



## 1. 施策の概要

- 英国のAISIは、自らAI安全性・セキュリティに係る評価ツールを開発。また米国は、企業が作ったリスクの高いAIを評価・分析する能力も持つべく強力に体制強化。**我が国のAISIは、令和6年2月に設置されて以降、人員・予算が脆弱なため、AIの安全性・セキュリティに係る機能が十分でない。**このような中、世界では、生成AIに加えて、AIEージェントやフィジカルAIといった高性能なAIシステムの利用・開発が急速に進んできているが、我が国ではこうした最先端のAIシステムにおいて原因究明が困難なインシデント発生時に対応する体制が十分に整備されておらず、AI利用・開発の促進を阻害する可能性。
- 令和7年9月のAI法の全面施行を受けて策定された**AI基本計画でも、AIの利用・開発を支えるべく、AIモデルの適正性に係る評価機能の構築を含むAISIの抜本的な機能強化が求められている**ほか、この度の**経済対策に係る総理大臣指示でもAI分野の官民連携投資が求められており、早急に民間企業等が安心して信頼できるAIを利用・開発できるよう機能の構築が不可欠。**
- そこで、我が国の民間企業等のAI利用・開発を促進しAIへの民間投資を増やすため、**AISIの抜本的な機能強化を前倒しで進め、以下の2つを実施する。**

### (1) AIEージェントの利用・開発に係るガイドやツール等の開発

- 各国が参画するAISI国際ネットワーク会合を通じて、日本AISIが、AIEエージェントをはじめとした最新のAIシステムの開発・利用におけるセキュリティ確保に係る課題等を把握。
- これらの結果をもとに、AIEエージェントをはじめとした最新のAIシステムが連動して動作するシステムにおけるセキュリティ確保のための調査（技術動向・事例・課題等）を実施。調査結果等に基づき、分野\*ごとのAI開発・利用に係るガイドの作成、及びガイドに基づくAIシステムに係る安全性検証を行うためのツールを開発。

※ガイド作成及びツール開発を行うドメインの想定例：分野共通、ヘルスケア分野、ロボット分野等のAIEエージェントの活用が想定される業界

### (2) AI評価に必要なベンチマークの開発

- 信頼できるAIの開発に不可欠な、学習用（質問と答えのセット）と評価用（よく似た質問と答えのセット）それぞれのデータセット作成及び評価手法の確立  
※を行うとともに、妥当性検証のためのAI評価プログラムの開発等を行いベンチマークを開発。

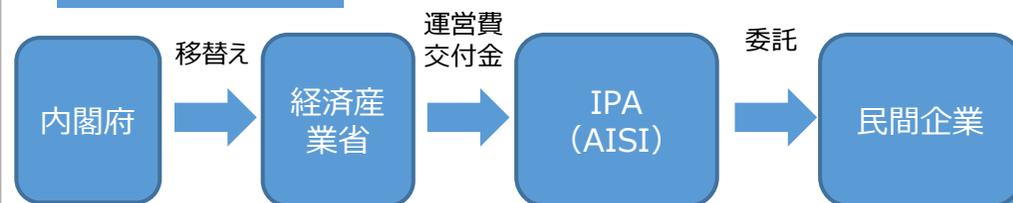
※ベンチマークを開発するカテゴリー（①安全性（バイアス、偽情報・誤情報、情報漏洩等）、②分野依存（ヘルスケア、法律、行政、経済等）、③ Jailbreak、④セキュリティ、⑤Eージェントモデル、⑥マルチモーダル、⑦ロボティクス、⑧評価プラットフォーム

- 国内外の類似の仕組みや認証スキームの検討に係る調査等を行い、これらの結果も参照しつつ、上記のベンチマーク等を用いてAIの安全性を認証するための枠組みを検討。

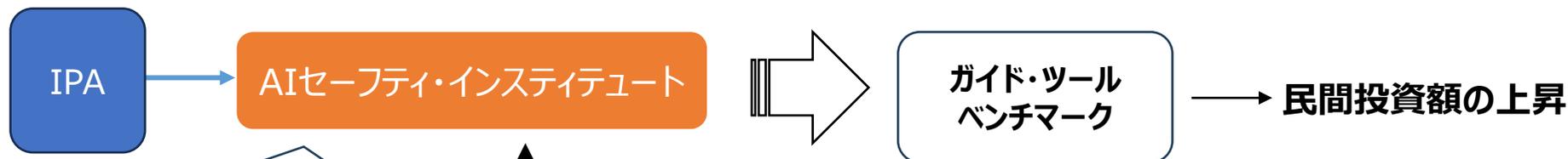
## 2. 施策の対象・成果イメージ

生成AIのみならず、AIEエージェント、フィジカルAIを含めた、**各ドメイン別のAIサーフェティ・セキュリティのガイド・ツール・評価基盤（ベンチマーク、コンソーシアムによる検証制度）を構築し、AIへの民間投資額を上昇させる。**

