**AACR** American Association for Cancer Research®

ANNUAL MEETING

2024 • SAN DIEGO

APRIL 5-10

#AACR24
AACR.ORG/AACR24

# Year of Open Science: Impact of the Cancer Research Data Commons

Erika Kim, PhD

Ina Felau, MS

Esmeralda Casas-Silva, PhD

Anthony Kerlavage, PhD

Center for Biomedical Informatics and Information Technology

National Cancer Institute, NIH, US

CANCER RESEARCH

The Foundational Cancer Journal
Driving Transformative Science

AACRJournals.org
@CR_AACR

**AACR** American Association for Cancer Research®

2302007fn

# Overview: NCI Cancer Research Data Commons (CRDC)

Ina Felau, MS

Health Science Administrator

Informatics and Data Science Program

Center for Biomedical Informatics and Information Technology

National Cancer Institute, NIH, US

Ina Felau, MS

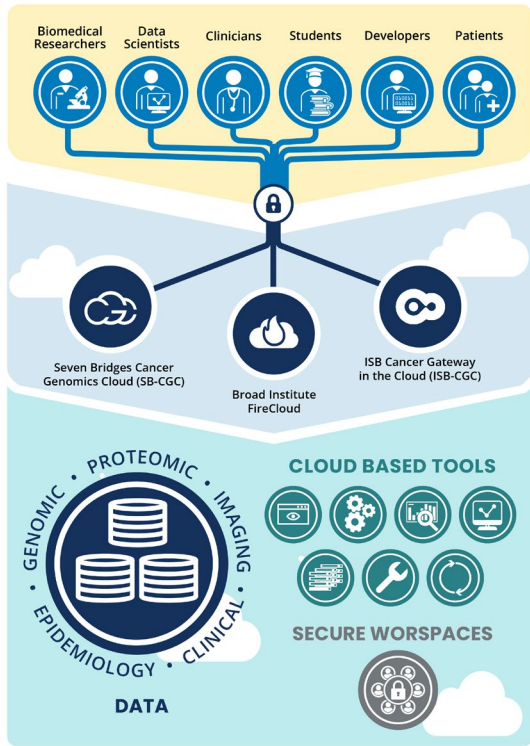I am a full-time paid employee of the NIH/NCI.
I have no financial relationships to disclose.

# Agenda

- **Overview: CRDC Vision**

- CRDC Ecosystem
  - Data Commons
  - Cloud Resources

- Interoperability Initiatives
  - Initiatives within CRDC
  - Trans-NIH initiative (NCPI)

- Resources for Researchers

NIH NATIONAL CANCER INSTITUTE
Cancer Research Data Commons

datacommons.cancer.gov

# Vision for the CRDC

Biomedical Researchers • Data Scientists • Clinicians • Students • Developers • Patients

Seven Bridges Cancer Genomics Cloud (SB-CGC)
Broad Institute FireCloud
ISB Cancer Gateway in the Cloud (ISB-CGC)

GENOMIC • PROTEOMIC • IMAGING • CLINICAL • EPIDEMIOLOGY

CLOUD BASED TOOLS

SECURE WORSPACES

DATA

= Secure data storage and access

## Mission

- Empower researchers by providing a secure, accessible cancer data ecosystem
- Provide state-of-the-art visualization, analysis, and interoperability tools in a flexible, cloud-based computational environment

## Lower Barriers

- Data submission
- FAIR data access, search, retrieval
- Integration of data for cross-domain analysis
- Analysis platforms, tools, and workflows

## Infrastructure & Sustainability

- Security and appropriate access for sensitive data
- Sustainable, reusable, and uniform architecture
- Comprehensive plan for long term data storage and accessibility to tools

## Stakeholder Focus

- Include all scientists and clinicians (of all technical abilities) using the data
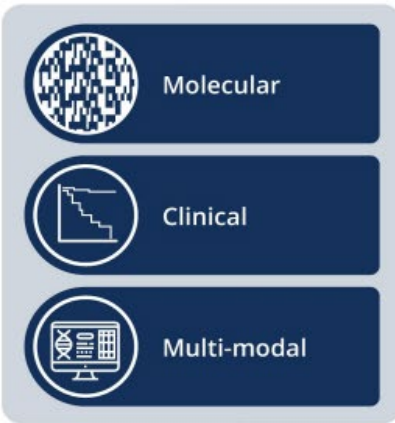
# FAIR Principles



