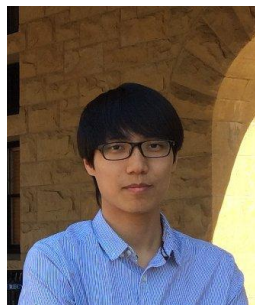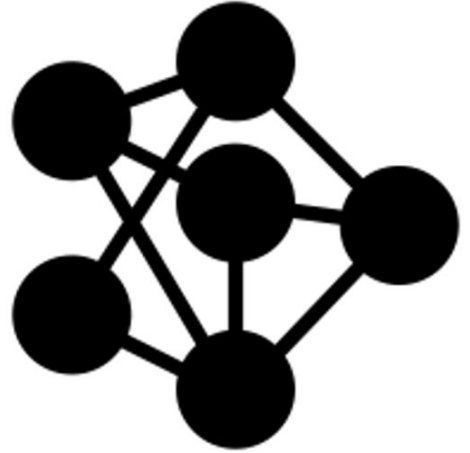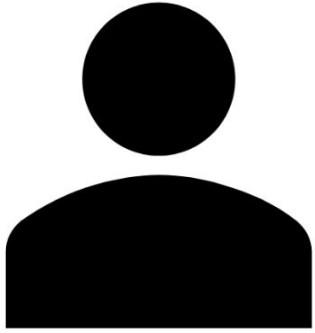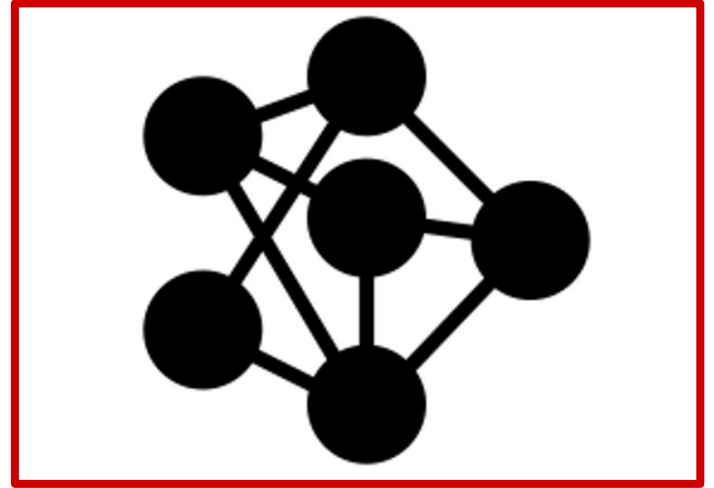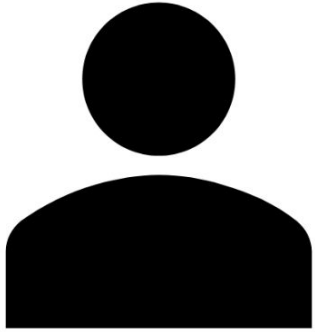# Human-in-the-Loop *Mixup*

Katie Collins, Umang Bhatt, Weiyang Liu, Vihari Piratla, Ilia Sucholutsky
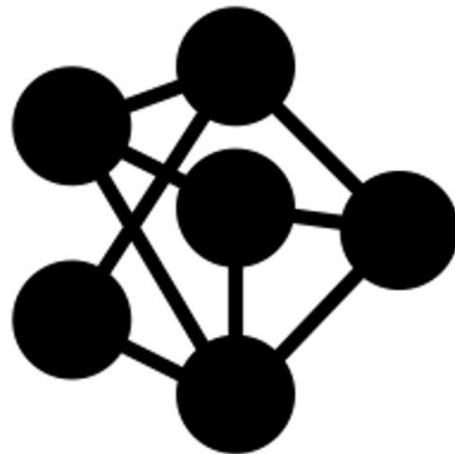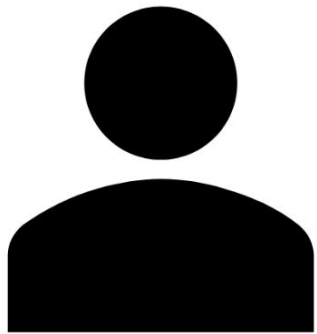Bradley Love, Adrian Weller

# Human-in-the-Loop *Mixup*

Katie Collins, Umang Bhatt, Weiyang Liu, Vihari Piratla, Ilia Sucholutsky
Bradley Love, Adrian Weller

Peterson*, Battleday*, et al 2019; Uma et al, 2020; Collins*, Bhatt*, Weller, 2022; Steyvers et al, 2022; Fel et al, 2022; Sucholutsky & Griffiths, 2023; Collins et al, 2023; Suchulotsky, Battleday, Collins et al, 2023, ... and several more!

Peterson*, Battleday*, et al 2019; Uma et al, 2020; Collins*, Bhatt*, Weller, 2022; Steyvers et al, 2022; Fel et al, 2022; Sucholutsky & Griffiths, 2023; Collins et al, 2023; Suchulotsky, Battleday, Collins et al, 2023, … and several more!

On the Informativeness of Supervision Signals

Ilia Sucholutsky[1]  Ruairidh M. Battleday[1]  Katherine M. Collins[2]  Raja Marjieh[3]  Joshua C. Peterson[1]
Pulkit Singh[1]  Umang Bhatt[2,4]  Nori Jacoby[5]  Adrian Weller[2,4]  Thomas L. Griffiths[1,2]

