

March 13th

Computational approaches to Assessing Radiation Exposure Across the Lifecourse: Risks, Insights, and Innovations.

### Heidi Hanson

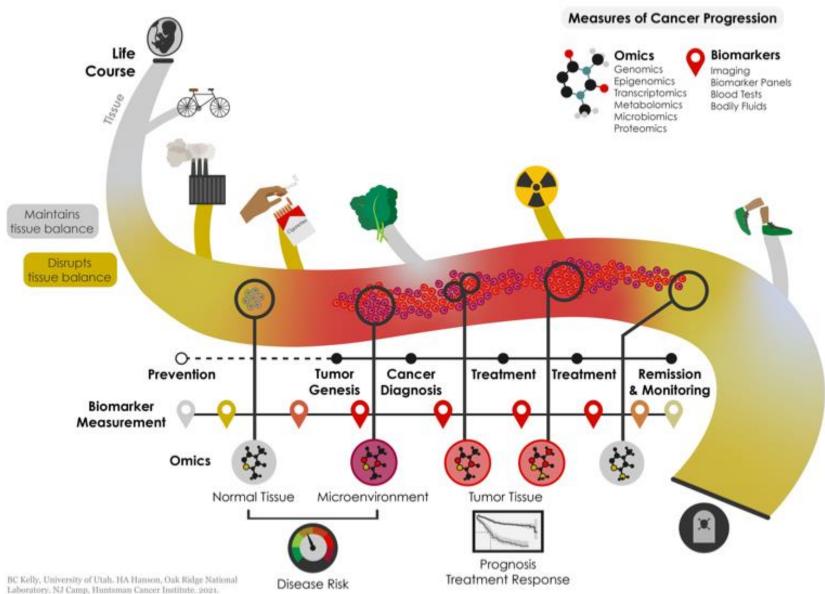
Computational Sciences and Engineering Division



ORNL IS MANAGED BY UT-BATTELLE LLC FOR THE US DEPARTMENT OF ENERGY



# **Life Course Approach to Radiation Exposure**





# Charting the Life Course: An Interdisciplinary Blueprint

5 Pathways Forward











1 Integrate datasets across the lifespan.



Use transdisciplinary principles and methods, borrowing from:



Engineering



Statistics



Communication



Data science cores focus on curation, cleaning, analysis, and modeling of data.









3 New ways to model highdimensional data.

4 Expand traditional epidemiological methods to include systems and network modeling.



Improve reproducibility with standard operationalizing measures.



Heidi Hanson and Shari Barkin

# Charting the Life Course: An Interdisciplinary Blueprint 5 Pathways Forward

# **Workflows for AI Ready Data and AI for Data Harmonization**





Integrate datasets across the lifespan.



Use transdisciplinary principles and methods, borrowing from:



Engineering





Statistics

Communication





Mathematics Epidemiology

Data science cores focus on curation, cleaning, analysis, and modeling of data.









New ways to model highdimensional data.

Expand traditional epidemiological methods to include systems and network modeling.



Improve reproducibility with standard operationalizing measures.



# Integrated Electronic Health, Administrative, and Population Data

Computational tools enabling decision making for risk assessment, prevention, and treatment

- Easy management of large amounts of heterogenous and siloed sensitive information
- Assessment of longitudinal trajectories of radiation exposure from residential radon, medical treatments, and occupational exposures.
- Identification of geographical areas at risk for high residential radon exposure
- Data-driven tools to guide decision making in the clinic and community

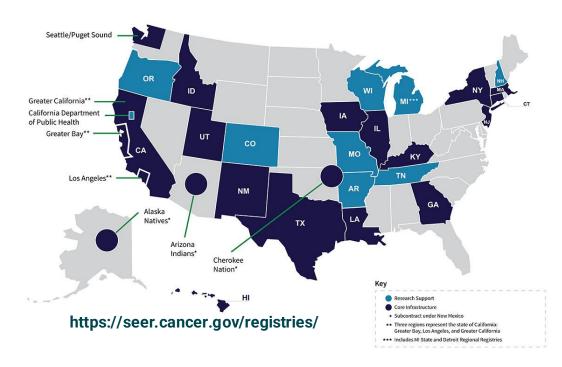






# Real World Evidence: Al for Near Real-Time Health Surveillance Covering 48% of the US Population

Surveillance Epidemiology End-Results (SEER) Registries > 850,000 Diagnoses Annually



Decision Making Tool Built into SEER Data Management System (2021)

# Auto-Extraction from Pathology Reports:

Accuracy: Auto-coding of 23-27% of path reports (N ~ 230,000) with > 98% accuracy across all data elements.

