# EE-559
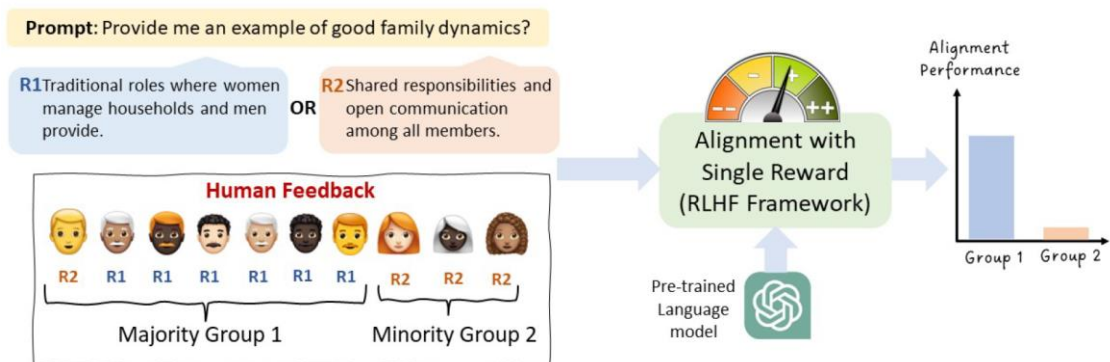# Deep Learning

## What's on today?

- Alignment with diverse preferences: on AI respecting the values of all
- Reinforcement learning from AI feedback: on AI training AI
- Interpretability: on identifying decision-defining features and paths
- Causal mediation analysis: on unpacking causal effects
- Datasheets for datasets: on responsible data collection and use

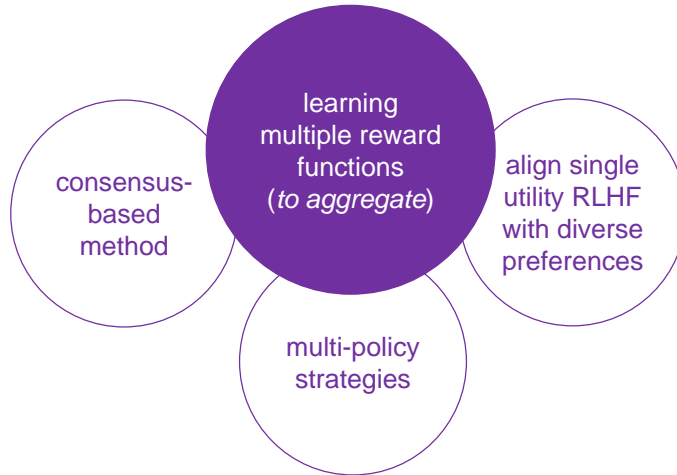# Alignment with diverse preferences

## Majority vs minority user groups



most RLHF approaches ignore diversity in human preference feedback by aligning the model with a **single reward function**

arXiv:2402.08925

# Diversity in opinions and preferences



consensus-based method

learning multiple reward functions (*to aggregate*)

align single utility RLHF with diverse preferences

multi-policy strategies

arXiv:2402.08925

# Mixture of preference distributions

$$P_u^*(\boldsymbol{y}_1 \succsim \boldsymbol{y}_2 \mid \boldsymbol{x}) = \mathbb{E}_{h \in H_u}[I(h \text{ prefers } \boldsymbol{y}_1 \text{ over } \boldsymbol{y}_2 \mid \boldsymbol{x})]$$  for all groups in $U$

$$U = \{H_1, H_2, \dots, H_{|U|}\} \qquad H = \bigcup_{u=1}^{|U|} H_u \qquad u: \text{human subpopulation index}$$

$$P^*(\boldsymbol{y}_1 \succsim \boldsymbol{y}_2 \mid \boldsymbol{x}) = \sum_{u=1}^{|U|} \left[ \sum_{h \in H_u} I_h(\boldsymbol{z}) q(h|u) \right] \eta(u) = \sum_{u=1}^{|U|} p_u^*(\boldsymbol{z}) \eta(u)$$

