

Tentative  
Translation

# **Overview of Draft 2019 Report**

June 2019

The Conference toward AI Network Society

## “2019 Report” of the conference toward AI Network Society

- Preface
- Chapter 1: Recent Trends in AI Networking
  - 1. Domestic trends (Trends in Japan)
  - 2. Overseas trends (Trends in the world)
  - 3. Trends in international discussions
- Chapter 2: Concept of the formulation of AI Utilization Guidelines
  - 1. Background and history
  - 2. Positioning of AI Utilization Guidelines
  - 3. Overview of AI Utilization Guidelines
  - 4. Issues to be studied in the future
- Chapter 3: Future Challenges
- Conclusion

### <Attachment 1> AI Utilization Guidelines

- Purpose and basic philosophies
- Classification of Related Entities
- AI Utilization Principles
- General flow of AI utilization
- Commentary on “AI Utilization Principles”
- Utilization phases that the Guidelines should be taken into consideration for

### <Attachment 2> Comparison of AI guidelines

### Domestic trends (Trends in Japan)

#### ➤ **Announcement of the “Social Principles of Human-Centric AI” (March 29, 2019)**

The government established the "Council for Social Principles of Human-centric AI," under the AI Strategy Expert Meeting for Strength and Promotion of the Innovation, for the purpose of formulating the basic principles for implementing and sharing AI in a better way and for usage in international discussions, and so released the “Social Principles of Human-Centric AI.” The social principles consist of seven principles: (1) Human-Centric; (2) Education/Literacy; (3) Privacy Protection; (4) Ensuring Security; (5) Fair Competition; (6) Fairness, Accountability and Transparency; and (7) Innovation.

### Overseas trends (Trends in the world)

#### ➤ **The Institute of Electrical and Electronics Engineers (IEEE) announced “Ethically Aligned Design, 1<sup>st</sup> edition” (March 25, 2019)**

The “Ethically Aligned Design” describes that the ethical and values-based design, development, and implementation of autonomous and intelligent systems(A/IS) should be guided by the following General Principles: (1) Human Rights; (2) Well-being; (3) Data Agency; (4) Effectiveness; (5) Transparency; (6) Accountability; (7) Awareness of Misuse; and (8) Competence. In addition, Ethically Aligned Design focuses on “From principle to practice” that is, IEEE launched projects of the IEEE P7000™ series of standards that explicitly focus on societal and ethical issues associated with a certain field of technology, and created an A/IS Ethics Glossary.

#### ➤ **HLEG on AI set up by European Commission announced “Ethics Guidelines for Trustworthy AI” (April 8, 2019)**

The High-Level Expert Group (HLEG) on AI set up by the European Commission presented their “Ethics guidelines for trustworthy artificial intelligence”. According to the guidelines, trustworthy AI should be lawful, ethical and robust, and it identifies four principles based on fundamental rights (respect for human autonomy, prevention of harm, fairness and explicability) and seven requirements for trustworthy AI: (1) human agency and oversight; (2) technical robustness and safety; (3) privacy and data governance; (4) transparency; (5) diversity, non-discrimination and fairness; (6) environmental and societal well-being; and (7) accountability. In addition, it has provided an assessment list aimed at operationalising the requirements.

### Trends in international discussions

#### ➤ **AI Expert Group at the OECD (AIGO) (France, etc., From September 2018 to February 2019)**

The Committee on Digital Economy Policy (CDEP) established the AI Expert Group at the OECD (AIGO) as part of a discussion on the need for the OECD Council to adopt policy principles that foster trust in, and adoption of, AI. It completed its proposal regarding principles for responsible stewardship of trustworthy AI and recommended national policy priorities for trustworthy AI.

#### ➤ **OECD and partner countries adopted the “Recommendation of the Council on Artificial Intelligence” (May 22, 2019)**