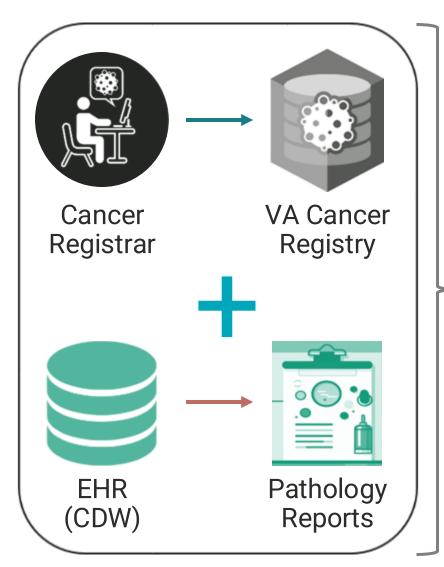
- Phenotype classifications:
  - Site = 70 categories
  - Sub-site = 324 categories
  - Histology = 626 categories
  - Laterality = 7 categories
  - Behavior = 4 categories

Production implementation Hierarchical Self Attention Model (HiSAN) with Deep Abstention:

- ■Total 16 registries (additional 4 since March 2022; ~31% of US population)
- ■Default as part of any new DMS installation, regardless of SEER affiliation
- •5+ new registries anticipated in 2023/2024
- Testing phase with the Veteran's Health Administration



# Real World Evidence: Clinical Decision Support with EHR Data





**FrESCO** 

- Classification (Reportability)
- Information Extraction
- Histology
- Behavior
- Site
- Laterality



#### Near Real-time Case Ascertainment

- National Quality of Care (Systems of Excellence)
- Tele-Oncology
- Cancer Registry



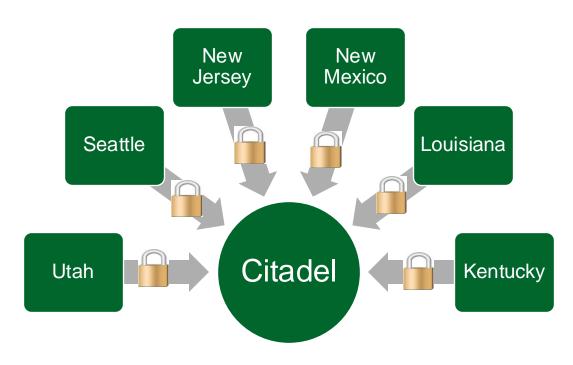






## Interoperable Al Systems at Scale:

Bringing Compute Closer to the Source Rapid Data Integration while Protecting Private Information



The National AI Research Resource (NAAIR) Task Force National Childhood Cancer Registry DOE ASCR Biopreparedness (BRaVE) Funding

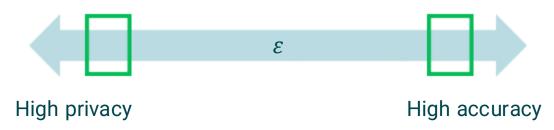
# Our current implementation

- Cross Silo
- Horizontal
- Model Centric Need for Cooperative Agreement between institutions/ participants

# Trusted host Differential privacy

"Epsilon Indistinguishability"

Privacy/accuracy tradeoff





**Innovation**: SEER data make us completely innovative in this space. We can design and test solutions for <u>real-world application</u> at population scale.

