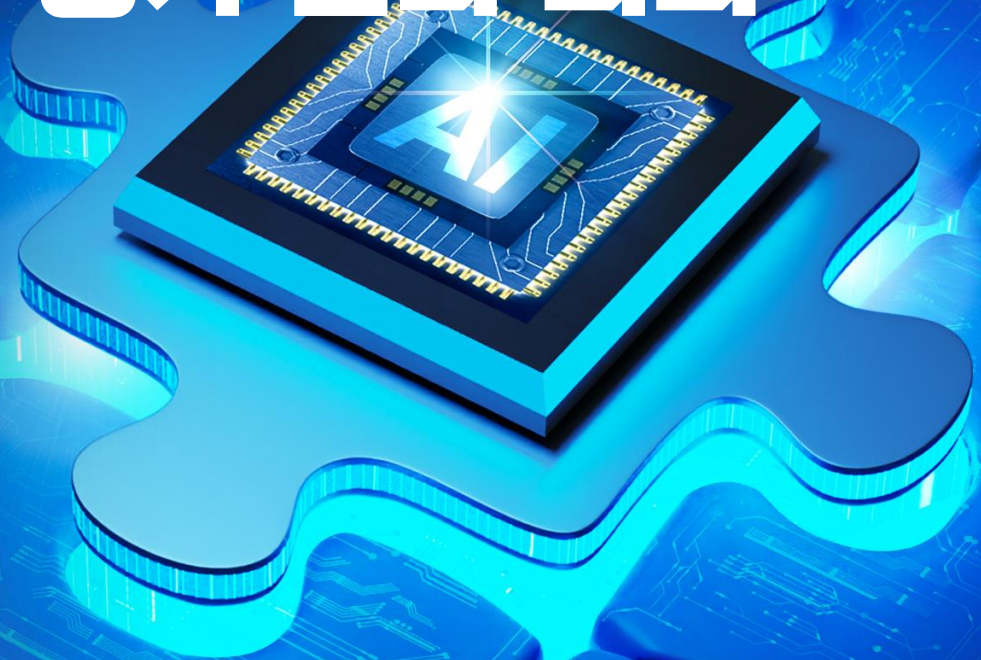


한국 AI안전연구소

: AI안전컨소시엄 정책 및 평가 분과 회의

2025. 4.22

ETRI **AISI** AI Safety Institute



진행 순서

- ▶ **SI안전연구소 업무 현황 소개**
- ▶ **SI안전컨소시엄 안내**
- ▶ **분과 의제**
- ▶ **기타 의제 제안 및 제언**
- ▶ **향후 일정 안내**

AI안전연구소 '25년도 주요 활동 목표

법·제도 정책 및 대외

과학 기반의 AI 안전 정책 지원과,
국내외 법·제도 및 글로벌 규범 논의와 선도를 위한 기반 구축

- ➡ AI기본법 제 32조 가이드라인 및 하위법령 제정 지원
- ➡ EU AI Act 행동 강령(Code of Practice) 안내서 마련
- ➡ AI 위험지도 개발
- ➡ 첨단 AI 위험관리 방법에 관한 가이드라인
- ➡ 해외 AI안전연구소 및 주요 기관과의 협력 강화 및 확대
- ➡ 첨단 AI 안전평가 및 테스트 분야 사실 표준 대응
- ➡ AI안전컨소시엄 운영
- ➡ AI안전연구소 운영 및 자문위원회 구성/운영

