# CDAO
## Chief Digital & Artificial Intelligence Office

## The DoD's Approach to Responsible AI & The Responsible AI Toolkit

**Drew Brooks**
**Lead Scientist for Responsible AI**
**Tools U.S. Department of Defense**
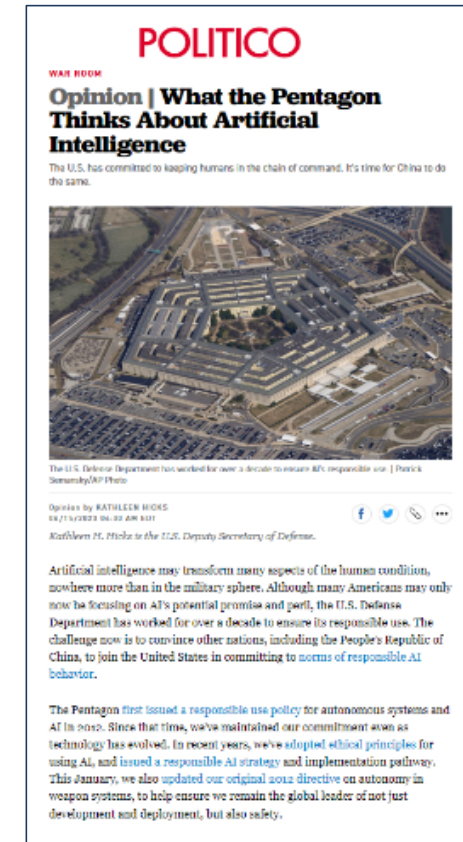
ResponsibleAI | ToolKit

# Winning Because of Our Values

> *"America and China are competing to shape the future of the 21st century, technologically and otherwise. That competition is one which we intend to win-not in spite of our values, but because of them."*
>
> –Deputy Secretary of Defense Kathleen Hicks

What does it mean to "win because of our values?"

CDAO

# What Is Responsible Artificial Intelligence (RAI)?

- RAI translates high-level values and Artificial Intelligence (AI) ethical principles into concrete actions, processes, metrics, and benchmarks to fit the use case at hand–and navigates any tradeoffs.

- RAI removes barriers to innovation and adoption through risk identification and reduction.

- RAI contributes to mission success through decision advantage and assurance.

CDAO

# Value Proposition of RAI: Assurance

**RAI increases assurance, thereby sustaining our tactical edge:**

- **Assurance for the Warfighter, Operational Commanders, and DoD Personnel to Reduce Cognitive Load:**
  - Provides assurance that technology has been developed to reduce risks of failure, unintended consequences, and dangerous or difficult ethical situations and choices for operational users.
  - Reduces cognitive load, allowing greater focus on contributors to mission success.

- **Assurance for the Department to Aid Adoption/Innovation:**
  - Provides assurance process to remove barriers to adoption and support effective innovation.

- **Assurance for Industry to Maintain Competitive Advantage:**
  - Ensures industry's trust that the U.S. Department of Defense (DoD) will responsibly steward its technologies.

- **Assurance for American Public:**
  - Ensures public's trust that AI-enabled capabilities employed by the DoD are aligned with our values.

- **Assurance for Allies to Increase Interoperability:**
  - Provides systems, tools, and processes grounded in shared values.
  - Is crucial, given the increasing need for interoperability.
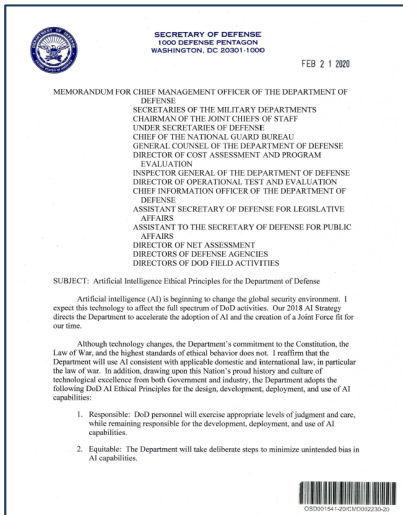
CDAO

# Background on the RAI Division

## RAI at DoD

- The DoD defines RAI as a dynamic approach to the design, development, deployment, and use of AI systems that implements the DoD AI ethical principles to advance the trustworthiness of such systems.

- RAI at the DoD emphasizes technical maturity, organizational change, modernized governance structures, and an understanding of sociotechnical risk.

## RAI Division's Role

- Is the primary technical advisor to the DoD on RAI.

