

Microsoft Responsible AI Impact Assessment Guide

FOR EXTERNAL RELEASE

June 2022

This guide accompanies the Microsoft Responsible AI Impact Assessment Template. We are releasing it externally to share what we have learned, invite feedback from others, and contribute to the discussion about building better norms and practices around AI.

We invite your feedback on our approach:

<https://aka.ms/ResponsibleAIQuestions>



Contents

Introduction

Impact Assessment: Guiding principles	3
Case study.....	4

01 Project Overview

1.1 System profile.....	6
1.2 System lifecycle stage.....	6
1.3 System description.....	6
1.4 System purpose	7
1.5 System features.....	7
1.6 Geographic areas and languages.....	8
1.7 Deployment mode.....	8
1.8 Identifying system uses.....	8

02 Intended uses

2.1 Assessment of fitness for purpose.....	14
2.2a Stakeholders, potential benefits, and potential harms.....	15
2.2b Brainstorming potential benefits and harms	18
2.3 Stakeholders for Goal-driven requirements from the Responsible AI Standard	24
2.4 Fairness considerations	26
2.5-2.8 Technology readiness assessment, task complexity, role of humans, and deployment environment complexity	27

03 Adverse Impact

3.1 Restricted Uses.....	29
3.2 Unsupported uses.....	29
3.3 Known limitations.....	30
3.4 Potential impact of failure on stakeholders	31
3.5 Potential impact of misuse on stakeholders.....	32
3.6 Sensitive Uses	32

04 Data Requirements

4.1 Data requirements	34
4.2 Existing data sets.....	34

05 Summary of Impact

5.1 Potential harms and preliminary mitigations	36
5.2 Goal applicability	42
5.3 Signing off on the Impact Assessment	42

INTRODUCTION

Responsible AI Impact Assessment Guide

This resource is intended to support your team as you work through the Responsible AI Impact Assessment, as required by the v2 Responsible AI Standard. It will help frame your conversations about Responsible AI with:

- Guidance
- FAQs
- Examples
- Activities
- A case study
- Prompts

Discuss the core concepts of the Impact Assessment with your team and document the highlights in the Impact Assessment Template.

Preparing for activities

Schedule at least one session for each of the activities. Invite members of your team with different expertise to discuss the core concepts. Then, document the highlights of your conversation in the Impact Assessment Template.

Guiding principles

FOR COMPLETING AN IMPACT ASSESSMENT

Exploration & Innovation

It is important that your team uses the Impact Assessment as an opportunity to pause and explore responsible AI challenges. Aim for a thorough investigation – the more you document, the easier it will be for potential reviewers to engage with your challenges, help you find mitigations, and work with you to develop innovative solutions in the future.

Collaborative work

Some parts of the Impact Assessment require teamwork. Review the following sections and plan your team's collaboration:

- 1.8** Defining intended uses
- 2.2** Creating a stakeholder list, identifying potential benefits and harms
- 2.4** Fairness considerations
- Section 3** Adverse Impact
- 5.1** Potential harms and preliminary mitigations



Imagine an AI system that optimizes healthcare resources

Case Study

This guide uses a case study to illustrate how teams might use the activities to complete the Impact Assessment Template.

Consider an AI system that optimizes healthcare resources such as the allocation of hospital beds or employee scheduling. The system makes predictions about how long a patient will stay in the hospital to inform both bed allocation and staff scheduling. The system's input includes both patient medical data and resource constraints like scheduling parameters. The system was trained on data from a specific hospital that included: 1) patients' medical history; 2) the type of surgery; 3) how long they stayed in the hospital; and 4) historic staffing and scheduling data.

The system has two intended uses. First, the system can be used by hospital staff to manage the allocation of hospital beds. Second, the system allows hospital administrators to automate scheduling for nurse shifts. For the rest of the guidance and activities, we focus on the first intended use, hospital staff using the system to manage the allocation of hospital beds. We call this system the Hospital Employee and Resource Optimization System (HEROS).

Section 1

Project overview



System profile and system lifecycle stage



System description, purpose and features



Geographic areas, languages and deployment mode



Intended uses

1.1

System Profile

Guidance

This section describes basic information about the system being evaluated in the Impact Assessment; tracks authors and updates; and establishes who will review the Impact Assessment when it is completed. Throughout the Impact Assessment when we refer to “the system,” we are referring to the system described in this section.

1.2

System lifestyle stage

Guidance

Your responses here will help potential reviewers understand the overall release timeline for the system.

1.3

System description

Guidance

Your response here should help potential reviewers understand what, exactly, you’re building. Describe what kind of AI capabilities the system has.

Use simple language and be specific, avoiding vague concepts as much as possible.

Write for an audience that has a basic understanding of AI systems but no understanding of this specific system.

Prompts

- What are you building?
- What does it do?
- How does it work?

1.4

System purpose

Guidance

In 1.3 you described what you are building, here your response should help potential reviewers understand why you're building this system.

Prompts

Focus on the why.

This statement should include:

1. the end user or primary customer of the system,
 2. how they complete this task today, or their current situation,
 3. the value that the system is intended to deliver,
 4. how it improves on today's situation.
-

1.5

System features

Guidance

This section should help potential reviewers understand the specific features (capabilities) of the system and how the system being evaluated in this impact assessment relates to existing systems or features.

In the section, please describe the system's features and capabilities overall, not the features of specific models that the system may use.

FAQs

What do 'system features' mean in this context?

System features are the functionalities or capabilities you use within a system to complete a set of tasks or actions.

What if there are no existing features because the system is completely new?

Do not complete 'existing features' section.

1.6

Geographic areas and languages

Guidance

When planning the system, knowing the geographic areas where it will or won't be deployed and supported languages helps to uphold fairness by aiding the identification of relevant demographic groups, languages, and other contextual details requiring consideration to ensure fairness concerns are identified and mitigated.

When describing supported languages, include both language and region (e.g., British English).

FAQs

What if the system does not use natural language processing?

In this case you should leave the languages section blank.

1.7

Deployment mode

Guidance

This section should help potential reviewers understand how the system will be deployed to users or customers.

Examples

Some possible responses could include:

- Online Service
 - Platform Service
 - Code
 - On Premises
 - Container
-

1.8

Identifying system uses: Activity

Steps

1) Brainstorm possible uses

Start by listing as many potential ways someone could use the system, no matter how outlandish they seem at first.

2) Categorize possible uses: Intended, unsupported or misuse

In the Impact Assessment we ask about several different types of uses. For each possible use you listed in step one, determine if the use is an intended use, unsupported use, or misuse of the system.

3) Check if any uses are also Sensitive Uses or Restricted Uses

Check all uses against the definitions for Restricted Uses and Sensitive Uses. Follow the guidance for any Restricted Uses. Report any Sensitive Uses to the Office of Responsible AI.

