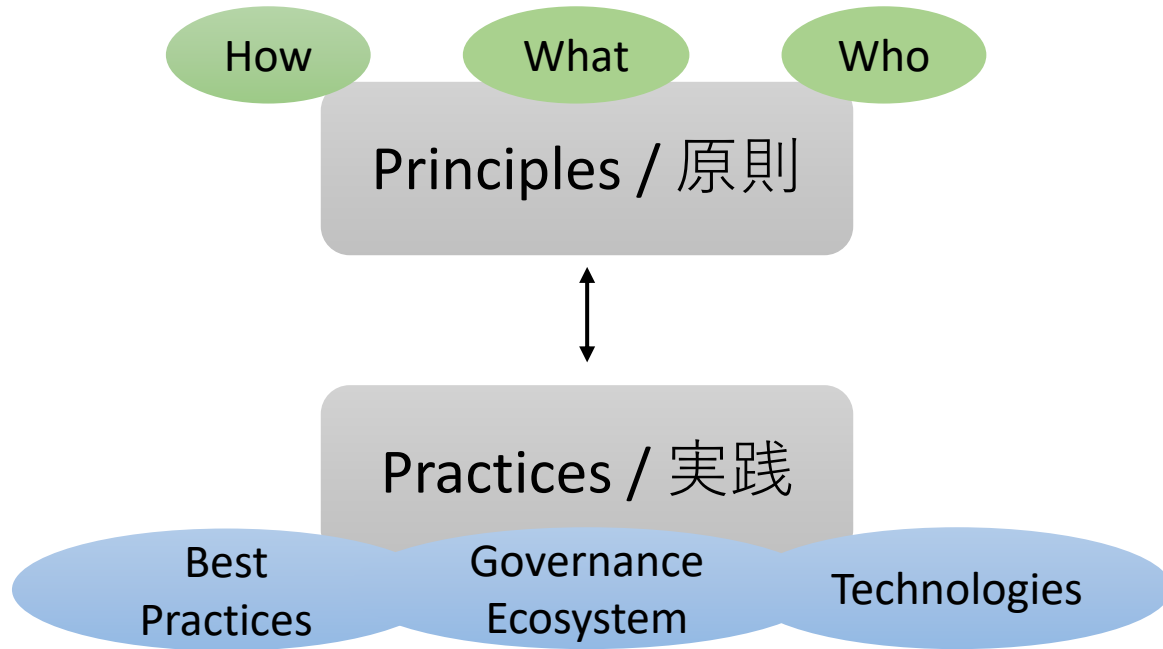
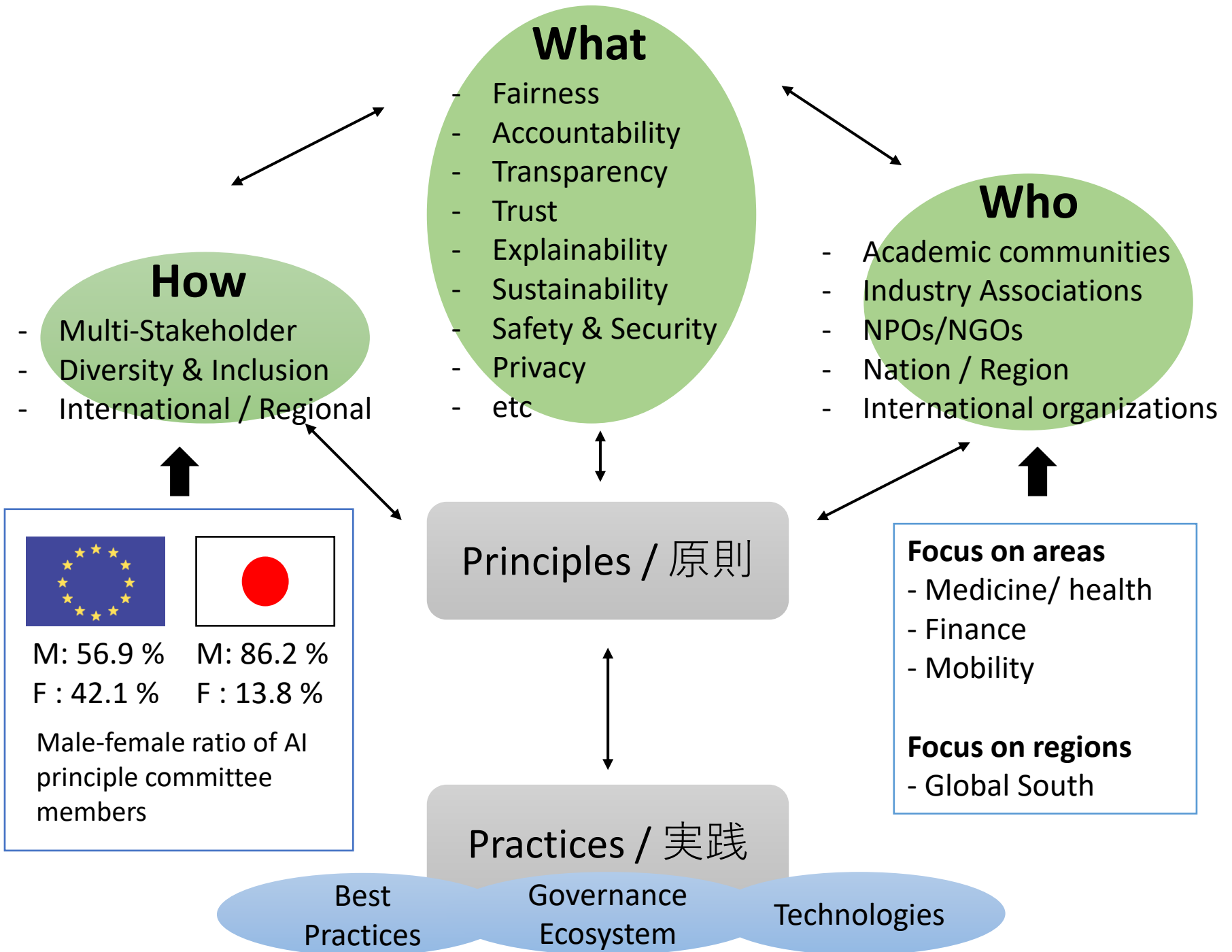