- Oversees execution of the RAI Strategy and Implementation Pathway.

- Coordinates development and implementation of RAI tools, guidance, and other resources.

- Convenes DoD components to develop and recommend RAI best practices governing the creation, development, and use of AI within the DoD.

CDAO

# DoD AI Ethical Principles



**February 2020: AI Ethical Principles Memorandum**
The DoD formally adopts five AI ethical principles and designates the Joint Artificial Intelligence Center (now CDAO) as DoD's lead for coordination and implementation of the principles.

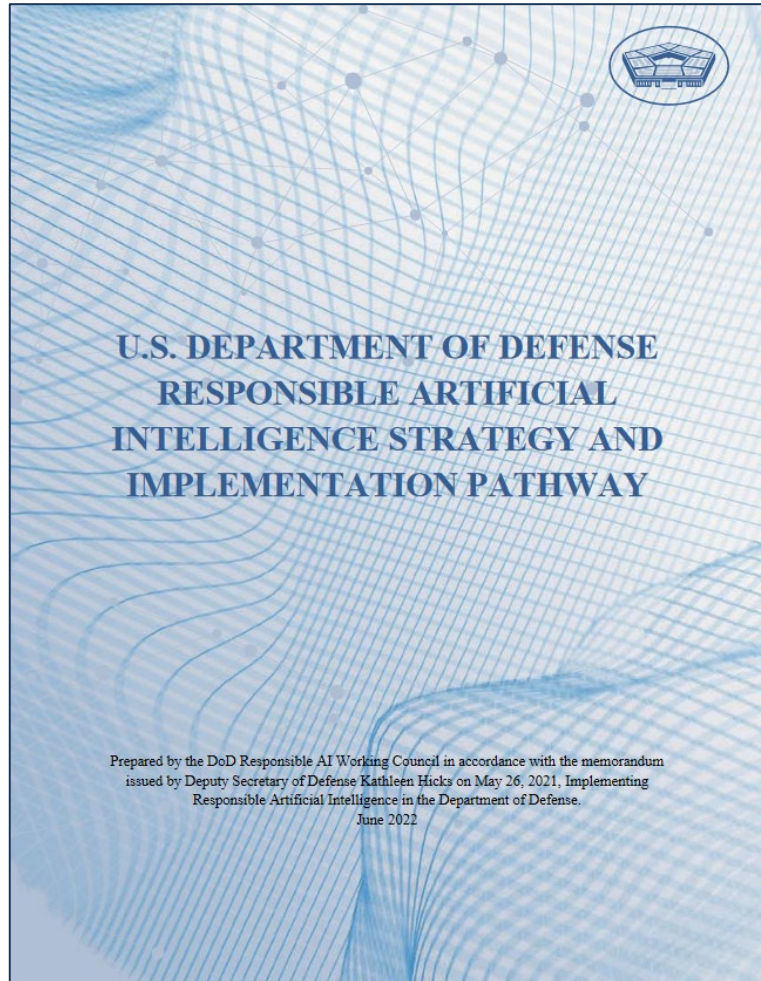| Principle | Description |
|---|---|
| **Responsible** | DoD personnel will exercise **appropriate levels of judgment and care,** while remaining responsible for the development, deployment, and use of AI capabilities. |
| **Equitable** | The department will take deliberate steps to **minimize unintended bias** in AI capabilities. |
| **Traceable** | The department's AI capabilities will be developed and deployed such that relevant personnel possess an appropriate understanding of the technology, development processes, and operational methods applicable to AI capabilities, including with **transparent and auditable methodologies, data sources, and design procedure and documentation.** |
| **Reliable** | The department's AI capabilities will have explicit, well-defined uses, and the safety, security, and effectiveness of such capabilities will be subject to **testing and assurance** within those defined uses across their entire life cycles. |
| **Governable** | The department will design and engineer AI capabilities to fulfill its intended functions while possessing the ability to **detect and avoid unintended consequences** and the ability to **disengage or deactivate deployed systems** that demonstrate unintended behavior. |

CDAO

# RAI Strategy and Implementation Pathway

**U.S. DEPARTMENT OF DEFENSE RESPONSIBLE ARTIFICIAL INTELLIGENCE STRATEGY AND IMPLEMENTATION PATHWAY**

Prepared by the DoD Responsible AI Working Council in accordance with the memorandum issued by Deputy Secretary of Defense Kathleen Hicks on May 26, 2021, Implementing Responsible Artificial Intelligence in the Department of Defense.
June 2022

**June 2022**

**Outlines the Department's Strategy for Operationalizing the Ethical Principles**

Six Tenets:

1. **RAI Governance:** Modernize structures for continuous oversight.

