# Practitioner Guidance for Human Rights-Based AI Governance

May 2025

BSR®

How can you take a human rights-based approach to responsible AI?

BSR

# Responsible AI Practitioner Guides

**1**    **Fundamentals of a Human Rights-Based Approach to Generative AI**

**2**    **A Human Rights-Based Approach to Governance and Management**

**3**    **A Human Rights-Based Approach to Impact Assessment**

**4**    **A Human Rights-Based Approach to Risk Mitigation**

**5**    **Conducting Stakeholder Engagement**

**6**    **A Human Rights-Based Approach to Policies and Enforcement**

**7**    **Aligning Transparency and Disclosure Practices with Human Rights Responsibilities**

**8**    **Remedy for Generative AI-Related Harms**

BSR

# BSR Speakers

**Lindsey Andersen**
**Associate Director**
San Francisco

**Hannah Darnton**
**Director**
San Francisco

**J.Y. Hoh**
**Associate Director**
Singapore

**Samone Nigam**
**Manager**
San Francisco

BSR

# Agenda

1. Review of key points from each practitioner guide

2. Key Takeaways

3. Q&A

# Guide 1:

# Fundamentals of a Human Rights-Based Approach to AI

BSR

A human rights-based approach means embedding respect for human rights into the development and deployment of AI

BSR

# Why a human rights-based approach?

**1**

**International standards for governments and companies**

- International human rights instruments
- The UN Guiding Principles on Business and Human Rights
- The OECD Guidelines for Multinational Enterprises on Responsible Business conduct

**2**

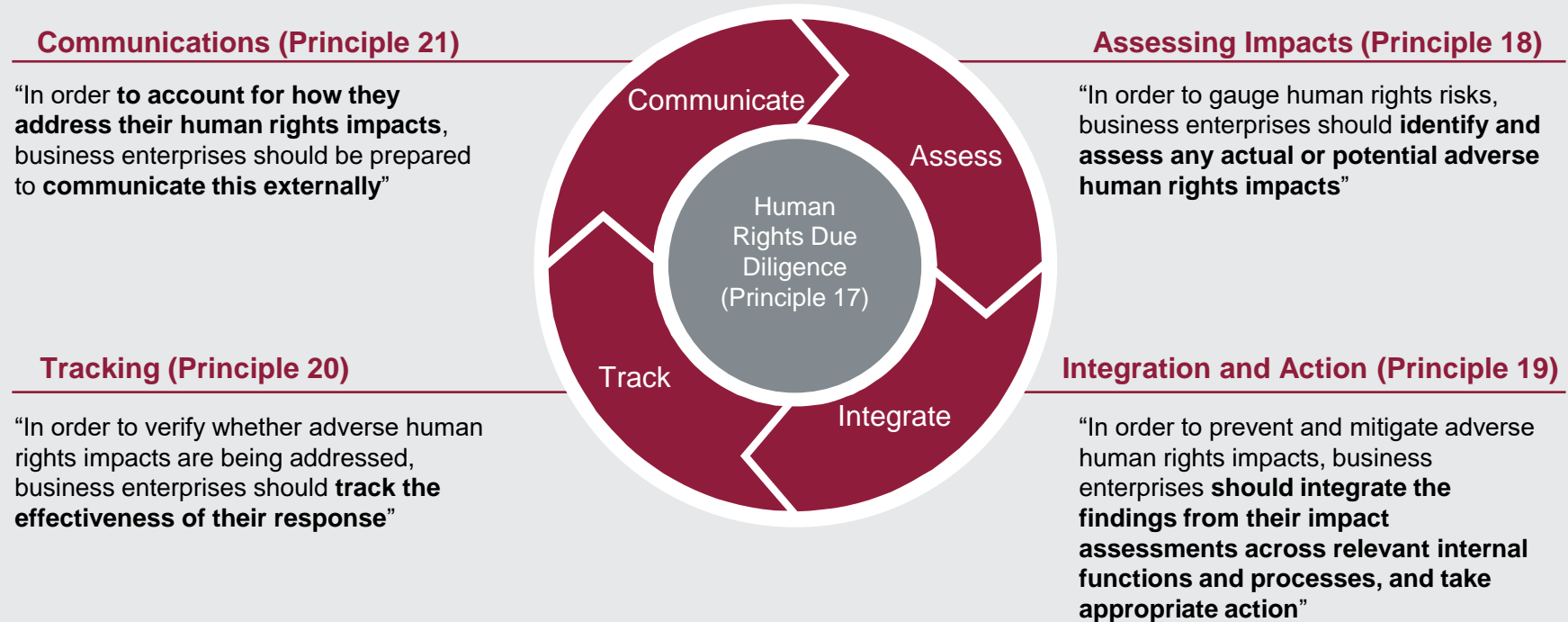**Established framework and methodology**

- The UNGPs provide an approach to identifying, assessing, prioritizing, and addressing risks to people

**3**

**Human rights integration into AI regulation**

- EU Digital Service Act
- EU AI Act
- Corporate Sustainability Due Diligence Directive
- Corporate Sustainability Reporting Directive

# The UNGPs require companies to conduct human rights due diligence

## Communications (Principle 21)

"In order **to account for how they address their human rights impacts**, business enterprises should be prepared to **communicate this externally**"

## Tracking (Principle 20)

"In order to verify whether adverse human rights impacts are being addressed, business enterprises should **track the effectiveness of their response**"

## Assessing Impacts (Principle 18)

"In order to gauge human rights risks, business enterprises should **identify and assess any actual or potential adverse human rights impacts**"

## Integration and Action (Principle 19)

"In order to prevent and mitigate adverse human rights impacts, business enterprises **should integrate the findings from their impact assessments across relevant internal functions and processes, and take appropriate action**"

Communicate

Assess

Track

Integrate

Human Rights Due Diligence (Principle 17)

