

# AI事業者ガイドライン案 (第1.0版) 概要

---

総務省  
経済産業省  
(令和6年4月)

# 「AI事業者ガイドライン」の背景・目的

- 生成AIに代表されるように、AI関連技術は日々発展をみせ、利用機会と可能性は拡大の一途をたどり、産業におけるイノベーション創出や社会課題の解決に向けても活用されている
- 我が国においては、Society 5.0の実現に向け、AIの高度な活用に対する期待が高まっている
- 我が国は、G7におけるAI開発原則に向けた提案を先駆けとし、G7・G20やOECD等の国際機関での議論をリードし、多くの貢献をしてきた
- このような背景を受け、我が国におけるAIガバナンスの統一的な指針を示すことで、イノベーションの促進及びライフサイクルにわたるリスクの緩和を両立する枠組みを関係者と連携しながら共創していくことを目指す

技術革新

Society 5.0の実現

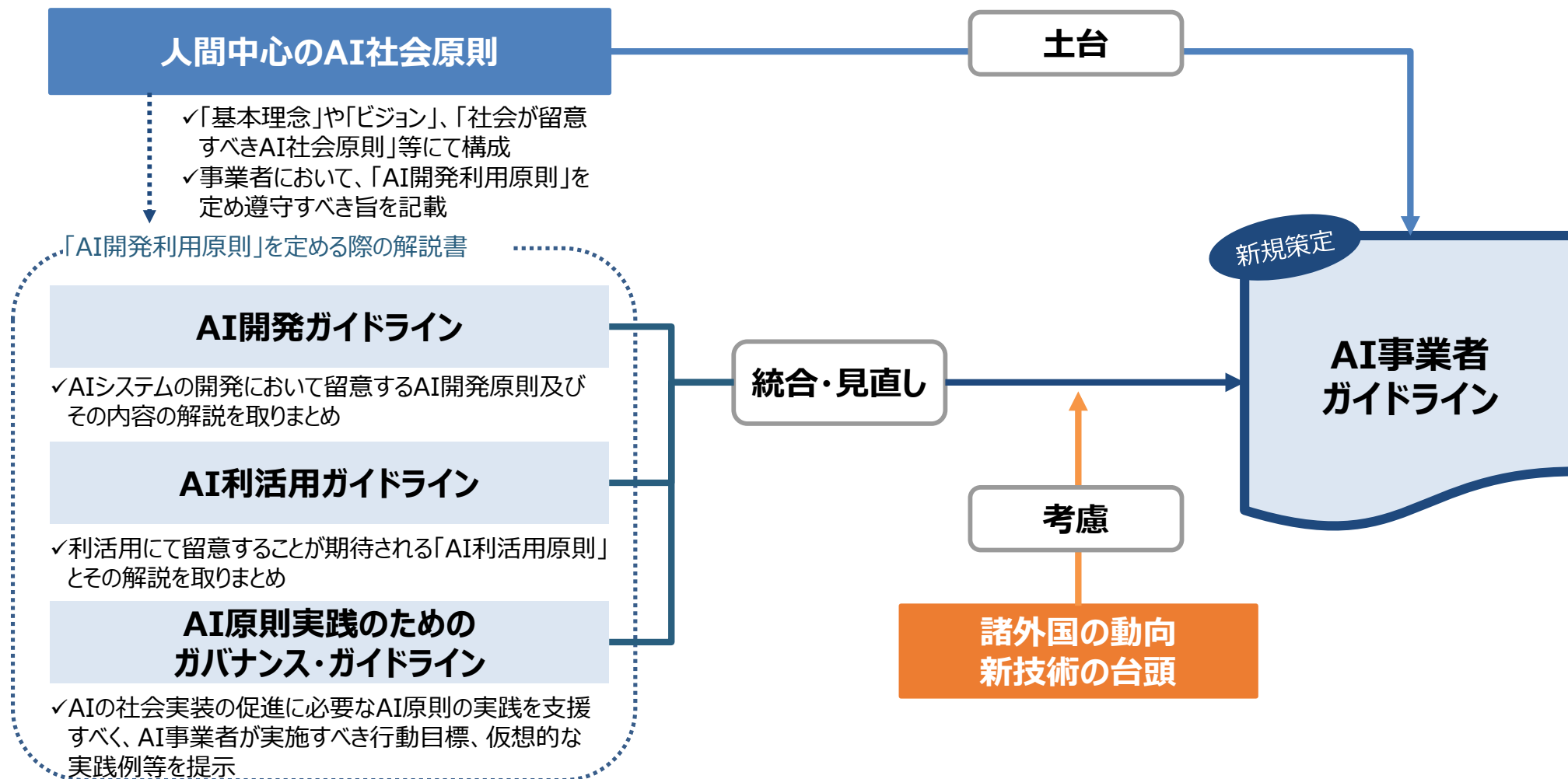
国際的な議論

## 「AI事業者ガイドライン」を策定

AIに関係する者が、国際的な動向及びステークホルダーの懸念を踏まえたAIのリスクを正しく認識し、必要となる対策をライフサイクル全体で自主的に実行できるように後押しし、イノベーションの促進及びライフサイクルにわたるリスクの緩和を両立する枠組みを関係者と連携しながら積極的に共創していくことを目指す

