

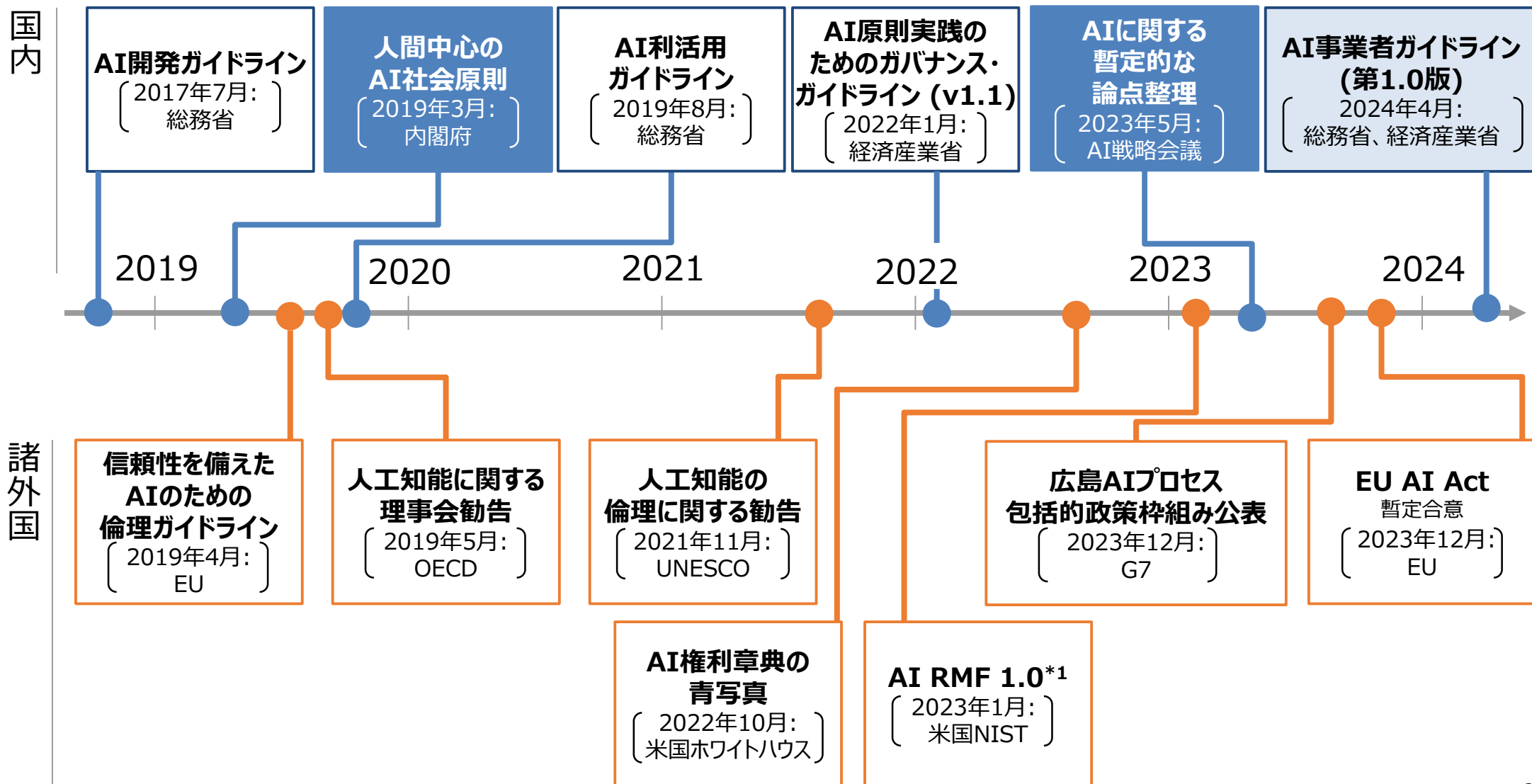
# AI事業者ガイドライン(第1.0版) について

---

経済産業省 商務情報政策局 情報経済課  
ガバナンス戦略国際調整官 飯野 悠介  
令和6年8月23日

# 【参考】AI事業者ガイドラインの経緯

- 「人間中心のAI社会原則」において、AIは、「社会に多大なる便益をもたらす一方で、その社会への影響が大きいゆえに、適切な開発と社会実装が求められる」ことが指摘。
- 「AIに関する暫定的な論点整理」において、AIの開発・提供・利用を促進しつつも、リスクへの適切な対処「ガードレール」が必要であると言及。

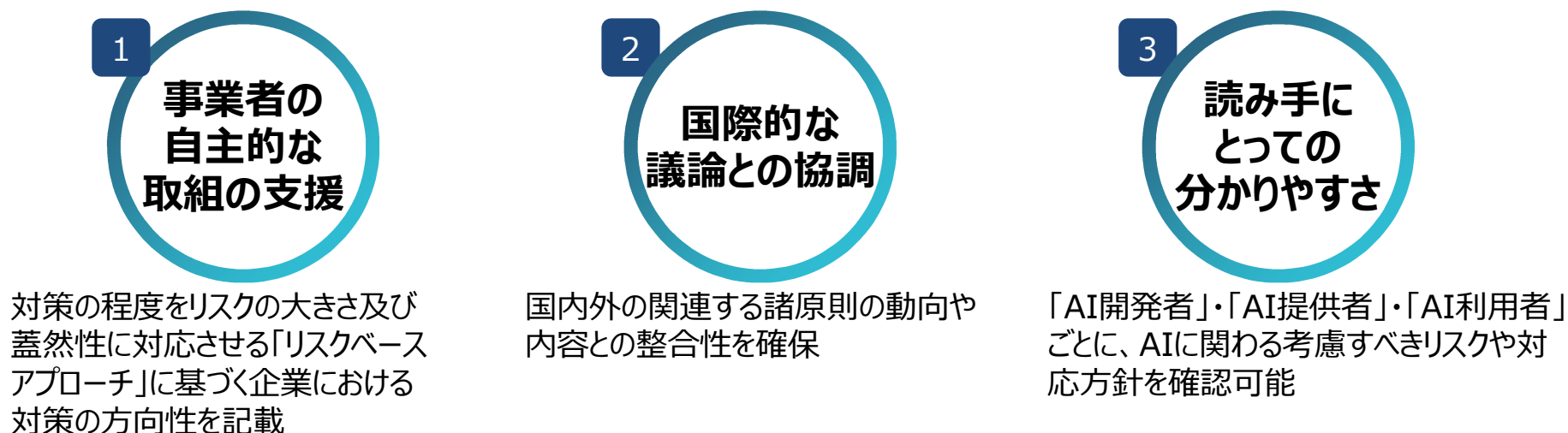


\*1: AI Risk Management Framework 1.0

# 「AI事業者ガイドライン」の基本的な考え方

- 本ガイドラインは、「**1** 事業者の自主的な取組の支援」、「**2** 国際的な議論との協調」、「**3** 読み手にとっての分かりやすさ」を基本的な考え方としている
- 加えて、「マルチステークホルダー」で検討を重ね実効性・正当性を重視するとともに、「Living Document」として今後も更新を重ねていく

