

さきがけ
「信頼されるAIの基盤技術」
研究総括説明

2022年4月27日

研究総括 有村 博紀(北海道大学・教授)



科学技術振興機構

戦略目標（概要）

信頼されるAIの基盤技術

- 人工知能(AI)技術を用いたシステムやサービス等のAI技術が社会の中に浸透しつつある。現在、「人間中心のAI社会」*1の実現のために「信頼される高品質AI」*2(Trusted Quality AI)の実現が課題。
- 一方で、深層学習に代表される現在のAI技術がもつさまざまな弱点や限界、社会との整合のために不足する要件が明らかに（脆弱性、予測のブラックボックス性、説明性や納得性、不公平性等）。
- そのため、現在のAIの限界を超えるAI技術そのものを革新的に発展させ、社会からの要請に応え得る根本的な信頼性確保を実現するための研究開発が必要となっている。

戦略目標（達成目標）

人間中心のAI社会の実現のため「信頼される高品質なAI」の創出に向けた研究開発の推進をはかる

(1) 現在のAI技術の限界を克服する新技術の創出

(2) AIシステムの信頼性・安全性を確保する技術の創出

(3) データの信頼性確保と意思決定・合意形成支援技術の創出

次世代の人工知能技術

基本的な問い：21世紀の社会で、
コンピュータが人間と働いていくために
必要な人工知能技術とは何か？

意思決定
タスク

法律

評価タスク

試験

運用タスク

生産

難しいタスク

医療

自動運転

さまざまな応用

介護

採用

社会インフラ

教育

社会的要請

- 社会的責任
- 人間との協働

- 説明可能性
- 信頼性・有意性
- 公平性・AIシステム

次世代の人工知能技術

- 現在「AI」としてどちらの方向で議論しているのか

2. “AIシステム”

シミュレーション?

制御システム?

- 複雑巨大システムへの機械学習アルゴリズムの組み込み
- 例: 生産プラント、交通インフラ、社会インフラ等

1. 機械学習アルゴリズム

2018年現在、盛んに研究されている

3. “汎用AI”

機械学習アルゴリズムを中核として、人間の知能に匹敵する汎用AIの実現を目指す

制御システムの専門家へのインタビュー







- 従来の複雑システムの場合。
- 理論に基づく制御コンポーネントを、非線形のシステムに組み込んで運用する場合がある。
- その場合には、制御理論と現場の経験を総合して、システムを設計し、運用する。

例: The Partnership on AI



- 機械学習技術の社会的受容を目指すコンソーシアム。
- 取り組むべき課題として次の6つをあげている

Our Work

-  安全・安心なAI (Safety-Critical AI)
-  公正性、透明性、責任有るAI
(Fair, Transparent, and Accountable AI)
-  AIと、労働者、経済
-  人々とAIシステムの協調
-  人々と社会へAIが与える影響
(Social and Societal Influences of AI)
-  AIと社会利益 (AI and Social Good*)

KDD 2017で同名のWorkshop

2016年創設。

創設メンバー

Amazon, Facebook, Google, DeepMind, Microsoft, IBM, Apple, ...