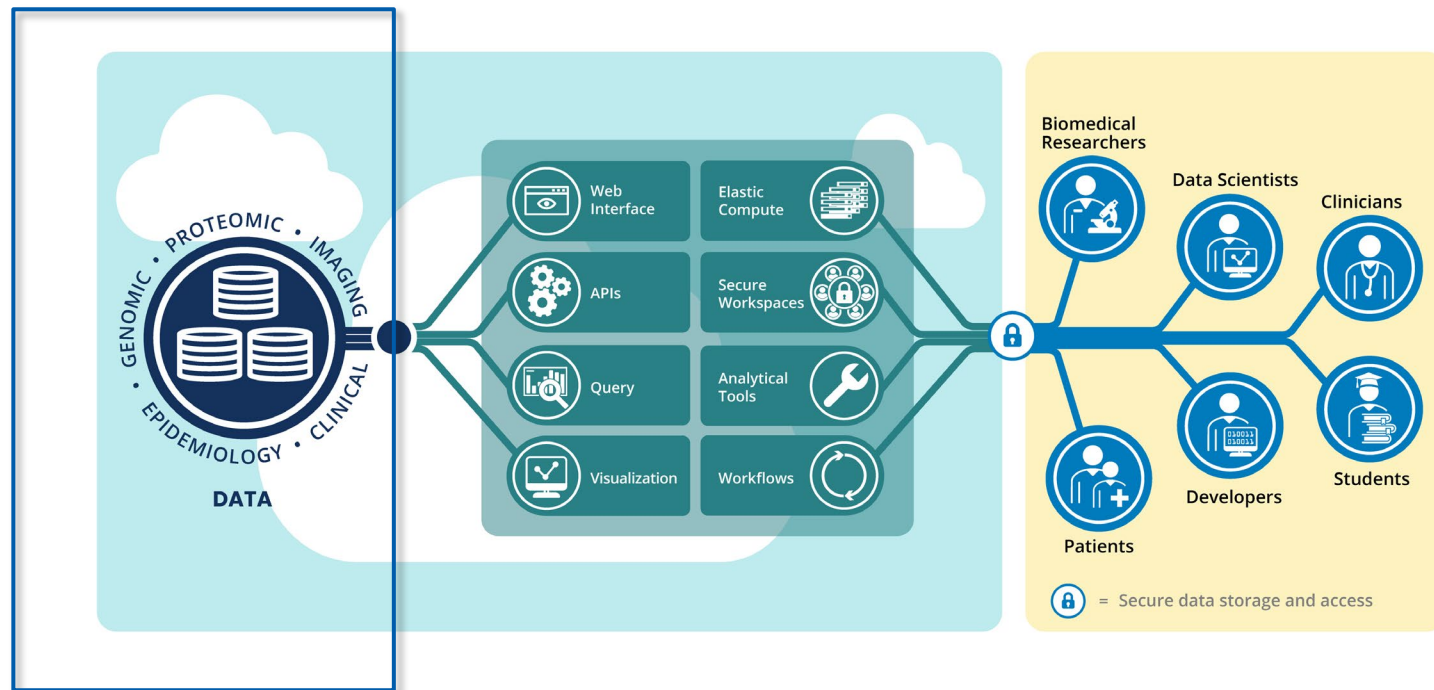## FINDABLE
Faceted Search and Key Word Search

Cases | Clinical | Genes | Mutations
∨ Search Cases ❓
🔍 e.g. TCGA-A5-A0G2, 432fe4a9-2...
Upload Case Set ▾
> Primary Site 🔍
> Program
> Project 🔍
> Disease Type 🔍

## ACCESSIBLE
Online Analysis and Visualization

- Molecular
- Clinical
- Multi-modal

## INTEROPERABLE
APIs and Standardized Metadata

- NCPI — NCPI & dbGaP
- Cancer Data Aggregator
- R — 3rd party R packages

## REUSABLE
Rich Metadata and Harmonized Scientific Data

- Raw, secondary and tertiary data
- Harmonized by uniform pipelines
- Customizable Cohorts

At the bottom are six specialized data commons (GDC, PDC, ICDC, CDS, IDC and CTDC). Selected CRDC features are used to demonstrate the implementation of FAIR principles.

# Agenda

- Overview: CRDC Vision

- **CRDC Ecosystem**
  - **Data Commons**
  - **Cloud Resources**

- Interoperability Initiatives
  - Initiatives within CRDC
  - Trans-NIH initiative (NCPI)