# 「AI事業者ガイドライン」の策定方針

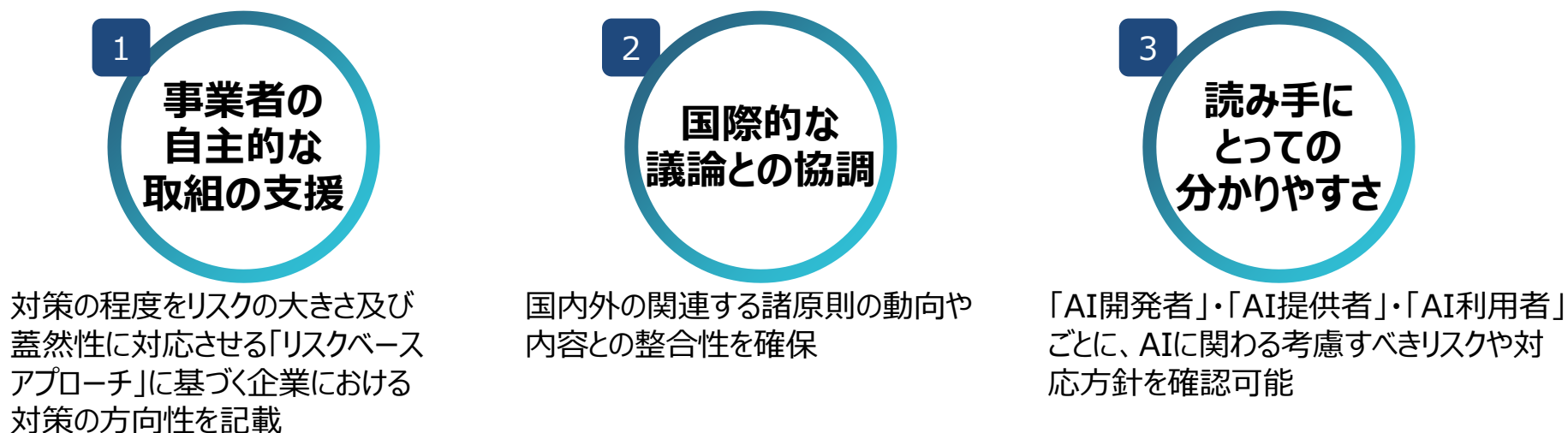
- 「AI事業者ガイドライン」は、「人間中心のAI社会原則」を土台としつつ、我が国における3つのガイドラインを統合し、諸外国の動向や新技術の台頭を考慮して策定する
- これまでのガイドラインとの整合性を担保することで、事業活動を支えるAIガバナンスの仕組みとして、連続性がある発展を遂げていくことが期待される



# 「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

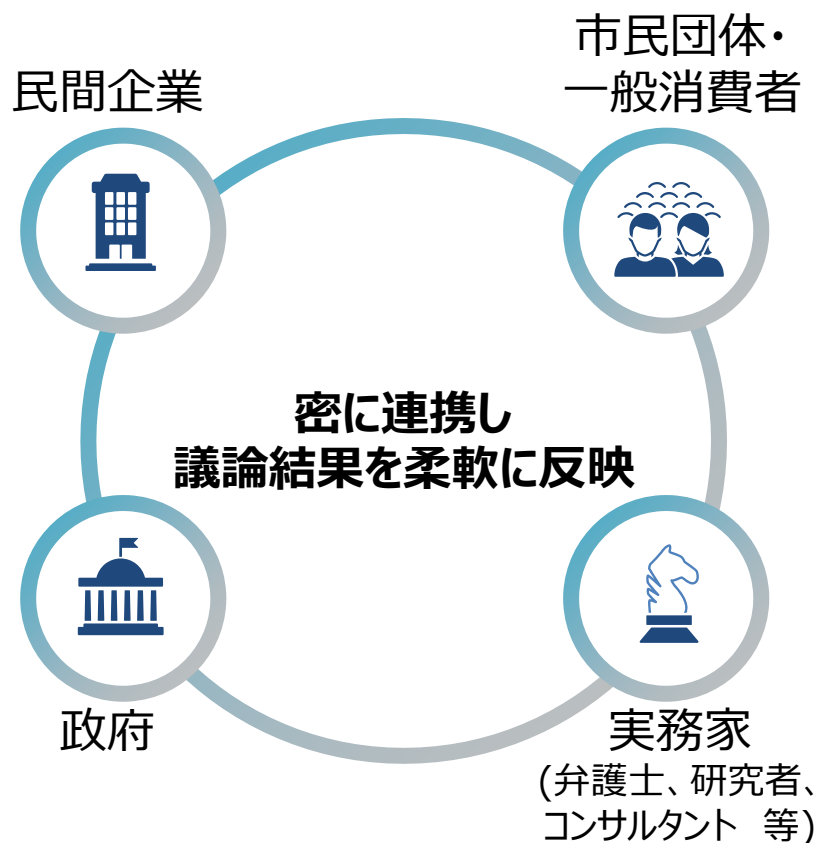
### Living Document

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

## 参考) マルチステークホルダーとの連携

- 政府単独ではなく、教育・研究機関、一般消費者を含む市民社会、民間企業等、多様なステークホルダー（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して策定する

