

国際的な議論のための A I 開発ガイドライン案

平成 29 年 7 月 28 日

A I ネットワーク社会推進会議

目次

序文

A I 開発ガイドライン案

- 1) 目的
- 2) 基本理念
- 3) 用語の定義及び対象範囲
- 4) 開発原則
- 5) 開発原則の解説

【別添】 関係するステークホルダに期待される役割

【参考】 A I ネットワーク化と智連社会

序文

AIに関する研究開発や利活用は今後飛躍的に発展することが期待されている。こうした中、2016年4月に日本で開催されたG7情報通信大臣会合において、ホスト国である日本はAI開発原則のたたき台を紹介し、各国関係閣僚による議論が行われた。その結果、G7において「AI開発原則」及びその内容の解説からなる「AI開発ガイドライン」の策定に向け、引き続きG7各国が中心となり、OECD等国际機関の協力も得て議論していくことで合意した。

少子高齢化などの課題を抱える日本は、AIを積極的に開発し利活用することにより様々な課題を解決するとともに、その知見を活かしつつAIの開発について留意することが期待される事項等を国際的に発信することにより、国際社会に大きく貢献することができる。

本ガイドライン案は、AIの便益の増進及びリスクの抑制のため研究開発において留意することが期待される事項に関するG7やOECDにおける国際的な議論のための基礎となる文書として作成されたものである。AIに関する技術がその発展の途上にあることに鑑みれば、国際的な「AI開発原則」及びその内容の解説からなる「AI開発ガイドライン」は、規制の導入を目指すものとすることは適当ではない。そこで、本ガイドライン案は、非規制的で非拘束的なソフトローとして国際的に共有される指針の案として作成されたものである。こうした指針に盛り込むことが適当な内容を議論することを通じ、

- 「AI開発ガイドライン」及び「AI利活用ガイドライン」の策定に向け、AIに関する研究開発や利活用に関係する多様なステークホルダ（開発者、サービス提供者、市民社会を含む利用者、各国政府、国際機関など）を含む国内及び国際的な議論を加速化すること
- AIの研究開発及び利活用におけるベストプラクティスの国際的な共有を促すことにより、AIへの利用者や社会の信頼を獲得し、AIの研究開発及び利活用の円滑化に資することが期待される。

A I 開発ガイドライン案¹

1. 目的

A I の研究開発や利活用は、今後急速に進展することが期待されているところであり、A I ネットワーク化（A I システムがインターネットその他の情報通信ネットワークと接続され、A I システム相互間又はA I システムと他の種類のシステムとの間のネットワーク（以下において「A I ネットワーク」という場合がある。）が形成されるようになることをいう。以下同じ。）が進展していく過程で、個人、地域社会、各国、国際社会²の抱える様々な課題の解決に大きく貢献するなど、人間及びその社会や経済に多大な便益を広範にもたらすことが期待される。このような方向に向けて、A I の研究開発や利活用を加速化していくことが求められる。

その一環として、A I システムが社会や経済にもたらす便益の増進を図るとともに、不透明化や制御喪失などA I システムに関するリスクの抑制を図る観点から、関連する社会的・経済的・倫理的・法的な課題に対応することが必要となる。特に、A I システムを利活用するサービスは、他の情報通信サービス同様、ネットワークを通じて国境を越えて提供されるものであることから、オープンな議論を通じ、国際的なコンセンサスを醸成し、非規制的で非拘束的なソフトローたるガイドラインやそのベストプラクティスをステークホルダ（開発者、サービス提供者、市民社会を含む利用者、各国政府、国際機関など）間で国際的に共有することにより、A I システムの便益の増進とリスクの抑制を図ることが求められる。

以上の問題意識に鑑み、本ガイドラインは、A I ネットワーク化の健全な進展を通じてA I システムの便益の増進とリスクの抑制を図ることにより、利用者の利益を保護するとともにリスクの波及を抑止し、人間中心の智連社会³を