$\boldsymbol{z} = (\boldsymbol{y}_1 \succsim \boldsymbol{y}_2 \mid \boldsymbol{x})$

distribution over the humans $H$

marginal probability distribution of *subpopulation $H_u$*

subpopulation with *specific preference distribution*

arXiv:2402.08925

3

# Mixture of preference distributions

$$p(\boldsymbol{z}') = \sum_{u=1}^{|U|} p_{\phi_u}^*(\boldsymbol{z}')\,\eta(u)$$  preference distribution  $\phi_u$ reward model parameter

$$\boldsymbol{z}' = (\boldsymbol{y}_w \succcurlyeq \boldsymbol{y}_l|\, \boldsymbol{x})$$

$\boldsymbol{y}_w$  chosen response by the human sub-population group $H_u$

$\boldsymbol{y}_l$  rejected response by the human sub-population group $H_u$

$$L(\phi) = \sum_{\boldsymbol{z}' \in D} \log \sum_{u=1}^{|U|} p_{\phi_u}(\boldsymbol{z}')\,\eta(u)$$

$$= \sum_{\boldsymbol{z}' \in D} \log \sum_{u=1}^{|U|} \frac{e^{r_{\phi_u}(\boldsymbol{y}_w, \boldsymbol{x})}}{e^{r_{\phi_u}(\boldsymbol{y}_w, \boldsymbol{x})} + e^{r_{\phi_u}(\boldsymbol{y}_l, \boldsymbol{x})}}\,\eta(u)$$  maximization of the log likelihood

arXiv:2402.08925

# Maximizing the minimum utility

**Alignment objective** with diverse human preferences and with KL-regularization

$$\underset{p}{\mathrm{argmax}} \left( \min_u \mathbb{E}_{\boldsymbol{x} \sim P, \boldsymbol{y} \sim p(.\,|\boldsymbol{x})}[r_{\phi_u^*}(\boldsymbol{y}, \boldsymbol{x})] \right) - \beta D_{KL}[p(.\,|\boldsymbol{x})||p_{\mathrm{REF}}(.\,|\boldsymbol{x})]$$

$\phi_u^*$ reward model parameter
  *for each human subpopulation* in $U$

$r_{\phi_u^*}$ reward model

$\beta > 0$  controls the deviation from the
  base reference policy $p_{\mathrm{REF}}$

arXiv:2402.08925

# **RL from AI feedback**

## Constitutional AI

**Reinforcement learning from AI feedback (RLAIF)**

input from humans:
list of **high-level principles**
(a "*constitution*" to guide the LLM training process)

deal with
potentially diverging
input from humans

aggregate the input into
consistent data about
"*collective*" preferences

separate LLM to generate **artificial preferences**
and **instruction data** for model fine-tuning

arXiv:2404.10271

# RL**AIF**

| | |
|---|---|
| RL**HF** distils *human preferences* into a single preference model | RL**AIF** distils *LM 'interpretations' of a set of principles* back into a hybrid human/AI preference model |

*use human labels for* **helpfulness**
*but only AI labels for* **harmlessness**

arXiv:2404.10271

# Constitutional AI

scale
supervision

reduce tension
between
**helpfulness** &
**harmlessness**

transparency

AI systems to help
*supervise* other AIs

eliminate
*evasive responses*
&
*explain objections*
to harmful requests

*declare principles*
governing AI behaviour

arXiv:2212.08073

# Two-stage process

| | | | |
|---|---|---|---|
| **Stage 1** | critique | revision | supervised learning |
| **Stage 2** | AI comparison evaluations | preference model | reinforcement learning |

arXiv:2212.08073

# Supervised stage

| generate responses to harmful prompts using a helpful-only AI assistant | ask the model to critique its response according to a principle in the constitution | randomly draw principles from the constitution at each step |
|---|---|---|
| *responses will be **harmful** & **toxic*** | *then ask to **revise response** in light of the critique* | *revise responses **repeatedly*** |

once process is complete,
fine-tune a pre-trained language model
with **supervised learning on the final revised responses**

arXiv:2212.08073

# Reinforcement learning stage

use model trained in Stage 1 to generate a pair of responses to each prompt in a dataset of harmful prompts

formulate each prompt & pair into a multiple choice question, ask which response is best according to a constitutional principle