- Resources for Researchers

# CRDC Ecosystem



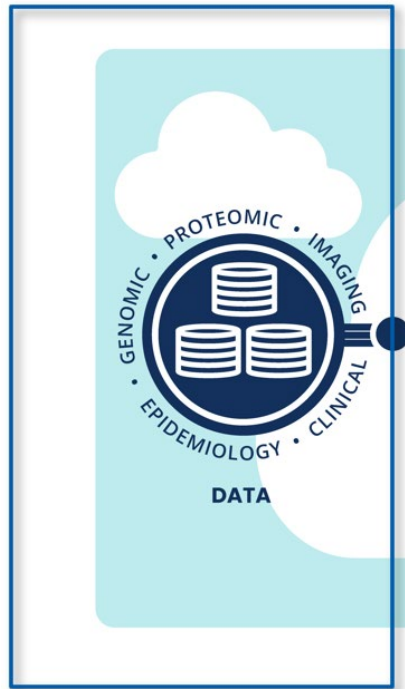https://datacommons.cancer.gov/

# CRDC Ecosystem



Data Commons

# CRDC Data Commons



Data Commons

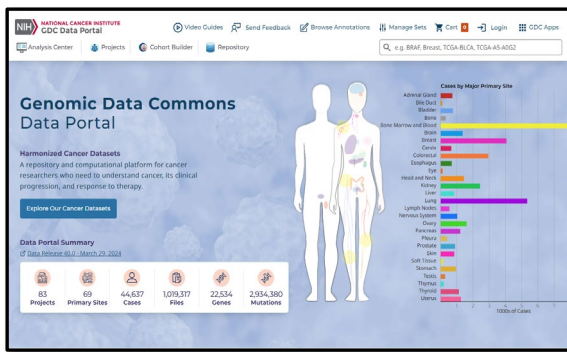Genomic Data Commons

Proteomic Data Commons

Imaging Data Commons
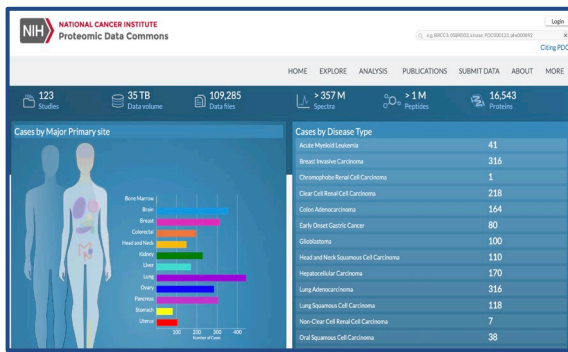
Integrated Canine Data Commons

Cancer Data Service

# CRDC Data Commons



## Genomic Data Commons

- Share, analyze, and visualize genomic data
- Harmonized to the same genome standard and variant calling pipeline

https://portal.gdc.cancer.gov/

## Proteomic Data Commons

- Filter, query, search, visualize and download proteomic data and metadata
- Data harmonization pipeline to uniformly analyze all PDC data

https://pdc.cancer.gov/

## Imaging Data Commons

- Share, analyze, and visualize de-identified multi-modal imaging data, as medical images (MRI, PET, CT)
- Uses DICOM standard

https://imaging.datacommons.cancer.gov/

# CRDC Data Commons

## Integrated Canine Data Commons

- Share data from canine clinical trials
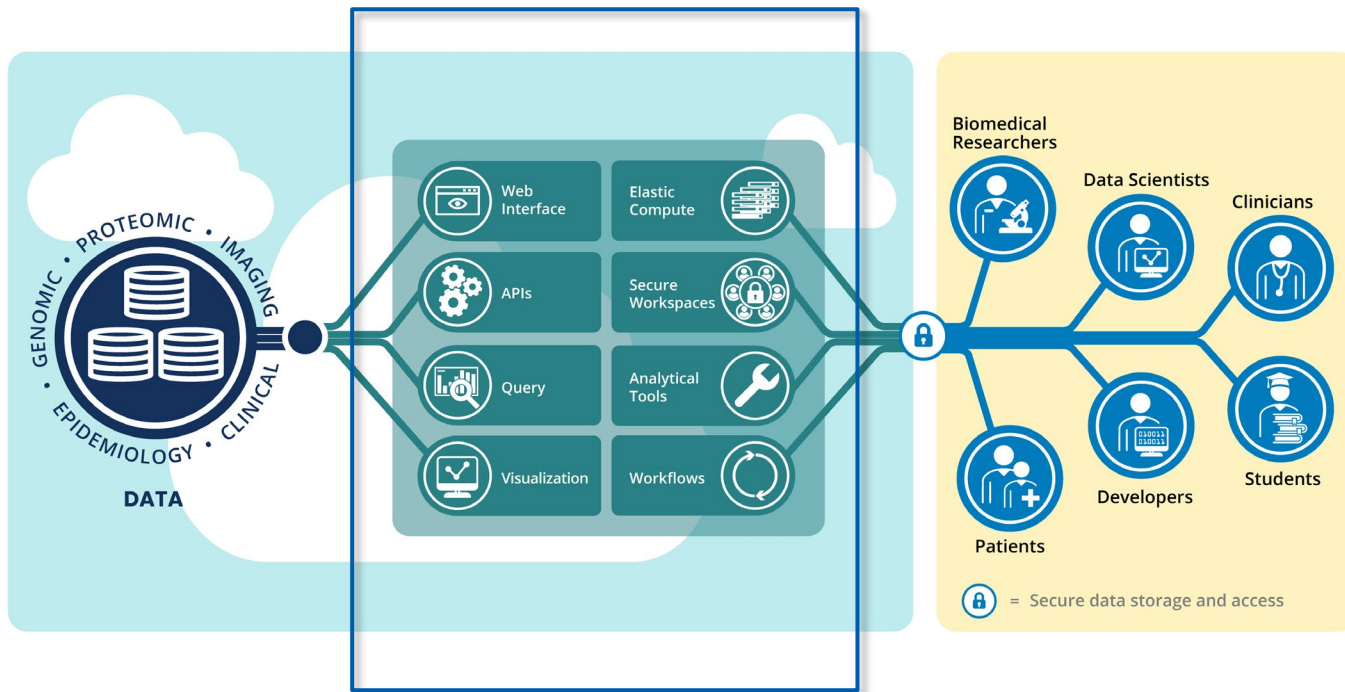- All data (including raw sequence data) are open access.

https://caninecommons.cancer.gov/



## Cancer Data Service

- Access NCI-funded data currently not hosted by other CRDC data commons
- All datatypes accepted

https://dataservice.datacommons.cancer.gov/

*Data types available across data commons: WGS, WXS, RNAseq, miRNA-seq, scRNAseq, ATAC-seq, DNA methylation, mass spectrometry-based proteomic data, DICOM.*

# CRDC Ecosystem



Cloud Resources

https://datacommons.cancer.gov/

# CRDC Cloud Resources

**Broad FireCloud (FC), powered by Terra**
- Based on the Google Cloud Platform (GCP)
- Offers extensive repositories of pre-built tools and workflows in the Workflow Definition Language (WDL).