現在、連携先は

>50%非営利、
13カ国からの、
80団体

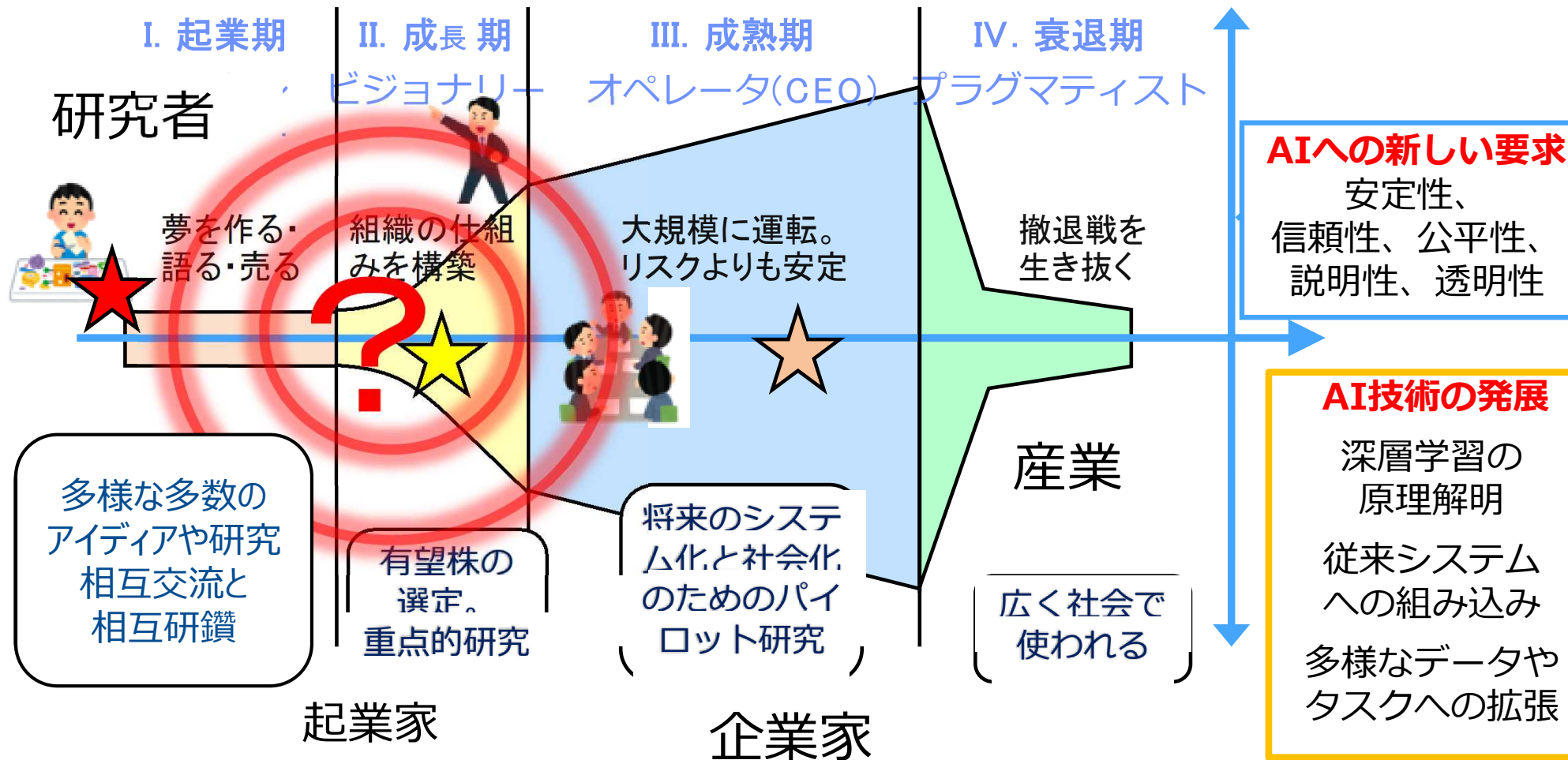
次世代の人工知能技術は今どこに？

- 現在のAI研究は、発展のどの段階か？

現在の人工知能技術

- 現在は爆発的な発展の変革期

人間中心のAI社会



(*) "The four: The Hidden DNA of Amazon, Apple, Facebook, and Google," Scott Galloway, 2017 (訳 渡会圭子、東洋経済新報社、2018)

さきがけ次世代IoT研究領域 (概要)

