Partnership
on Technology
Innovation and the
Environment

# WORKSHOP ON ENVIRONMENTAL PROTECTION, ARTIFICIAL INTELLIGENCE, AND MACHINE LEARNING

# SUMMARY REPORT

**SEPTEMBER 21, 2018**

Contents

This document provides a summary of the discussion and conclusions of the Workshop on Environmental Protection, Artificial Intelligence and Machine Learning that was held at American University on September 21, 2018. This workshop was organized under the auspices of the Partnership for Technology Innovation and the Environment (PTIE) and is the fifth such workshop convened by the Partnership.

The purpose of the workshop was to bring together experts from many fields to explore, evaluate, and propose ways of using artificial intelligence (such as machine learning and natural language programming) in environmental protection. The workshop focused on two specific areas of environmental decision making as case studies: water quality (especially nutrient) management and chemical risk assessment.

In June 2018, American University's Center for Environmental Policy (Center), the PTIE Secretariat, in partnership with EPA, launched a planning process for the workshop. Dan Fiorino, Director of the Center, and his team reached out to individuals from government, industry and academia to discuss the environmental protection challenges that are already benefitting from AI applications and those that may benefit in the future. The team held phone interviews with approximately 20 individuals to understand their perspectives and obtain their input on the topics that would be relevant and to get names of other individuals who could contribute substantively to the process.

PTIE convened a steering committee to plan the September workshop (Attachment A). The committee was comprised of EPA officials (Office of the Chief Financial Officer, Office of Water, Office of Information Management, Office of Research and Development, Office of Chemical Safety and Pollution Prevention), American University staff members, and representatives from IBM Watson, Intel, the Environmental Defense Fund and the Environmental Law Institute. Steering committee meetings were held on June 27 and July 18 to identify the specific programmatic themes, to develop a process for diving deeply into the selected topics, and to engage a wider group of stakeholders in the working groups.

Subsequently, two workgroups, were formed -- one on chemical risk assessment and one on water quality/nutrient monitoring. These groups met regularly by phone to discuss the priorities for their work and develop issue papers to frame the issues for future discussion (Attachment B).
This summary presents the results of: (A) the water quality breakout and plenary discussion; (B) the risk assessment breakout and plenary discussion; (C) the cross-cutting issues discussed; and (D) follow-up actions. Attachment C includes the list of participants and their contact information.

**I.      Summary of the Water Quality Breakout**

The water quality breakout considered five questions:
1. What are some use cases where there is a demand for AI applications?
2. What about the supply or availability of AI tools?
3. Why are AI tools not used more widely in water quality?
4. What needs to happen to scale up the use of AI water quality and nutrient management?
5. How might an innovation challenge advance the use of AI tools to better manage water quality?

*1. What are some "use cases"[1] where there is a demand for AI applications?*
- **State of Maryland**: AI could help in projects such as stormwater pond automation to optimize overflow storage and water quality potential; overflow prediction; development of land cover sets using satellite images; and optimization of winter road salt use to limit the use of chlorides.
- **Chesapeake Bay**: Despite the large amount of Chesapeake Bay data available, there is a lag between nutrient reductions and estuary impacts data. In addition, the relationship of nutrient load to algal blooms is not clear. AI tools may support the development of predictive models.
- **Water Quality/Nutrient Trading:** The use of continuous nutrient sensors to improve treatment processes, generate credits, and obtain real-time information can reduce uncertainty in monitoring and optimize sampling and sensor placement. The major economic and environmental benefits of nutrient trading come from non-point sources; however, significant uncertainty remains about the reliability and timing of claimed reductions. There are benefits of monitoring water quality not only in streams but in rivulets, and there may be a role for AI to optimize monitoring, improve data analyses, and minimize such uncertainties.
- **Metropolitan Sewer District for Greater Cincinnati (MSDGC):** Given the vast distances and 300 discrete overflow points that are connected, AI would optimize the system further to minimize adverse environmental and health impacts without increasing the resources required for new infrastructure.
- **Water Treatment Facilities in general:** Regulation or other arrangements could draw upon technologies and AI techniques for their daily work. From a return on investment (ROI) perspective, these projects deliver returns to the enterprise in terms of existing CAPEX vs OPEX.
- **Lake George:** New York has gathered about 70,000 data points in 30 years. With the deployment of sensors this year, it will have far more data. With all of these data, there is a need for machine learning to effectively use those data. Further, there is a need to deploy commercial versions of AI to understand other water bodies.