*produces an AI-generated* **preference dataset** *for* **harmlessness**, *which is mixed with human feedback* **helpfulness** *dataset*

train a **preference model** on this comparison data

fine-tune the model from Stage 1 via RL against this **preference model** resulting in a policy trained by RLAIF

arXiv:2212.08073

# Interpretability

# Human-understandable explanations
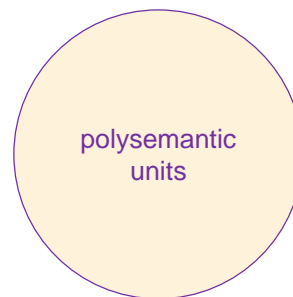
*How do neural networks calculate their outputs*?

*What are the internal processes*?

*Can we make targeted changes to these processes*?

# Neurons activate for multiple contexts

superposition

polysemantic units

*over-complete set of directions*
in activation space

*difficult-to-interpret* units

# Explaining behavior of models

| dictionary learning | → | sparse autoencoders | → | causal subnetwork |
|---|---|---|---|---|

<div style="text-align:center">

semantically meaningful
decomposition of the
activation space

</div>

arXiv:2403.19647

# Computational subgraphs

$$x, \epsilon(x), b \in \mathbb{R}^D$$

feature activations    bias

$$x = \hat{x} + \epsilon(x) = \sum_{i=1}^{S} f_i(x) v_i + b + \epsilon(x)$$

features
(*unit vectors*)    error term

identify directions in a latent
space that represent $S$
*human-interpretable features*

$$S = 64 \times D$$

**Concepts**: Feature disentanglement with SAEs;
SAE trained to minimize L2 *reconstruction error* and L1 *regularization term* (to promote sparsity).

# Attribution patching

$$\boldsymbol{a} \in \mathbb{R}^D$$
node in the graph

*c: clean*
*p: patch*

$$\boldsymbol{x}_c \in \mathbb{R}^D$$
$$\boldsymbol{x}_p \in \mathbb{R}^D$$
inputs

$$m(.)$$
real-valued metric

➡️

measure
the importance of $\boldsymbol{a}$
on a pair of inputs $(\boldsymbol{x}_c, \boldsymbol{x}_p)$

contrastive pair

**Concept**:
Inferring the causal role of the *patched* component (*contrastive input*) in producing the original behavior.

# Attributing causal effect

$$F\big(m; \boldsymbol{a}; \boldsymbol{x}_c, \boldsymbol{x}_p\big) = m(\boldsymbol{x}_c | r(\boldsymbol{a} = \boldsymbol{a}_p)) - m(\boldsymbol{x}_c)$$ 
indirect effect

**Attribution patching**

$$F_a\big(m; \boldsymbol{a}; \boldsymbol{x}_c, \boldsymbol{x}_p\big) = \nabla_{\boldsymbol{a}} m_{|\boldsymbol{a}_c}(\boldsymbol{a}_p - \boldsymbol{a}_c)$$ 
first-order Taylor expansion
(linear approximation)

