

第7講 未来への飛翔 – 生成A I、自動運転、 A I 倫理が紡ぐ社会の未来

【学習到達目標】

- ・日本のA I 戦略は教育から始まることを説明できる。
- ・教育に利用される生成A I や自動運転について、事例を挙げて説明できる。
- ・国際的なA I のルール作りとA I 倫理について考えることができる。

1. 日本のA I 戦略

政府見解として公開されたのが、2019年4月で「我が国のA I 原則/A I 戦略について」（平井卓也・内閣府特命担当大臣）です。HirAI PitchにおけるA I 関連テーマの例として「人材育成」、「社会実装」、「データ戦略」などがあり、A I 戦略に関連する幅広いテーマについて産官学の関係者が意見交換いたしました。

具体的には、「A I の中小企業向けサービス」、「A I 検定試験」、「データ流通市場」や「衛星写真を組み合わせた都市の3Dマップ作成技術」がA I 関連テーマの例です。今日でいうデジタルツインというようなものが当時出てきたわけです。

・「A I 戦略」の概要

2022年のA I 戦略では、「人間尊重」、「多様性」と「持続可能」の3つの理念と、ソサイアティ5.0を実現しSDGsに貢献するということを目指しています。この3つの理念の実装念頭において、5つの戦略目標、「人材」、「産業競争力」、「技術体系」及び「国際」、それに加えて5つ目として、戦略目標0ということで「差し迫った危機への対処」というのを設定するようになりました。教育と関係するのは、「未来への基盤作り」の中の

「教育改革」です。

「産業・社会の基盤作り」項目では、「社会実装」が非常に重要されています。A I戦略2022での社会実装というのは、「他国の先進的な事例」との比較に基づいて、新たな目標を設定して推進するということです。

・人間中心のA I社会原則の概要

A I戦略2019で言う「人間中心のA I社会原則の概要」では、「人間中心の社会」を中心に周りに「人間の尊厳」、「多様性」と「持続可能性」というものがあります。

「人間中心の原則」というのは、A Iは人間の能力を拡張するが、A I利用に関わる最終判断は人が行うものです。教育と関係するのは「教育・リテラシーの原則」ですが、リテラシーを育む教育環境を全ての人々に平等に提供するとなっております。その他にも、「公正競争確保の原則」、「公平性、説明責任及び透明性の原則」、「イノベーションの原則」、「セキュリティ確保の原則」や「プライバシー確保の原則」があります。

・A I戦略の進捗状況と新たな戦略

A I戦略のこれまでの進捗状況と新たな戦略の必要性を示したものです。

A I戦略の主な成果の第1番目は、「教育改革」です。教育改革の1つ目は、「数理、データサイエンス、A I教育プログラム認定制度」が開始されたことです。これはリテラシーレベル（読み書きレベル）ですが、2021年までに78件を認定するという成果を挙げています。

教育改革の2つ目は、皆さんも関わっておられる小学校、中学校における「G I G Aスクール構想」で、1人1台端末を前倒して実現したということです。

他にもあの「人工知能研究開発ネットワーク」を設立したこととか、スマート農業実証プロジェクトとかが成果として現れてきました。

進捗としては、171件中154件、90%が計画通りに進捗しています。

国外においては、A I技術は国家安全保障、民主主義保全などの社会の根本機能維持の必須技術となっています。

例えば、アメリカでは国家安全保障の観点から、A I予算（非国防）ということで年間320億ドル、プロ野球の大谷選手の契約金1000億円超の約50倍を、国家予算として軍事でない面だけでこれだけ投入しています。

一方、中国では、軍事に使うということで5か年計画を発表しています。他方、EUはAI利用に関する包括規制案を公表しています。

日本では、これまでの「社会・経済システムの変革」に加え、地震とか水害とか火山噴火という「大規模災害」への備えや、コビッド19のような「パンデミックへの備え」に対応するような新たなAI戦略が要求されます。