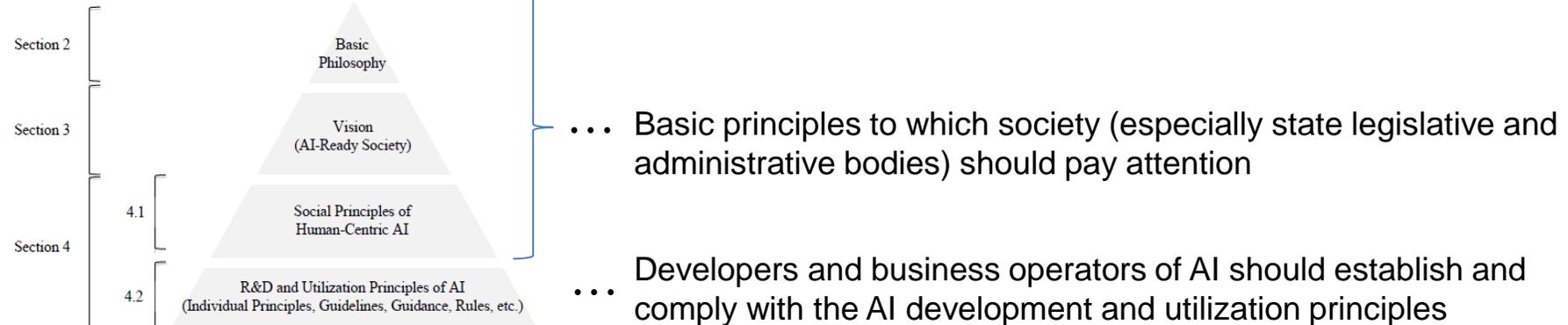
The OECD’s 36 member countries along with six partner countries signed up to the OECD Principles on Artificial intelligence at the Organisation’s annual Ministerial Council Meeting. The recommendations identify five principles for responsible stewardship of trustworthy AI, namely: (1) inclusive and sustainable growth and well-being; (2) human-centred values and fairness; (3) transparency and explainability; (4) robustness and safety; and (5) accountability. They also include recommended national policy priorities for trustworthy AI. The recommendation consists of the high level principles and specific measures to be taken will be considered continuously at the CDEP meeting after the formulation of the recommendation.

#### ➤ **G7 Multistakeholder Conference on Artificial Intelligence (Canada, December 6, 2018)**

Several key AI experts and G7 focal point organizations collaboratively drafted discussion papers for each topic: (1) AI for Society; (2) Unleashing Innovation; (3) Accountability in AI; and (4) the Future of Work, and directed breakout sessions during the multistakeholder conference. Japan together with Canada was in charge of (3) Accountability in AI.

## Reference when private sectors, etc. formulate their own principles

### Structure of “Social Principles of Human-Centric AI [1]”



Each developer or business operator is encouraged to establish AI development and utilization principles.

Practical guidance to be referred to is needed.

## Contribution to the international discussions

Consensus on AI principles is being reached. Hereafter, the focus of the discussion will be on **how to implement the AI principles**. Japan will continue to contribute to the international discussions, and foster sharing the recognitions.

(Examples of the discussion on how to implement AI principles)

[European Commission]

- **Assessment list** in the “Ethics Guidelines for Trustworthy AI”
- A revised version of the list, taking into account the feedback gathered through the piloting phase, will be presented to the European Commission in early 2020.

[OECD]

- Practical measures to realize “Recommendation of the Council on Artificial Intelligence” (= **Practical Guidance**)
- Discussion at OECD/CDEP will start from summer of 2019.

[1] <https://www.cas.go.jp/jp/seisaku/jinkouchinou/pdf/humancentricai.pdf>



## Purpose

To **facilitate AI utilization and social implementation** by way of increasing the benefits, and mitigating the risks of AI, as well as fostering trust in AI through the sound progress of AI networks

## Basic Philosophies

- To achieve a **human-centred society**
- To respect the **diversity** of users and advance **inclusion** of people with diverse backgrounds
- To achieve a **sustainable society** that can solve various problems faced by individuals, local communities, countries, and the international community
- To ensure an appropriate **balance between the benefits and risks** of AI
- To **realize appropriate role assignment among stakeholders** with consideration for **ability and knowledge** on AI that each user is expected to have
- To **share the Guidelines and their best practices** internationally among stakeholders
- To constantly **review** the Guidelines and **flexibly revise** them as necessary

## Developers

Those who conduct the R&D of AI systems

## Users

Those who use AI systems, AI services, or AI-accompanying services.

## Data providers

Those who provide data for learning or other methods of AI systems used by others.

## Third parties

Those whose rights and interests are affected due to AI systems or AI services used by others.

## AI service providers

Users who provide AI services or AI-accompanying services on business basis to others.

## End users

Users who use AI systems or AI services without providing AI services or AI-accompanying services on business basis to others.

## Business users (including non-profit specialists and administrative bodies)

End users who use AI systems or AI services on business basis.

(Note) It is expected that business users who use AI systems or AI services without operating AI systems themselves won't be able to pay the same level of attention as other business users. However, even in such cases, they are expected to promptly request appropriate measures from developers or AI service providers, etc.

## Consumer users

End users who use AI systems or AI services (excluding business users)

(Note) Consumer users who operate AI systems themselves may be required to pay the same level of attention as developers and AI service providers, etc.