**The ISB Cancer Gateway in the Cloud       (ISB-CGC)**
- Offers Google Cloud Platform (GCP) native tools and Google BigQuery for big data analytics and Google Compute Engine for complex workflow execution.
- Designed for users looking to use derived data.

**The Seven Bridges Cancer Genomics Cloud (SB-CGC), powered by Velsera**
- Based on the Amazon Web Services (AWS) platform
- Offers a curated library of over 850 tools and workflows optimized for the cloud using the Common Workflow Language (CWL).

| | |
|---|---|
| Eliminate the need to download data | Access to workspaces, analysis tools, workflows & pipelines |
| Bring your own data and tools: collaborative pre-publication workspaces | Integrate your data with other CRDC data and tools in the cloud |

# Agenda

- Overview: CRDC Vision

- CRDC Ecosystem
  - Data Commons
  - Cloud Resources

- **Interoperability Initiatives**
  - **Initiatives within CRDC**
  - **Trans-NIH initiative (NCPI)**

- Resources for Researchers

# CRDC: Interoperability Needs for Cancer Data

**Challenge:** Access comprehensive datasets like TCGA and CPTAC from multiple repositories for integrative analysis

- Discover relevant datasets *across multiple resources using common standards*

- *Aggregate and analyze data* housed in separate data repositories using latest analytical tools

# CRDC: Internal Interoperability Projects

## CRDC Data Standards Services (DSS)

- Semantic harmonization across CRDC datasets
- Shared data models for submission & search
- Leverage existing standards, eg NCIt



## CRDC Cancer Data Aggregator (CDA)

- Search by harmonized, common language terms to aggregate data distributed across CRDC repositories
- Get information about subjects, files or specimens in a standard tsv format that can be opened in Excel, integrated into a pipeline or uploaded to a cloud resource
- cdapython available via interactive browser, notebooks or local install

# Agenda

- Overview: CRDC Vision

- CRDC Ecosystem
  - Data Commons
  - Cloud Resources

- Interoperability Initiatives
  - Initiatives within CRDC
  - Trans-NIH initiative (NCPI)

- **Resources for Researchers**

# CRDC Website

https://datacommons.cancer.gov/

# Support for Researchers

https://datacommons.cancer.gov/support-for-researchers

| CRDC COMPONENT | RESOURCE | AVAILABLE SUPPORT |
|---|---|---|
| Cloud Resources | Broad FireCloud Powered by Terra (FC) | • How to Set Started on FC⧉<br>• Broad Institute FireCloud Workshop Tutorials⧉<br>• Terra Self-Service Learning Resources⧉<br>• Broad Institute FireCloud FAQs⧉ |
| | ISB Cancer Gateway in the Cloud (ISB-CGC) | • How to Get Started on ISB-CGC⧉<br>• ISB-CGC FAQs⧉<br><br>The ISB-CGC team offers virtual office hours through Google Meet. Note that the link is different for each of the days.<br><br>• Tuesdays at 2:00 pm ET; Link: http://meet.google.com/jkg-cxke-yzs⧉<br>• Thursdays at 11:00 am, ET; Link: http://meet.google.com/jai-kgkg-sii⧉<br><br>Find tutorials and user guides on the ISB Cancer Gateway website⧉. |
| | Seven Bridges Cancer Genomics Cloud (CGC), powered by Velsera (SB-CGC) | • How to Get Started on SB-CGC⧉<br>• SB-CGC Introduction to the CGC Webinar⧉<br>• SB-CGC Scaling Single-Cell Research⧉<br>• SB-CGC Troubleshooting Tutorial⧉<br><br>The Seven Bridges team offers virtual office hours through Google Meet at https://meet.google.com/kbs-ojnj-dcg⧉ at the following times:<br><br>• Tuesdays at 10:00 am ET<br>• Thursdays at 2:00 pm ET<br><br>Learn more about SB-CGC through their user guides, video tutorials, and webinars⧉. |

| CRDC COMPONENT | RESOURCE | AVAILABLE SUPPORT |
|---|---|---|
| Data Commons | Genomic Data Commons (GDC) | • How to get started on GDC⧉<br>• GDC Webinars⧉<br>• GDC FAQs⧉<br>• NGS Studies of Familial Data Using Cloud Computing⧉ |
| | Proteomic Data Commons (PDC) | • PDC FAQs⧉<br>• NCI OCCPR Webinar on PDC⧉ |
| | Imaging Data Commons (IDC) | The IDC offers community office hours every week through Google Meet at https://meet.google.com/xyt-vody-tvb⧉.<br><br>• Tuesdays, 16:30 – 17:30 (ET/New York)<br>• Wednesdays, 10:30-11:30 (ET/New York)<br><br>Learn more from the IDC user guide and white papers⧉.<br><br>If you have questions about the IDC, email the team at support@canceridc.dev or start a thread in their online forum⧉. |
| | Integrated Canine Data Commons (ICDC) | If you have questions, please email the ICDC team at: ICDCHelpDesk@mail.nih.gov. |
| | Cancer Data Service (CDS) | If you have questions, please email the CDS team at: CDSHelpDesk@mail.nih.gov. |

| CRDC COMPONENT | RESOURCE | AVAILABLE SUPPORT |
|---|---|---|
| Infrastructure | Cancer Data Aggregator (CDA) | If you have questions, contact the CDA team through the CDA Helpdesk⧉ |

# CRDC Insights: External Newsletter (quarterly)

https://datacommons.cancer.gov/crdc-insights

# Overview:
# NCI Cancer Research Data Commons (CRDC)



datacommons.cancer.gov

# CRDC Impact and Success Stories

Esmeralda Casas-Silva, Ph.D.