¹ 本ガイドライン案は、国際的に共有される「A I 開発ガイドライン」の策定に向けた国際的な議論に用いられることを念頭に作成されたものである。

² 国際社会の抱える課題については、国連の「持続可能な開発目標」(SDGs) (http://www.un.org/ga/search/view_doc.asp?symbol=A/70/L.1)などを参照。

³ 智連社会 (Wisdom Network Society) とは、A I ネットワーク化の進展の結果として、人間がA I ネットワークと共生し、データ・情報・知識を自由かつ安全に創造・流

実現することを目的とする。

そこで、本ガイドラインは、上記の目的を達成する観点から、今後のA Iシステムの開発において留意することが期待される事項であるA I開発原則及びその内容の解説を取りまとめたものである。

なお、A Iシステムの研究開発は多様な利活用の分野に及ぶものであり、分野ごとにA Iシステムのもたらす便益やリスクは異なる可能性がある。このため、本ガイドラインでは、A Iの開発に関し留意することが期待される事項のうち、利活用の分野に共通する事項又は分野間の連携に関し留意することが期待される事項について定めることとし、分野ごとの事情に応じて留意することが期待される事項については、本ガイドラインとは別に、各々の分野ごとに当該分野に係るガイドラインの在り方に関し、策定の要否を含め、各分野の国際機関を含む関係するステークホルダによる議論が行われることが期待される。

また、A Iシステムは利活用の過程を通じて学習等により出力やプログラムが継続的に変化する可能性があることから、開発者が留意することが期待される事項のみならず、利用者が留意することが期待される事項も想定される。このため、本ガイドラインとは別にA I利活用ガイドラインを策定することについて、国際的な議論を進めていくことが期待される。

2. 基本理念

本ガイドラインの目的に鑑み、次に掲げる理念を一体的なものとして本ガイドラインの基本理念とする。

1. 人間がA Iネットワークと共生することにより、その恵沢がすべての人によってあまねく享受され、人間の尊厳と個人の自律が尊重される**人間中心の社会を実現すること**。
2. A Iの研究開発と利活用が今後急速に発展し、ネットワーク化されたA Iシステムが国境を越えて人間及び社会に広範かつ多大な影響を及ぼすものと見込まれることから、A Iシステムの研究開発の在り方について、非拘束的なソフトローたる**指針やそのベストプラクティスをステークホルダ間で国際**

通・連結して「智のネットワーク」(Wisdom Network)を構築することにより、あらゆる分野におけるヒト・モノ・コト相互間の空間を越えた協調が進展し、もって創造的かつ活力ある発展が可能となる社会である (【参考】「A Iネットワーク化と智連社会」参照)。

的に共有すること。

3. イノベーティブでオープンな研究開発と公正な競争を通じ、A I ネットワークの便益を増進するとともに、学問の自由や表現の自由といった民主主義社会の価値を最大限尊重しつつ、A I ネットワークにより権利利益が侵害されるリスクを抑制するため、**便益とリスクの適正なバランスを確保**すること。
4. A I 関連の技術が引き続き急速に進展していくことが見込まれる中、**技術的中立性を確保**する観点から特定の技術や手法に基づくA I の研究開発を阻害しないよう配慮するとともに、**開発者にとって過度の負担とならないものとするよう留意**すること。
5. A I 関連技術やA I の利活用が今後とも飛躍的に発展することが期待されることから、A I ネットワーク化の進展等を踏まえ、国際的な議論を通じて、**本ガイドラインを不断に見直し、必要に応じて柔軟に改定**すること。また、本ガイドラインの見直しに際しては関係するステークホルダの参画を得るなど、広範で柔軟な議論に努めること。

3. 用語の定義及び対象範囲

3-1 用語の定義

2. 基本理念に鑑み、本ガイドラインにおける「A I」については、以下のとおり定義する。

「A I」とは、「A I ソフト及びA I システムを総称する概念」をいう⁴。

- 「A I ソフト」とは、データ・情報・知識の学習等⁵により、利活用の過程を通じて自らの出力やプログラムを変化させる機能を有するソフトウエ

⁴ 本ガイドラインにおけるA I の定義は、現在既に実用化されている特化型A I を主たる対象として想定しているが、自律性を有するA I や汎用A I (Artificial General Intelligence) の開発など今後予想されるA I に関する急速な技術発展を見据え、今後開発される多種多様なA I についても、学習等により自らの出力やプログラムを変化させる機能を有するものである場合には、含み得るものとしている。