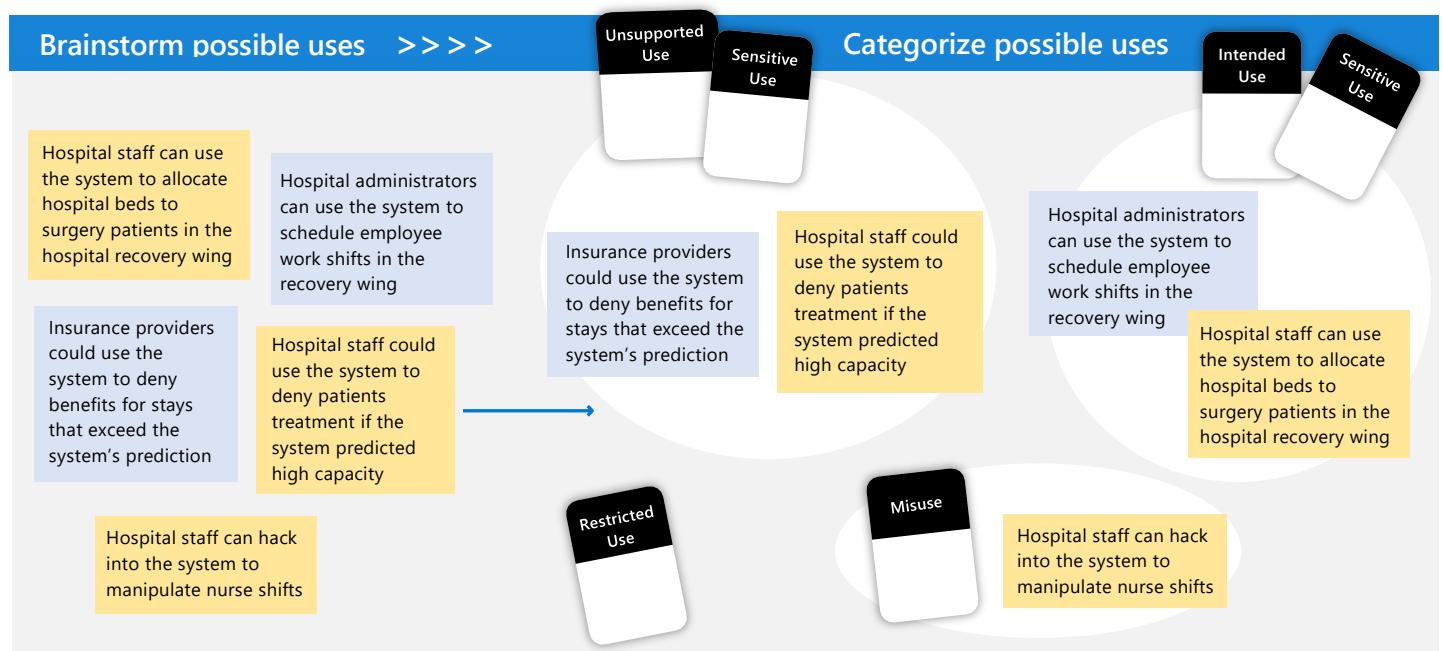
Case Study: Hospital Employee and Resource Optimization System (HEROS)

Focusing on the system's anticipated features and models, the team brainstormed some of the possible ways the system could be used. After listing several different uses, the team looked at the prompts below to determine if the uses are intended uses, unsupported uses, or misuses.

Two of the uses (hospital staff denying patients treatment due to system predictions and insurance providers denying claims based on system predictions) are uses the team is not designing or testing the system for, so they are considered unsupported uses.

Two of the uses (hospital staff using the system to allocate hospital beds and hospital administrators using the system to schedule work shifts) are uses the team is intentionally designing and evaluating the system for, so they are intended uses.

One of the uses (hospital staff can hack into the system to manipulate nurse shifts) is a misuse of the system. Lastly the team compared all uses with the Restricted Uses and the triggers for Sensitive Uses. None of the uses are Restricted Uses according to the current guidance. But, after reviewing the definition of Sensitive Uses, the team determined that all uses are Sensitive Uses. They submitted the system for review.



Categories of use

Intended use

Uses for which the system is designed and tested. An intended use is a statement of who is using the system, for what purpose or task, and where they are when they're using the system.

E.g., HEROS has two intended uses: 1) for hospital staff to allocate beds to surgery patients in the hospital and 2) for hospital administrators to schedule employee shifts in the hospital based on predicted occupancy.

For more on intended uses, see 1.8.

Misuse

All systems have the potential to be intentionally or unintentionally misused. Consider ways that someone could hack the system, use it to cause harm, or mistakenly use the system incorrectly.

E.g., Hospital staff can hack into HEROS to manipulate nurse shifts.

For more on intended uses, see 3.6

Unsupported use

Reasonably foreseeable uses for which the system was not designed or evaluated or that we recommend customers avoid.

E.g., An unsupported use of HEROS includes hospital staff using the system to deny patients treatment if the system predicted high capacity.

For more on unsupported uses, see 3.2

Restricted uses and sensitive uses

Restricted Use

Restricted Uses are uses of AI which are subject to specific restrictions (typically on development and deployment). They are defined by the Office of Responsible AI and updated periodically. If any of the uses of the system are Restricted Uses, follow the guidance for that Restricted Use.

Sensitive Use

Some uses of AI systems are particularly sensitive and impactful on individuals and society. If any uses or misuses of the system match any of the triggers for Sensitive Uses, report them to the Office of Responsible AI.

Intended uses

Guidance

We begin an Impact Assessment from an intended use to understand the system's socio-technical context. The impact of the system on people and society will depend on its intended uses (the uses that your team plans to design and test for).

Intended uses may also be referred to as use cases or scenarios.

Intended uses are statements that include:

1. who the end user or primary customer is,
2. what they will achieve with the system, and
3. where they are when they use it.

Systems that have a single end user and fulfill a single purpose might only have one intended use.

Other systems may have multiple intended uses (refer to the case study as an example).

FAQs

What is the difference between intended uses and system features?

Intended uses are not system features – features focus on what the system's functionalities are, while intended uses focus on how people will use them to accomplish their goals in particular settings. A single system feature may have multiple intended uses.

What if my system can be used in many different settings?

Remember the focus is on the uses your team is designing and testing for, not all the places someone could use a system. Think about the settings your team talks about most often. It's also okay to use a high-level setting like "at work", "at home", or "online."

What if the system is a platform technology?

If the system is a platform technology that could be used in many different settings, list one or a few examples to help you think through stakeholders, harms, and benefits in this section.

What if the intended uses are very similar?