BSR

# How is a human rights-based approach different from others?

## Ethics

- Framework for decision-making in situations where right and wrong, good and bad, are not clearly defined
- Different "schools of thought" and standards that support different approaches and choices
- Different traditions, cultures, countries, and religions may choose different outcomes
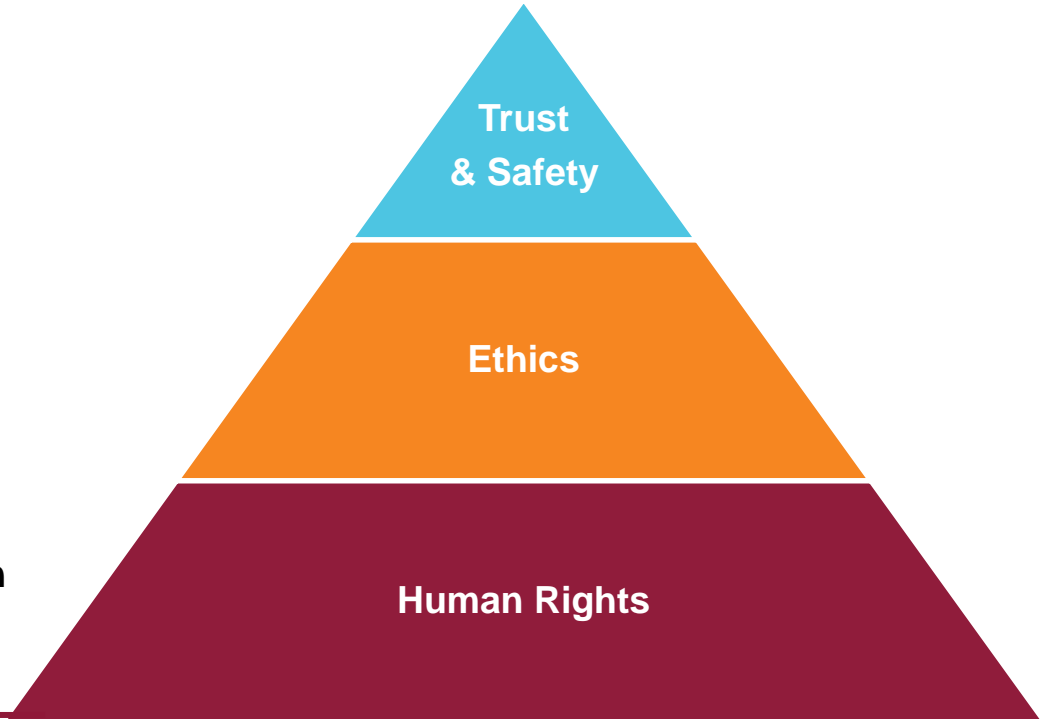
## Human Rights

- Internationally defined framework for government and company responsibilities
- List of rights that should always be protected / respected regardless of context or culture
- Focus on the experiences of the most vulnerable
- Establish a floor rather than a ceiling—e.g., respect for human rights is a minimum requirement

## Trust & Safety

- Function within companies to operationalize ethics / human rights / safety efforts
- Focused on practically ensuring genAI tools are "safe"
- Brings approaches and lessons learned from online platform content governance
- Anchored on pre-established risk taxonomies

# How does it all fit together?

**Human rights should be the <u>foundation</u> for responsible AI upon which other approaches can be integrated.**

Trust & Safety

Ethics

Human Rights

# Guide 2:

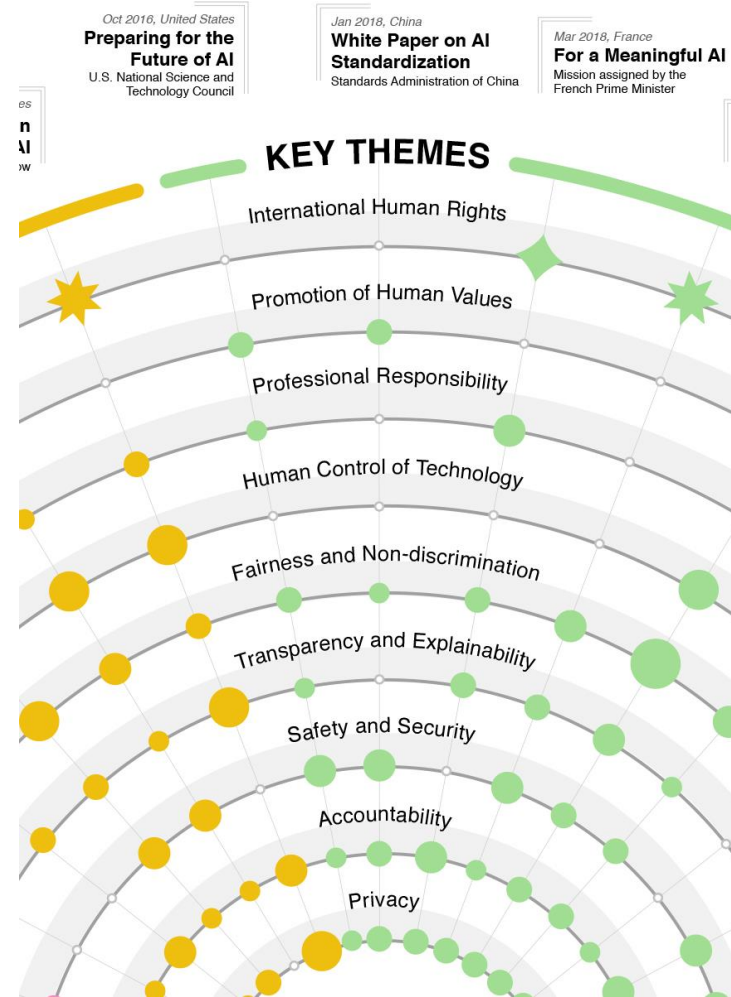# A Human Rights-Based Approach Governance and Management

BSR

# What are AI principles?

- A set of values that guide an organization's AI development / deployment

# Why are they important?

- AI principles provide a **foundation for embedding responsible AI** across an organization (e.g. provide remit to responsible AI teams)

- Including human rights in AI principles provides a **foundation for a human rights-based approach** to responsible AI

- Including human rights provides **clarity and consistency** to nebulous principles

Source: Principled Artificial Intelligence: A Map of Ethical and Rights-Based Approaches to Principles for AI, https://cyber.harvard.edu/publication/2020/principled-ai



Oct 2016, United States
**Preparing for the Future of AI**
U.S. National Science and Technology Council