領域名: 信頼されるAIの基盤技術

「人間中心のAI社会」の実現には、社会における多様で重要なタスクをこなし、人間と共に働く「信頼されるAI」が必要。そのために、従来のAI技術の単なる延長でなく、現在のAI技術/AIシステムがもつ本質的な問題点に取り組み、新たな原理・技術・社会実装の発見・発明を通じて、それらを解決する「信頼されるAI」のための革新的なAI基盤技術の創生を目指す。さらに、社会のさまざまな局面で、重要なタスクを遂行するための「AIシステム」の構成原理と実現技術の確立を目指す。

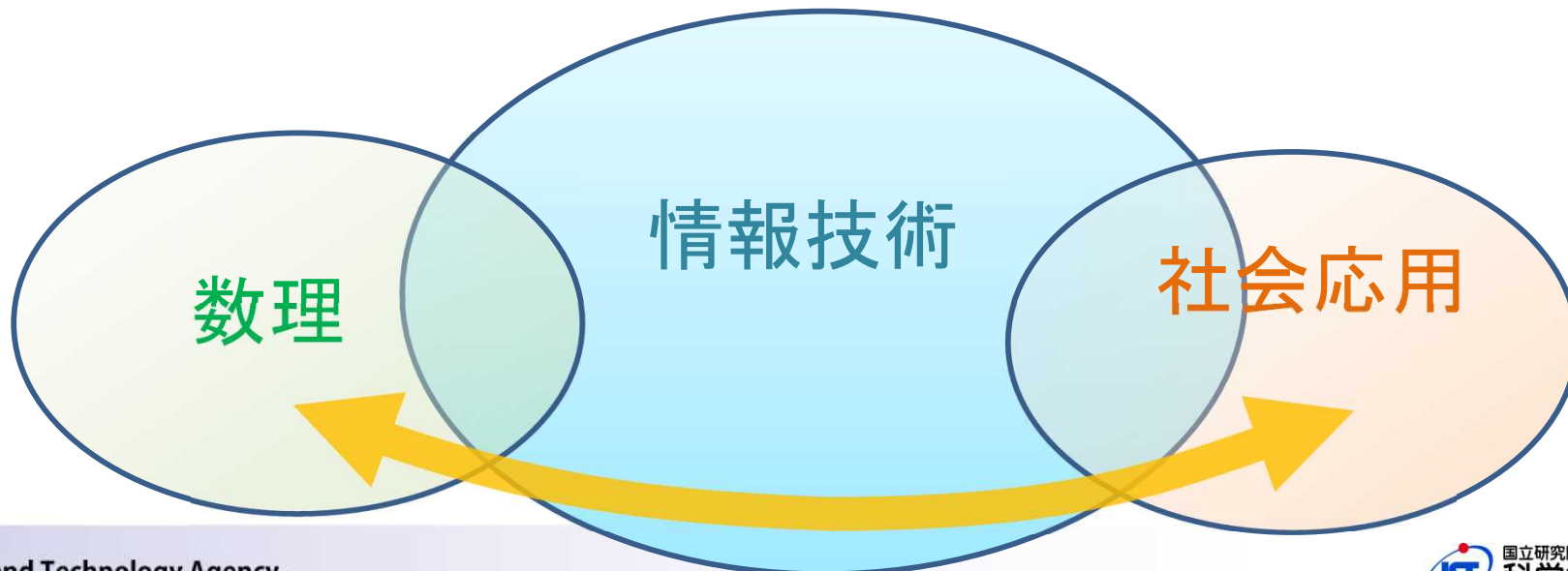
対象とする技術項目

- (1) 現在のAI技術の限界を克服するための新しい数理・計算・解析の基盤技術の研究、
- (2) AIシステムの信頼性・頑健性・透明性・公平性等、社会における新たなAI応用タスクの概念化・定式化、新しい構成原理・実現技術、(3) AIシステムを支えるデータや情報の信頼性保証や、AI技術を用いた意思決定や合意形成の技術、(5) AIシステムの設計・開発・運用の

方法論等

研究総括方針 (1/4)