AI안전연구소 '25년도 주요 활동 목표

평가 기법 개발 및 실증

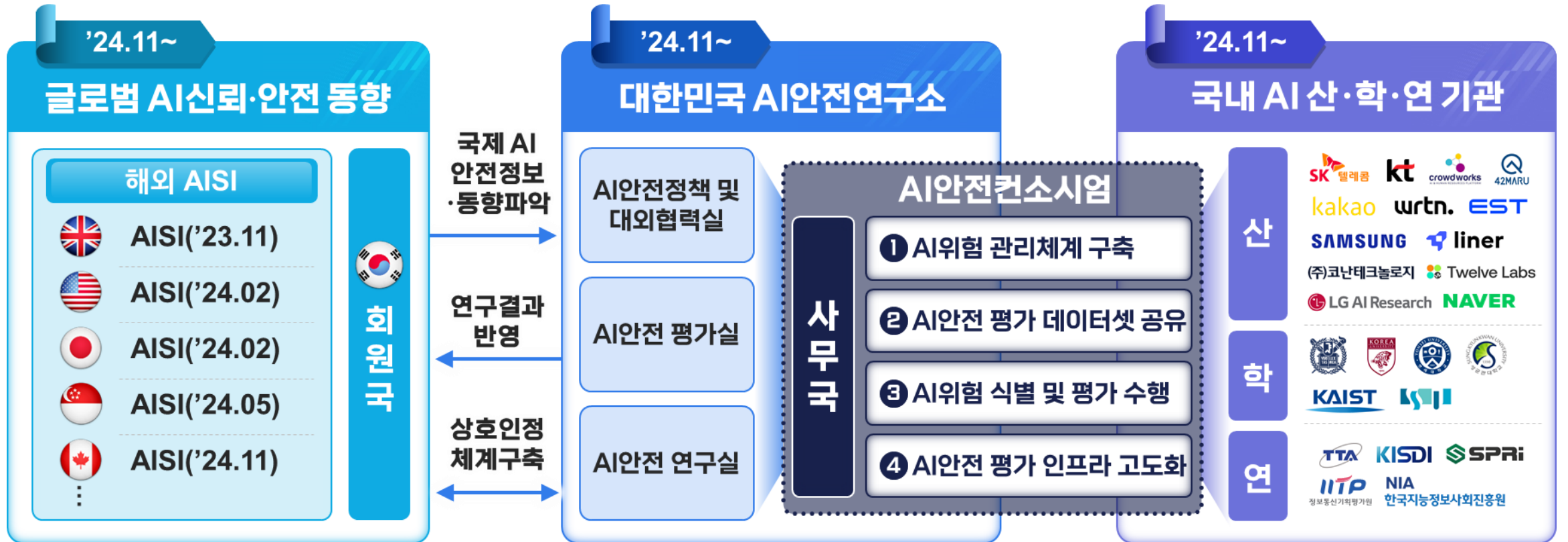
첨단 AI 대상의 안전 확보를 목표로,
안전성 평가를 위한 방법론, 데이터셋, 도구 등의 기술 개발

- ▶ 첨단 AI 평가 기법 및 안전 사례 선행 연구
- ▶ AI 에이전트의 안전성 평가 기술 개발
- ▶ 첨단 AI 안전성 평가를 위한 벤치마크 데이터셋 개발
- ▶ 국산 첨단 AI 시스템에 대한 안전성 평가 수행
- ▶ 첨단 AI의 위험 발굴을 위한 국제 레드티밍 행사 개최
- ▶ 첨단 AI의 안전성 평가를 위한 프롬프트 상시 수집
- ▶ 평가 신속 대응 및 기밀 보장을 위한 컴퓨팅 인프라 구축

AI안전컨소시엄 소개: 개요

AISI를
매개점으로

국내 산학연 24개 기업 및 기관 참여,
AI안전연구소 업무현황과 현안을 공유하고, AI안전 연구 협력과 담론 형성



AI안전컨소시엄 소개: 분과별 소개

정책 분과

국내외 법, 정책 및 거버넌스에 대한 대응과, 적극적인 글로벌 협력 전략 마련, 이를 위하여 우리나라가 최우선적으로 대응해야 할 위험의 정의

삼성전자
이스트소프트
뤼튼테크놀로지스
라이너
서울대학교
연세대학교
소프트웨어정책연구소
정보통신정책연구원
한국지능정보사회진흥원
네이버
케이티
포티투마루
고려대학교
성균관대학교
한국정보통신기술협회