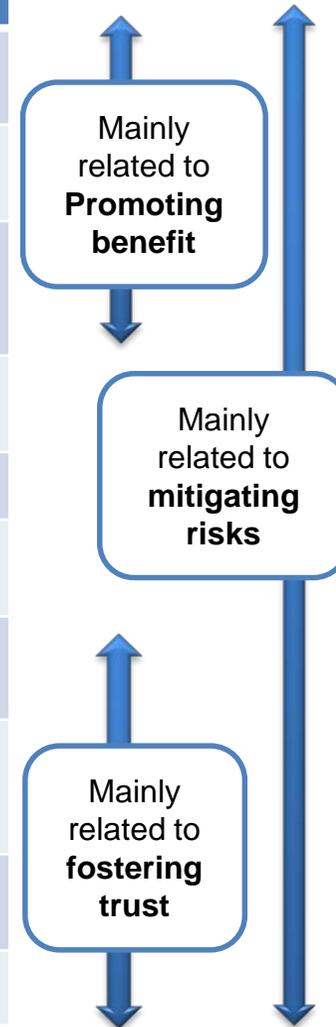
- AI systems : Systems that incorporate AI software as a component.
- AI services : Services that provide the functions of AI systems
- AI-accompanying services : AI-systems-update services or additional learning services, etc.

(Note) One individual or enterprise may be included in multiple entities.

# AI Utilization Principles

Compiled AI Utilization Principles that are expected to be referred to by AI service providers, business users, etc. and are to be shared internationally

Principle of	Description
<b>1. Proper Utilization</b>	Users should make efforts to utilize AI systems or AI services in a proper scope and manner, under the proper assignment of roles between humans and AI systems, or among users.
<b>2. Data quality</b>	Users and data providers should pay attention to the quality of data used for learning or other methods of AI systems.
<b>3. Collaboration</b>	AI service providers, business users, and data providers should pay attention to the collaboration of AI systems or AI services. Users should take into consideration that risks might occur and even be amplified when AI systems are to be networked.
<b>4. Safety</b>	Users should take into consideration that AI systems or AI services in use will not harm the life, body, or property of users or third parties through the actuators or other devices.
<b>5. Security</b>	Users and data providers should pay attention to the security of AI systems or AI services.
<b>6. Privacy</b>	Users and data providers should take into consideration that the utilization of AI systems or AI services will not infringe on the privacy of users' or others.
<b>7. Human dignity and individual autonomy</b>	Users should respect human dignity and individual autonomy in the utilization of AI systems or AI services.
<b>8. Fairness<sup>1</sup></b>	AI service providers, business users, and data providers should pay attention to the possibility of bias inherent in the judgements of AI systems or AI services, and take into consideration that individuals will not be discriminated unfairly by their judgments.
<b>9. Transparency<sup>2</sup></b>	AI service providers and business users should pay attention to the verifiability of inputs/outputs of AI systems or AI services and the explainability of their judgments.
<b>10. Accountability<sup>3</sup></b>	Users should make efforts to fulfill their accountability to the stakeholders.

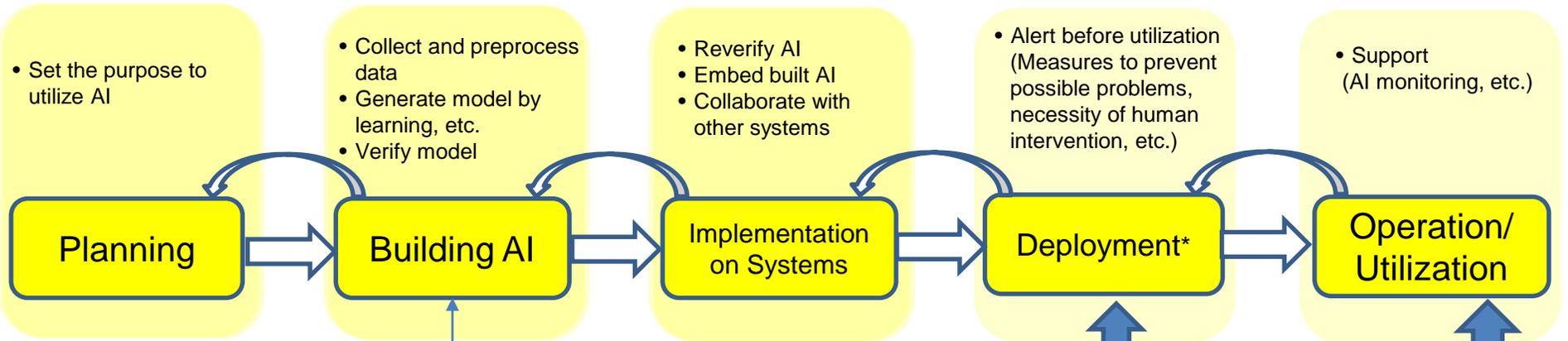


1) It should be noted that there are multiple definitions and criteria for “fairness.”