数理・情報技術・社会応用等の多様な立場から、
人間中心のAI社会の実現に資する、
次世代AI技術と周辺技術に関して、
若手研究者による新しい発想に基づいた、
独創的かつ挑戦的な研究を推進する。



研究総括方針 (2/4)

■ 例1: 現在のAI技術の弱点の克服

提案募集する研究の例

- ・現在のAI技術の原理・限界の解明。新しいAIの構成原理。多様なデータやタスクを対象とする新しい数理モデルや、アルゴリズム、システムの提案。
- ・帰納的情報処理(例: 深層学習)と演繹的情報処理(例: 記号推論・知識処理・シミュレーション等)を融合したAIの研究。高次の機械学習タスク(因果関係の発見や、反事実説明、統計的有意な予測等)の研究
- ・人間の脳情報処理や認知に関する知見に基づく新しいAI原理。

これら従来の研究トピックにとどまらず、将来の人間中心のAI社会の基盤技術として、革新をもたらすようなあらゆる研究テーマ。

■ 例2: AIシステム・サービスの信頼性・安全性を確保するための研究

- ・AI技術を中核とする情報システム(以降、AIシステム)の設計・開発・運用や、品質評価、安全設計、説明性、透明性、公平性等に関わる新しい数理・方法論の研究。
- ・未知・想定外ケースや環境変化、データのバイアス、悪意ある攻撃等に対して頑健なAIシステム・サービスを実現するための研究

■ 例3: 社会におけるAIシステムのための周辺技術

- 多数の参加者・利害関係者が参加するAI技術を用いたサービスやシステムの数理・アルゴリズム・制度設計等の研究(例: マッチングや、オークション、クラウドソーシング等)。納得性や、公平性、頑健性等の望ましい性質をもった意思決定や合意形成の数理や技術。フェイク、データ改ざん等の不正(AIによるものを含む)を検知し、対処する技術の研究。

さきがけ「信頼されるAI」

1期・2期生 (R2・R3年度採択) 分野マッピング

1期生

2期生

データ信頼性・ 信憑性の問題

(フェイク・改竄、思考誘導等)

AI技術の信頼性・ 安全性の問題

(脆弱性・解釈性、差別・品質問題等)

現在のAIの 限界

(大量学習、常識、脳機構等)

五十嵐 歩美
(NII)「信頼される資源配分メカニズムの構築」

小林 泰介
(NAIST)「頑健性と安全性の性能限界を明らかにする深層強化学習」

西野 正彬
(NTT)「誤りがないことを保証する検証器つき機械学習の研究」

佐々木 勇和
(阪大)「グラフデータの説明可能なバイアスに関する基盤技術の創出」

松原 崇
(阪大)「望まれる性質を設計段階で保証する幾何学的深層学習の構築」

HOLLAND Matthew・Japmes
(阪大)「学習過程における価値観の多様化と性能保証の両立」

原 聡
(阪大)「機械学習モデルとユーザのコミュニケーション：モデルの説明と修正」

竹内 孝
(京大)「リアルな意思決定のための時空間因果推論モデルの研究」

飯塚 里志
(筑波大)「実応用に向けた動画画像コンテンツ加工のためのユーザ制御可能な例ベース深層学習フレームワークの確立」

谷中 瞳
(東大)「診療の意思決定を支援する言語・非言語時間情報検索」

横川 大輔
(東大)「化学的知見を生かした転送性の高い特徴量の抽出と利用」

藤井 慶輔
(名大)「生物集団移動の専門家が利用可能な説明・意思決定のための基盤技術」

吉井 和佳
(京大)「人とAIの同化に基づく能力拡張型音楽理解・創作基盤」

菅原 朔
(NII)「説明性の高い自然言語理解ベンチマークの構築」

小野 峻佑
(東工大)「センシングと知識発見の間に橋をかける数理的データ解析基盤」

栗田 修平
(理研AIP)「与えられた指示文章に従い言語で判断を説明するAI」

日高 昇平
(JAIST)「機械理解の創成に向けた随伴関手の統計的推定理論の構築」

大関 洋平
(東京大学)「認知・脳情報処理による人間らしい言語処理モデルの開発」

西田 知史
(NICT)「脳情報に基づいたAIの信頼性評価技術の開発」

岡田 謙介
(東京大学)「透明性の高い達成度テスト運用基盤の開発」

統計・数理

AI・機械学習

画像診断・処理

ロボティクス
・脳科学・心理学

科学・工学応用

自然言語処理

現在、20名のさきがけ研究者!

