



Teacher : Dr. Cécile Hardebolle CS-290: Responsible software

 $\frac{28}{10}/\frac{2024}{90}$ minutes

Student 1

 $\mathrm{SCIPER} \colon 999000$

Do not turn the page before the start of the exam. This document is double-sided, has 17 pages, the last ones possibly blank. Do not unstaple.

- Place your student card on your table.
- No paper materials other than one (1) A4 sheet of notes recto-verso is allowed to be used during the exam.
- Using a **calculator** or any electronic device is not permitted during the exam.
- First part : single choice questions (12 questions, 12 points)
 - for the singles choice questions, we give :
 - +1 points if your answer is correct,
 - 0 points if you give no answer or your answer is incorrect.,
- Second Part: true/false questions (4 questions, 4 points)
 - for the true/false questions, we give :
 - +1 points if your answer is correct,
 - 0 points if you give no answer or your answer is incorrect.
- Third part : case studies (3 questions, 20 points)
 - for each question, the number of points is noted above each question. Leave the checkbox empty.
- Use a black or dark blue ballpen and clearly erase with correction fluid if necessary.
- If a question is wrong, the teacher may decide to nullify it.

Respectez les consignes suiva	ntes Observe this guidelines Beachten Sie bitte	e die unten stehenden Richtlinien				
choisir une réponse select an answer Antwort auswählen	ne PAS choisir une réponse NOT select an answer NICHT Antwort auswählen	Corriger une réponse Correct an answe Antwort korrigieren				
ce qu'il ne	ICHT tun sollte					





First part: single choice questions

For each question, mark the box corresponding to the single correct answer.

O .	-4
Question	- 1
& account	_

Question 1
A group of computer scientists with similar background, all experts in AI-based software development, is aware of cognitive biases. They are starting a new AI project for healthcare and they aim to minimize the
impact of these biases when making design decisions.
Select the strategy they should use (best answer):
slow down the decision-making processes
systematically include all members of their group to increase heterogeneity
choose one or two of them to play the devil's advocate
\square systematically include all members of their group to apply a participatory design method
Question 2 The CEO of a tech company stated in the media: "In the past, we've invested in technology to positively impact people's lives, and we have no intention of changing that strategy in the future - technology remains the best alternative." We may interpret this as: (select the best answer)
Sunk cost fallacy
Source cues
System 1 thinking
☐ Illusory truth
Question 3 A start-up developed a machine learning model designed to connect people based on their personal interests. A big company has then bought the start-up and is currently using the algorithm to connect jobseekers with employers. It is a case of
deployment bias
aggregation bias
measurement bias
intersectional bias
Question 4 A company has developed a complex algorithm to predict whether athletes suspected of doping actually do it. A positive result means that the algorithm classifies the athlete as at risk of doping, while a negative result means no risk of doping. The system has been used for 5 years and we have access to data about athletes that were indeed caught for doping. We found that the proportion of athletes predicted to dope amongst all predictions is higher for men rather than for women. The fairness metric we have used is:
error rate balance
conditional use accuracy equality
equal accuracy
demographic parity



... ethical sensitivity

You develop a software that analyzes the weather forecast to send the population a notification in case of upcoming extreme rain (positive result). In this context: True Positive = rain is predicted and the prediction is correct True Negative = no rain is predicted, and the prediction is incorrect False Positive = rain is predicted and the prediction is correct False Negative = no rain is predicted, and the prediction is correct Question 6 You develop a software that analyzes the weather forecast to send the population a notification in case of upcoming extreme rain (positive result). The False Negative Rate (FNR) is: The number of times rain is predicted among all times it actually didn't rained The number of times no rain is predicted among all times it actually didn't rained The number of times no rain is predicted among all times it actually rained The number of times rain is predicted among all times it actually rained Question 7 A dilemma is... ... a situation in which you have to decide between two alternatives using a coin flip (better to leave things to chance) ... a situation in which you have to weigh the pros and cons of each decision (and their consequences), with no decision 100% perfect or 100% unperfect ... a situation in which you have to weigh the pros and cons of each decision (and their consequences) and choose the one with the higher number of pros. ... a situation in which you should escalate the decision to your management line. You work on a chatbot to provide students assistance on campus questions. At evaluation Question 8 time it generates plausible nonsense with a 15% rate. You probably face what we called an... ... ethical issue ... ethical dilemma ... ethical blindness