2. **Warfighter Trust:** Achieve justified confidence through training and education and test and evaluation and verification and validation.

3. **AI Product and Acquisition Life Cycle:** Identify and mitigate risks throughout life cyle.

4. **Requirements Validation:** Ensure AI systems are aligned with operational needs.

5. **Responsible AI Ecosystem:** Promote shared understanding through domestic and international engagements.

6. **AI Workforce:** Build, train, equip, and retain an RAI-ready workforce.

CDAO

# Examples of RAI Tools and Capabilities

RAI tools function in a number of ways to support the operationalization of DoD's AI ethical principles for capability developers, RAI practitioners, and senior leaders.

| What | Function | Example Tools |
|---|---|---|
| **Technical or Software Based** | Helps developers and testers assess factors such as bias, reliability, and safety | Data Bias Detection Tools<br>Explainability Tools<br>T&E Harness |
| **Documentation and Artifacts** | Provides traceability of data sources, model limitations, risk identification, and mitigation efforts | Use Case/Harms Analysis<br>Data Cards<br>Model Cards |
| **Frameworks and Checklists** | Provides prompts to guide users in creating muscle memory around new processes for risk assessment and ethical considerations | Common Failure/Mishap List<br>Algorithmic Impact Assessments<br>Ethics Maturity Assessments<br>User Research and Design Tools |
| **Knowledge Sharing** | Provides centralization for information sharing, learning, and common lexicon, practices, etc. | Use Case Repositories<br>Information Management Systems |
| **Executive Dashboards** | Provides visibility into organizational compliance, status, and risk | Key Performance Metrics<br>Progress Tracking |

CDAO

# RAI Toolkit

- **The Responsible AI toolkit is our organizing framework to make the capabilities being built out under the RAI Strategy & Implementation Pathway:**

  o Findable

  o Usable

  o Interoperable

- **Living document and web application (currently in minimum viable product form) building upon and incorporating:**

  o Industry best practices and tools (currently 70+ listed in the toolkit) and academic innovations

  o DIU RAI Guidelines and Worksheets, NIST AI RMF + Playbook, IEEE 7000, etc.

  o Tools being built through the RAI Strategy and Implementation Pathway

CDAO

# RAI Toolkit Priorities

- **Provides a process for demonstrating consistency/alignment with the DoD AI ethical principles**

- **Enables traceability and promotes assurance**

- **Provides a mechanism for collecting lessons learned that can serve as inputs to policy**

  - Enables empirical tracking of how RAI influences mission success

- **Provides common framework for partners and allies to develop shared assurance cases**

  - Aids interoperability and trust

  - Is co-developed NATO version of Toolkit

  - Is developing collaborations over the toolkit with IC, interagency, other allies and partners

CDAO

# Approach to Toolkit

**Top-Down Approach:**

Identifies the classes of tools that would be needed to align with the U.S. Constitution, executive orders, DoD AI ethical principles, other RAI Frameworks, long-standing international norms and values, etc.

**Bottom-Up Approach:**

Draws from market research studies of COTS/GOTS/OS RAI Tools, AI Ethical Frameworks, RAI Processes, and Standards (e.g., NIST AI RMF and Playbook, IEEE 7000, DIU Responsible AI guidelines, etc.)
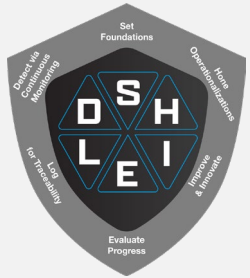
Ensures Coverage

Establishes Common Foundation & Enables Tailorability of Toolkit for Partners

Identifies Gaps

Updates as New Policy/ Tools/ Frameworks Emerge

CDAO

# Design Challenges and Principles

**The RAI Toolkit aims to seamlessly assist users to plan and execute the necessary RAI activities and select appropriate supporting artifacts and tools.**

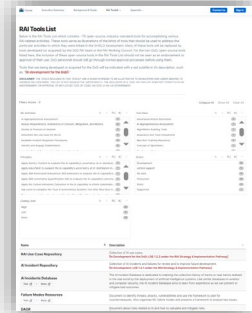| Challenges | Principles |
|---|---|
| Wide diversity of use cases and priorities across the DoD. | Modular and Tailorable |
| Demonstrate alignment with DoD AI ethical principles. | Traceable |
| Existing assessment processes can overwhelm a small team. | Lightweight |
| RAI processes require coordination among diverse team roles and stakeholder considerations. | Holistic |
| RAI activities should take place during all phases of AI development. | Integrated |
| Existing approaches assume expert RAI knowledge. | Upskilling |
| RAI research and practice is still evolving. | Iterative (Living Document) |