・ 国家規模の危機への対処

「国家規模の危機への対処」ということで、日本の場合は少子化ということで人口減少に伴う、我が国の体力の低下とか、デジタル化の遅れがあります。これまでの閉塞を破る起爆剤として、AIを大きく活用すべきというように捉えられているわけです。それには、2つの流れがあって、1つは最大速度でのデジタル化・AI化ということです。基幹インフラのAI化前提でデジタルツイン、特に災害の場合に実際の街で水害とか火山の噴火できないので、デジタルツインの上でシミュレーションするというものです。

例えば、渋谷の街をそのまま再現して、デジタルで作った世界で、例えば人口がどう動くとか？を見るわけです。

それからもう1つの流れは、強靱な社会システムへの転換ということで、グローバルネットワークの強化ということがあります。

・ 地球規模の危機への対処/強靱な基盤作り

規模の危機への対処」と「強靱な基盤作り」です。強靱な基盤作りということで、「責任あるAIの概念」を構築する「説明可能なAI」等の技術は、情報基盤の信頼性を担保します。そして、我が国は高品種と安心安全という競争上の利点を生み出し、「責任あるAI」に向けた取り組み行っていくことになります。地球規模の危機への対処の方は、SDGsで言うサステイナビリティ分野でのAIの応用を考えていくというわけです。

・ AIに関する「思い込み」

AIの社会実装の推進に臨む姿勢には、3つの思い込みを捨てる必要があります。

1つ目は「AIは人の仕事を代替する？」という思い込みです。AIが我々の仕事を奪うのではと言うのは、思い込みであって、実はAIは人と協調します。人がAIと協調することで労力を最小化して利益を最大化することができます。

2030年雇用の大崩壊があり、AIが人の仕事を奪うと言われています。

澤井は、日本の少子化を解決するため、むしろA Iをどんどん使い込んで、生産技術では日本が世界で一番優れているとなったなら、日本の国際競争力が維持できると考えています。そんな訳で、人はA Iと協調していく方が良いのではないかと思います。

2つ目の思い込みは「技術者だけがA Iを深く理解できる？」です。これも思い込みです。最近のチャットGPTとかBARDとかBING AI等を見ると、「結構だれでも使える」感じですが、ビジネスケースからA Iは理解できるということです。自らA Iを構築しなくても、既存のA Iを利用し、他の部分で差別化していくことも一つの有効な手段です。

3つ目は、「データが全て？」です。今のA Iはデータがなければ何も出来ないという思い込みがあります。実はそうではなく、使いながらどんどん増やしていけば良いというのが「ループの形成」です。A Iによるサービスの提供を通じてデータを取得するというようなループを形成し、サービスの提供とデータの収集を同時に行うことが重要です。

・A Iデータの課題

「データ」について、A I戦略2019に出た話ですが、ノイズや偏ったデータによってはA Iが信頼できる結果を出すことができない可能性があります。例えば採用試験を、A I採用システムがこれまでずっと過去の履歴では男性を取ることが多いとした場合、学習データのバイアスによって、女性を不適切な判断で不採用にするのではないかと心配です。

また、少しのノイズでA Iシステムが誤認識してしまいます。例えば最高時速50と書いてある標識が泥で汚れてしまうと、速度制限50以上出しても良いと読み間違えてしまうというような誤動作が発生します。

その意味でA I製品のサービスの信頼性を担保する仕組みが必要です。その1つは国際標準化で、いま1つは第三者評価プロセス、ガイドラインの整備です。

A I戦略2022では社会実装に向けた4つの取り組みがあります。

ディープラーニングを重要分野として位置づけ、企業による実装を念頭に置き、「A Iの信頼性の向上」、「人材確保等の環境整備」、「A I利活用データの充実」と「政府におけるA I利活用の推進」に取り組みます。そして、物理、科学、機械など日本が強みを有する分野とA Iの融合により、競争力の高い製品やサービスを生み出すことが今後大事になってきます。

・日本のA I戦略のまとめ

日本のA I戦略は3つにまとめることができます。1つは「人間中心の社会原則の必要性」です。2つ目は、「A I社会原則を踏まえたA I社会実装」で人間中心という意味を真に理解した上で実現すべき未来社会に向かうためのA I社会実装を推進します。3つ目は、「我が国としての方向性」です。

人間中心に基づく健全なA I社会実装で世界に先行し、広島サミットの宣言をおこなっていることです。2023年末A I戦略会議にて生成A Iに関する暫定的な論理を整理しています。これがまとまると次のA I戦略になるというわけです。