Our goal with defining intended uses is to assess the impact the system may have on people and society. When thinking about whether or not an intended use should be distinct, consider the following:

- Similar end users or customers who use the system for similar tasks and purposes in similar settings can be grouped into a single intended use.
- Different tasks, purposes, and goals should be separate intended uses as they will likely result in different types of potential harms.
- While a system could potentially be used in many different settings, you may not need a separate intended use for every possible location. Consider the locations your team is designing and testing for, these are the locations you should include explicitly in an intended use.

Prompts

Elements of an intended use:

- End user (who)
- Purpose/task/goal (for what?)
- Setting (where)

Examples

An AI system that can identify an individual using a biometric scan of the individual's face.

Intended uses:

1. Allow an organization to grant physical access to a controlled space such as in an office.
2. Provide a personalized experience for people in a physical space like a bank or hotel lobby.

HEROS (case study)

Intended uses:

1. Allows hospital staff to use predictions of how long a patient will stay in the hospital to manage the allocation of beds.
2. Allows administrators to automate scheduling for nurse shifts in the hospital.

Repeating the intended uses section

Guidance

For each intended use of the system, copy and paste the section titled Intended Use #1 [Name of Intended Use].

Repeat questions 2.1-2.8 for each intended use.

Examples

HEROS has two intended uses, so in the Impact Assessment template we repeat this section two times. The first time, we answer questions 2.1-2.8 for intended use #1. The second time, we answer questions 2.1-2.8 for intended use #2.

Throughout the rest of this guide we provide examples for the core concepts, often using intended use #1 of HEROS:

Hospital staff can use the system to allocate hospital beds to surgery patients in the hospital recovery wing.

Section 2

Intended uses



Assessment of fitness for purpose



Stakeholders, potential benefits & potential harms



Stakeholders for goal-driven requirements



Fairness considerations



Technology readiness assessment, task complexity, role of humans, and deployment environment complexity

2.1

Assessment of fitness for purpose

Guidance

Your responses here will help potential reviewers understand how the system effectively solves the intended problem posed by each intended use, recognizing that there may be multiple valid ways in which to solve the problem.

Examples

A solution that is not fit for purpose might:

- Be trained on data in one language and then deployed to other languages without retraining. For example, a system trained to detect spelling errors on American English may not be appropriate for other varieties of English.
- Have outputs that do not represent the target phenomena. For example, images of people smiling or not smiling do not realistically represent emotion. Emotion is an internal state, while smiling is not.
- Be developed under the premise that it will make a process more efficient, but without evidence showing that it actually solves the problem it's meant to, or that it improves efficiency.

Prompts

- What is the problem to be solved for with this intended use? What evidence will establish that the system is fit for this intended use?
- How is the problem currently solved?
- Why will this system be a better solution than other approaches to solving the same problem?
- What evidence will demonstrate that the system is a better solution than other approaches?
- What are the system outputs intended to represent?
- What are the limitations to their realistic representation?
- What is the justification for using these outputs despite these limitations?
- What evidence will demonstrate the validity of this representation?

FAQs

What types of evidence are considered acceptable?

Some examples of acceptable evidence include:

- Stakeholder conversations, such as input from domain experts on proposed concepts and designs.
- Results of systematic stakeholder research, e.g.:
 - Focus groups or user panels
 - Ethnographic studies or shadowing
 - Iterative concept testing
 - Surveys
 - Experiments

- Marketing reports or trend analyses
 - System analytics
 - Other data sources
-

2.2a

Stakeholders, potential benefits, and potential harms

Guidance

The stakeholder benefits and harms table is one of the most important elements of your Impact Assessment. It is key to understanding the potential impact of the system on people.

In this section we have prompts to identify stakeholders from two broad categories: direct and indirect stakeholders. The specific category the stakeholder belongs to is not necessarily important. These categories are useful for identifying a broad range of stakeholders who may be impacted by the system.

Direct stakeholders include people who interact with the system directly. They can be system owners, primary users, secondary users, decision subjects or data subjects and malicious actors.

Indirect stakeholders are affected by the system but, unlike direct stakeholders, do not engage with the system itself. Indirect stakeholders can include bystanders, people responsible for decision subjects or data subjects (such as parents), society at large, or communities who may be affected by the system but don't use it.

FAQs

How many stakeholders?

Try to find at least one stakeholder per prompt. Consider whether each stakeholder prompt applies to your system even if you don't think it does at first glance.

How specific should the stakeholders be?

Try to be as specific as possible when describing your stakeholders. If you are working with a platform technology and you find it hard to generate specific stakeholders, you can stay at a higher level - e.g., 'end user'.

What if I have a lot of stakeholders?

The stakeholder exercise is meant to be explorative so the more stakeholders you have the better your team can explore the potential impact the system may have. Make sure that the stakeholders are relevant and different. If the stakeholders are very similar, in terms of the role that they have in the system, try to group them.

What if one stakeholder has multiple roles?

In some cases, a specific individual or group may hold more than one stakeholder role. In this case, include the stakeholder only once in the table, but make a note of the different roles. This will help you think about benefits and harms in the next section.

Generating stakeholders: [Activity](#)

Steps

1) Intended use

Start by thinking about one specific intended use (as defined by the end user, purpose, and context). Who are the stakeholders in this scenario?

2) Stakeholder prompts

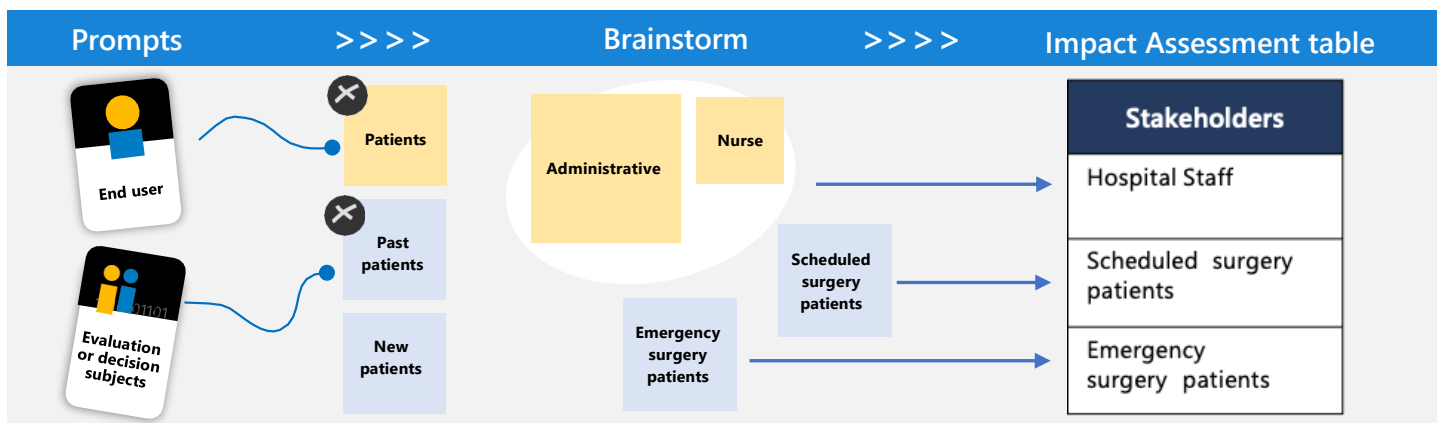
Use the prompts on the next page to generate stakeholders for the intended use.