CDAO

# RAI Toolkit Components

## Currently Available in Toolkit MVP



### Planning Tools
Identify and document potential risks and plan RAI activities for mitigation

### Tools and Resource Database
Provide resources for implementing RAI plan

### Software Tools, Guidance and Best Practices, Checklists, Metrics

## In Development



### Evaluation Tools
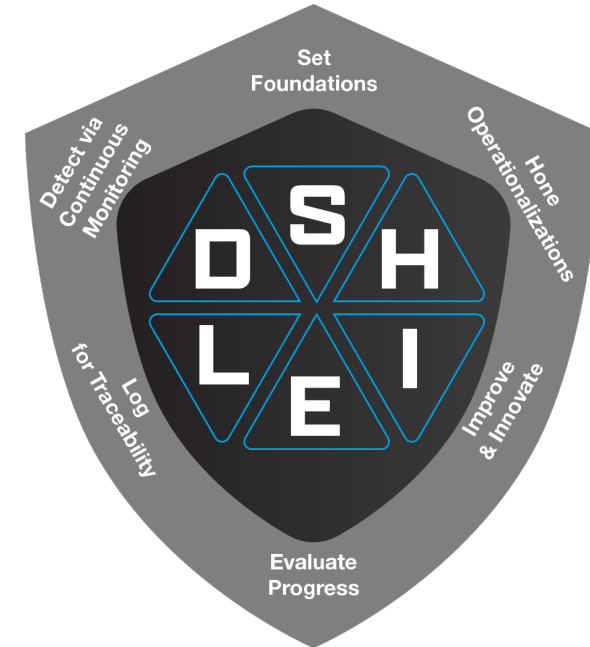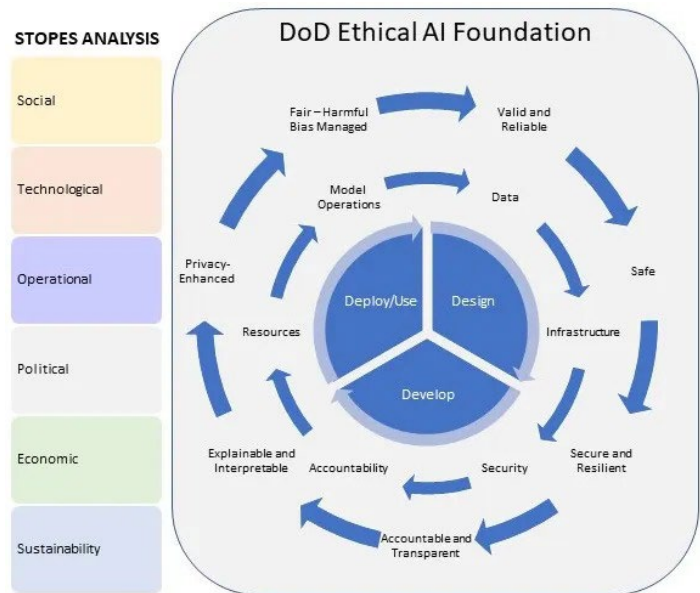Evaluate progress against RAI plan

### Oversight Dashboard
Monitor RAI progress and risk profile across programs and portfolios

**The RAI toolkit assists users to plan and execute RAI activities, and select appropriate supporting artifacts and tools.**

CDAO

# How: RAI Planning and Assessment



## Defense AI Guide on Risk Assessment

- Risk management guidance for DoD is aligned to NIST AI Risk Management Framework

- Risk management process initiates a SHIELD Assessment

- Supporting tools in development

## SHIELD Planning Process

- A series of six sequential classes of activities identifies RAI-related issues for tracking and mitigation

- List of issues are tracked throughout the life cycle via statements of concern

- Elements in the SHIELD assessment route the user to relevant tools within the tools database

# Tools and Resource Database MVP

- **Searchable Database (70+ items) of COTS/GOTS/open-source RAI Tools:**
  - Informed by CDAO market research and RAI FY22 tool survey
    - — Industry best practices and tools (70+)
    - — Academic methodologies
    - — DIU Responsible AI guidelines & worksheets
    - — NIST AI RMF + Playbook
    - — IEEE 7000

