

**EPFL**

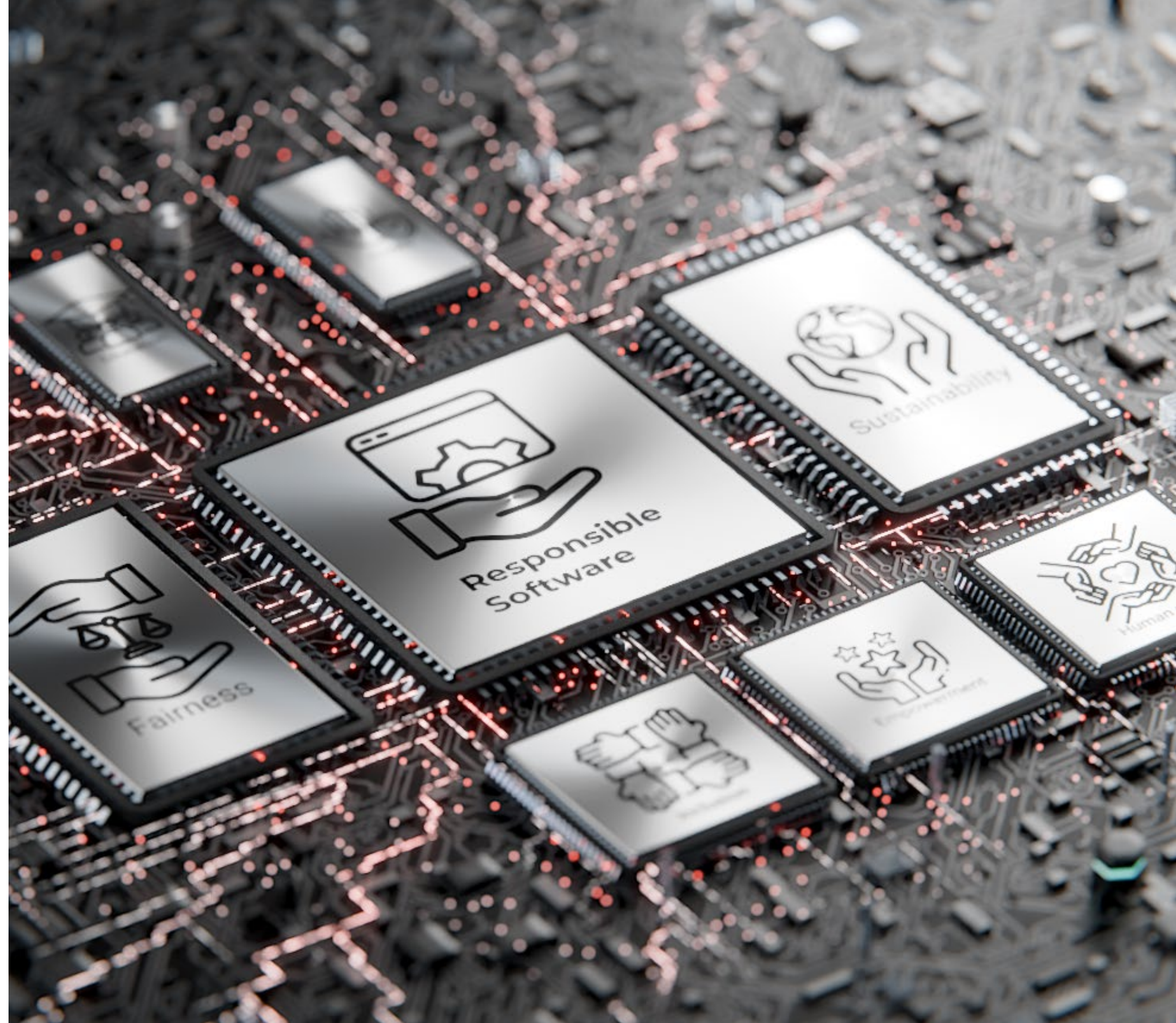
# **Introduction**

## **Session**

**15 sept.**

Cécile Hardebolle

**Responsible  
Software**



# Responsibility

---

In this course, we consider that being responsible as a software engineer means:

- a. Making sure a liability clause is included in the software license agreement.
- b. Reacting rapidly to correct software bugs when they are reported.
- c. Being accountable for the decisions made by the development team.
- d. Anticipating the potential negative impacts of the software on others.

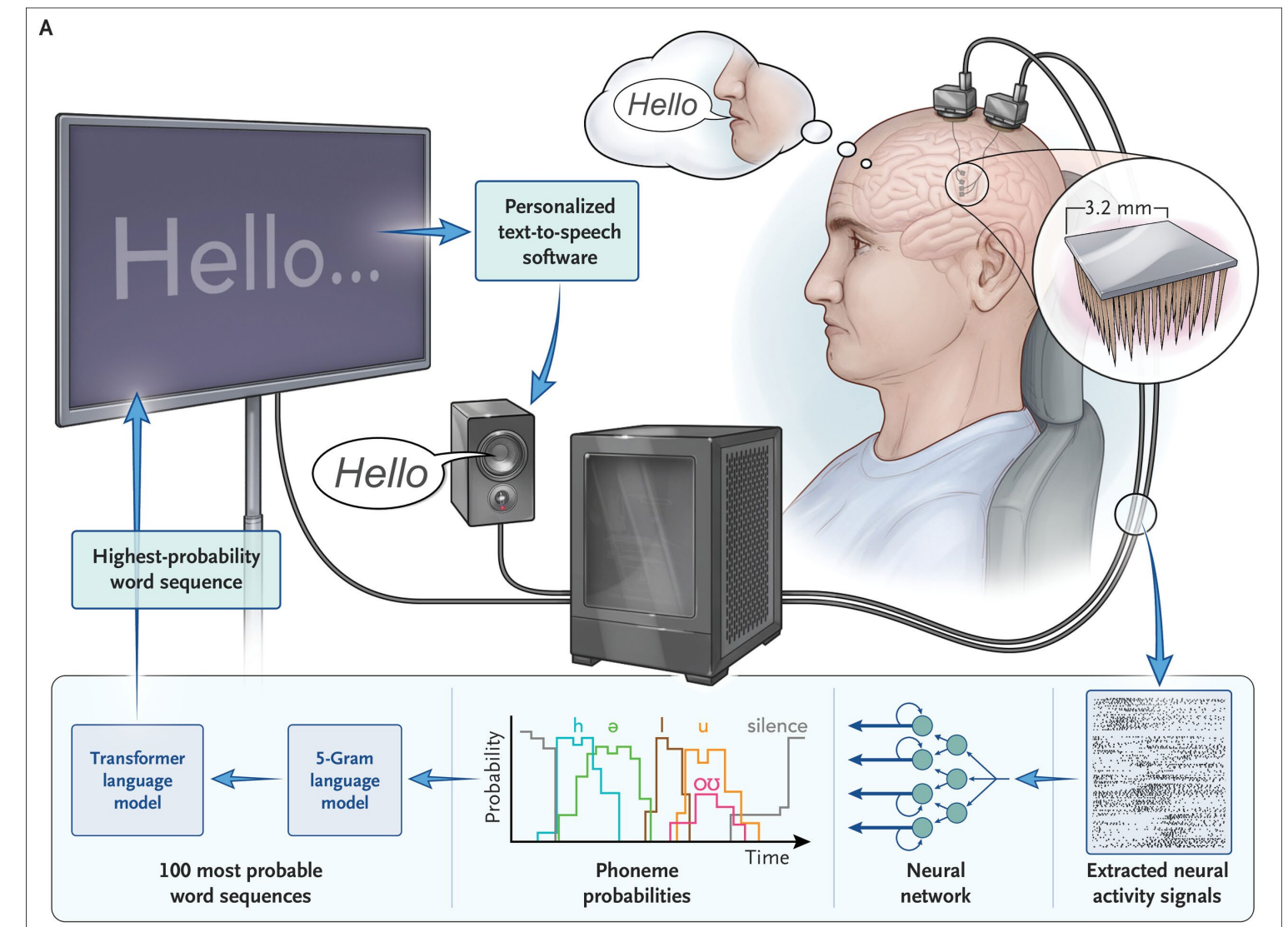
# Types of issues (1/2)

(Card et al., 2024)

Brain-to-speech software can translate neural activity associated with attempted speech into spoken words. A key challenge is ensuring that only intentional communication is captured, not private inner thoughts.

This is:

- A technical issue
- An ethical issue
- An ethical dilemma



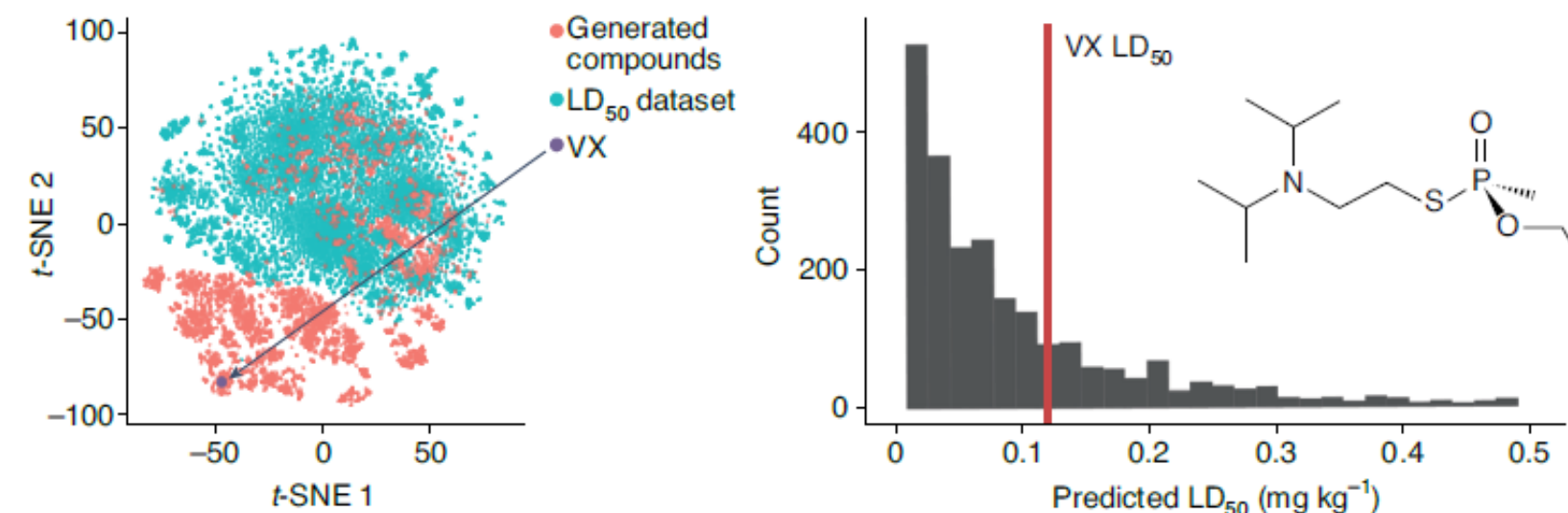
# Types of issues (2/2)

(Urbina et al., 2022)

A software company has developed a Machine Learning model that is able to discover new chemical compounds for medicine development. They identify that the model can also discover new chemical weapons.

This is:

- A technical issue
- An ethical issue
- An ethical dilemma



**Fig. 1 |** A t-SNE plot visualization of the LD<sub>50</sub> dataset and top 2,000 MegaSyn AI-generated and predicted toxic molecules illustrating VX. Many of the molecules generated are predicted to be more toxic in vivo in the animal model than VX (histogram at right shows cut-off for VX LD<sub>50</sub>). The 2D chemical structure of VX is shown on the right.

# Normative ethical theories (1/2)

---

A software engineer refuses to hide a critical bug in a released product and tells you: “I believe that it is always wrong to lie, even if telling the truth might result in harm to some people.”  
Which ethical theory does this engineer follow?

- a. Utilitarianism
- b. Deontology
- c. Virtue
- d. Care

# Normative ethical theories (2/2)

---

A software engineer refuses to hide a critical bug in a released product and tells you: “If I do not report this bug, I am not being a trustworthy and courageous person.” Which ethical theory does this engineer follow?

- a. Utilitarianism
- b. Deontology
- c. Virtue
- d. Care

# Stakeholders analysis

---

Which of the following statements are true about stakeholders?

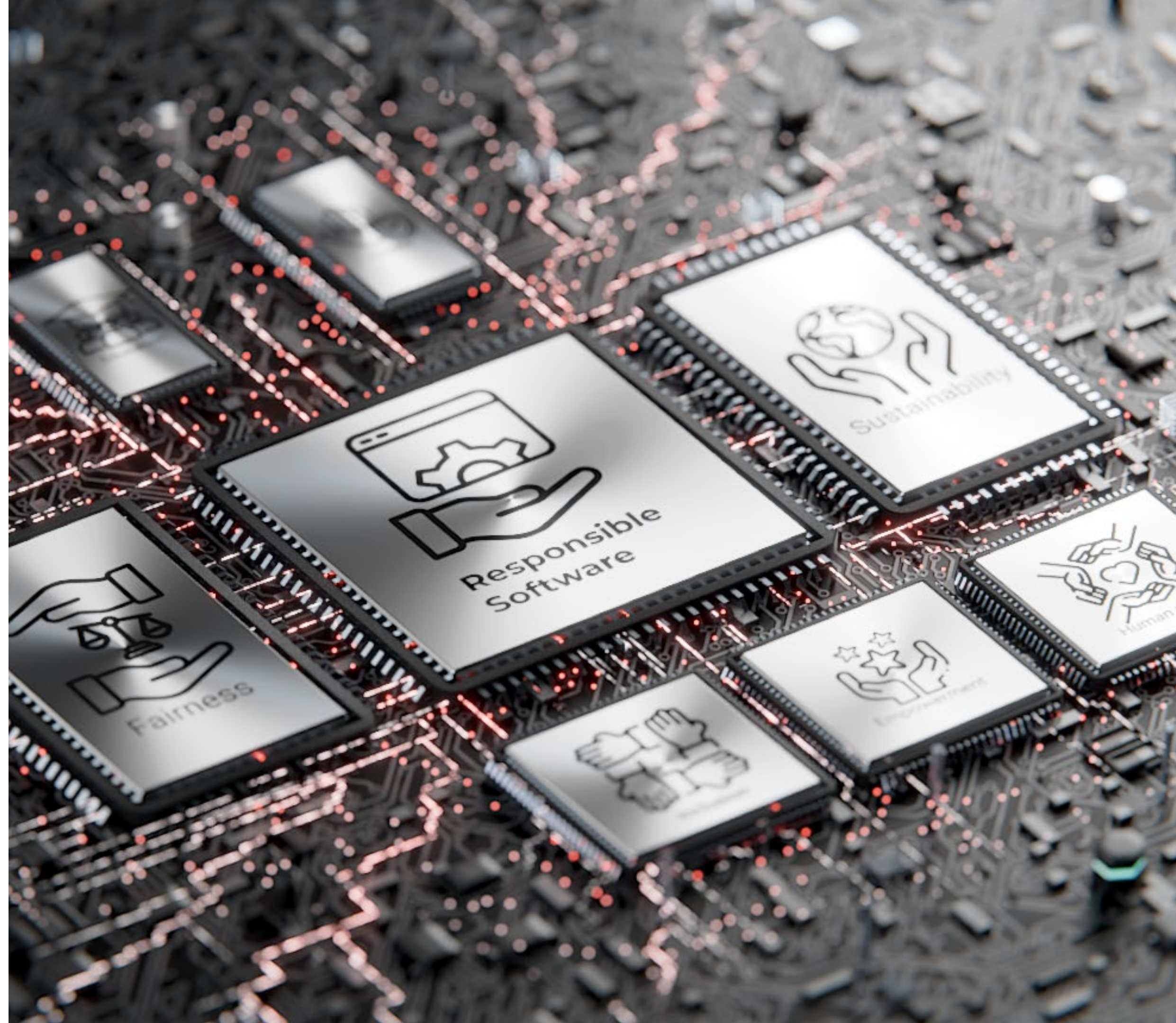
- a. Can be persons
- b. Can be non-humans
- c. Can be affected positively
- d. Can be affected negatively
- e. Are in contact with the software
- f. Do not interact with the software but are affected by it

**EPFL**

**Safety 1  
Review &  
Case Studies  
22 sept.**

Cécile Hardebolle

**Responsible  
Software**



# Autonomous car software - 1

---

The software of an autonomous car has a 10% error rate in recognizing traffic signs correctly.

We are in the presence of (select all that apply):

- 0% a. A safety threat
- 0% b. A security threat
- 0% c. A safety hazard
- 0% d. A security hazard

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

# Autonomous car software - 2

---

Stickers placed on a stop sign lead the software of an autonomous car to misclassify it as a speed limit sign.

We are in the presence of (select all that apply):

- 0% a. A safety threat
- 0% b. A security threat
- 0% c. A safety hazard
- 0% d. A security hazard

URL: ttpoll.eu

Session ID: cs290

# Worldwide “CrowdStrike” outage in 2024

This event is an example of:

- 0% a. Malfunction
- 0% b. Misuse, abuse
- 0% c. Unintended use
- 0% d. Intended use

URL: [ttpoll.eu](https://ttpoll.eu)  
Session ID: cs290

CrowdStrike IT outage affected 8.5 million Windows devices, Microsoft says

20 July 2024

Share ↵ Save +

Joe Tidy  
Cyber correspondent, BBC News



The New York Times

## Stranded in the CrowdStrike Meltdown: ‘No Hotel, No Food, No Assistance’

Airlines pledged assistance, refunds and reimbursements to passengers whose travel had been disrupted by this summer’s software outage. Instead, passengers told us, they were on their own.

# Bad actors, safety and security

---

- 0% a. Bad actors generate safety issues only
- 0% b. Bad actors generate security issues only
- 0% c. Bad actors generate both security and safety issues

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

# Bad actors and the 4 scenarios

---

Bad actors can be involved in (select all that apply):

- 0% a. Malfunction
- 0% b. Misuse, abuse
- 0% c. Unintended use
- 0% d. Intended use

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

# The “confusing” matrix - 1

---

We use software to detect fissures in concrete walls before they become visible to the naked eye.

A positive result means presence of fissure.

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

Select all the correct statements:

- 0% a. TN = actual absence of fissure, correct prediction
- 0% b. TP = actual absence of fissure, correct prediction
- 0% c. FN = actual presence of fissure, incorrect prediction
- 0% d. FP = actual presence of fissure, incorrect prediction

# The “confusing” matrix - 2

---

We use software to detect fissures in concrete walls before they become visible to the naked eye.

A positive result means presence of fissure.

From a safety perspective, the indicator we should pay most attention to is:

0% a. TN

0% b. TP

0% c. FN

0% d. FP

URL: ttpoll.eu  
Session ID: cs290

# The "confusing" matrix - 3

We use software to detect fissures in concrete walls before they become visible to the naked eye.

A positive result means presence of fissure.

Here is the confusion matrix you get 🙌

What is the False Negative Rate (FNR)?