## 考え方



## プロセス

### マルチステークホルダー

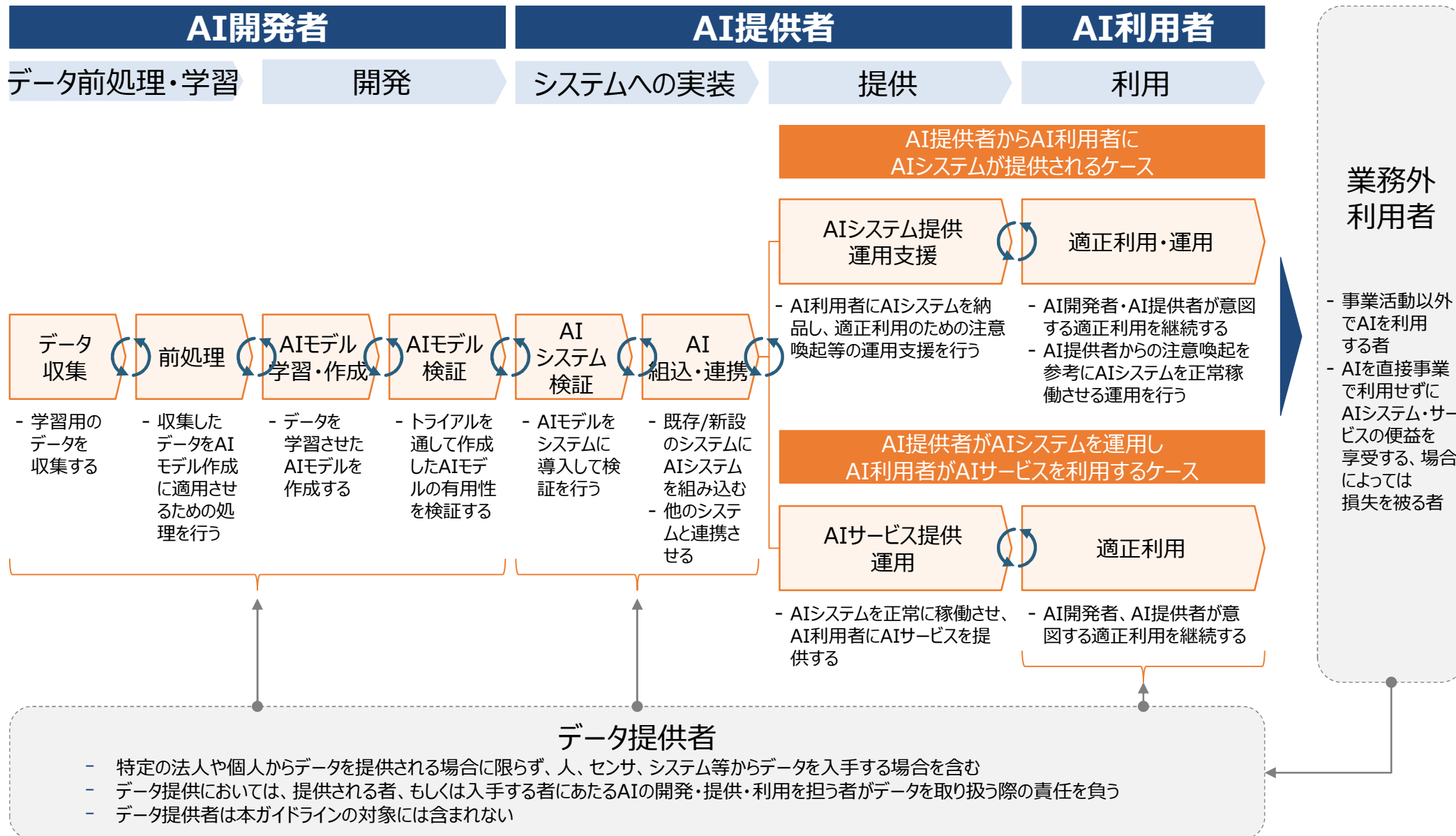
教育・研究機関、一般消費者を含む市民社会、民間企業等で構成されるマルチステークホルダーで検討を重ねることで、実効性・正当性を重視したものととして策定

### Living Document

AIガバナンスの継続的な改善に向け、アジャイル・ガバナンスの思想を参考にしながら適宜、更新

# 一般的なAIの事業活動を担う主体

- AIライフサイクルにおける具体的な役割を考慮し、AIの事業活動を担う立場として、「AI開発者」、「AI提供者」、「AI利用者」の3つに大別して整理する ※「データ提供者」、「業務外利用者」は対象外とする



# 「AI事業者ガイドライン」本編、別添の位置づけ

- 本編では、事業者がAIの安全安心な活用を行い、AIの便益を最大化するために重要な「どのような社会を目指すのか（基本理念=why）」及び「どのような取組を行うか（指針=what）」を示した
- 別添（付属資料）では、「具体的にどのようなアプローチで取り組むか（実践=how）」を示すことで、事業者の具体的な行動へとつなげることを想定している

## 本編（why, what）

## 別添（付属資料）（how）



どのような社会を  
目指すのか  
(基本理念=why)



どのような取組を  
行うか  
(指針=what)



どのようなアプローチで  
取り組むか  
(実践=how)



# 【参考】別添1 B. AIによる便益/リスク（AIによる便益）

## 主な記載内容

- 便益を享受する最終利用者に焦点を当ててAIによる便益を整理し、その理解促進につなげる

	開発	マーケティング	販売	物流・流通	顧客対応	法務	ファイナンス	人事
従来から存在する便益の例	コード検証、ドキュメント作成の自動化	広告用メールの自動配信	受注後の対応メール等の自動発信	需要予測に基づく生産・在庫数最適化	チャットボットによる自動対応	翻訳	財務諸表の自動作成	給与計算等の自動化
（生成AIで更に向上）								
	類似コード・データの抽出・検証	データに基づいたパーソナライゼーション広告	チャネル別、ニーズ別の売上予測	配送ルート最適化	過去の問合せ内容に基づいたFAQ作成	法務文章のレビュー	過去実績にもとづいた将来予測、不正検知	職務経歴書等に基づいた人材需要マッチング
生成AI特有の便益の例	学習データの生成、コーディングアシスタント、新製品のブレインストーミング	販売促進（マーケティング素材・キャッチコピー等）の自動作成	営業トークスクリプトの自動作成	物流条件交渉のアシスタント	対応内容の自動生成、要約	規定に基づいた契約書ドラフトの自動生成	文脈を踏まえた上での社内問合せ対応	文脈を踏まえた上での人事面接の対応

# 【参考】別添1 B. AIによる便益/リスク（AIによるリスク）

## 主な記載内容

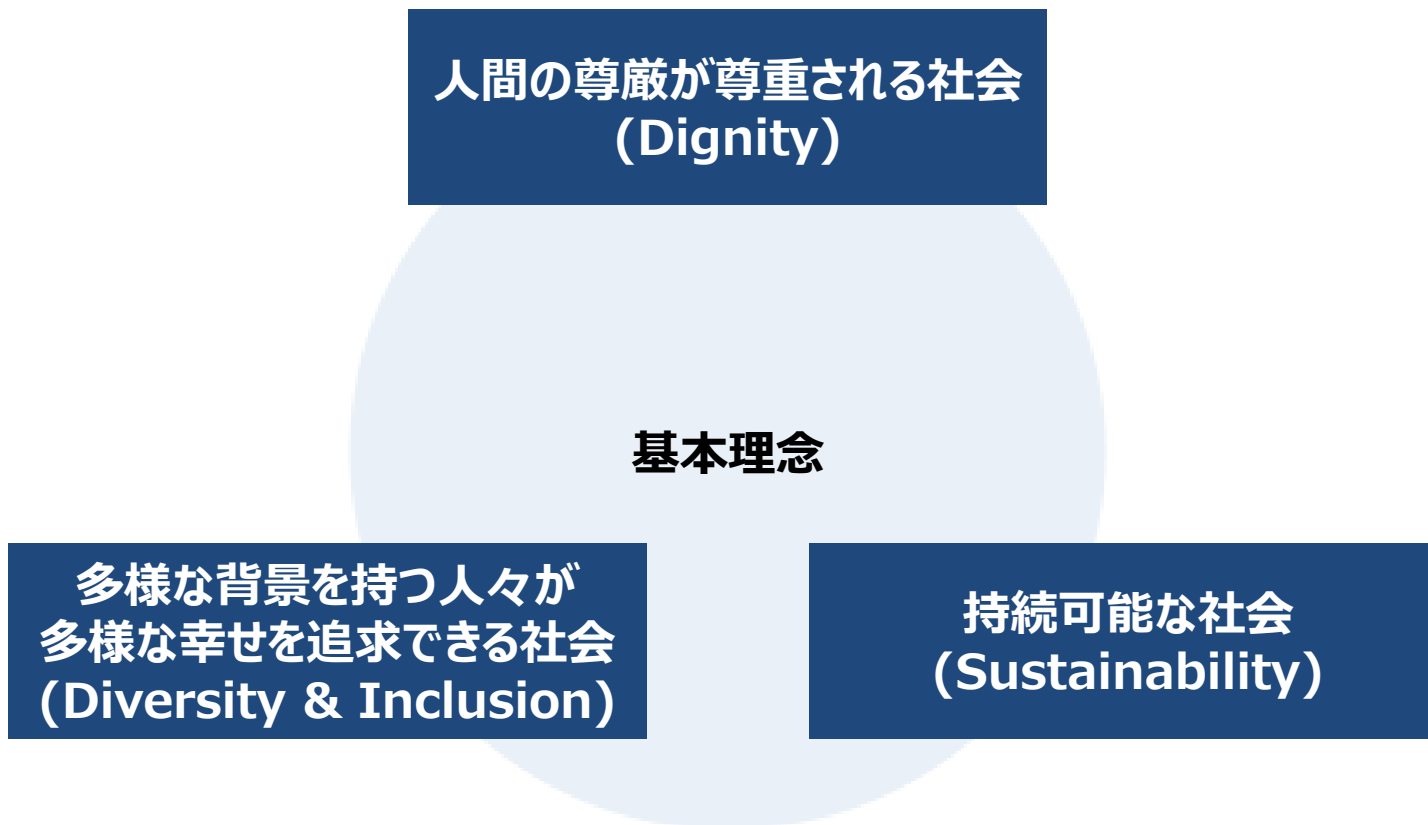
- 従来型のAIからのリスク及び生成AIで特に顕在化したリスクについて、主なものを整理
- AIによるリスクを知ること、AIガバナンスの必要性の理解を深める