You will repeat this for each intended use. While stakeholders may be similar across intended uses, potential harms and benefits are likely different.

Case Study: Hospital Employee and Resource Optimization System (HEROS)

Focusing on the intended use, the team used the stakeholder prompts (see next page) to brainstorm different stakeholders, then discussed how to record them in the stakeholder table in the Impact Assessment Template. The team started with the prompt for end user. One team member suggested that patients are the end user. A colleague disagreed because according to the prompt the end user is someone 'who is directly involved in using or operating the system' and this would be staff and not the patients. The rest of the team agreed and suggested 'Administrative' and 'Nurse' as the specific end users for this intended use. Since these two stakeholders perform the same role in the system for this intended use - interpreting hospital stay predictions and allocating beds - they were grouped into a single stakeholder 'Hospital Staff' and added to the stakeholder table in the Impact Assessment Template.

For the stakeholder prompt for evaluation subjects, the team brainstormed 'new patients', 'past patients', "emergency patients" and "scheduled surgery patients". The team decided not to include past patients because their data has been anonymized to train the system and they cannot be impacted by the system when it is in production. 'Emergency surgery patients' and 'scheduled surgery patients' were added as two different stakeholders because the context and model performance is different for these stakeholders. For emergency surgeries, staff may be more stressed (context), and the model may make less accurate predictions due to the diversity of illnesses and injuries and the lack of medical history for some emergency patients (performance). On the other hand, staff dealing with scheduled surgery patients enjoy a more routine environment (context) and the model makes more accurate predictions due to the consistent stream of patients undergoing scheduled surgeries (performance).



Generating stakeholders: Prompts

<div>End user (DIRECT STAKEHOLDER)</div> <div>Who will be most directly involved in using or operating the system? Who will have to interpret system outputs in order to make decisions? <i>E.g., marketing team, students</i></div>	<div>Evaluation or decision subjects (DIRECT STAKEHOLDER)</div> <div>Who will be evaluated or monitored by the system, whether or not by choice? Who will the system make predictions or recommendations about? <i>E.g., registered customer,</i></div>	<div>Oversight and control team (DIRECT STAKEHOLDER)</div> <div>Who will troubleshoot, manage, operate, oversee or control the system during and after deployment? Who can discontinue the system? <i>E.g., Microsoft, consumer customer, enterprise customer, B2B, B2C</i></div>	
<div>System owner or deployer (DIRECT STAKEHOLDER)</div> <div>Who will own and make decisions about whether to employ a system for particular tasks? Who develops and deploys systems that integrate with this system? <i>E.g., enterprise customer, Microsoft, hospital administrators</i></div>	<div>System builders or developers (DIRECT STAKEHOLDER)</div> <div>Who will be involved in the system design and development? <i>E.g., your team, customer dev team</i></div>	<div>Malicious Actors (DIRECT STAKEHOLDER)</div> <div>Who may intentionally misuse the system? <i>E.g., hackers</i></div>	
<div>Bystanders (INDIRECT STAKEHOLDER)</div> <div>Who in the vicinity of the deployed system may be impacted by its use? <i>E.g., passers-by</i></div>	<div>Regulators and civil society organizations (INDIRECT STAKEHOLDER)</div> <div>Who may advocate for regulation of this system or be concerned about compliance? <i>E.g., government health entities</i></div>	<div>Communities (INDIRECT STAKEHOLDER)</div> <div>Which communities may be affected by the short- or long-term use of the system? <i>E.g., communities with low digital literacy</i></div>	<div>Associated Parties (INDIRECT STAKEHOLDER)</div> <div>Who may have substantial interest in the system based on their relationship to other stakeholders? <i>E.g., company partners, family members</i></div>

2.2b

Brainstorming potential benefits and harms

FAQs

How many potential benefits and harms should I have per stakeholder?

For some stakeholders you might have a single benefit or harm and for another stakeholder you might have four or five. Other times you may not have any benefits to a stakeholder but several harms. Some stakeholders may have no benefits or harms.

What if a benefit or a harm applies to more than one stakeholder?

You may identify similar benefits or harms for different stakeholders. Think about any potential differences from the perspective of each stakeholder. If there aren't any significant differences in how the stakeholders may experience the benefit or harm, you can list the same harm or benefit for multiple stakeholders.

Why do we need to document potential benefits and harms?

This is an initial exploration that will help your team anticipate the impact that the system might have. A thorough exploration of potential harms can help you prevent problems before they happen.

Brainstorming potential benefits and potential harms: [Activity](#)

Steps

1) Stakeholder

Think about one stakeholder at a time.

2) How might they benefit?

Ideate potential benefits for a stakeholder by asking: How might they benefit directly or indirectly from interacting with the system?

3) Harms prompts




It is often difficult to think of ways the system might negatively impact stakeholders. Use the prompts on the following pages to help you brainstorm potential harms related to the Microsoft AI principles.

[Case Study: Hospital Employee and Resource Optimization System \(HEROS\)](#)

In the example below the team brainstormed potential benefits and harms for the stakeholder 'scheduled surgery patient', which is the evaluation or decision subject. The team started by asking how this stakeholder could benefit directly or indirectly from using the system. After brainstorming, they summarized the benefits as: better understanding of the length of hospital stay and better able to plan for things like childcare or house sitting while in recovery.

To think about the potential harms of the system the team looked at the harms prompts on the following pages. For the sake of this case study, only two examples of potential harms are listed below (but there may be more). The accountability prompt for human oversight and control led the team to think about how automatically allocating beds based on system predictions could cause harm. The transparency prompt for

system intelligibility started an interesting conversation about how knowledge of the outputs could affect hospital operations. The team discussed how not understanding the system well enough could result in patient care decisions that are not in the patient's best interest. The team documented each of these harms, aligned to the respective AI principles, in the stakeholders, potential benefits, and potential harms table as shown below.