領域アドバイザー

氏名(敬称略)	所属	役職
石川 冬樹	情報・システム研究機構 国立情報学研究所 アーキテクチャ科学系	准教授
宇野 毅明	情報・システム研究機構 国立情報学研究所 情報学原理研究系	教授
浦本 直彦	(株)三菱ケミカルホールディングス データ&先端技術部	部長
大野 和則	東北大学 未来科学技術共同研究センター	教授
岡崎 直観	東京工業大学 情報理工学院	教授
鹿島 久嗣	京都大学 情報学研究科	教授
佐久間 淳	筑波大学 システム情報系	教授
櫻井 祐子	名古屋工業大学 大学院工学研究科	教授
佐倉 統	東京大学 大学院情報学環／理化学研究所 革新知能統合研究センター	教授／チームリーダー
谷口 忠大	立命館大学 情報理工学部	教授
長井 志江	東京大学 ニューロインテリジェンス国際研究機構	特任教授
藤吉 弘亘	中部大学 工学部	教授
松井 知子	統計数理研究所 モデリング研究系	教授
持橋 大地	統計数理研究所 数理・推論研究系	准教授
森永 聡	日本電気(株) データサイエンス研究所	上席主席研究員

関連分野(順不同) : 人工知能・機械学習(理論・先端技術開発・応用)・数理統計・自然言語処理・画像処理・信号処理・ロボティクス・情報知覚・創発システム・認知科学・神経脳科学・ソフトウェア工学・データマイニング・アルゴリズム・エージェント(メカニズムデザイン)・集合知・科学技術と社会・自然科学応用など(応募分野は上記に限られません)

運営の基本方針

年2回の領域会議を中心に

さががけは自由な挑戦の場:

- 短期的成果でなく、失敗をおそれず、挑戦的な研究を

さががけは切磋琢磨の場:

- 未来のリーダーになる研究者同士のネットワークづくり

さががけは成長の場:

- 領域アドバイザーからの助言（研究構想、計画、キャリアアップ等）。
- 関係する人文・企業・行政等の多様な関係者との交流。
CREST等の領域との連携。

研究者・アドバイザーの有志による研究会（
つながり会）なども行っています
（2022年4月 東京、京都）

応募の留意点

- 短期的な成果でなく、それぞれの研究分野に**大きなインパクト**を与えるような**成果**を目指してほしい
- **将来の「信頼できるAI」**のために、どのように現在の**問題点・困難**を解決し、どのような**新しい概念や技術**を作り出したいかを、**申請者自身の言葉**で語ってほしい
- ご自身の**目標実現**に際して、どのような**領域参加者**や**他領域**と、**どのような連携**をしたいかを自由に記載いただけるとうれしい

研究期間と予算

- 研究期間: 3.5年以内
- 研究費の規模: 4,000万円(直接経費)上限
- 採択数: 10件程度 X 3期(令和2、3、4年予定)

詳しくはJSTの募集要項をご覧ください

公募スケジュール

研究提案募集

募集
締切

5月31日 (火) 正午

※切後は提案を一切受理しませんのでご留意下さい

書類選考会

7月12日(火)

書類選考通過者への連絡期限

7月19日(火)

面接選考会

8月1日(月)、2日(火)

※ 具体的な面接日時についてはJSTから指定させていただきます。あらかじめご了承ください。

研究開始

10月以降