	リスク	事例	対応する「共通の指針」
従来型AIから存在するリスク	バイアスのある結果及び差別的な結果の出力	<ul style="list-style-type: none"> <li>IT企業が自社で開発したAI人材採用システムが女性を差別するという機械学習面の欠陥を持ち合わせていた</li> </ul>	1) 人間中心 3) 公平性
	フィルターバブル及びエコーチェンバー現象	<ul style="list-style-type: none"> <li>SNS等によるレコメンドを通じた社会の分断が生じている</li> </ul>	1) 人間中心
	多様性の喪失	<ul style="list-style-type: none"> <li>社会全体が同じモデルを、同じ温度感で使った場合、導かれる意見及び回答がLLMによって収束してしまい、多様性が失われる可能性がある</li> </ul>	1) 人間中心
	不適切な個人情報の取扱い	<ul style="list-style-type: none"> <li>透明性を欠く個人情報の利用及び個人情報の政治利用も問題視されている</li> </ul>	1) 人間中心 4) プライバシー保護
	生命、身体、財産の侵害	<ul style="list-style-type: none"> <li>AIが不適切な判断を下すことで、自動運転車が事故を引き起こし、生命や財産に深刻な損害を与える可能性がある</li> <li>トリアージにおいては、AIが順位を決定する際に倫理的なバイアスを持つことで、公平性の喪失等が生じる可能性がある</li> </ul>	2) 安全性 3) 公平性
	データ汚染攻撃	<ul style="list-style-type: none"> <li>AIの学習実施時及びサービス運用時には学習データへの不正データ混入、サービス運用時ではアプリケーション自体を狙ったサイバー攻撃等のリスクが存在する</li> </ul>	5) セキュリティ確保
	ブラックボックス化、判断に関する説明の要求	<ul style="list-style-type: none"> <li>AIの判断のブラックボックス化に起因する問題も生じている</li> <li>AIの判断に関する透明性を求める動きも上がっている</li> </ul>	6) 透明性 7) アカウンタビリティ
	エネルギー使用量及び環境の負荷	<ul style="list-style-type: none"> <li>AIの利用拡大により、計算リソースの需要も拡大しており、結果として、データセンターが増大しエネルギー使用量の増加が懸念されている</li> </ul>	1) 人間中心
生成AIで特に顕在化したリスク	悪用	<ul style="list-style-type: none"> <li>AIの詐欺目的での利用も問題視されている</li> </ul>	2) 安全性 8) 教育・リテラシー
	機密情報の流出	<ul style="list-style-type: none"> <li>AIの利用においては、個人情報や機密情報がプロンプトとして入力され、そのAIからの出力等を通じて流出してしまうリスクがある</li> </ul>	5) セキュリティ確保 8) 教育・リテラシー
	ハルシネーション	<ul style="list-style-type: none"> <li>生成AIが事実と異なることをもつもらしく回答する「ハルシネーション」に関してはAI開発者・提供者への訴訟も起きている</li> </ul>	2) 安全性 8) 教育・リテラシー
	偽情報、誤情報を鵜呑みにすること	<ul style="list-style-type: none"> <li>生成AIが生み出す誤情報を鵜呑みにすることがリスクとなりうる</li> <li>ディープフェイクは、各国で悪用例が相次いでいる</li> </ul>	1) 人間中心 8) 教育・リテラシー
	著作権との関係	<ul style="list-style-type: none"> <li>知的財産権の取扱いへの議論が提起されている</li> </ul>	2) 安全性
	資格等との関係	<ul style="list-style-type: none"> <li>生成AIの活用を通じた業法免許や資格等の侵害リスクも考えうる</li> </ul>	2) 安全性
	バイアスの再生成	<ul style="list-style-type: none"> <li>生成AIは既存の情報に基づいて回答を作るため既存の情報に含まれる偏見を増幅し、不公平や差別的な出力が継続/拡大する可能性がある</li> </ul>	3) 公平性