〈 정책 분과 Agenda 〉

구분	안건	성과
~6월	단기 AI 위험지도 개발 및 검토	AI위험 정의 및 지도 개발
~8월	중기 EU GPAI CoP 안내서 제작 참여 및 검토	CoP 안내서
수시	글로벌 AI안전 정책 및 규범 의견수렴	AI안전컨소시엄 의견서

AI안전컨소시엄 소개: 분과별 소개

평가 분과

국내 상황에 맞는 위험 요소에 대한 평가 연구와 함께,
AI안전연구소 네트워크를 비롯, 국제사회와의 논의를 통한 글로벌 호환성

카카오
네이버
엘지에이아이연구원
케이티
에스케이텔레콤
코난테크놀로지
포티투마루
클라우드웍스
고려대학교
성균관대학교
한국정보통신기술협회
뤼튼테크놀로지스
서울대학교
소프트웨어정책연구소
정보통신정책연구원
한국지능정보사회진흥원

< 평가 분과 Agenda >

구분	안건	성과
~6월 단기	국제AI안전보고서 집필	한글버전 보고서
~11월 중기	안전성 평가 관점 수립 및 검토	AI안전성 평가 로드맵
장기	국제 공동 평가 수행	국제AISINetwork Joint Testing 보고서

AI안전컨소시엄 Agenda 소개

평가 분과

국제AI안전보고서 한글버전 집필 및 국제AISINetwork Joint Testing 개요

1. 국제AI안전보고서 한글버전 집필

- International AI Safety Report
 - 의장: 요슈아 벤지오
 - 전문가 자문단: 전세계 30여개국 96명 전문가 (한국: 이경무 교수)
 - 2023년 영국 블레츨리 선언에 따라 작성 시작
 - 2025년 2월 파리시행동정상회의에서 v1.0 발표
 - 한국 버전 협력 제안
- 주요 내용
 - GPAI의 세 가지 위험
 - 악의적 사용, 기술적 결함, 시스템적 위험
 - 위험 완화 방안
 - 2025년 초, Agentic AI에 대한 이슈 추가
 - 기존 위험 + 자율성, 통제 불능, Hijacking, AI협업 등



AI안전컨소시엄 Agenda 소개

평가 분과

국제AI안전보고서 한글버전 집필 및 국제AISINetwork Joint Testing 개요

2. 국제 AISI 네트워크 공동 테스트

- International Network of AI Safety Institutes
 - 출범일: 2024.11.20 (미국 샌프란시스코 회의)
 - 의장국: 미국
 - 참여국(10개국): 한국, 영국, 미국, 일본, 싱가포르, EU, 프랑스, 케냐, 캐나다, 호주
- 목표
 - AI안전에 대한 위험과 그 대응 방안에 대해 과학적인 공동 이해 구축
- 구성
 - Track 1: 합성 콘텐츠 위험 대응
 - Track 2: AI 시스템 테스트
 - Track 3: 첨단 AI 시스템의 위험 평가
- 변화
 - 미국 행정부 변화, 영국 AISI 명칭 변경
 - 2월 파리 AI행동 정상회의 이후 7월까지 의장단 미정

AI안전컨소시엄 Agenda 소개

평가 분과

국제AI안전보고서 한글버전 집필 및 국제AISI네트워크 Joint Testing 개요

2. 국제 AISI 네트워크 공동 테스트

- 개요
 - 주도국: US, UK, SG, JP (2024.11)
 - 목적: AI 시스템 테스트를 위한 최선의 방법론 도출
 - 평가도구: Moonshot(SG), Inspect(UK)
- 1차 조인트 테스트
 - 대상모델: LLaMA 3.1 405B
 - 벤치마크: GSM8K, SQuAD 2.0, MMMLU
 - **결과**: 다국어 테스트 가능성 확인, 방법론 일치의 필요성
- 2차 조인트 테스트
 - 대상모델: Mistral Large, Gemma 2 (27B)
 - 벤치마크: AILuminat (MLC), AnswerCarefully (다문화/일본), CyberSecEval, Cybench
 - **결과**
 - 다국어 테스트 실행 - 10개 언어
 - 채점 과정의 본질적인 주관성과 번역의 불완전성
 - 자동 평가의 한계