Prompts >> Brainstorm >> Impact assessment table	Stakeholders	Potential system benefits	Potential system harms
	Refer to stakeholder activity	<i>How might this stakeholder benefit from the system?</i>	  <p>Not enough knowledge of systems limitations</p>
	 <p>Scheduled surgery Patient (evaluation or decision subject)</p>	<p>The patient will have information about how long they might remain in the hospital</p> <p>They can plan better around their recovery</p> <ul style="list-style-type: none"> Better understanding of length of hospital stay Better able to plan for things like childcare or house sitting while they're recovering 	<p>Automatic allocation = compromised patient care</p> <p>It is hard for staff to make decisions</p> <p>Poor decisions = compromised patient care</p> <p>Accountability - If the system automatically allocates hospital beds, without human oversight and approval, there may be instances where an inaccurate prediction compromises patient (decision subject) recovery.</p> <p>Transparency - If the hospital staff (decision maker) does not understand how to interpret system outputs, a patient (decision subject) could be allocated a bed for a shorter amount of time impacting their ability to fully recover.</p>

Brainstorming potential benefits and harms: Prompts

Accountability

POTENTIAL HARM

Significant adverse impacts

- Could the system impact legal position or life opportunities?
- Could the system uphold or become a threat to human rights?
- Could the system result in a risk of physical or psychological injury?
- Does the system meet the definition of a Restricted Use?

Accountability

POTENTIAL HARM

Fit for purpose

What harms might this stakeholder experience if the system does not effectively solve the intended problem?

E.g., If the system is unable to accurately predict the length of hospital stays for scheduled surgery patients (the intended problem), then decision makers will either make poor decisions based on the system outputs or stop using the system.

Accountability

POTENTIAL HARM

Data governance and management

What harms might this stakeholder experience if the data used to train the system have not been sufficiently managed or evaluated in relation to the system's intended use(s)?

E.g., If the system is trained using data from all types of hospital stays it may not accurately represent hospital stays specifically for scheduled surgery patients.

Accountability

POTENTIAL HARM

Human oversight and control

What harms might this stakeholder experience if the system is not subject to appropriate human oversight and control?

If the system automatically allocates hospital beds, without human oversight and approval, there may be instances where an inaccurate prediction compromises patient (decision subject) recovery.

Transparency

POTENTIAL HARM

Significant intelligibility

What harms might this stakeholder experience if there is not enough information to make appropriate decisions about people, using the system's outputs?

E.g., If the hospital staff (decision maker) do not understand how to interpret system outputs, a patient (decision subject) could be allocated a bed for a shorter amount of time impacting their ability to fully recover.

Transparency

POTENTIAL HARM

Communication to stakeholders

What harms might this stakeholder experience if they are unable to understand what the system can or cannot do?

E.g., If the 'hospital management' (system deployer) doesn't understand what the system was designed to do, they might use the system in a way it was not intended to be used. An unsupported use of this system would be denying patients (decision subject) treatment based on system predictions.

Transparency

POTENTIAL HARM

Disclosure of AI interaction

What harms might this stakeholder experience if they are unaware that they are interacting with an AI system when that system impersonates human interaction or generates or manipulates image, audio or video content that could falsely appear to be authentic?

E.g., While this harm does not apply to HEROS, it may apply to the system you are working on. It can be deceptive for an AI chatbot to handle conversations with a user without any indication of it being an AI system. It could also be deceptive for a system to create synthetic voice from input text with no indication that this is an AI system.

Fairness

POTENTIAL HARM

Quality of service

How might the system perform better or worse for different demographic group(s) this stakeholder might identify as?

While this harm does not apply to HEROS, it may apply to the system you're working on. AI systems may perform differently for different demographic groups. By evaluating system performance for different demographic groups, we can identify fairness harms to mitigate.

Fairness

POTENTIAL HARM

Allocation

Could the system recommend the allocation of resources or opportunities to a stakeholder differently based on their demographic group(s)?

E.g., patients from 'low-income backgrounds' (demographic group) might have a history of short hospital stays for financial reasons. This may result in predictions for shorter stays even when they have severe health needs, thereby impacting the patient's access to a hospital bed.

Fairness

POTENTIAL HARM

Minimization of stereotyping, demeaning, or erasing outputs

How might the system represent this stakeholder in ways that stereotype, erase, or demean them based on their demographic group(s)?

E.g., if the system uses a binary gender classification (male/female) it would erase other gender identities.

Reliability & Safety

POTENTIAL HARM

Reliability and safety guidance

What harms might this stakeholder experience if the system performs unreliably or unsafely (e.g., models are incompletely trained, system operates outside of the acceptable ranges, or performs with unacceptable error rates)?

E.g., If the system predicts the length of patient stay unreliably, then patients (decision subjects) might be in the hospital for longer than necessary, impacting their job security, or not long enough, impacting their recovery.

Reliability & Safety

POTENTIAL HARM

Failures and remediations

What harms might this stakeholder experience due to a predictable failure (e.g., false positives, false negatives, or other types of failures), or inadequately managing unknown failures once the system is in use?
E.g., In the event of an unexpected major system outage, the system may not work at all. This would be particularly problematic for the hospital in emergency situations.

Reliability & Safety

POTENTIAL HARM

Ongoing monitoring, feedback, and evaluation

What harms might this stakeholder experience related to system changes and operation after release, especially related to identification of issues, maintenance, and improvement over time?
E.g., It's possible that practices within the hospital shift over time, and a model trained on the original data set could become less accurate over time. Predictions would be less reliable, potentially compromising patient (decision subjects) care.

Inclusiveness

POTENTIAL HARM

Accessibility Standards

Refer to the standards for compliance.

Privacy and Security

POTENTIAL HARM

Microsoft Privacy Standard

Refer to the standard for compliance.

Microsoft Security Policy

Refer to the policy for compliance.

2.3

Stakeholders for Goal-driven requirements from the Responsible AI Standard

Guidance

Certain Goals in the Responsible AI Standard require you to identify certain types of stakeholders. When a Goal applies to the system, list the specific stakeholder(s) associated with that Goal to meet the requirement.

FAQs

What if a Goal does not apply to the system?

In this case you can enter "N/A" in the response area.

What if I already listed these stakeholders in the previous section?

If you used the stakeholder prompts earlier in this guide you likely already identified the Goal-based stakeholders. These specific questions are included in the Impact Assessment to ensure you meet the Goal-driven requirements for identifying stakeholders. It's okay to copy and paste the appropriate stakeholder from the previous section.

Examples

Goal A5: Human oversight and control

Stakeholders for Goal A5 include those who troubleshoot, manage, operate, oversee, or control the system during and after deployment. For these stakeholders, you should also describe their oversight and control responsibilities.

E.g., for intended use #1 of HEROS, the A5 stakeholders include the hospital administrators and their development team. Their responsibilities include overseeing the day-to-day operations of the system and troubleshooting issues that arise.

Goal T1: System intelligibility for decision-making

This Goal only applies to AI systems when the intended use of the generated outputs is to inform decision making by or about people.

Stakeholders for Goal T1 include those who use the outputs to make decisions and those about whom decisions are made.

E.g., for intended use #1 of HEROS, the T1 stakeholders include the hospital staff and patients.