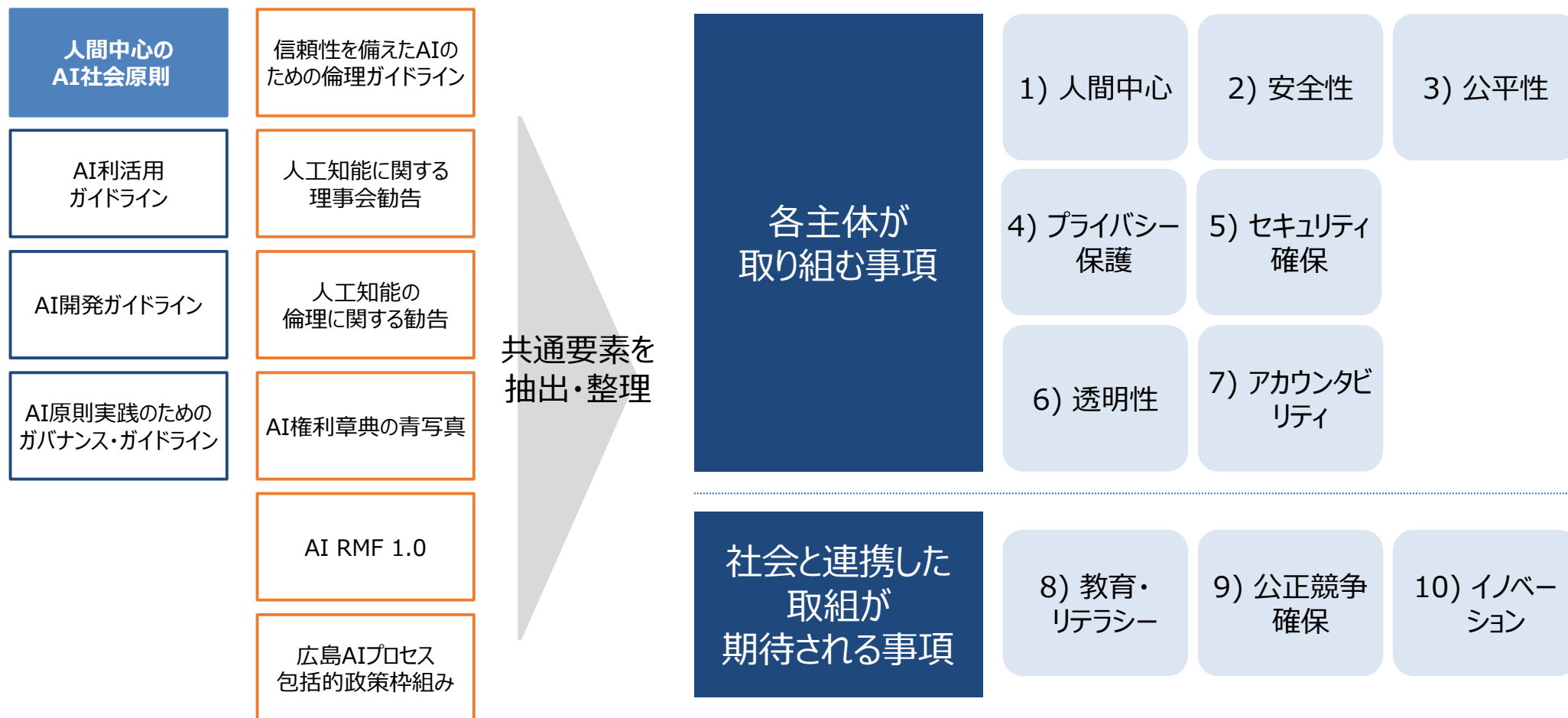
# 基本理念

- 「人間中心のAI社会原則」において、AIがSociety 5.0の実現に貢献し、AIを人類の公共財として活用することで、社会の在り方の質的变化や真のイノベーションを通じて地球規模の持続可能性へとつなげることが重要であると述べられている
- 加えて、下記の3つの価値を「基本理念」として尊重し、「その実現を追求する社会を構築していくべき」としており、この普遍的な考え方は、今後も目指すべき理念であり続けている



# 各主体に共通の指針

- AIの活用による目指すべき社会の実現のために各主体が連携して取り組む内容を原則としてまとめた上で、「共通の指針」として整理する
- 「共通の指針」は、「人間中心のAI社会原則」を土台としつつ、諸外国における議論状況や、新技術の台頭に伴い生じるリスクへの対応等を反映している
- その結果、各主体が取り組む事項、及び社会と連携して取り組むことが期待される事項に分類される



## 【参考】各主体に共通の指針 [1/2]

- 各主体は、1) 人間中心に照らし、法の支配、人権、民主主義、多様性及び公平公正な社会を尊重する
- 憲法、知的財産関連法令及び個人情報保護法をはじめとする関連法令、AIに係る個別分野の既存法令等を遵守すべきであり、国際的な指針等の検討状況についても留意することが重要
- AIガバナンスを構築し継続的に運用（AIのもたらすリスクの程度や各主体の資源制約に配慮しつつ実施）

## 指針

## 内容（主な項目の抜粋）

指針	内容（主な項目の抜粋）
各主体が 取り組む事項	1) 人間中心 <ul style="list-style-type: none"> <li>✓ AI が人々の能力を拡張し、多様な人々の多様な幸せ（well-being）の追求が可能となるよう行動する</li> <li>✓ AI が生成した<b>偽情報・誤情報・偏向情報</b>が社会を不安定化・混乱させるリスクが高まっていることを認識した上で必要な対策を講じる</li> <li>✓ より多くの人々がAIの恩恵を享受できるよう<b>社会的弱者によるAIの活用</b>を容易にするよう注意を払う</li> </ul>
	2) 安全性 <ul style="list-style-type: none"> <li>✓ 適切なリスク分析を実施し、<b>リスクへの対策</b>を講じる</li> <li>✓ 主体のコントロールが及ぶ範囲で本来の利用目的を逸脱した提供・利用により危害が発生することを避ける</li> <li>✓ AIシステム・サービスの特性及び用途を踏まえ、学習等に用いるデータの正確性等を検討するとともに、<b>データの透明性の支援、法的枠組みの遵守</b>、AIモデルの更新等を合理的な範囲で適切に実施する</li> </ul>
	3) 公平性 <ul style="list-style-type: none"> <li>✓ 特定の個人ないし集団へのその人種、性別、国籍、年齢、政治的信念、宗教等の多様な背景を理由とした<b>不当で有害な偏見及び差別をなくす</b>よう努める</li> <li>✓ AIの出力結果が公平性を欠くことがないよう、AIに単独で判断させるだけでなく、適切なタイミングで人間の判断を介在させる利用を検討した上で、無意識や潜在的な<b>バイアスに留意</b>し、AIの開発・提供・利用を行う</li> </ul>
	4) プライバシー保護 <ul style="list-style-type: none"> <li>✓ 個人情報保護法等の<b>関連法令の遵守</b>、<b>各主体のプライバシーポリシーの策定・公表</b>により、社会的文脈及び人々の合理的な期待を踏まえ、ステークホルダーのプライバシーが尊重され、保護されるよう、その重要性に応じた対応を取る</li> </ul>
	5) セキュリティ確保 <ul style="list-style-type: none"> <li>✓ AI システム・サービスの<b>機密性・完全性・可用性を維持</b>し、常時、AIの安全な活用を確保するため、その時点での技術水準に照らして合理的な対策を講じる</li> <li>✓ AIシステム・サービスに対する外部からの攻撃は日々新たな手法が生まれており、これらの<b>リスクに対応するための留意事項を確認</b>する</li> </ul>

## 【参考】各主体に共通の指針 [2/2]

- 各主体は、1) 人間中心に照らし、法の支配、人権、民主主義、多様性及び公平公正な社会を尊重する
- 憲法、知的財産関連法令及び個人情報保護法をはじめとする関連法令、AIに係る個別分野の既存法令等を遵守すべきであり、国際的な指針等の検討状況についても留意することが重要
- AIガバナンスを構築し継続的に運用（AIのもたらすリスクの程度や各主体の資源制約に配慮しつつ実施）

## 指針

## 内容（主な項目の抜粋）

指針	内容（主な項目の抜粋）
各主体が 取り組む事項 (続き)	6) 透明性 <ul style="list-style-type: none"> <li>✓ AIを活用する際の社会的文脈を踏まえ、AIシステム・サービスの検証可能性を確保しながら、必要かつ技術的に可能な範囲で、<b>ステークホルダーに対し合理的な範囲で適切な情報を提供</b>する（AIを利用しているという事実、活用している範囲、データ収集及びアノテーションの手法、AIシステム・サービスの能力、限界、提供先における適切/不適切な利用方法、等）</li> </ul>
	7) アカウンタビリティ <ul style="list-style-type: none"> <li>✓ トレーサビリティの確保や共通の指針の対応状況等について、ステークホルダーに対して情報の提供と説明を行う</li> <li>✓ 各主体の<b>AIガバナンスに関するポリシー、プライバシーポリシー等の方針を策定</b>し、公表する</li> <li>✓ 関係する情報を文書化して一定期間保管し、必要なときに、必要なところで、入手可能かつ利用に適した形で参照可能な状態とする</li> </ul>
社会と 連携した 取組が 期待される 事項	8) 教育・リテラシー <ul style="list-style-type: none"> <li>✓ AIに関わる者が、その関わりにおいて<b>十分なレベルのAIリテラシーを確保</b>するために必要な措置を講じる</li> <li>✓ AIの複雑性、誤情報といった特性及び意図的な悪用の可能性もあることを勘案して、<b>ステークホルダーに対しても教育を行う</b>ことが期待される。</li> </ul>
	9) 公正競争確保 <ul style="list-style-type: none"> <li>✓ AIを活用した新たなビジネス・サービスが創出され、持続的な経済成長の維持及び社会課題の解決策の提示がなされるよう、<b>AIをめぐる公正な競争環境が維持</b>に努めることが期待される</li> </ul>
	10) イノベーション <ul style="list-style-type: none"> <li>✓ 国際化・多様化、<b>産学官連携</b>及びオープンイノベーションを推進する</li> <li>✓ 自らのAIシステム・サービスと他のAIシステム・サービスとの相互接続性及び相互運用性を確保する</li> <li>✓ 標準仕様がある場合には、それに準拠する</li> </ul>