International Network
of AI Safety Institutes

FEBRUARY 11-12, 2025 | PARIS, FRANCE

**International Network of AI Safety
Institutes Joint Testing Exercise:**

Improving Methodologies for AI Model
Evaluations Across Global Languages

Contributors: Members of the International Network of AI Safety Institutes

AI안전컨소시엄 Agenda 소개

평가 분과

국제AI안전보고서 한글버전 집필 및 국제AISI네트워크 Joint Testing 개요

2. 국제 AISI 네트워크 공동 테스트

- 3차 조인트 테스트
 - 보다 체계적인 접근: 위험 우선순위, 테스트 설계 및 Threat modelling
 - 기간: 4월 ~ 7월 (ICML)
 - 후보 주제
 - 다국어/다문화, 유해성/편향
 - 사이버보안
 - 에이전트 테스트 - multistep reasoning

공지 및 향후 일정

결과 공개

분과 회의 결과 및 관련 자료는 모두 AISI 홈페이지에 게시
<https://aisi.re.kr>

멤버 추가

SI안전컨소시엄 멤버 추가 가입의 경우,
별도의 요건 혹은 기존 멤버의 동의 여부에 대한 정책 마련 수요 조사

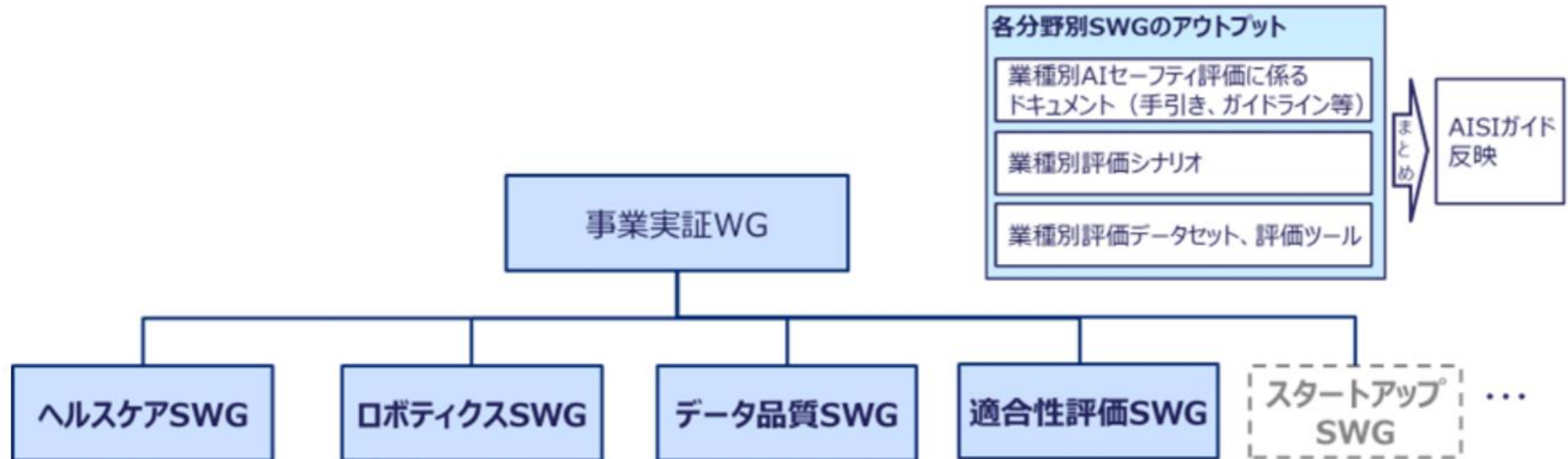
2차 회의

차기 회의는 6월 중순(합동, 오프라인)으로 예정,
회의 후 작업 참여 및 상황을 봐서 이후 일정 확정 및 공지

[참고] 글로벌 AI안전연구소의 컨소시엄 및 파트너십 활동

일본 AISI

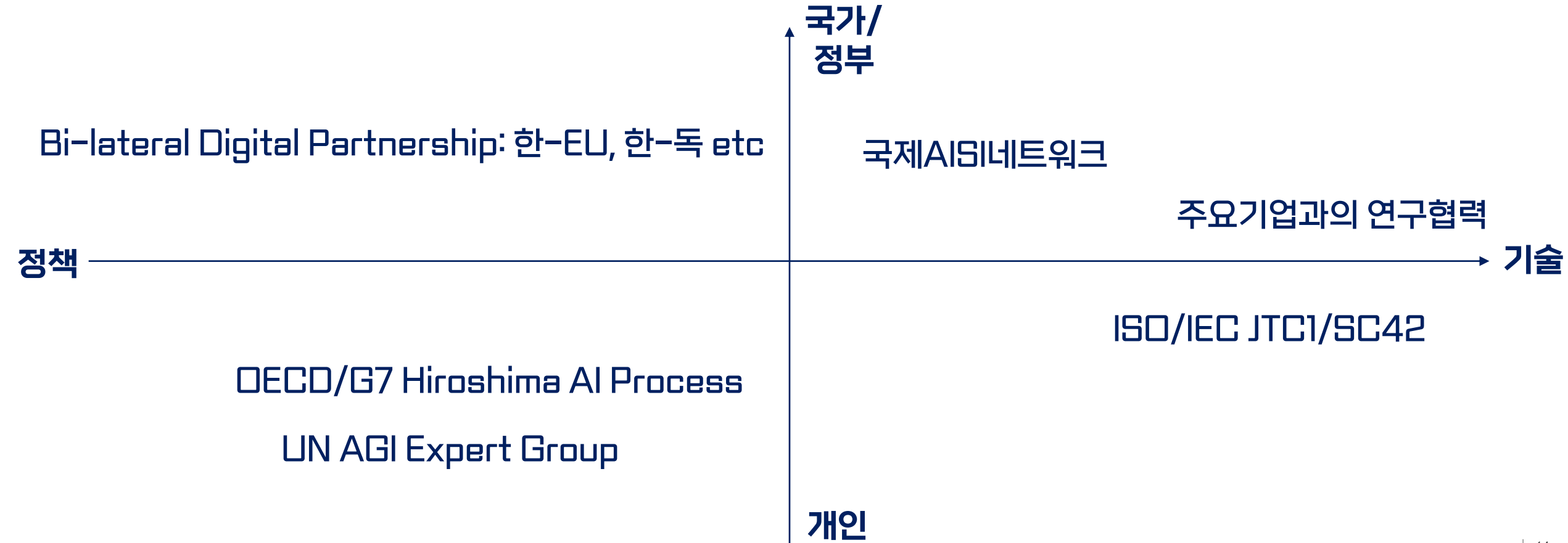
일본 AISI는 파트너십을 체결한 산학연 기업 및 기관과 연구협력을 통해, 각 산업 분야별 안전평가 실증 사업을 추진('25.3.24.)



[참고2] 글로벌 AI 거버넌스 및 협력 추진 현황

글로벌 참여

글로벌 AI 거버넌스 형성과 국제표준화 등
AI안전 분야 협력 활동 활발



Creating a Sustainable Future with Safe AI

Thank you for your time and attention

AISI AI Safety Institute