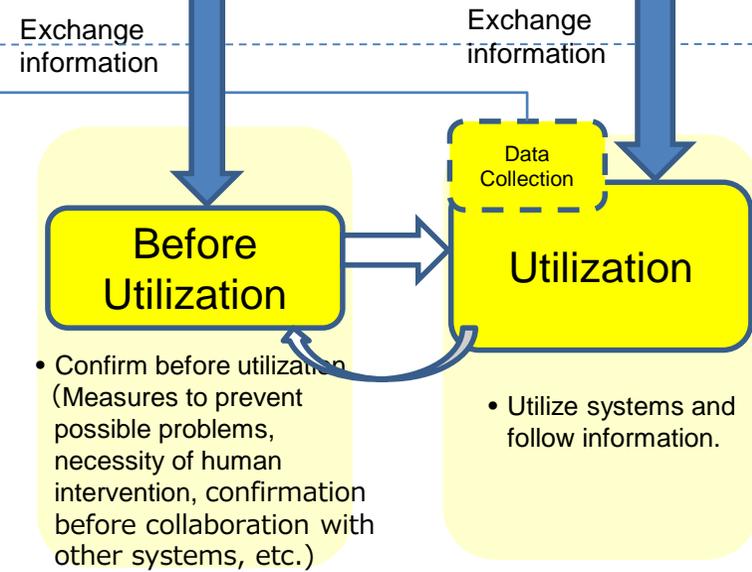
2) This principle is not intended to ask for the disclosure of algorithm, source code, or learning data. In interpreting this principle, privacy of individuals and trade secrets of enterprises are also taken into account.

3) Accountability: In order to convince the people who are affected by the judgments of AI systems or AI services, taking measures such as explaining the proper meanings and reasons for the judgments and if necessary, compensating them

### Case where the user operates AI Systems



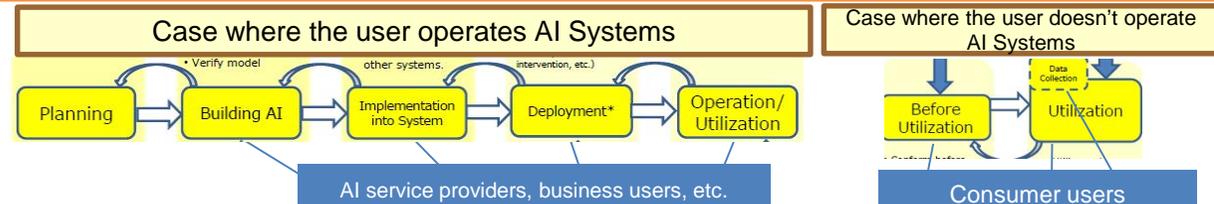
### Case where the user doesn't operate AI Systems



\*Deployment: To make AI software/systems available

Note: This general flow of AI utilization describes a typical case in order to clarify the phase in which each principle is to be considered. Keep in mind that there are various cases of AI utilization, such as in the case that development and operation are performed simultaneously (ex. DevOps).