**Integrated gradient**

$$F_g\big(m; \boldsymbol{a}; \boldsymbol{x}_c, \boldsymbol{x}_p\big) = \left( \sum_{\alpha} \nabla_{\boldsymbol{a}} m_{|\alpha \boldsymbol{a}_c + (1-\alpha)\boldsymbol{a}_p} \right) (\boldsymbol{a}_p - \boldsymbol{a}_c)$$ 
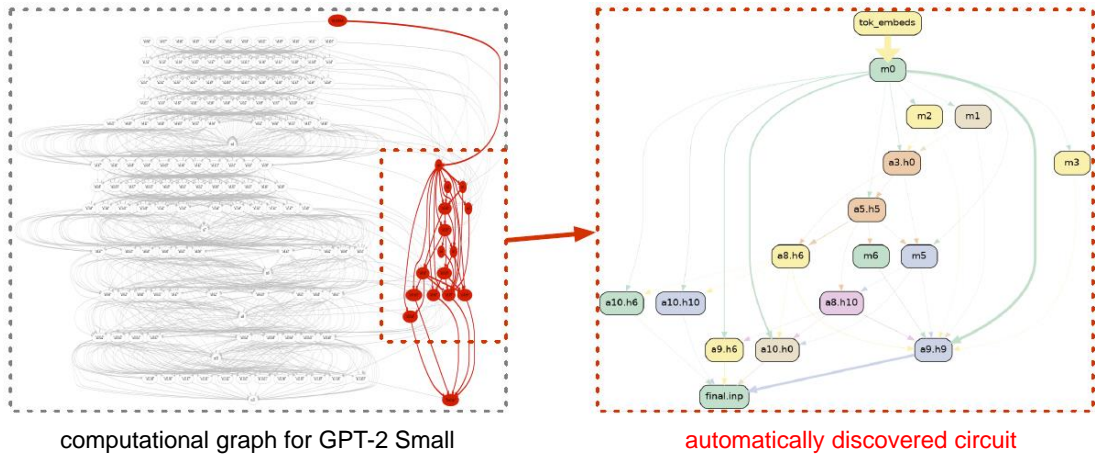more accurate approx.

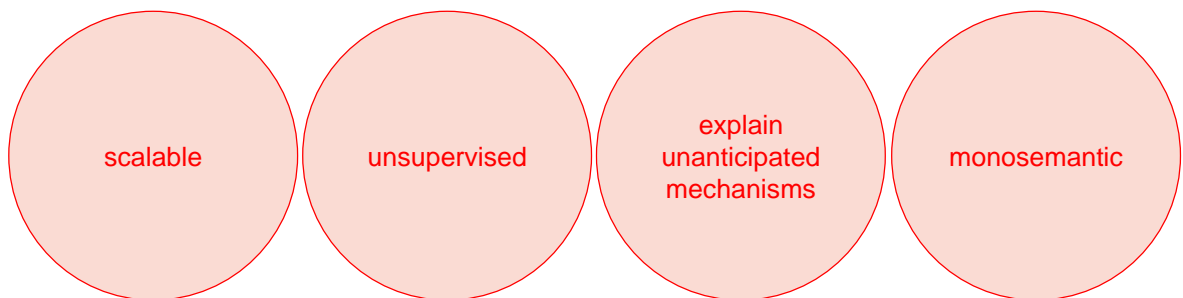$$\alpha \in \left\{ 0, \frac{1}{N}, \dots, \frac{N-1}{N} \right\} \quad N = 10$$