# MOSSAIC Computational Workflows Rapid Response for Health Surveillance

Standard data processing workflows across siloed systems for preparing unstructured text

DuckDB

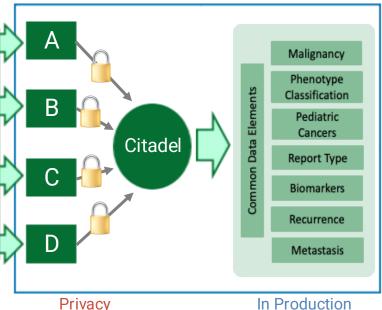
**Portable Data Engineering** 

Process NAACCR formatted XMLs and

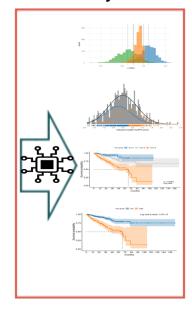
HL7s + other source files

In Production

Federated Learning for extracting common data elements



Downstream Analysis



Privacy Preserving FL: In R&D

In Progress:

>>> ARROW

Multimodal/longitudinal
DNN, Phrase Level Attention
DNN, Tools for Bias
Detection, Exposomic Data
Linkage, Privacy Preserving
FL, Synthetic Clinical
Reports

bardi

**NLP Pre-Processing** 

In Production

Available Deep Learning
Tools
MT-CNN, MT-HiSAN,
Clinical BigBird,
PathBigBird, Deep

**Abstaining Classifier** 

FrESCO

https://github.com/DOE-NCI-MOSSAIC



Silo A

Silo B

Silo C

Silo D

<xml />

CSV



Heidi Hanson, Josh Grant, Maggie Davis,

#### Raster Data

- Satellite images, aerial photographs, and derived raster products

#### Vector Data

 Geographic features such as boundaries, roads, and points of interest

#### Tabular Data

Meta data, user information, analytical results, modeled data

#### Large Source Files

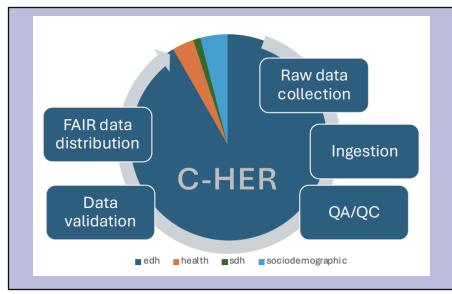
Raw data files sored in MinIO

#### Data Ingestion Code

The code used to ingest the data is also stored in the database, allowing us to quickly refresh data and identify issues

**OAK RIDGE**National Laboratory

Metadata and data storage: Metadata is captured with Prefect, which has been integrated with the Centralized Health and Environmental Repository (C-HER) tools. This provides dynamic metadata extraction and ensures robust error handling, making metadata documentation reliable and scalable. Data security is ensured with encryption for data at rest and in transit. We implement strict access controls and authentication measures to safeguard our data.



The C-HER Ecosystem utilizes
Open-Source software such as
MinIO, PostgreSQL, Prefect, and
Python. The goal of the project is
to optimize the reuse of data, thus
the C-HER ecosystem adheres to
principles of findability,
accessibility, interoperability, and
reusability (FAIR).



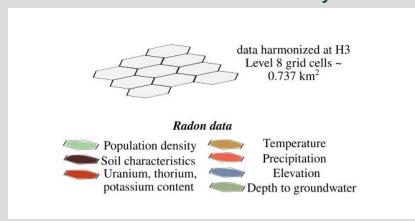
<u>Data Quality Assurance</u>: Regularly validate data using automated tools and manual checks to ensure accuracy and reliability. Document all QA-QC procedures and maintain logs of all data checks.

<u>Ethical Concerns</u>: We consider and address privacy concerns for data that is potentially identifiable. We ensure that data usage complies with all applicable laws and ethical guidelines.