[1]Dept. of Computer Science, Princeton University,
[2]Dept. of Engineering, University of Cambridge,
[3]Dept. of Psychology, Princeton University,
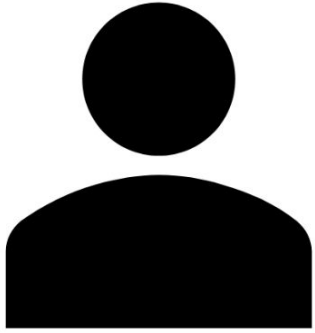[4]Alan Turing Institute,
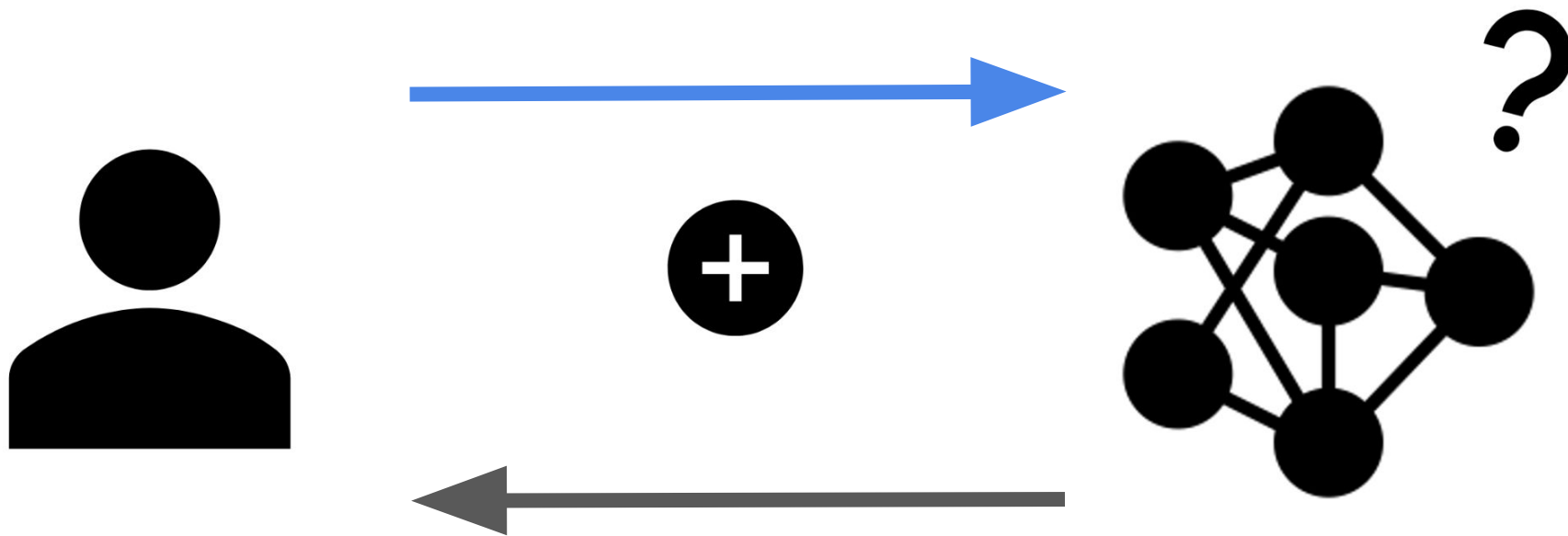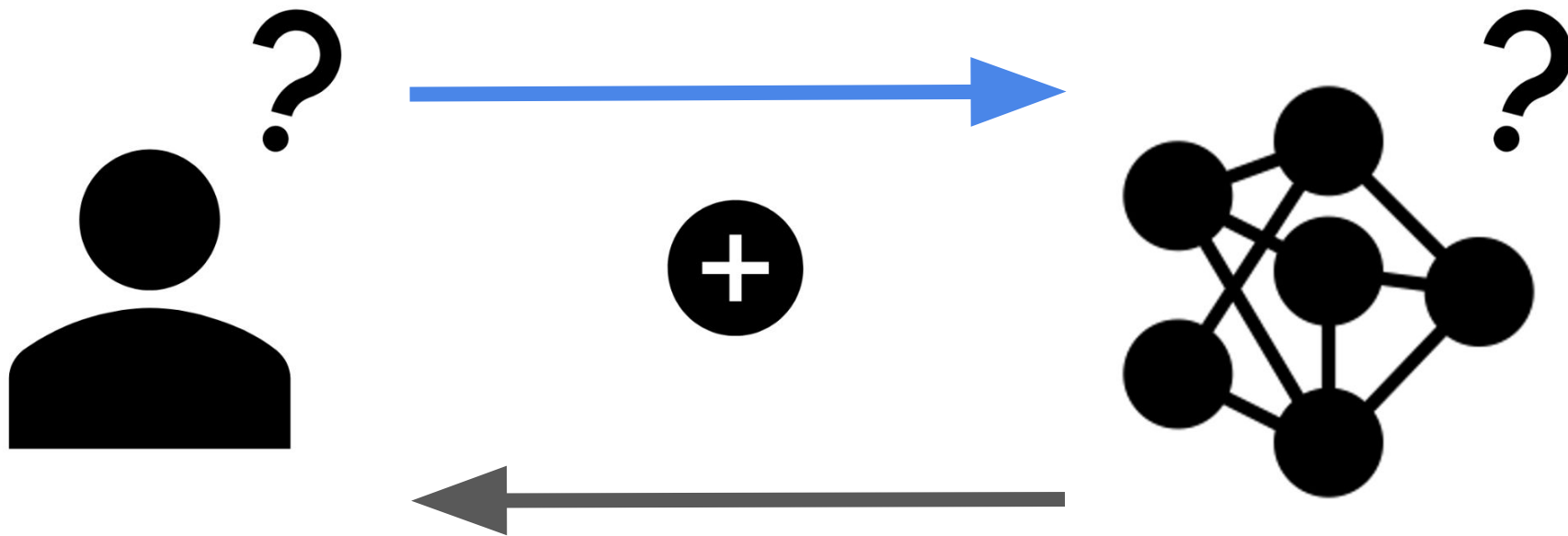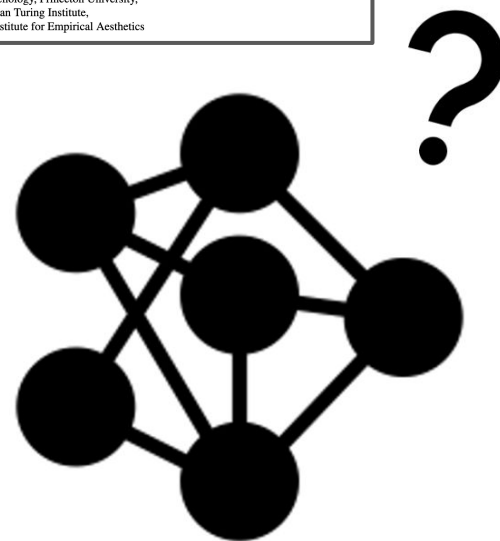[5]Max Planck Institute for Empirical Aesthetics

Peterson*, Battleday*, et al 2019; Uma et al, 2020; Collins*, Bhatt*, Weller, 2022; Steyvers et al, 2022; Fel et al, 2022; Sucholutsky & Griffiths, 2023; Collins et al, 2023; Suchulotsky, Battleday, Collins et al, 2023, ... and several more!
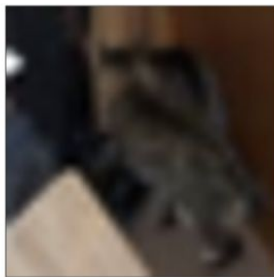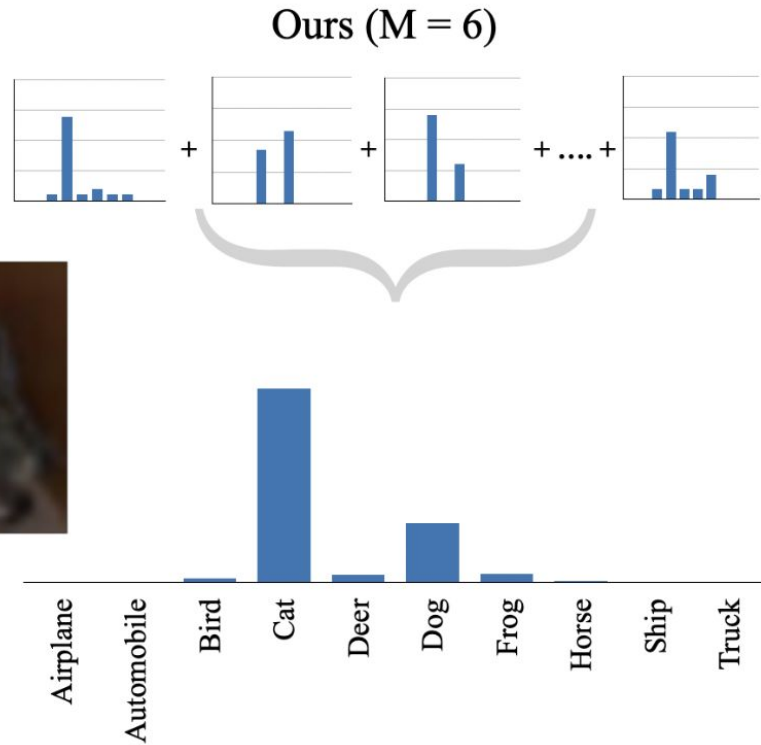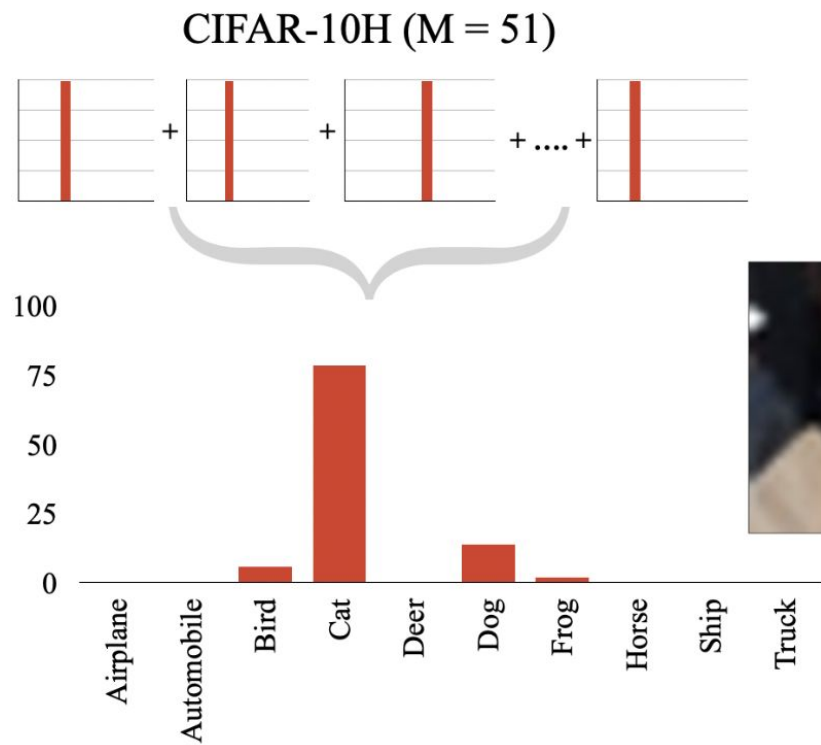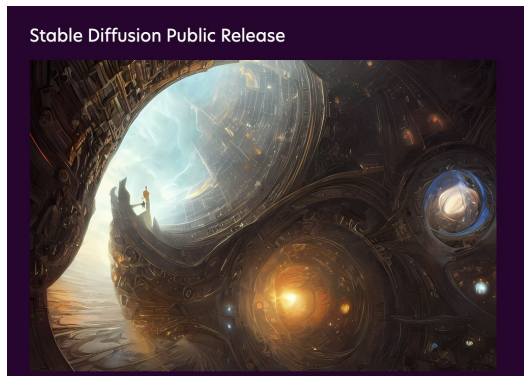
CIFAR-10H (M = 51)