### 連携主体



### 連携方法

#### 意見交換、議論の場を多数設定

- 左記連携主体で構成された検討会
- 実務家を中心としたワーキンググループ
- 民間企業との意見交換会

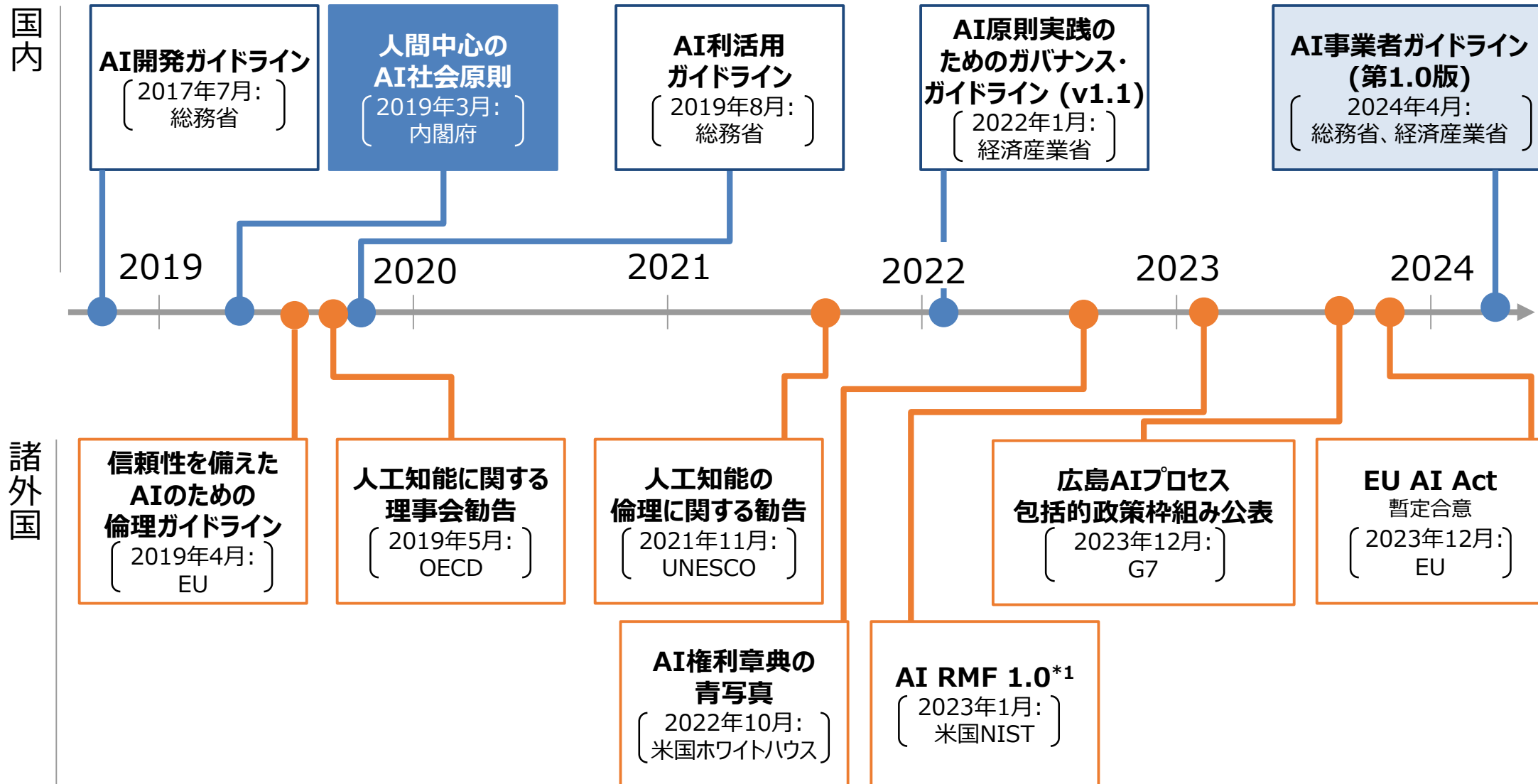
#### 意見照会を通じて広く知見を収集

- 100名程度の有識者
  - 民間企業担当者