*2. What about the supply or availability of AI tools?*
- The water quality group is focused on utilizing currently disparate structured data sources and in leveraging real-time data or Internet of Things (IoT) content from sensors for use in decision making. In the case of Lake George, for instance, there is the need to better predict the impact

---

[1]A use case is a description of all the ways an end-user wants to make use of a system. These are requests made of a system; use cases describe what that system does in response to such requests. Use cases provide a structure for gathering customer requirements and setting the project scope. They are also useful for end users in testing as it is designed, which leads to quicker development and a more useable system.
(https://www.pmi.org/learning/library/use-cases-project-manager-know-8262)

of salination by combining various datasets of waterbody, weather, and environmental attribute data, through the sensors deployed in particular areas. In the case of the Idaho National Lab, weather, transportation, and vehicle information is harvested and combined to predict conditions and stresses at any time of day to help inform decisions.

- It is possible to look not only at EPA's own structured datasets, but also to pull information more broadly. For instance, videos, voice, and text information is readily available; these could be curated and analyzed to inform decisions. One example is New York's subway system, which developed a model including internet information about local events or rush hour issues in the city to help operators decide about loading additional trains and on what line.
- Manufacturing companies need to get permits, and complain about cycle time. AI may make it easier for regulatory communities to understand their obligations, and how to streamline resources in the permitting process.
- Although blockchain is not a AI or ML tool, it may apply when there is a need to track many different systems and institutions–authorities that substantiate data accuracy or verify whether changes are authorized or not. Blockchain technology could be useful for water quality trading (and other) applications by enhancing transparency, verification, and documentation.

### 3. Why are AI tools not used more widely in water quality?

- *Issues associated with trusting the models (the black box):* The machine builds a model using training data fed into it to fit the parameters; data scientists need to ensure that the model is trained with appropriate data. A model developed for one city cannot necessarily be applied directly to other cities/locations. Specific parameters need to be adjusted to avoid inappropriate results.
- *Issues regarding generalizability and transferability of AI algorithms:* Training has three components: data, algorithm, and decision-making. During training, the algorithms learn basic information and patterns in the data. If the model works on a new dataset in the same way it did during training, then the model is generalizable to decision-making. Generalization is related to the scalability of the product. However, the models react to specific situations as the data changes, meaning that the model may be transferable but needs calibration. Oversimplifying this concept, there are two steps. In the first step, the model is trained at the lowest layers. It will learn how to discover basic information from the data set that will be reused during step 2. In step 2, the highest layers need updates to fine-tune the model with the task of the specific watershed decision-making. The lower level layers are unmodified as their functions are universal and less specific. When the decision making is similar (e.g. storm water treatment) data collection process and dataset labels are mostly shared and the techniques and methods are mostly standardized.
- The AI problem formulation needs to be set up first. Any good project needs to approach the problem formulation very well at the outset.

### 4. What needs to happen to scale up the use of AI water quality and nutrient management?

- Advancing the effective use of AI models in nutrient and water quality management may require rulemaking changes.
- Job losses from automation are often a concern. A major challenge is to overcome the resistance and train the workforce to acquire new skills.
  There is no shortage of technological solutions. The major problem is to find the right drivers that will promote demand; these may include funds, policy, and regulatory incentives.

*5. How might an innovation challenge advance the use of AI tools to better manage water quality?*

- Innovation challenges could apply to a water body, community, organizations, treatment facilities, or any other potential users. They may compete with each other to access opportunities to demonstrate AI/ML use. This opportunity could be distinct depending on the specific problem the agencies/facilities/organizations are going to address.