- 0% a. 13%
- 0% b. 17%
- 0% c. 20%
- 0% d. 25%

		Predicted	
		Fissure	No Fissure
Actual	Fissure	60	15
	No Fissure	20	100

# Harm categories - 1

---

A user sees their post unfairly censored.  
This harm is in the category (select one):

- 0% a. Physical injury
- 0% b. Emotional or psychological injury
- 0% c. Opportunity loss
- 0% d. Economic loss
- 0% e. Dignity loss
- 0% f. Liberty loss
- 0% g. Privacy loss
- 0% h. Environmental impact
- 0% i. Manipulation
- 0% j. Social detriment

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

# Harm categories - 2

---

A fitness app leaks GPS location data on social media.

This harm is in the category (select one):

- 0% a. Physical injury
- 0% b. Emotional or psychological injury
- 0% c. Opportunity loss
- 0% d. Economic loss
- 0% e. Dignity loss
- 0% f. Liberty loss
- 0% g. Privacy loss
- 0% h. Environmental impact
- 0% i. Manipulation
- 0% j. Social detriment

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

# Harm categories - 3

---

Online ads lead a compulsive shopper to additional purchases.

This harm is in the category (select one):

- 0% a. Physical injury
- 0% b. Emotional or psychological injury
- 0% c. Opportunity loss
- 0% d. Economic loss
- 0% e. Dignity loss
- 0% f. Liberty loss
- 0% g. Privacy loss
- 0% h. Environmental impact
- 0% i. Manipulation
- 0% j. Social detriment

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

# Harm categories - 4

---

A recruitment software indirectly discriminates based on people's name.

This harm is in the category (select one):

- 0% a. Physical injury
- 0% b. Emotional or psychological injury
- 0% c. Opportunity loss
- 0% d. Economic loss
- 0% e. Dignity loss
- 0% f. Liberty loss
- 0% g. Privacy loss
- 0% h. Environmental impact
- 0% i. Manipulation
- 0% j. Social detriment

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

# Harm categories - 5

---

The results of an image search engine for “Nurse” show only women.  
This harm is in the category (select one):

- 0% a. Physical injury
- 0% b. Emotional or psychological injury
- 0% c. Opportunity loss
- 0% d. Economic loss
- 0% e. Dignity loss
- 0% f. Liberty loss
- 0% g. Privacy loss
- 0% h. Environmental impact
- 0% i. Manipulation
- 0% j. Social detriment

URL: [ttpoll.eu](http://ttpoll.eu)

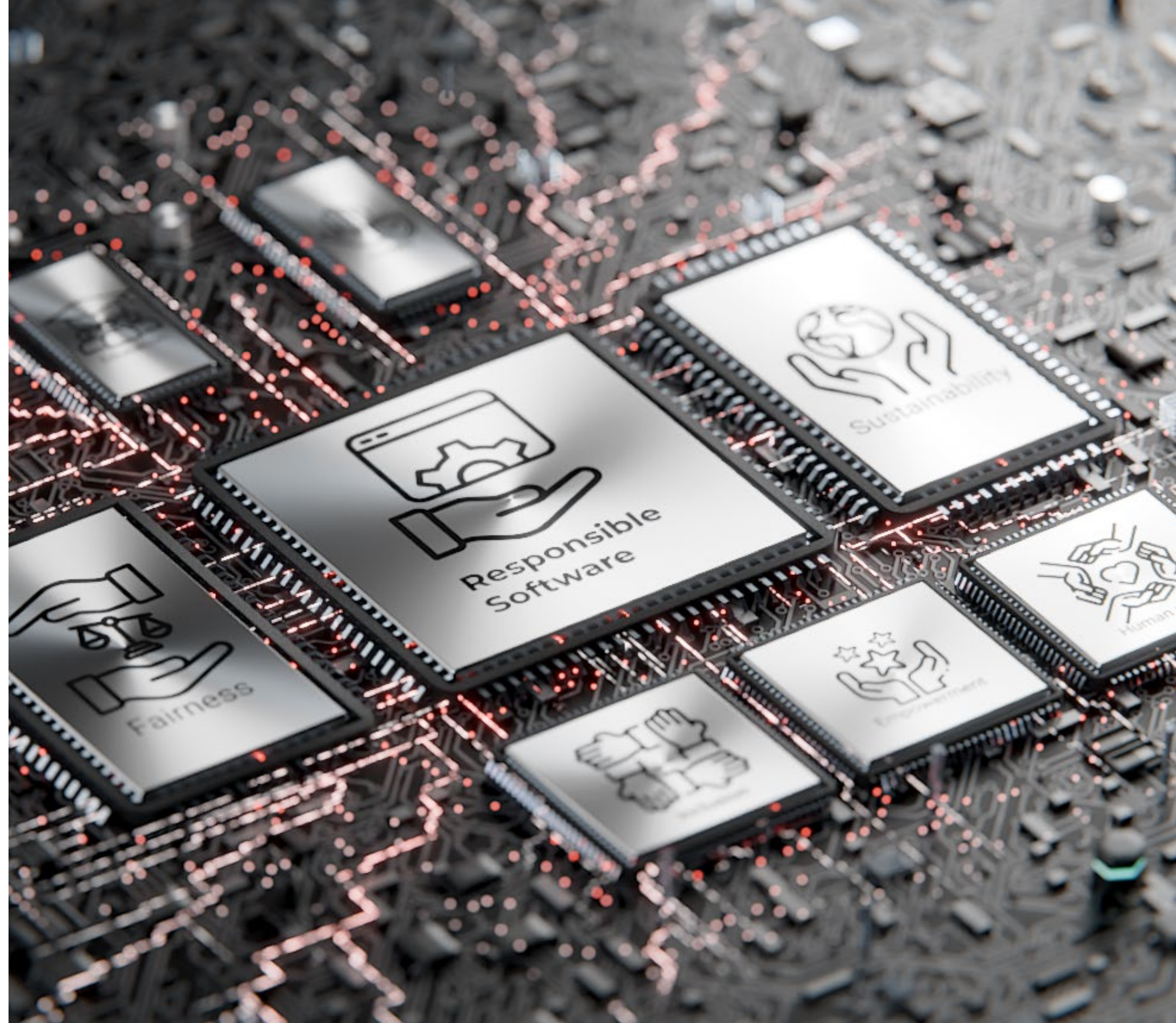
Session ID: cs290

**EPFL**

**Safety 2  
Review &  
Case studies  
29 sept.**

Cécile Hardebolle

**Responsible  
Software**



# Macro-level perspective

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

A macro-level perspective is useful (select all correct statements):

- a. When software is under design
- b. After software is deployed
- c. After an analysis with a meso-level perspective
- d. When considering expanding to new countries
- e. When software is used by public institutions

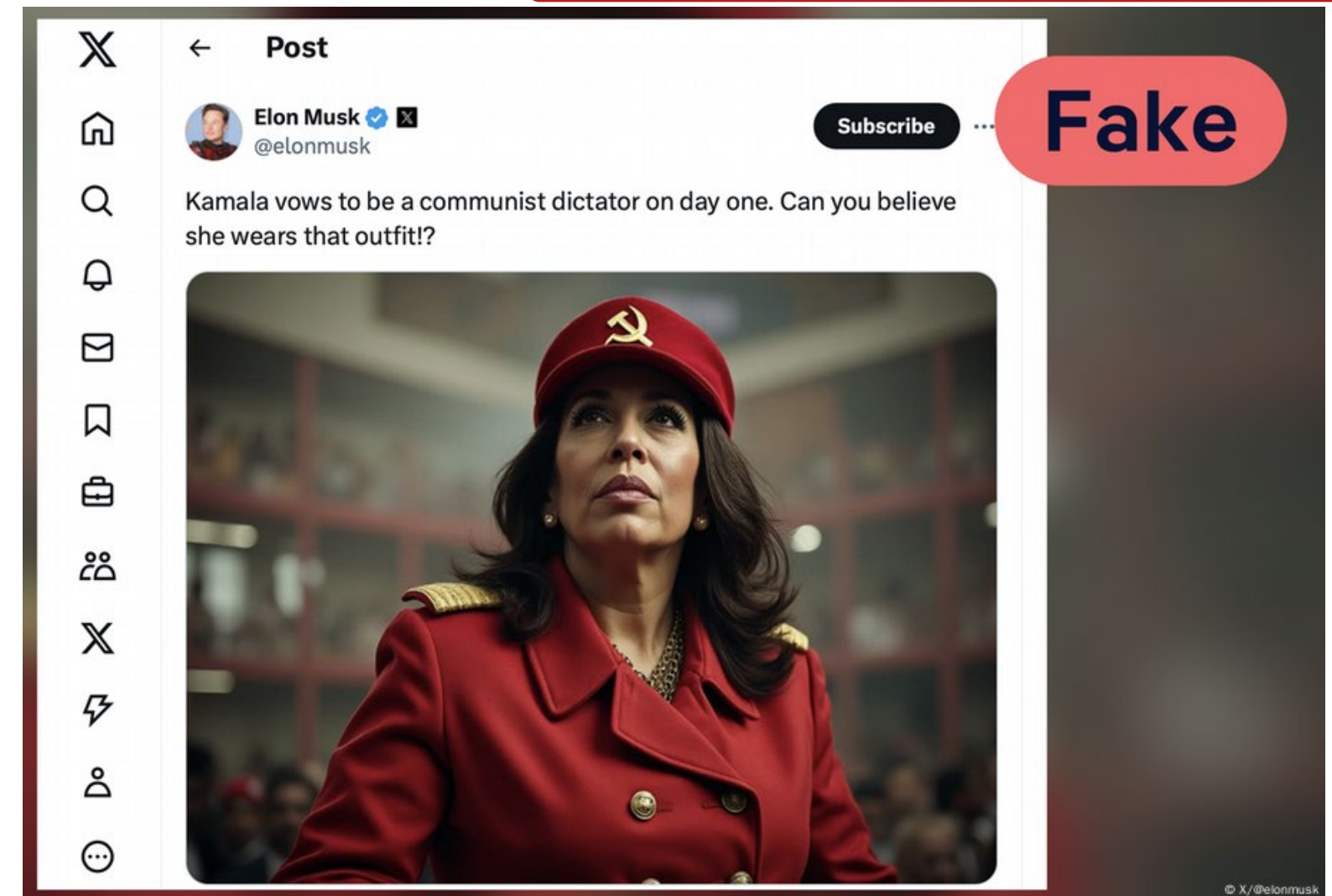
# False beliefs

URL: ttpoll.eu  
Session ID: cs290

One dis-/mis-information post by Elon Musk appears in your Twitter timeline.

You are more likely to believe it because of (choose one):

- a. System 2
- b. Illusory truth
- c. Source cues
- d. Prebunking



Fact check: Elon Musk spreads US election lies. (2024, February 11).  
Dw.Com. <https://www.dw.com/en/fact-check-how-elon-musk-is-spreading-us-election-lies/a-70663408>

Exam  
type

# Dis/Mis-information

URL: ttpoll.eu

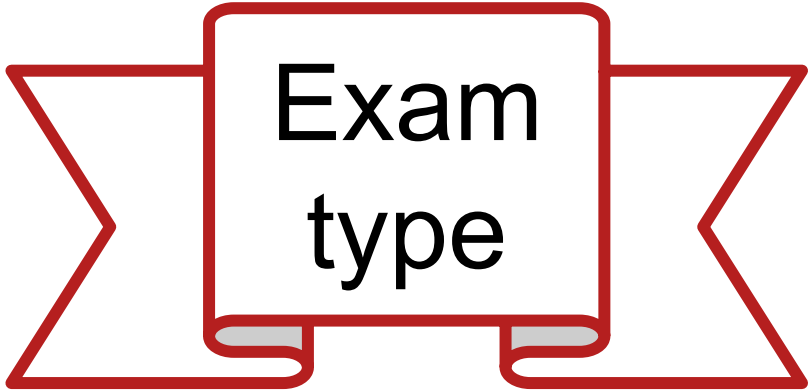
Session ID: cs290

Your friend tells you:

“Eating carrots will drastically improve your night vision.”

This is (choose one):

- a. Misinformation
- b. Disinformation
- c. Malinformation
- d. Fake news



Exam  
type

# Software & disinformation

URL: ttpoll.eu

Session ID: cs290

Software playing a role in disinformation can be (select all that apply):

- a. Generative AI
- b. Bots
- c. Content moderation systems
- d. Content recommendation systems

# Humans & disinformation

URL: [ttpoll.eu](http://ttpoll.eu)

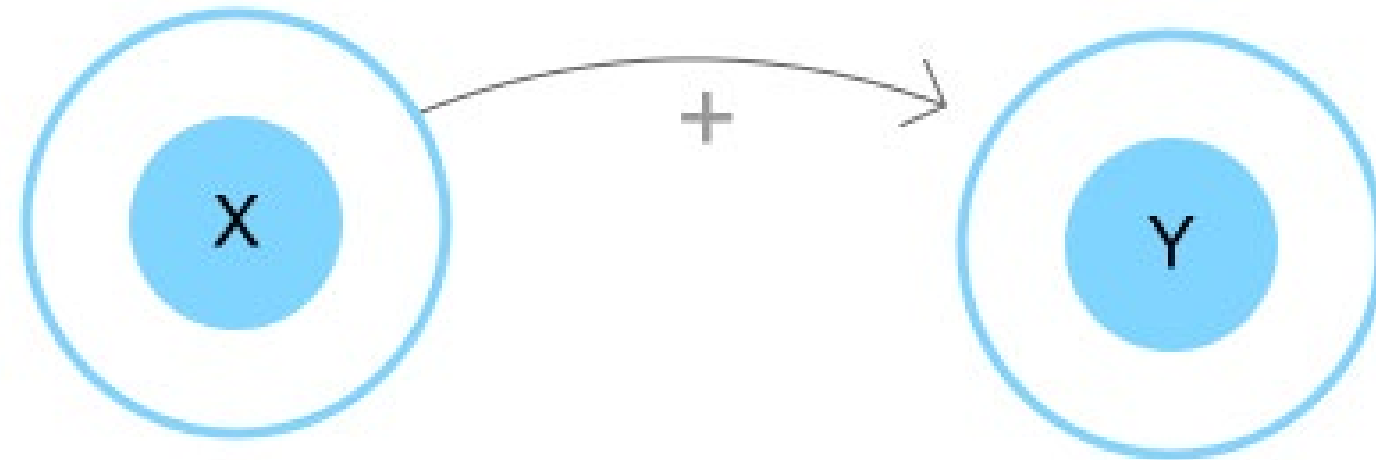
Session ID: cs290

Humans playing a role in disinformation do it (select all that apply):

- a. Because it's their work
- b. By inattention
- c. For political reasons
- d. To please other people
- e. To spark emotions

# Causal Loop Diagrams

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

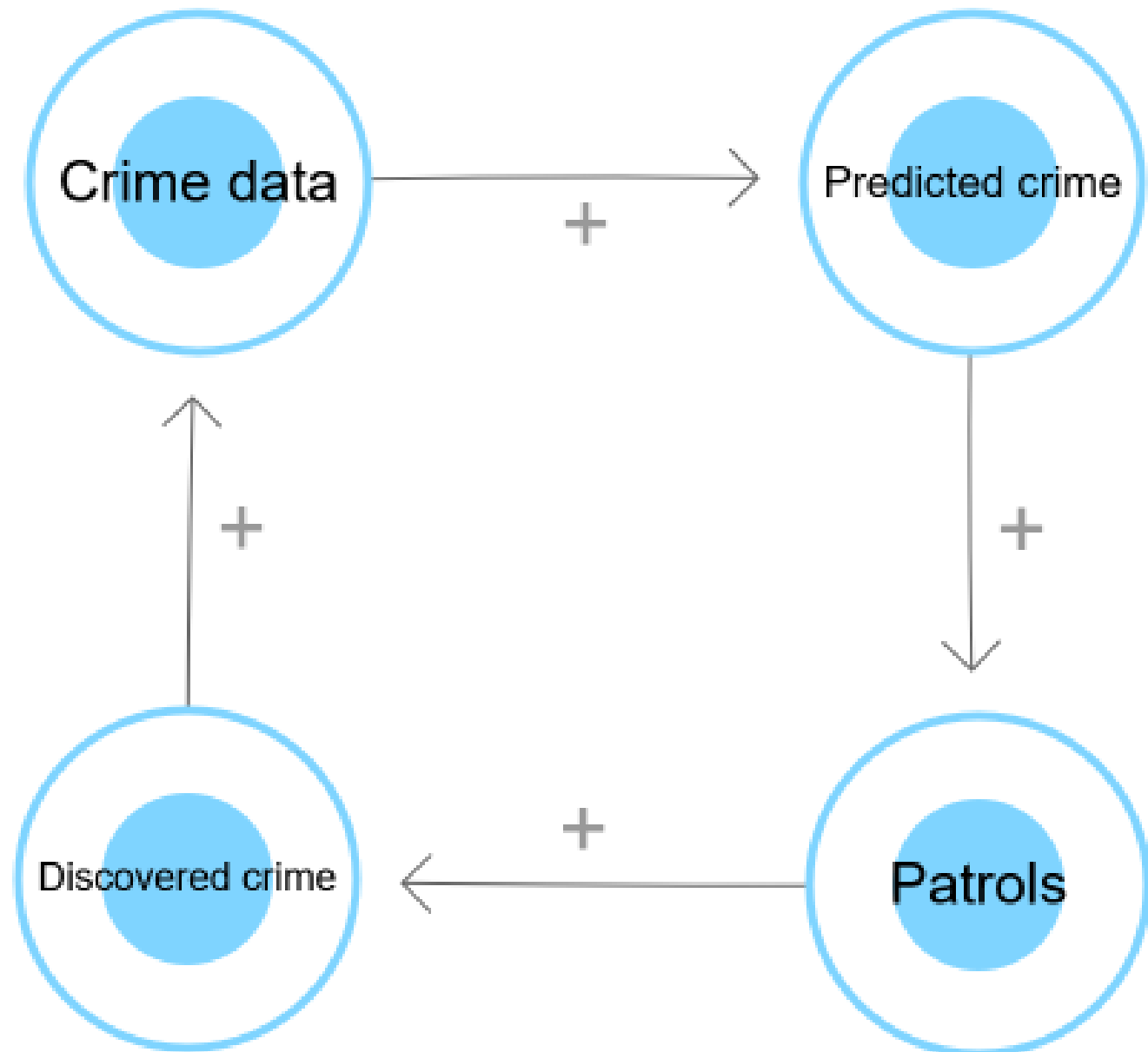


The arrow with label “+” means:

- There's a transition from state X to state Y on token “+”
- The quantity in X is added to the quantity in Y
- X and Y both change in an increasing direction
- Y changes in the same direction as X