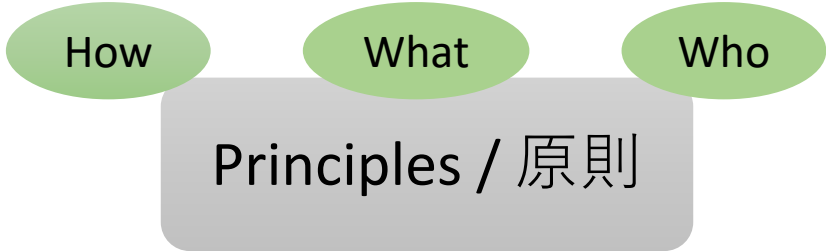
AIガバナンス 原則から実践へ

東京大学未来ビジョン研究センター准教授
理化学研究所AIPセンター客員研究員

江間 有沙







Best Practices

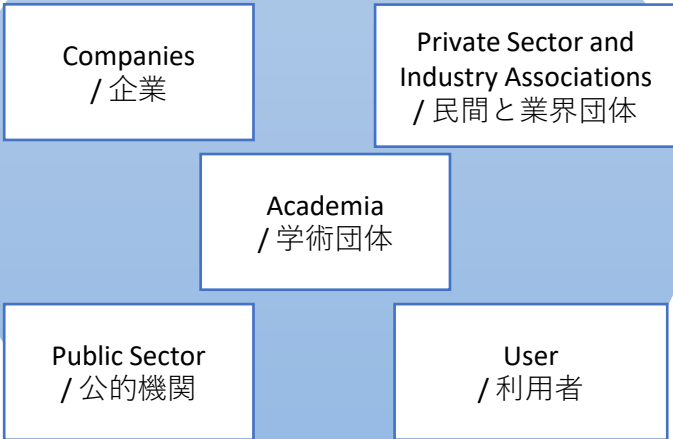
Utilizations / 利活用

- Medicine / 医療
- Finance / 金融
- Risk Prevention / 防災
- Agriculture / 農業
- Military / 軍事