Health Science Administrator

Center for Biomedical Informatics and Information Technology

National Cancer Institute, NIH, US

# Esmeralda Casas-Silva, Ph.D.

I am a full-time paid employee of the NIH/NCI.
I have no financial relationships to disclose.

**Problem:**

- Higher recurrence and drug resistance in PTC subset
- No reliable way to predict which PTCs will progress

**Goal:**

- Identify molecular risk factors for papillary thyroid cancer progression
- Develop PTC genomic classifiers, stratify patients based on recurrence risk
  - Improve prognosis predictions

frontiers | Frontiers in Endocrinology

A clinically useful and biologically informative genomic classifier for papillary thyroid cancer

Steven Craig[1,2†], Cynthia Stretch[3†], Farshad Farshidfar[3], Dropen Sheka[4], Nikolay Alabi[4], Ashar Siddiqui[5], Karen Kopciuk[3,6], Young Joo Park[7,8], Moosa Khalil[9], Faisal Khan[9,10], Adrian Harvey[11] and Oliver F. Bathe[3,11,12*]

# The Cancer Genome Atlas (TCGA)



33 Tumor types
10 Rare cancers
>11,000 patients

2.5 Petabytes of data

Genomic Data Commons
Proteomic Data Commons
Imaging Data Commons

- ✓ Genomic
- ✓ Transcriptomic
- ✓ Epigenomic
- ✓ Imaging
- ✓ Clinical
- ✓ More…

Map genetic mutations across cancers

**CRDC Success Story 1:** Using GDC to Advance Papillary Thyroid Cancer (PTC) Care Through Genomic Classifiers

AACR
American Association for Cancer Research®
ANNUAL MEETING
2024 • SAN DIEGO
APRIL 5-10 • AACR.ORG/AACR24 • #AACR24

**Approach:**

- Leverage GDC to access cases from landmark TCGA study detailing genomic profile of PTC
  - Transcriptional data
  - Copy Number Variation
  - Methylation status

- Apply machine learning to 500+ cases
  - Stratify into molecular subtypes based on recurrence risk

- Used TCGA methylation data in GDC to explore epigenetic differences between cancer subtypes

**Key findings:**

- 3 unique molecular subtypes
  - Subtype 1: Lowest recurrence rate
    - Lower BRAFV600E, higher RAS mutations
  - Subtype 2: Moderate recurrence
    - Higher BRAFV600E, inflammatory, EMT pathways
  - Subtype 3: High recurrence rate
    - immunosuppressive microenvironment, high EZH2-HOTAIR pathway, BRAFV600E and TERT promoter mutations

**Impact:**

- Genomic classifiers outperformed the American Thyroid Association's clinical risk stratification system

**Problem:**
- Immunotherapy only successful in a small proportion of cancer cases

**Goal:**
- Develop comprehensive understanding TME across cancers
- Reveal immune cell surveillance and tumor immune evasion mechanisms

**Cell**

Resource

**Pan-cancer proteogenomics characterization of tumor immunity**

Francesca Petralia,[1,36,*] Weiping Ma,[1,36] Tomer M. Yaron,[2,3,4,36] Francesca Pia Caruso,[5,33,36] Nicole Tignor,[1,36] Joshua M. Wang,[6,7,36] Daniel Charytonowicz,[1,37] Jared L. Johnson,[2,8,9,37] Emily M. Huntsman,[2,3,37] Giacomo B. Marino,[10,37] Anna Calinawan,[1,37] John Erol Evangelista,[10] Myvizhi Esai Selvan,[1,12] Shrabanti Chowdhury,[1] Dmitry Rykunov,[1] Azra Krek,[1] Xiaoyu Song,[11,12] Berk Turhan,[1] Karen E. Christianson,[13] David A. Lewis,[10] Eden Z. Deng,[10]

# NCI **C**linical **P**roteomic **T**umor **A**nalysis **C**onsortium (CPTAC)

- Government, academia, industry partnership
- 1500+ patients, 10 tumor types
- Mass spectrometry proteomic data (PDC)
  - Protein expression
  - Post-translational modifications
  - Protein-protein interactions
- Genomic Data (GDC, CDS)
  - WGS, WXS
  - RNA-seq
  - CNV

- Clinical Data (GDC, PDC)
  - Treatment outcomes
  - More
- Imaging Data (IDC)



Image: Clinical Proteomic Tumor Analysis Consortium, National Cancer Institute.
https://proteomics.cancer.gov/programs/cptac

## Approach:

- Combine CPTAC data from across CRDC

- Analyze genomic, epigenetic, transcriptomic, and proteomic alterations across tumors