# Part 1: behavior

URL: ttpoll.eu  
Session ID: cs290

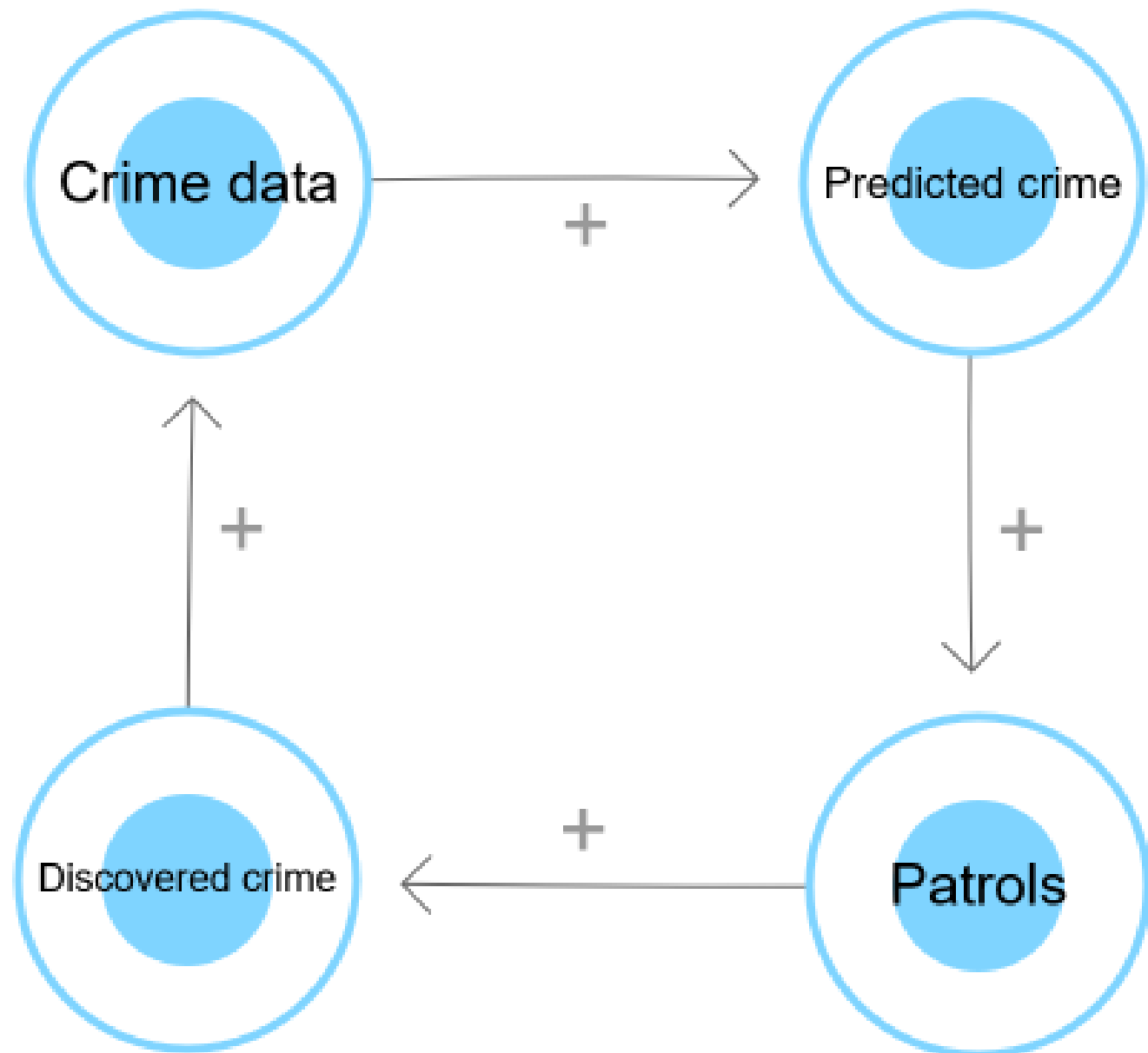


Over time, the quantities in this system will:

- a. Stabilize
- b. Increase
- c. Decrease
- d. It depends

# Part 1: type of feedback loop

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290



The feedback loop in this diagram is:

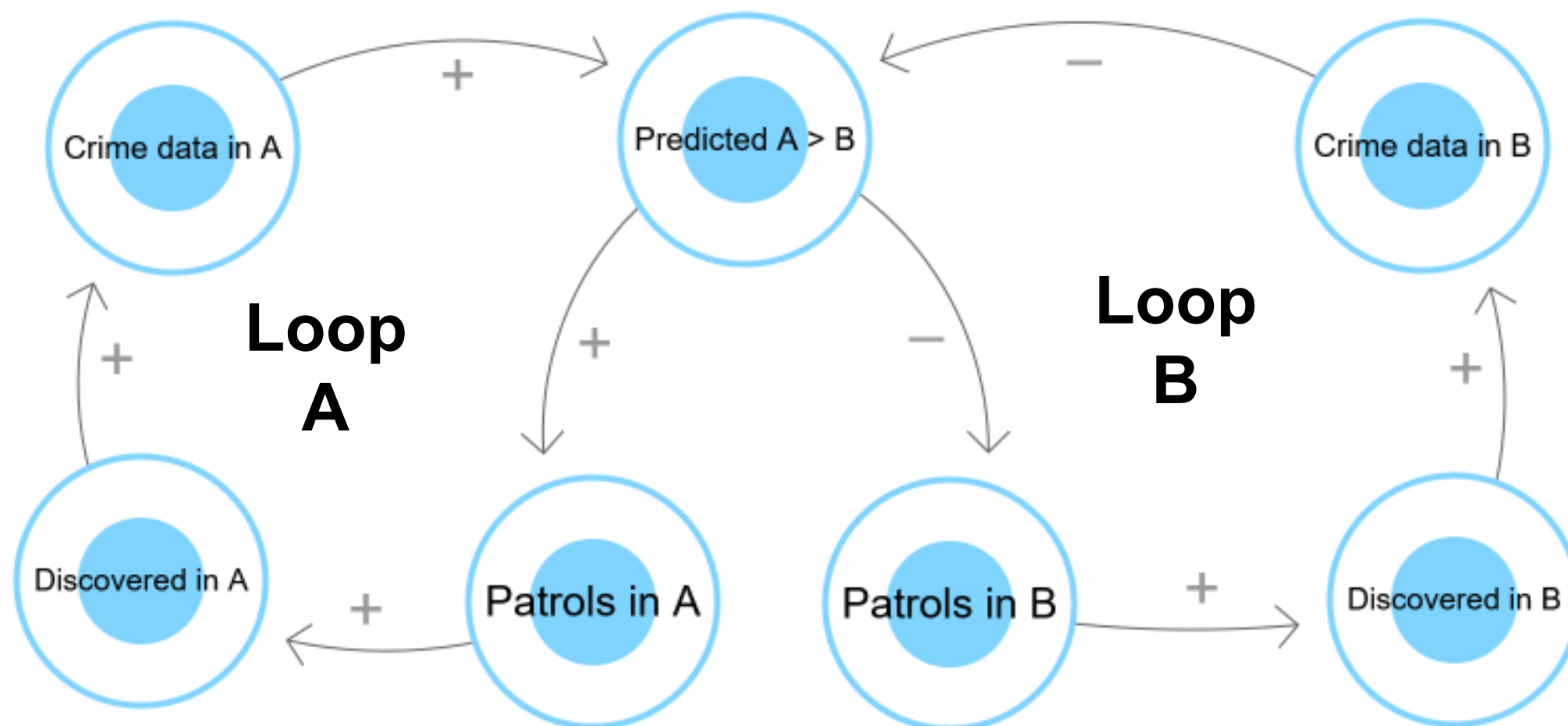
- a. Balancing
- b. Reinforcing

# Part 2: types of feedback loops

URL: ttpoll.eu  
Session ID: cs290

What is the type of loops A and B? (select 2 answers):

- a. Loop A is balancing
- b. Loop A is reinforcing
- c. Loop B is balancing
- d. Loop B is reinforcing

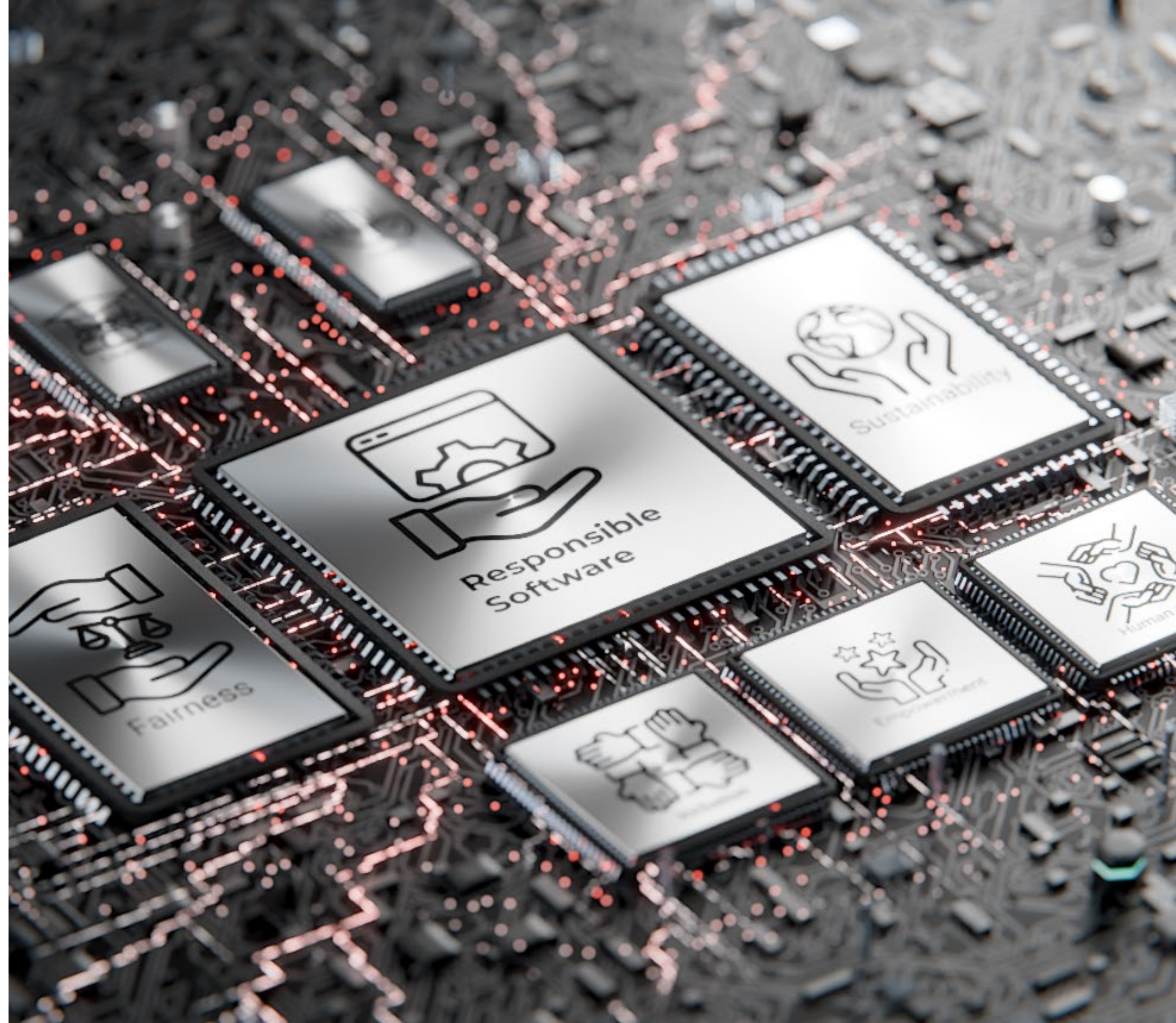


**EPFL**

**Fairness 1  
Review &  
Case studies  
7 oct.**

Cécile Hardebolle

**Responsible  
Software**



# Attributes - 1

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

**What are the characteristics of hair color as an attribute to represent people? (select all that apply)**

- a. Not sensitive
- b. Sensitive
- c. Observed
- d. Latent
- e. Objective
- f. Subjective

# Attributes - 2

URL: ttpoll.eu

Session ID: cs290

Let's imagine a software that relies on SAT scores (standardized test for university admission in the US) to make recommendations of when to approve study loans.

**What are the characteristics of the SAT score?**

- a. Not sensitive
- b. Sensitive
- c. Private
- d. Public
- e. Proxy
- f. System

# Bias - 1

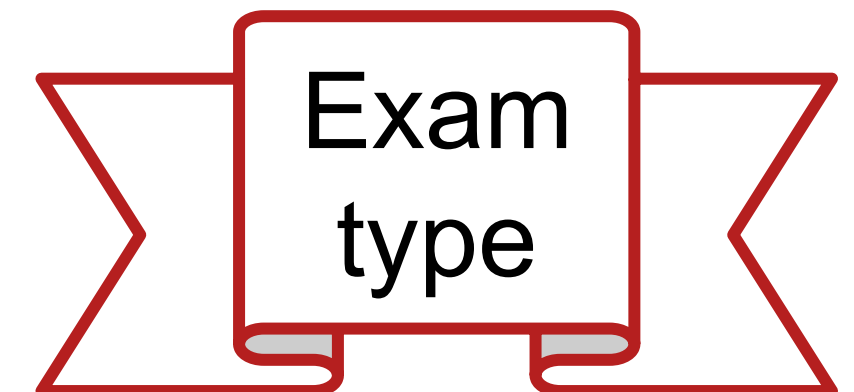
URL: ttpoll.eu

Session ID: cs290

The city of Lozhann decides to deploy a smartphone app that allows residents to report potholes throughout the city to help with the identification of repair needs.

**What bias will the data collected by the app probably exhibit?**  
(select one answer)

- a. Confirmation bias
- b. Representation bias
- c. Measurement bias
- d. Automation bias



# Bias - 2

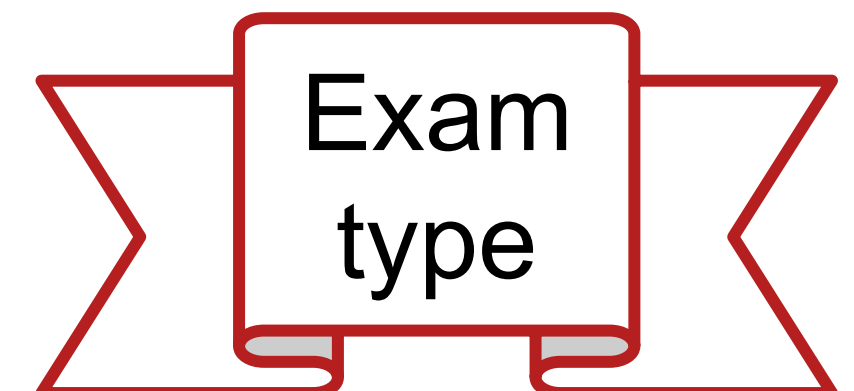
URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

In the new ArcFit fitness tracker, the calory burn feature uses the "metabolic equivalent of task" formula, which estimates the energy a body uses during a specific activity. The same calculation is used during walking and running.

**What type of bias will the calory burn variable probably have?**  
(select one answer)

- a. Confirmation bias
- b. Representation bias
- c. Measurement bias
- d. Automation bias



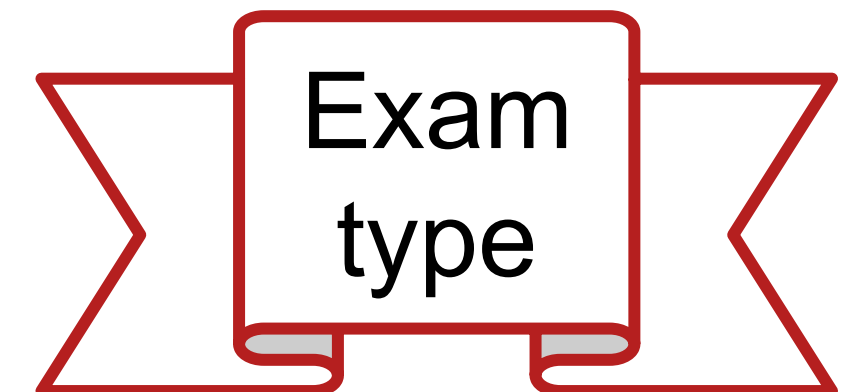
# Bias - 3

URL: ttpoll.eu  
Session ID: cs290

A group of computer scientists with similar background, all experts in software development, are starting a new software project in the healthcare domain.

**What type of bias will these scientists probably have?**  
(select one answer)

- a. Confirmation bias
- b. Automation bias
- c. Pre-existing bias
- d. Sunk cost fallacy



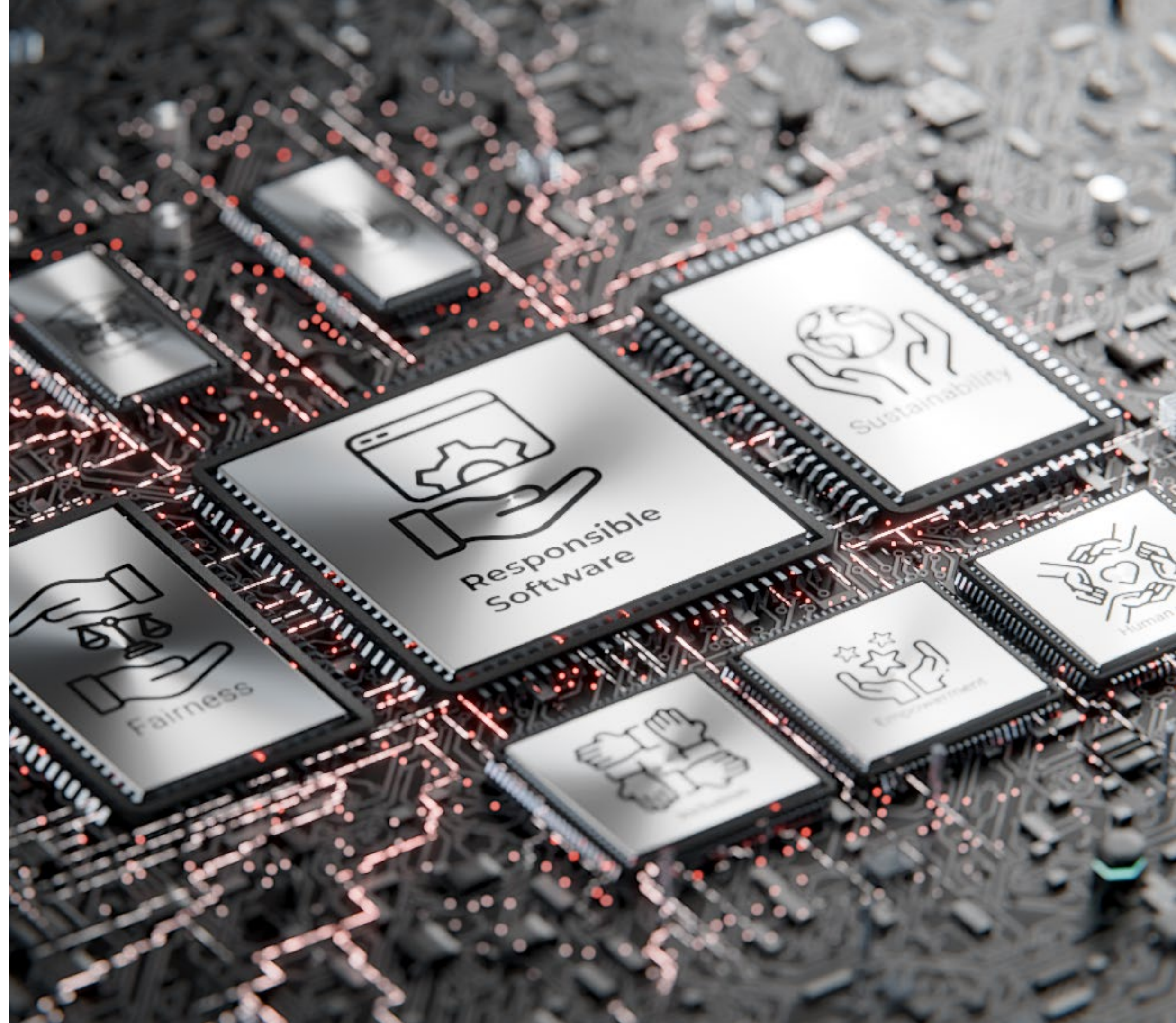
**EPFL**

# **Fairness 2 Review & Case studies**

**13 oct.**

Cécile Hardebolle

**Responsible  
Software**



# Biases in the ML lifecycle - 1

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

Simpson's paradox is when the patterns observed at the level of the full sample and at the level of subgroups are opposed.