2 生成A Iと教育

・A Iの1丁目1番は教育

A I戦略2019では、「A Iの1丁目1番は教育」ですと言っています。登壇者は、座長の安西祐一郎・慶応義塾大学元学長、北野ソニー・コンピュータサイエンス研究所社長と内閣府の専門官です。

A I戦略2019での教育改革に関する主な取り組みです。狙いは、デジタル社会の「読み・書き・そろばん」である「数理・データサイエンス・A I」の基礎などの必要な力を全国民が育み、あらゆる分野で活躍できるようにすることです。

・教育改革に関する主な取り組み

G I G Aスクール構想の前倒しで生徒1人1台端末もI C T環境これは実現されたので、今後小中学校では4校に1人以上、高校では1校に1人以上の多様なI C T人材を登用することになります。

「学習内容の強化」としては、高校におけるA Iの基礎となる実習授業の充実で、大学でのM O O C等を活用した標準カリキュラムの開発と展開です。2025年の育成目標は、小中学生全員と高校卒業生全員ということで毎年100万人にしていこう、大学・高専卒業者は全員ということで毎年50万人、合わせて150万です。

「応用基礎」は、高校の一部と大学の50%で2025年には25万人育成していこうというわけです。「認定制度・資格の活用」では、「大学等の優れた教育プログラムを政府が認定します。

「制度構築」し、全学生が受けれる講座を作ることです。この講座は著者だけじゃなくて何人かの先生方で一緒にすることです。更に「A I応用力の習得」ということで、A Iと専門分野、例えば文化系の文化創造であれば、A Iと文化創造とのダブルメジャーになるような人の育成を促進していくことになっています。

「先鋭的な人材を発掘・伸ばす環境整備」ということでエキスパートを2025年には2000人を育成する状況にして行きたいということです。

・生成A I / 教育

2022年11月にオープンA Iが、ご存知のようにC h a t G P Tを発表いたしました。爆発的な広がりを持って、2ヶ月で、1億人に達しました。C h a t G P Tは、事前学習した大規模なニューラルネットワーク（大規模言語モデル）を対話形式で操作する方法を採用し、「生成A I」として世界的に話題を呼びました。

生成A Iとは、全く新しいオリジナルのアウトプットを生み出すA Iのことで、具体的には新しいデジタルの画像／動画、音声／音楽、文章／コード等を生成するA I、もしくはこれらを組み合わせて生成するA Iのことを指す、とされています。

生成A Iは、非常に高性能で、コンピュータが我々の住む世界のことを良く理解し、人間と自然な対話（チャット）ができるような新しい能力を獲得しているかのように振る舞います。

生成A Iは、数億パラメータ以上の膨大な知識量を有し、答えを瞬時に提示することを得意としており、正解を提示する能力では人間を超えています。生成A Iは使い次第で人間の創造的な作業を支援します。

・未来への基盤作り

教育改革に関する主な取り組みと、研究開発における主な取り組みです。「初等中等教育における環境整備」は、「1人1台端末」で実現されました。

リテラシーとして高等学校でのA Iの基礎となる実習授業の充実と、大学等の優れた教育プログラムを政府が認定する制度で、毎年50万人育成します。

応用基礎で毎年25万人、エキスパートと言われる人材は毎年2000人育成していく計画です。

研究開発における主な取り組みでは、「人工知能研究開発ネットワークを理研、産総研とN I C Tを中核にし、大学、海外の研究機関や国の研究機関とつぐネットワークを充実していきましょう」という計画です。

・生成A Iに関する暫定的な論点整理

以下の論点整理は、最近の技術の急激な変化や広島A Iプロセスを踏まえて、A I戦略会議構成員がA I関連の論点を整理したものです。

○生成A Iの可能性

生成A Iの登場は、内燃機関の発明・IT革命と同じく、幅広く生活の質を向上させる「歴史の転機」となる可能性があります。生産性の向上・情報アクセスの改善など、諸課題の解消も期待されています。