Goal T2: Communication to stakeholders

Stakeholders for Goal T2 include those who will make decisions about whether to employ the system for particular tasks and those who develop or deploy systems that integrate with this system.

E.g., for intended use #1 of HEROS, the T2 stakeholders include hospital administrators.

Goal T3: Disclosure of AI interaction

This Goal only applies to AI systems that impersonate interactions with humans, unless it is obvious from the circumstances or context of use that an AI system is in use and AI systems that generate or manipulate image, audio, or video content that could falsely appear to be authentic.

Stakeholders for Goal T3 include those who will use or be exposed to the system.

This Goal does not apply to HEROS. But, for an AI system that impersonates a human in a chat bot, the T3 stakeholders might include the end user. For an AI system that creates synthetic audio from text for announcements in public spaces, the T3 stakeholders might include the general public.

2.4

Fairness considerations

Guidance

This section requires you to assess which of the Fairness Goals from the Responsible AI Standard apply to the system and identify which stakeholders should be considered for this Goal. After identifying the affected stakeholders you're asked to identify which demographic groups, especially marginalized groups, would be most at risk of experiencing a fairness harm.

To complete this section, please follow the process below:

1. Identify the relevant stakeholder(s) (e.g., end user, person impacted by the system, etc.).
2. Identify any demographic groups, including marginalized groups, that may require fairness considerations.
3. Prioritize these groups for fairness consideration and explain how the fairness consideration applies.

Demographic groups can refer to any population group that shares one or more particular demographic characteristics. Depending on the AI system and context of deployment, the list of identified demographic groups will change.

Marginalized groups are demographic groups who may have an atypical or even unfair experience with the system if their needs and context are not considered. May include minorities, stigmatized groups, or other particularly vulnerable groups. Additionally, marginalized groups can also include children, the elderly, indigenous peoples, and religious minorities. Groups to include for consideration will depend in part on the geographic areas and intended uses of your system.

Goal F1: Quality of Service

Applies to: AI systems when system users or people impacted by the system with different demographic characteristics might experience differences in quality of service that Microsoft can remedy by building the system differently.

E.g., a system that uses natural language processing may perform differently for users who speak supported languages as a second language or who speak less common varieties of a language.

Goal F2: Allocation of resources and opportunities

Applies to: AI systems that generate outputs that directly affect the allocation of resources or opportunities relating to finance, education, employment, healthcare, housing, insurance, or social welfare.

E.g., an automated hiring system that exhibits bias against hiring certain demographic groups (e.g., women).

Goal F3: Minimization of stereotyping, demeaning, and erasing outputs

Applies to: AI systems when system outputs include descriptions, depictions, or other representations of people, cultures, or society.

E.g., a system that uses natural language processing to generate text for images may under-represent particular groups. Such an example is if the system generated the caption of "CEO" only for white males.

2.5 – 2.8

Technology readiness assessment, task complexity, role of humans, and deployment environment complexity

Guidance

Your responses to these questions will help potential reviewers understand important details about how the system has been evaluated to date, what type of tasks the system is designed to execute, how humans interact with the system, and how you plan to deploy the system.

Examples

See Impact Assessment Template for inline examples.

Section 3

Adverse impact



Restricted Uses



Unsupported uses



Known limitations



Potential impact of failure on stakeholders



Potential impact of misuse on stakeholders



Sensitive uses

Thinking through adverse impact

Guidance

Even the best systems have limitations, fail sometimes, and can be misused. Consider known limitations of the system, the potential impact of failure on stakeholders, and the potential impact of misuse.

Prompts

- Try thinking from a hacker's perspective.
 - Consider what a non-expert might assume about the system.
 - Imagine a very negative news story about the system. What does it say?
-

3.1

Restricted uses

Guidance

Certain uses of AI technology at Microsoft are restricted. The list of Restricted Uses evolves periodically, so it's necessary to consult the list of Restricted Uses and their definitions each time you complete or review an Impact Assessment and follow the guidance for any Restricted Use.

3.2

Unsupported uses

Guidance

Some of the potential uses of the system fall outside of the scope of an intended use.

Unsupported uses can include:

- Reasonably foreseeable uses for which the system was not designed or evaluated
- Uses that we recommend customers avoid

Examples

A system that uses computer vision to recognize handwritten text was not designed or tested to verify the authenticity of signatures on forms.

A recommendation system that can be tailored to a customer's specific needs was not designed to make recommendations in sensitive domains like healthcare or finance.

3.3

Known limitations

Guidance

Every system will have limitations. Describing those limitations will ensure that the system is used for its intended purposes.

Examples

A system that translates speech to text will perform poorly in a noisy environment where several people in the proximity of the user are speaking.

A system that uses natural language processing may perform poorly for non-native speakers of a supported language.

Prompts

- Are there conditions in the deployment environment that would affect the system's performance?
- Are there types or ranges of input that are not suited for the system?

3.4

Potential impact of failure on stakeholders

Guidance

We should anticipate known failures and understand how they might impact stakeholders.

Define and document the predictable failures, including false positive and false negative results, and how they would impact stakeholders. Consider how system failures would manifest in each of the identified intended uses and for the system as a whole. Consider how reliability, accuracy, scope of impact, and failure rates of components and the overall system may impact appropriate use. Identify and document whether the likelihood of failure or consequences of failure differ for any marginalized groups. When serious impacts of failure are identified, note them in the summary of impact as a potential harm.

Prompts

- What are the predictable failures of this system?
- How would a false positive impact stakeholders?
- How would a false negative impact stakeholders?
- Does the likelihood or consequence of failure differ for any marginalized groups?

Examples

Example for: A system that scans web content and blocks suspicious websites:

What are the predictable failures of this system?	<i>False positives, false negatives, and major system outages.</i>
How would a false positive impact stakeholders?	<i>A website is labeled suspicious and blocked incorrectly. The end user can't access legitimate content, and the website owner will lose traffic.</i>
How would a false negative impact stakeholders?	<i>A website that should have been blocked is not and end users are exposed to malicious content.</i>
Does the likelihood or consequence of failure differ for any marginalized groups?	<i>In the case of false negatives, the effects could be more harmful to minors.</i>

3.5

Potential impact of misuse on stakeholders

Guidance

Every system could be intentionally or unintentionally misused. It is important to understand what misuse could be for the system and how that misuse may impact stakeholders.

Prompts

- How could someone misuse the system?
- How would misuse impact stakeholders?
- Do the consequences of misuse differ for any marginalized groups?

Examples

Example for: A system that scans web content and blocks suspicious websites:

	Misuse #1	Misuse #2
How could someone misuse the system?	Malicious actors could figure out how to evade the detection system.	Website owners could figure out how to trigger warnings on competitors' websites.
How would misuse impact stakeholders?	End users would be exposed to malicious content more frequently.	Competitors would be improperly flagged by the system and receive fewer visits and less traffic.
Do the consequences of misuse differ for any marginalized groups?	The effects could be more harmful to minors.	