Question 9

Here are three variables:

- Disinformation spread
- Public trust in information
- ullet Development of disinformation software

We know that:

- As the spread of disinformation increases, the public trust in information decreases
- As the public trust in information decreases, bad actors see a growing opportunity to develop disinformation software exploiting this mistrust

In a causal loop diagram representing the dynamics between these variables, we would have (select the correct answer):

correct answer):
☐ The arrow between "Public trust in information" and "Development of disinformation software" has a positive sign and the arrow between "Development of disinformation software" and "Disinformation spread" has a negative sign.
The arrow between "Public trust in information" and "Development of disinformation software" has a negative sign and the arrow between "Development of disinformation software" and "Disinformation spread" has a positive sign.
The arrow between "Public trust in information" and "Development of disinformation software" has a positive sign and the arrow between "Development of disinformation software" and "Disinformation spread" has a positive sign.
☐ The arrow between "Public trust in information" and "Development of disinformation software" has a negative sign and the arrow between "Development of disinformation software" and "Disinformation spread" has a negative sign.
Question 10 Fill the blanks: If a piece of software behaves in a way at first glance, but puts people of at, then it
is a case of discrimination.
negative / several groups / an advantage / direct
positive / identified groups / an advantage / inverse
neutral / specific groups / a disadvantage / indirect
neutral / several groups / a disadvantage / direct
negative / specific groups / a disadvantage / indirect
Question 11 A bad actor launched a phishing attack on employees of Swiss public institutions to steal their login credentials. An online media outlet reported on it, with the most upvoted comments on the article criticizing the institutions for their inability to counter online threats, harming their reputation. The harm to reputation can be classified as an impact that is:
Direct
☐ Both direct and indirect
☐ Neither direct nor indirect
☐ Indirect





Question 12

Question 12	
Imagine that you develop software for people from a single country.	If you nonetheless envision cultural
differences you are probably using the (select the best answer)	
edge case strategy	
stride strategy	
bad actor strategy	
people behind the data strategy	





Second part: true/false questions

For each question, mark the box (without erasing) TRUE if the statement is **always true** and the box FALSE if it is **not always true** (i.e., it is sometimes false).

You found a dataset with 5 variables, all self-reported by participants: eye-color, extraversion and 3 health-related variables. When analyzing the data you identify that:

- there are positive and substantial correlations among the 3 health variables
- there is a positive and substantial correlation between eye-color and extraversion
- there is no correlation between eye-color and the health variables

Question 13	Eye-color is a latent variable	
	☐ TRUE	FALSE
Question 14	Eye-color is a proxy for health	
	TRUE	FALSE
Question 15	Eye-color is a sensitive attribute	
	TRUE	FALSE
Question 16	Extraversion is a latent variable	
	TRUE	☐ FALSE





Third part: case Studies

Answer in the empty space below. Your answer should be carefully justified, and all the steps of your argument should be discussed in details. Leave the check-boxes empty, they are used for the grading.

Question 17: Case 1: Harms modeling - Social assistant chatbot This question is worth 5 points.