## 3. 資金の流れ



#### 4. 取組内容（システム概念図）



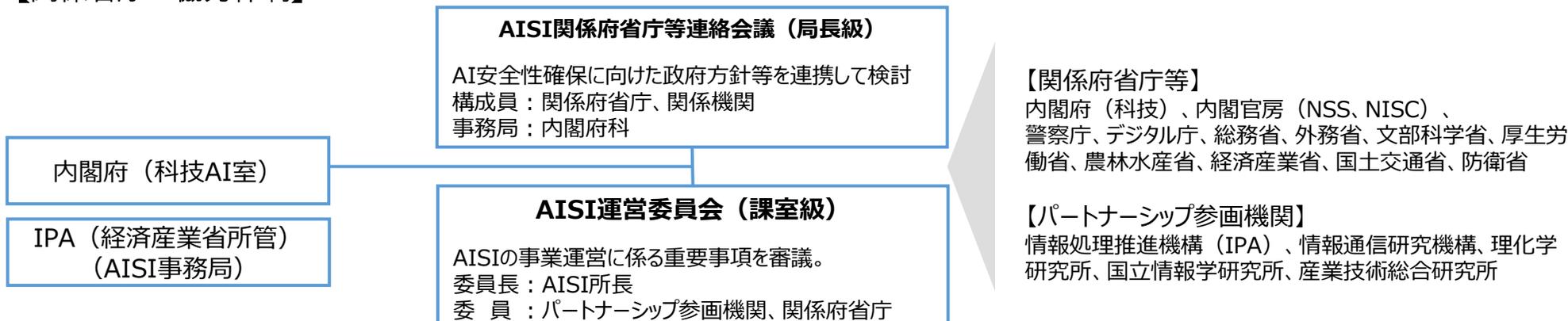
##### 【これまでのAISIの業務】

- 安全性評価に係る調査、基準等の検討
  - ✓ 安全性に係る標準、チェックツール、偽情報対策技術、AIとサイバーセキュリティに関する調査
  - ✓ 安全性に係る基準、ガイダンス等の検討
  - ✓ 上記に関するAIのテスト環境の検討
- 安全性評価の実施手法に関する検討
- 他国の関係機関（英米のAISI等）との国際連携に関する業務

##### 本事業で実施する新たな取組

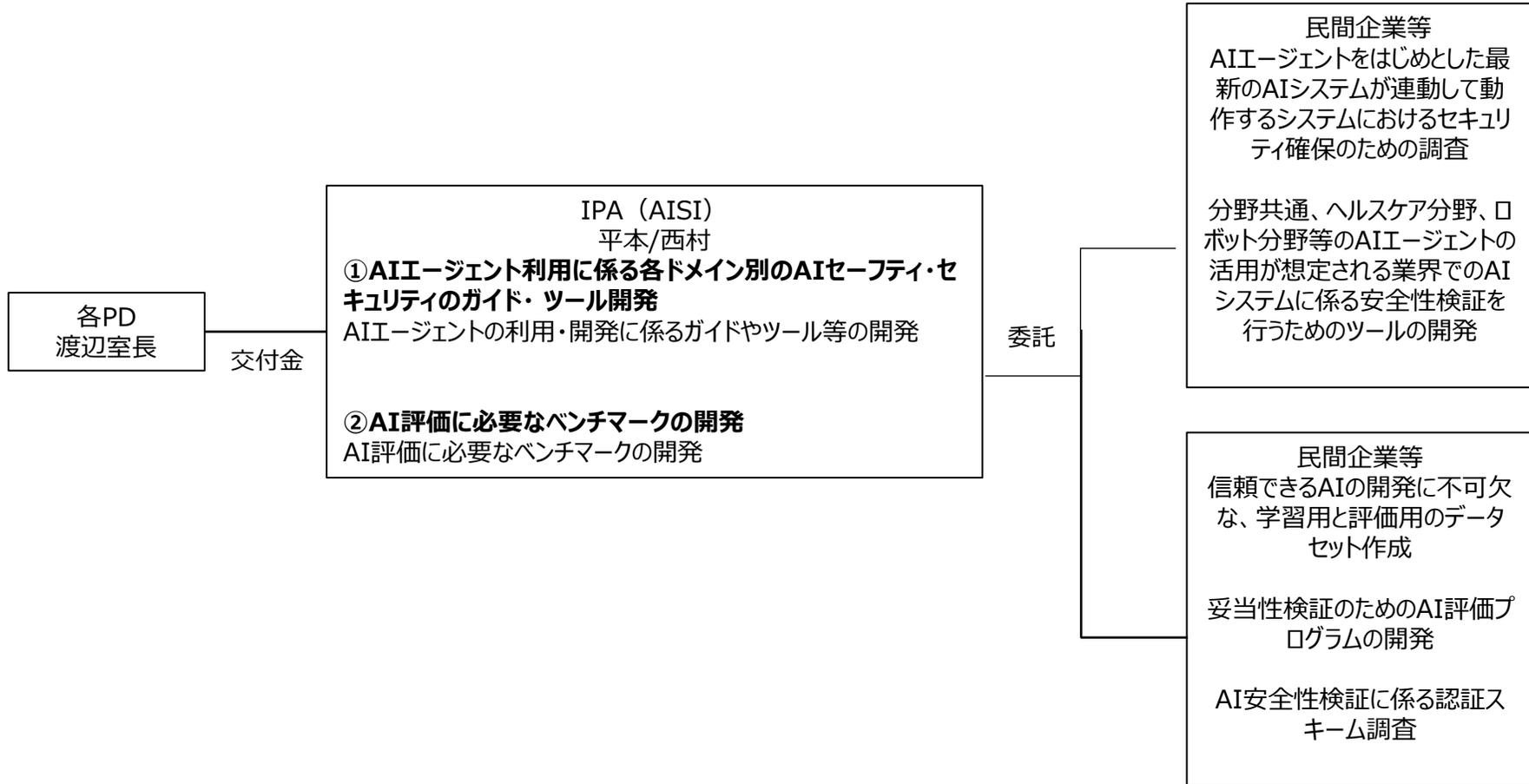
- ① AIエージェントの利用・開発に係るガイドやツール等の開発
- ② AI評価に必要なベンチマークの開発

