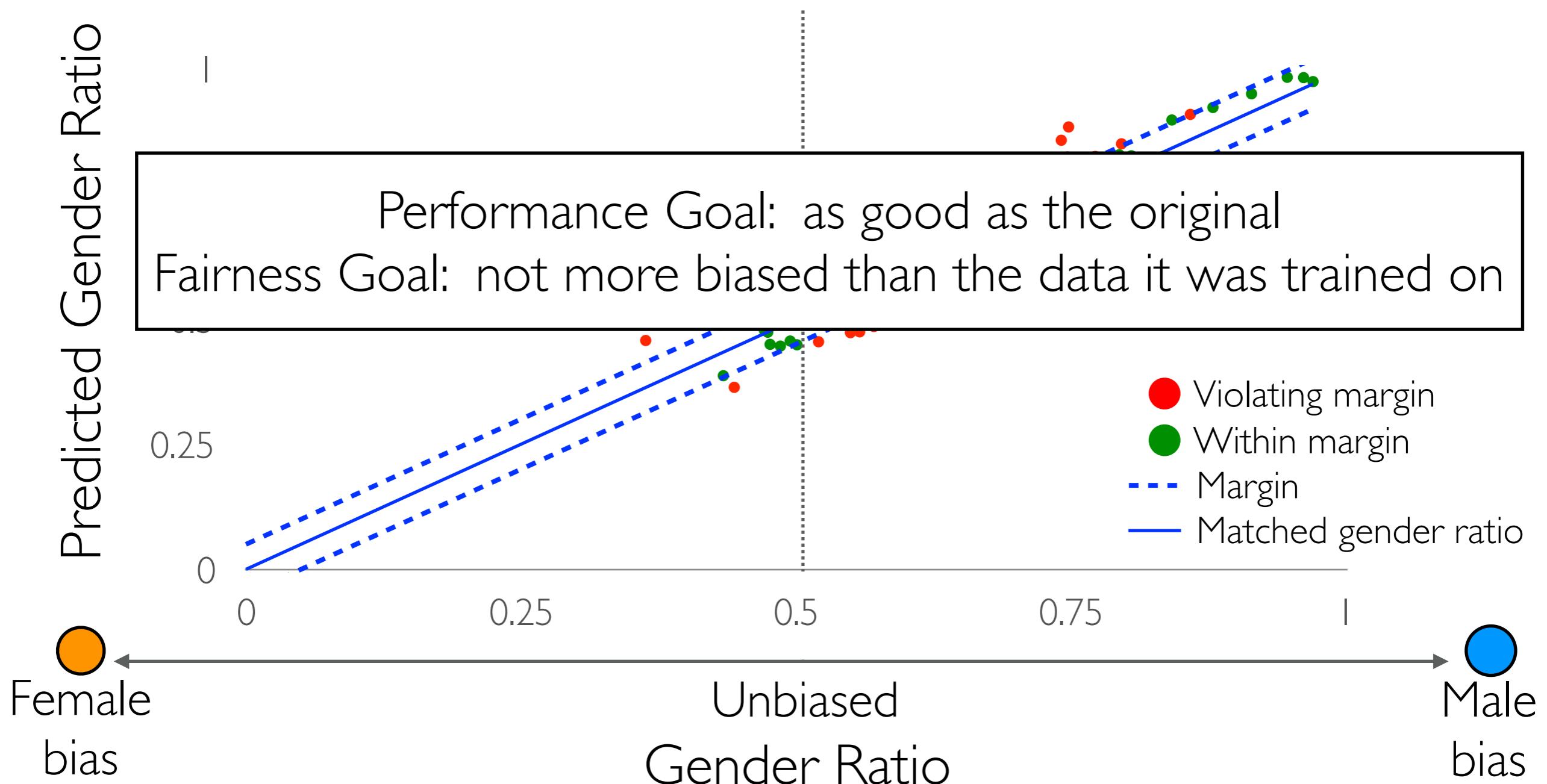
Gender Bias De-amplification in COCO

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



Gender Bias De-amplification in COCO

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



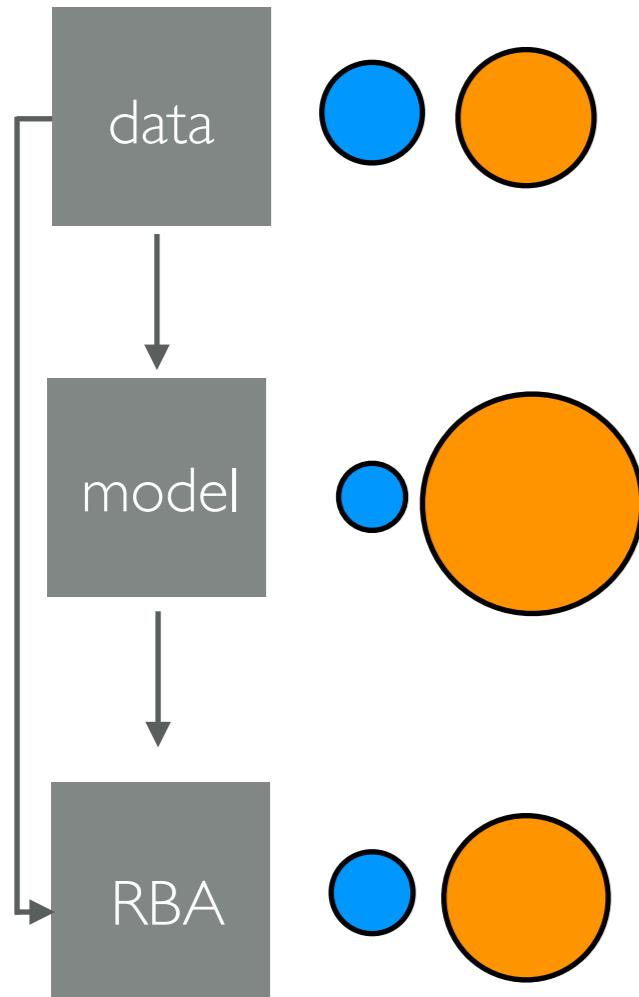
Contributions



imSitu vSRL
(events)



COCO MLC
(objects)

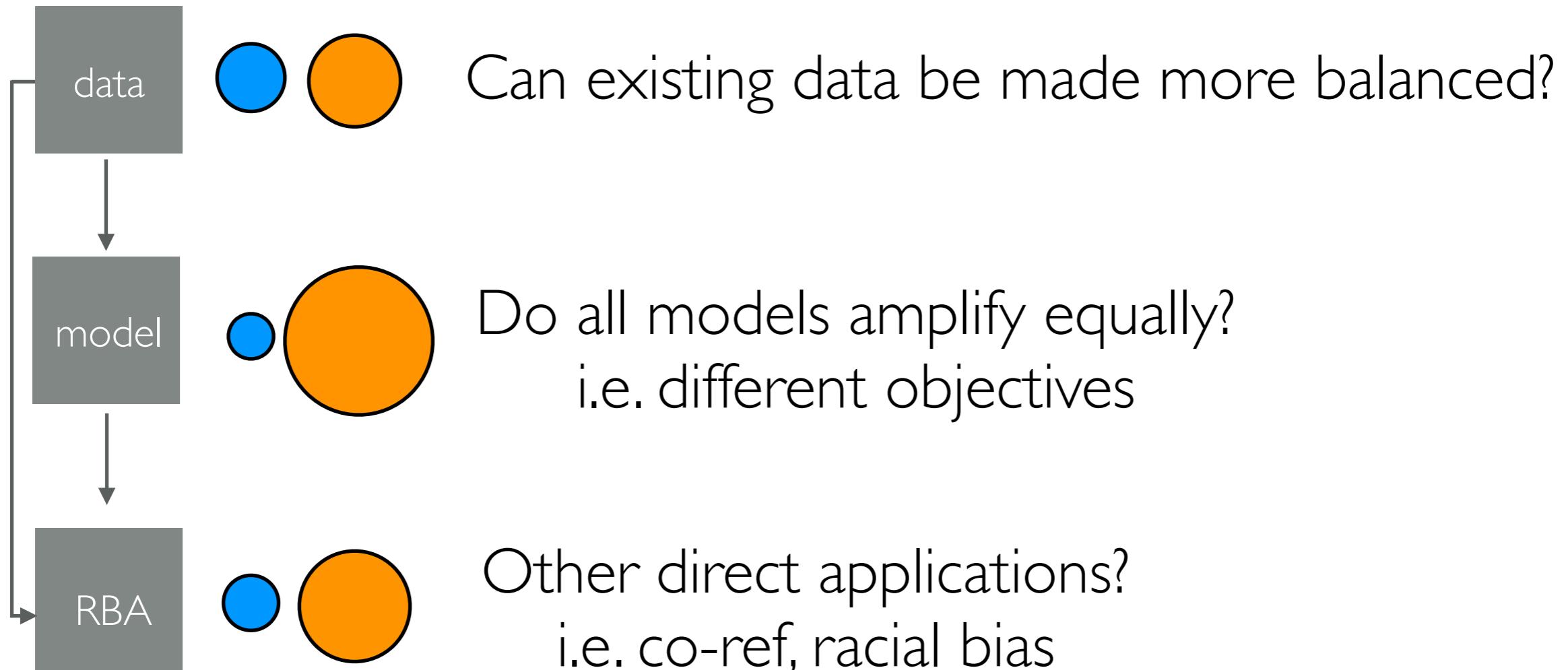


High dataset gender bias
38% (objects) 47% (events) exhibit strong bias

Models amplify existing gender bias
~70% objects and events have bias amplification

Reducing bias amplification
~50% reduction in amplification
Insignificant loss in performance

Future Work



Questions?

<https://github.com/uclanlp/reducingbias>



imSitu vSRL
(events)



COCO MLC
(objects)