Jan 2018, China
**White Paper on AI Standardization**
Standards Administration of China

Mar 2018, France
**For a Meaningful AI**
Mission assigned by the French Prime Minister

**KEY THEMES**

International Human Rights
Promotion of Human Values
Professional Responsibility
Human Control of Technology
Fairness and Non-discrimination
Transparency and Explainability
Safety and Security
Accountability
Privacy

# Example AI principles that include human rights

Google ———————————————————————————————— AI Principles

## 2 Responsible development and deployment

Because we understand that AI, as a still-emerging transformative technology, poses evolving complexities and risks, we pursue AI responsibly throughout the AI development and deployment lifecycle, from design to testing to deployment to iteration, learning as AI advances and uses evolve. This means:

- Implementing appropriate human oversight, due diligence, and feedback mechanisms to align with user goals, social responsibility, and widely accepted principles of international law and human rights.
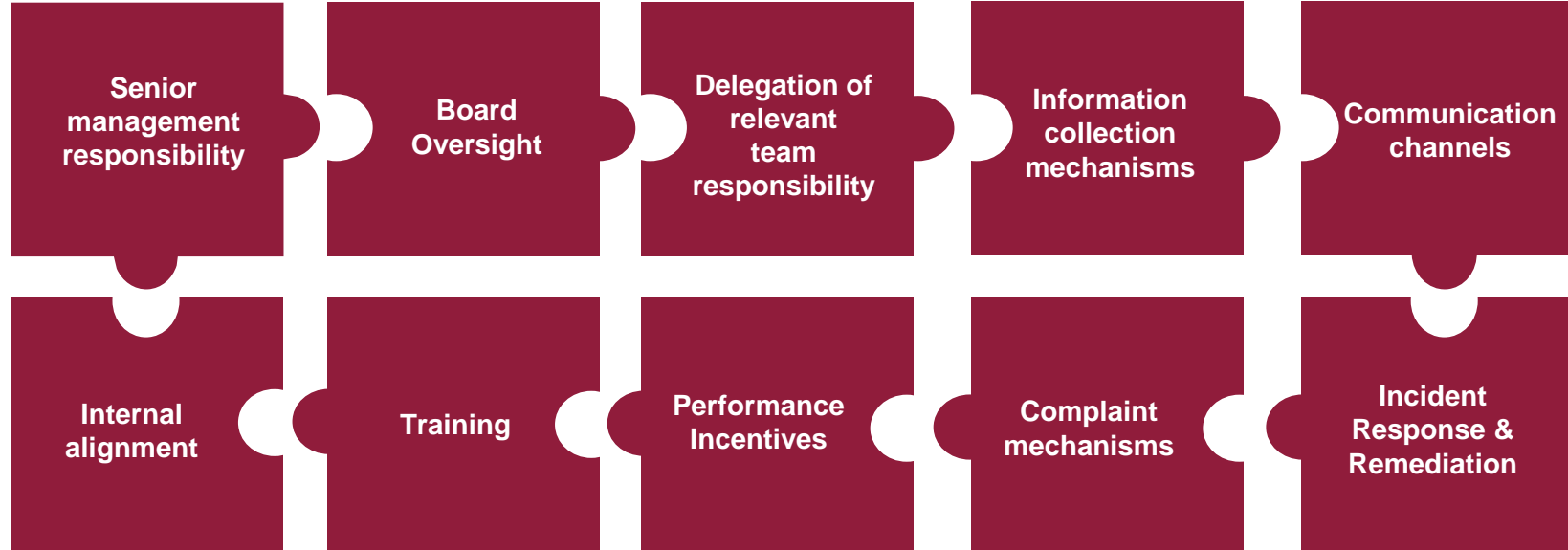
BSR

# Example AI principles that include human rights

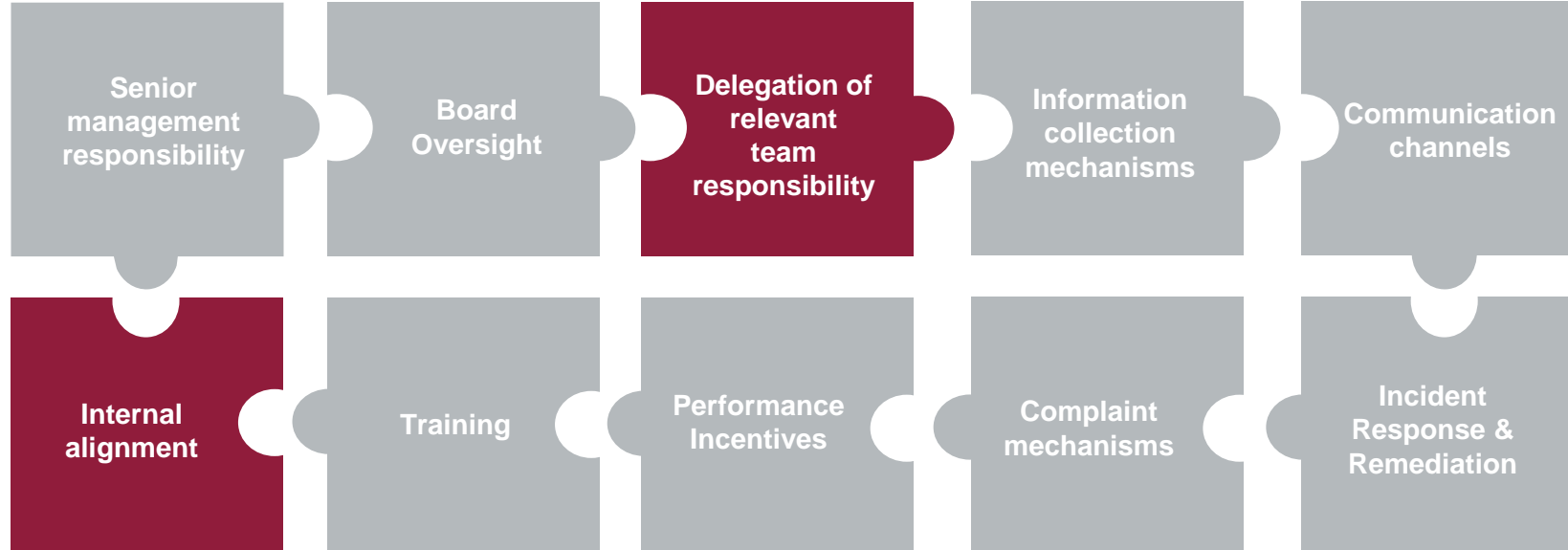## Meet Salesforce's Trusted AI Principles

### Responsible

We strive to safeguard human rights, to protect the data we are trusted with, observe scientific standards and enforce policies against abuse. We expect our customers to use our AI responsibly, and in compliance with their agreements with us, including our Acceptable Use Policy.