## 【参考】高度なAIシステムに関する事業者に通の指針

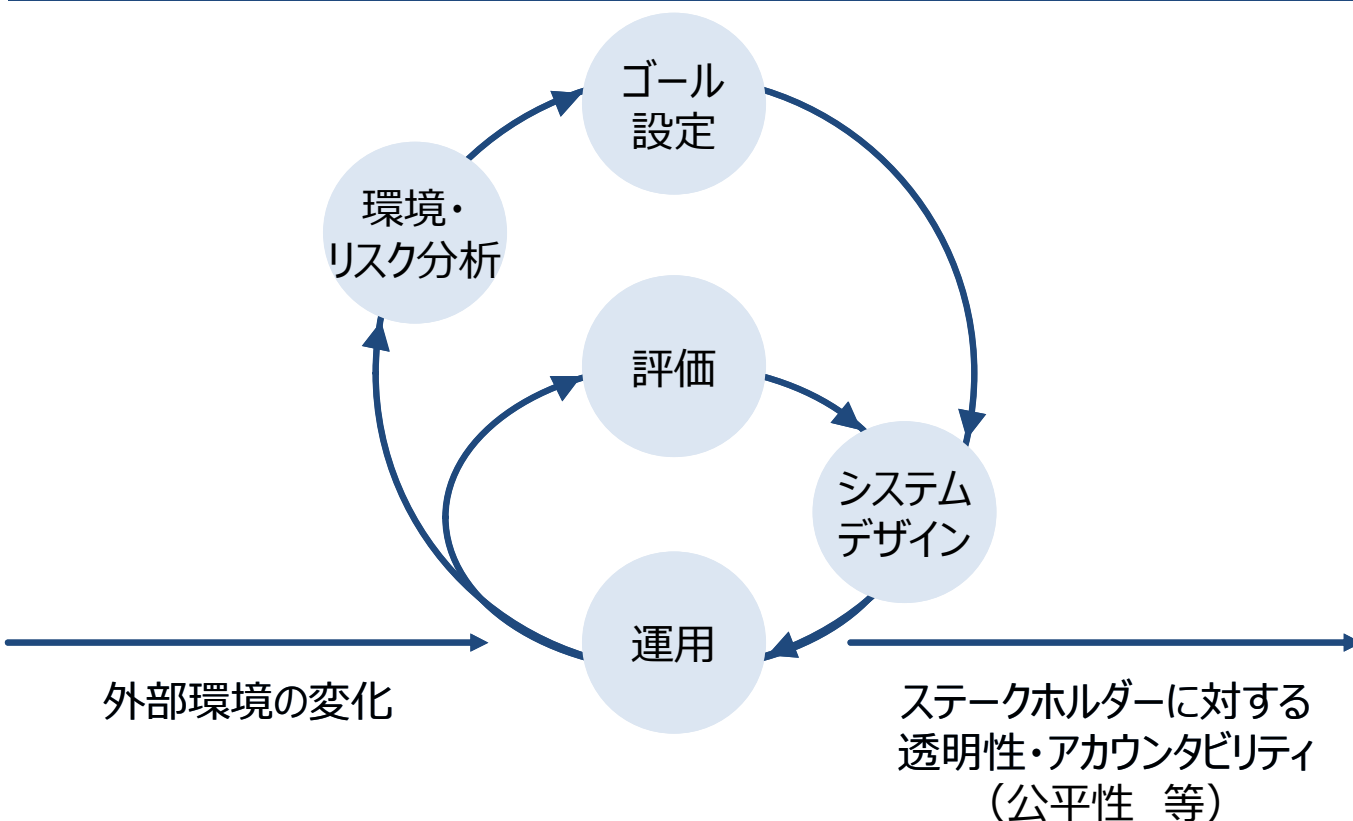
- 「共通の指針」に加え、以下を遵守すべきである\*1。ただし、I) ～ XI) は高度なAIシステムを開発するAI開発者にのみ適用される内容もあるため、各主体は適切な範囲で遵守することが求められる。
  - I. AIライフサイクル全体にわたるリスクを特定、評価、軽減するために、高度なAIシステムの開発全体を通じて、その導入前及び市場投入前も含め、適切な措置を講じる
  - II. 市場投入を含む導入後、脆弱性、及び必要に応じて悪用されたインシデントやパターンを特定し、緩和する
  - III. 高度なAIシステムの能力、限界、適切・不適切な使用領域を公表し、十分な透明性の確保を支援することで、アカウントビリティの向上に貢献する
  - IV. 産業界、政府、市民社会、学界を含む、高度なAIシステムを開発する組織間での責任ある情報共有とインシデントの報告に向けて取り組む
  - V. 特に高度なAIシステム開発者に向けた、個人情報保護方針及び緩和策を含む、リスクベースのアプローチに基づくAIガバナンス及びリスク管理方針を策定し、実施し、開示する
  - VI. AIのライフサイクル全体にわたり、物理的セキュリティ、サイバーセキュリティ、内部脅威に対する安全対策を含む、強固なセキュリティ管理に投資し、実施する
  - VII. 技術的に可能な場合は、電子透かしやその他の技術等、AI利用者及び業務外利用者等が、AIが生成したコンテンツを識別できるようにするための、信頼できるコンテンツ認証及び来歴のメカニズムを開発し、導入する
  - VIII. 社会的、安全、セキュリティ上のリスクを軽減するための研究を優先し、効果的な軽減策への投資を優先する
  - IX. 世界の最大の課題、特に気候危機、世界保健、教育等（ただしこれらに限定されない）に対処する対処するため、高度なAIシステムの開発を優先する
  - X. 国際的な技術規格の開発を推進し、適切な場合にはその採用を推進する
  - XI. 適切なデータインプット対策を実施し、個人データ及び知的財産を保護する
  - XII. 高度な AI システムの信頼でき責任ある利用を促進し、貢献する

\*1: 詳細は、G7デジタル・技術大臣会合（2023年12月）で採択された「広島AIプロセスG7デジタル・技術閣僚声明」における「広島AIプロセス包括的政策枠組み」の「II. 全てのAI関係者向け及び高度なAIシステムを開発する組織向けの広島プロセス国際指針」を参照。

# AIガバナンスの構築

- AIを安全安心に活用していくために、経営層のリーダーシップのもと、下記に留意しながら適切なAIガバナンスを構築することで、リスクをマネジメントしていくことが重要となる
  - 複数主体に跨る論点について、バリューチェーン/リスクチェーンの観点で主体間の連携確保
  - 上記が複数国にわたる場合、データの自由な越境移転の確保のための適切なAIガバナンスの検討
  - 経営層のコミットメントによる、各組織の戦略や企業体制への落とし込み/文化としての浸透

## 適切なAIガバナンスの構築



### 複数主体間の連携確保

バリューチェーン/リスクチェーンの観点  
で主体間の連携を確保



### 適切なデータ流通

複数国にまたがる想定の場合の適切な  
リスク管理/AIガバナンスの実施



### 経営層のコミットメント

戦略/体制への落とし込み、  
各組織の文化としての浸透



# 【参考】別添2. 第2部 [E.AIガバナンスの構築]関連

- 別添2では、AIガバナンスの構築のための「行動目標」、「実践のポイント」及びそれに対応する仮想的な「実践例」、実際の企業の取組事例を掲載
- 各事業者にとって、自身のAIガバナンスの構築のための検討を、具体例を交えつつ行うことを可能とする

## A.経営層によるAIガバナンスの構築及びモニタリング

- 行動目標
  - 一般的かつ客観的な目標を記載
- 実践のポイント
  - 「AI原則実践のためのガバナンスガイドライン Ver. 1.1」を基に、国内外ガイドラインやISO等の要素も取り入れ
- 実践例
  - 仮想事例に基づく事例を記載
  - 生成AI等の最新動向への対応事例も織り込み