arXiv:2403.19647

# Subgraph with distinct functionality: example

iterative patching experiments to remove unnecessary components and connections



computational graph for GPT-2 Small

automatically discovered circuit

arXiv:2304.14997

# Mechanistic interpretability



scalable

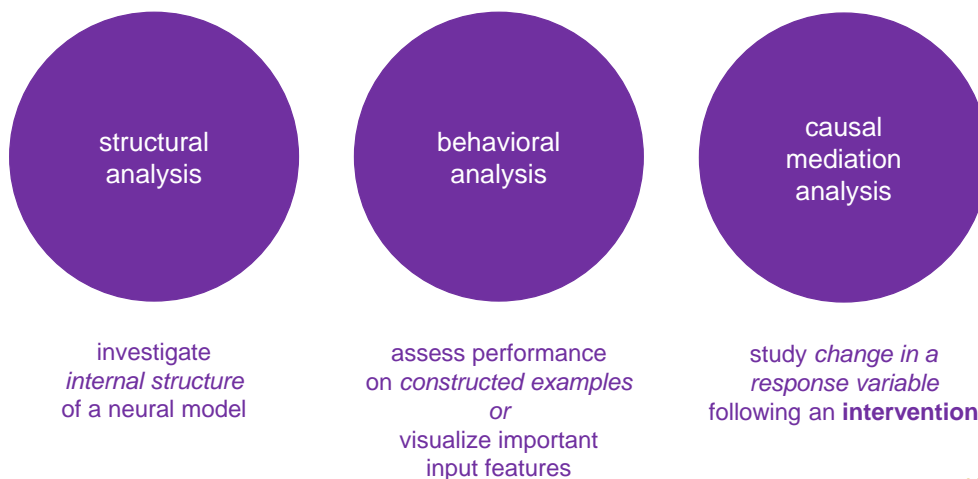unsupervised

explain unanticipated mechanisms

monosemantic

**Concept**:
Features and decision paths that are *uniquely responsible* for a decision.

# Causal mediation analysis

## Interpreting neural models

structural analysis

behavioral analysis

causal mediation analysis

investigate *internal structure* of a neural model

assess performance on *constructed examples* *or* visualize important input features

study *change in a response variable* following an **intervention**

arXiv:2004.12265

# Indirect effect of mediators



causal graph

the mediator decouples the total effect
into direct & indirect effect

# Structural-behavioral analysis: example

causal mediation analysis yields insights
on the role of model components in mediating gender bias



**structural analysis**

**behavioral analysis**

highlight internal
model components
responsible for gender bias

components causally
implicated in how gender bias
manifests in the model outputs

# Structural-behavioral analysis

$p_\theta(x_t | x_1, \ldots, x_{t-1})$  pre-trained language model

$\boldsymbol{h}_{l,i} \in \mathbb{R}^K$  contextual representation of word $i$ in layer $l$

$\boldsymbol{h}_{l,i,k}$  $1 \leq k \leq K$  neural activations

$\alpha_{l,h,i,j} \geq 0$  attention directed from word $i$ to word $j$ by head $h$ in layer $l$

$$\sum_j \alpha_{l,h,i,j} = 1$$

arXiv:2004.12265

# Measure of gender bias

prompt $\boldsymbol{x}$: *The nurse said that …*

stereotypical candidate: *she*

anti-stereotypical candidate: *he*

$$p_\theta(she|\boldsymbol{x}) > p_\theta(he|\boldsymbol{x})$$

**societal bias** associating *nurses* with *women* more than *men*

measure of grammatical gender bias in the model

$$y(\boldsymbol{x}) = \frac{p_\theta(\text{antistereotypical} | \boldsymbol{x})}{p_\theta(\text{stereotypical} | \boldsymbol{x})}$$

$$y(\boldsymbol{x}) = \frac{p_\theta(\text{he} | \textit{The nurse said that})}{p_\theta(\text{she} | \textit{The nurse said that})}$$
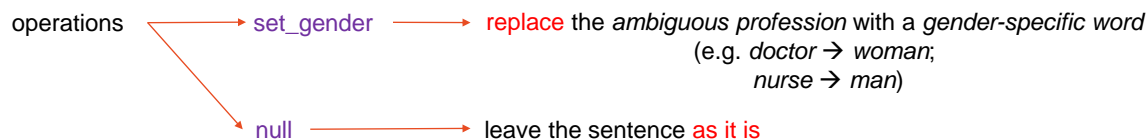
perfectly **unbiased model**  $y(\boldsymbol{x}) = 1$

arXiv:2004.12265

# Understanding the role of components

targeted interventions on the input text ⟶ measure the effect on gender bias

operations → set_gender ⟶ replace the *ambiguous profession* with a *gender-specific word*
(e.g. *doctor* → *woman*;
*nurse* → *man*)

→ null ⟶ leave the sentence as it is

population of prompts $\quad y_o(x) \quad$ for operation (intervention) o

arXiv:2004.12265

# Unit-level total effect

**Total Effect**
(of the intervention)

$$TE(\text{set\_gender}, \text{null}; y, x) = \frac{y_{\text{set\_gender}}(x) - y_{\text{null}}(x)}{y_{\text{null}}(x)} = \frac{y_{\text{set\_gender}}(x)}{y_{\text{null}}(x)} - 1$$