BSR

# Applying Human Rights Governance Principles to Responsible AI



Source: OECD Due Diligence Guidance for Responsible Business Conduct

BSR

# Applying Human Rights Governance Principles to Responsible AI

Senior management responsibility

Board Oversight

Delegation of relevant team responsibility

Information collection mechanisms

Communication channels

Internal alignment

Training

Performance Incentives

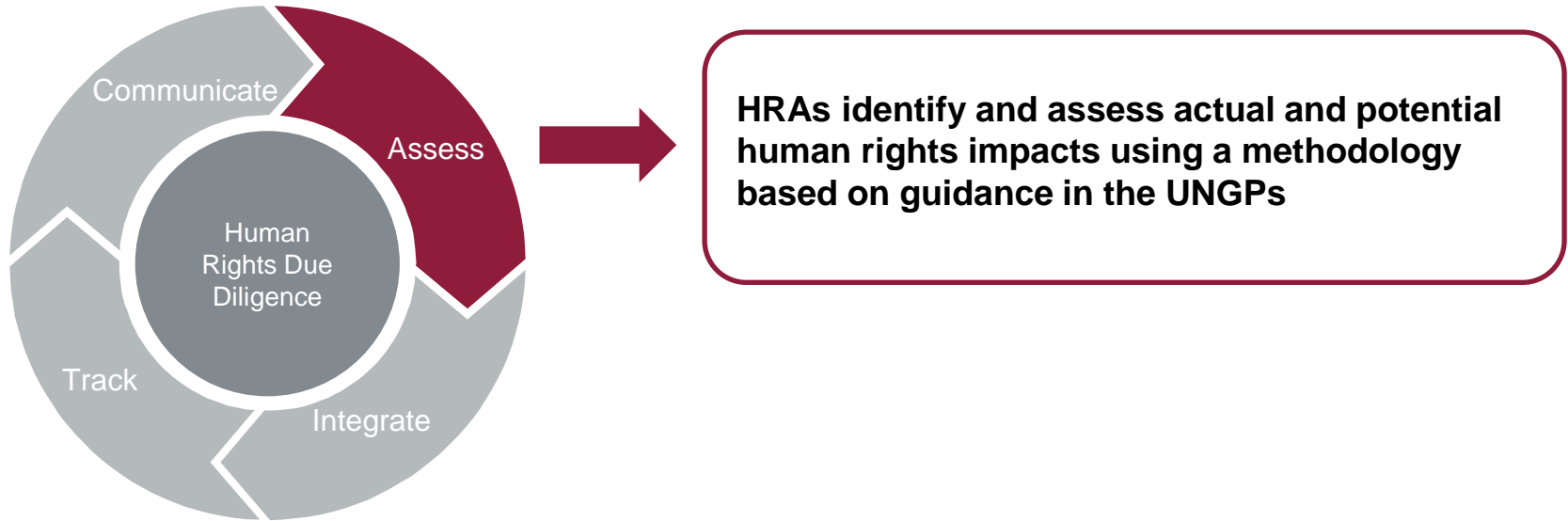Complaint mechanisms

Incident Response & Remediation

Source: OECD Due Diligence Guidance for Responsible Business Conduct

BSR

# Guide 3:

# A Human Rights-Based Approach to Impact Assessment

# What are human rights assessments?

Communicate

Assess

Human
Rights Due
Diligence

Track

Integrate

**HRAs identify and assess actual and potential human rights impacts using a methodology based on guidance in the UNGPs**

# Core Elements of a Human Rights Assessment

- Identifying impacts human rights impacts using all internationally recognized human rights as a reference point

BSR

# List of Internationally Recognized Human Rights

- Right to equality and non-discrimination
- Right to life, liberty, and personal security
- Freedom from slavery
- Freedom from torture and degrading treatment
- Due process and fair trial rights
- Freedom from arbitrary arrest and exile
- Right to privacy
- Freedom of movement
- Right to asylum
- Right to a nationality and the freedom to change nationality
- Right to marriage and family
- Right to own property
- Freedom of thought
- Freedom of religion and belief
- Right to remedy

- Freedom of opinion, expression, and access to information
- Right of peaceful assembly and association
- Right to political participation
- Right to social security
- Labor Rights (e.g. safe working conditions, adequate remuneration, right to join unions)
- Right to rest and leisure
- Right to adequate living standards
- Right to health
- Right to education
- Right to participate in the cultural life of the community
- Right to benefit from scientific advancement
- Right to internet access
- Right to a healthy environment
- Disability rights (e.g. right to accessibility)
- Child Rights

BSR

# Core Elements of a Human Rights Assessment

- Identifying impacts human rights impacts using all internationally recognized human rights as a reference point

- Assessing and prioritizing impacts based on severity to people

- Emphasis on vulnerable and marginalized groups; stakeholder engagement

- Considering interconnectivity between rights

- Accounting for context

BSR

# BSR's Human Rights Assessment Process

**Step 1:**
Define the scope of the assessment

**Step 2:**
Secure relevant buy-in

**Step 3:**
Gather and analyze relevant information and data

**Step 4:**
Engage internal and external stakeholders

**Step 5:**
Identify actual and potential human rights impacts

**Step 6:**
Assess and prioritize impacts

**Step 7:**
Identify appropriate action to address impacts

BSR

# Why HRAs for AI?

| Benefits of HRAs | Limitations of HRAs |
|---|---|
| • Focus on impacts to people | |
| • Comprehensiveness of risk / impact identification | |
| • An approach to prioritizing impacts | |
| • An established, internationally accepted methodology | |
| • Adaptability to a variety of contexts | |
| • Assistance with regulatory compliance | |