**When training a ML model, Simpson's paradox can lead to**  
(select 1 answer):

- a. Evaluation bias
- b. Aggregation bias
- c. Optimization choices
- d. Deployment bias

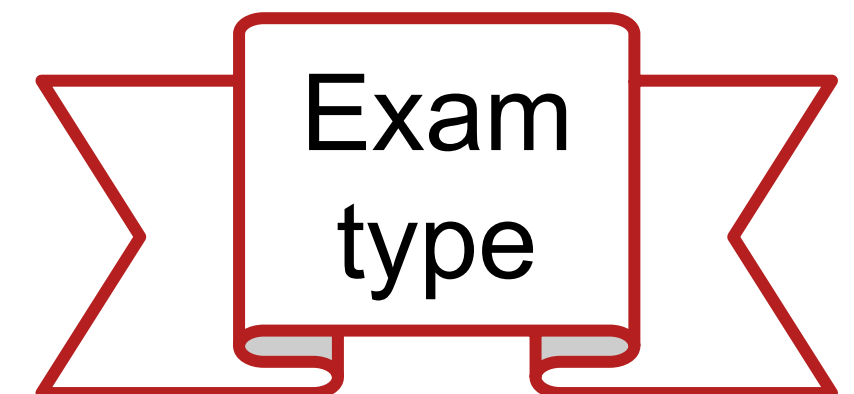
# Biases in the ML lifecycle - 2

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

The society RetailProtect has developed a ML model to identify instances of shoplifting in retail shops. For evaluating their model, they use a benchmark in which actors from diverse ethnicities simulate a range of shoplifting actions.

**This can lead to (select 1 answer):**

- 0% a. Evaluation bias
- 0% b. Aggregation bias
- 0% c. Optimization choices
- 0% d. Deployment bias



# Fairness metrics - 1

URL: [ttpoll.eu](http://ttpoll.eu)

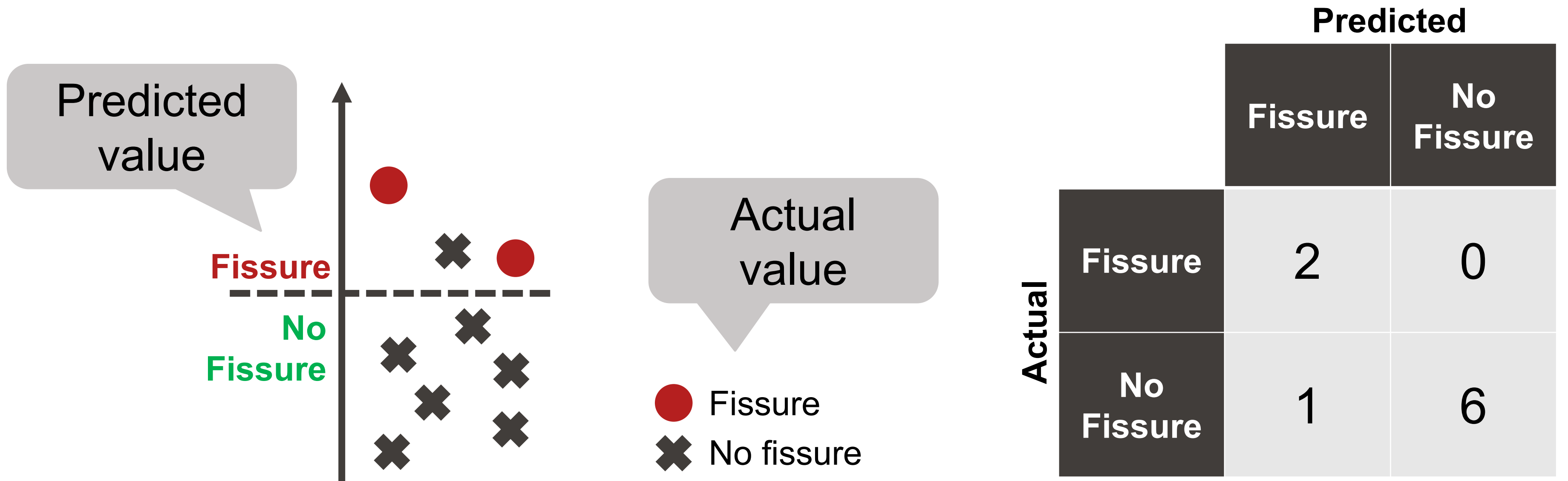
Session ID: cs290

Among the metrics below, **which can be used to assess the fairness** of a piece of software? (select all that apply)

- 0% a. Accuracy
- 0% b. False Positive Rate
- 0% c. False Negative Rate
- 0% d. False Discovery Rate
- 0% e. False Omission Rate
- 0% f. Positive Predictive Value
- 0% g. Negative Predictive Value
- 0% h. Proportion of positive prediction (also called acceptance rate)

# Fairness metrics – 2

The company SuperCrack has developed a model to detect fissures in concrete before they become visible. They evaluate their model against a benchmark. The results look like this:



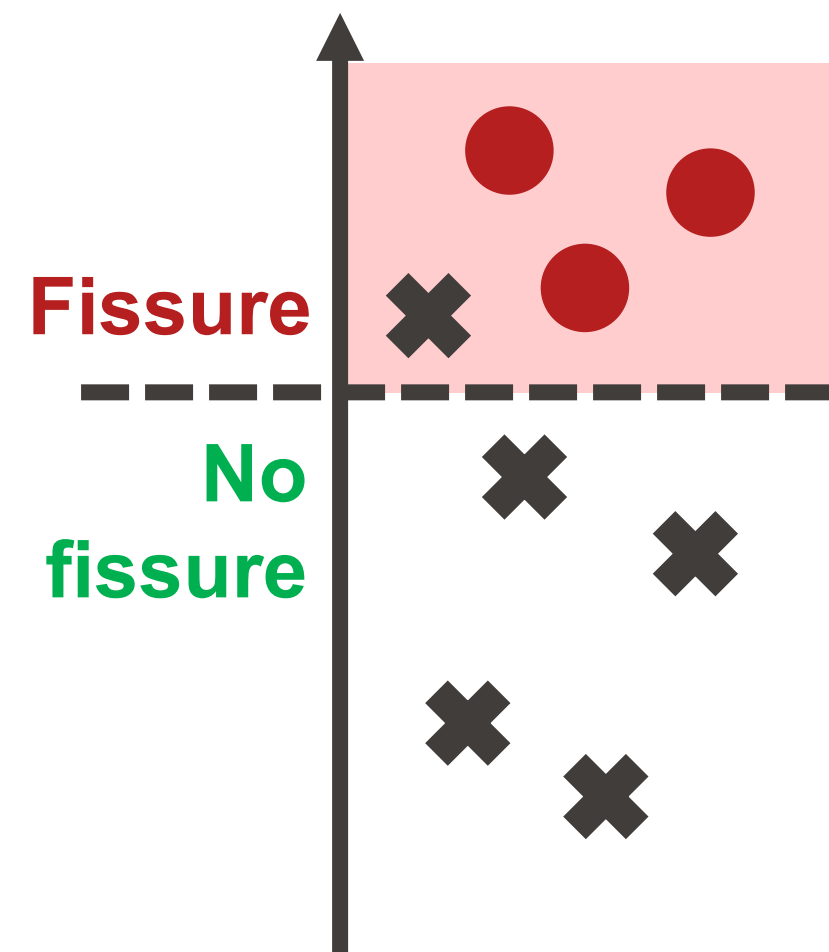
# Fairness metrics – 2a

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

They want to know whether their model performs equally well for plain concrete and for reinforced concrete. Here are the results:

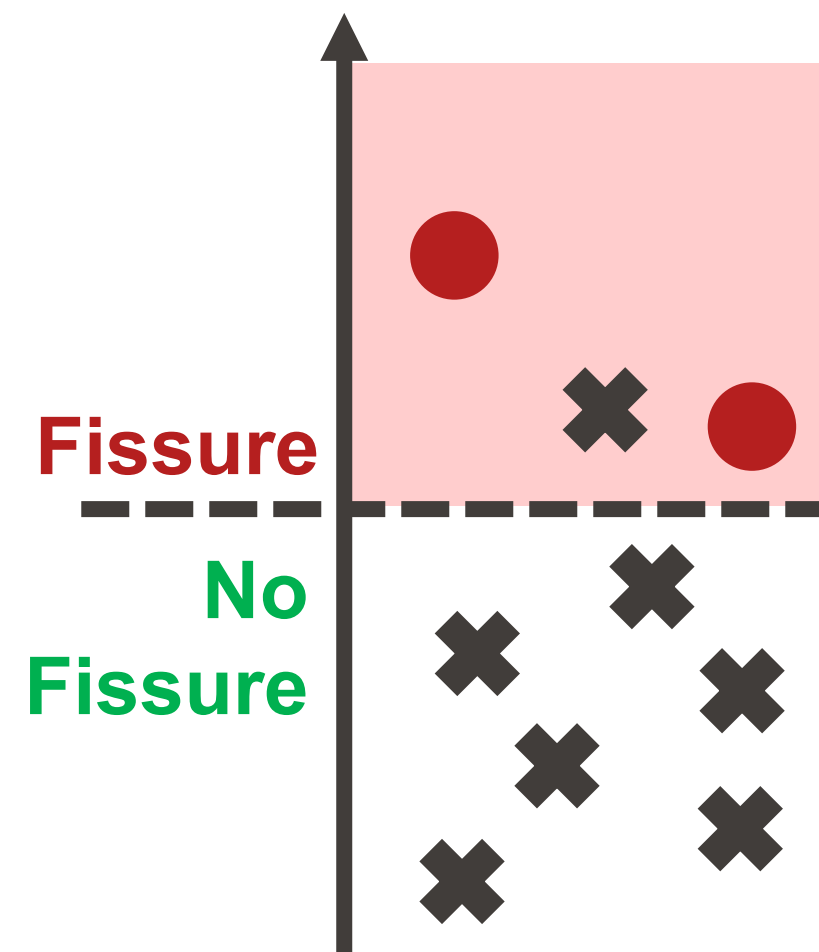
Metric = 4 / 8

Plain  
Concrete



Metric = 3 / 9

Reinforced  
Concrete



**Which notion of fairness are they using?**  
(select 1 answer)

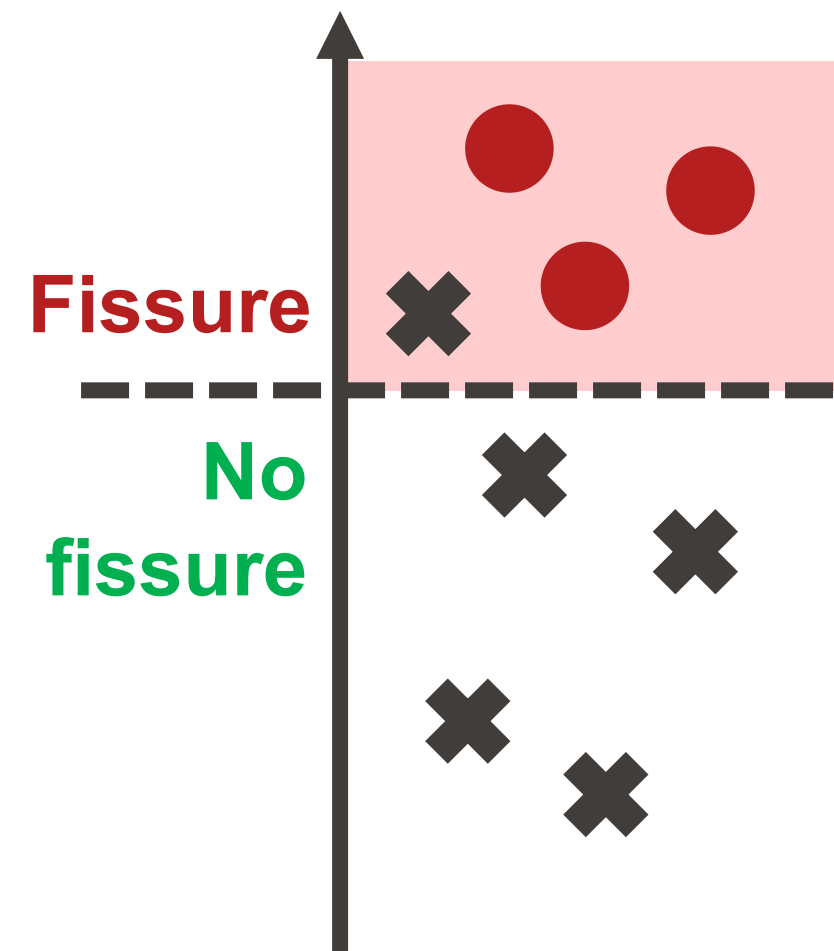
- 0% a. Equal accuracy
- 0% b. Error rate balance
- 0% c. Error parity
- 0% d. Demographic parity

# Fairness metrics – 2b

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

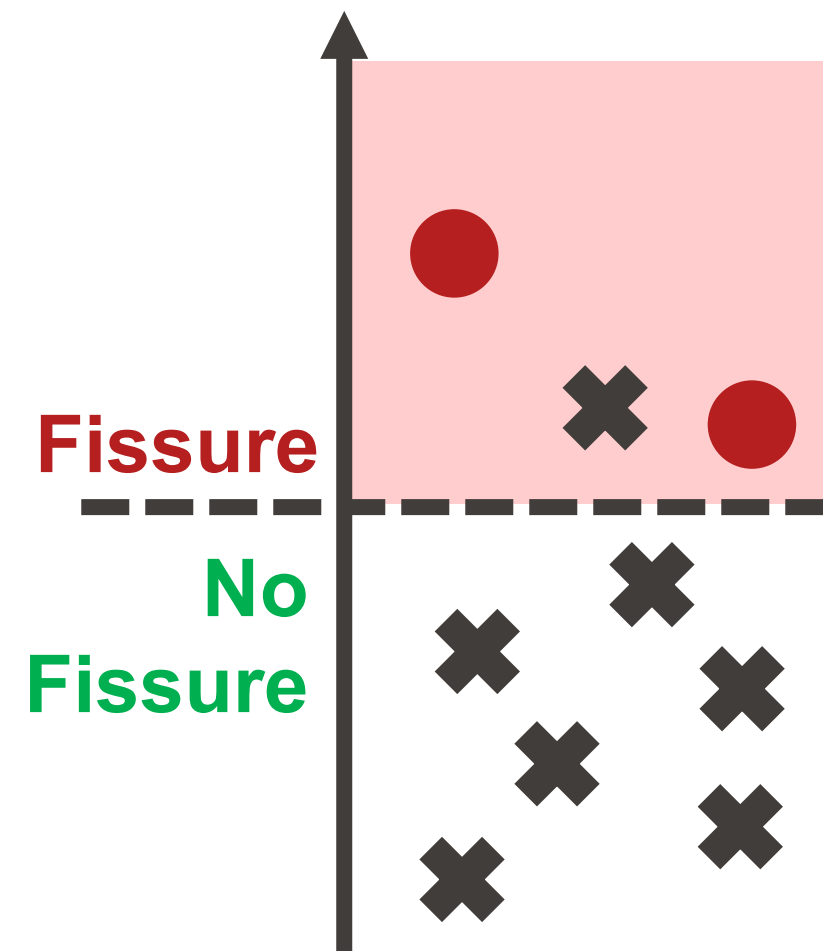
Metric = 4 / 8

Plain  
Concrete



Metric = 3 / 9

Reinforced  
Concrete



**According to this metric,  
is their model fair?**  
(select 1 answer)

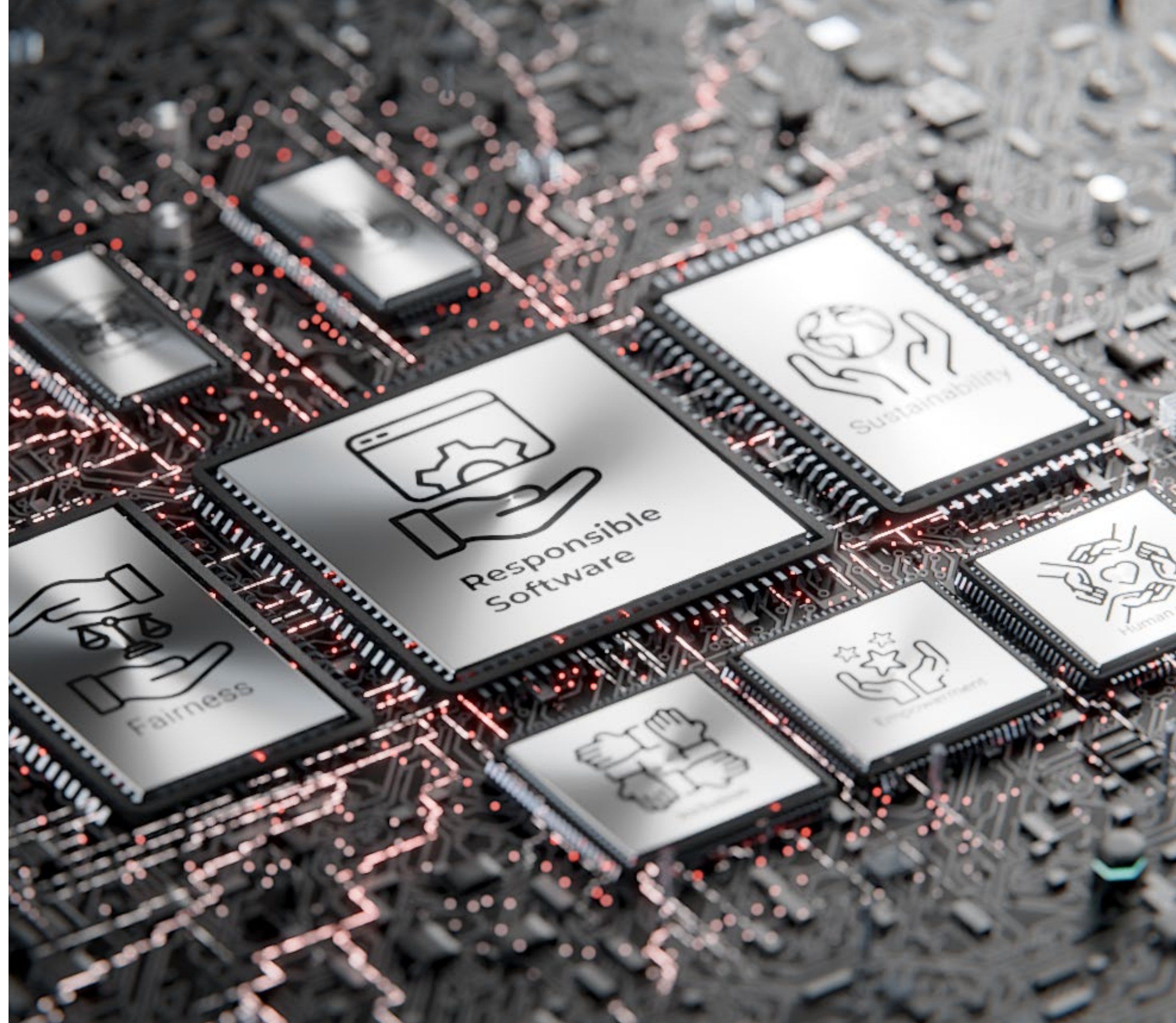
- 0% a. Yes
- 0% b. No
- 0% c. Other option

**EPFL**

**Sustainability 1  
Review & Case  
studies  
10 nov.**

Cécile Hardebolle

**Responsible  
Software**



# Carbon footprint factors

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

What are (some of) the factors in the carbon footprint of software?  
(select all that apply)

- a. The programming language
- b. The computational complexity of the code
- c. The type of hardware
- d. The carbon intensity of the electricity mix
- e. The location where software is hosted
- f. The time at which software runs

# CO<sub>2</sub> equivalent

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

An electricity production facility reports the following emissions per kWh produced:

- 250 g of carbon dioxide (CO<sub>2</sub>)
- 8 g of fossil methane (CH<sub>4</sub>)

What are the carbon emissions of the facility in g CO<sub>2</sub> eq / kWh (considering the GWP-100)?