本ガイドラインにおいては、基本理念4. に掲げる技術的中立性の確保の見地などから、上述のとおりA I を定義しており、今後開発される多種多様なA I についてもその機能次第で含み得るものとしている。なお、本ガイドラインにおけるA I の定義の在り方については、A I の技術発展の動向等を踏まえ、今後継続的に議論を行っていくことが必要である。

⁵ 学習以外の方法によりA I ソフトが自らの出力やプログラムを変化させる要因としては、例えば、データ・情報・知識に基づく推論や、センサやアクチュエータ等を通じた環境とのインタラクションなどが考えられる。

アをいう。例えば、機械学習ソフトウェアはこれに含まれる。

- 「**AIシステム**」とは、AIソフトを構成要素として含むシステムをいう。例えば、AIソフトを実装したロボットやクラウドシステムはこれに含まれる。

AIシステムの「開発者」及び「利用者」については、以下のとおり定義する。ただし、「開発者」及び「利用者」は場面に応じて個別に決まる相対的な概念であることに留意する必要がある。

- 「**開発者**」とは、AIシステムの研究開発（AIシステムを利用しながら行う研究開発を含む。）を行う者（自らが開発したAIシステムを用いてAIネットワークサービスを他者に提供するプロバイダを含む。）をいう。
- 「**利用者**」とは、AIシステムを利用する者（最終利用者（エンドユーザ）のほか、他者が開発したAIネットワークサービスを第三者に提供するプロバイダを含む。）をいう。

3-2 対象範囲

本ガイドラインの対象とする**AIシステムの範囲**は、AIシステムがネットワークを通じて国境を越えて利用され、広く人間及び社会に便益やリスクをもたらす可能性があることから、ネットワーク化され得るAIシステム（ネットワークに接続可能なAIシステム）とする。

本ガイドラインの対象とする**開発者の範囲**は、本ガイドラインが非拘束的なソフトローとしての指針であることに鑑み、3-1で定義された開発者すべてとする。

本ガイドラインの対象とする**開発の範囲**は、学問の自由の尊重、社会に与える影響の大きさ等に鑑み、閉鎖された空間（実験室、セキュリティが十分に確保されたサンドボックス等）内での開発は対象とせず、ネットワークに接続して行う段階とする。

4. AI開発原則

（主にAIネットワーク化の健全な進展及びAIシステムの便益の増進に関する原則）

① **連携の原則**-----開発者は、A I システムの相互接続性と相互運用性に留意する。

(主にA I システムのリスクの抑制に関する原則)

② **透明性の原則**-----開発者は、A I システムの入出力の検証可能性及び判断結果の説明可能性に留意する。

③ **制御可能性の原則**-----開発者は、A I システムの制御可能性に留意する。

④ **安全の原則**-----開発者は、A I システムがアクチュエータ等を通じて利用者及び第三者の生命・身体・財産に危害を及ぼすことがないよう配慮する。

⑤ **セキュリティの原則**-----開発者は、A I システムのセキュリティに留意する。

⑥ **プライバシーの原則**-----開発者は、A I システムにより利用者及び第三者のプライバシーが侵害されないよう配慮する。

⑦ **倫理の原則**-----開発者は、A I システムの開発において、人間の尊厳と個人の自律を尊重する。

(主に利用者等の受容性の向上に関する原則)

⑧ **利用者支援の原則**-----開発者は、A I システムが利用者を支援し、利用者を選択の機会を適切に提供することが可能となるよう配慮する。

⑨ **アカウントビリティの原則**-----開発者は、利用者を含むステークホルダに対しアカウントビリティを果たすよう努める。

5. 開発原則の解説

① **連携の原則**-----開発者は、A I システムの相互接続性と相互運用性に留意する。

(解説)