- 1,056 tumor samples,10 cancers

- Classify tumors into immune subtypes

- Correlate with clinical outcomes



Image: Pan-cancer Immune Landscape Infographic. Derived from Petralia, F., et al., Cell 187, 2024

**Key findings:**

- 7 distinct immune subtypes
  - Common immune reactions, evasion mechanisms independent of cancer type
- Correlations between PFS and immune subtypes, TME immune cell load
- Specific kinases activated in subtypes
  - Immune evasion, pathogenesis, and host immunity

**Impact:**

- Multi-dimensional view of tumor biology
- Novel patient stratification, therapeutic strategies
- New interactive web portals: PhosNet Vis, ProKap
  - Leverage PDC's CPTAC pan cancer kinase and transcription factor activity score data to explore relationships with immune subtypes
  - New avenues for research and target discovery

## CCDI- Exemplar for building a learning health care system for cancer

- Community of researchers, advocates, hospitals and networks committed to sharing pediatric cancer data to accelerate research on childhood cancers.

- Federated Pediatric Cancer Data Ecosystem
  - Childhood cancer data and resources from across the nation
    - Research repositories
    - Patient registries
    - Hospitals



**PIECE IT TOGETHER**
The Childhood Cancer Data Initiative is completing the puzzle to learn from and help heal children, teens, and young adults with cancer.

**BUILD A STRONG BASE**
Progress requires data from many sources that is connected and easy to access.

**MAKE DATA EASY TO USE**
More thoughtful tools for analyzing data will help answer important questions.

**ASSEMBLE BETTER DATA**
Complete data sets are needed to understand each type of cancer.

**IMPROVE TREATMENTS**
Data is the foundation that informs new treatments and improves lives faster.

Image: Childhood Cancer Data Initiative. National Cancer Institute.
https://www.cancer.gov/research/areas/childhood/childhood-cancer-data-initiative/about

## New CCDI Hub

- Entry point for researchers looking to use and connect with CCDI-related data and resources



Image: Childhood Cancer Data Initiative Hub. National Cancer Institute. https://ccdi.cancer.gov/

- Cancer Data Service (CDS) and Imaging Data Commons (IDC)
  - CCDI data from 1400+ participants
  - 70,000+ files
    - Genetic testing data
    - WGS and WXS
    - Full transcriptome sequencing
    - Single cell analysis
    - Imaging data

- Hosts Data for CCDI project
  - CCDI's Molecular Characterization Initiative (MCI)
    - Detailed clinical and molecular information, patients treated at academic and medical institutions around the country

## CRDC Infrastructure supports new CCDI Data Hub

- Bento framework
  - GUI for Data Exploration
- Authentication & Authorization controls
  - Integration of dbGaP A&A workflows for controlled data access

# CRDC Infrastructure supports new CCDI Data Hub

- Access to SB-CGC
  - Cloud-based environment
  - 500+ bioinformatics tools, workflows
  - Combine with own datasets or data across CRDC, NCI Cloud Platform Interoperability (NCPI)
  - Collaborative workspaces



**NCPI**



Image: NIH Cloud Platform Interoperability Effort. NIH Office of Data Science Strategy. https://datascience.nih.gov/nih-cloud-platform-interoperability-effort."

AACR American Association for Cancer Research®

ANNUAL MEETING
2024 • SAN DIEGO

APRIL 5-10
#AACR24
AACR.ORG/AACR24

# CRDC Lessons Learned and Future State

Anthony Kerlavage, Ph.D.

Director

Center for Biomedical Informatics and Information Technology

National Cancer Institute, NIH, US

NIH NATIONAL CANCER INSTITUTE

# Disclosure Information

Anthony Kerlavage, PhD

I am a full-time paid employee of the NIH/NCI.
I have no financial relationships to disclose.

# The Year of Open Science

AACR
American Association
for Cancer Research®

ANNUAL MEETING
2024 • SAN DIEGO
APRIL 5-10 • AACR.ORG/AACR24 • #AACR24

The White House Office of Science & Technology Policy

**YEAR OF OPEN SCIENCE RECOGNITION CHALLENGE**

- Advancing national open science policy
- Providing access to the results of the nation's taxpayer-supported research
- Accelerating discovery and innovation
- Promoting public trust
- Driving more equitable outcomes

NIH's Data Management and Sharing Policy went into effect on January 25, 2023, fulfilling the memorandum's provisions around public access to scientific data.

# CRDC: Celebrating 10 Years of Data Sharing

# CRDC: Celebrating 10 Years of Data Sharing

## Data Infrastructure & Analysis

**DCC, DSS & Sustainability Awards**
Feb

**Cloud Resources Awards**
Sep

**Cloud Resources Live**
Jan

**DCF Award**
Sep

**DCF Live**
Mar

**CDA Award**
May

**CDA Live**
April

2014  2015  2016  2017  2018  2019  2020  2021  2022  2023

**DATA INFRASTRUCTURE AND ANALYSIS**

**Cloud Resources:** NCI Cloud Resources
**DCF:** Data Commons Framework
**CDA:** Cancer Data Aggregator

**DCC:** CRDC Data Hub
**DSS:** CRDC Data Standards Services
**Sustainability:** CRDC Sustainability Study

# Commemorating 10 Years of the CRDC

A four-part invited series in AACR *Cancer Research* journal showcasing how the CRDC empowers the cancer research community.