- a. 240 g / kWh
- b. 258 g / kWh
- c. 490 g / kWh
- d. 656 g / kWh
- e. 906 g / kWh

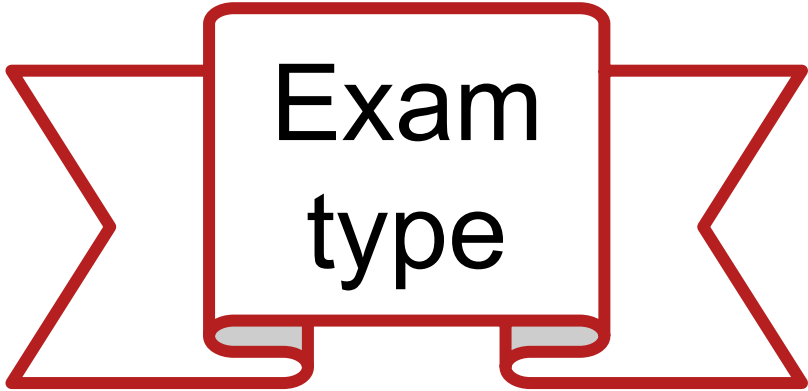
# Power Usage Effectiveness

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

The GreenDC datacenter consumes an average of 1 MW.  
This means annually a total of 8 760 MWh of electricity.  
50% of this electricity is used to power the IT equipment.  
What is the PUE of GreenDC?

- a. 0.5
- b. 1
- c. 1.5
- d. 2

A red-outlined graphic consisting of a central rectangle with the text "Exam type" inside, flanked by two chevron-like shapes pointing outwards.

Exam  
type

# Scopes in the GHG protocol

URL: ttpoll.eu

Session ID: cs290

For a software development company, the electricity consumed by software during the development phase falls into:

- a. Scope 1 (direct)
- b. Scope 2 (indirect, energy)
- c. Scope 3 (indirect, value chain)
- d. It depends

# Direct stakeholders

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

Which of the following stakeholders can be considered **direct** stakeholders in the case:

- a. Internal IT employees working on IT infrastructure
- b. Corporate clients using the center to provide applications
- c. Users of applications hosted by the center
- d. Companies providing energy to the center
- e. Local population in the area of the center
- f. Local ecosystems in the area of the center

# Rebound effect

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

The rebound effect is when **higher energy efficiency** in a product (i.e. lower energy consumption from use) leads to an **increase of total energy consumption** because:

(select all that apply)

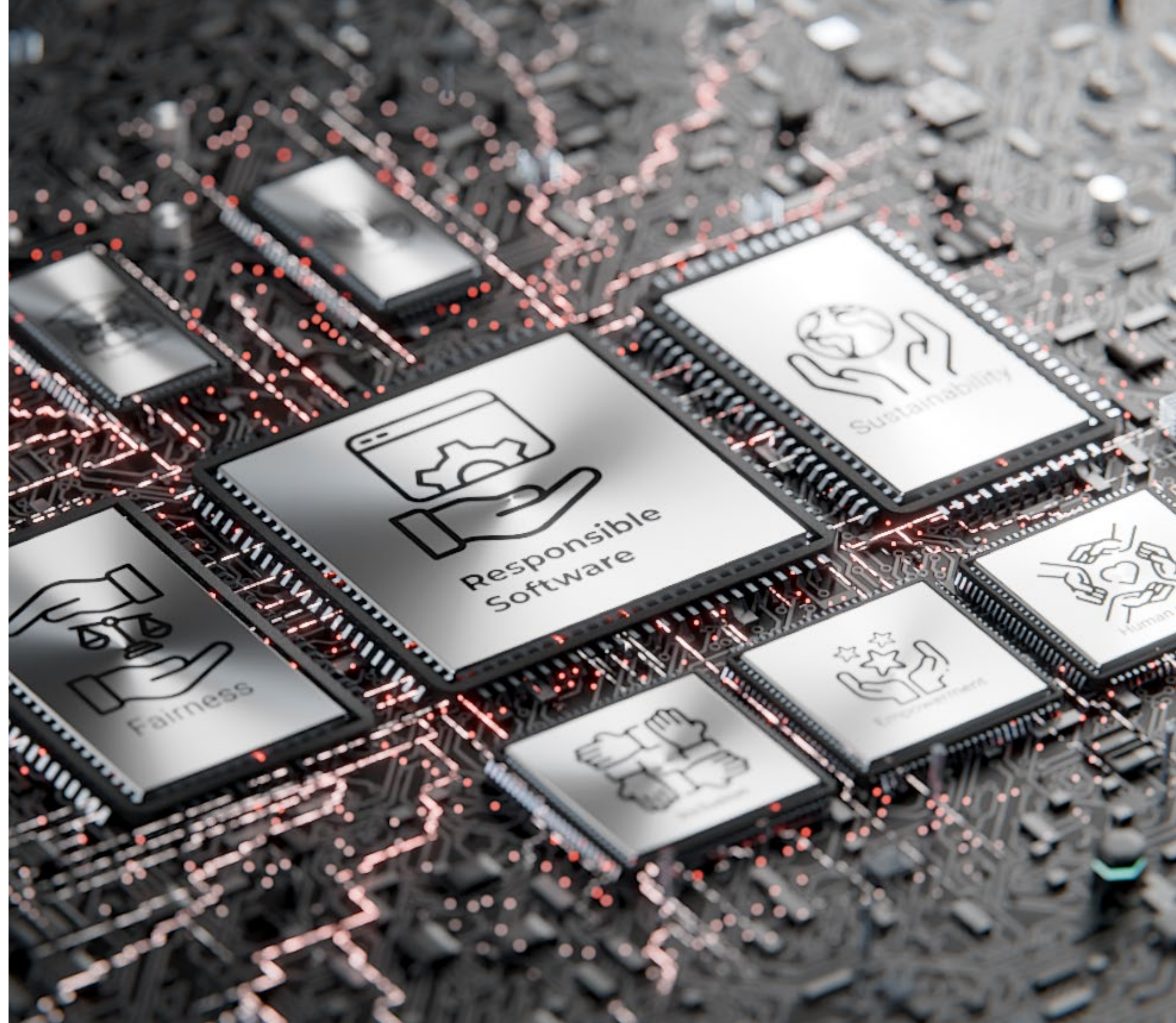
- a. The demand for the product increases
- b. The demand for the product decreases
- c. The product is used more often
- d. The product is used less often
- e. The consumption of other products increases
- f. The consumption of other products decreases

**EPFL**

**Sustainability 2  
Review & Case  
studies  
17 nov.**

Cécile Hardebolle

**Responsible  
Software**



# The footprint of training - 1

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

What are the 3 most important elements in the carbon footprint of ML training?

Rank them **by decreasing impact** (i.e. most impactful first) :

- a. The training time
- b. The power consumption of the CPU
- c. The power consumption of the GPU
- d. The PUE of the datacenter
- e. The carbon intensity of the electricity

# The footprint of training - 2

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

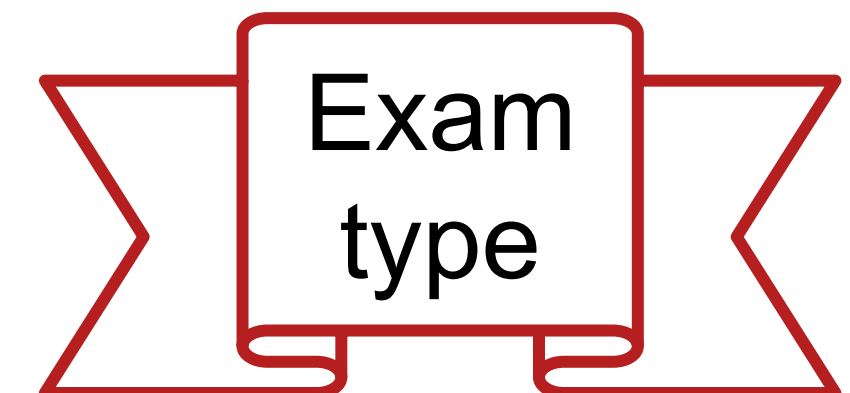
Let's consider the training of the model SupChat-7B. The computing node has 2 GPUs of the model Nvidia A100 80GB, which consume 400W each. Our datacenter, which has a PUE of 1.2, is located in Germany (carbon intensity: 381g CO<sub>2</sub>e / kWh).

The training time is 80 000 hours of total GPU computation time.

**What is the carbon footprint for the training of SUPMOD-7B?**

- a. 14,63 tons CO<sub>2</sub>e
- b. 29,26 tons CO<sub>2</sub>e
- c. 14 630,4 tons CO<sub>2</sub>e
- d. 29 260,8 tons CO<sub>2</sub>e

Exam type but with calculations that can be done by hand (i.e. simpler than here)



# The footprint of inference - 1

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

What are the 3 most important elements in the carbon footprint of ML inference?

Rank them **by decreasing impact** (i.e. most impactful first) :

- a. The number of user queries
- b. The electricity consumed per query
- c. The PUE of the datacenter
- d. The carbon intensity of the electricity

# The footprint of inference - 2

---

The model SupChat-7B is now deployed in production. It is hosted on the same computing node with 2 GPUs of the model Nvidia A100 80GB, which consume 400W each. Our datacenter, which has a PUE of 1.2, is located in Germany (carbon intensity: 381g CO<sub>2</sub>e / kWh). Our model is able to serve 120 tokens per second.

It has an average of 2000 users daily and generates an average of 5000 tokens per user per day.

## **What is the carbon footprint of 1 day of inference?**

1. What is the total GPU computation time used over 1 day (in h)?
2. What is the power consumed by the model for inference (in W)?
3. What is the total electricity consumed over 1 day (in kWh)?
4. What is the carbon footprint over 1 day (in kg CO<sub>2</sub>e)?

# Total carbon footprint

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

We have obtained the carbon footprint of SupChat-7B at training and at inference time. What is its total carbon footprint?

- a. Training
- b. Inference
- c. Training x Inference
- d. Inference – Training
- e. Training + Inference
- f. Other

# Hardware renewal

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

We want to optimize the energy consumption of SupChat-7B at inference time. We decide to upgrade our hardware platform and to replace our A100 GPUS with H100 GPUs. The H100 are 4 times more performant than the A100 in terms of computation speed. Their power consumption is 700W at maximum use.

What effect(s) are we likely to observe (select all that apply)?

- a. A decrease in the energy consumption
- b. An increase in the energy consumption
- c. A decrease in the overall carbon footprint
- d. An increase in the overall carbon footprint

# Water Usage Effectiveness

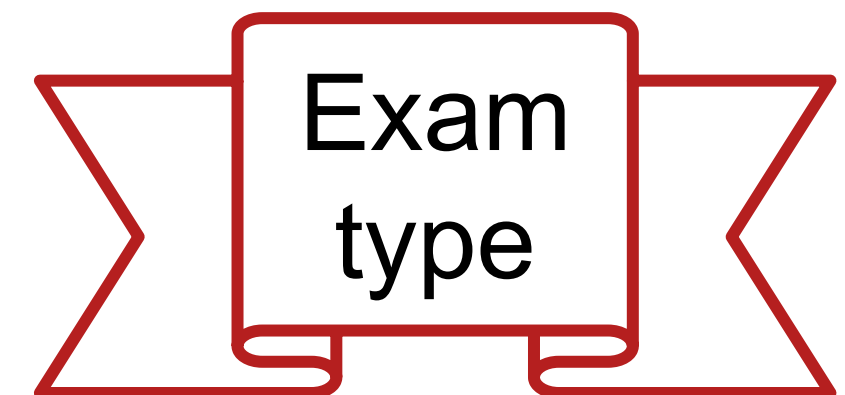
URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

The datacenter hosting SupChat-7B consumes an average of 1 MW. This means annually a total of 8 760 MWh of electricity. It consumes approximately 15.8 million liters of water each year. What is the WUE of the datacenter (onsite only)?

- a. 0,18
- b. 0,55
- c. 1,8
- d. 18,03
- e. 55,44

Exam type but with calculations that can be done by hand (i.e. simpler than here)

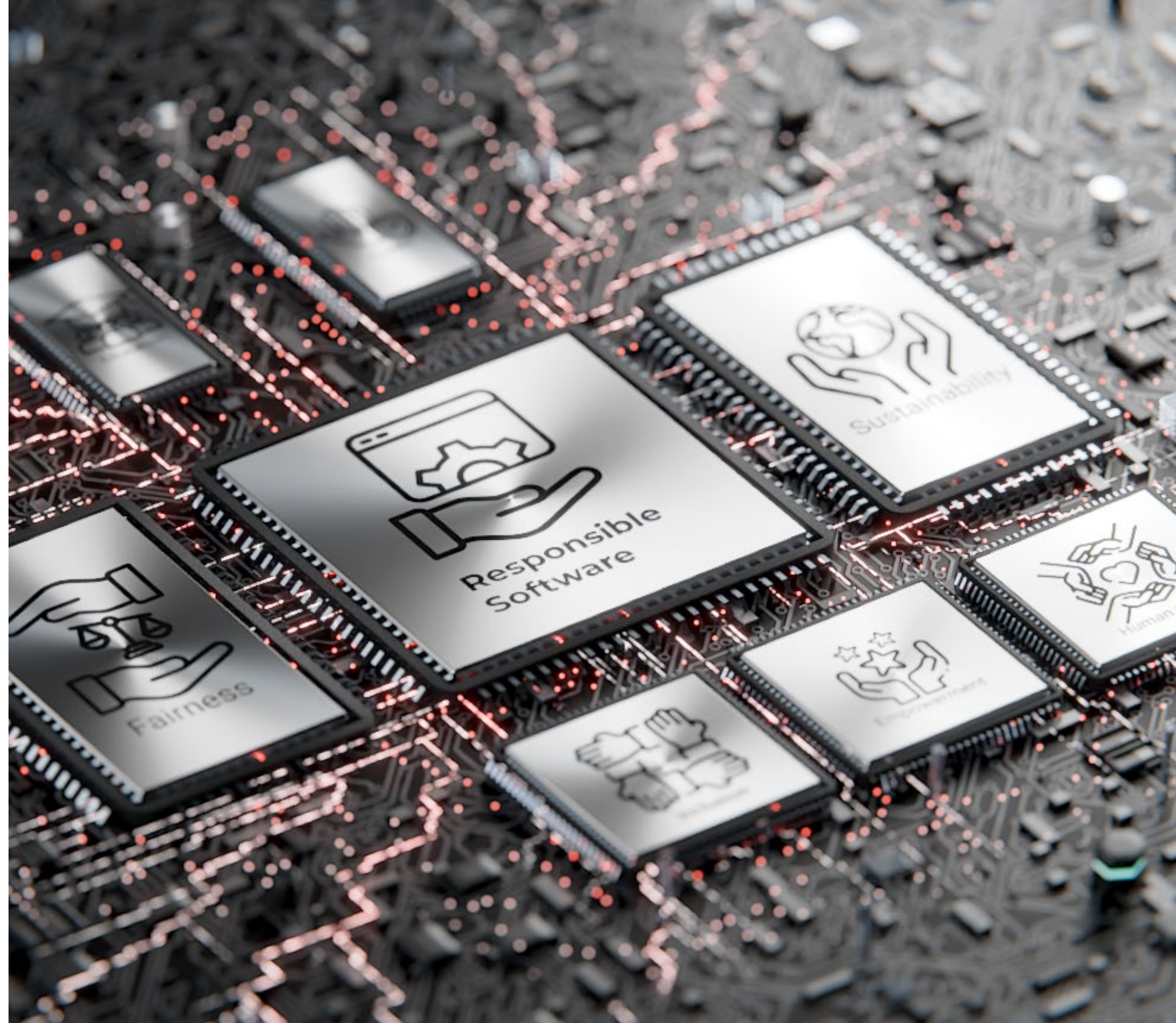


**EPFL**

**Empowerment 1  
Review & Case  
studies  
24 nov.**

Cécile Hardebolle

**Responsible  
Software**



# Meditation app

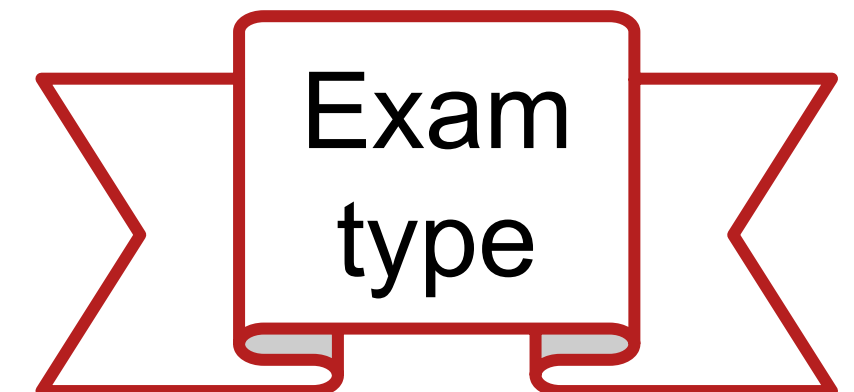
URL: ttpoll.eu

Session ID: cs290

ZenPath is an app dedicated to mental well-being that offers guided meditation sessions online. To reduce user dropout, they decide to display a popup after a user skips two sessions where the “Resume Today!” button is preselected.