開発者は、A I ネットワーク化の健全な進展を通じて、A I システムの便益を増進するとともに、リスクの抑制に関する複数の開発者の取組が相互に調和して有効に機能するよう、A I システムの多様性を踏まえつつ、自らの開発するA I システムと他のA I システム等との相互接続性と相互運用性⁶に留意

⁶ ここで相互運用性と相互接続性としては、自らの開発するA I システムが情報通信ネッ

することが望ましい。そのため、開発者は、以下の事項について留意することが望ましい。

- 相互接続性と相互運用性を確保するために有効な関連情報の共有に向けて協力するよう努めること。
- 国際的な標準や規格がある場合には、これに準拠したA Iシステムを開発するよう努めること。
- データ形式の標準化、A P Iを含むインターフェースやプロトコルのオープン化に対応するよう努めること。
- 自らの開発するA Iシステムと他のA Iシステム等との相互接続・相互運用により意図しない事象が生ずるリスクに留意すること。
- A Iの開発に関連する知的財産に関し、保護と利活用のバランスに配慮しつつ、標準必須特許などA Iシステムと他のA Iシステム等との相互接続性・相互運用性の確保に資する知的財産権のライセンス契約及びその条件についてオープンかつ公平な取扱いを図るよう努めること。

②透明性の原則——開発者は、A Iシステムの入出力の検証可能性及び判断結果の説明可能性に留意する⁷。

(解説)

本原則の対象となるA Iシステムとしては、利用者及び第三者の生命、身体、自由、プライバシー、財産などに影響を及ぼす可能性のあるA Iシステムが想定される。

開発者は、A Iシステムに対する利用者を含む社会の理解と信頼が得られるよう、採用する技術の特性や用途に照らし合理的な範囲で、A Iシステムの入出力の検証可能性及び判断結果の説明可能性に留意することが望ましい。

③制御可能性の原則——開発者は、A Iシステムの制御可能性に留意する。

(解説)

開発者は、A Iシステムの制御可能性に関するリスクを評価するため、あら

トワークと接続され、他のA Iシステム等と相互に適切に協調して運用することが可能であることを念頭に置いている。

⁷ 本原則は、開発者によるアルゴリズム、ソースコード、学習データの開示を想定するものではない。また、本原則の解釈に当たっては、プライバシーや営業秘密への配慮も求められる。

かじめ検証及び妥当性の確認⁸を行うよう努めることが望ましい⁹。こうしたリスク評価の手法としては、社会において実用化される前の段階において、実験室内やセキュリティが確保されたサンドボックスなどの閉鎖空間において実験を行うことが考えられる。

また、開発者は、制御可能性を確保するため、採用する技術の特性に照らして可能な範囲において、人間や信頼できる他のAIによる監督（監視、警告など）や対処（AIシステムの停止、ネットワークからの切断、修理など）の実効性に留意することが望ましい。

④安全の原則——開発者は、AIシステムがアクチュエータ等を通じて利用者及び第三者の生命・身体・財産に危害を及ぼすことがないよう配慮する。

（解説）

本原則の対象となるAIシステムとしては、アクチュエータ等を通じて利用者及び第三者の生命、身体、財産に危害を及ぼす可能性のあるAIシステムが想定される。

開発者は、本原則に関する国際標準等を参照するとともに、特にAIシステムが学習等によって出力やプログラムが変化する可能性を踏まえ、以下の事項について留意することが望ましい。

- AIシステムの安全性に関するリスクを評価・抑制するため、あらかじめ検証及び妥当性の確認を行うよう努めること。
- AIシステムがアクチュエータ等を通じて稼動する際の本質安全（アクチュエータの運動エネルギーなど本質的な危険要因の低減）や機能安全（自動ブレーキなど付加的な制御装置の作動によるリスクの抑制）に資するよう、AIシステムの開発の過程を通じて、採用する技術の特性に照らし可能な範囲で措置を講ずるよう努めること。
- AIシステムを利用する際の利用者及び第三者の生命、身体、財産の安全に関する判断（例えば、AIを搭載したロボットの事故発生時に、優先的

⁸ 検証 (verification) 及び妥当性の確認 (validation) は、あらかじめリスクを評価し抑制するための手法であるが、前者は形式的な整合性の確認を意味して用いられるのに対し、後者は実質的な妥当性の確認を意味して用いられることが一般的である (See e.g., The Future of Life Institute (FLI), *Research Priorities for Robust and Beneficial Artificial Intelligence* (2015)).