## 1 LESSONS LEARNED AND FUTURE STATE

Traces the history of the CRDC over the past 10 years, noting its progress in providing access to data and tools along with training and outreach to support the cancer research community. This review also provides an assessment of the CRDC's impact, lessons learned, and future plans to promote data sharing, data accessibility, interoperability, and reuse.

**Read Part One** 📖

## 2 RESOURCES TO SHARE KEY CANCER DATA

Describes each of the CRDC's data commons, including their unique and shared features, accomplishments, and challenges. This paper also details how the CRDC data commons implement Findable, Accessible, Interoperable, Reusable (FAIR) principles and promote data sharing in support of the NIH Data Management and Sharing Policy.

**Read Part Two** 📖

## 3 CLOUD-BASED ANALYTICAL RESOURCES

Details how the three Cloud Resources (CRs), including the Broad Institute FireCloud, Institute for Systems Biology Cancer Gateway in the Cloud (ISB-CGC), and Seven Bridges' Cancer Genomics Cloud powered by Velsera (SB-CGC) provide access to large, cloud-hosted multi-modal cancer datasets, as well as offer tools and workspaces for performing data analysis where the data resides. Included is a review of publicly available analytical tools.

**Read Part Three** 📖

## 4 CORE STANDARDS AND SERVICES

Outlines core CRDC services to aggregate descriptive information from multiple studies for findability via a single interface. These standards and services aggregate and semantically harmonize multiple data types making the CRDC a single point of discovery and access for cancer research data originating from multiple sources. They also facilitate the evolution of the CRDC as one hub for managing, storing, and sharing diverse types of data.

**Read Part Four** 📖

https://datacommons.cancer.gov/publications/aacr-cancer-research

# CRDC Impact to Date

**354** STUDIES

**134K** SUBJECTS

**9.4PB+** DATA

**2.4K+** YEARS OF COMPUTE

**2K+** PUBLIC TOOLS AND WORKFLOWS

**82.3K+** UNIQUE USERS / YEAR

**30K+** CRDC DATA CITATIONS

# Lessons Learned: Community Engagement

- Increase in training and educational resources
  - Community engagement
    - Challenges, seminars, conferences, training sessions
  - Training videos, tools, and documentation
    - cloud use, cloud cost prediction tools, multi-modal analysis
  - CRDC Insights: Quarterly Newsletter

# Lessons Learned: Sustainability Study

- Planning for long-term sustainability of CRDC resources

  - CRDC supports the democratization of cancer research by providing **cloud-based, secure storage and analytic tools** for cancer data.

  - The Sustainability Study **supports the CRDC in planning for the financial sustainability** of its future work, **ensuring it operates efficiently and delivers value for money to the cancer research community** for years to come.

**RECOMMENDATIONS FOR COST SAVINGS & OPERATIONAL OPTIMIZATION**

**BASELINE FOR CURRENT FINANCIAL REQUIREMENTS TO OPERATE CRDC; DEVELOP FUTURE FINANCIAL PROJECTIONS UNDER VARIOUS SCENARIOS**

**IDENTIFY BEST PRACTICES FROM EXTERNAL ENVIRONMENTS**

Poster #3568: "Managing large-scale cancer research data programs" on April 8th, 1:30 – 5:00 pm

# Lessons Learned: Lowering Barriers

# ARPA-H Biomedical Data Fabric* Toolbox

NCI in partnership with ARPA-H will advance the next-generation of tools to synthesize and speed use of health research data, starting with cancer

**Make biomedical research data easier to use**

**Reduce effort for data integration**

**Develop new data fabric capabilities & tools**

**Build health data science models that can be applied across disciplines**

\* A data fabric provides a unified, consistent layer of data services that can work across many different systems and environments.

# BDF Toolbox - Technical Areas

**TA1: Automated Data Collection**

Lower barriers to high-fidelity, timely, and automated data collection of research data across labs and health record systems



**TA2: Machine-Assisted Curation**

Prepare, connect, and harmonize multi-source data for analysis at scale



**TA3: Intuitive Exploration**

Enable advanced, human-centered data exploration and dashboards for use by diverse stakeholders and decision-makers



**TA4: User Engagement**

Evaluate data fabric tools across researchers, clinicians, and patients to create tools that will be enthusiastically adopted.



**TA5: Cross-Domain Generalization**

Leverage tools and platforms to generalize data across biomedical domains and disease types.

# CRDC Roadmap

## Data Commons

- ✓ CDS 1.0 Website
- ✓ GDC 2.0 Website
- ✓ CTDC 1.0 Website

➢ Common clinical data model & data dictionaries available
➢ Enclave 1.0 API
➢ Population Science DC Pilot

➢ PDC & ICDC 2.0 Website
➢ Image automated de-ID & conversion
➢ Support U, P grants

➢ Transition BDF tools
➢ PDC Metabolomics 1.0
➢ Support R Grants

➢ Clinical/EHR integration
➢ Cancer Models

**2024**     **2025**     **2026**

## Infrastructure & Analysis

- ✓ DCF RAS A&A
- ✓ 1st AIDR Challenge
- ➢ Automated file compression/archiving

➢ Data Submission & Discovery Dashboard MVP
➢ CDA full integration
➢ Helpdesk & Concierge 1.0
➢ Trans-NIH use cases