**II.      Summary of the Water Quality Plenary Session**

- The demand for AI applications is clear; connecting the supply and demand is essential.
- In trying to help move forward with AI, it is important to show success. It is also important to create the right incentives to encourage innovative thinking.
- Procurement and technology approval are still large impediments for AI applications. One way of overcoming that is through a partnership; organizations focused on a particular area could share upfront costs, knowledge, experiences, and similar datasets to demonstrate the value of partnerships.
- Organizations need to formulate the problem, guide criteria for determining the AI tools in particular application and develop action plans. In addition, social and behavioral aspects are important considerations on what AI can do in an organization and its management structure.
- The Partnership for Public Service has provided a forum for AI experts from DOD and CIA, among others, to share their experiences about what is important to decision-makers.
- There is a need to encourage citizen participation with sensors and to engage citizens in resource conservation and protection efforts.
  There may be value in producing a position paper on the significance of AI to EPA and identifying where AI can help with advances in different programs and elements of EPA decision making.

**III.      MSDGC Case Study Presentation**

Reese Johnson, from the Metropolitan Sewer District for Greater Cincinnati (MSDGC) Watershed Operations delivered a presentation on the utility's approaches to reduce combined sewage overflow pollution. Below are the summary points.

- Approximately 11 billion gallons/year are discharged from Cincinnati's hundreds of combined sewer overflow points. A $ 3.2 billion consent decree was approved to address overflow pollution ten years ago. With scarce local resources, MSDGC is using several approaches to address the problem, including: improving current gray infrastructure with green ones; keeping stormwater out of combined sewers during heavy rains through strategic separation; taking out pipes running directly into waterways; upsizing pipes; getting information and prioritizing operations of closing or opening gates from real-time sensors; and by improved monitoring system (SCADA).
- Regarding the real-time sensors, there can provide better stream data. The needed sensors would cost approximately $7 million, more than budgeted for this purpose. Hence, a partnership was established with EPA to pilot a sensor challenge. The challenge was a way for the industry understand what the utility needed; and helped them to consider multiple designs.
- This challenge was successful and may identify a roadmap to address the difficulties facing wastewater treatment operators. The public entity opens a competition, bidding, and RFP. The utility gets a particular answer that is mostly standardized through the use of multiple different type of sensors. These sensors generate data that are transferred to the SCADA system, and all

data are gathered and assessed altogether. As a result, MSDGC is able to operate and manage the process more remotely, dynamically, and flexibly to minimize combined sewer overflows.

- MSDGC has achieved economic, social, and environmental benefits. It improved treatment operations, significantly reduced costs, reduced overflows, and improved watershed protection.
- MSDGC is seeking AI tools to further optimize this system. 300 overflow points are located across a large spatial area with variation in rainfall levels that require a better prediction model.
- As a matter of comparison, the cost generally in the industry is about $.87/gallon for treatment, a storage system costs about $.40/gallon, source control costs about $ .23/gallon, real-time control with adjustable gates costs $.03/gallon, which is the current MSDGC cost. With AI further optimization, according to the cost-benefit analysis, this expense could drop to $.01/gallon.
- What if better sensors to estimate pollutant load for water quality are available at a reasonable cost? MSDGC's operation costs and ongoing retrofitting are conditioned on volumetric measurement. The reason is that current legislation requirements are based on water quantity as opposed to water quality. Is it possible to change this with a water quality sensor challenge?

<p style="text-align:center"><u>RISK ASSESSMENT WORKGROUP</u></p>

**I.     Summary of the Risk Assessment**

- The risk assessment workgroup focused on identifying incremental, short-term opportunities for streamlining the systematic review (SR) process for chemical risk assessments. While there are many areas for improvement in the systematic review process, the members felt that to achieve a useful outcome from this discussion they would need to focus on a manageable goal.
- The SR process is the gold standard in chemical risk assessment. SRs are designed to be clear, transparent, and trackable. SRs are extremely labor intensive and expensive, however, making them a logical starting point for using AI in the chemical risk assessment (RA) process.
- There are a variety of opportunities to bring AI into SR. Chemicals are categorized based on chemical structure and biological activity, both of which are used to build models. Currently, efforts are underway to build models that can mine text (e.g., from scientific journals) using annotation.
- The workgroup walked through the Modular Risk Assessment Workflow from the issue paper (Attachment A). SR is based on a modular and iterative approach. Typically, it starts with a broad literature search to scope, plan, and formulate the problem associated with a particular chemical. Machine learning (ML) and evidence mapping are used to enhance problem formulation. The company *Sciome* has developed machine learning tools to be used during title abstract screening.
- Challenges remain in the manual process of extracting, labeling, and organizing data. These can be addressed by reducing the burdens of literature reviews. This is possible by limiting which parts of a study need to be reviewed, which may be eliminated, and assessing the quality of studies.
- A significant bottleneck in the data extraction/analytics process is language choices. There are multiple common and scientific words for things as basic as the definition of a "rat" or a "tumor" to different strains of chemicals (i.e., based on purity or composition). Ontologies and best data management practices are critical to organizing the evidence and promoting interoperability.