Ours (M = 6)

Figure from Collins*, Bhatt*, Weller, 2022
Peterson*, Battleday*, et al 2019
Krizhevsky, 2009

# What about Synthetic Data?

# What about Synthetic Data?


Stability AI, 2022

**Review**

Next-generation deep learning based on simulators and synthetic data

Celso M. de Melo [1,*] Antonio Torralba,[2] Leonidas Guibas,[3] James DiCarlo,[4] Rama Chellappa,[5] and Jessica Hodgins[6]



Engine        Fire        Rain

Figure from Girdhar*, El-Nouby* et al, 2023


Zhang et al, 2017; Krizhevsky, 2009

**nature**

Explore content ∨    About the journal ∨    Publish with us ∨    Subscribe

nature > outlook > article

OUTLOOK | 27 April 2023

## Synthetic data could be better than real data

Machine-generated data sets have the potential to improve privacy and representation in artificial intelligence, if researchers can find the right balance between accuracy and fakery.

Neil Savage

### Synthetic Data - what, why and how?

James Jordon
jjordon@turing.ac.uk

Lukasz Szpruch
l.szpruch@ed.ac.uk

Florimond Houssiau
fhoussiau@turing.ac.uk

Mirko Bottarelli
mirko.bottarelli@warwick.ac.uk

Giovanni Cherubin
gcherubin@turing.ac.uk

Carsten Maple
cm@warwick.ac.uk

Samuel N. Cohen
scohen@turing.ac.uk

Adrian Weller
aweller@turing.ac.uk

**The Alan Turing Institute**    **THE ROYAL SOCIETY**

# What about Synthetic Data?



Stable Diffusion Public Release

Stability AI, 2022
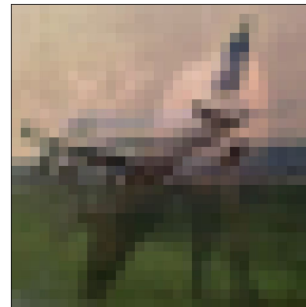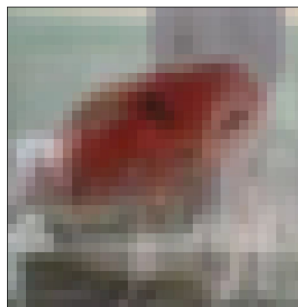


Engine    Fire    Rain

Figure from Girdhar*, El-Nouby* et al, 2023

## Review
# Next-generation deep learning based on simulators and synthetic data

Celso M. de Melo [1,*] Antonio Torralba, [2] Leonidas Guibas, [3] James DiCarlo, [4] Rama Chellappa, [5] and Jessica Hodgins [6]



Zhang et al, 2017; Krizhevsky, 2009

### Synthetic Data - what, why and how?

James Jordon
jjordon@turing.ac.uk

Lukasz Szpruch
l.szpruch@ed.ac.uk

Florimond Houssiau
fhoussiau@turing.ac.uk

Mirko Bottarelli
mirko.bottarelli@warwick.ac.uk

Giovanni Cherubin
gcherubin@turing.ac.uk

Carsten Maple
cm@warwick.ac.uk

Samuel N. Cohen
scohen@turing.ac.uk

Adrian Weller
aweller@turing.ac.uk

The Alan Turing Institute        THE ROYAL SOCIETY

## How perceptually-sensible are synthetic examples?
## Aligning model + human reprs?

# What about Synthetic Data?



Stable Diffusion Public Release

Stability AI, 2022



Engine  Fire  Rain

Figure from Girdhar*, El-Nouby* et al, 2023
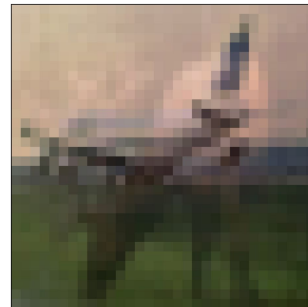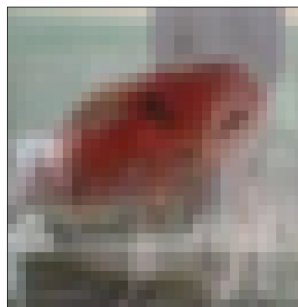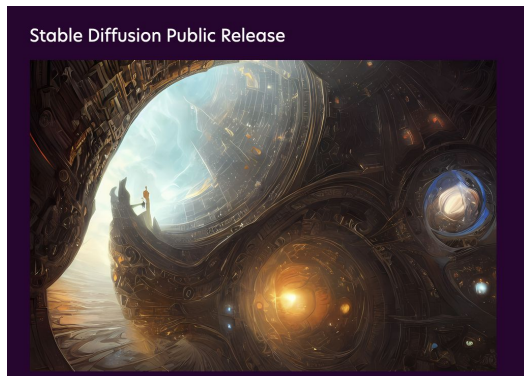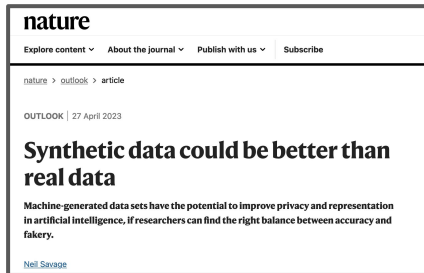


Review

Next-generation deep learning based on simulators and synthetic data

Celso M. de Melo [1,*] Antonio Torralba,[2] Leonidas Guibas,[3] James DiCarlo,[4] Rama Chellappa,[5] and Jessica Hodgins[6]



Zhang et al, 2017; Krizhevsky, 2009

**How perceptually-sensible are synthetic examples? Aligning model + human reprs?**

nature

Explore content ∨  About the journal ∨  Publish with us ∨  Subscribe

nature > outlook > article

OUTLOOK | 27 April 2023

## Synthetic data could be better than real data

Machine-generated data sets have the potential to improve privacy and representation in artificial intelligence, if researchers can find the right balance between accuracy and fakery.

Neil Savage

### Synthetic Data - what, why and how?

| | |
|---|---|
| James Jordon | Lukasz Szpruch |
| jjordon@turing.ac.uk | l.szpruch@ed.ac.uk |
| Florimond Houssiau | Mirko Bottarelli |
| fhoussiau@turing.ac.uk | mirko.bottarelli@warwick.ac.uk |
| Giovanni Cherubin | Carsten Maple |
| gcherubin@turing.ac.uk | cm@warwick.ac.uk |
| Samuel N. Cohen | Adrian Weller |
| scohen@turing.ac.uk | aweller@turing.ac.uk |

**The Alan Turing Institute**

**THE ROYAL SOCIETY**

# This Talk

- Why *Mixup*?
- Overview of *Mixup* Data Generation
- HMix and HILL MixE Suite
- Learning with Human Relabelings
- Taking Stock and Looking Ahead
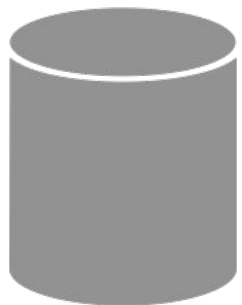
# Why *Mixup*?