例) 行動目標1-1【便益/リスクの理解】【実践例v:生成AIへの対応】  
 「自社の業務に活用するにあたり、JDLA【※】の発行する「生成AIの利用ガイドライン」に沿って、自社内での利用ガイドラインの策定を行っている。併せて、個人情報保護委員会の「生成AIサービスの利用に関する注意喚起等」等の行政からの情報発信も確認している。また、ニュース及びSNS等を通じ生成AIに関する情報の収集も欠かせない。これらを通じて、便益及びリスクの最新動向の把握に努めている。」【※一般社団法人日本ディープラーニング協会】

## B.AIガバナンスの構築に関する実際の取組事例

- 実際の取組事例
  - 「AI原則実践のためのガバナンスガイドライン Ver. 1.1」に基づくAIガバナンスの取組事例について、5社分をコラム化
  - AIガバナンスを推進するにあたり、多くの企業がつまづく観点を記載



# 【参考】別添3～5. 各主体向け 構成

- 各主体の重要事項に対して、「ポイント」、「具体的な手法」、「参考資料」を以下の構成で掲載

## 別添 記載内容

## 解説

### A. 本編「第3部 AI 開発者に関する事項」の解説

[本編の記載内容（再掲）]

データ前処理・学習時

#### D-2) i. 適切なデータの学習

Article I. プライバシー・バイ・デザイン等を通じて、学習時のデータについて、適正に収集するとともに、第三者の個人情報、知的財産権に留意が必要なもの等が含まれている場合には、法令に従って適切に扱うことを、AI のライフサイクル全体を通じて確保する（「2」安全性、「4」プライバシー保護、「5」セキュリティ確保）

Article II. 学習前・学習全体を通じて、データのアクセスを管理するデータ管理・制限機能の導入検討を行う等、適切な保護措置を実施する（「2」安全性、「5」セキュリティ確保）

#### [ポイント]

- AI モデルの質の向上のために、AI 開発者は、AI の学習等に用いるデータの質に留意することが重要となる。
- 利用する AI の特性及び用途を踏まえ、AI の学習等に用いるデータの質（正確性及び完全性等）に留意する
- また、AI によりなされる判断は、事後的に精度が損なわれたり、低下したりすることが想定されるため、想定される権利侵害の規模、権利侵害の生じる頻度、適用できる技術水準、精度を維持するためのコスト等を踏まえ、あらかじめ精度に関する基準を定めておくことが期待される。精度が当該基準を下回った場合には、データの質に留意して改めて学習させる
- なお、ここで言う「精度」には、AI が倫理的に正しい判断を行っているか（例えば、AI が暴力的な表現を行っていないか、ヘイトスピーチ等を行っていないか 等）も含まれる

#### [具体的な手法]

1. データに個人情報、機密情報、著作権等の権利又は法律上保護される利益に関するものが含まれていないか、確認を実施
  - 1.1. 固有表現抽出

#### [参考文献]

1. 国立研究開発法人産業技術総合研究所「機械学習品質マネジメントガイドライン 第4版」（2023年12月）
2. NIST, “AI Risk Management Framework Playbook”（2023年1月）

本編の記載内容を再掲

### ポイント

- 本編記載事項に加え、重要観点を補足

### 具体的な手法

- 他のガイドライン等を参照しつつ具体的に解説

### 参考文献

- 「ポイント」や「具体的な手法」の参照元を記載

本編の  
第3～5部に  
対応する解説

本編 第2部に  
対応する解説

「共通の指針」に  
関する各主体の  
具体的なアプローチ



- AI開発者は、AIモデルを直接的に設計・変更ができるため、AIが提供/利用された際にどのような影響を与えるか、事前に可能な限り検討し、対応策を講じておくことが特に重要

データ前処理  
学習時

- D-2) i. 適切なデータの学習
- プライバシー・バイ・デザイン等を通じて、個人情報、知的財産権に留意が必要なもの等が含まれている場合には、法令に則って適切に扱う
  - データ管理・制限機能の導入検討を行う等、**適切な保護措置を実施**する
- D-3) i. データに含まれるバイアス等への配慮
- 学習データ、モデルの学習過程でバイアスが含まれうることに留意し、**データの質を管理するための相当の措置**を講じる
  - バイアスを完全に排除できないことを踏まえ、**AIモデルが代表的なデータセットで学習され、AIシステムに不公正なバイアスがないか点検**されることを確保する

## AI開発時

- D-2) ii. 人間の生命・身体・財産、精神及び環境に配慮した開発
- 予期しない環境を含む様々な状況下での利用に耐えうる性能の要求
  - **リスクを最小限に抑える**方法の要求
- D-2) iii. 適正利用に資する開発
- **AIを安全に利用可能な使い方について明確な方針・ガイダンスを設定**する
  - AIモデルに対する事後学習を行う場合に、**学習済AIモデルを適切に選択**する
- D-3) ii. AIモデルのアルゴリズム等に含まれるバイアスへの配慮
- AIモデルを構成する**各技術要素によってバイアスが含まれうる**ことまで検討する
  - AIモデルが代表的なデータセットで学習され、AIシステムに不公正なバイアスがないか点検する
- D-5) i. セキュリティ対策のための仕組みの導入
- 採用する技術の特性に照らし適切に**セキュリティ対策を講ずる**（セキュリティ・バイ・デザイン）
- D-6) i. 検証可能性の確保
- AIの予測性能及び出力の品質が、活用開始後に大きく変動する可能性又は想定する精度に達しないこともある特性を踏まえ、**事後検証のための作業記録を保存**しつつ、その品質の維持・向上を行う

- AI開発者は、AIモデルを直接的に設計・変更ができるため、AIが提供/利用された際にどのような影響を与えるか、事前に可能な限り検討し、対応策を講じておくことが特に重要

## 開発後

- |                                     |   |
|-------------------------------------|---|
| D-5) ii. 最新動向への留意                   | - AIシステムに対する攻撃手法は日々新たなものが生まれており、これらのリスクに対応するため、 <b>開発の各工程で留意すべき点を確認</b> する  |
| D-6) ii. 関連する<br>ステークホルダーへの<br>情報提供 | - AIシステムの技術的特性、安全性確保の仕組み、予見可能なリスク及びその緩和策、不具合の原因及び対応状況等に関する <b>情報提供</b> を行う  |
| D-7) i. AI提供者への共通の<br>指針の対応状況の説明    | - AI提供者に対して、AIに活用開始後に品質が変動する可能性及び、その結果として <b>生じるリスク等の情報提供及び説明</b> を行う   |
| D-7) ii. 開発関連情報の<br>文書化             | - AIシステムの開発過程、意思決定に影響を与えるデータ収集及びラベリング、使用されたアルゴリズム等について <b>文書化</b> する  |
| D-10) i. イノベーションの<br>機会創造への貢献       | - AIの <b>品質・信頼性、開発の方法論等の研究開発</b> を行う<br>- <b>持続的な経済成長の維持及び社会課題解決</b> につながるよう貢献する<br>- DFFT等の国際議論の動向の参照、AI開発者コミュニティ又は学会への参加等の取組を行う等、国際化・多様化及び産学官連携を行う<br>- <b>社会全体への情報提供</b> を行う |