# Linking the points of each Principle and utilization phases



Principle	Points of each principle	AI service providers, business users, etc.				Consumer users			
		Building AI	Implementation into system	Deployment	Operation/Utilization	Before utilization	Utilization	Data Collection	
<b>1. Proper Utilization</b>	a) Utilization in the proper scope and manner	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
	b) Human intervention	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
	c) Cooperation among stakeholders			<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
<b>2. Data quality</b>	a) Attention to the quality of data used for learning or other methods of AI	<input type="radio"/>						<input type="radio"/>	
	b) Attention to security vulnerabilities of AI by learning inaccurate or inappropriate data	<input type="radio"/>		<input type="radio"/>			<input type="radio"/>	<input type="radio"/>	
<b>3. Collaboration</b>	a) Attention to the interconnectivity and interoperability of AI systems		<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
	b) Address the standardization of data formats, protocols, etc.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	
	c) Attention to problems caused and amplified by AI networking		<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
<b>4. Safety</b>	a) Consideration for the life, body, or property		<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
<b>5. Security</b>	a) Implementation of security measures		<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
	b) Service provision, etc. for security measures			<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
	c) Attention to security vulnerabilities of learned models	<input type="radio"/>		<input type="radio"/>			<input type="radio"/>	<input type="radio"/>	
<b>6. Privacy</b>	a) Respect for the privacy of end users and others		<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
	b) Respect for the privacy of others in the collection, preprocessing, provision, etc. of personal data	<input type="radio"/>		<input type="radio"/>		<input type="radio"/>		<input type="radio"/>	
	c) Attention to the infringement of the privacy of users' or others and prevention of personal data leakage		<input type="radio"/>				<input type="radio"/>		
<b>7. Human dignity and individual autonomy</b>	a) Respect for human dignity and individual autonomy			<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
	b) Attention to the manipulation of human decision making, emotions, etc. by AI			<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
	c) Reference to the discussion of bioethics, etc. in the case of linking AI systems with the human brain and body		<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
	d) Consideration for prejudice against the subject in profiling which uses AI	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>		
<b>8. Fairness</b>	a) Attention to the representativeness of data used for learning or other methods of AI	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	
	b) Attention to unfair discrimination by algorithm	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	
	c) Human intervention (viewpoint of ensuring fairness)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>				
<b>9. Transparency</b>	a) Recording and preserving logs such as inputs/outputs, etc. of AI		<input type="radio"/>	<input type="radio"/>	<input type="radio"/>				
	b) Ensuring explainability	<input type="radio"/>							
	c) Ensuring transparency when AI is used in administrative bodies	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>				
<b>10. Accountability</b>	a) Efforts to fulfill accountability	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	
	b) Notification and publication of usage policy on AI systems or AI services	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	

※ This table assumes that AI service providers, business users, etc., operate AI by themselves and that consumer users do not operate AI by themselves.

## Concept of organization of detailed explanation

- The detailed explanation for each point, of each principle, is organized as a matter to be noted from the following perspectives (the lower left text (in Japanese)) :
  - AI service providers, business users and data providers
  - Consumer users
- The explanation is supplemented by a concept diagram on “machine learning”, which is used as major AI technologies in recent years, on the assumption that “the AI Utilization Guidelines” are reviewed regularly (the lower right figure)

## Examples of detailed explanation

AIサービスプロバイダ、ビジネス利用者とデータ提供者は、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（正確性や完全性など）に留意することが期待される。特に機械学習においては、以下の方法によりデータの質を確保することが考えられる。

【データ収集時の対策（例）】

- ・ 収集するデータがAIの利用目的に適ったものかを確認する。
- ・ 社会的に信用の高い者が公開するデータを活用する。
- ・ データの作成履歴を確認した上で収集する。
- ・ 自らデータを収集する際には、データに付随する権利に留意する。

【データ前処理時の対策（例）】

- ・ 人間でも判定が困難と考えられるデータは、学習等の対象から除外する<sup>1</sup>。
- ・ 機械（学習器）が誤認識しやすいと考えられるデータは積極的に学習の対象とする<sup>2</sup>。
- ・ （特に教師あり学習等で）アノテーション（ラベル付与）を行う際には、誤って行わないよう留意する。
- ・ 利用時に利用（入力）されるデータの形式を意識してデータセットを作成する。
- ・ 前処理をどのように行ったのか（データ前処理に関する来歴）について、ログを取得し保存する。

【学習時の対策（例）】

- ・ 既存の学習モデルを利用して転移学習<sup>3</sup>等を行う。
- ・ 学習の精度を上げるため、特定のデータを拡張した上で学習を行う。
- ・ 一過性のある時系列データを学習する場合などは、学習対象とすべきデータの範囲を適切に画定する。

1) 例えば、画像認識などで、対象となるオブジェクトが人間の目で見てても画定できない場合など。  
2) 例えば、画像認識などで、対象となるオブジェクトが稀にあるなど。  
3) 転移学習(Transfer Learning)とは、深層学習を含む機械学習で用いられる技術の1つで、特定の領域(ドメイン)で学習させたモデルを別の領域に適用する技術である。少ないデータで高精度な学習結果を得ることが出来る可能性があるのがメリットである。  
4) 「データの拡張(Data Augmentation)とは、特定の学習データが少ない際、汎化性能(未知のデータに対する性能)を高めることにより、データの正確性を確保するために用いられる手段の1つである。学習に用いるデータを拡張し(例えば画像データであれば、反転、拡大、縮小を適用し)それぞれを別のソースに基づくデータとして用いることにより、汎化性能が改善されることがある。

また、AIによりなされる判断は、事後的に精度が異なったり、低下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、技術水準、精度を維持するためのコスト等<sup>1</sup>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ことが期待される。

なお、消費者的利用者が提供するデータを用いることが予定されている場合には、AIの特性及び用途を踏まえ、データ提供の手段、形式等について、あらかじめ消費者的利用者に情報を提供することが期待される。