Zhang et al, 2017

# Why *Mixup*?

- Simple generative process

Data Mixing Policy: $f(x_i, x_j, \lambda_f) = \lambda_f x_i + (1 - \lambda_f)x_j = \tilde{x}$
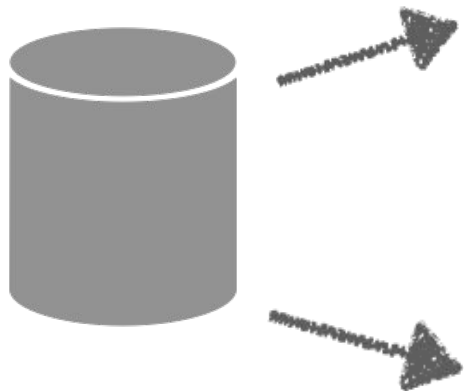
Label Mixing Policy: $g(y_i, y_j, \lambda_g) = \lambda_g y_i + (1 - \lambda_g)y_j = \tilde{y}$

- Powerful and popular regularizer + calibrator
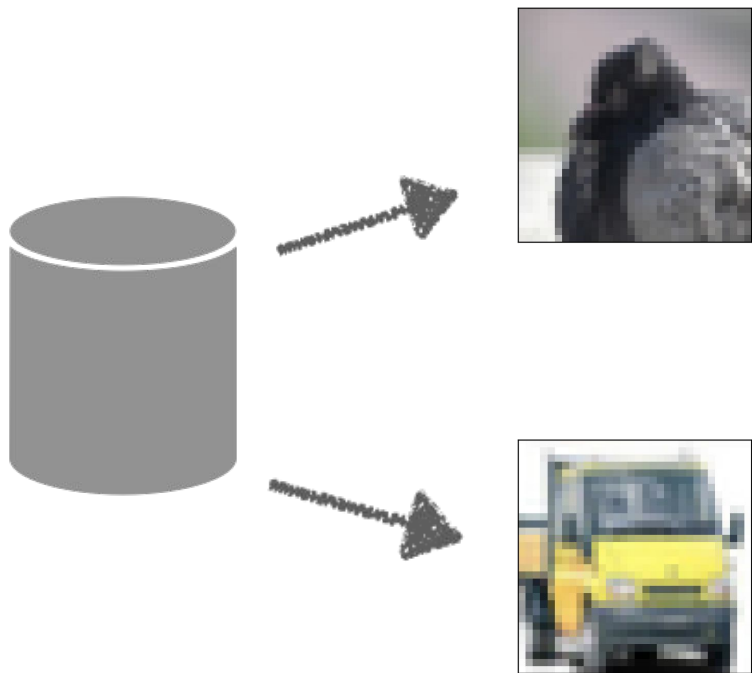- Cognitive neuroscience suggests misalignment

Zhang et al, 2017

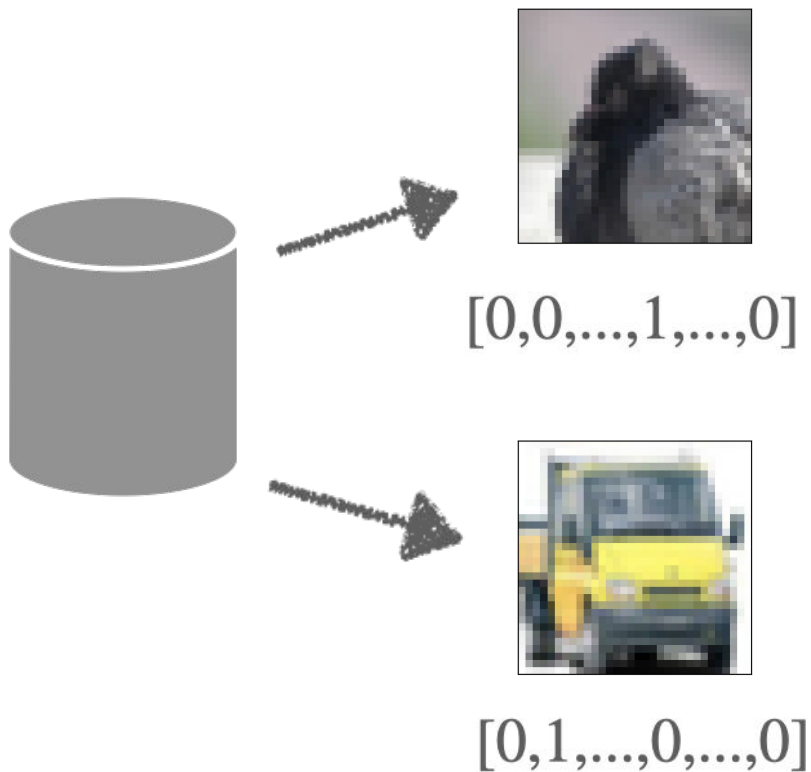# *Mixup* Generative Process

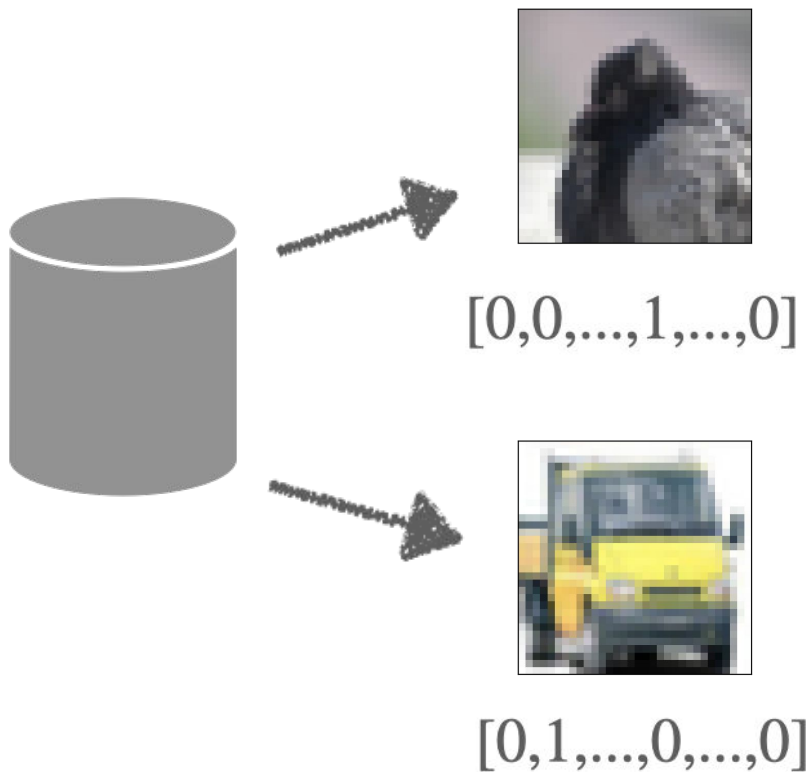# *Mixup* Generative Process

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$[0,0,...,1,...,0]$



$[0,1,...,0,...,0]$

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$$\lambda_f = 0.1$$

[0,0,...,1,...,0]

[0,1,...,0,...,0]

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$[0,0,...,1,...,0]$

$\lambda_f = 0.1$

$[0,1,...,0,...,0]$

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$$\lambda_f = 0.1$$

$[0,0,...,1,...,0]$

$[0,0.9,...,0.1,...,0]$

$[0,1,...,0,...,0]$

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$$\lambda_f = 0.1 \qquad \lambda_f = 0.5$$

$[0,0,...,1,...,0]$

$[0,0.9,...,0.1,...,0]$

$[0,1,...,0,...,0]$

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$\lambda_f = 0.1$      $\lambda_f = 0.5$

$[0,0,...,1,...,0]$

$[0,0.9,...,0.1,...,0]$

$[0,1,...,0,...,0]$

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$$\lambda_f = 0.1 \qquad \lambda_f = 0.5$$

[0,0,...,1,...,0]

[0,1,...,0,...,0]

[0,0.9,...,0.1,...,0]    [0,0.5,...,0.5,...,0]

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$[0,0,...,1,...,0]$

$\lambda_f = 0.1$  $\lambda_f = 0.5$  $\lambda_f = 0.7$

$[0,0.9,...,0.1,...,0]$  $[0,0.5,...,0.5,...,0]$

$[0,1,...,0,...,0]$

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$[0,0,...,1,...,0]$

$[0,1,...,0,...,0]$

$\lambda_f = 0.1$    $\lambda_f = 0.5$    $\lambda_f = 0.7$

$[0,0.9,...,0.1,...,0]$    $[0,0.5,...,0.5,...,0]$

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$\lambda_f = 0.1$     $\lambda_f = 0.5$     $\lambda_f = 0.7$