- AI提供者は、AIの稼働と適正な利用を前提としたAIシステム・サービスの提供を実現することが重要

AIシステム  
実装時

- |          |  |  |
|----------|--|--|
| P-2) i.  | 人間の生命・身体・<br>財産、精神及び環境に<br>配慮したリスク対策     | - 様々な状況下でAIシステムがパフォーマンスレベルを維持できるようにし、 <b>リスクを最小限に抑える</b> 方法を検討する   |
| P-2) ii. | 適正利用に資する提供                               | - AI開発者が設定した範囲でAIを活用する<br>- AIシステム・サービスの正確性等を担保すると同時に、 <b>AI開発者の想定利用環境とAI利用者の利用環境に違い等がないか検討</b> する   |
| P-3) i.  | AIシステム・サービスの<br>構成及びデータに含まれ<br>るバイアスへの配慮 | - データの公平性を担保し、参照する情報、外部サービス等の <b>バイアスを検討</b> する<br>- AIモデルの入出力及び <b>判断根拠を定期的に評価</b> し、バイアスの発生をモニタリングする<br>- AIモデルの出力結果を受け取るAIシステム等において、利用者の判断を恣意的に制限するようなバイアスが含まれる可能性を検討する |
| P-4) i.  | プライバシー保護のため<br>の仕組み及び対策の導<br>入           | - 採用する技術の特性に照らし適切に個人情報へのアクセスを管理・制限する仕組みの導入等の <b>プライバシー保護対策を講ずる</b> （プライバシー・バイ・デザイン）  |
| P-5) i.  | セキュリティ対策のため<br>の仕組みの導入                   | - 採用する技術の特性に照らし適切に <b>セキュリティ対策を講ずる</b> （セキュリティ・バイ・デザイン）  |
| P-6) i.  | システムアーキテクチャ等<br>の文書化                     | - AIシステムの意思決定に影響を与えるシステムアーキテクチャ、データの処理プロセス等について <b>文書化</b> する  |

- AI提供者は、AIの稼働と適正な利用を前提としたAIシステム・サービスの提供を実現することが重要

AIシステム・  
サービス  
提供後

- |          |                          |   |   |
|----------|--------------------------|---|---|
| P-2) ii. | 適正利用に資する提供               | - | <b>適切な目的</b> でAIシステム・サービスが利用されているかを定期的に検証する   |
| P-4) ii. | プライバシー侵害への<br>対策         | - | AIシステム・サービスにおけるプライバシー侵害に関して <b>適宜情報収集し、侵害を認識した場合等は適切に対処するとともに、再発の防止</b> を検討する   |
| P-5) ii. | 脆弱性への対応                  | - | 最新のリスクに対応するために提供の各工程で気を付けるべき点の動向を確認し、 <b>脆弱性に対応することを検討する</b>  |
| P-6) ii. | 関連するステークホルダー<br>への情報提供   | - | AIシステムの・サービスの技術的特性、予見可能なリスク、緩和策、出力又はプログラムの変化の可能性、不具合の原因と対応状況、インシデント事例、学習データの収集ポリシー、その学習方法等に関する情報を説明できるようにする<br>- AIの性質及び利用目的等に照らして、 <b>AIを利用しているという事実や適切/不適切な使用方法、更新内容とその理由等の情報提供や説明の実施</b> |
| P-7) i.  | AI利用者への共通の<br>指針の対応状況の説明 | - | AI利用者に <b>適正利用を促し</b> 、正確性・必要に応じて最新性等が担保されたデータの利用やコンテキスト内学習による不適切なモデルの学習に対する注意喚起、 <b>個人情報を入力する際の留意点についての情報を提供する</b><br>- AIシステム・サービスへの個人情報の不適切入力について注意喚起する                                  |
| P-7) ii. | サービス規約等の<br>文書化          | - | AI利用者に向けた <b>サービス規約を作成するとともにプライバシーポリシーを明示する</b>   |

- AI利用者は、AI提供者が意図した範囲内で継続的に適正利用、必要に応じたAIシステムの運用を行うことが重要であり、より効果的なAI利用のために必要な知見を習得することが期待される

AIシステム  
サービス  
利用時

- |   |  |
|---|--|
| U-2) i. 安全を考慮した<br>適正利用                   | <ul style="list-style-type: none"> <li>- AI提供者が定めた利用上の留意点を遵守して、<b>AI提供者が設計において想定した範囲内で利用</b>する</li> <li>- AIの出力について精度及びリスクの程度を理解し、<b>様々なリスク要因を確認した上で利用</b>する</li> </ul>  |
| U-3) i. 入力データ又はプロンプ<br>トに含まれるバイアスへの<br>配慮 | <ul style="list-style-type: none"> <li>- 公平性が担保されたデータの入力を行い、プロンプトに含まれるバイアスに留意して、責任をもって<b>AI出力結果の事業利用判断を行う</b></li> </ul>   |
| U-4) i. 個人情報の不適切<br>入力及びプライバシー侵<br>害への対策  | <ul style="list-style-type: none"> <li>- AIシステム・サービスへ個人情報を不適切に入力しないよう注意を払う</li> <li>- AIシステム・サービスにおける<b>プライバシー侵害に関して適宜情報収集</b>し、防止を検討する</li> </ul>   |
| U-5) i. セキュリティ対策の実施                       | <ul style="list-style-type: none"> <li>- AI提供者による<b>セキュリティ上の留意点を遵守</b>する</li> <li>- AIシステム・サービスに機密情報等を不適切に入力しないよう注意を払う</li> </ul>  |
| U-6) i. 関連するステー<br>クホルダーへの情報提供            | <ul style="list-style-type: none"> <li>- 公平性が担保されたデータの入力を行い、プロンプトに含まれるバイアスに留意して、<b>出力結果を取得し、結果を事業判断に活用した際は、その結果を関連するステークホルダーに合理的な範囲で情報提供</b>する</li> </ul>  |
| U-7) i. 関連するステー<br>クホルダーへの説明              | <ul style="list-style-type: none"> <li>- AIの特性や用途、データの提供元となる関連するステークホルダーとの接点、プライバシーポリシー等を踏まえ、データ提供の手段、形式等について、あらかじめ<b>当該ステークホルダーに平易かつアクセスしやすい方法で情報提供</b>する</li> <li>- AIの出力結果を特定の個人又は集団に対する評価の参考とする場合は、人間による合理的な判断のもと、説明責任を果たす</li> <li>- <b>関連するステークホルダーからの問合せに対応する窓口を合理的な範囲で設置</b>し、AI提供者とも連携の上説明及び要望の受付を行う</li> </ul> |
| U-7) ii. 提供された文書の活用<br>と規約の遵守             | <ul style="list-style-type: none"> <li>- AI提供者から提供されたAIシステム・サービスについての<b>文書を保管・活用</b>する</li> <li>- AI提供者が定めた<b>サービス規約を遵守</b>する</li> </ul>   |

# 別添7. チェックリスト 概要

- 別添7では、AIによるリスクを抑えつつ便益を享受する取組の立案、実践を確実に推進するための「チェックリスト」及び「具体的なアプローチのためのワークシート」を用意している

本編・  
別添1～5（付属資料）

チェックリスト

具体的なアプローチ検討の  
ためのワークシート



本編・別添を読んでAIガバナンスの重要性や、各事業者に期待されることを理解する

「チェックリスト」を活用し、本編・別添についての各主体の取組(What)を確認する

「具体的なアプローチ検討のためのワークシート」を使用し、各事業者の具体的なアプローチ(How)を検討する

# 別添7. チェックリスト 活用方法

## 別添7

- 全事業者は、「別添7 Aチェックリスト[全主体向け]」を活用し、各自の取組状況の概観を確認する
- 高度なAIシステムに関係する事業者に該当する場合には、「別添7 Bチェックリスト[高度なAIシステムに関係する事業者向け]」も実施する