  - 専門家、研究者
  - 市民団体、消費者団体 等

#### パブリックコメントを通じ、幅広い意見を収集

# 参考) AIに関連する主な諸原則等

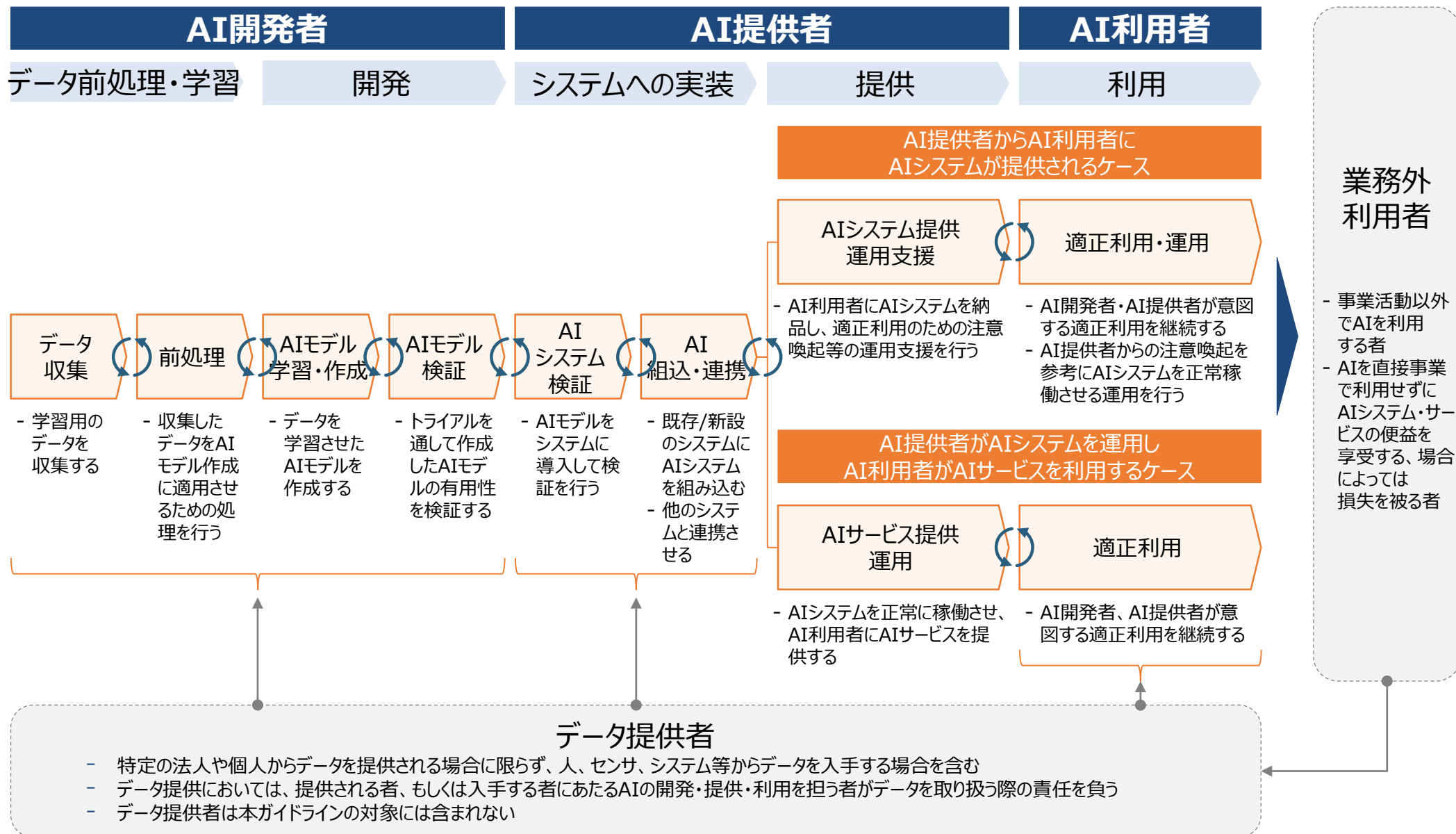
- 諸外国において、各種規制及びガイドラインの策定が積極的に議論されているため、本ガイドラインにおいても、諸原則や規制動向等との整合を意識する