Work / 働き方

- Human-machine interaction / 人と機械の関係性

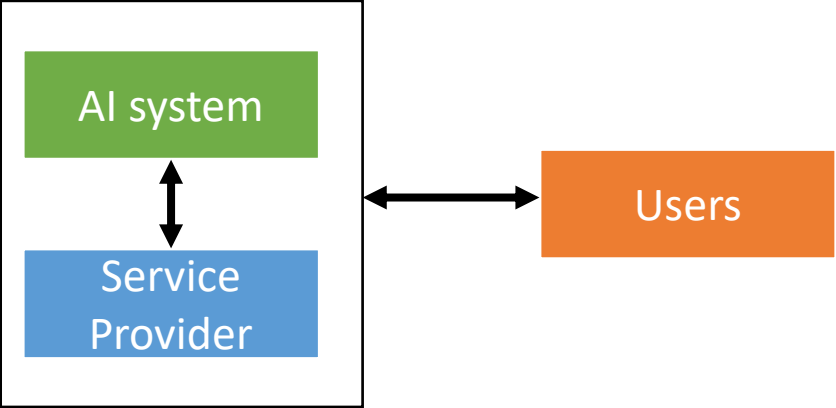
Governance Ecosystems



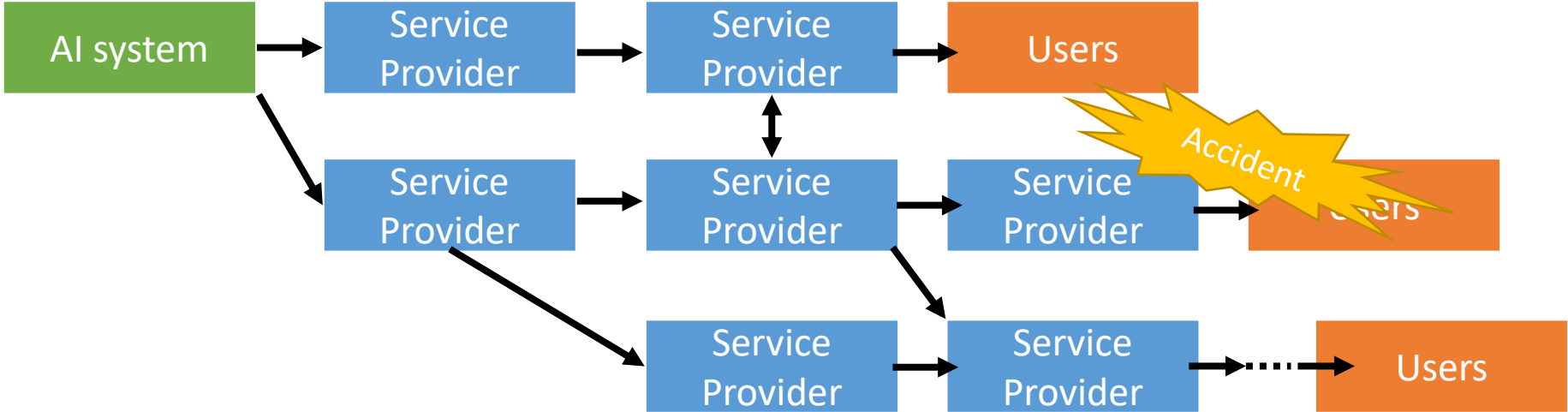
Technologies

- Trustworthy AI / 信頼されるAI
 - Fair algorithm and Data / 公平なアルゴリズムとデータ
 - Explainable AI / 説明可能AI
 - Sustainable AI / 持続可能AI
 - Robust AI / 頑健なAI
- Next AI & Environment / AIの次と環境
 - 5G Network / 通信網
 - Neuroscience & Quantum / 脳科学・量子

B2C Company



B2B2C supply chain



Ecosystem of AI governance / AIガバナンスのエコシステム

Companies / 企業

- Corporate vision / ビジョン
- AI Principles / AI原則
- Risk Assessment / リスク評価
- Risk Control / リスクコントロール
- Employee/ Employer Education / 教育

Private Sector and Industry Association / 民間と業界団体

- Guidelines / ガイドライン
- Standards / 標準
- Audit / 監査
- Insurance / 保険
- Fact Check / ファクトチェック

Academia / 学術団体

Public Sector / 公的機関

- Hard Law, Soft Law / 法や規制
- Third-Party Incident Committee / 事故調査制度
- Whistleblower System / 内部告発制度

Users / 利用者

- Education / 教育
 - Engineer / エンジニア
 - Experts / 専門家
 - Policy makers / 政策関係者
 - General publics / 一般