⁹ リスク評価の要素としては、例えば、AIシステムが与えられた目標を形式的に達成するために開発者の意図に実質的に反する動作（報酬ハッキング）を行うリスクやAIシステムが学習等による利活用の過程を通じた変化に伴い開発者の意図しない動作を行うリスク等に配慮することが考えられる。See e.g., Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman & Dan Mané, *Concrete Problems in AI Safety*, arXiv:1606.06565 [cs.AI] (2016).

に保護される生命、身体、財産の順位などに関する判断)を行うA Iシステムを開発する場合において、利用者等ステークホルダに対して当該A Iシステムの設計の趣旨及びその理由を説明するよう努めること。

⑤セキュリティの原則——開発者は、A Iシステムのセキュリティに留意する。

(解説)

開発者は、OECDセキュリティガイドラインなどセキュリティに関する国際的な指針を踏まえるほか、A Iシステムが学習等によって出力やプログラムが変化する可能性があることを踏まえ、以下の事項について留意することが望ましい。

- A Iシステムの情報セキュリティについては、通常、情報の機密性、完全性及び可用性が確保されることが求められるが、必要に応じて、A Iシステムの信頼性（意図したとおりに動作が行われ、権限を有しない第三者による操作を受けないこと）や頑健性（物理的な攻撃や事故への耐性）にも留意すること。
- A Iシステムのセキュリティに関するリスクを評価・抑制するため、あらかじめ検証や妥当性確認を行うよう努めること。
- A Iシステムの開発の過程を通じて、採用する技術の特性に照らし可能な範囲でセキュリティ対策を講ずるよう努めること（セキュリティ・バイ・デザイン）。

⑥プライバシーの原則——開発者は、A Iシステムにより利用者及び第三者のプライバシーが侵害されないよう配慮する。

(解説)

本原則にいうプライバシーの範囲には、空間に係るプライバシー（私生活の平穏）、情報に係るプライバシー（個人データ）及び通信の秘密が含まれる。開発者は、OECDプライバシーガイドラインなどプライバシーに関する国際的な指針を踏まえるとともに、A Iシステムが学習等によって出力やプログラムが変化する可能性があることを踏まえ、以下の事項について留意することが望ましい。

- プライバシー侵害のリスクを評価するため、あらかじめプライバシー影響評価を行うよう努めること。
- A Iシステムの利活用時におけるプライバシー侵害を回避するため、当該システムの開発の過程を通じて、採用する技術の特性に照らし可能な範囲で措置を講ずるよう努めること（プライバシー・バイ・デザイン）。

⑦倫理の原則——開発者は、A I システムの開発において、人間の尊厳と個人の自律を尊重する。

(解説)

開発者は、人間の尊厳と個人の自律を尊重するに当たり、人間の脳や身体と連携するA I システムを開発する場合は、生命倫理に関する議論などを参照しつつ、特に慎重に配慮することが望ましい。

開発者は、採用する技術の特性に照らし可能な範囲で、A I システムの学習データに含まれる偏見などに起因して不当な差別が生じないよう所要の措置を講ずるよう努めることが望ましい。

開発者は、国際人権法や国際人道法を踏まえ、A I システムが人間性の価値を不当に毀損することがないよう留意することが望ましい。

**⑧利用者支援の原則——開発者は、A I システムが利用者を支援し、利用者
に選択の機会を適切に提供することが可能となるよう配慮する。**

(解説)