- **Customizable User Interface:**
  - Tailorable labels for ethical principles, development lifecycle phases, category names, roles, and disciplines
  - Interactive search and exploration

# Who: RASCI and Personas List

## RAI Role**

| RAI Role** |
|---|
| Users/Stakeholders |
| Mission Commanders |
| Senior Leader/AI Innovation Leader |
| Functional Requirements Owner |
| Program Manager |
| AI Ethics and Risk Specialist |
| Relevant Legal, Ethical, or Policy Expert |
| UX/Design/HMT/AI Adoption Specialist |
| AI Development Team<br>  System Architect<br>  Data Architect<br>  Data Operations Specialist<br>  Data Analyst<br>  Data Scientist<br>  Data Officer<br>  AI Engineer/AI/ML Specialist<br>  Data Steward |
| AI Test and Evaluation Specialist |
| IT/Cyber Expert |

| Role* | Definition |
|---|---|
| **Responsible** | The person who does the work to complete the task or create the deliverable. |
| **Accountable** | The person ultimately accountable for the work or decision being made; this person gives final approval. |
| **Supporting** | The person who provides support for those who are responsible or accountable; participates in doing the work of a task. |
| **Consulted** | Anyone who must be consulted with or add input prior to a decision being made and/or the task being completed. |
| **Informed** | The people who need to be updated on project status or informed when a decision is made or work completed. |

**Individuals or teams may be dual hatted;
Roles map to DoD Cyber Workforce (DCWF) roles – **BLUE text indicates relevant DCWF Role**

CDAO

# When: RAI Development Life Cycle

**Design**

**Develop**

Apply metrics and standards to evaluate system performance (on an ongoing basis, if necessary)

**Intake**

Identify use case and its relationship with existing system/capabilities

**Ideation**

Determine model parameters and establish AI model performance baseline against existing systems

**Assessment**

Develop business case for product and conduct data audit

**Acquisition / Development**

Run RFI process and select vendor; secure resources for internal systems

**TEV&V**

**Deploy**

**Use**

**Use**

Trained end user monitors outputs

**Integration & Deployment**

**Framework introduces the concept of "Gates" to the AI Development Lifecycle\***

- "Gates" indicate recommended considerations for progression to the next phase of development

*\*DoD RAI Strategy & Implementation Pathway (p13)*

CDAO

# RAI Toolkit Current Features



Navigation by
AI Product Lifecycle Stage

Export/Import function
to save and share progress

"GATE" filter
(displays most
essential
assessment
questions)

Export as PDF

Navigation by
Type of RAI Activity

SHIELD Assessment
identifies risk and
opportunities

Links to tools to
address identified
risks and
opportunities

Filters assessment
questions by
persona/project
role, AI ethical
principle, discipline,
etc.

# RAI Toolkit Web Application



RAI Toolkit Web App:

https://rai.tradewindai.com/

CDAO

# Toolkit Way Ahead

- **Develop versions of toolkit to support approvals and reviews for various use cases**
  - Deconflict with other required processes and documentation to support creation of integrated template or documentation process.
  - Use Tabletop Exercises/'Mock Reviews' to refine documentation process.

- **Develop versions of toolkit focused on generative AI/LLMs**

- **Pilot on other use cases throughout DoD, interagency, international partners**
  - Codevelop shared versions of the toolkit with partners.
  - Collect, organize, and share lessons learned.

- **RAI Toolkit tabletop exercises & technical exchanges with key allies & partners to aid interoperability**

- **Develop acquisitions-focused version of toolkit**

- **Integrate into AI training courses**

- **Continue to add functionality**
  - Develop user interface.
  - Add further tailorability features (data & model type, use case, risk profile, etc.).
  - Continue dashboard development.
  - Integrate with other tools (T&E, cyber) and with platforms.
  - Integrate feedback.

CDAO

# Closing Thoughts

"...ultimately, AI systems only work when they are based in trust. We have a principled approach to AI that anchors everything that this Department does. We call this Responsible AI, and it's the only kind of AI that we do. Responsible AI is the place where cutting-edge tech meets timeless values." - General Lloyd J. Austin III, Secretary of Defense

# Thank You and  Questions

Drew Brooks

*Lead Scientist for Responsible AI Tools*

U.S. Department of Defense

andrew.l.brooks.civ@mail.mil

RAI Toolkit Web App:
https://rai.tradewindai.com/