3.6

Sensitive uses

Guidance

If you are designing, developing, or deploying AI systems that could be applied in Sensitive Use scenarios, report those to the Office of Responsible AI as early as possible.

Compare the uses and misuses of the system with the triggers for Sensitive Uses. If any of the uses meet any of the triggers for a Sensitive Use, report them to the Office of Responsible AI.

Section 4

Data requirements



Data requirements



Pre-defined data sets

4.1

Data requirements

Guidance

Define and document data requirements for training, validating, and testing the system with respect to the system's intended uses, stakeholders, and the geographic areas where the system will be deployed.

4.2

Existing data sets

Guidance

If you plan to use existing data sets to train, validate, or test the system, assess the quantity and suitability of available data sets that will be needed by the system in relation to the data requirements defined in 4.1.

Section 5

Summary of impact



Potential harm and preliminary mitigations



Goal applicability



Signing off on the Impact Assessment

5.1

Potential harms and preliminary mitigations

Guidance

This table will help reviewers understand how the system's potential harms will be addressed (either through complying with the Responsible AI Standard or through other mitigations).

Mitigating potential harms: [Activity](#)

Steps

- 1)** Collect the harms you previously identified throughout the Impact Assessment in this table (check your stakeholder table, fairness considerations, adverse impact section, and any other place you described potential harms).
- 2)** For each potential harm, use the mitigations prompts in this section to understand if one or more of the Goals from the Responsible AI Standard can serve as a mitigation.
 - If a Goal can serve as a mitigation to a harm, copy and paste the mitigation text in the second column of the table. Then, document the team's plan for how to implement this mitigation in the design of the system in the third column of the table.
- 3)** After completing step 2 for each harm, discuss the harms that remain unmitigated with your team and develop mitigations for those in the table.

Case Study: Hospital Employee and Resource Optimization System (HEROS)

The team started by listing all the relevant potential harms identified in the Impact Assessment in the first column of the table.

The first harm in the list was from the fairness evaluation: “Patients from low-income backgrounds might have a history of short hospital stays for financial reasons. This may result in predictions for shorter stays and impact the patient’s access to a hospital bed.” The team read the prompts (on the following pages) evaluating what kind of harm this was and what would be appropriate mitigations. The team identified two prompts for harms mitigations that aligned well with their harm: F1 which describes differences in system performance for specific demographic groups, and F2 which refers to allocation and differences across stakeholder groups. The team read both sets of mitigations and discussed how to apply them to the system. Both prompts refer to “evaluating the datasets and the system then modifying the system to avoid differences in performance and allocation.” Since the system’s intended use is to allocate healthcare resources, F2 is the best fit for this potential harm. They recorded the appropriate mitigation and strategy in the Impact Assessment.

	Describe the potential harm	Corresponding Goal from the Responsible AI Standard (if applicable)	Describe your initial ideas for mitigations or explain how you might implement the corresponding Goal in the design of the system
Prompts and discussion >>>> Impact assessment table	Potential Harm: Patients from ‘low-income backgrounds’ might have a history of short hospital stays for financial reasons.	<div> <div>F1</div> <div>OR</div> <div>F2</div> </div>	Mitigation strategy for this system: Evaluate fairness for this demographic group by comparing the system’s performance for this group to the system’s performance for other groups.
	Patients from ‘low-income backgrounds’ might have a history of short hospital stays for financial reasons. This may result in predictions for shorter stays and impact the user’s access to a hospital bed.	F2 - This harm is mitigated by evaluating the data sets and the system then modifying the system to minimize differences in the allocation of resources and opportunities between identified demographic groups.	Evaluate the dataset and system performance for the demographic group of patients from low-income communities. If we detect disparities in hospital stay length predictions or the allocation of resources, we will reassess the system design to minimize the disparity.

Harms mitigation: Accountability prompts

The requirements from each respective Goal are the expected mitigation for the relevant harm.

Prompt

Is this harm the result of a consequential impact on legal position or life opportunities; risk of physical or psychological injury; a threat to human rights; or a Restricted Use?

Could this harm be mitigated by clarifying the problem to be solved by the system and communicating evidence that the system is fit for purpose to stakeholders?

Is this harm the result of data that has not been sufficiently managed or evaluated in relation to the system's intended use(s)?

Could this harm be mitigated if the system had adequate human oversight and control?

If you answer yes to any of the prompts, enter the corresponding mitigation(s) below the prompt in the second column of table 5.1 in the Impact Assessment Template.

Mitigation

Goal A2: Oversight of significant adverse impacts

This Goal applies to all AI systems.

Harms that result from Sensitive Uses must be mitigated by guidance received from the Office of Responsible AI's Sensitive Uses team. Please report your system as Sensitive Use. For Restricted Uses, see guidance [here](#).

Goal A3: Fit for purpose

This Goal applies to all AI systems.

This harm is mitigated by assessing whether the system is fit for purpose for this intended use by providing evidence, recognizing that there may be many valid ways in which to solve the problem.

Goal A4: Data governance and management

This Goal applies to all AI systems.

This harm is mitigated by ensuring that data used to train the system is correctly processed and appropriate based on the intended use, stakeholders, and geographic areas.

Goal A5: Human oversight and control

This Goal applies to all AI systems.

This harm can be mitigated by modifying system elements (like system UX, features, educational materials, etc.) so that the relevant stakeholders can effectively understand and fulfill their oversight responsibilities.

Example

For example, a system that recommends diagnoses based on a patient's medical history triggers the sensitive use category: risk of physical or psychological injury. To mitigate this potential harm, we must seek guidance from the Sensitive Uses team to understand how best to design the system to minimize such errors.

For example, a facial recognition system that estimates an individual's age is not suitable for verifying the legal age of an individual for the sale of controlled substances.

For example, a system that predicts the length of hospital stays was trained on data from UK hospitals. To use this system in other geographic areas the system should be retrained with data from hospitals in the appropriate geographic locations.

For example, a system that diagnoses a medical condition based on video recordings of a patient in a specific orientation. If the patient is not in the required orientation, the system may misdiagnose the patient. To mitigate this potential harm the system can ensure the patient is in the correct orientation during the image preprocessing stage and alert a technician when the patient needs to be reoriented.

Harms mitigation: Transparency prompts

The requirements from each respective Goal are the expected mitigation for the relevant harm.