- The model for data extraction/analytics needs to be scalable to address the challenge of variable risk assessment standards. Some assessments have more stringent, legally-mandated standards than others.
- The group relied on questions from Booz Allen's *Data Science Playbook* as a discussion guide: https://www.boozallen.com/content/dam/boozallen_site/sig/pdf/publications/data-science-playbook.pdf)

## II.    Issues Related to Data

- Data extraction requires several steps: labeling the data; identifying its location; extracting free text fields; extracting the quantitative data; and extracting the qualitative data to support the evaluation. The first incremental step is to improve data extraction that will inform decision-making
- Confidentiality surfaced as a broad discussion point. In the natural language processing (NLP) space, a significant obstacle is dealing with copyrighted data (subscription journals versus open source). Unlike medical research, only a small percentage of the toxicology literature provides full open access. While all studies with public funding are required to be public after one year, the journals do not often allow for the data to be extracted, even after the licenses expire.
- Medical research overcomes the proprietary data issue by masking the author and individual details. This issue is often addressed in toxicology research by extracting data sets from their source. This is labor intensive, however, and the data have to be recreated for every review.
- EPA is doing three SR-data extractions and study evaluation to develop training sets in partnership with NIEHS. EPA found 3,000 relevant, high quality references, although only 200-300 were both relevant *and* publicly available, and thousands are needed to complete a thorough assessment. They are looking to combine this data set with human health data to build a training set that can be used to build algorithms.
- One issue discussed is whether we are applying a higher standard to AI than we apply to humans. An example is that of self-driving cars that use human experience to establish rules. Many such cars trained in this way failed. Experience doesn't work for AI in the same way that it does for humans. AI models need to be built with a lot of distractions/extraneous information so that systems can learn to ignore it.
- Relative to human health risk assessments, ecotoxicity scientists are light years behind in implementation of SR and even further in AI. Ecotoxicity does not have the issue of Confidential Business Information (CBI).

## III.    Issues Related to People and Organizational Culture

- To address its shortage of data scientists, NIH created a data science strategic plan that other government agencies could use. NIH recruits data scientists to start at an entry or mid-level and work their way up. Someone who knows baseline AI but also understands the agency's problems and needs is preferred.
- It also is important to communicate with various vendors to gain different perspectives and to build partnerships with academics at all levels of the organization doing the assessments. There also is a need to also establish partnerships with other government agencies. International collaboration and outreach, such as through the International Collaboration for the Automation of Systemic Reviews (ICASR), also is important.
- A major challenge is how to recruit data scientists, who are expensive and hard to find. There are significant gaps in skills within EPA and in overseeing contractors. Internal training is

important -- beyond recruiting externally -- to allow EPA staff to better manage contracts and allow access to expertise available outside the agency. A data scientist today would be able to use tools to process data. The answer is not just more data scientists. Mid-level staff who can evaluate whether the approach is correct for the problem at hand is needed. One approach is to send the same problem to multiple contractors to see what comes back; this will allow EPA managers to improve their understanding and to compare contractors against one another.

- Partnering with universities for additional technical support was recommended. IBM works closely with universities. EPA works closely with other agencies such as NIEHS. EPA has a variety of vehicles for gaining access the necessary skill sets, such as ORISE Fellows.