\*1: AI Risk Management Framework 1.0

# 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)



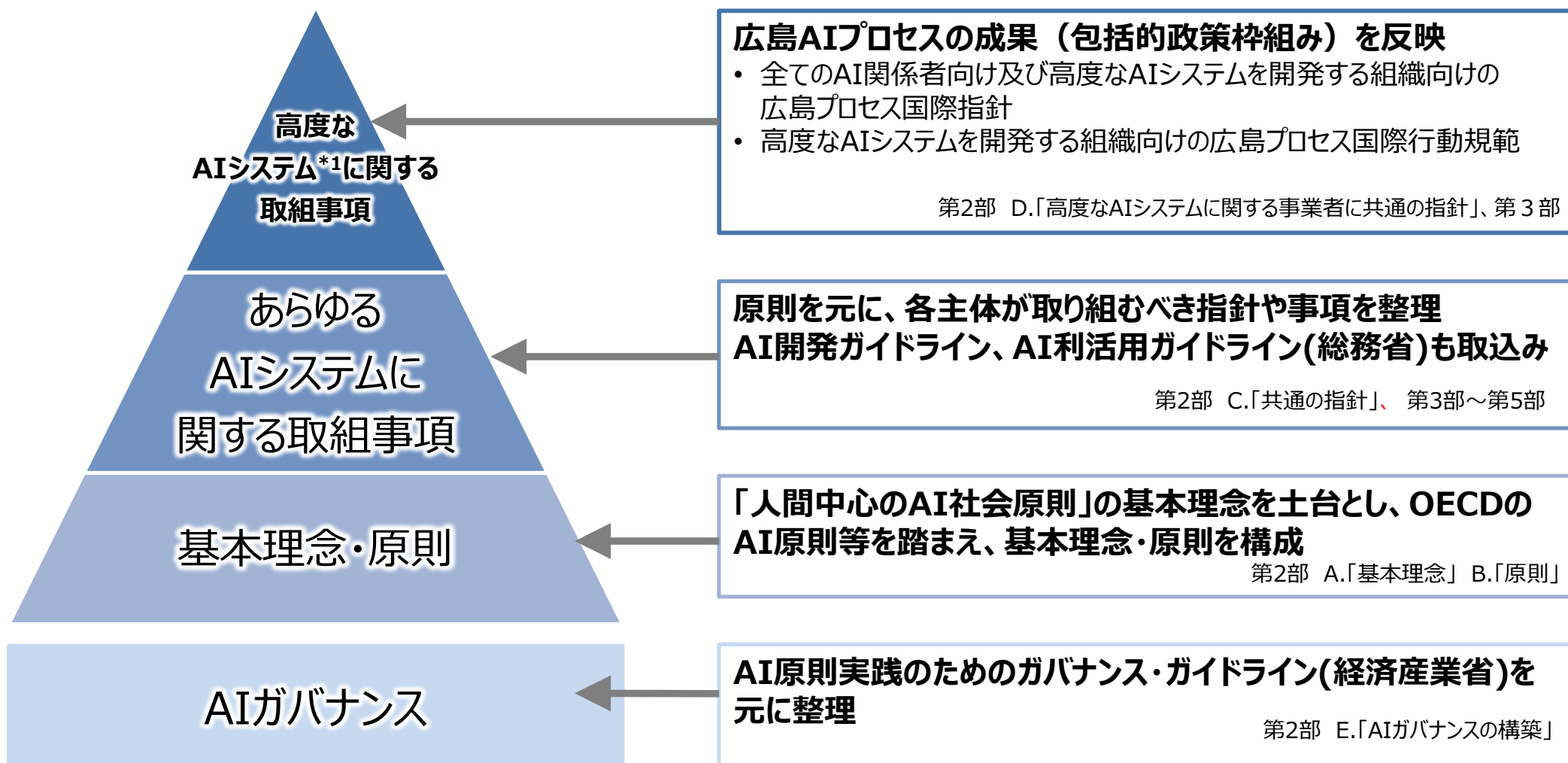
# 「AI事業者ガイドライン」の構成

- 別添の記載内容は本編と対応しており、本編の読解及びそれに基づく検討や行動をサポートする解説書としての役割を果たす

	本編 (why, what)	別添 (付属資料) (how)
主体共通	第1部 AIとは	1. 第1部関連 [AIについて] A. AIに関する前提 B. AIによる便益/リスク
	第2部 AIにより 目指すべき社会 及び 各主体が取り組む 事項 A.「基本理念」 B.「原則」 C.「共通の指針」 D.「高度なAIシステムに関する 事業者に通の指針」 E.「AIガバナンスの構築」	2. 第2部関連 [E.AIガバナンスの 構築] A. 経営層によるAIガバナンスの構築及び モニタリング B. AIガバナンスの事業者取組事例
主体別	第3部 AI開発者に 関する事項 ※「高度なAIシステムを開発する組織向けの 広島プロセス国際行動規範」における 追加的な記載事項 も含む	3. 第3部関連 [AI開発者向け] A. 「第3部 AI開発者に関する事項」の解説 B. 「第2部」の「共通の指針」の解説 C. 高度なAIシステムの開発にあたって遵守 すべき事項
	第4部 AI提供者に 関する事項	4. 第4部関連 [AI提供者向け] A. 「第4部 AI提供者に関する事項」の解説 B. 「第2部」の「共通の指針」の解説
	第5部 AI利用者に 関する事項	5. 第5部関連 [AI利用者向け] A. 「第5部 AI利用者に関する事項」の解説 B. 「第2部」の「共通の指針」の解説
その他 参考資料		6. 「AI・データの利用に関する契約ガイドライン」を参照 する際の主な留意事項について 7. チェックリスト 8. 主体横断的な仮想事例 9. 海外ガイドライン等の参照先