$[0,0,...,1,...,0]$

$[0,1,...,0,...,0]$

$[0,0.9,...,0.1,...,0]$     $[0,0.5,...,0.5,...,0]$     $[0,0.3,...,0.7,...,0]$

Zhang et al, 2017; Krizhevsky, 2009

# *Mixup* Generative Process



$$[0,0,...,1,...,0]$$

$$[0,1,...,0,...,0]$$

$\lambda_f = 0.1$  $\lambda_f = 0.5$  $\lambda_f = 0.7$

$[0,0.9,...,0.1,...,0]$  $[0,0.5,...,0.5,...,0]$  $[0,0.3,...,0.7,...,0]$

Zhang et al, 2017; Krizhevsky, 2009
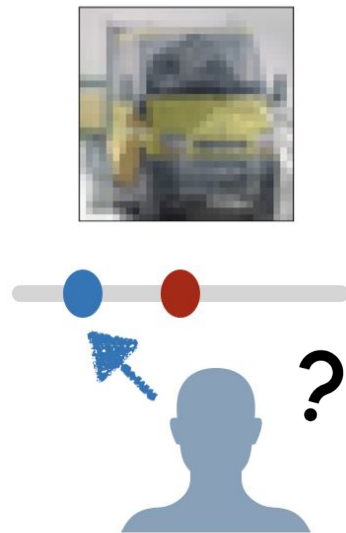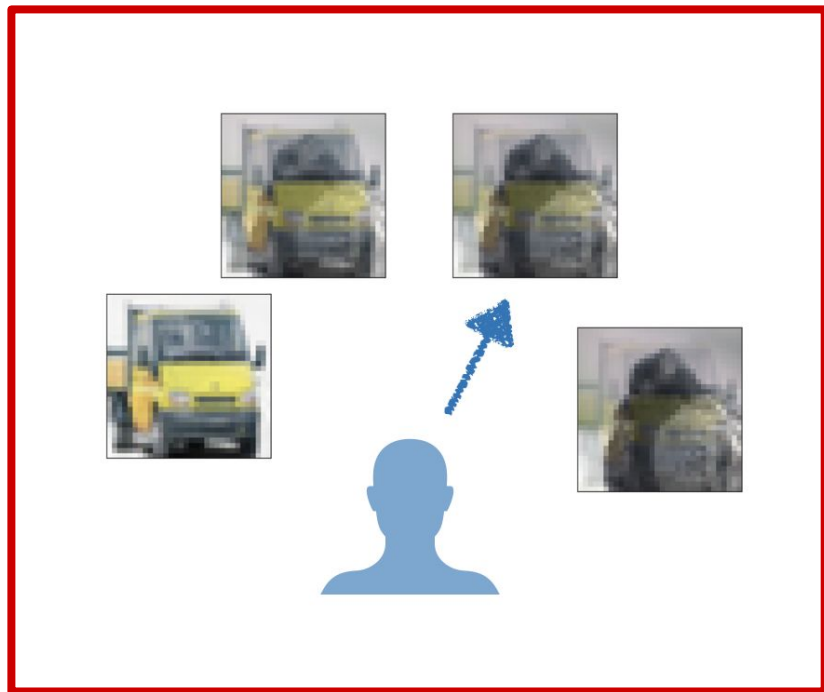
# Eliciting Human Percepts of Synthetic Examples

# Eliciting Human Percepts of Synthetic Examples

# Eliciting Human Percepts of Synthetic Examples

# Eliciting Human Percepts of Synthetic Examples

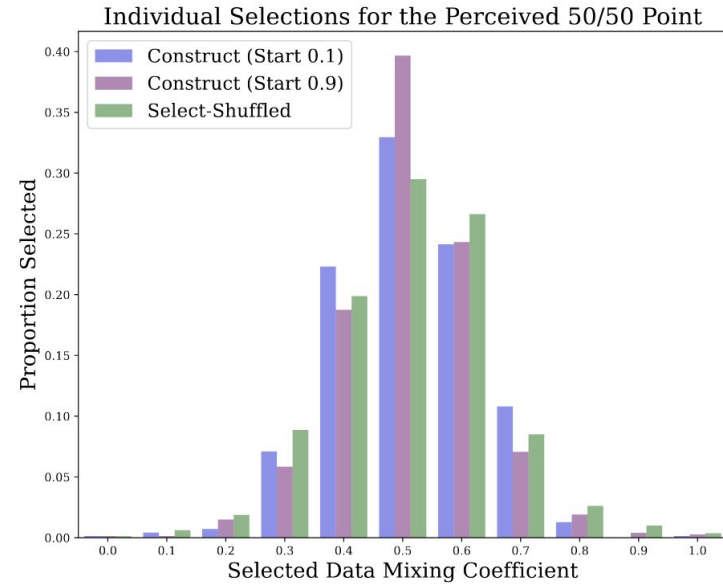# Eliciting Human Percepts of Synthetic Examples

# Selecting a Matching Midpoint

- 249 mixed images
- 70 participants
- 2 interface types
  - Construct
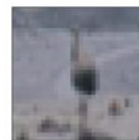  - Select-Shuffled

# Selecting a Matching Midpoint

- 249 mixed images
- 70 participants
- 2 interface types
  - Construct
  - Select-Shuffled
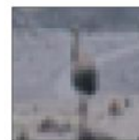


Individual Selections for the Perceived 50/50 Point

# Selecting a Matching Midpoint



$\lambda_f = 0.0$



$\lambda_f = 1.0$

# Selecting a Matching Midpoint



$\lambda_f = 0.0$

$\lambda_f = 0.5$
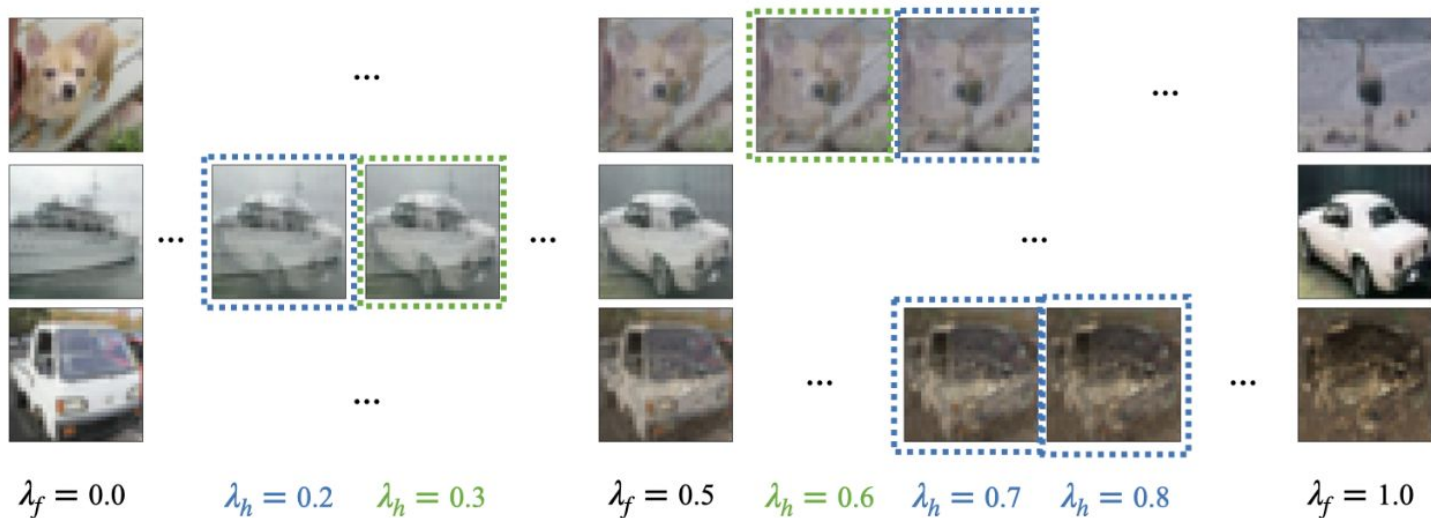
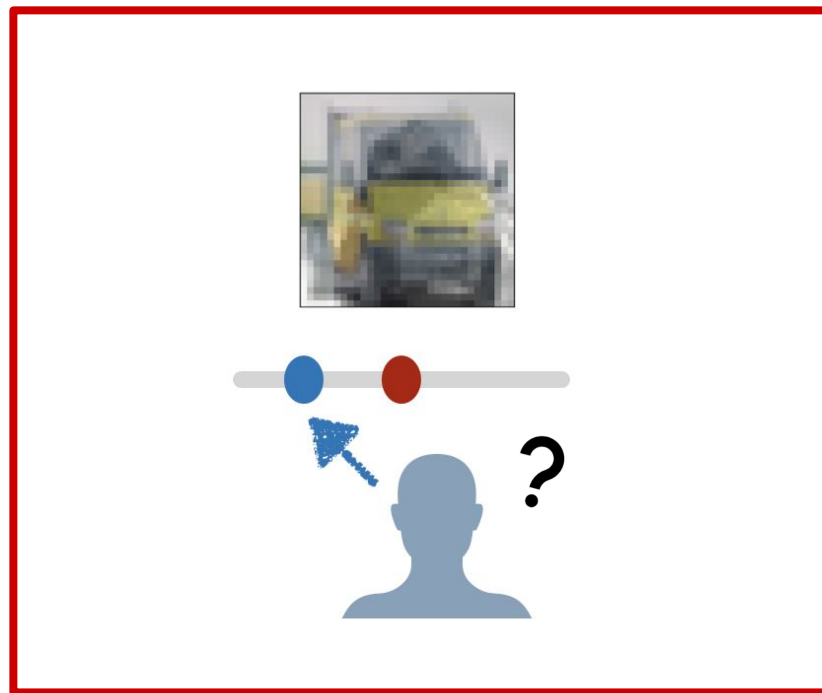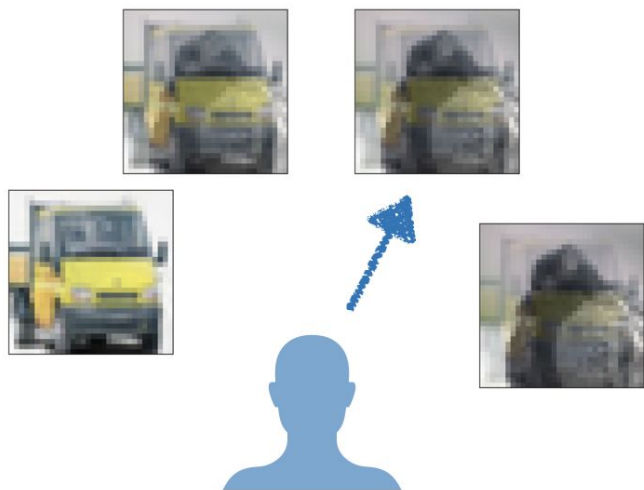$\lambda_f = 1.0$