Companies
/ 企業

Corporate vision and values /
ビジョンや実現したい価値

↔ AI Principles

Consider AI Service Requirements &
Technologies / AIサービス要件と技術把握

↔

- Facial Recognition
- Building Entrance

Create Risk Scenarios
/ リスクシナリオ作成

↔

- Data bias & wrong feedback
- Influence of the noise

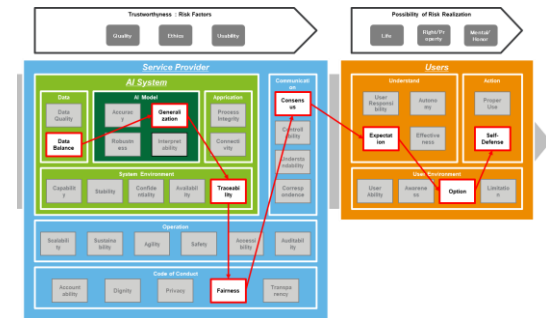
Draw Risk Chains
/ リスク要因の関係性作成

↔

- Usage Policy
- Education
- Data modification

Consider Risk Control
/ リスク管理実行

↔



リスクアセスメント&コントロール (Case2 : 採用AI)

- リスクの複雑性評価 → リスクコントロールを検討 -

実現すべき価値・目的 (リスクの影響先)	リスク No.	リスクシナリオ	技術的 難易度	環境変 化	利用者 起因	R C	コントロールのサマリー		
							AIシステム	サービスプロバイダ	ユーザー
1 人材採用レベル の維持・向上	R001	適切な評価	○			●	システム環境の確保 個別モデルの開発	目標精度の設定・見直し 再学習の指示	正確なフィードバック
	R002	予測性能の維持	○			●	十分な正解率の確保 検証可能性の確保	予測性能の検証 再学習	代替運用
	R003	ノイズによる影響	○			●	モデルの頑健性 判断根拠の出力	判断根拠の分かりやすさ 判断根拠の検証	判断根拠の検討
	R004	虚偽の申込	○			●	十分な正解率の確保 判断根拠の出力	過去の異常事例の検証 利用部門との連携	最終選考プロセス
	R005	過度なAI依存	○		○	●	判断根拠の出力	判断精度の開示 判断根拠の分かりやすさ	予測性能・リスクの認識 最終判断プロセス
	R006	誤ったフィードバック	○		○	●	学習データの検証 学習時の情報の保管	判断精度・異常値の検証 再学習	正確なフィードバック 人事システムからの反映
	R007	人材トレンドの変化	○	○		●	データ分布変化の認識 汎化性能の見直し	分布／正解率の監視 再学習	人材トレンド変化の認識
	R008	新たな職種	○	○		●	開発環境の準備 学習データ確保 モデルの開発	開発体制の整備 開発したモデルの検証	要求仕様の定義
2 採用活動に係る コストの削減	R009	コスト超過					適切な価格設定	コスト管理	
3 海外グループを含 めたサービス提供	R010	地域の会社への対応	○	○		●	システム環境の確保 個別モデルの開発	個別の目標精度の設定 モデルの性能監視 開発体制の確保	
	R011	不十分な開発スピード					再学習環境の準備 モデル開発者の確保	PJ体制の確保	
4 企業の社会的責 任(公平性のある 採用活動)	R012	判断根拠情報の不正販売	○		○	●	操作ログの記録	不正利用の監視 外部弁護士との連携	内部牽制
	R013	公平性	○		○	●	データの偏り モデルの汎化性	公平性ポリシーの検討 ネガティブな判断の開示	AIの判断傾向の理解 公平な最終判断
	R014	予測結果の目的外利用			○		データ保護	アクセス管理 目的内利用	データの取扱
	R015	風評被害			○		データ保護	職業倫理の教育	職業倫理の教育
	R016	プライバシー保護			○		データ保護	法令順守の教育	法令順守の教育 データの取扱