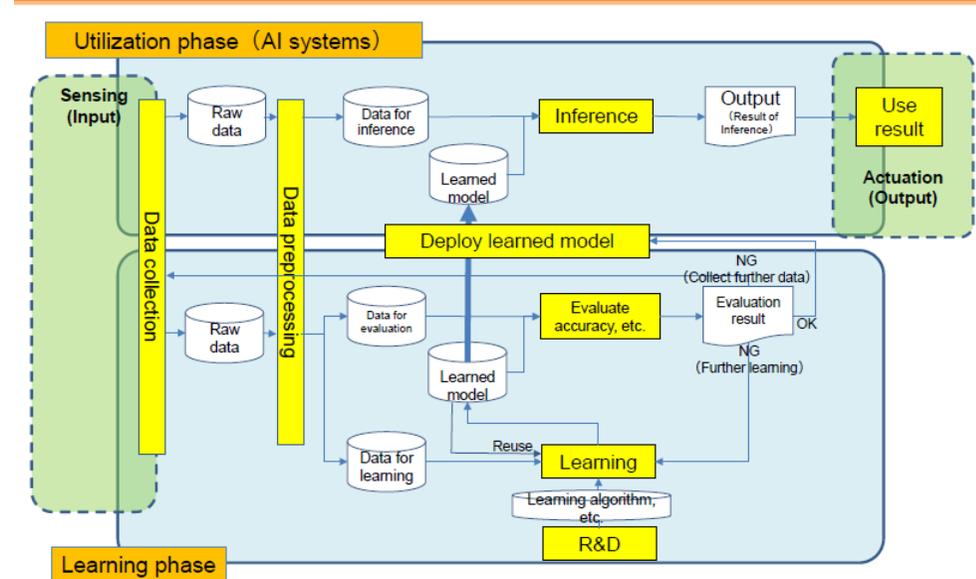
1) 例えば、機械学習を中心としたAIは帰納的な処理を行うため、当該AI単体では、判断結果に基づき原則的に100%の精度を担保できないこと等が挙げられる。

<参考>

消費者的利用者は、自らデータを収集し、利用するAIの学習等を行うことが予定されている場合には、データの形式について、開発者、AIサービスプロバイダ等から提供された情報を踏まえた上でデータの収集、保存を行うことが望ましい。

Detailed explanation example (Focusing on notes on AI service providers, business users and data providers (In addition, notes on consumer users as reference))

## Flow of Learning and Utilization for Machine Learning



Illustrated flow example

	issues	overview
<b>1. Matters related to the sound development of AI networking</b>		
(1)	Dissemination and development of “AI R&D/Utilization Guidelines”	Holding symposiums to disseminate “AI R&D/Utilization Guidelines”; Dissemination of detailed explanations to realize the principles in international frameworks, etc.
(2)	Following-up for discussions regarding AI development/ utilization principles/guidelines	Following-up and continuously reviewing international discussions on AI development / utilization principles/guidelines.
(3)	Issues related to environmental improvements addressed by related stakeholders	Cooperation among stakeholders, sharing of best practices, and studies on the ideal state of legal systems.
(4)	Securement of the smooth collaboration of AI systems or AI services	Studies on the range of related information expected to be shared among stakeholders and how to share it.
(5)	Securement of a competitive ecosystem	Keeping watch on the trends of related markets.
(6)	Protection of the interests of users	Studies on the manner of voluntary information provision from developers to users and on the ideal state of protecting users (e.g., by insurance).
<b>2. Matters related to the evaluation of the socioeconomic impact of AI networking</b>		
(1)	Scenario analysis on the socioeconomic impact of AI networking	Ongoing implementation of scenario analysis and international sharing.
(2)	Establishment of evaluation indicators for the impact of the progress of AI networking, and for richness and happiness	Study on setting indicators.
(3)	Fostering social acceptability on the utilization of AI systems	Keeping watch on the degree of social acceptance for the utilization of AI systems.
<b>3. Matters related to issues over a human who is in a society under the ongoing progress of AI networking</b>		
(1)	Study on the ideal state of relationship between humans and AI systems	Studies on the ideal state of role assignments of professionals (doctors, lawyers, accountants, etc.) and AI systems.
(2)	Study on the ideal state of relationship among stakeholders	Studies on the ideal state of responsibility sharing in case of AI’s risks becoming apparent.
(3)	Safety net development	Keeping watch on labor market trends and prevention of the unfair redistribution of income accompanying the progress of AI networking.