**Average Total Effect**

$$TE(\text{set\_gender}, \text{null}; y) = \mathbb{E}_x \left[ \frac{y_{\text{set\_gender}}(x)}{y_{\text{null}}(x)} - 1 \right] \quad \textit{expectation over the population}$$

arXiv:2004.12265

# Example

compute relative probabilities of the baseline

$p_\theta(he|x) = p_\theta(he|\textit{The nurse said that}) \approx 3.1\%$

$p_\theta(she|x) = p_\theta(she|\textit{The nurse said that}) \approx 22.4\%$

$y_{\text{null}}(x) = 3.1/22.3 \approx 0.14$

set $x$ to an anti-stereotypical case

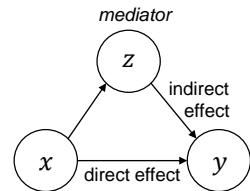$p_\theta(he|x, \text{set\_gender}) = p_\theta(he|\textit{The \textbf{man} said that}) \approx 31.5\%$

$p_\theta(she|x, \text{set\_gender}) = p_\theta(she|\textit{The \textbf{man} said that}) \approx 2.4\%$

$y_{\text{set\_gender}}(x) = 31.5/2.4 \approx 13.1$

**Total Effect**
(of the intervention)

$$TE(\text{null}, \text{set\_gender}, y, x) = \frac{13.1}{0.14} - 1 \approx 92.57$$

arXiv:2004.12265

# Natural direct and indirect effect



*mediator*

$z$

indirect effect

$x$   direct effect   $y$

**Natural Direct Effect**

$$NDE(\text{set\_gender}, \text{null}; y) = \mathbb{E}_x \left[ \frac{y_{\text{set\_gender}, z_{null}(x)}(x)}{y_{\text{null}}(x)} - 1 \right]$$

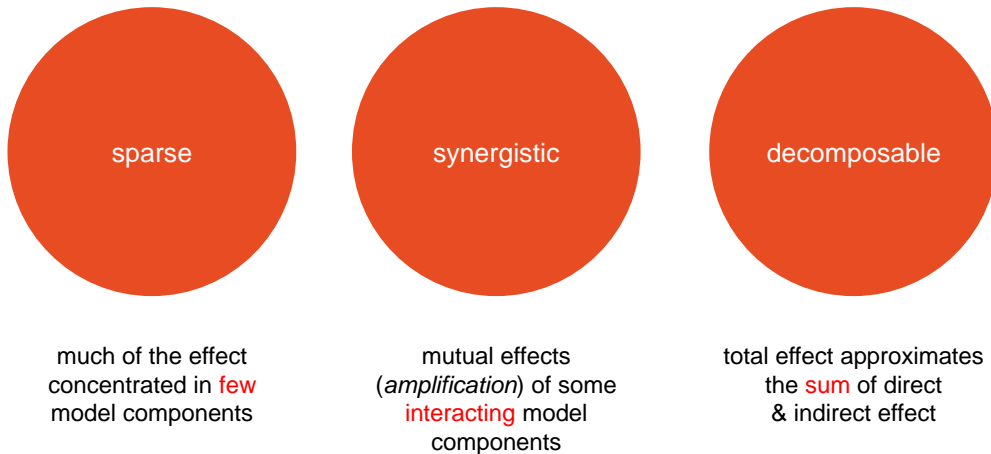measures **direct effect** on gender bias that does not pass through mediator

**Natural Indirect Effect**

$$NIE(\text{set\_gender}, \text{null}; y) = \mathbb{E}_x \left[ \frac{y_{\text{null}, z_{\text{set\_gender}}(x)}(x)}{y_{\text{null}}(x)} - 1 \right]$$

measures **indirect effect** flowing from $x$ to $y$ through mediator $z$

arXiv:2004.12265

# Model's sensitivity to grammatical gender

sparse

synergistic

decomposable

much of the effect concentrated in few model components

mutual effects (*amplification*) of some interacting model components

total effect approximates the sum of direct & indirect effect

arXiv:2004.12265

# **Datasheets for datasets**

# Documenting datasets

| | | | |
|---|---|---|---|
| **motivation** | **composition** | **collection process** | **preprocessing, cleaning, labeling** |
| **uses** | **distribution** | **maintenance** | |

transparency & accountability

# Motivation

- For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled?

- Who created the dataset and on behalf of which entity?

- Who funded the creation of the dataset?