# Selecting a Matching Midpoint



$\lambda_f = 0.0$  $\lambda_h = 0.2$  $\lambda_h = 0.3$  $\lambda_f = 0.5$  $\lambda_h = 0.6$  $\lambda_h = 0.7$  $\lambda_h = 0.8$  $\lambda_f = 1.0$

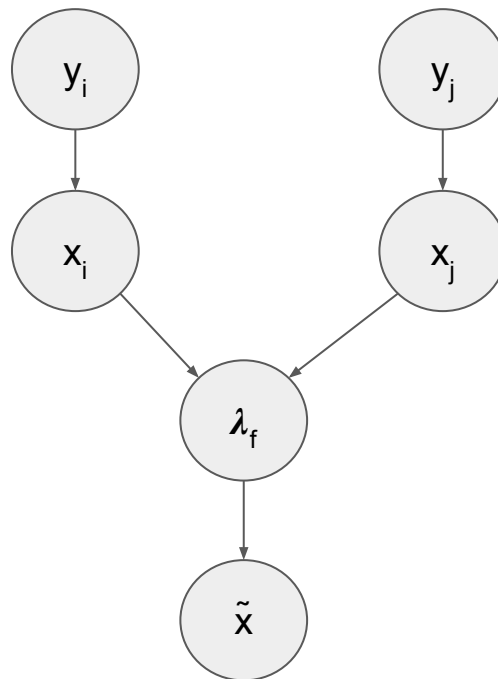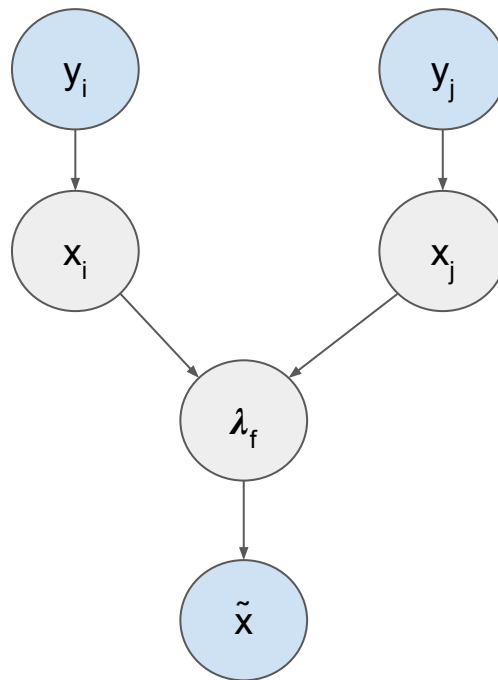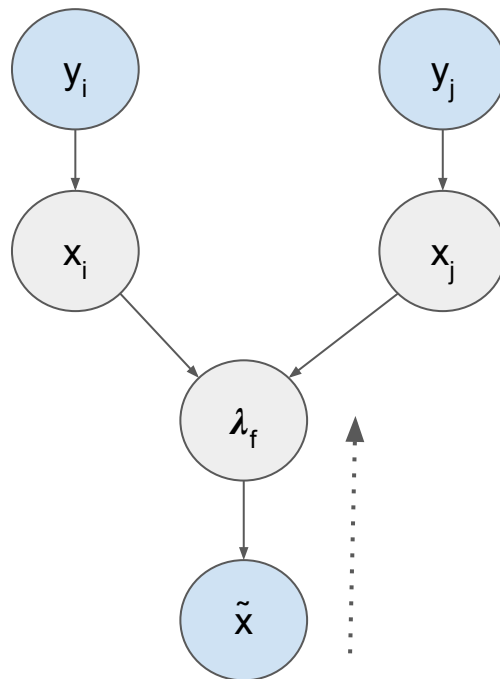# Eliciting Human Percepts of Synthetic Examples

# Inferring the Data Mixing Coefficient

- 2070 mixed images
- 81 participants

# Inferring the Data Mixing Coefficient

- 2070 mixed images
- 81 participants

# Inferring the Data Mixing Coefficient

- 2070 mixed images
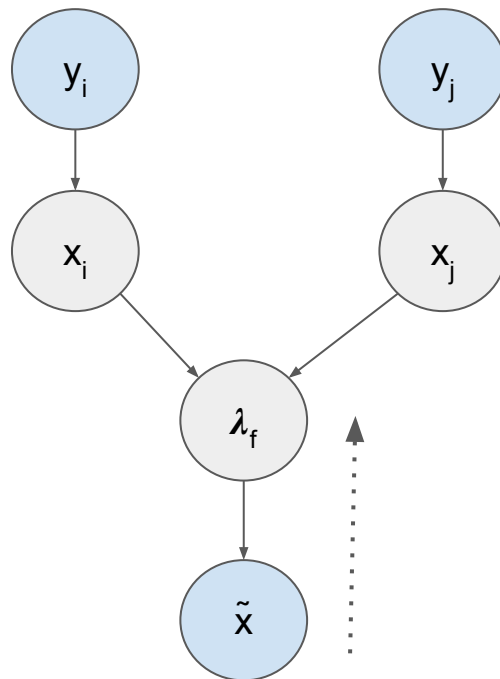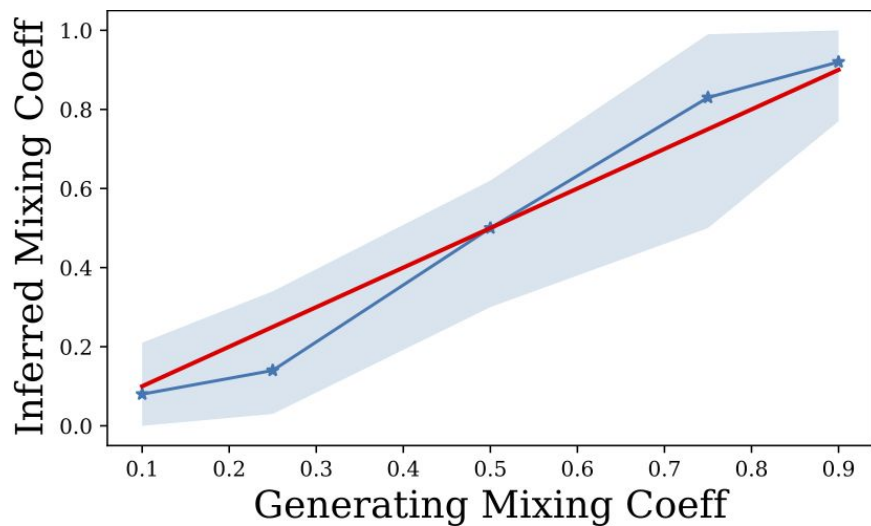- 81 participants