BSR

# Why HRAs for AI?

| Benefits of HRAs | Limitations of HRAs |
|---|---|
| • Focus on impacts to people | • May not cover all relevant impacts |
| • Comprehensiveness of risk / impact identification | • Are more qualitative than quantitative |
| • An approach to prioritizing impacts | • Are not technical assessments |
| • An established, internationally accepted methodology | |
| • Adaptability to a variety of contexts | |
| • Assistance with regulatory compliance | |

BSR

# Integrating human rights into other AI impact assessments

| Assessment Type | Description | How to Integrate Human Rights |
|---|---|---|
| **Algorithmic Impact Assessments / Audits** | Systematic examination of the algorithms and data used in an AI system to assess their fairness, accountability, transparency, and ethical implications. | Utilize the list of internationally recognized human rights (see the appendix) as a foundation for brainstorming to help identify impacts or create a risk/harm taxonomy. Consider severity when assessing impacts. |
| **Model / Application Evaluations** | Empirical assessments of an AI system's performance or impact on people and society. | Utilize human rights as a foundation for identifying impacts/harms to evaluate. |
| **Fairness Testing** | Assessment of whether an AI system exhibits biases or discrimination against certain groups of individuals based on protected characteristics such as race, gender, ethnicity, or age. Often includes model/application evaluation. | Utilize the vulnerable groups framework to help identify groups for the basis of testing.<br><br>Consider how additional human rights may be impacted as a result of identified fairness issues. |
| **Data Quality Reviews** | Examination of the data used to train AI models to look for issues such as incorrect labels, representativeness, accuracy, and bias, that may lead to inaccurate or problematic outputs. | Consider how different data quality issues could lead to human rights impacts, and consider the severity of those impacts to help prioritize corrective actions / mitigation of related impact. |
| **Red Teaming** | A range of assessment methods for AI systems that involves using adversarial techniques and approaches to test the security, robustness, and resilience of AI systems. | Identify pathways to human rights impacts as part of the red-teaming process.<br><br>Include red teamers with a background suited to identifying risks to people, as well as people representative of, or familiar with, risks and needs of vulnerable groups. |

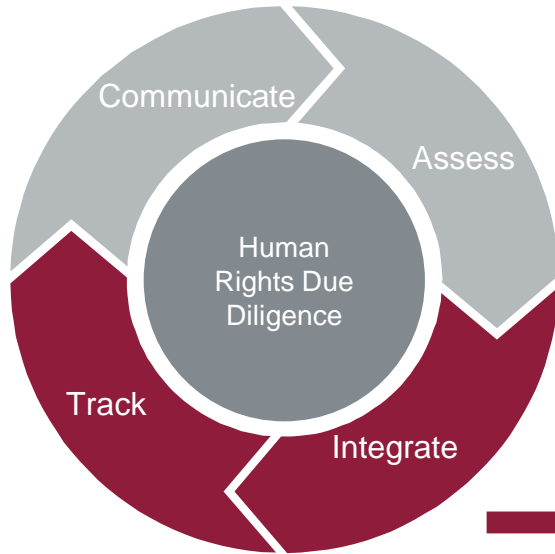# Integrating human rights into other AI impact assessments

| Assessment Type | Description | How to Integrate Human Rights |
|---|---|---|
| **Algorithmic Impact Assessments / Audits** | Systematic examination of the algorithms and data used in an AI system to assess their fairness, accountability, transparency, and ethical implications. | Utilize the list of internationally recognized human rights (see the appendix) as a foundation for brainstorming to help identify impacts or create a risk/harm taxonomy. Consider severity when assessing impacts. |
| **Model / Application Evaluations** | Empirical assessments of an AI system's performance or impact on people and society. | Utilize human rights as a foundation for identifying impacts/harms to evaluate. |
| **Fairness Testing** | Assessment of whether an AI system exhibits biases or discrimination against certain groups of individuals based on protected characteristics such as race, gender, ethnicity, or age. Often includes model/application evaluation. | Utilize the vulnerable groups framework to help identify groups for the basis of testing.<br><br>Consider how additional human rights may be impacted as a result of identified fairness issues. |
| **Data Quality Reviews** | Examination of the data used to train AI models to look for issues such as incorrect labels, representativeness, accuracy, and bias, that may lead to inaccurate or problematic outputs. | Consider how different data quality issues could lead to human rights impacts, and consider the severity of those impacts to help prioritize corrective actions / mitigation of related impact. |
| **Red Teaming** | A range of assessment methods for AI systems that involves using adversarial techniques and approaches to test the security, robustness, and resilience of AI systems. | Identify pathways to human rights impacts as part of the red-teaming process.<br><br>Include red teamers with a background suited to identifying risks to people, as well as people representative of, or familiar with, risks and needs of vulnerable groups. |

# Guide 4:

# A Human Rights-Based Approach to Risk Mitigation

BSR

# Where does risk mitigation fit?



**Risk mitigation entails:**

- **Integrating and acting upon the findings from impact assessments**

- **Tracking the effectiveness of mitigation measures**

BSR

# What human rights concepts can inform AI risk mitigation?

**1 — Attribution**

How closely connected an entity is to an impact

**+**

**2 — Leverage**

Ability to address the harm / influence the entity causing the harm

**Cause**

Entity's actions alone lead to the impact

- Cease or prevent the impact
- Provide remedy if it occurs

**Contribute**

Facilitates, enables, incentivizes a third party to cause the impact

- Cease or prevent the contribution
- Use leverage to mitigate remaining impacts
- Provide remedy if impact occurs

**Linked**

Impact is directly linked to its operations, products, or services by a business relationship

- Utilize leverage to address the impact

# Other human rights principles for AI risk mitigation

- Risk tolerance and "offsetting" risks and benefits is not acceptable in a human rights context

- Tensions and trade-offs (aka competing equities) can be worked through using "counterbalancing"

- Think broadly about risk mitigation—technical / product actions, policies, process, transparency / communications

- Identify risks that arise from mitigations

- Integrate risk mitigation across the organization

- Track the effectiveness of risk mitigation over time

BSR

# Other human rights principles for AI risk mitigation

- **Risk tolerance and "offsetting" risks and benefits is not acceptable in a human rights context**