# Composition

- What do the instances that comprise the dataset represent (e.g. documents, photos, people, countries)? Are there multiple types of instances (e.g. movies, users, and ratings)?

- How many instances are there in total (of each type, if appropriate)?

- Does the dataset contain all possible instances or is it a sample from a larger set?

- What data does each instance consist of? Raw data or features?

- Is there a label or target associated with each instance?

- Is any information missing from individual instances?

arXiv:1803.09010

# Composition

- Are relationships between individual instances made explicit (e.g. social network links)?

- Are there recommended data splits (e.g. training, development/validation, testing)?

- Are there any errors, sources of noise, or redundancies in the dataset?

- Is the dataset self-contained, or does it link to or otherwise rely on external resources?

- Does the dataset contain data that might be considered confidential?

- Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?

arXiv:1803.09010

# Composition (people)

- Does the dataset identify any subpopulations (e.g. by age, gender)?

- Is it possible to identify individuals (i.e. one or more natural persons), either directly or indirectly (i.e. in combination with other data) from the dataset?

- Does the dataset contain data that might be considered sensitive in any way (e.g. data that reveals race or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)?

arXiv:1803.09010

# Collection process

- How was the data associated with each instance acquired?

- If the data was reported by subjects or indirectly inferred from other data, was the data validated?

- What procedures were used to collect the data? How were these procedures validated?

- If the dataset is a sample from a larger set, what was the sampling strategy?

- Who was involved in the data collection process and how were they compensated?

- Over what timeframe was the data collected?

arXiv:1803.09010

# Collection process (people)

- Were any ethical review processes conducted?

- Did you collect the data from the individuals in question directly, or obtain it via third parties?

- Were the individuals in question notified about the data collection? Did the individuals in question consent to the collection and use of their data?

- If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses?

- Has an analysis of the potential impact of the dataset and its use on data subjects (e.g. a data protection impact analysis) been conducted?

arXiv:1803.09010

# Preprocessing, cleaning, labeling

- Was any preprocessing, cleaning and/or labeling of the data done (e.g. discretization or bucketing, tokenization, part-of-speech tagging, feature extraction, removal of instances, processing of missing values)?

- Was the raw data saved in addition to the preprocessed/cleaned/labelled data (e.g. to support unanticipated future uses)?

- Is the software that was used to preprocess, clean and/or label the data available?

arXiv:1803.09010

# Uses

- Has the dataset been used for any tasks already?

- Is there a <span style="color:red">repository</span> that links to any or all papers or systems that use the dataset?

- What <span style="color:red">(other) tasks</span> could the dataset be used for?

- Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might <span style="color:red">impact future uses</span>?
  E.g., is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other risks or harms (e.g., legal risks, financial harms)? Is there anything a dataset consumer could do to mitigate these risks or harms?

- Are there tasks for which the dataset should <span style="color:red">not be used</span>?

arXiv:1803.09010

# Distribution

- Will the dataset be <span style="color:red">distributed</span> to third parties?

- How will the dataset will be distributed? Does the dataset have a digital object identifier (<span style="color:red">DOI</span>)? When will the dataset be distributed?

- Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable <span style="color:red">terms of use</span>?

- Have any third parties imposed IP-based or other <span style="color:red">restrictions</span> on the data associated with the instances?

- Do any export controls or <span style="color:red">regulatory restrictions</span> apply to the dataset or to individual instances?

arXiv:1803.09010

# Maintenance

- Who will be supporting, hosting and maintaining the dataset? How can the owner, curator, manager of the dataset be contacted?

- Is there an erratum? Will the dataset be updated (e.g. to correct labeling errors, delete/add instances)? Will older versions of the dataset continue to be supported/hosted/maintained?

- If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances?

- If others want to extend, augment, build on, contribute to the dataset, is there a mechanism for them to do so?

arXiv:1803.09010

# What did we learn today?

- Alignment with diverse preferences
- Reinforcement learning from AI feedback
- Interpretability
- Causal mediation analysis
- Datasheets for datasets

# **EE-559**
# **Deep Learning**

andrea.cavallaro@epfl.ch