開発者は、A I システムの利用者を支援するため、以下の事項について留意することが望ましい。

- 利用者の判断に資する情報を適時適切に提供し、かつ、利用者にとって操作しやすいインターフェースが利用可能であることに配慮するよう努めること。
- 利用者に選択の機会を適時適切に提供する機能（例えば、デフォルトの設定、理解しやすい選択肢の提示、フィードバックの提供、緊急時の警告、エラーへの対処など）が利用可能であることに配慮するよう努めること。
- ユニバーサルデザインなど社会的弱者の利用を容易にするための取組に努めること。

また、開発者は、A I システムの学習等による出力又はプログラムの変化の可能性を踏まえ、利用者に対し適切な情報提供を行うよう努めることが望ましい。

**⑨アカウントビリティの原則——開発者は、利用者を含むステークホルダに
対しアカウントビリティを果たすよう努める。**

(解説)

開発者は、A I システムへの利用者や社会の信頼を得られるよう、自らの開発するA I システムについてアカウントビリティを果たすことが期待される。

具体的には、利用者に対しA I システムの選択及び利活用に資する情報を提供す

るとともに、利用者を含む社会によるA Iシステムの受容性を向上するため、開発者は、本ガイドラインの掲げる開発原則①～⑧の趣旨に鑑み、利用者等に対し自らの開発するA Iシステムの技術的特性について情報提供と説明を行うほか、多様なステークホルダとの対話を通じて様々な意見を聴取するなど、ステークホルダの積極的な関与（フィードバック）を得るよう努めることが望ましい。

また、開発者は、自らの開発するA Iシステムによってサービスを提供するプロバイダ等と情報を共有し、協力するよう努めることが望ましい。

【別添】関係するステークホルダに期待される役割

本ガイドラインの目的に鑑み、関係する産学民官のステークホルダは、例えば以下のような役割を果たすことが期待される。

1. 各国政府及び国際機関は、本ガイドラインの運用や見直しにおいて各国政府、国際機関、開発者、市民社会を含む利用者など**多様なステークホルダ間の対話の促進**に向けた環境整備に努めることが期待される。
2. 開発者、市民社会を含む利用者など関係するステークホルダは、上記1の対話に参加するとともに、本ガイドラインに適う**ベストプラクティスを共有**し、AIをめぐる議論の多様性を確保しつつ、AIの便益の増進及びリスクの抑制について、認識の共有を図り、相互に協力するよう努めることが期待される。
3. 標準化団体等は、本ガイドラインに適う**推奨モデルを作成し公表**することが期待される。
4. 各国政府は、**AIの開発者コミュニティを支援**し、本ガイドラインに掲げるAIのもたらす便益の増進やリスクの抑制といった課題の解決に向けた取組と併せ、**AIに関する研究開発を支援する政策**を積極的に推進することが期待される。

AI ネットワーク化と智連社会

AIネットワーク化

1. AIシステムが、他のAIシステムとは連携せずに、インターネットその他の情報通信ネットワークを介して単独で機能。

2. 複数のAIシステム相互間のネットワークが形成され、ネットワーク上のAIシステムが相互に連携して協調。

3. センサやアクチュエータを構成要素として含むAIネットワークが人間の身体又は脳と連携することを通じて、人間の潜在的な能力が拡張。

4. 人間とAIネットワークが共生し、人間社会のあらゆる場面においてシームレスに連携。

智連社会

智連社会 (Wisdom Network Society 【WINS】)は、人間がAIネットワークと共生し、データ・情報・知識を自由かつ安全に創造・流通・連結して「智のネットワーク」(Wisdom Network)を形成することにより、あらゆる分野におけるヒト・モノ・コト相互間の空間を越えた協調が進展し、もって創造的かつ活力ある発展が可能となる人間中心の社会像。

人間がAIネットワークと共生し

データ・情報・知識を自由かつ安全に創造・流通・連結して「智のネットワーク」を形成することにより

あらゆる分野におけるヒト・モノ・コト相互間の空間を越えた協調が進展し

人機共生

総智連環

協調遍在

創造的かつ活力ある発展が可能となる社会