# C-HER Serves as the foundation for modeling exposures

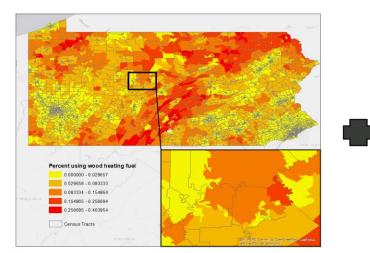
**Example: Indoor Radon** 

### C-HER Protocol for AI Ready Data

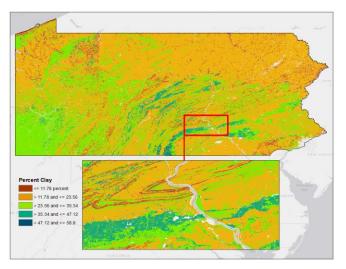




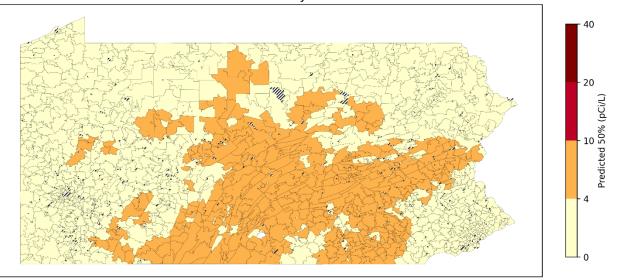
#### ZCTA level housing Characteristics Individual test results with ZCTA information



Low spatial resolution geologic, soil, and other data



#### Predicted 50% by ZCTA



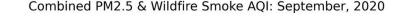
# 25+ years of multimodal exposure data linked to cancer records

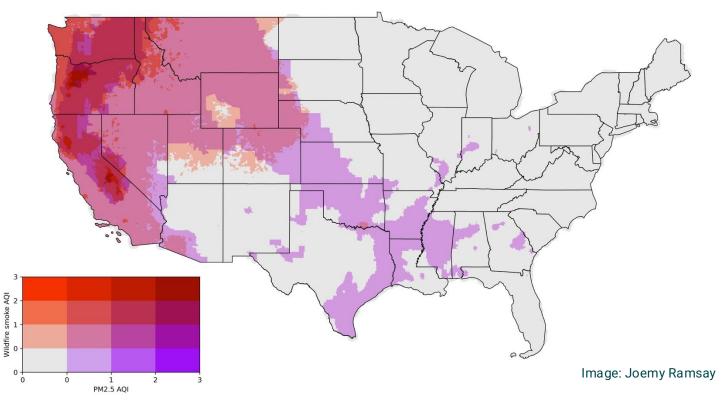
# M O S S A I C Modeling Outcomes Using Surveillance Data & Scalable Artificial Intelligence for Cancer



### LexisNexis Residential History data

- 11 SEER registries have been linked (3.2 million individuals diagnosed from 2005 2022)
- 15 SEER registries should be linked by the end of 2023
- Residential history constructed for Louisiana and expected to be expanded to all linked registries this year
- High quality data from 1995 2020
- 83% are geocoded to the point location



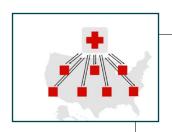




"Electronic Health Record informed Lagrangian method for precision public Health (EHRLICH)" will be a set of new computational capabilities developed to enable the rapid assimilation of real-world data for digital twins of population level biological threats.



Scalable, Trustworthy and Responsible Al



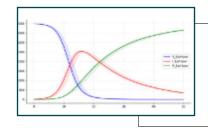
Privacy-Aware Simulation for Clinical AI and Learning (PASCAL)



Al Ready Social and Environmental Determinants of Health Data



High-spatial resolution and high-fidelity synthetic populations that include demographic and biologic characteristics



ENABLE: An agent-based simulation framework

















# Interoperable Al Systems at Scale:

Bringing Compute Closer to the Source Rapid Data Integration while Protecting Private Information

# Alternative approach: synthetic data generation

#### Our Solution:

- Leverage generative Al
- Minimize leakage of PHI & other identifiable info
- Evaluate for utility, similarity, and privacy