- Tensions and trade-offs (aka competing equities) can be worked through using "counterbalancing"

- Think broadly about risk mitigation—technical / product actions, policies, process, transparency / communications, collaborative actions

- Identify risks that arise from mitigations

- Integrate risk mitigation across the organization

- Track the effectiveness of risk mitigation over time

# Other human rights principles for AI risk mitigation

- Risk tolerance and "offsetting" risks and benefits is not acceptable in a human rights context

- **Tensions and trade-offs (aka competing equities) can be worked through using "counterbalancing"**

- Think broadly about risk mitigation—technical / product actions, policies, process, transparency / communications, collaborative actions

- Identify risks that arise from mitigations

- Integrate risk mitigation across the organization

- Track the effectiveness of risk mitigation over time

BSR

# Other human rights principles for AI risk mitigation

- Risk tolerance and "offsetting" risks and benefits is not acceptable in a human rights context

- Tensions and trade-offs (aka competing equities) can be worked through using "counterbalancing"

- **Think broadly about risk mitigation—technical / product actions, policies, process, transparency / communications, collaborative actions**

- Identify risks that arise from mitigations

- Integrate risk mitigation across the organization

- Track the effectiveness of risk mitigation over time

BSR

# Example risk mitigations across the generative AI value chain

## Foundation Model Developers

| | |
|---|---|
| **Example Mitigation** | **Managed model rollout** |
| **Description** | The process of making informed and calculated decisions about model release on a gradient from closed to open source, which may include receiving feedback from relevant stakeholders to [...] tive model fine-tuning |
| **Risks It Addresses** | **Upstream:** gaps in training data that lead to model perf[...] issues, such as unequal performance across demographic[...] communities or outputs that are not representative of or[...] the user base<br><br>**Downstream:** applications of foundation models associa[...] adverse impacts on people |
| **How to Integrate a Human Rights-Based Approach** | • Consider how human rights could be adversely impa[...] the application of foundation models into technolog[...] range of use cases<br><br>• Consider how human rights could be adversely impa[...] limiting access and availability of the foundation mo[...]<br><br>• Identify any tensions in the above scenarios (i.e., adv[...] from closed vs. open source)<br><br>• Consider what leverage the foundation model deve[...] to mitigate and provide remedy for harms resulting f[...] stream applications of the model<br><br>• Establish decision-making processes for navigating h[...] trade-offs related to types of model release that sec[...] possible expression of the competing rights without [...] limiting them |

## Deployers

| | |
|---|---|
| **Example Mitigation** | **Human oversight** |
| **Description** | The process of human verification and/or approval of AI-generated predictions, decisions, or other outputs |
| **Risks It Addresses** | **Upstream:** Upstream: instances in which programmed safeguards fail to prevent inaccurate, inappropriate, biased, or otherwise harmful system outputs<br><br>**Downstream:** unintended adverse impacts on people associated with AI outputs |
| **How to Integrate a Human Rights-Based Approach** | • Consider how the integration of genAI tools into operations, processes, or workflows could cause, contribute to, or be directly linked to potential adverse human rights impacts<br><br>• Be aware of issues that could increase the likelihood of adverse human rights impacts (e.g., inaccurate or biased outputs, harmful behavior specific to the use case and application domain, etc.)<br><br>• Ensure users/operators are adequately trained on how genAI works, aware of potentially harmful outputs and behavior, and review outputs accordingly<br><br>• Establish processes for flagging and escalating harmful outputs<br><br>• Establish decision-making processes for navigating human rights trade-offs that secure the fullest possible expression of the competing rights without unnecessarily limiting them |

BSR

# Guide 5:

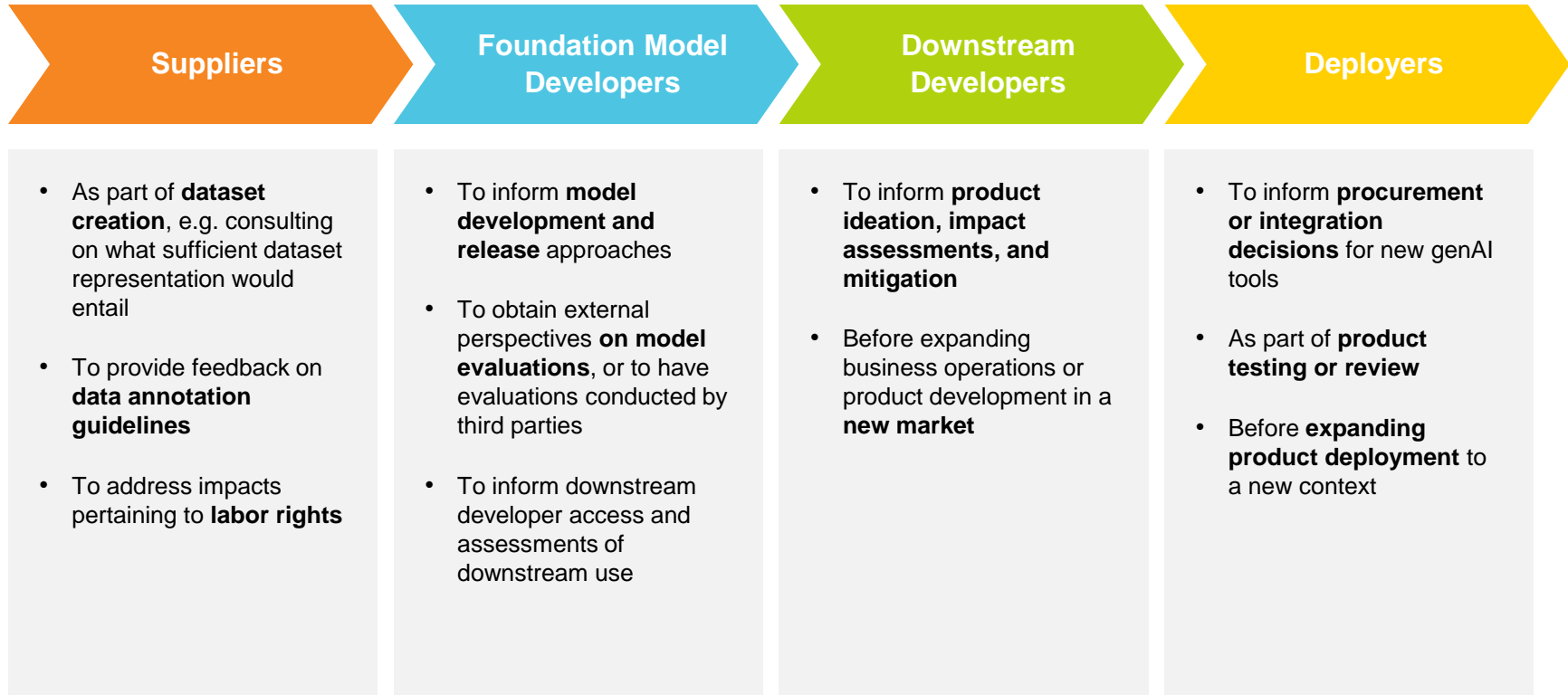# Conducting Stakeholder Engagement aka "Participatory Approaches"