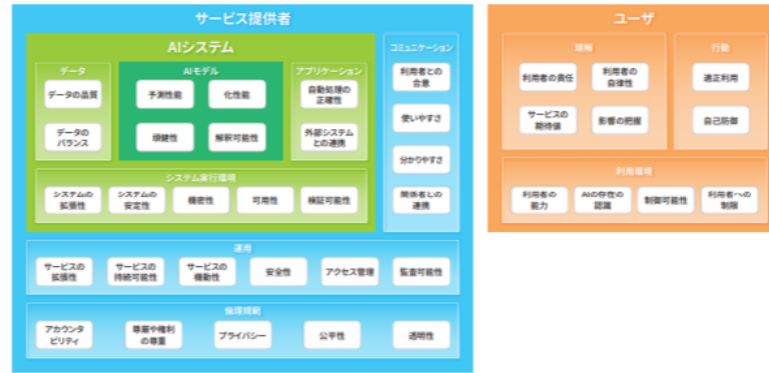


2021年6月
**リスクチェーンモデル (RCModel)
 ガイド Ver 1.0**

東京大学未来ビジョン研究センター 技術ガバナンス研究ユニット
 AIガバナンスプロジェクト

ケース事例も公開中！

- Case01.採用AI (2021/07)
- Case02.無人コンビニ(2021/07)
- Case03.送電線の外観検査ドローン (2021/07)
- Case04.不良品検知AI (2021/07)
- Case05.道案内ロボット(2021/07)
- Case06.再犯可能性の検証AI (2021/07)



AIシステム		サービス提供者		ユーザ		
内容	主なコントロールの例	検証対象	内容	検証対象	内容	主なコントロールの例
データ		信頼性		理解		
データの品質	データの品質確保	データ信頼性の検証	説明責任、利用者の保護	責任範囲の明確化 適切な期待値の維持	利用者の責任	利用者の責任の理解 AIシステムの概要への合意
データのバランス	データの分布	汎化性能の検証	尊重や権利の尊重	利用者の優先順位の検討	利用者の自律性	自律的に判断する内容の理解 AIの判断の訂正
アプリケーション		公平性	プライバシー	サービス全体での公平性	サービスの期待値	AIサービスに対する期待値の理解 期待値と実態等の理解
自動処理の正確性	補助的な自動処理	差別のない自動処理	透明性	AIサービスに関する必要情報の明示	影響の把握	AIサービスが利用者の理解 AI利用のワークスタイルの変化の理解
外部システムとの連携	外部システム等との連携	外部システム等との連携	適用	AIサービスに依存する必要な情報の明示	利用環境	利用環境
システム実行環境		システム実行環境	サービスの拡張性	利用環境の増加等へ対応できる サービス設計	利用者の能力	サービス利用において 必要な知識・スキルの習得
システムの拡張性	利用環境の増加等へ対応できる システム環境	機能モデルを実装する環境	サービスの持続可能性	持続的なサービスの維持	AIの存在の認識	AIの存在の認識 利用者の告知
システムの安定性	安定稼働するシステム環境	システム環境の確保 AIのメンテナン	安全性	透明の対応	制御可能性	利用者の制約 AI判断の訂正機能
機密性	機密や権利が保護されたシステム環境	システム環境の確保 アクセスコントロールの実装	アクセス管理	サービス全体での安全性の確保	利用者の制約	不正利用等の防止に向けた 利用者の制約
可用性	必要な時に使用できるシステム環境	システム環境の確保 システム稼働時間の確保	監査可能性	適切なシステム 利用履歴・経過の管理	行動	適切な目的に沿ったサービス利用 期待値への訂正・修正
検証可能性	AIサービスを事後検証できるための性質	十分なログ・モニタリングの確保 期待値や権利侵害への対応 期待値や権利侵害への対応 期待値や権利侵害への対応	コミュニケーション	第三者を含む必要も監視	自己防衛	不測が発生した際のクレーン 環境変化による事前対応
			利用者と			
			使いやすさ			
			分かりやすさ			
			関係者との連携			

責任ある研究・イノベーションとは？

- 2000年代前半から欧米で議論され始める
- 欧州連合のHorizon 2020（2014-2020）でも言及

応答性と変化への適応性

変化する状況、知識や展望に対応して思考や行動様式、組織構造を包括的に変更できること

先見性と省察性

研究やイノベーションがどのように未来を形成するかをよりよく理解するための前提、価値観や目的を熟考し、影響を想定すること

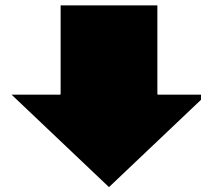
多様性と包摂性

研究・イノベーションの実践、普及と意思決定において、科学技術の発展の早い段階から多様な人を巻き込むこと

公開性と透明性

人々が情報を精査し対話できるように、方法、結果や影響についてバランスよく伝達すること

- コリングリッジのジレンマ（1980）
 - 情報の問題：技術が社会で使われる前にその影響力を予測することは難しい
 - 力の問題：一度普及してしまった技術は制御するのが難しい



- 設計の段階から影響の想定が求められる
- 「社会実験」ではなく「実験社会」
- 私たち一人一人の役割と責任とは？