○生成A Iと日本の親和性

我が国は、①研究・技術水準の高さ、②ロボット・A Iへの肯定的イメージ、③労働人口急減、④デジタル化への高いニーズ、⑤きめこまやかさ・創造性など、生成A Iとの親和性が高く、大きなチャンスです。アトムとかドラえもんとかロボットA Iへの肯定的イメージが強くて、そこへ持ってきて、少子化による労働人口の急減があり、デジタル化の必要性が出てきています。きめ細やかさ創造性など、日本人の特性ともあって、生成A Iとの信和性が大きく、大きなチャンスになる可能性があるというわけです。

○いま戦略を検討することの重要性

我が国に、A Iの勃興とともに再び成長の機運が見えており、諸外国の後塵を拝さないよう、今こそ大胆な戦略が必要です。政府は、人々がA Iがもたらす社会変化に対して安心感を持ち、各プレイヤー、つまり各産業界の人たちが予見可能性を持てるようリスクに対応すべきという訳です。また、企業・研究者が存分に活動できるインフラ整備を行うべきです。例えばスパコン、A I用のスパコン整えるとかそういうような話です。

○懸念されるリスクの具体例と対応

懸念されるリスクの具体例と対応については、次の6点があります：1) 機密情報の漏洩や個人情報の不適正な利用のリスク、2) 犯罪の巧妙化・容易化につながるリスク、3) 偽情報などが社会を不安定化・混乱させるリスク、4) サイバー攻撃が巧妙化するリスク、5) 教育現場における生成A Iの扱い、6) 著作権侵害のリスク、7) A Iによって失業者が増えています。

澤井は、日本の場合むしろチャンスにすべきで、A Iによって生産力を上げ、逆に職業が増えていくという風にすべきだと思っています。

○A Iの利用

A Iの利用については、次の4点が考えられます：

1) 生成A Iは、デジタル化を加速させ、我が国全体の生産性向上のみならず、様々な社会課題解決に資する可能性があります。2) A I利用を加速するため、医療や介護・行政・教育・金融・製造等のデータ連携基盤の構築・D F F T (Data Free Flow with Trust) 構想の具体化・人材育成・スタートアップの事業環境整備を進めるべきです。3) 政府機関が一体となって、機密情報漏洩のリスクなどに配慮しつつ、率先して生成A Iの利用可能性を追求することが重要です。4) 幅広い世代が生成A Iの恩恵を享受できるよう、スキル・リテラシーを身に付けることが大切です。

A Iの利用は、デジタル化を加速し、我国全体の生産性向上のみならず様々な社会課題の解決する可能性があるので、教育に関しては幅広い世代が生成A Iの恩恵を享受できるスキル・リテラシーを身につけることは大切ではないかということです。

・生成A Iを活用した授業の実践例

【事例1】千葉県印西市立原山小学校（松本博幸校長）の生成A Iを活用した授業の実践例です。授業では「文章生成A Iの仕組みを知り、これからどのようにかかわっていけばよいか考える」ことを狙いとし学習活動に取り組んでいます。保護者の理解及び協力、そして専門的な知見を有する外部人材の助言によって、全国的に見ても先進的な生成A Iについての授業が実践されていました。

この授業実践例では、専門家の「N P O法人 みんなのコード」が「対話A Iは、ネット上にある大量のデータを使って、それをまずA Iに学習させて、人間の質問や指示に答えるようにするのだよ。悪いことに使わなければうまく使うことができるよ。A Iの言っていることは100%正しいということはないということを理解して使うのだよ。」とわかりやすく教えています。

【事例2】同志社中学校・高等学校の反田 任先生のA I発音チェックアプリ“ELSA Speak”を活用した英語授業デザインの事例です。目的はELSA Speakを活用し「話すこと（発話）」の力を向上させることです。ELSA Speakでは英語の音素単位で発音判定をするため、単語や英文の正確な発音が身に付きます。

・生成 A I が脅威的な進化

生成 A I が脅威的な進化をしています。最近「A I アライアンス」が出てきました。A I アライアンスとは、Microsoft の「being A I」、ChatGPT、GPT4 や GPT4 ターボに対抗する国際的な連合グループです。