## CROSS CUTTING ISSUES

**Summary of the Cross-Cutting Plenary Session**

- *Ethics:* For instance, choosing neighborhoods for which water overflows will be allowed.
- *Legal/Liability:* Can AI be questioned as would a human being, such as in the case of vehicle crash? This brings public participation, legal, transparency, accessibility, liability altogether.
- *Organizational:*
  - IT capacity: IT data science infrastructure may not be suitable for collaboration across organizations.
  - Budget constraints: There is a scarcity of funds to recruit talent, research new approaches.
  - Socializing the use of AI: Cultural change is needed to deploy AI in organizations.
  - Workforce implication: New skills may be needed for cross-domain knowledge of chemo-bio-informatics and other expertise.
  - Governance: Stakeholder and leadership acceptance of models is needed; there is asymmetry in capturing the benefits of AI.
- *Safety/Accountability:* The quality of data used to make informed decision can be questioned. Criteria and standards for data quality may be necessary.
- *Public Participation:*
  - Accuracy: If the question is what type of information to collect, there is the issue of bias or accuracy. For instance, AI recognizes white male faces much better than other gender and skin tones because of the quantity of white male faces used in a dataset to train the model.
  - Public acceptance: In the sensors case, the public participation includes specialized organizations or regulatory community that may provide the standards. Having the public involved to get acceptance of the system is necessary.
  - Trading off public participation and transparency: this happens when organizations develop an AI model (as for a chemical) that is fully validated but has intellectual property issues.

These issues will not be resolved soon. One approach is to apply AI in a promising area and employ data in a practical purpose under existing constraints. The next step is to look at cross-cutting issues to tackle. These cross-cutting issues should be considered in any further work on water quality or risk assessment.

**I.    Water Quality Next Steps**

- The workgroup participants -- Environmental Protection Agency, Maryland Department of the Environment, The Metropolitan Sewer District of Greater Cincinnati, Environmental Defense Fund, Intel, IBM -- are going to look into conducting an innovation challenge. It will involve people from different states and possibly from different watersheds.
- There are two relevant categories of AI: automation and modeling. The group needs to recognize the categories and decide which one of those will be included in the challenge.
- The application of AI needs to be explored as a group, walking through the steps from the beginning to the end, and maybe use it as a pilot to develop a framework for how to implement AI in water quality.
- The approach to be taken is through small steps: understand the problem and design a goal, then determine what incremental steps that can be taken to advance on the use of AI. If that is accepted by an organization, and if it saves time, money, and effort, then work toward a more robust model. These steps show small improvements and are more likely to be accepted.

**II.    Risk Assessment Next Steps**

- Complete training sets for the EPA/NIEHS Data Extraction Challenge for Systematic Review (The deadline is October 13-14). Share findings among workgroup members and the public.
- Continue the conversation among workgroup members and other interested stakeholders to develop a bank for extracted data so that it can be reused for other studies.
- Explore opportunities for setting standards for data extraction/labeling/organization processes, including definitions of the roles and responsibilities of various parties.
- Identify mechanisms for consulting with the publishing community regarding making publicly-funded data easily accessible after one year behind a payment wall. While this is currently the agency's policy, the publishers do not regularly make it easy to access the datasets.

**Partnership for Technology Innovation and the Environment: Workshop on Environmental Protection, Artificial Intelligence, and Machine Learning**

**Background Paper for the Water Quality Work Group**

This work group is assessing the use of AI/machine learning tools to improve water quality, with the goal of using AI tools and real-time data to manage the impacts of nutrients on water quality.

Many technologies and methods are in use to monitor and model large volumes of diverse data relevant to sources, fate, transport, and effects of nutrients in water. AI may improve the capacity to organize, assess, integrate and interpret complex data for improving nutrient management. The literature provides studies that document opportunity for better prediction accuracy and results using AI tools in some cases as opposed to traditional methods (attached for reference).

In preparing for the workshop, the work group has done the following:
- Developed illustrative use cases defining needs and opportunities for nutrient issues
- Considered AI tools that have application to water quality and nutrient management
- Decided to develop an innovation challenge based on a model EPA has used before
- Examined AI needs and applications in the pretreatment community, where data availability creates a more immediate environment for using AI tools
- Considered the Chesapeake Bay as a possible focus for an innovation challenge, based on the availability of extensive monitoring data and contacts with possible stakeholders
- Identified likely barriers to the application of AI to nutrient management issues
- Identified a range of potential stakeholders for further work in developing AI solutions

The work group has identified several illustrative use cases for matching up water quality and nutrient management needs with AI tools and expertise. Among these (as illustrations) are using cloud based automated technology to open and close pond outfalls to optimize flood storage and water quality; suing tools for optimizing the use of road salt during the winter season and minimizing environmental impacts from salt storage, handling, and application while protecting safety; and using methods for predicting low dissolved oxygen levels and emergence of algal blooms having ecological and health implications.