Prompt	Is this harm the result of inadequate intelligibility of system outputs?	Could this harm be mitigated by a better understanding of what the system can or cannot do?	Is this harm the result of users not understanding that they are interacting with an AI system or AI-generated content?
Mitigation	<p>Goal T1: System intelligibility for decision making</p> <p><i>This Goal applies to AI systems when the intended use of the generated outputs is to inform decision making by or about people.</i></p> <p>This harm is mitigated by designing system elements (like system UX, features, educational materials, etc.) so that the affected stakeholders can interpret system behavior effectively.</p>	<p>Goal T2: Communication to stakeholders</p> <p><i>This Goal applies to all AI systems.</i></p> <p>This harm is mitigated by providing stakeholders with relevant information about the system to inform decisions about when to employ the system or platform.</p>	<p>Goal T3: Disclosure of AI interaction</p> <p><i>This Goal applies to AI systems that impersonate interactions with humans, unless it is obvious from the circumstances or context of use that an AI system is in use; and AI systems that generate or manipulate image, audio, or video content that could falsely appear to be authentic.</i></p> <p>This harm is mitigated by modifying system elements (like system UX, features, educational materials, etc.) so that the relevant stakeholders will understand the type of AI system they are interacting with or that the content they are exposed to is AI-generated.</p>
Example	<p><i>For example, a system that predicts whether the user would get approved for admission to college can offer factors that could have led to an approval such as 'if you had volunteered with a nonprofit' or 'if your GPA had been higher.'</i></p>	<p><i>For example, a voice transcription system to provide healthcare to people that speak a wide variety of languages doesn't work well in the emergency room admissions area because it is often a noisy environment. To mitigate potential harms, we can communicate this system limitation to stakeholders.</i></p>	<p><i>For example, an AI-powered chatbot begins each conversation by clearly stating it is an autonomous chat agent.</i></p>

Harms mitigation: Fairness prompts

The requirements from each respective Goal are the expected mitigation for the relevant harm.

Prompt

Is this harm the result of the system providing a worse quality of service for some demographic groups?

Is the harm the result of the system affecting the allocation resources or opportunities relating to finance, education, employment, healthcare, housing, insurance, or social welfare, differently for different demographic groups?

Is this harm the result of outputs of the system that stereotype, demean, or erase some demographic groups?

If you answer yes to any of the prompts, enter the corresponding mitigation(s) below the prompt in the second column of table 5.1 in the Impact Assessment Template.

Mitigation

Goal F1: Quality of Service

This Goal applies to AI systems when system users or people impacted by the system with different demographic characteristics might experience differences in quality of service that Microsoft can remedy by building the system differently.

This harm is mitigated by evaluating the data sets and the system then modifying the system to improve system performance for affected demographic groups while minimizing performance differences between identified demographic groups.

Goal F2: Allocation of resources and opportunities

This Goal applies to AI systems that generate outputs that directly affect the allocation of resources or opportunities relating to finance, education, employment, healthcare, housing, insurance, or social welfare.

This harm is mitigated by evaluating the data sets and the system then modifying the system to minimize differences in the allocation of resources and opportunities between identified demographic groups.

Goal F3: Minimization of stereotyping, demeaning, and erasing outputs

This Goal applies to AI systems when system outputs include descriptions, depictions, or other representations of people, cultures, or society.

This harm is mitigated by a rigorous understanding of how different demographic groups are represented within the AI system and modifying the system to minimize harmful outputs.

Example

For example, people who speak language varieties that are underrepresented in the training data may experience worse quality of service for a voice transcription system. The system can be evaluated and modified to improve quality of service for these demographic groups.

For example, a hiring system that scans resumes and recommends candidates for hiring trained on historical data tends to be biased toward male candidates. The system can be evaluated and modified to reduce unfair allocation of opportunities.

For example, an image search system that predominantly returns images of men in response to the query "chief executive officer" may underrepresent non-male chief executive officers. To mitigate this, the system can be modified to provide more representative outputs.

Harms mitigation: Reliability & Safety prompts

The requirements from each respective Goal are the expected mitigation for the relevant harm.

Prompt	<p>Could this harm be mitigated by defining and documenting reliable and safe performance of the system and providing documentation to customers?</p>	<p>Is this harm the result of a predictable failure, or inadequately managing unknown failures once the system is in use?</p>	<p>Could this harm be mitigated by monitoring and evaluating the system in an ongoing manner?</p>
<p>If you answer yes to any of the prompts, enter the corresponding mitigation(s) below the prompt in the second column of table 5.1 in the Impact Assessment Template.</p>			
Mitigation	<p>Goal RS1: Reliability and safety guidance <i>This Goal applies to all AI systems.</i></p> <p>This harm is mitigated by defining safe and reliable behavior for the system, ensuring that datasets include representation of key intended uses, defining operational factors and ranges that are important for safe & reliable behavior for the system, and communicating information about reliability and safety to stakeholders.</p>	<p>Goal RS2: Failures and remediations <i>This Goal applies to all AI systems</i></p> <p>This harm is mitigated by establishing failure management approaches for each predictable failure.</p>	<p>Goal RS3: Ongoing monitoring, feedback, and evaluation <i>This Goal applies to all AI systems.</i></p> <p>This harm is mitigated by establishing system monitoring methods that allow the team to identify and review new uses, identify and troubleshoot issues, manage and maintain the system, and improve the system over time.</p>
Example	<p>For example, a computer vision system not optimized for use in varying lighting situations may compromise performance. To mitigate this potential harm, we can define the lighting conditions in which the system performs well.</p>	<p>For example, a synthetic voice service could be used in ways that are harmful, for example, for Deepfake videos. To mitigate these potential harms, the system can produce a watermark to signal that the content is AI generated.</p>	<p>For example, monitoring revealed that a system that uses infrared cameras to verify an individual logging into a device fails five times more frequently on weekend days. Specifically, this is because the weekend logins take place in different environmental conditions with more infrared light. To improve performance, the system can include secondary verification templates created in a variety of lighting conditions.</p>

5.2

Goal Applicability

Guidance

The tables in section 5.2 will help you assess which of the Goals in the Responsible AI Standard apply to the system. Some of the Goals in the Standard apply to *all* AI systems, while other Goals apply to only specific types of AI systems. When a Goal applies to the system being evaluated in the Impact Assessment, you *must* complete the requirements associated with that Goal. When a Goal does *not* apply to the system, we ask that you explain why in a textbox below the tables.

5.3

Signing off on the Impact Assessment

Guidance

Signing off is the final step in completing your Impact Assessment. In this section of the Impact Assessment, ensure that your Impact Assessment has been reviewed and approved by the reviewers named in 1.1.

Scan this code to access responsible AI resources from Microsoft:



© 2022 Microsoft Corporation. All rights reserved. This document is provided "as-is." It has been edited for external release to remove internal links, references, and examples. Information and views expressed in this document may change without notice. You bear the risk of using it. Some examples are for illustration only and are fictitious. No real association is intended or inferred. This document does not provide you with any legal rights to any intellectual property in any Microsoft product. You may copy and use this document for your internal, reference purposes.