・雇用の未来

オズボーン氏の論文「雇用の未来」の中では、「2030年には、現在ある職業の47%がなくなる」と言われています。そのため、A I 時代を生き抜く人材を育成する教育が重要になる。消える職業、なくなる仕事で筆頭に上がっているのは銀行の融資担当者とかスポーツの審判、不動産のブローカー、レストランの案内係や、保険の審査担当者というようなことが上がっています。

3 自動運転

自動車産業領域では、ディープラーニングが強みを持つ「認識」領域を中心に、数年前から様々な実証実験が行われています。安全性を最重視しつつ、将来の実用化に向けた取り組みが進められています。

・自動運転

カメラ・センサーなどから「信号・道路標識・障害物」等の外部情報を把握し、アクセスやブレーキ、ステアリング等の各種操作につなげていく自動運転の工程内で、ディープラーニング技術が活用されているわけです。

実際の動きを見てみますと、車の上にセンサーが付いていて、他にもセンサーついていると思うのですが道路の障害物とかを認識して動いています。

内閣官房 I T 総合戦略室等では、高遠道路でのトラックの隊列走行を早ければ2022年に商業化することを目指しています。高度自動運転の市場化・サービス化にあたり、様々な走行環境における実証実験の実施」が不可欠となって来ました。

成果としては、まず自家用車の分野では2021年3月にホンダより、自動運転レベル3に適合する技術であるトラフィックジャムパイロット（渋滞運転機能）を搭載した車両が市場化されました。

ホンダの自動運転レベル3は、高速道路渋滞時など一定の条件化で、システムがドライバーに代わって運転操作を行うことが可能になりました。

このことにより、2020年の実現を目指していた高速道路での自動運転レベル3が実現されることになりました。また「運転支援システムの高度化」について、高精度3次元地図を利用した速道路でのハンズオフの運転支援が日産などで実現されました。

自動運転の国内外の動向を見てみましょう。

アメリカの場合、民間主導でレベル4を実用化しています。具体的には、IT企業が自動車を時期デジタル領域と見据えて、無人モビリティサービスの事業開発と促進しています。GOOGLEなどのIT企業を中心に、WaymoタクシーやGakic無人配送トラックといったレベル4サービスカーが一部公道で実用化しています。

欧州では、欧州OEMのディーゼル不正やHEV技術力不足を補う競争力獲得のため、自動化・電動化を政府指導で推進しています。官民ともに、レベル2のとかレベル3が中心で、官主導でレベル4の試みが行われています。レベル3までは車の運転を基本的に人間が主導で行っています。レベル4以降では基本的に人間はアシストしません。レベル5に至っては人間のアシストが全くありません。今見る限りレベル5は、未だ出てきていません。

中国の場合、2025年までに製造強国の仲間入りを、政府目標としており自動化技術開発を政府指導で推進しています。

新興国では、民間主導でレベル4を目指しており、公共交通の未整備地やスマートシティ内の移動手段確保を目的に、民間を中心に自動運転を実証しています。

日本における自動運転の導入地域とAD/ADAS市場の拡大です。ADは「自動運転システム (Autonomous driving system)」で、ADASは「先進運転支援システム (Advanced Driver-Assistance Systems)」です。

自家用車については2022年にレベル1からレベル2までに普及しています。

2025年には運転支援の進化により、レベル2からレベル3になりAD/ADAS導入が促進されます。そして高度な運転支援車とサービスカーの基盤を共用し、AD/ADAS ECU、センサーやアルゴリズムなどの量産効果によりコスト低減されたシステムの活用を通して、商用車レベル4サービスカーの本格普及に寄与します。

商用車については2022年にレベル4のサービスカーを導入し、自家用車との協調領域における自動運転開発を行いながら、2025年にはレベル4サービスカーの導入を40箇所を増やし促進して行きます。2030年以降に自家用車の技術を活用しながら、レベル4サービスカーの本格普及を目標にしています。

自家用車のレベル4自動運転サービスの実現は、限定空間の遠隔監視のみで自動運転サービスの実現に向けた実証事業の推進という形で行われています。事業性向上に向けて、4台までの車両を1人が同時監視します。2023年5月21日より、福井県永平寺町において、レベル4での自動運転移動サービスが開始されました。西村経済産業大臣（当時）が乗って現地で記念式典に臨んでいます。町民の方々による乗車や、遠隔監視室でレベル4運行を監視して様子が見れます。遠隔監視室では、1人の人が3台の自動運転を監視し、操作しています。