**Vision**

Our vision is to create an innovation challenge that links sources of AI technology and expertise with the needs defined by regulators, water quality managers, communities, treatment systems, and others. This innovation challenge is a model for applications in other water quality settings.

The larger vision is that AI tools and expertise will support water quality managers in making decisions based on diverse, reliable, real-time data and lead to better permitting and regulatory decisions.

**Challenges**

- *Availability and quality of data:* AI applications depend on having extensive data of a suitable quality. The needed level of quality will depend on the uses we have in mind. We discussed the Chesapeake Bay and the availability of data over time.
- *Experience and knowledge base of users of AI tools:* The likely users will require some degree of understanding of available tools, their strengths and limits, and the data needs.
- *Understanding by AI providers of water quality needs:* Water quality managers have identified a range of possible applications or use cases. They will need to educate AI providers to find the appropriate match between the tools and the water quality needs.
- *Regulatory and legal issues:* Different uses of AI will call for different levels of data quality. Permitting and compliance will demand a higher level of reliability than use of AI in identifying priorities and management setting priorities.
- *Developing the needed algorithms:* There will need to be enough quality data to be able to train machines and develop algorithms for applying AI tools to water quality settings.

**Actions and Synergy**

The work group plans to address the following questions in the workshop:

- What can we learn from any existing applications of AI to water quality management?
- What has to happen for AI to play a role in nutrient management within the next decade?
- What are the more significant barriers to using AI for water quality, especially nutrients?
- What are problems that AI could be used to improve water quality and nutrient quality?
- How might an innovation challenge begin to advance the state of AI use for nutrients?
- What would next steps be in creating an innovation challenge around AI and nutrients?

We also are eager to hear from the risk assessment work group to draw on their experience on methods, opportunities, barriers, and needs. Although the subject is different there are common lessons for each of the groups to draw upon. These lessons may include opportunities for using AI; a sense of the limits and role of AI in decision making; awareness of AI tools and methods; and of strategies for incorporating AI into diverse aspects of environmental decision making.

The work group provided several opportunities for synergy that we can use both at the workshop and in the follow-up activities. Among the perspectives were those of state and federal water

quality managers and data experts; technology providers and expertise; nutrient experts; and data scientists and managers. The workshop is an opportunity to develop these discussions in depth.

The work group agreed to approach the workshop as an opportunity to assess opportunities for and barriers to expanded use or AI tools in water quality and especially nutrient management. The goal at the workshop is not to design the innovation challenge but define the issues and opportunities and lay the groundwork for additional activity to prepare and launch a challenge.

**Suggestions**

Our advice is to use the workshop to create a platform for further discussions of the use of AI tools for managing water quality, especially in managing nutrients, which constitute the principal threat to water quality across the country. AI tools have many possible uses so support improved water quality but that many barriers, such as having sufficient quality, reliable data, do exist.

**Participants in the Water Quality and AI Work Group** (took part in at least one of the calls)

Dan Fiorino and Norie Ogata, American University
Denice Shaw, EPA Office of Research and Development
Robin Thottungal, EPA Office of Environmental Information
Kirsten Schroeder, IBM
Claude Yusti, IBM
Joe Rudek, Environmental Defense Fund
Lee Curry, Maryland Department of the Environment
Dinorah Dalmasy, Maryland Department of the Environment
Jennifer Molloy, EPA Office of Water
Dwane Young, EPA Office of Water
Steve Harper, Intel
Joe Greenblott, EPA Office of the Chief Financial Officer

**Partnership for Technology Innovation and the Environment: Workshop on Environmental Protection, Artificial Intelligence, and Machine Learning Risk Assessment Track**

# Overview

Environmental and occupational chemical exposures have the potential to negatively impact human health and environmental outcomes. Assessments of multiple studies or types of evidence are conducted to reach conclusions on potential hazards and risks due to chemical exposures. Approaches and tools for chemical risk assessment depend on the decision-making context and can range from screening-level assessments based on in silico evidence (often on data poor substances) to complex literature-based assessments that need to consider the results of hundreds of studies conducted in humans, animals, and cell-based systems.