##### 【関係省庁の協力体制】



## 5. 取組スケジュール

テーマ名	令和7年度			令和8年度												
	1月	2月	3月	4月	5月	6月	7月	8月	9月	10月	11月	12月	1月	2月	3月	
①AIエージェントの利用・開発に係るガイドやツール等の開発	契約・手続き			国内外における、AIエージェントをはじめとした最新のAIシステムの開発・利用技術調査・分析				調査分析結果を踏まえた、国内外における、AIエージェントをはじめとした最新のAIシステムの開発・利用におけるガイドの作成								
									国内外における、AIエージェントをはじめとした最新のAIシステムの特有のリスクを評価するツールの開発							
		AIエージェントに係るガイドやツール開発ワーキンググループ（仮称）														
②AI評価に必要なベンチマークの開発				信頼できるAIの開発に不可欠な学習用及び評価用データセットの作成												
				AI安全性ベンチマーク、AI評価プログラムの作成												
				ベンチマーク開発ワーキンググループ（仮称）												



## 7. 実施内容・到達目標（KPI）

テーマ名	実施内容の概要 到達目標（KPI）
①AIエージェントの利用・開発に係るガイドやツール等の開発	<p>【実施内容】</p> <ul style="list-style-type: none"> <li>国際AISIネットワークの構築と拡大の推進。AIエージェント等の最新のAIシステムに係る議論を行い必要な情報を把握</li> <li>国内外における、AIエージェントをはじめとした最新のAIシステムの開発・利用技術調査・分析</li> <li>調査分析結果を踏まえた、国内外における、AIエージェントをはじめとした最新のAIシステムの開発・利用におけるガイドの作成</li> <li>国内外における、AIエージェントをはじめとした最新のAIシステムの特有のリスクを評価するツールの開発</li> </ul> <p>【到達目標（KPI）】</p> <ul style="list-style-type: none"> <li>AIシステムの開発・利用におけるガイドの作成・提供 1件</li> <li>最新のAIシステムの特有のリスクを評価するツールの開発・提供 1件</li> </ul>
②AI評価に必要なベンチマークの開発	<p>【実施内容】</p> <ul style="list-style-type: none"> <li>信頼できるAIの開発に不可欠な学習用及び評価用データセットの作成</li> <li>妥当性検証のためのAI評価プログラムの開発</li> <li>AI安全性検証に係る認証スキーム調査</li> </ul> <p>【到達目標（KPI）】</p> <ul style="list-style-type: none"> <li>信頼できるAIの開発に不可欠な学習用及び評価用データセットの作成・公開 3件</li> <li>妥当性検証のためのAI評価プログラムの開発・公開 1件</li> <li>AI安全性検証に係る認証スキームの調査結果をもとにリスク分類や評価ガイドラインをもとにAIシステムを認証するためのスキームを構築</li> </ul>