・タクシー需要予測

タクシーの需要予測にもAIが使われるようになりました。タクシードライバーが抱える一般的な課題に、「需要がある場所や時間帯を予測するのが困難」というものがあります。これから紹介するケースでは、タクシー需要を予測するために各種データ(リアルタイム人口統計データ、気象データ、過去のタクシー運行データ等)を活用して、どこに車を配車したら1番効率がいいか正解を予測する本事例では、「抽象的で複雑な特徴を獲得する」ことを狙い、オートエンコーダ系技術を採用しています。オートエンコーダとは、3層ニューラルネットにおいて、入力層と出力層に同じデータを用いて教師あり学習させたものです。

【事例】株式会社NTTドコモの「AIタクシー」は現在から30分後までの未来のタクシー乗車需要の予測結果などのデータをオンラインで配信するサービスです。需要を予測するために、発展させたオートエンコーダ(Stacked denoising Autoencoder、SdA)を使用しています。タクシー需要を92.9%の精度で予測しました。実証実験を経て、2018年2月よりサービス提供開始し2022年6月に終了しました。

・未来の自動運転

2023年8月13日、米国ミシガン州政府は、デトロイトからアナーバーの間に全長40マイル(約64km)に及ぶ米国初の自動走行レーンを建設すると発表しました。レベル5の自動運転は難しすぎるというので、アメリカでは道路になんか仕掛けをしようというようになってきているようです。デトロイトというとアメリカの自動車の発祥の地なので、デトロイトで行うことが将来の自動運転に相当のインパクトを与えるわけです。

これまでGoogleは、インフラに依存しない「自立型」の自動走行の実現を目指してきたことで有名です。

2000万マイル（約3200万km）以上の公道走行を続けながら開発してきた技術です。2016年に設立された自動運転専門会社「Waymo」が、アリゾナ州で提供している有料の自動走行ライドシェアサービスWaymo Oneのいずれも自立型の自動走行を前提としてきました。

グーグルは、関連会社Caviumがインフラ協調の必要性を主張したことにより自動運転の考えが変化しつつあるようです。

オレンジに引いた線のところが、自動運転車レベル5が走るイメージです。混雑している道路の中でオレンジのところだけは自動運転で走っています。

4 国際的なAIのルール作り

上図は、国際的なAIのルール作りです。政府は「AI戦略2022」や「人間中心のAI社会原則」などをこれまで定めてきました。

論点整理は、従来の基本戦略・理念は維持しつつ、急激な技術変化やG7広島サミットで合意されたビジョンと目標（「我々が共有する民主的課に沿った、信頼できるAI」）等を踏まえ、AI戦略会議の構成員が有識者として、生成AIを中心に課題と方向性などを整理しています。

・広島プロセス

広島プロセスには3つの動きがあります。

- 1) 2023年5月G7広島サミットにおいて岸田総理が提唱したことです。
- 2) 広島AIプロセスのG7デジタル技術・閣僚声明です。広島AIプロセス包括的政策枠組みです。①生成AIに関するG7の理解に向けたOECDレポート、②全てのAI関係者に向け及び高度な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、③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、④プロジェクトベースの協力の4つあります。広島AIプロセスを更に前進させるための作業計画の策定です。
- 3) 2023年12月6日G7首脳による最終合意で、AIアラスシンポジウム2023で公表された指針です。①全てのAI関係者向けの広島プロセス国際指針として、導入前及び市場投入前を含め適切な措置を講じます。②導入後の話です。③十分な透明性の確保を支援します。④関連組織間での責任ある情報共有します。⑤AIガバナンス及びリスクの管理方針を策定し、実施し、開示します。⑥強固なセキュリティ管理に投資し実施します。

⑦ユーザーがAIを生成したコンテンツを識別できるようにします。⑧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ます。⑨世界の最大の課題、特に気候危機、世界保護、教育等に対するために、高度なAIシステムの開発を優先します。⑩国際的な技術企画の開発を推進することです。⑪適切なデータインプット対策を実現し、それによって個人データ及び知的財産を保護します。⑫高度なAIシステムの信頼でき責任ある利用を促進し、デジタルリテラシー、訓練及び認識を向上させる機会を求めるべきであるということです。