Across the chemical risk assessment space, there are many opportunities to implement artificial intelligence (AI) and automation including:

- Machine Learning (ML) and Natural Language Processing (NLP) to semi-automate the process of searching the literature, screening studies for relevance, evaluating the quality of individual studies, and summarize study methodology and findings ("data extraction").
- Labelling and organizing of extracted data for enhanced querying across data sources according to biological process to inform risk assessments.
- AI approaches to use data from one chemical (or group of chemicals) to inform the assessment on another substance including quantitative substance activity relationships (QSAR) and read-across approaches.
- In addition to enhancing the speed and accuracy of an assessment itself, AI can contribute to the prioritization of substances and chemical testing at the front and back end of a risk assessment.

# Vision

The overall vision of the Risk Assessment Track is to create an interoperable chemical risk assessment platform for performing sophisticated data querying simultaneously from multiple data sources (e.g. AOPkb, CEBS, HAWC, HERO, IUCLID, RapidTox, ToxRefDB, etc.). Inputs include published scientific studies and gray literature in both structured and unstructured formats.

We anticipate the modules depicted in Figure 1 can be addressed by case studies that include:

- ML and evidence mapping to enhance problem formulation, scoping, and planning based on title and abstract screening. Note that AI tools exist and are in use.
    Implementation of ML, NLP, and evidence mapping during full text screening;
    Summarizing study methods and findings from full text "data extraction";
    Automating point-of-departure assessments (benchmark dose modeling);
- NLP to label data matched to ontology concepts;
- Interoperability promoted through ontology concepts coded by key, source, and date;
    Automation to increase data labeling efficiency;
    NLP and ML to enhance study quality evaluation efficiency and consistency;
    Clustering around chemical and biological inventories using unsupervised learning approaches;
    Read-across approaches (e.g., building QSARs);
    Identifying next steps for testing (prioritizing); and
    Identifying safer alternatives based on chemical/biological structure.

# Challenges

Deployment of AI across the regulatory and research chemical risk assessment space faces technological, policy, organizational and many other barriers including:
- *Identification of appropriate data sets;
- Common principles for representing chemical, toxicological, and exposure data;
- AI standards (metrics of acceptable performance), interoperability, and evolution;
- Expert (cross-domain knowledge of chemo-/bio-informatics, AI, and chemical risk assessment) engagement;
- IT data science infrastructure suitable for collaborating across the EPA, other US federal agencies, and globally;
- Data accessibility — many pieces of a company, EPA, and other organizational data are controlled by trade secret statutes or organizational silos;
- Ownership and rights to underlying data that underpin models;
    Stakeholder/political/scientific acceptance of statistical models;
    Scarce dollars to recruit talent, research new approaches, and implement research in regulatory contexts.

# Actions

*For the sake of focusing the workshop discussions, the issue of AI to help in the data extraction process and to label and organize the extracted data will be the priority topics. We will focus on experiences and challenges encountered to date in developing appropriate data sets that may also include:*
- Continued collaboration to develop training and test data sets for model development;
- Developing an approach for leveraging staff resources;
- Leveraging public and private funding models;
- Task based workgroups with access to resources.

# Contributors (alphabetical)

Michelle Angrish, EPA
Dan Fiorino, American University
Joe Greenblott, EPA
Kristan Markey, EPA
Norie Ogata, American University
Dave Rajewski, Environmental Law Institute
Seema Schappelle, EPA
Michele Taylor, EPA
Kris Thayer, EPA
Danielle Wagner, American University
Claude Yusti, IBM

# References/Resources

Luechtefeld, T., et al. (2018). "Machine Learning of Toxicological Big Data Enables Read-Across Structure Activity Relationships (RASAR) Outperforming Animal Test Reproducibility." Toxicological Sciences **165**(1): 198-212.

Wilkinson, M. D. and M. Dumontier (2016). "The FAIR Guiding Principles for scientific data management and stewardship." **3**: 160018.

Jaspers., et al. (2018). "Machine Learning Techniques for the automation of literature reviews and systematic reviews in EFSA." European Food Safety Authority.