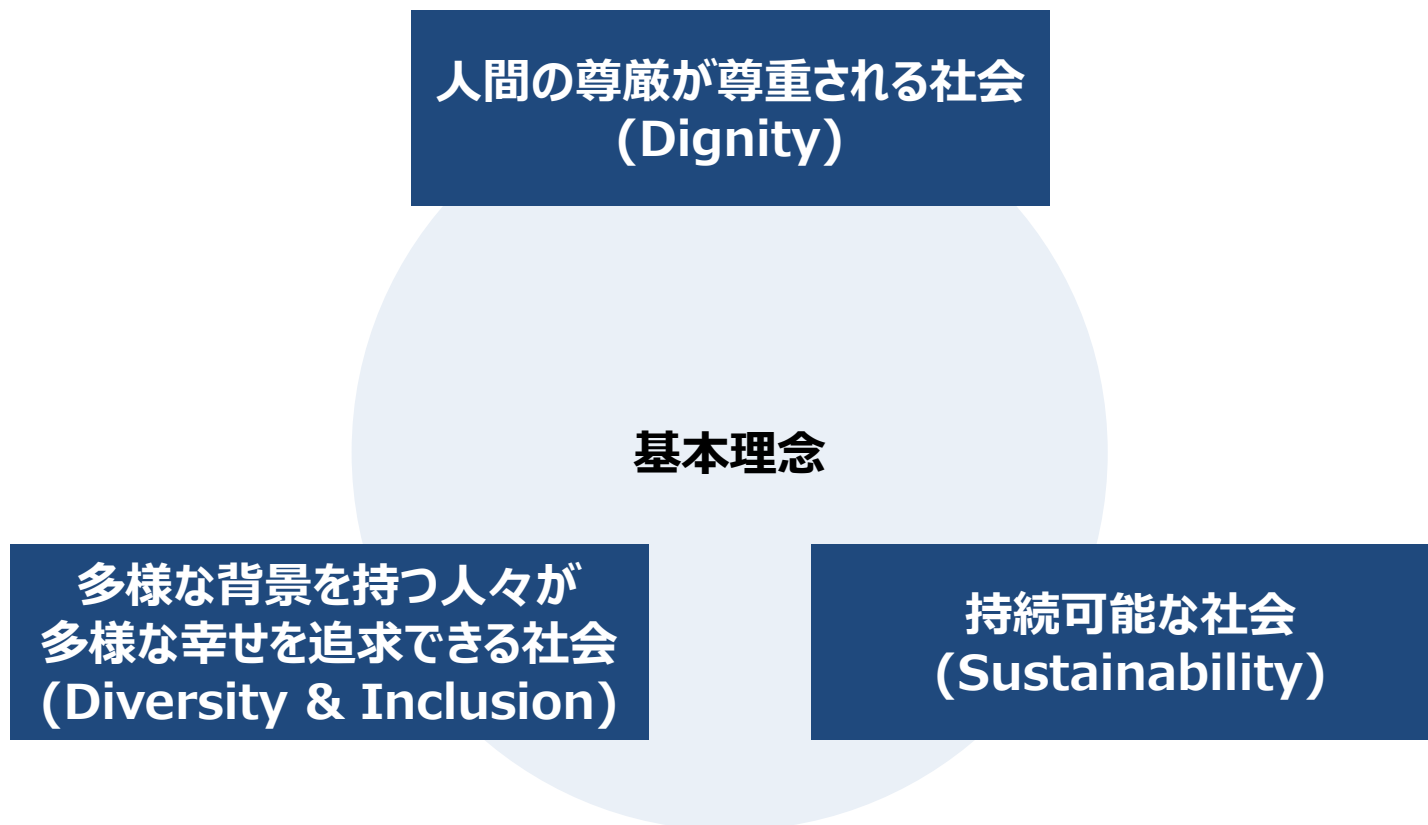
# 「AI事業者ガイドライン」の対象範囲

- ・ 広島AIプロセスで取りまとめられた高度なAIシステムに関する国際指針及び国際行動規範を反映しつつ、**一般的なAIを含む（想定され得る全ての）AIシステム・サービスを広範に対象**
- ・ 実際のAI開発・提供・利用においては、本ガイドラインを参照し、**各事業者が指針遵守のために適切なAIガバナンスを構築するなど、具体的な取組を自主的に推進することが重要**