・リスクの対応として

中央省庁が生成AIを利用する場合の手続きで、生成AIの業務利用に関します。2023年9月オープンAIに対する注意喚起がありました。

AIと知的財産権とAIと著作権問題に関して、研究会・審議会における議論が行われています。教育とも絡む面では学習用データの整備で、政府等保有データの提供促進があります。加えて、新たなモデルの開発、基礎的・先進的な研究、人材育成があります。注目したい点は次の3つあります。

1) ルールメイキング+エンフォースメント、2) 大規模汎用モデル+多様なモデル、ファインチューニング、3) 一部の地域での利用+世界各国での本格的な利用というものがあります。今後、安全性・信頼性、オープンクローズ、多様性/国際整合性の議論がますます重要になります。

5 AI倫理

現在の第4次AIブームでは、自動運転や画像診断など私たちの暮らしにAI技術が急速に入り込んで来ています。21世紀の基幹テクノロジーとされるAIとどう付き合い、その活用をどこまで許容していくのか？EUではAI倫理に基づく輸入規制を計画しており、日本のAI倫理が問われています。

・AI倫理の背景

2023年12月17日OpenAIの「超知性」誕生に備える研究チームがGPT-2（弱いAI）モデルでGPT-4のように強力なAI（強いAI）を制御する方法を説明しました。OpenAIは、人間よりもはるかに賢いAIである「超知性」が2033年までの10年間で開発されると推測しており、「スーパーアライメントチーム」を立ち上げ、超知性を制御するための研究が行われています。

AIの賢さを下回る人間ではAIの監視が困難になります。OpenAIのスーパーアライメントチームは、人間が超知性を適切に監視できるかを見る代わりに、大規模言語モデルのGPT-2がより強力なGPT-4を監督できるかテストを行いました。

注) GPT-2のパラメータは15億程度であるのに対し、GPT-4のパラメータは約1760億に上るといいます (IEEE Spectrum)。

・ AI倫理の定義

倫理とは、Websterによれば「a system of moral principle」となっており、AI倫理は「a system of moral principle for using AI」と定義できます。日本の文部科学省が推進する全国の児童・生徒1人に1台のコンピュータと高速ネットワークを整備する「GIGAスクール構想」に必要な「IoEEE」（倫理のインターネット、教育のインターネット、生きる力のインターネット）AI倫理チャットボット機能を試作・検証することにあります。

・ 1人1台端末に必要なAI倫理チャットボット

AI戦略の教育改革「1人1台端末」、GIGAスクール構想の最大の課題は「チャットによるいじめ問題」と言われています。AI倫理を探究するうちに「チャットによるいじめ問題」にAI倫理チャットボットが使えるのではないかと考えられています。

近年、人間と会話をすることができる対話システムへの注目が集まっています。例えば、音楽の再生やメールの確認などを行う Google Assistant や Siri、また顧客からの問い合わせ対応を代替するチャットボット（チャットのロボット）といった、何らかのタスク達成を目的としたタスク指向型対話システムが広く浸透してきています。今日、自動運転や画像診断など私たちの暮らしにAI技術が急速に入り込んで来ています。21世紀の基幹テクノロジーとされるAIとどう付き合い、その活用をどこまで許容していくのか？

「AI倫理」とでも呼ぶべき社会規範をきちんと議論しなくてはならないと言われています。

ウィーン大学の哲学者クーケルバーク氏は著書「AI倫理」で、AIを使うための「運転免許証」がない。と警鐘を鳴らしています。

ノーベル文学賞受賞者のカズオイシグロ氏は「クララとお日さま」と言う最新の小説の中で、人工知能を搭載したロボット「クララ」を登場させています。クララは、観察と学習への意欲と理解力を持つに至り、人間社会で生きていく力「生きる力」(Energy of Life)を得るようになります。

・教師あり学習

機械に学習させる「機械学習」には「教師あり学習」、「教師なし学習」、と「強化学習」の三つの学習の枠組みがある。図5は人間の脳のニューロンが層状に接続した構造を模擬した機械学習の三つの枠組みです。