Scenario:

In the realm of technological innovation, a revolutionary social-assistant chatbot emerges, designed to offer guidance on relationships. This cutting-edge human-centered AI, inspired by Snapchat's AI chatbot, aims to become an indispensable part of people's lives. Sarah and James are two individuals with contrasting lives. James, a young artist, craves genuine connections with like-minded people, while Sarah, a young consultant, struggles to balance her career with her personal life. Sarah and James turn to this chatbot for relationship advice. Its sophisticated algorithm analyzes their preferences, communication styles, and social behaviors to offer tailored suggestions for interactions. In addition to helping them identify others' emotions, it provides them with conversation starters and even helps plan memorable dates. As the chatbot gains traction and spreads throughout society, it becomes an integral part of society's social, economic, and political landscape. It reshapes how people approach dating and relationships, influencing not only their personal lives but also impacting the dating industry, advertising strategies, and even political campaign tactics. In addition, companies rely on the chatbot to predict the emotions of their staff and their clients to maximize their benefits. Yet, there are those who remain skeptical of the chatbot's far-reaching influence. Some individuals, wary of data privacy concerns and the potential for manipulation, opt to abstain from using the technology. They seek more traditional avenues for forming connections, believing in the value of genuine human interactions and the potential risks that come with relying on AI for personal advice.

Task:

Considering the following extract of the harms modeling table, describe what should go in the different cells:

- For cells A, B, C and E: describe 1 harm that corresponds to the category
- For cell D: indicate the corresponding harm category

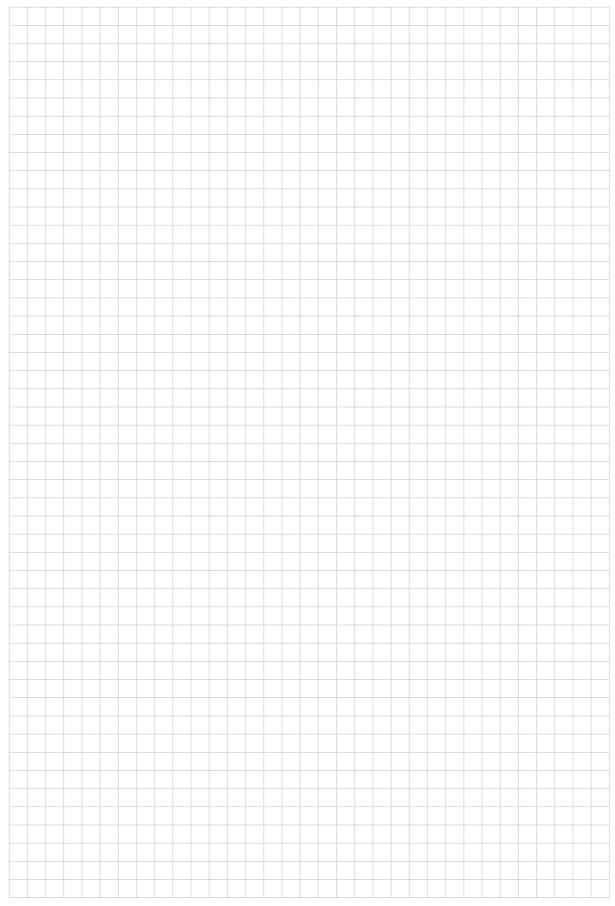
Make sure to identify your answers with the corresponding letters [1 point / answer].

Category	Type of harm	Social assistant chatbot
Humans	Physical injury	A)
Allocation of Resources	Opportunity loss	В)
Human Dighta	Liberty loss	C)
Human Rights	D)	Most intimate feelings are now "public"
Social System Harms	Social detriment	E)





+1/8/53+







Question 18: Case 2: Values analysis - Personalized deals This question is worth 5 points.

05152	

Scenario:

A webshop manager wants to offer interesting deals to the shop's customers, and thinks that it would be best to offer personalized deals to each one of them. As the customers provide their email address when registering, the manager creates the following script: for each user, the script finds some account linked to the mail address (Facebook, YouTube, Amazon, Retail stores, etc.) and buys the data related to that user. With that data, a personalized offer containing deals adapted to the centers of interest of the user is sent directly by email.