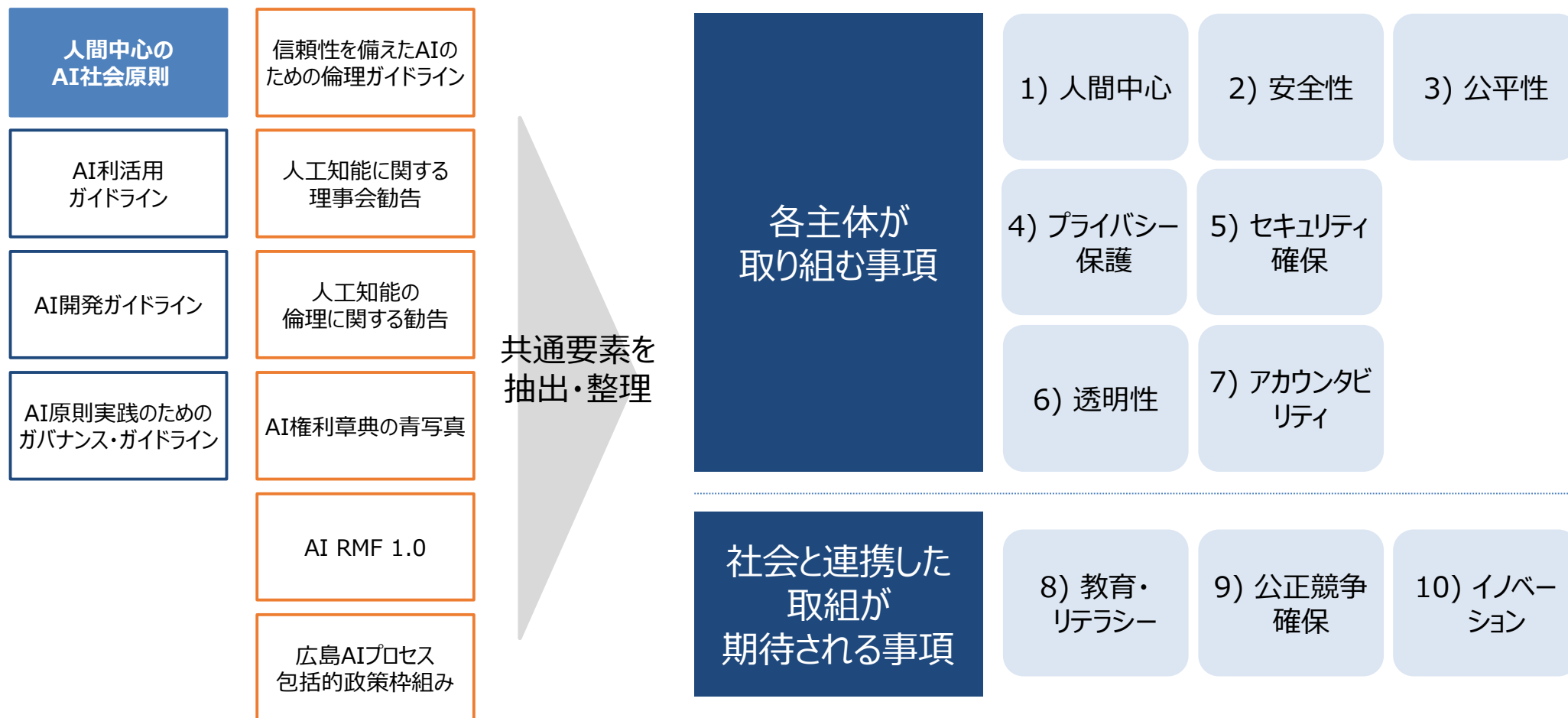
# 基本理念

- 「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



# 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>

# 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）

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、データ収集及びアナレーションの手法、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>