What type of nudging technique is most likely used here?

- a. Opt-in
- b. Social proof
- c. Scarcity
- d. Default



# Use of data

URL: ttpoll.eu  
Session ID: cs290

← Back

## Data for Generative AI Improvement

Can LinkedIn and its affiliates use your personal data and content you create on LinkedIn to train generative AI models that create content?

Use my data for training content creation AI models

On



This setting controls the training of generative AI models used to create content. When this setting is on LinkedIn and its affiliates may use your personal data and content you create on LinkedIn for that purpose. [Learn more.](#)

This is one of the settings on LinkedIn in the USA, set to its default value.

What is the most likely outcome?

- Most users will turn the setting off
- Most users will turn the setting on
- Most users will let the setting as is
- Most users will change the setting

# Navigation app

URL: ttpoll.eu

Session ID: cs290

In an effort towards more sustainability, the itinerary search in Noodle Maps now returns 2 itinerary options in the following order:

- 1) the most fuel-efficient but longest itinerary
- 2) the shortest but least fuel-efficient itinerary

What are the characteristics of this nudge? (select all that apply)

- a. Takes advantage of System 1
- b. Takes advantage of System 2
- c. Transparent to the user
- d. Covert
- e. Ethically fine
- f. Ethically problematic

# Deceptive patterns vs nudges

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

Which of the following are characteristics shared by nudges and deceptive patterns? (select all that apply)

- a. They modify the choice architecture
- b. They make users do things they didn't originally mean to
- c. They take advantage of how humans make decisions
- d. They intentionally bias user behavior
- e. They restrict choices
- f. They benefit users
- g. They benefit another party
- h. They make users lose track of time

# E-commerce platform

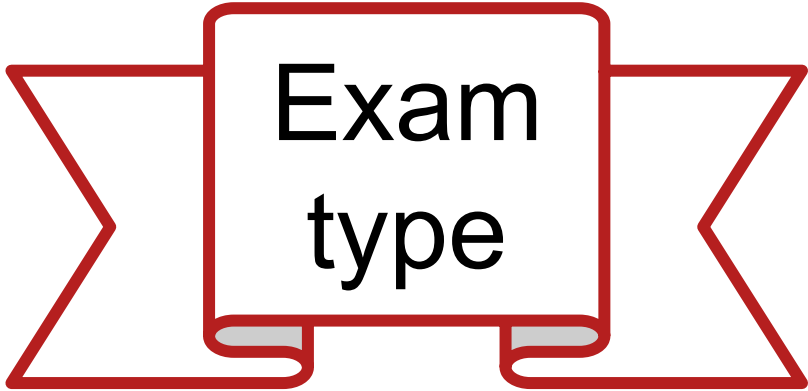
URL: ttpoll.eu

Session ID: cs290

The e-commerce platform Shine would like to implement new features to improve the experience of its various categories of users. Here is the list of envisaged features.

Which of them best matches the definition of a deceptive pattern?

- a. Personalize style recommendations based on past browsing
- b. Display user-provided past purchase data to recommend sizes
- c. Register users to a ShineClub membership trial on checkout
- d. Provide downloadable QR codes for the free return of items



Exam  
type

# Translation

URL: ttpoll.eu  
Session ID: cs290

Consider the following translation. What is the issue here?

French ▾ ↔ English (American) ▾ Glossary

Dans un souci de durabilité, la recherche d'itinéraire dans Noodle Maps renvoie désormais 2 options d'itinéraire dans l'ordre suivant :

- 1) l'itinéraire consommant le moins de carburant mais le plus long
- 2) l'itinéraire le plus court **mais consommant plus de carburant**

×

In the interests of sustainability, the route search in Noodle Maps now returns 2 route options in the following order:

- 1) the most fuel-efficient but longest route
- 2) the shortest **but most fuel-efficient route**

- a. Parity error
- b. Factuality error
- c. Measurement error
- d. Faithfulness error

# Evaluating the level of risk - 1

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

Consider the following Privacy risk: **“Tracks personal app usage”**  
How would you evaluate the level of this risk in terms of probability and severity of impacts?

(select 2 options: 1 for probability, 1 for severity)

- a. Probability: low
- b. Probability: medium
- c. Probability: high
- d. Severity: low
- e. Severity: medium
- f. Severity: high

# Evaluating the level of risk - 2

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

Consider the following Welfare risk: “**Excessive reminders could lead to stress or anxiety**”. How would you evaluate the level of this risk in terms of probability and severity of impacts?

(select 2 options: 1 for probability, 1 for severity)

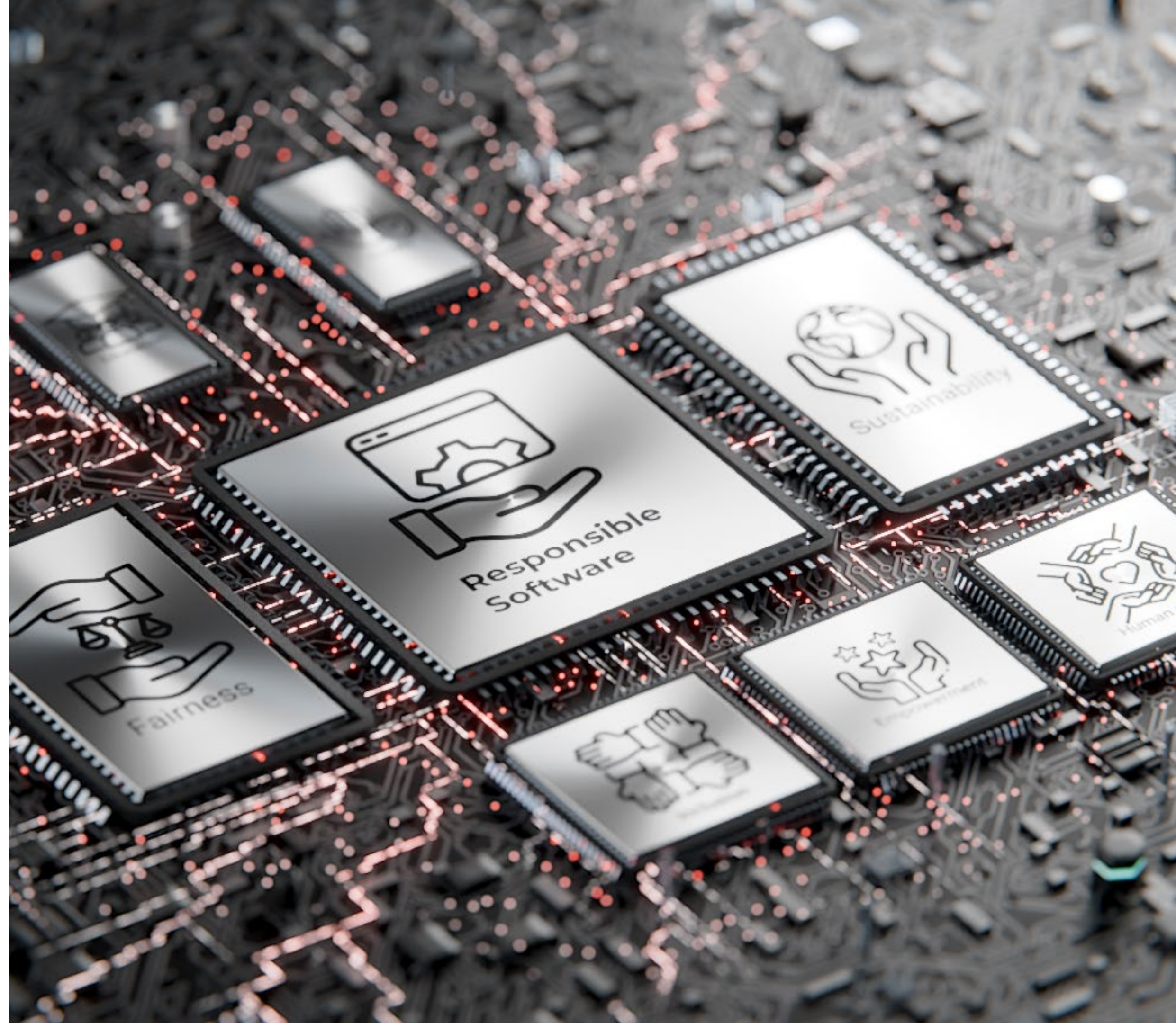
- a. Probability: low
- b. Probability: medium
- c. Probability: high
- d. Severity: low
- e. Severity: medium
- f. Severity: high

**EPFL**

**Empowerment 2  
Review & Case  
studies  
8 dec.**

Cécile Hardebolle

**Responsible  
Software**



# Privacy policies

URL: ttpoll.eu  
Session ID: cs290

Several studies have shown that the privacy policies of many online platforms and websites are extremely long (several thousand of words, taking in the 20 minutes to read on average), use legalistic terminology and are hard to navigate.

This can be said to be a transparency issue because (select all that apply):

- Information is not accessible
- Information is not understandable
- Information is not relevant



The New York Times

SHARE

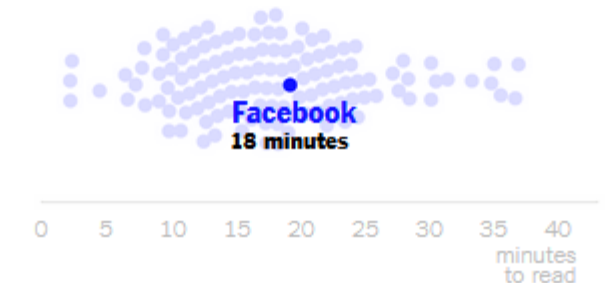
Opinion | THE PRIVACY PROJECT

## We Read 150 Privacy Policies. They Were an Incomprehensible Disaster.

By Kevin Litman-Navarro

In the background here are several privacy policies from major tech and media platforms. Like most privacy policies, they're verbose and full of legal jargon — and opaquely establish companies' justifications for collecting and selling your data. The data market has become the engine of the internet, and these privacy policies we agree to but don't fully understand help fuel it.

To see exactly how inscrutable they have become, I analyzed the length and readability of privacy policies from nearly 150 popular websites and apps. Facebook's privacy policy, for example, takes around 18 minutes to read in its entirety - slightly above average for the policies I tested.



(Sherman, 2024; Litman-Navarro, 2019)

# Beer brewing dataset - 1

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

One of the results of your Bachelor thesis is a very cool dataset which contains tasting profiles and consumer reviews for 3197 unique beers from 934 different breweries. This dataset can be used to train machine learning models for sentiment analysis and classification tasks.

You want to make the dataset public.

For ensuring transparency you should also publish with it:

(select all that apply):

- a. Composition of the data, including demographics
- b. Description of the collection process
- c. Description of the pre-processing performed
- d. Description of the purposes and intended use

# Beer brewing dataset - 2

URL: [ttpoll.eu](http://ttpoll.eu)

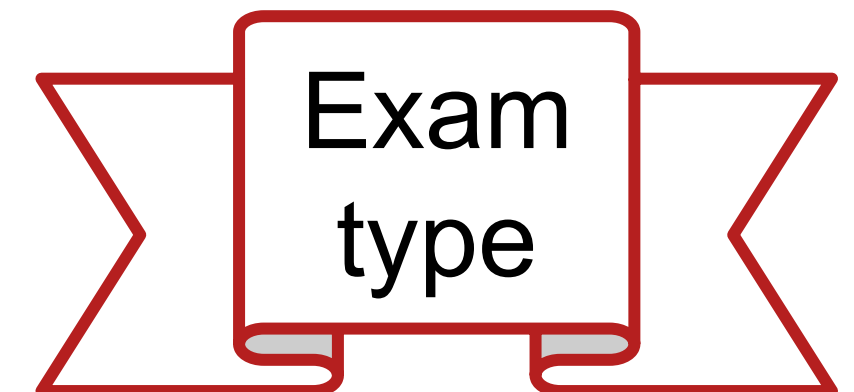
Session ID: cs290

One of the results of your Bachelor thesis is a very cool dataset which contains tasting profiles and consumer reviews for 3197 unique beers from 934 different breweries. This dataset can be used to train machine learning models for sentiment analysis and classification tasks.

You have created a datasheet for your dataset.

Which of the FAIR principles do you follow by providing a datasheet?

- a. Findable
- b. Accessible
- c. Interoperable
- d. Reusable



# Linear Regression Model

---

You have found on HuggingFace an open-source Linear Regression model that predicts the price of a house based on a range of features like lot area, construction year, number of rooms, etc. For recall, a linear regression model has the following mathematical form, where  $y'$  is the predicted price,  $x_i$  are the features and  $b$  and  $w_i$  are the final parameters of the model:

$$y' = b + w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + \dots$$

Let's imagine that you want to modify this model. How could you do it?

- a. Modify the values of the parameters
- b. Use a post-hoc interpretability method
- c. Retrain the model with a new dataset
- d. It is not possible to modify the model

# Logistic Regression Model

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

In the Fairness 2 notebook you have created a Logistic Regression model on the ProPublica dataset to try to reproduce how the COMPAS software predicts the risk of recidivism.

The Logistic Regression model you have created can be said to be (select all that apply):

- a. White-box
- b. Black-box
- c. Post-hoc interpretable
- d. Interpretable by design

# COMPAS

URL: [tppoll.eu](http://tppoll.eu)

Session ID: cs290

To have transparency on the ML model behind the COMPAS software would mean to have access to:

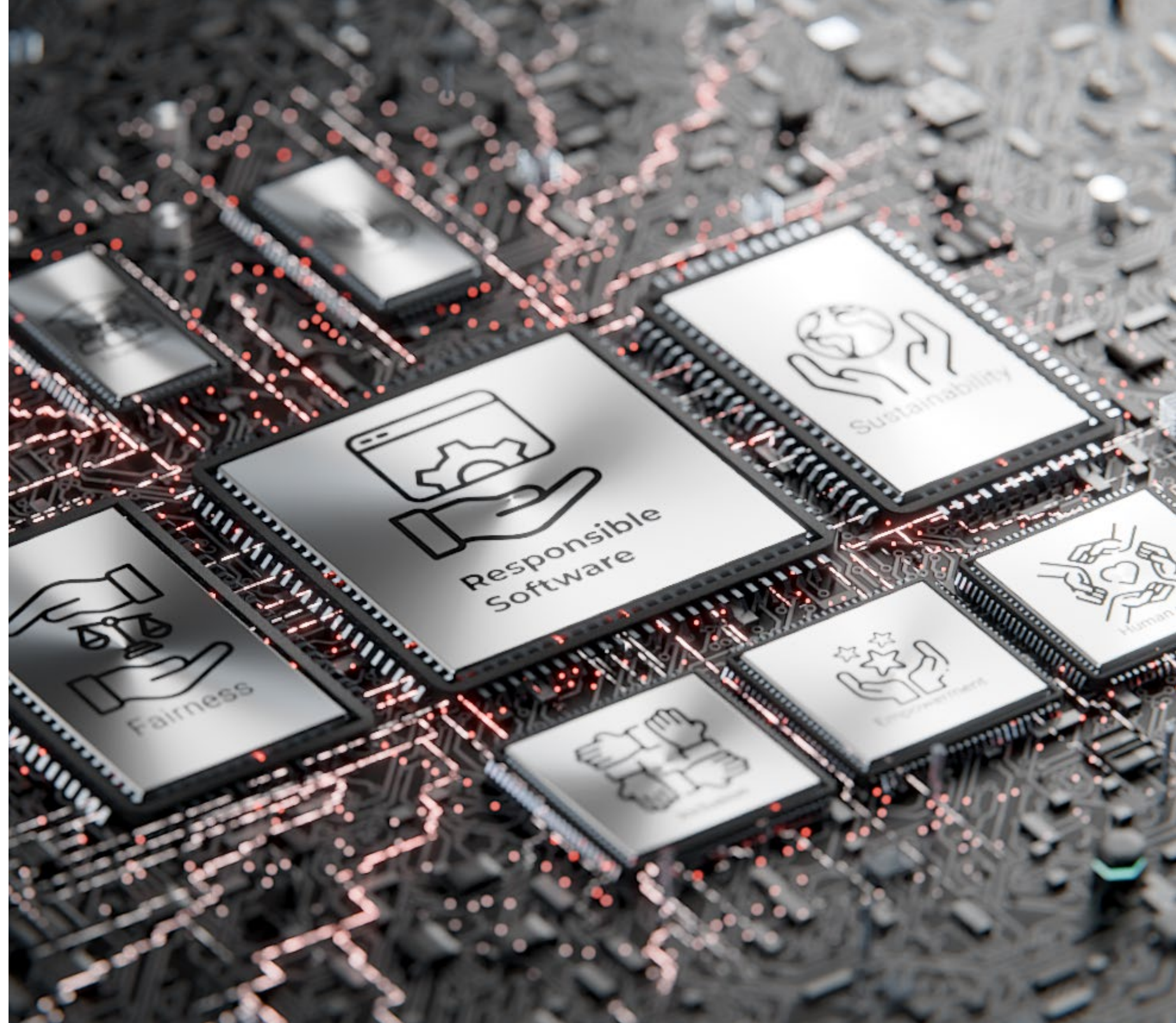
- a. The design documentation
- b. The user documentation
- c. The training code
- d. The training dataset
- e. A post-hoc interpretability method
- f. It depends

**EPFL**

**Conclusion  
Case studies  
+ Q&A  
15 dec.**

Cécile Hardebolle

**Responsible  
Software**



**Review questions**  
**“Whole Course”**

# Ethical sensitivity

New

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

What is ethical sensitivity?