# Inferring the Data Mixing Coefficient

- 2070 mixed images
- 81 participants

# Inferring the Data Mixing Coefficient

- 2070 mixed images
- 81 participants

# Aligning Model Representations with Human Percepts?

# Aligning Model Representations with Human Percepts?
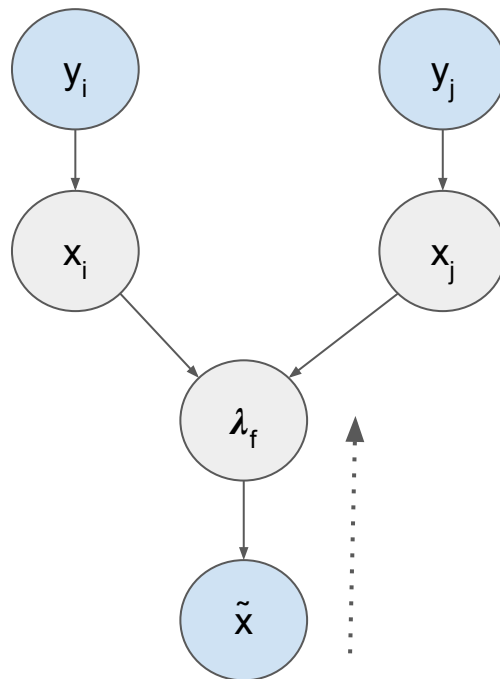
| Generalization | Calibration | Robustness |
|:---:|:---:|:---:|

Szegedy et al, 2014; Hendrycks and Dietterich, 2019; Bhatt et al, 2021; Thomas and Uminisky, 2022

# Relabeling with Human Perceptual Judgments

# Relabeling with Human Perceptual Judgments

| Label Type | CE | FGSM | Calib |
|---|---|---|---|
| Regular (No Aug) | $2.02\pm0.12$ | $13.12\pm2.65$ | $0.28\pm0.011$ |
| + Random | $2.11\pm0.13$ | $12.81\pm2.84$ | $0.24\pm0.014$ |
| + Uniform | $2.16\pm0.14$ | $12.71\pm2.79$ | $0.25\pm0.012$ |
| + *mixup* | $1.65\pm0.11$ | $10.62\pm2.44$ | $0.23\pm0.005$ |
| + Ours (Relabel) | $1.78\pm0.12$ | $11.69\pm2.90$ | $0.24\pm0.009$ |

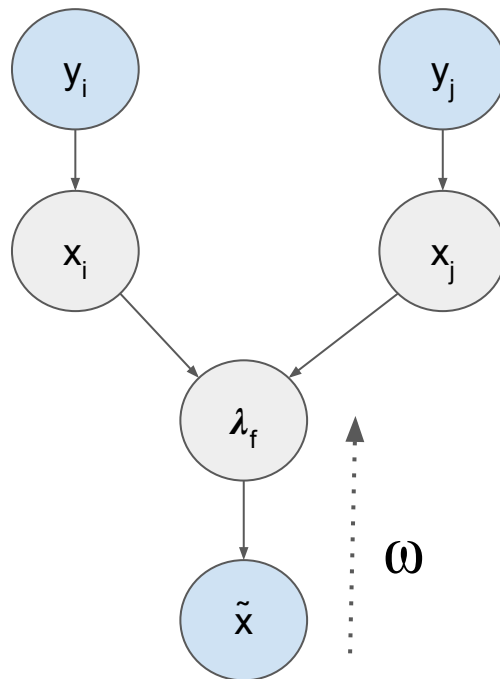# Human *Uncertainty* in Inference

# Human *Uncertainty* in Inference

# Human *Uncertainty* in Inference
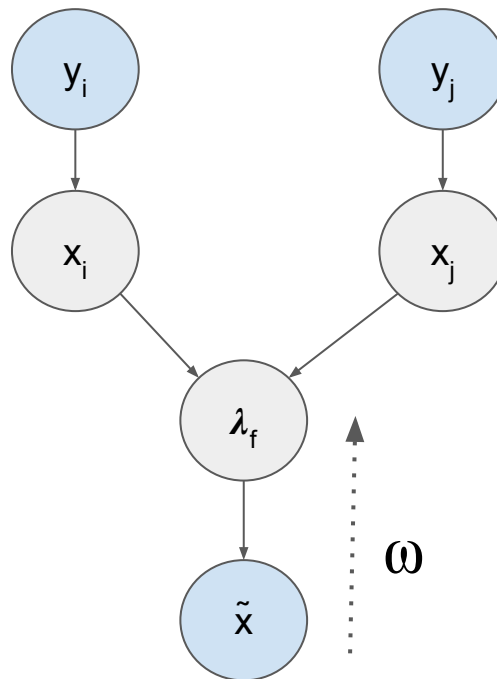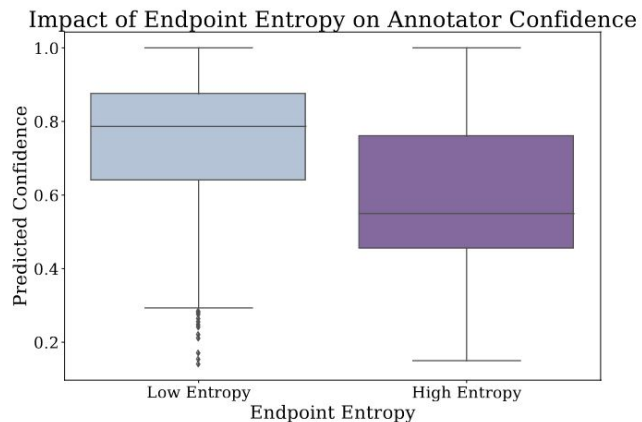
# Human *Uncertainty* in Inference

| Mixing Coefficient | Reported Confidence |
|---|---|
| 0.1 | $0.79 \pm 0.17$ |
| 0.25 | $0.72 \pm 0.20$ |
| 0.5 | $0.63 \pm 0.20$ |

# Human *Uncertainty* in Inference

| Mixing Coefficient | Reported Confidence |
|---|---|
| 0.1 | $0.79 \pm 0.17$ |
| 0.25 | $0.72 \pm 0.20$ |
| 0.5 | $0.63 \pm 0.20$ |



Impact of Endpoint Entropy on Annotator Confidence

# Relabeling with Human Perceptual Judgments

| Label Type | CE | FGSM | Calib |
|---|---|---|---|
| Regular (No Aug) | 2.02±0.12 | 13.12±2.65 | 0.28±0.011 |
| + Random | 2.11±0.13 | 12.81±2.84 | 0.24±0.014 |
| + Uniform | 2.16±0.14 | 12.71±2.79 | 0.25±0.012 |
| + *mixup* | 1.65±0.11 | 10.62±2.44 | 0.23±0.005 |
| + Ours (Relabel) | 1.78±0.12 | 11.69±2.90 | 0.24±0.009 |
| (Relabel & $\omega$) | **1.48±0.06** | **8.89±1.59** | **0.19±0.001** |

# Is human relabeling scalable?

# Relabeling with (In-Filled) Human Perceptual Judgments
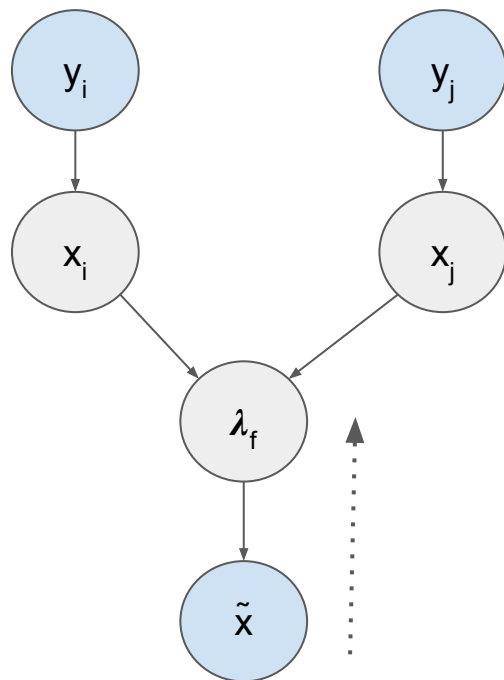
# Relabeling with (In-Filled) Human Perceptual Judgments