# References

## 4. Principle of Safety

AI service providers should take the following measures with consideration for potential damage, etc.:

- (A) Build a mechanism to ensure the safety of the entire system when implementing AI (fail-safe)
- (B) Inspect, repair and update AI

## 8. Principle of Fairness

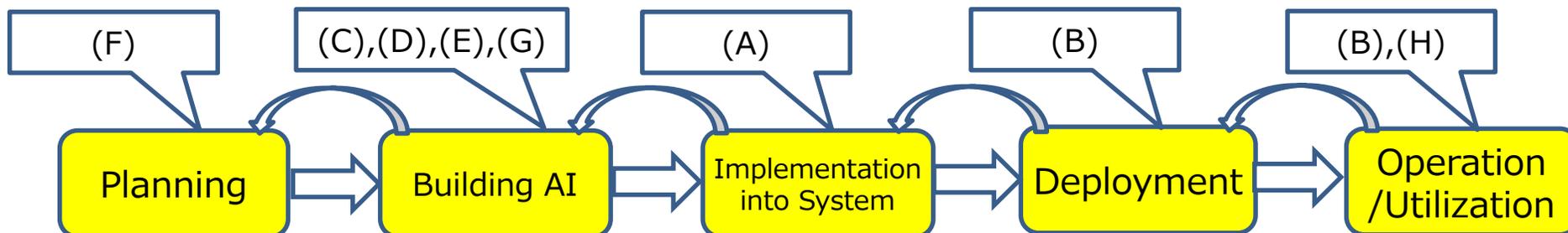
AI service providers should take the following actions according to the social context, etc.:

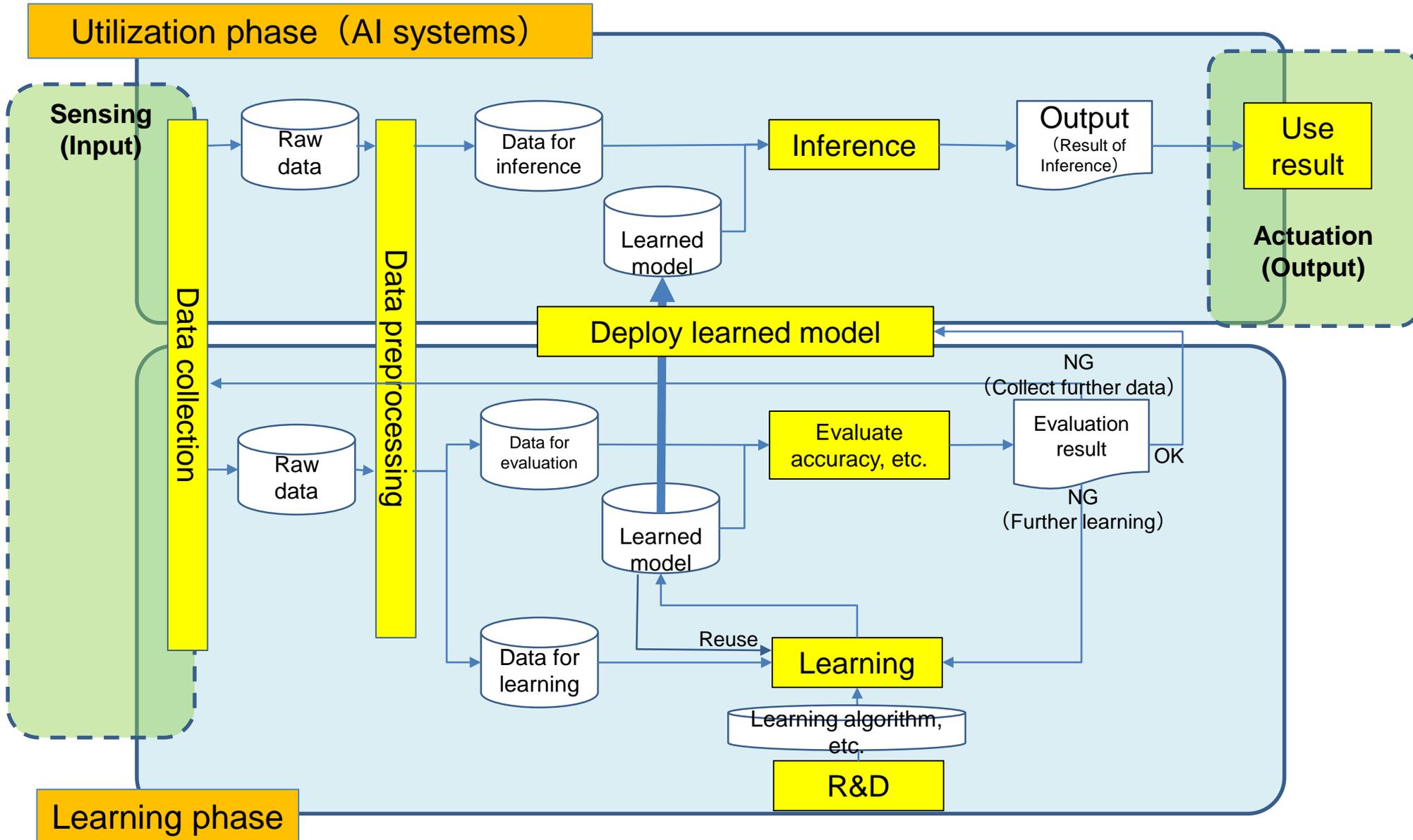
- (C) Ensure representativeness of learning data
- (D) Remove biases inherent in learning data
- (E) Eliminate biases of those who label learning data
- (F) Clarify matters to be fair and design an algorithm satisfying them