- a. The ability to predict all technical outcomes before deployment
- b. The capability to identify the impact of a situation on others
- c. The ability to act to benefit others even at your own expense
- d. The capacity to account for all ethical values simultaneously

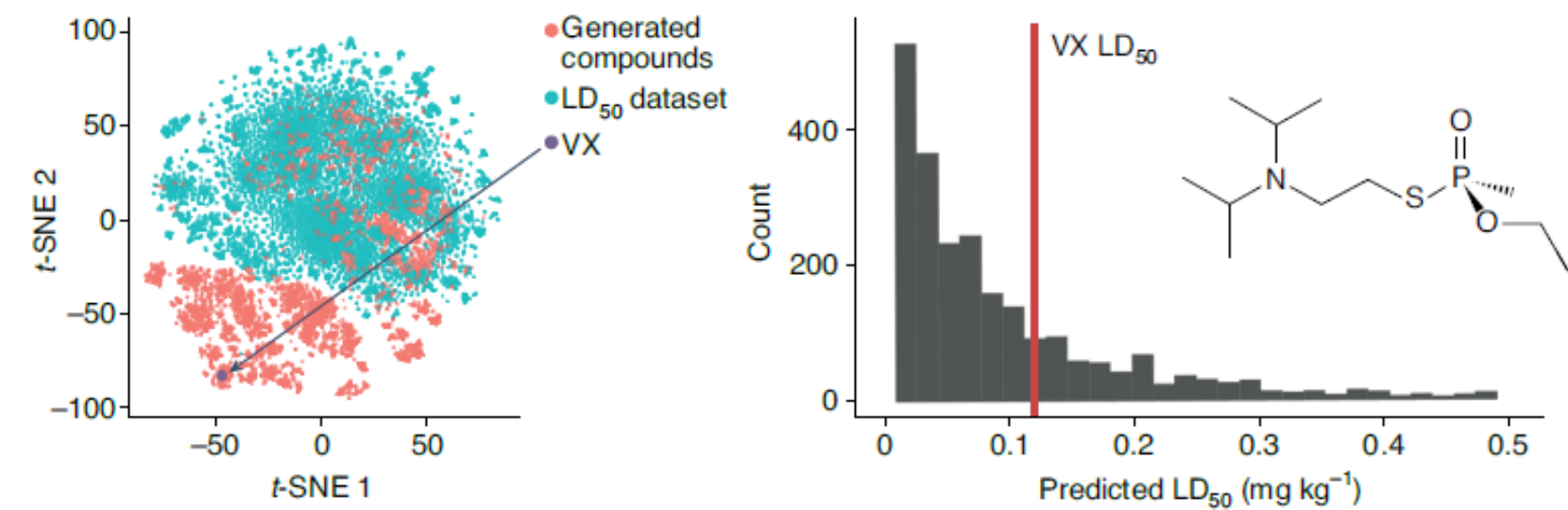
# Chemical discovery

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

A software company has developed a Machine Learning model that is able to discover new chemical compounds for medicine development. They identify that the model can also discover new chemical weapons.

What type issue is this?

- a. A technical issue
- b. An ethical issue
- c. An ethical dilemma



**Fig. 1 |** A t-SNE plot visualization of the LD<sub>50</sub> dataset and top 2,000 MegaSyn AI-generated and predicted toxic molecules illustrating VX. Many of the molecules generated are predicted to be more toxic in vivo in the animal model than VX (histogram at right shows cut-off for VX LD<sub>50</sub>). The 2D chemical structure of VX is shown on the right.

(Urbina et al., 2022)

# Vulnerabilities

New

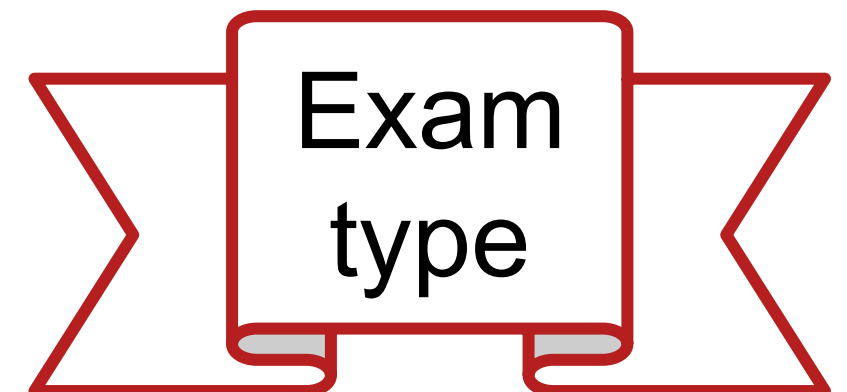
URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

A software engineer decides to postpone the launch of a new feature due to the late discovery of a security vulnerability and justifies: “The new feature would bring us some short-term benefits but would have serious negative consequences for all of our customers, our aim must be the greatest good for the greatest number.”

Which ethical theory does this engineer follow?

- a. Utilitarianism
- b. Deontology
- c. Virtue
- d. Care



# Food delivery

New

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

An online food delivery app experiences a data breach where customer payment details are stolen.

What type of risks are represented in this situation?

- a. Safety risks from misdiagnosed food allergies
- b. Safety risks from incorrect delivery scheduling
- c. Sociotechnical risks in app-driver communication
- d. Security risks from unauthorized system access

# Hospital

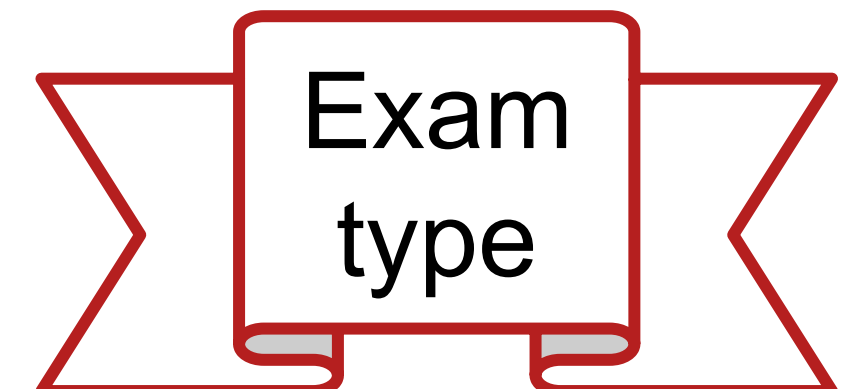
New

URL: ttpoll.eu  
Session ID: cs290

Patient records in a hospital have been encrypted by cybercriminals who demand payment to restore access, causing emergency services to halt and delay critical care for patients.

Which harm scenario does this represent?

- a. Unintended use
- b. Malfunction
- c. Misuse
- d. Intended use



# Fissures in concrete

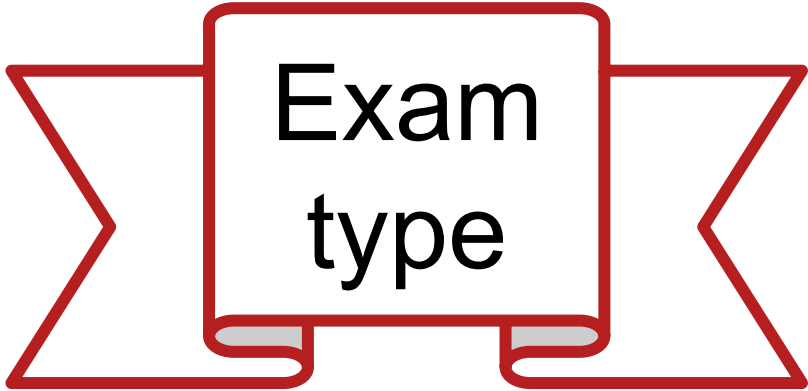
URL: ttpoll.eu

Session ID: cs290

The company SuperCrack has developed a model to detect fissures in concrete walls before they become visible to the naked eye. A positive result means presence of fissure.

Which of the statements below is correct?

- a. TN = actual absence of fissure, correct prediction
- b. TN = actual presence of fissure, incorrect prediction
- c. TP = actual presence of fissure, incorrect prediction
- d. TP = actual absence of fissure, correct prediction

A red-outlined graphic element consisting of a central rectangular box with the text "Exam type" inside, flanked by two chevron-like shapes pointing outwards.

Exam  
type

# Contagious disease

New

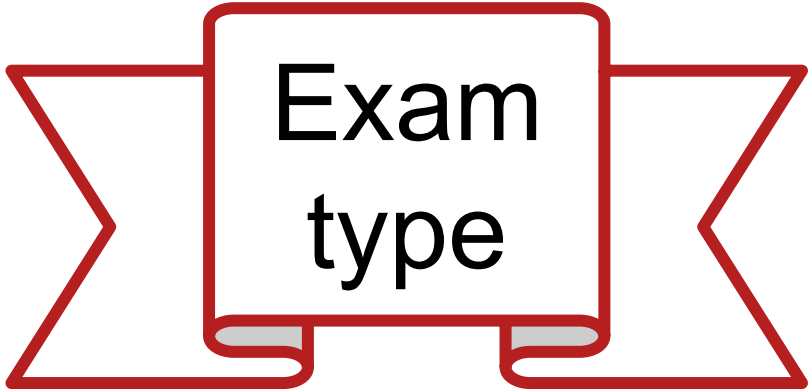
URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

A rapid test for a contagious disease (infected = positive result) shows a high number of false negatives.

What are the consequences of false negatives in terms of safety?

- a. Healthy people continue their daily activities as normal.
- b. Healthy people receive unnecessary quarantine.
- c. Infected individuals receive the appropriate medication.
- d. Infected individuals spread the disease unknowingly.

A red-outlined rectangular box with a ribbon-like shape on the left and right sides, containing the text "Exam type".

Exam  
type

# Political campaign

New

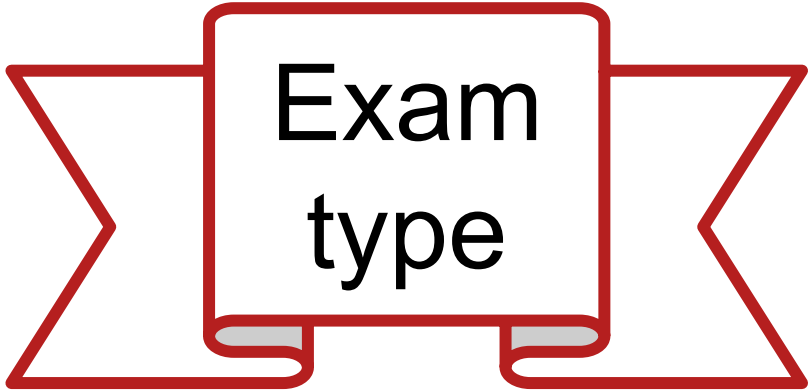
URL: ttpoll.eu

Session ID: cs290

A whistleblower releases authentic internal documents from the campaign of a political party with the goal of damaging the party's public image for the upcoming election.

What type of information is this?

- a. Misinformation
- b. Disinformation
- c. Malinformation
- d. Fake news

A red-outlined ribbon graphic containing the text "Exam type".

Exam  
type

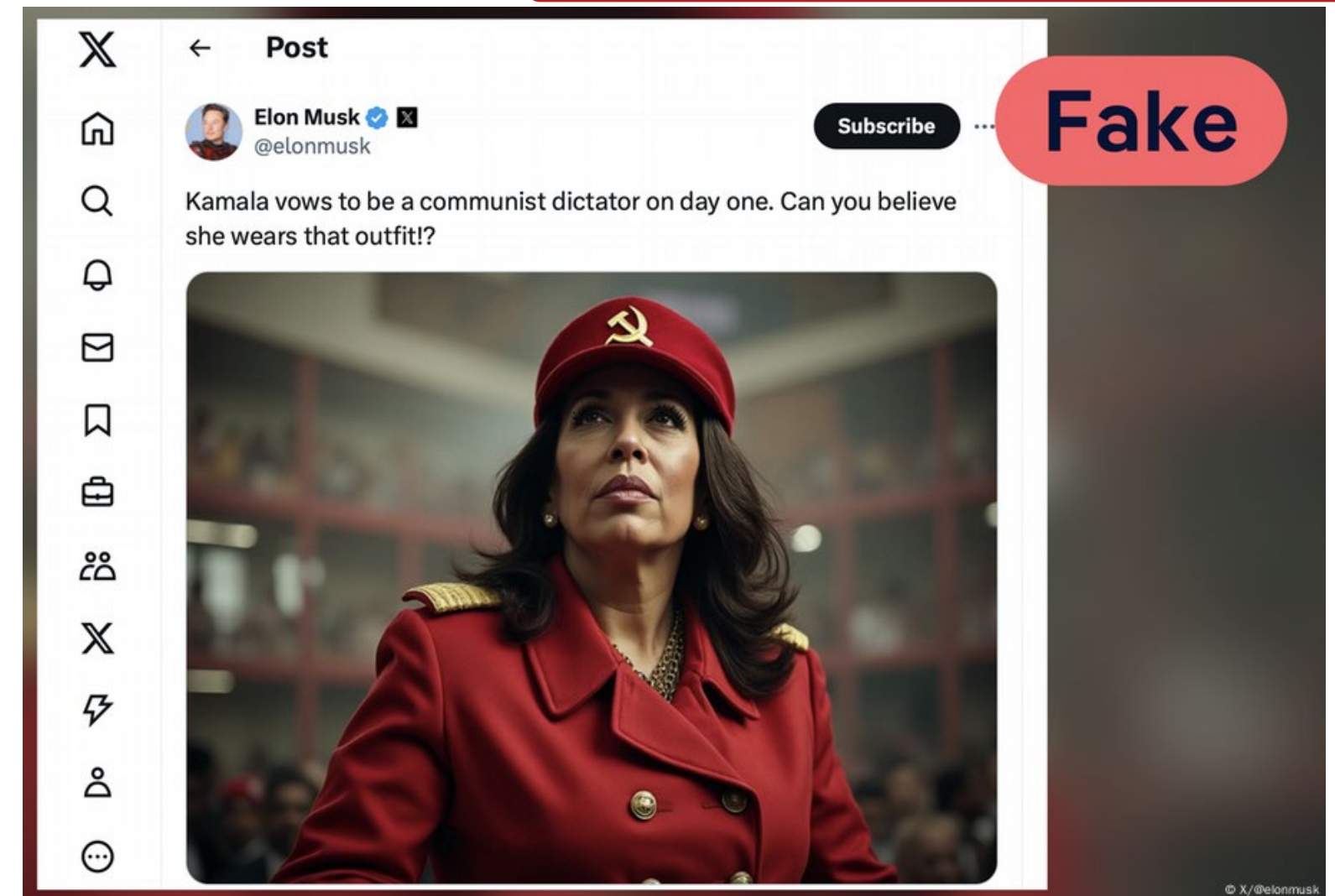
# Posts on Twitter

URL: ttpoll.eu  
Session ID: cs290

One dis-/mis-information post by Elon Musk appears in your Twitter timeline.

Why would you be more likely to believe it than other posts?

- System 2
- Illusory truth
- Source cues
- Prebunking



Fact check: Elon Musk spreads US election lies. (2024, February 11).  
Dw.Com. <https://www.dw.com/en/fact-check-how-elon-musk-is-spreading-us-election-lies/a-70663408>

Exam  
type

# Loans

New

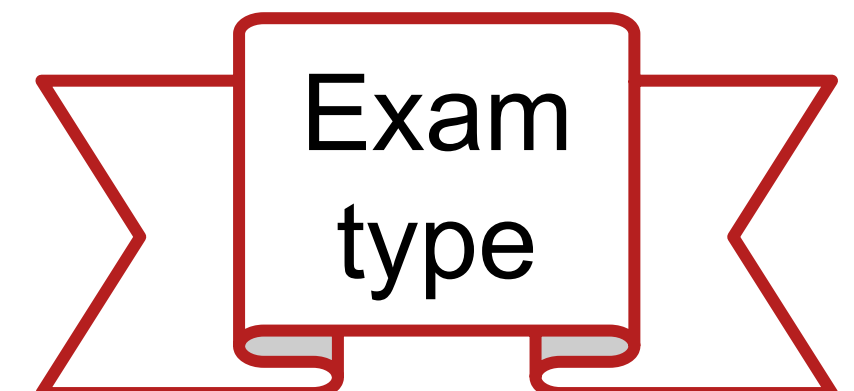
URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

A ML model for loan approval consistently denies loans to applicants from rural neighborhoods. The model has been trained on data from the bank covering all the loan decisions taken in the last 5 years for all the neighborhoods served by the bank.

Which type of bias is most likely present in the data from this scenario?

- a. Sampling bias
- b. Representation bias
- c. Measurement bias
- d. Preexisting bias



# Shoplifting

URL: ttpoll.eu

Session ID: cs290

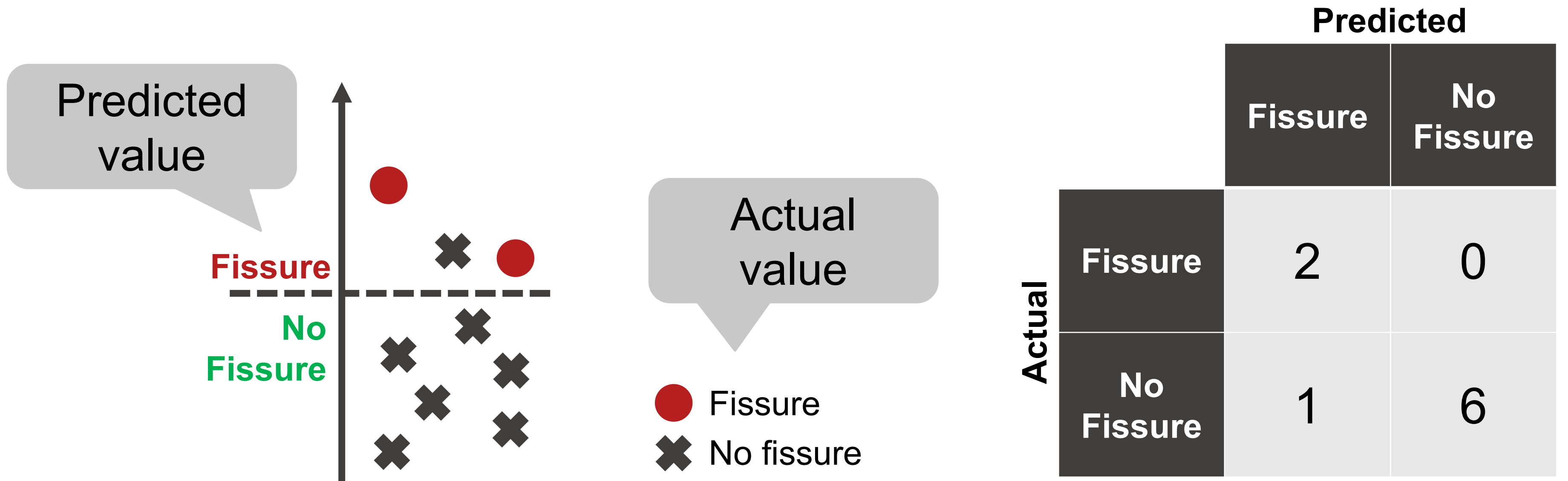
The society RetailProtect develops a ML model to identify instances of shoplifting in retail shops. They evaluate their model on a benchmark in which actors from diverse ethnicities simulate a range of shoplifting actions. They plan to deploy soon in shops.

What type of bias is present in this scenario?