Initial demonstration on SEER data

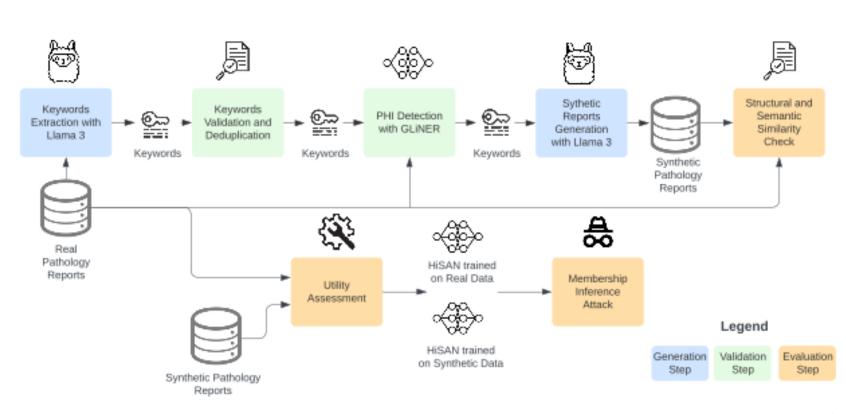


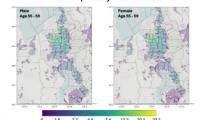
Image credit: Patrycja Krawczuk



# EHRLICH tools are designed to be agile and adaptable

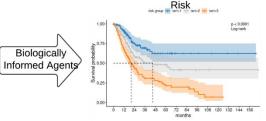
# Data-Driven Characterizations of Agents and their Environments

Comorbid Conditions that Increase Susceptibility to Disease



### High-Fidelity Synthetic Populations

Accute Identification of Populations at



Health's not Eventy Distributed in Space

Physical, Social, and Chemical Environments



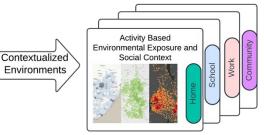


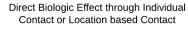


Physical, Social, and Chemical Environments Determine Risk



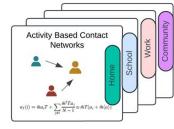










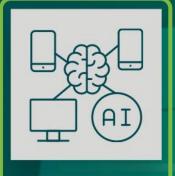








# PLACING A PULSE ON POPULATION HEALTH



Federated Learning



Identify existing patterns and anomoly detection



Synthetic Data Generation



Predict patterns relating to environmental exposures



Privacy Preserving Al



High-performance computing for modeling and simulating health outcomes



## **MOSSAIC**

#### **National Cancer Institute**

Lynne Penberthy (PI)
Elizabeth Hsu (Technical Lead)
Valentina Petkov
Serban Negoita
Ola Adeyemi
Sylkk Ansah
Sarah Bonds

#### **IMS**

Linda Coyle Jennifer Stevens Scott Depuy Rusty Sheilds Gary Beverungen

Los Alamos National Laboratory
Jamaludin Mohd Yusof
Sayera Dhaubhadel
Tanmoy Bhattacharya

#### Oak Ridge National Laboratory

Heidi Hanson (PI)
John Gounley (Technical Lead)
Shelaine Curd (Program Manager)

Georgia Tourassi

Joe Lake

Adam Spannaus

Dakota Murdock

**Zachary Fox** 

Patrycja Krawczuck

**Dakotah Maguire** 

Jordan Miller

Mayanka Chandra Shekar

Noah Schaefferkoetter

Sajal Dash

Isaac Lyngaas

Abhishek Shivanna

**Robert Bridges** 

**Christopher Stanley** 

Vandy Tombs

**Christoph Metzner** 





## **EHRLICH**



















#### **Oak Ridge National Laboratory**

Heidi Hanson (PI)
John Gounley (Technical Lead)
Shelaine Curd (Program Manager)
Adam Spannaus
Zachary Fox
Patrycja Krawczuck
Dakotah Maguire
Mayanka Chandra Shekar
Robert Bridges
Christopher Stanley
Vandy Tombs
Sudip Seal
James Nutaro
Sifat Moon

#### **Los Alamos National Laboratory**

Jamal Mohd-Yusof Cristina Garcia Cardona

**Christoph Metzner** 

#### **Argonne National Laboratory**

Rick Stevens Thomas Brettin

#### **University of Utah**

James VanDerslice Joemy Ramsay

#### **University of Arizona**

Nirav Merchant Ravi Tandon

#### **University of Southern California**

Rima Habre

#### **Duke**

Amanda Randles

## University of Chicago/Morehouse School of Medicine Lilly Immergluck