# 高度なAIシステムに関する事業者に通の指針

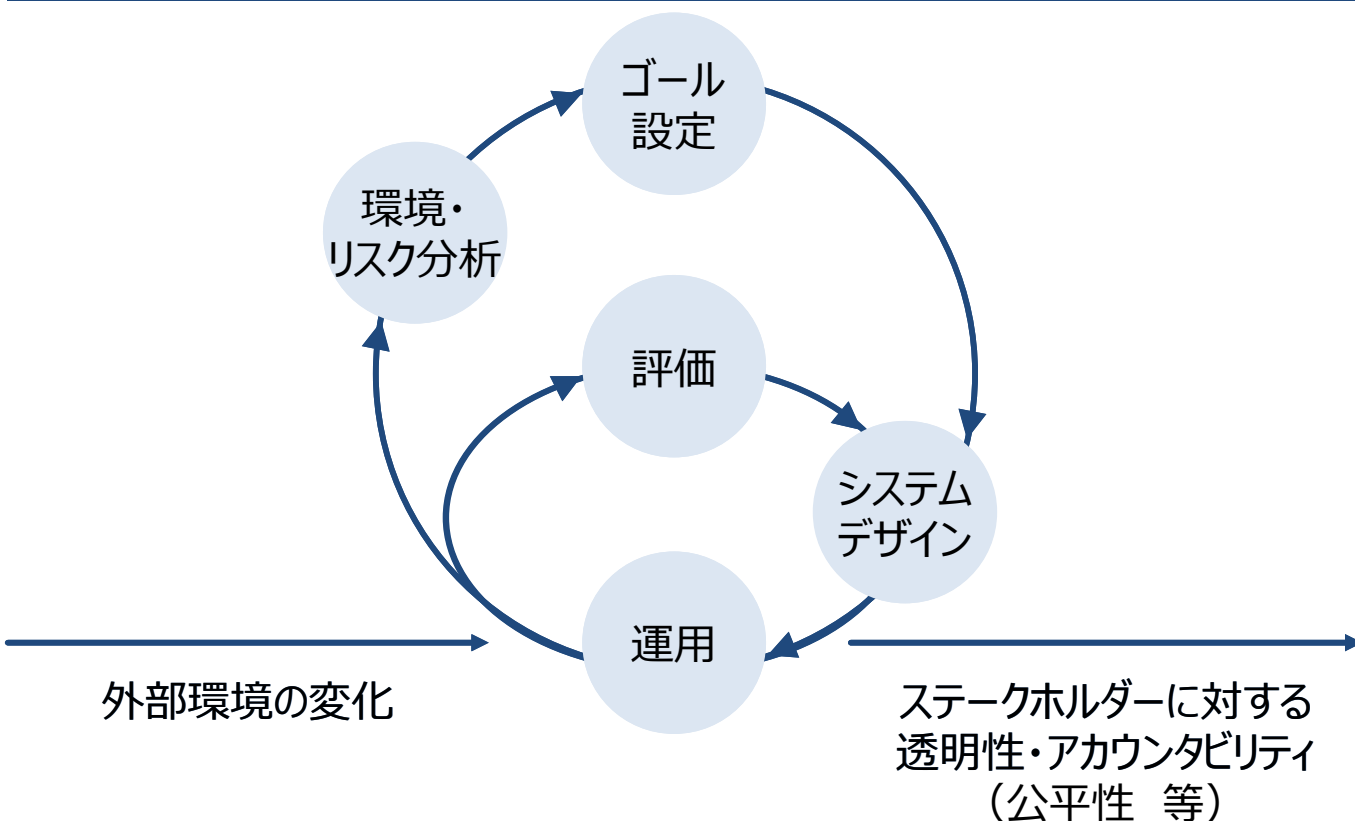
- 「共通の指針」に加え、以下を遵守すべきである\*1。ただし、I) ~ XI) は高度なAIシステムを開発するAI開発者にのみ適用される内容もあるため、各主体は適切な範囲で遵守することが求められる。
  - I. AIライフサイクル全体にわたるリスクを特定、評価、軽減するために、高度なAIシステムの開発全体を通じて、その導入前及び市場投入前も含め、適切な措置を講じる
  - II. 市場投入を含む導入後、脆弱性、及び必要に応じて悪用されたインシデントやパターンを特定し、緩和する
  - III. 高度なAIシステムの能力、限界、適切・不適切な使用領域を公表し、十分な透明性の確保を支援することで、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





- 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、**開発の各工程で留意すべき点を確認**する
- 
- D-6) ii. 関連する  
ステークホルダーへの  
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**情報提供**を行う
- 
- D-7) i. AI提供者への共通の  
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**生じるリスク等の情報提供及び説明**を行う
- 
- D-7) ii. 開発関連情報の  
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**文書化**する
- 
- D-10) i. イノベーションの  
機会創造への貢献
- AIの**品質・信頼性、開発の方法論等の研究開発**を行う
  - 持続的な経済成長の維持及び社会課題解決**につながるよう貢献する
  - 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
  - 社会全体への情報提供**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る<br>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の<br>仕組み及び対策の導入               | - 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の<br>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>   |

高度な AI システムを取り扱うAI提供者は、「第 2 部D. 高度なAIシステムに関する事業者に通の指針」について、I) ～XI) を適切な範囲で遵守し、XII)について遵守すべき

- 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

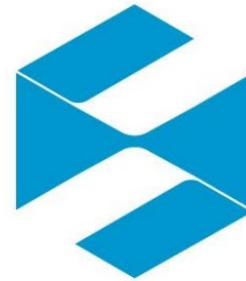
- |          |                           |  |
|----------|---------------------------|--|
| U-2) i.  | 安全を考慮した<br>適正利用           | <ul style="list-style-type: none"> <li>- AI提供者が定めた利用上の留意点を遵守して、<b>AI提供者が設計において想定した範囲内で利用</b>する</li> <li>- AIの出力について精度及びリスクの程度を理解し、<b>様々なリスク要因を確認した上で利用</b>する</li> </ul>  |
| U-3) i.  | 入力データ又はプロンプトに含まれるバイアスへの配慮 | <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.  | 個人情報の不適切入力及びプライバシー侵害への対策  | <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し、防止を検討する</b></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し、AI提供者とも連携の上説明及び要望の受付を行う</b></li> </ul> |
| U-7) ii. | 提供された文書の活用と規約の遵守          | <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>   |





総務省

Ministry of Internal Affairs  
and Communications



経済産業省

*Ministry of Economy, Trade and Industry*