- 0% a. Evaluation bias
- 0% b. Aggregation bias
- 0% c. Optimization bias
- 0% d. Deployment bias

# Fissures in concrete (again)

The company SuperCrack has developed a model to detect fissures in concrete before they become visible. They evaluate their model against a benchmark. The results look like this:



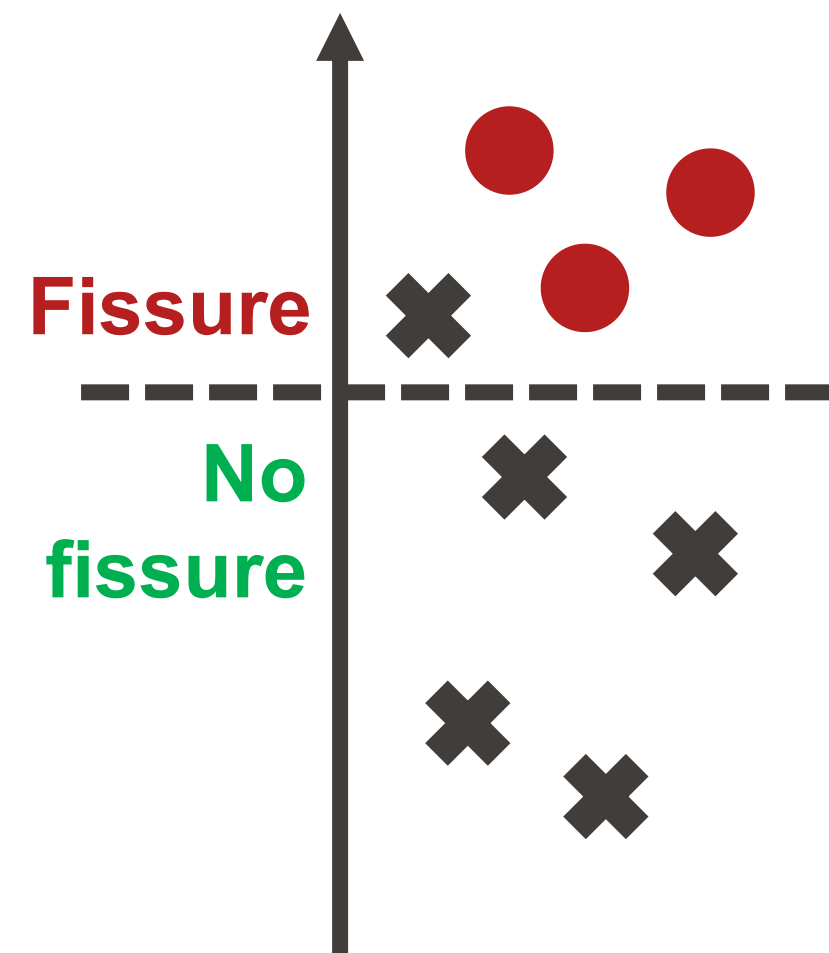
# Fissures in concrete (again) <sup>New</sup>

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

They want to know whether their model performs equally well for plain concrete and for reinforced concrete. Here are the results:

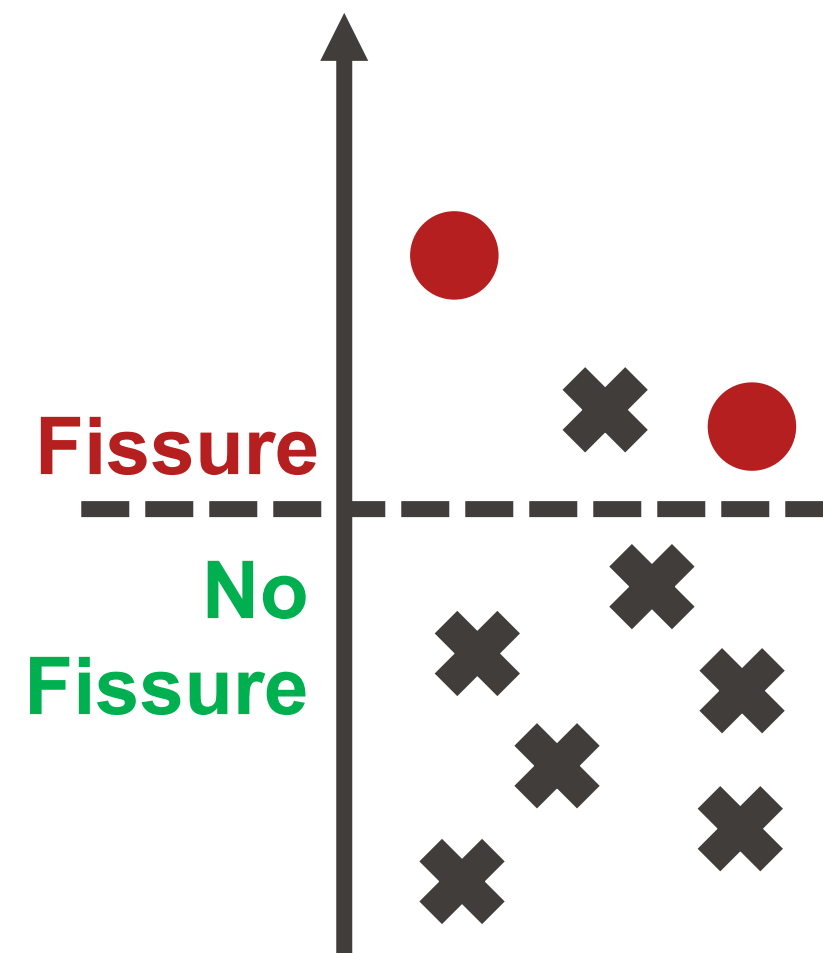
Metric = 1 / 5

Plain  
Concrete



Metric = 1 / 7

Reinforced  
Concrete



Which metric are they using? (select 1 answer)

- 0% a. Accuracy
- 0% b. FNR
- 0% c. FPR
- 0% d. Positive prediction rate

# University admissions

New

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

A model trained to help screen applications to university has an accuracy of 97% and the false positive rate (FPR) is 5% for group X and 6% for group Y. However, the Disparate Impact Ratio is 0,613 with group X having a higher admission rate.

What is most likely happening in this situation?

- 0% a. Differences in the FNR are causing the low DIR
- 0% b. The DIR indicates a higher error rate for group Y
- 0% c. The applicants from group X have stronger profiles
- 0% d. Group Y has a lower rate of actual positive labels

# Datacenter cooling

URL: [ttpoll.eu](http://ttpoll.eu)

Session ID: cs290

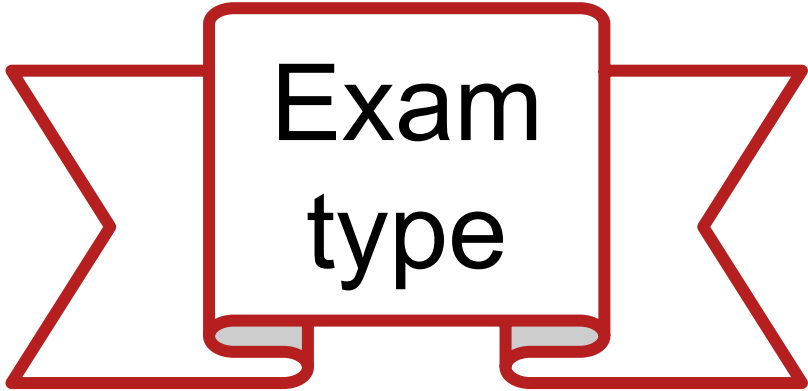
The GreenDC datacenter consumes an average of 1 MW.  
This means annually a total of 8 760 MWh of electricity.  
50% of this electricity is used to power the IT equipment.  
What is the PUE of GreenDC?

0% a. 0.5

0% b. 1

0% c. 1.5

0% d. 2



Exam  
type

# Datacenter water

URL: ttpoll.eu

Session ID: cs290

The TitanCore datacenter consumes a total of 24 000 MWh of electricity annually. It consumes approximately 16 million liters of water each year.

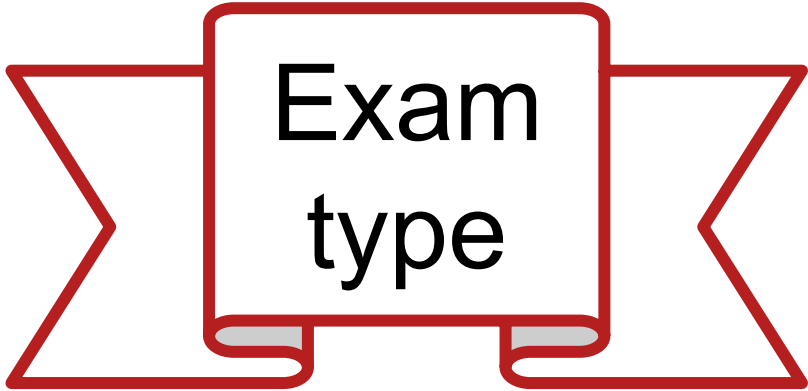
What is the WUE of the datacenter (onsite only)?

0% a. 0,000667

0% b. 0,667

0% c. 1,5

0% d. 1500

A red-outlined graphic consisting of a central rectangle with the text "Exam type" inside, flanked by two stylized, arrow-like shapes pointing outwards.

Exam  
type

# LLM training

New

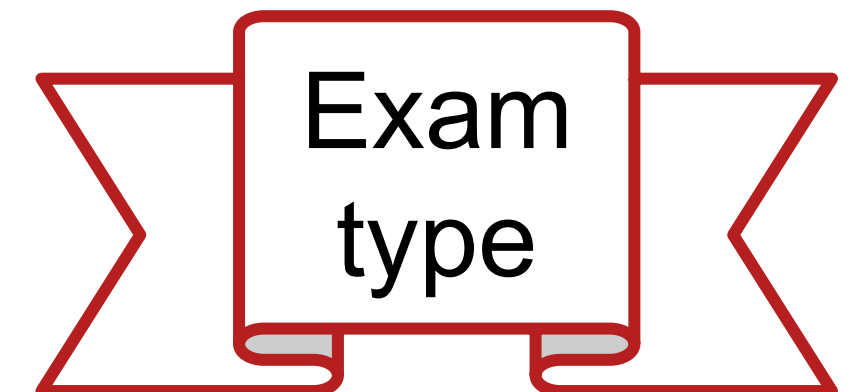
URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

The training of the LLM “BreezeTalk” took 3 months using 100% of the resources available on a 10-server cluster.

Each server has an embodied footprint of 1200 kg CO<sub>2</sub>e and a 3-year lifespan.

What share of embodied footprint should be allocated to BreezeTalk (training only), in kg CO<sub>2</sub>e?

- a. 100
- b. 1000
- c. 3000
- d. 12000



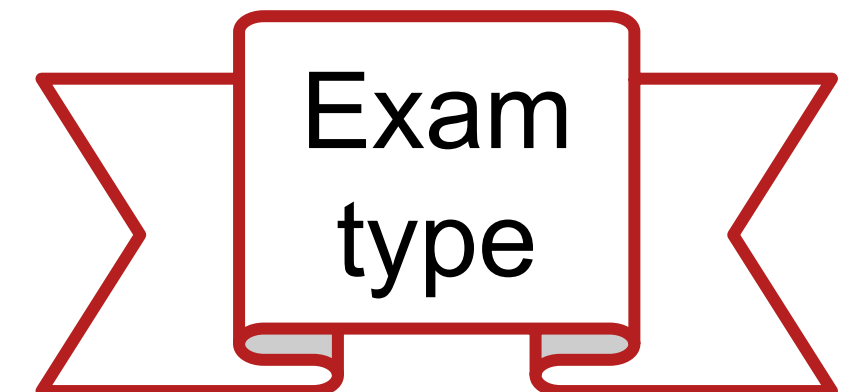
# Meditation app

URL: ttpoll.eu  
Session ID: cs290

ZenPath is an app dedicated to mental well-being that offers guided meditation sessions online. To reduce user dropout, they decide to display a popup after a user skips two sessions where the “Resume Today!” button is preselected.

What type of nudging technique is most likely used here?

- 0% a. Opt-in
- 0% b. Social proof
- 0% c. Scarcity
- 0% d. Default



# E-commerce platform

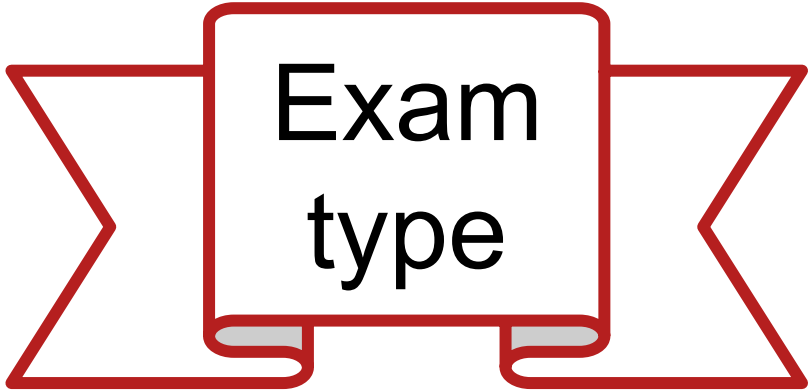
URL: ttpoll.eu

Session ID: cs290

The e-commerce platform Shine would like to implement new features to improve the experience of its various categories of users. Here is the list of envisaged features.

Which of them best matches the definition of a deceptive pattern?

- 0% a. Personalize style recommendations based on past browsing
- 0% b. Display user-provided past purchase data to recommend sizes
- 0% c. Register users to a ShineClub membership trial on checkout
- 0% d. Provide downloadable QR codes for the free return of items



Exam  
type

# Beer brewing dataset

URL: [ttpoll.eu](http://ttpoll.eu)

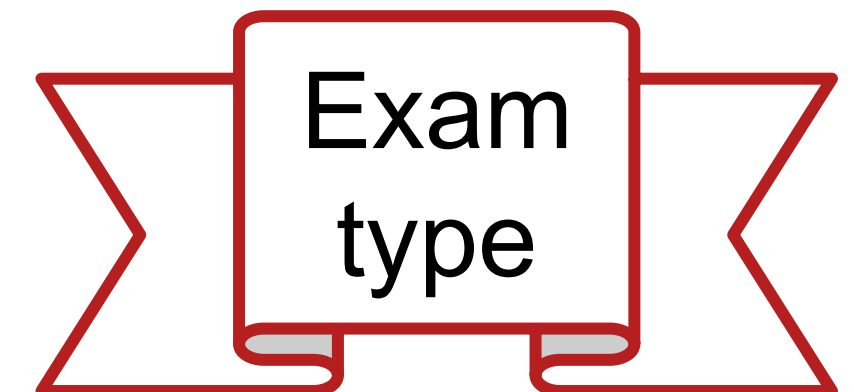
Session ID: cs290

One of the results of your Bachelor thesis is a very cool dataset which contains tasting profiles and consumer reviews for 3197 unique beers from 934 different breweries. This dataset can be used to train machine learning models for sentiment analysis and classification tasks.

You have created a datasheet for your dataset.

Which of the FAIR principles do you follow by providing a datasheet?

- 0% a. Findable
- 0% b. Accessible
- 0% c. Interoperable
- 0% d. Reusable



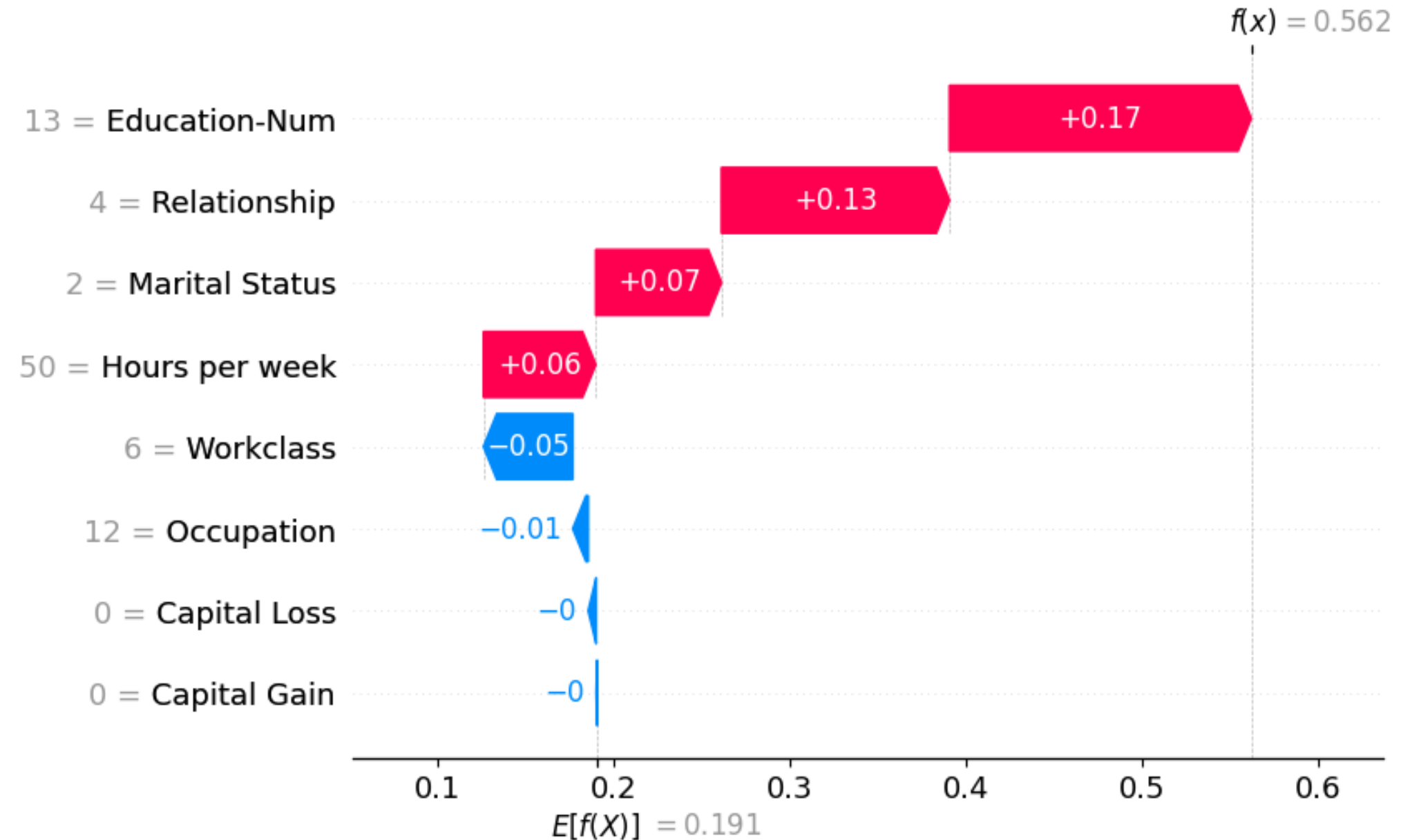
# Loans (again)

New

URL: [ttpoll.eu](http://ttpoll.eu)  
Session ID: cs290

The plot on the right displays the SHAP values obtained for the prediction generated by our ML model for customer 1113.

What does this plot represent in terms of interpretability method?



- a. A local explanation
- b. A global explanation
- c. A feature importance analysis
- d. A feature correlation analysis