➢ Data Submission & Discovery Dashboard 1.0
➢ Single NCI workflow engine
➢ Evaluate BDF Tools

➢ Trans-HHS use cases
➢ Federated workflow engines
➢ Molecular and clinical data integration program

➢ Integrate BDF tools
➢ CRDC as a ML testbed
➢ Consolidated CRDC architecture

**2024**     **2025**     **2026**

# Learning Health System for Cancer

## NCI

Durga  Addepalli
Jill Barnholtz-Sloan
Erin Beck
Emi Casas-Silva
Zhaoyi  Chen
Heather Creasy
Tanja Davidsen
Ina Felau
Emily Greenspan
Jaime Guidry Auvil
Toby Hecht
Tony Kerlavage
Erika Kim
Henry Rodriquez
Pothur Srinivas
Louis Staudt
David Sturgill
Granger Sutton
Zhining Wang
Xu Zhang

## FFRDC/FNL

Amanda Bell
Paula Darte
Hayley Dingerdissen
Sharon Gaheen
Naila Gulzar
Mark Jensen
Gina Kuffel
John Otridge
Sam Pathak
Todd Pihl
Sudha Venkatachari
Ulli Wagner
Mike Warfe

**All CRDC contractors**

**All partners throughout NCI/NIH and data contributors.**

## Essex

Shanthala Basavappa
Melissa Cook
Kim Gambini
Amal Ghannam
Shannon Lane
Luis Santana

## ARPA-H

Andrea Bild
Julie Bletz
Jennifer Roberts
Alastair Thomson

CANCER RESEARCH

The Foundational Cancer Journal Driving Transformative Science

AACRJournals.org
@CR_AACR

AACR American Association for Cancer Research

Questions & Discussion

datacommons.cancer.gov

# AACR Cancer Research Series

A four-part invited series published online in March 2024 highlighting the CRDC's accomplishments from the past 10 years.

▶ LESSONS LEARNED AND FUTURE STATE

▶ RESOURCES TO SHARE KEY CANCER DATA

▶ CLOUD-BASED ANALYTICAL RESOURCES

▶ CORE STANDARDS AND SERVICES

Learn more about the series on the CRDC Website.

## CANCER RESEARCH

**The Foundational Cancer Journal Driving Transformative Science**

AACRJournals.org
@CR_AACR

AACR American Association for Cancer Research®

NIH ▶ NATIONAL CANCER INSTITUTE

# 2024 AACR Annual Meeting: San Diego, CA

## Presentations

- **Impact of the Cancer Research Data Commons (CRDC)**
  - Sunday, April 7 - 1:00pm – 2:00pm

- **NCI Artificial Intelligence (AI) Programs and Resources for Advancing Cancer Research**
  - Wednesday, April 10 - 10:15am -11:15am

## Posters

- **ISB – Cancer Gateway in the Cloud**
  - Monday, April 8

- **CRDC – Sustainability Implementation Planning**
  - Monday, April 8
- **Velsera – Seven Bridges, Cancer Genomics Cloud**
  - Monday, April 8
  - Tuesday, April 9
  - Wednesday, April 10
- **Broad – FireCloud (Terra)**
  - Wednesday, April 10

View the **AACR Program** for more details.

NIH ❯ NATIONAL CANCER INSTITUTE

# 2024 AACR Annual Meeting: San Diego, CA

### Posters

**Sunday, April 7, 2024  / 1pm - 5pm**
Session PO.RSP01.01 – Regulatory Science and Policy 920/3 – Insights from the NCI request for information on existing data sharing processes for NIH-funded research

**Monday, April 8, 2024  / 9 am – 12:30 pm**
Session PO.SHP01.01 – Science and Health Policy 1303/17 – Harness the power of data to improve cancer care – understanding the complex landscape of health policies and regulations

**Tuesday, April 9, 2024  / 2:30 – 3:30 pm**
Session NIH10 – Building on the Power of Data and Community – CCDI data ecosystem: Tools and resources

**Wednesday, April 10, 2024  / 10:15 – 11:15 am**
Session NIH12 - NCI Artificial Intelligence (AI) Programs and Resources for Advancing Cancer Research – NCI funding opportunities, resources and activities

**NIH** NATIONAL CANCER INSTITUTE

# 2024 CRDC Fall Symposium: October 16-17, 2024

A one-and-a-half day event highlighting the 10th anniversary of the CRDC as well as plans for the future.

📍 NIH MASUR AUDITORIUM, BETHESDA MD (10/16)
NCI CAMPUS, ROCKVILLE MD (10/17)

**PRE-REGISTRATION REQUIRED**

REGISTER & MORE INFORMATION AT
DATACOMMONS.CANCER.GOV

**CRDC and ODS Collaboration Session**
Wednesday, October 16 @ 1:30 PM ET

- Data Sharing & Access within CRDC
- CRDC Symposium Kick-Off

*Immediately following NCI Office of Data Sharing Symposium (separate event registration)*

**CRDC Session**
Thursday, October 17 @ 9:00 AM ET

- CRDC History & Current State
- Success Stories & Impactful Programs
- Future Spotlight
- Fireside Chat

NIH NATIONAL CANCER INSTITUTE