# Relabeling with (In-Filled) Human Perceptual Judgments

- Fit logistic functions per category boundary

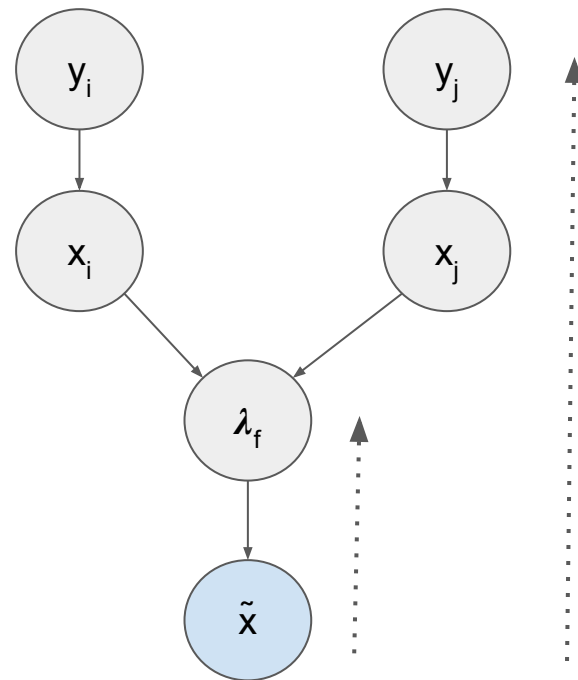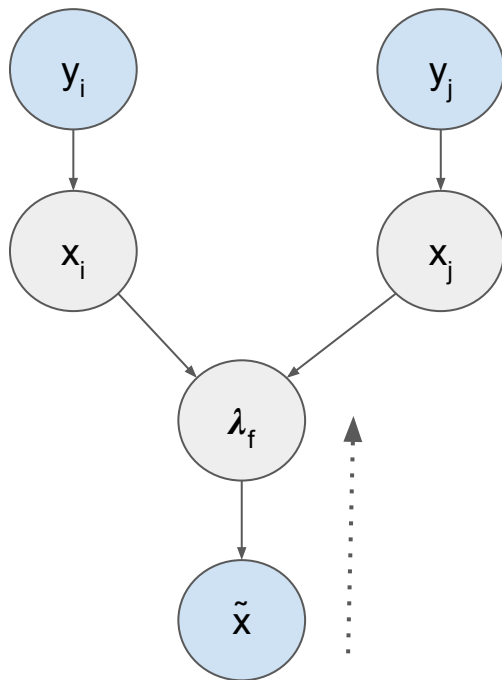| Label Policy | CE | FGSM | Calib |
|---|---|---|---|
| *mixup* | **1.15±0.08** | 7.46±2.40 | **0.10±0.01** |
| Human-Fits (Ours) | 1.16±0.08 | **7.32±2.27** | **0.10±0.01** |

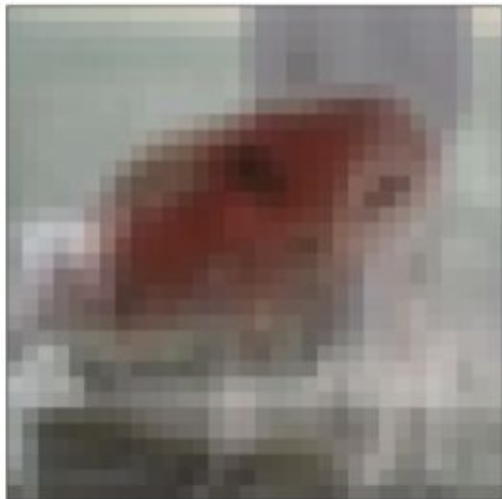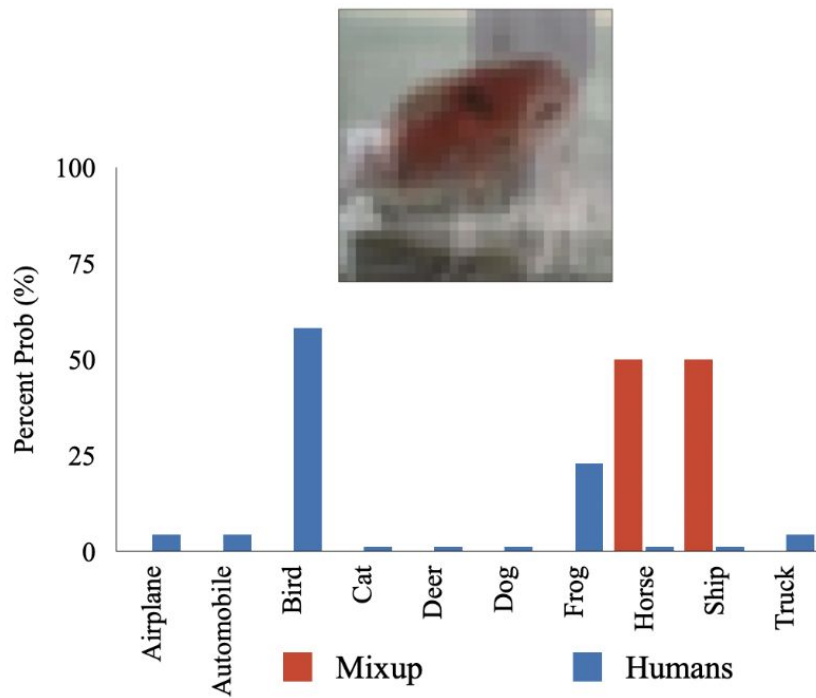# Richer Human Uncertainty

# Richer Human Uncertainty

# Richer Human Uncertainty

# Richer Human Uncertainty

# Richer Human Uncertainty

# Takeaways

- Synthetic examples generated in *mixup* likely differ in fundamental ways from human perception

# Takeaways

- Synthetic examples generated in *mixup* likely differ in fundamental ways from human perception
- Relabeling with human perceptual judgments — **espec accounting for human uncertainty** — has potential to possibly improve performance

# Takeaways

- Synthetic examples generated in *mixup* likely differ in fundamental ways from human perception
- Relabeling with human perceptual judgments — **espec accounting for human uncertainty** — has potential to possibly improve performance
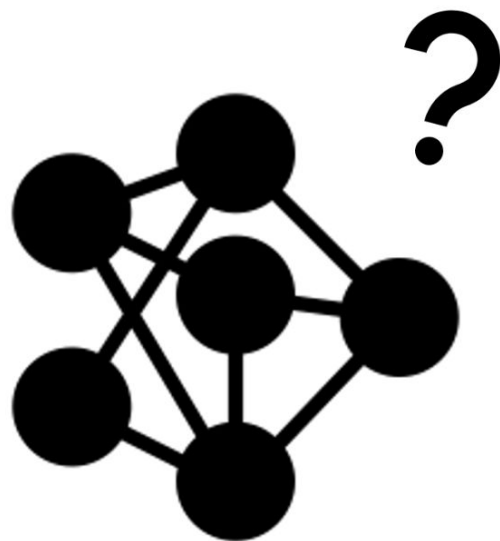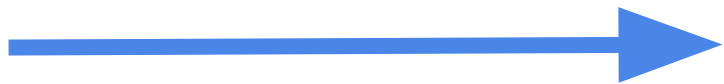- Scalability challenges

# Takeaways

- Synthetic examples generated in *mixup* likely differ in fundamental ways from human perception
- Relabeling with human perceptual judgments — **espec accounting for human uncertainty** — has potential to possibly improve performance
- Scalability challenges

HILL MixE Suite Interfaces

H-Mix Data

https://github.com/cambridge-mlg/hill-mixup

For more details,
please check out our paper + poster :)

H-Mix Data + HILL MixE Suite interfaces at our repo:
https://github.com/cambridge-mlg/hill-mixup

More questions? Thoughts?
kmc61@cam.ac.uk