Task:

Your overall task is to perform an analysis of the values and value tensions involved for the different stakeholders in the case. Follow the 2 steps below:

1 [4 points] Consider some stakeholders in the case and identify 2 values which are supported by the software (= 2 value-based benefits) and 2 values that are opposed by the software (= 2 value-based harms).

Consider the value-based benefit/harm table template below and describe what would go in each cell for each of the values you identified:

- (A) Describe the stakeholder
- (B) Name the value (you should use the names in Appendix 3.1) and explain in your own words what the value means for this stakeholder
- (C) Indicate if the value is supported (value-based benefit) or harmed (value-based harm) for this stakeholder
- (D) Justify why it is supported / harmed by the software

Make sure to identify your answers with the corresponding letters.

A list of Schwarts's values is provided in appendix 3.1.

2 [1 point] Draw a value-based tension map showing a value tension and provide an explanation of the tension.

Stakeholder	Key Value	Benefits	Harms	Justification
Stakeholder: (A)	Value name and description: (B)	Benefit or	Harm: (C)	It's a value-based benefit/harm for this stakeholder because: (D)





Appendix 3.1

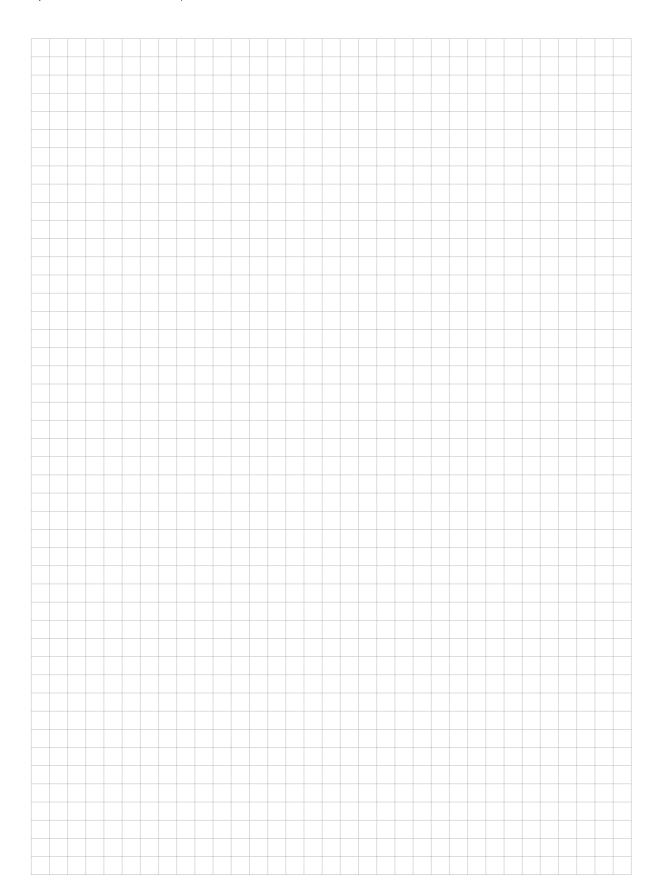
	Power Resources	Power through control of material and social resources					
Self- enhancement	Power Dominance	Power through exercising control over people					
	Achievement	Personal success through demonstrating competence according to social standards					
	Hedonism	Pleasure and sensuous gratification for oneself					
	Stimulation	Excitement, novelty, and challenge in life					
Openness to change	Self-direction Action	The freedom to determine one's own actions					
	Self-direction Thought	The freedom to cultivate one's own ideas and abilities					
	Universalism Tolerance	Acceptance and understanding of those who are different from oneself					
Self-	Universalism Concern	Commitment to equality, justice, and protection for all people					
transcendence	Universalism Nature	Preservation of the natural environment					
	Humility	Recognizing one's insignificance in the large scheme of things					
	Benevolence Dependability	Being a reliable and trustworthy member of the in-group					
	Benevolence Caring	Devotion to the welfare of in-group members					
	Tradition	Maintaining and preserving cultural, family, or religious traditions					
	Conformity Interpersonal	Avoidance of upsetting or harming other people					
Conservation	Conformity Rules	Compliance with rules, laws, and formal obligations					
	Security Societal	Safety and stability in the wider society					
	Security Personal	Safety in one's immediate environment					
	Face	Security and power through maintaining one's public image and avoiding humiliation					

Table 1: Source: Schwartz et al. (2012).