# Stakeholder Engagement Across the GenAI Value Chain

| Suppliers | Foundation Model Developers | Downstream Developers | Deployers |
|---|---|---|---|
| • As part of **dataset creation**, e.g. consulting on what sufficient dataset representation would entail<br><br>• To provide feedback on **data annotation guidelines**<br><br>• To address impacts pertaining to **labor rights** | • To inform **model development and release** approaches<br><br>• To obtain external perspectives **on model evaluations**, or to have evaluations conducted by third parties<br><br>• To inform downstream developer access and assessments of downstream use | • To inform **product ideation, impact assessments, and mitigation**<br><br>• Before expanding business operations or product development in a **new market** | • To inform **procurement or integration decisions** for new genAI tools<br><br>• As part of **product testing or review**<br><br>• Before **expanding product deployment** to a new context |

# Best Practices for Stakeholder Engagement

| | |
|---|---|
| **Engage at multiple levels** | Ongoing organizational engagement, one-off engagement on an issue, and product-level engagements based on risk |
| **Identify staff to own engagement** | Dedicated, trained, and resourced staff should own engagement processes |
| **Prepare well** | Preparation ensures engagement goes smoothly and positive relationships are developed |
| **Build and maintain relationships** | Invest in relationships over time that are built on a foundation of trust, respect, and provide mutual benefit |
| **Consider compensation** | Be creative about compensation models to avoid the perception of an extractive relationship |

BSR

# Guide 6:

# A Human Rights-Based Approach to Policies and Enforcement

# Policies and enforcement as AI risk mitigations

| Product Policies | Organizational Policies |
|---|---|
| Terms of service (ToS) | Privacy / data protection policies |
| Acceptable use policies (AUPs) | Security / data access policies |
| Model Policies | Quality assurance policies |
| Content Policies | Compliance policies |
| AI Principles | Procurement policies |

BSR

# Elements of a human rights-based approach to policies and enforcement

- Policies should reference human rights when relevant

- Policies should address impacts to all human rights

- Policies and enforcement approaches should adhere to human rights principles

- Policy development and enforcement should be informed by stakeholder engagement

- Policy development should address the most severe impacts first

- Policies and enforcement should honor human rights principles when faced with conflicting requirements

- Policies enforced in conflict-affected areas should receive enhanced attention

- Reporting and appeals channels should be established to help identify policy violations and correct enforcement errors

- Policies should be updated and enforcement reviewed on an ongoing basis

BSR

# Elements of a human rights-based approach to policies and enforcement

- **Policies should reference human rights when relevant**

- Policies and enforcement should honor human rights principles when faced with conflicting

- Policies should address impacts ... should

- Policies and enforcement approaches ... adhere to human rights principles ... and

- Policy development and enforcement ... informed by stakeholder engagement

- Policy development should address the ... severe impacts first

**"Usage Restrictions" in Microsoft Enterprise AI Services Code of Conduct**

---

"Customers, users, and applications built with Microsoft AI Services must NOT use the services […] to make decisions or take actions without appropriate human oversight as part of an application that may have a consequential impact on any individual's legal position, financial position, life opportunities, employment opportunities, or **human rights**, or may result in physical or psychological harm to an individual."

BSR

# Benefits of integrating human rights into policies and enforcement

**1**

**Global approaches across geographic borders**

**2**

**Enforcement that is thoughtful, consistent, and rights-informed**

**3**

**Flexibility to adapt to evolving rights over time**

BSR

# Guide 7:

# Aligning Transparency and Disclosure Practices with Human Rights Responsibilities

# Key Definitions

A **disclosure** is information that an entity reveals about itself that would not otherwise be available or easily discoverable

### Entity / Company-Level Disclosures

- Governance
- Strategy
- Risks and impacts
- Indicators / metrics / targets

### Model / Product-Level Disclosures

- Model / system cards
- Datasheets

BSR

# The human rights disclosure landscape

| Voluntary Standards | Mandatory Standards |
|---|---|
| The UN Guiding Principles on Business and Human Rights (Principle 21) | EU Corporate Sustainability Reporting Directive (CSRD) |
| OECD Due Diligence Guidelines for Responsible Business Conduct | European Sustainability Reporting Standards (ESRS) |
| The Global Reporting Initiative (GRI) | |
| International Financial Reporting Standards (IFRS) Foundation's Sustainability Disclosure Standards | |

BSR

# Benefits of AI Disclosures

- **The discipline of putting together disclosures that meet standards can improve responsible AI workflows**—e.g. help prioritize key issues, establish baselines, craft mitigations, track progress, and guide resource allocation

- **Creating disclosure workflows helps prepare for external audits or evaluations**, which are increasingly being adopted into regulation

- **Disclosures can address investor and other key stakeholder concerns**, reassuring stakeholders you are tracking key issues and taking steps to address risks

- **Disclosure about impacts to people spurs progress on disclosure across the AI industry** by motivating peers and

- **Disclosure improves knowledge/understanding of the public and policymakers** about nuanced and complex issues, informing public policy and regulation