## 9. Principle of Transparency

AI service providers should take the following measures in light of AI usage, etc.:

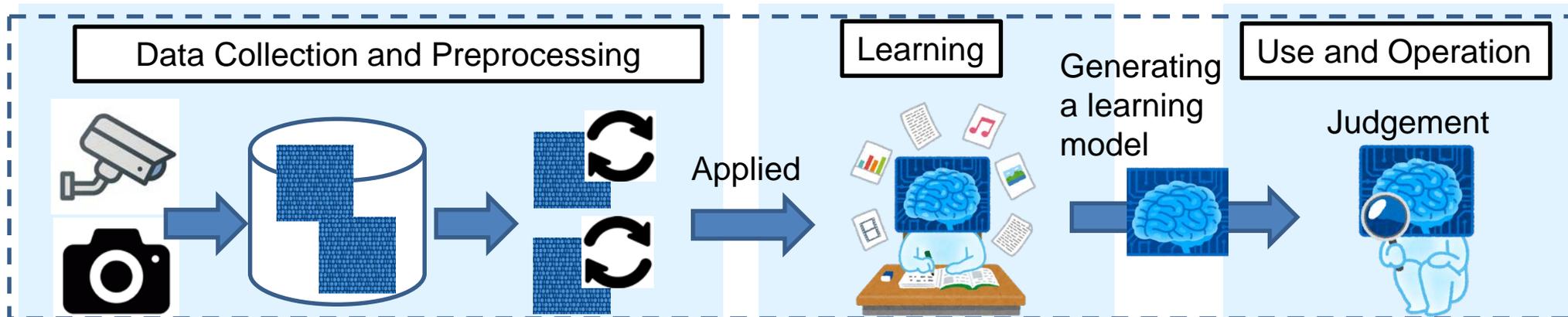
- (G) Manage the provenance of data used for learning, etc.
- (H) Save the related logs (input/output, etc.)





Indicate specific measures to implement the principles according to the process of AI utilization

Ex. : A system that judges whether there is a possibility of a crime being committed through human image input



[Fairness]

**Ensure data representativeness:**

Not to target persons living in specific regions

[Fairness]

**Eliminate bias in annotating data:**

Not to label data with personal intentions and impressions

[Transparency]

**Data provenance**

[Privacy]

**Respect for the privacy** of persons who were captured on the images

[Fairness]

**Attention to unfair discrimination by algorithm:**

Not to discriminate in the computing process

[Proper utilization]

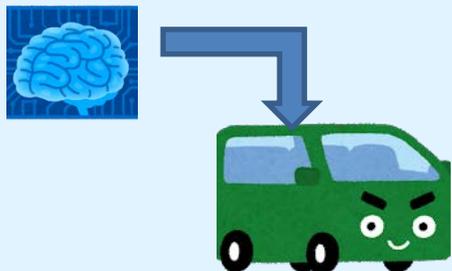
**Human intervention:**

Consider end user's right to be influenced by AI judgements

Indicate specific measures to implement the principles according to the process of AI utilization

## Ex. Autonomous driving

Build system with AI



Deployment



Use and Operation



[Safety]  
Ensure safety across the entire system (**Fail-safe**)

[Security]  
**Take reasonable measures corresponding to the current technology level** to prevent system hacking

[Collaboration]

- Negotiate and coordinate among autonomous vehicles, and **support data format / protocol**
- Address risks that a problem in one AI system spreads to the entire system

[Safety/Security]  
**Share information** about measures to be taken when infringement occurs

[Proper Utilization]  
**Human intervention:**  
Share conditions on switching control from AI to human

[Safety/Security]  
**Provide updates (information) for systems with AI**

[Transparency/Accountability]  
**Ensure explainability, and fulfill accountability**, when accident occurs

# Examples of Trade-offs

↔  
The relationship seems to be in a trade-off.