教師あり学習とは主に人間の小脳が担う学習機能で、代表的な統計手法は回帰と分類です。学習者に対し、教師が明示的に正解を教えたり、学習者の誤りを指摘したりすることで、学習者が正しい解を得ることを助けます。正しい入出力の組合せを与えて学習することで、新規の入力に対し適切に出力します。

「回帰」の代表的手法は誤差逆伝播法(Back Propagation)です。「分類」の手法として、正解、若しくは誤りを入力として、未経験入力に対する意志を決定する決定木(Decision Tree)や決定表(Decision Table)の作成などがあります。本研究では EXCEL 上の決定表で「倫理表」を試作しました。

・「A I 倫理」処理システムの試作

「教師あり学習」A I を使い、社会規範・倫理と、設計者の故意ではないA I の誤認識（機能不全、誤作動や機能低下を含む）を検証し、適切な処理を行う

「I o E E E」（Internet of Ethics, Internet of Education, Internet of Energy of Life）A I 倫理チャットボット機能の試作を行いました。

本研究では、「教師あり学習」を使い、教育禁止用語や放送禁止用語等のような社会規範・倫理とA I の誤認識が処理・説明できるシステム作りを目指しました。入力はA I 音声入力でもキーボード入力でもできます。デープラーニングによるA I 音声入力はiPhoneで行い、リモートマウスで接続したパソコン上でA I 倫理処理を行いました。「A I 音声入力では何故誤認識したか？」は言葉では説明できない、つまり暗黙知です。

社会規範・倫理とA I の誤認識の検出・修正（言換え）処理は EXCEL の VBA プログラムで瞬時に処理され、修正した音声入力文と修正理由を説明した説明文はそれぞれ EXCEL ファイルに保存される。学習データは、社会規範・倫理例、A I の誤認識と学習済みの TensorFlow.js モデル・デーモン（システム）等で、インターネットとブロックチェーンで参照します。

例えば入力文「Slave is a bad word」（奴隷は良くない言葉です）を入力すると、正しい表現に言い換え、その理由を説明する社会規範・倫理例の放送禁止用語は教育禁止用語としてウェブ検索する具体的にはアイヌ系からロンパリに始まりブスとかチビといった誹謗中傷の類からジョンやアメ公といった人種差別用語まで教育上使わない方が良いとかがえられる用語は網羅されています。

・「教師あり学習」モデルを使った検証

音声入力文に、①アイデンティティベースの憎悪、②侮辱、③わいせつ、④重度の毒性、⑤性的に露骨、⑥脅威、⑦毒性などの有毒なコンテンツが含まれているかどうかを、約 200 万件を事前に「教師あり学習」した学習済みの TensorFlow.js モデル・デーモンを使い検出しグラフ化し、「I o E E E」AI 倫理チャットボット機能の検証を行いました。例えば G I G A 端末の入力文「馬鹿!消えてしまえ!」を、学習済み TensorFlow.js モデル・デーモンに入力し分類すると、②侮辱)かつ⑦毒性が「TRUE (きわめて有害)」、及び①アイデンティティ攻撃、③卑猥、④重度の毒性、⑤性的な露骨及び、⑥威嚇は「FALSE (無害)」と分類します。

G I G A 端末でトラブルを引き起こす入力文例 37 件中 26 件(10%)が TRUE (きわめて有害)、53 件(20%)が NULL (要注意)、残りは FALSE (無害)と検出できました。「I o E E E」AI 倫理チャットボット機能は、30%を超える抽出率で、倫理テーブルで確認できなかった未定義の誹謗中傷を検出できました。本チャットボット機能は非常に効果的でした。

結果、G I G A 端末やケータイに人工知能を搭載した人間に親切な「I o E E E」AI 倫理チャットボット機能を搭載し、1) 社会規範・倫理と AI の誤認識の修正処理を行い、その後 2) 「教師あり学習」モデルを使った検証を行うと、抜けが少ない有効な AI 倫理処理ができると分かりました。

課題

自動運転車が引き起した事故は、誰が責任を負うべきか？について考察し、あなたの考えを 800 字で説明しなさい。