**別添7 A チェックリスト [全主体向け]** 令和6年4月19日

本チェックリストは、全事業者が行うべき取組と、共通の留意事項を記載したものです。事業者ごとの具体的な取組事項は別添7 Bをご活用ください。

高度なAIシステムに関係する事業者は、「別添7 B 高度なAIシステムに関係する事業者向け」も実施ください。

**別添7 B チェックリスト [高度なAIシステムに関係する事業者向け]** 令和6年4月19日

本チェックリストは、全事業者が行うべき取組と、高度なAIシステムに関係する事業者に共通の留意事項を記載したものです。

高度なAIシステムに関係する事業者は、本編と別添7 Bに示された取組事項の両方について実施する必要があります。本編と別添7 Bの両方ともご活用ください。

チェック項目	チェック項目
<input type="checkbox"/> 人間中心の考え方を基に、画法が保障される又は法的的に認められた人権を侵害することがないようしているか？	<input type="checkbox"/> 高度なAIシステムの市場導入前及び開発全体を通じて、AIのバイアスや差別に十分注意を払う。評価、検証するにあたって適切な評価を実施しているか？
<input type="checkbox"/> AIに関わる全ての者の生命・身体・財産、精神及び環境に危害を及ぼすことがないよう安全性を確保しているか？	<input type="checkbox"/> 市町村の人権・倫理性、インクルーシブ、異文化理解を促進し、実践しているか？
<input type="checkbox"/> 潜在的なバイアスをなくすよう留意し、それと対応できないバイアスがあることを認識しつつ、回避できないバイアスが人権及び多様な文化を尊重する取組の観点から許容可能か評価しているか？	<input type="checkbox"/> 十分な透明性の確保やアウトプットの向上のため高度なAIシステムの能力、限界、通知・干渉的な影響等を留意しているか？
<input type="checkbox"/> プライバシーを尊重・保護し、関係法令を遵守しているか？	<input type="checkbox"/> 国・自治体、民間、市民社会、学術界との関係構築等、責任ある情報共有に十分コミットしているか？
<input type="checkbox"/> 不正操作によってAIの振る舞いが意図せぬ変更又は停止が生じることのないよう、セキュリティを確保しているか？	<input type="checkbox"/> リスクベースのアプローチにもとづいてAIのリスク及びリスク管理方針を策定、実施、評価しているか？
<input type="checkbox"/> 透明性を確保するため、AI自体やAIシステム・サービスが提供されるユーザーに対し合理的で技術的に可能な範囲で提供しているか？	<input type="checkbox"/> AIのバイアスや差別を回避するため、管理プロセスを強化し、バイアスや差別を回避する、検閲プロセスや監査を実施し、実施しているか？
<input type="checkbox"/> データの出し所、AIの意思決定等のプロセスやAIに関する情報やリスクの対応状況等について、関係するステークホルーに対し合理的な範囲でアウトプットやフィードバックを実施しているか？	<input type="checkbox"/> 法的に可能な場合は、AIが生成したコンテンツを識別するために、電子透視しその他の技術等、信頼性の高いコンテンツ検証及び実装のメカニズムを構築、導入しているか？
<input type="checkbox"/> AIのバランスやプライバシーに関するポリシー等を策定しているか？	<input type="checkbox"/> 社会的、倫理的、セキュリティ上のリスクを軽減するための研究を推進し、効果的な実践に導入に活用しているか？
<input type="checkbox"/> 上記の実現のため、各事業者の状況に応じた具体的なアプローチは構築しているか？	<input type="checkbox"/> 倫理規範、検証・評価などが世界の標準や国際標準等と一致し、高度なAIシステムの開発を促進しているか？
	<input type="checkbox"/> 国際的な倫理規範の取組を推進しているか？
	<input type="checkbox"/> 適切なデータ入力確保と、偏りデータ及び怪しいデータの取組を実施しているか？
	<input type="checkbox"/> 国・自治体の取組等との連携や協力に関するプロセスやアウトプットの向上や透明性の確保への協力と情報共有等、高度なAIシステムの提供と責任ある取組の促進、実施しているか？
	<input type="checkbox"/> 上記の実現のため、各事業者の状況に応じた具体的なアプローチは構築しているか？

• チェック項目は本編の要約を記載

• チェックすることで、各自の取組状況を概観

• 具体的な実践内容の検討に、「別添7C 具体的なアプローチのためのワークシート」を活用 (活用方法次頁)

※別添7 Bには「具体的なアプローチ検討のためのワークシート」もご活用ください

# 別添7. 具体的なアプローチ検討のためのワークシート

## 別添7

### 活用方法

- 本ガイドラインの記載内容に関して、具体的なアプローチを検討する際に重要となる事項を記載
- 事業者の事業内容及び置かれた状況等に応じ、各自でカスタマイズして活用することを前提としている

別添7C. 具体的なアプローチ検討のためのワークシート (共通の指針関連)

別添7C.1. 具体的なアプローチ検討のためのワークシート (共通の指針関連)

別添7C.2. 具体的なアプローチ検討のためのワークシート (共通の指針関連)

事業者の事業内容及び置かれた状況等	具体的なアプローチ	実施の順序	実施の時期	実施の場所	実施の方法	実施の担当者	実施の進捗状況
1. 事業者の事業内容及び置かれた状況等	2. 具体的なアプローチ	3. 実施の順序	4. 実施の時期	5. 実施の場所	6. 実施の方法	7. 実施の担当者	8. 実施の進捗状況

• 活用の前に、**実施責任者等を決定**する

• 事業者の事業内容や置かれた状況等に応じ、**各自でカスタマイズして活用**する

• 「見直し日」も検討・記載することで**定期的な更新**を行う

• 具体の**アプローチ**を検討する際に重要になる事項が記載されており、事業者が**各自でカスタマイズ**する際の**リファレンス**となる



# 【参考】別添8. 主体横断的な仮想事例

## 概要

- 別添8.では、本ガイドラインに沿って、AI開発者、AI提供者、AI利用者が重要事項の検討を行った場合の「主体横断的な仮想事例」を掲載
- 各主体が本ガイドラインの内容を実際に落とし込む際の具体的なイメージの想起や、各主体間での連携が重要になるポイントの明確化が可能
- 現在「採用AI」を扱う事業者を例として取り上げているが、今後、他の事例も追加していく予定

Case 採用AI		AI開発部門	人材採用部門（採用AIチーム）	人材採用担当者
		AI開発者	AI提供者	AI利用者
No.	分類	共通の前提/各主体が実施している活動	本ICにおける主体が実施している活動	本ICにおいて主体が実施している活動
<b>1) 人間中心</b>				
本主体は、AIシステム・サービスの開発・提供・利用において、後述する各事項を含む全ての取り扱われる事項が満たされる土壌として、少なくとも実法が定める又は慣習的に認められる人権を尊重することがないよう努めてある。また、AIが人々の能力を拡張し、手塚				
<b>① 人間の尊厳及び個人の自由</b>				
1	目的	AIが活用される際の社会的文脈を踏まえ、人間の尊厳及び個人の自由を尊重する	AIシステムの開発において、学習データの収集やラベリング、モデルの性能評価等は、AI開発者だけで完結せず、AI提供者側で	AIサービスの提供において、AI利用者が最終判断(応募者の合否)を行えるようになっている(Human-in-the-loop)
2	目的	特に、AIを人間の脳・身体と連携させる場合には、その周辺技術に関する情報を踏まえつつ、諸外国及び研究施設における生命	脳・身体と連携するケースではないため対象外	脳・身体と連携するケースではないため対象外
3	目的	個人の権利・利益に重要な影響を及ぼす可能性のある分野においてAIを利用したプロファイリングを行う場合、個人の尊厳を尊重し、アウトプットの正確性を可能な限り確保を努め、AIの予測、推奨、判断等の結果を理解して利用し、かつ生じうる不利益等を慎重に検討した上で、不適切な目的に利用しない	AIシステムの開発において、実際の予測結果を学習データに用いる際には個人情報取扱いに関する契約書の締結やアクセス権管理等を実施している。 ※公平性とプライバシーについては、「3) 公平性」/「4) プライバシー保護」を参照	AIシステムの開発において、実際の予測結果を学習データに用いる際には個人情報取扱いに関する契約書の締結やアクセス権管理等を実施している。 ※公平性とプライバシーについては、「3) 公平性」/「4) プライバシー保護」を参照
<b>② AIによる意思決定・感情の操作等への留意</b>				
1	目的	人間の意思決定、認知等、感情を不当に操作することを目的とした、又は意図的に知覚できないレベルでの操作を前提としたAIシステム・サービスの開発・提供・利用は行わない	本ケースに関しては、2)①-3と同じ観点になる	本ケースに関しては、2)①-3と同じ観点になる
2	目的	AIシステムの開発・提供・利用において、自動化/バイアス等のAIに留意して存在するリスクに注意を払い、必要な対策を講じる	本ケースに関しては、2)①-3と同じ観点になる	本ケースに関しては、2)①-3と同じ観点になる

**ご清聴ありがとうございました**