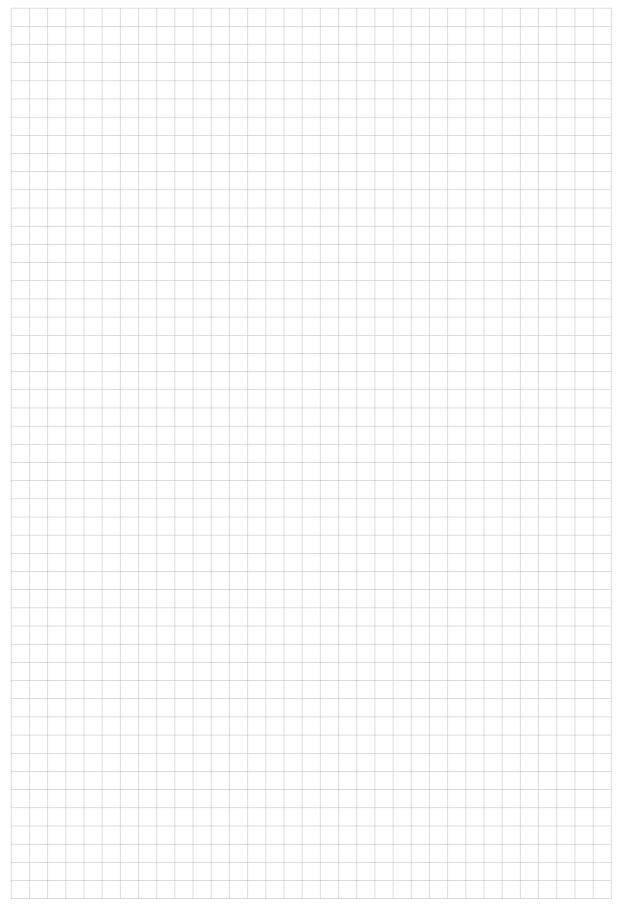
1) Value-based benefits/harms







+1/12/50+







2) Value-based tension map

-																
-																
ŀ																
-																





Question 19: Case 3: Universal Digital Identity Platform This question is worth 10 points.

|--|--|

Scenario:

Consider this made up very futuristic case study: You are part of an international team tasked with developing a Universal Digital Identity Platform (UDIP). This platform is intended to be the ultimate authentication system and replace the paradigm of having one account for each service we use. UDIP provides every individual a unique digital identity, which can be used globally for accessing various services such as banking, healthcare, education, and government services. The platform will utilize biometric data, including facial recognition and fingerprints, to ensure secure and accurate identification and authentication. The aim is to streamline access to services, reduce fraud, and enhance global connectivity. Governments and private companies worldwide are eager to adopt this system to improve efficiency and security. The platform has the potential to become a foundational technology, potentially affecting billions of people.

Task:

As the ethics referee of the team, you are asked to anticipate potential consequences of the deployment of the platform in terms of safety and fairness. Follow the 3 steps below:

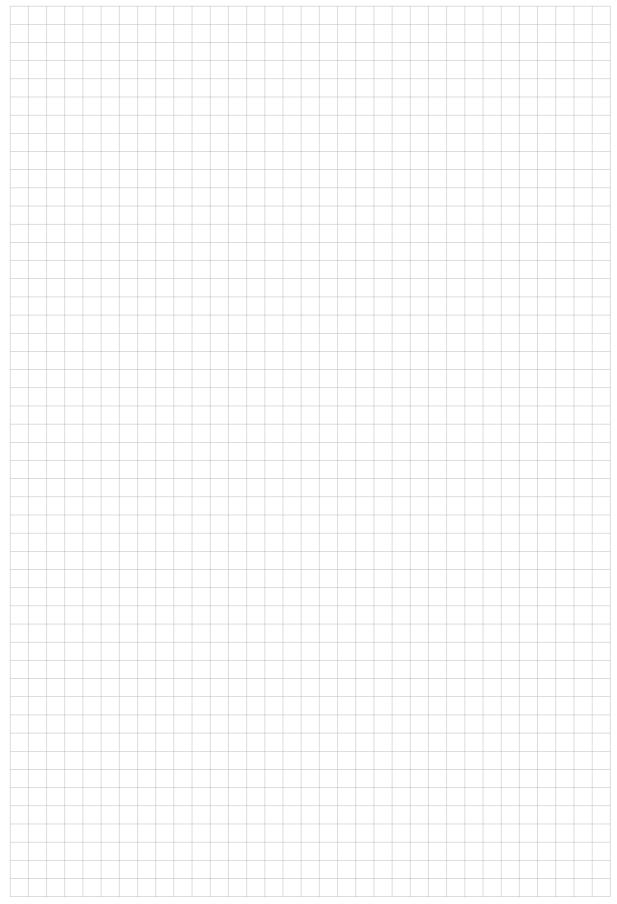
- 1 [1 point] Name one strategy seen in the course that you can apply for this task. Warning: you cannot use "Harm Modeling" for this case.
- 2 [3 points] Explain the strategy: Justify why this strategy is appropriate for this task.
 - (a) Describe briefly how to apply this strategy.
 - (b) Describe the result of applying the strategy.
- 3 [3 points] Present one safety issue you identify in the case
- 4 [3 points] Present one fairness issue you identify in the case
 Specify any assumption you make (that is not clearly stated in the scenario) about the system and its
 stakeholders.



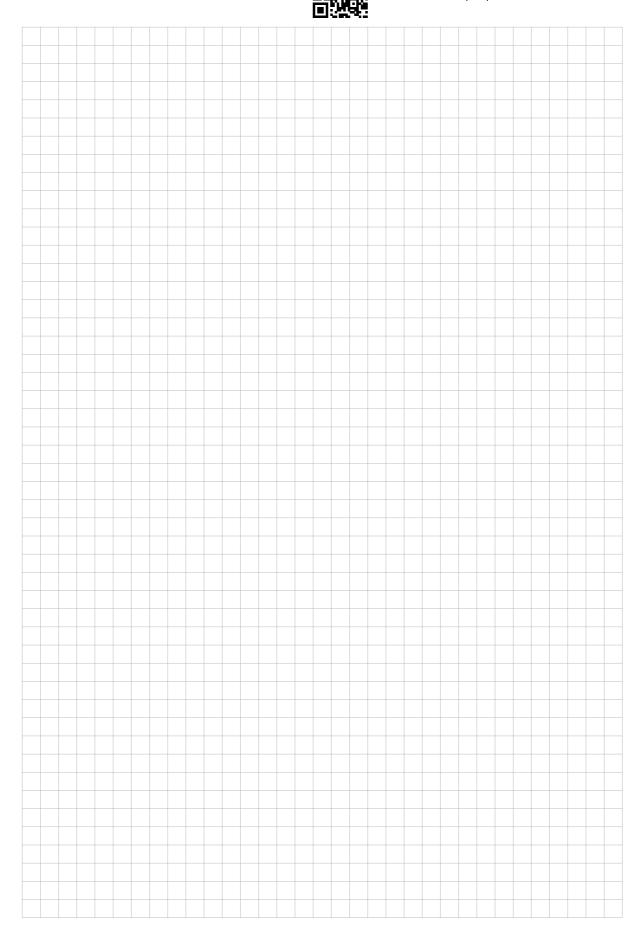




+1/15/46+











+1/17/44+