BSR

# Applying human rights and sustainability disclosure best practices to AI

## Report Content

| | | | | |
|---|---|---|---|---|
| **PRINCIPLE 1**<br>MATERIALITY AND CONCISENESS | **PRINCIPLE 2**<br>STRATEGIC AND FORWARD LOOKING | **PRINCIPLE 3**<br>SUSTAINABILITY CONTEXT | **PRINCIPLE 4**<br>KEY PERFORMANCES INDICATORS AND NARRATIVE | **PRINCIPLE 5**<br>COMPLETENESS |

## Report Quality

| | | | | |
|---|---|---|---|---|
| **PRINCIPLE 6**<br>STAKEHOLDER ENGAGEMENT | **PRINCIPLE 7**<br>BALANCE | **PRINCIPLE 8**<br>ASSURANCE | **PRINCIPLE 9**<br>CONSISTENCY AND COMPARBILITY | **PRINCIPLE 10**<br>CONNECTIVITY OF INFORMATION |

BSR

# Disclosure Examples Across the Value Chain

| Suppliers | Foundation Model Developers | Downstream Developers | Deployers |
|---|---|---|---|
| • Dataset documentation | • Model and system cards<br><br>• Responsible use guides<br><br>• Non-technical disclosures (e.g. blog posts) | • Non-technical disclosures for users (e.g. infographics)<br><br>• Non-technical disclosures for other stakeholders (e.g., hallucination rates, summary of stakeholder feedback) | • Non-technical disclosures for users (e.g. product info) |

BSR

# Guide 8:

# Remedy for AI-Related Harms

BSR

# Why is remedy important?

- **Mitigating business risks –** Remedy mitigates business risks (e.g. lawsuits, public shaming campaigns)

- **Source of business intelligence** – Remedy mechanisms can be a useful source of business intelligence (e.g. how products are impacting people in the world)

# Remedy in the UNGPs

- Individuals whose rights have been harmed by businesses must have access to remedy

- Companies should provide or cooperate in remediation for impacts they cause or contribute to.

- Companies should establish or participate in effective operational grievance mechanisms

BSR

## Remedy in the UNGPs

- Individuals whose rights have been harmed by businesses must have access to remedy

- Companies should provide or cooperate in remediation for impacts they cause or contribute to.

- Companies should establish or participate in effective operational grievance mechanisms

## Why is remedy important?

- **Mitigating business risks –** Remedy mitigates business risks (e.g. lawsuits, public shaming campaigns)

- **Source of business intelligence** – Remedy mechanisms can be a useful source of business intelligence (e.g. how products are impacting people in the world)

# Five Categories of Remedy

| | |
|---|---|
| **Satisfaction** | Ceasing the violation, acknowledging the harm, disclosing the truth, providing an apology, and sanctioning those responsible |
| **Restitution** | Restoring, to the extent possible, whatever has been lost and returning the rightsholder to the state before the harm occurred |
| **Guarantees of non-repetition** | Changes to policies and procedures to prevent future harms, or the taking of disciplinary action |
| **Rehabilitation** | Medical, psychological, legal, social, or other services to restore the victim |
| **Compensation** | Money or other benefits, where damage can be financially assessed |

BSR

# Remedy Across the GenAI Value Chain

| Value Chain Actor | Remedy Examples |
|---|---|
| **1**    **Suppliers** | • Channels for removing personal data from datasets<br>• Compensates victims of copyright violations<br>• Creates and publicizes new data procurement process |
| **2**    **Foundation Model Developers** | • Fine-tunes model to prevent recurrence of rights violation<br>• Publicizes new fine-tuning method that reduces odds of further rights violations<br>• Publishes blog post acknowledging error and detailing mitigations |
| **3**    **Downstream Developers** | • Create new product feature that manages risk<br>• Provides free access to healthcare services to remedy impacts on rights to health<br>• Reporting channel feedback informs product development |
| **4**    **Deployers** | • Apology letters sent to users<br>• Compensation to supplier for lost revenue<br>• GenAI tool no longer used |

# Remedy Across the GenAI Value Chain

| Value Chain Actor | Remedy Examples |
|---|---|
| **1**   **Suppliers** | • Channels for removing personal data from datasets<br>• Compensates victims of copyright violations<br>• Creates and publicizes new data procurement process |
| **2**   **Foundation Model Developers** | • Fine-tunes model to prevent recurrence of rights violation<br>• Publicizes new fine-tuning method that reduces odds of further rights violations<br>• Publishes blog post acknowledging error and detailing mitigations |
| **3**   **Downstream Developers** | • Create new product feature that manages risk<br>• Provides free access to healthcare services to remedy impacts on rights to health<br>• Reporting channel feedback informs product development |
| **4**   **Deployers** | • Apology letters sent to users<br>• Compensation to supplier for lost revenue<br>• GenAI tool no longer used |

# Remedy Across the GenAI Value Chain

| Value Chain Actor | Remedy Examples |
|---|---|
| **1** **Suppliers** | • Channels for removing personal data from datasets<br>• Compensates victims of copyright violations<br>• Creates and publicizes new data procurement process |
| **2** **Foundation Model Developers** | • Fine-tunes model to prevent recurrence of rights violation<br>• Publicizes new fine-tuning method that reduces odds of further rights violations<br>• Publishes blog post acknowledging error and detailing mitigations |
| **3** **Downstream Developers** | • Create new product feature that manages risk<br>• Provides free access to healthcare services to remedy impacts on rights to health<br>• Reporting channel feedback informs product development |
| **4** **Deployers** | • Apology letters sent to users<br>• Compensation to supplier for lost revenue<br>• GenAI tool no longer used |

BSR

# Remedy and the single point of contact

**The duty to coordinate remedy for harms that require business action should lie with a single point of contact.**

- The single point of contact will usually be **the value chain entity that is directly interfacing with the affected stakeholder**.

- There will be cases, where entities from the broader remedy ecosystem, such as law enforcement or social services organizations, must play a role in effective remedy.

BSR

# Key Takeaways

**1**

**Human rights should be the foundation for other approaches**

**2**

**Human rights concepts and frameworks can augment responsible AI workflows**

**3**

**You don't need to be a human rights expert to take a human rights-based approach**

BSR

# Q&A

BSR