

## Frederick National Laboratory for Cancer Research



### JDACS4C – Joint Design of Advanced Computing Solutions for Cancer

Frederick National Laboratory Advisory Committee Meeting

Eric Stahlberg, PhD

May 8, 2017

The Frederick National Laboratory is a Federally Funded Research and Development Center operated by Leidos Biomedical Research, Inc., for the National Cancer Institute  
DEPARTMENT OF HEALTH AND HUMAN SERVICES • National Institutes of Health • National Cancer Institute

## JDACS4C – NCI-DOE Collaboration Overview

Frederick  
National  
Laboratory  
for Cancer Research

- **Shared Interests**

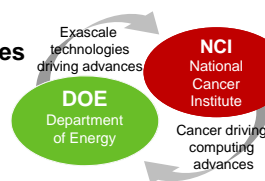
- Cancer scientific challenges driving advances in computing
- Exascale technologies driving cancer advances

- **Supports Two Primary Executive Office Initiatives**

- Precision Medicine Initiative (Jan 2015)
- National Strategic Computing Initiative (July 2015)

- **Three Pilot Efforts**

- Molecular domain pilot
  - NCI Lead: **Frank McCormick** (FNL/UCSF), **Dwight Nissley** (FNL)
  - Lead DOE Leads: **Fred Streitz** (LLNL), **Felice Lightstone** (LLNL)
- Pre-clinical domain pilot
  - NCI Lead: **Jim Doroshov** (DCTD) and **Yvonne Evrard** (FNL)
  - Lead DOE Leads: **Rick Stevens** (ANL)
- Population/clinical domain pilot
  - NCI Lead: **Lynne Penberthy** and **Paul Fearn** (DCCPS)
  - Lead DOE Labs: **Gina Tourassi** (ORNL), **Gil Weigand** (ORNL)



## JDACS4C Collaboration Pilots: Capabilities to Accelerate Precision Oncology

**NCI Mission Impact:**  
*Accelerating development of new treatment options for precision cohorts*

**Pilot 2:**  
Biological Models  
*Multi-scale computational biological models*

**Pilot 1:**  
Pre-clinical Models  
*Predictive patient drug response models with advanced computing*


**Pilot 3:**  
Cancer Surveillance  
*Computational insight into factors impacting clinical response*

## Accelerating Computational and Data-driven Cancer Research

**Exascale in a nutshell:**

- Millions of CPU cores contributing to a single task
- Nearly 1000 times faster than fastest computer today
- Focus of DOE Advanced Strategic Computing

## Joint Design of Advanced Computing Solutions for Cancer



**JDACS4C**


Exascale technologies driving advances


DOE  
Department of Energy


NCI  
National Cancer Institute


Cancer driving computing advances


**Initiatives Supported NSCI and PMI**

 NATIONAL CANCER INSTITUTE

 ARGONNE




 OAK RIDGE

 LAWRENCE LIVERMORE NATIONAL LABORATORY


 LOS ALAMOS


Frederick National Lab for Cancer Research

### Integrated Precision Oncology

	Molecular	Pre-clinical	Population
	<b>Pre-clinical Domain – Improved predictive models</b> <i>Computational/hybrid predictive models of drug response</i> <i>Improved experimental design</i>		
	<b>Clinical Domain – Precision oncology surveillance</b> <i>Expanded SEER database information capture</i> <i>Modeling patient health trajectories</i>		
	<b>Molecular Domain – Multiscale biological models</b> <i>Models for RAS-RAS complex interactions</i> <i>Insight into RAS related cancers</i>		
<b>CANcer Distributed Learning Environment (CANDLE)</b> Scalable Deep Learning for Cancer			

JDACS4C established June 27, 2016 with signed MOU between NCI and DOE

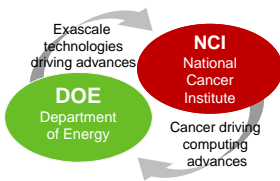

5

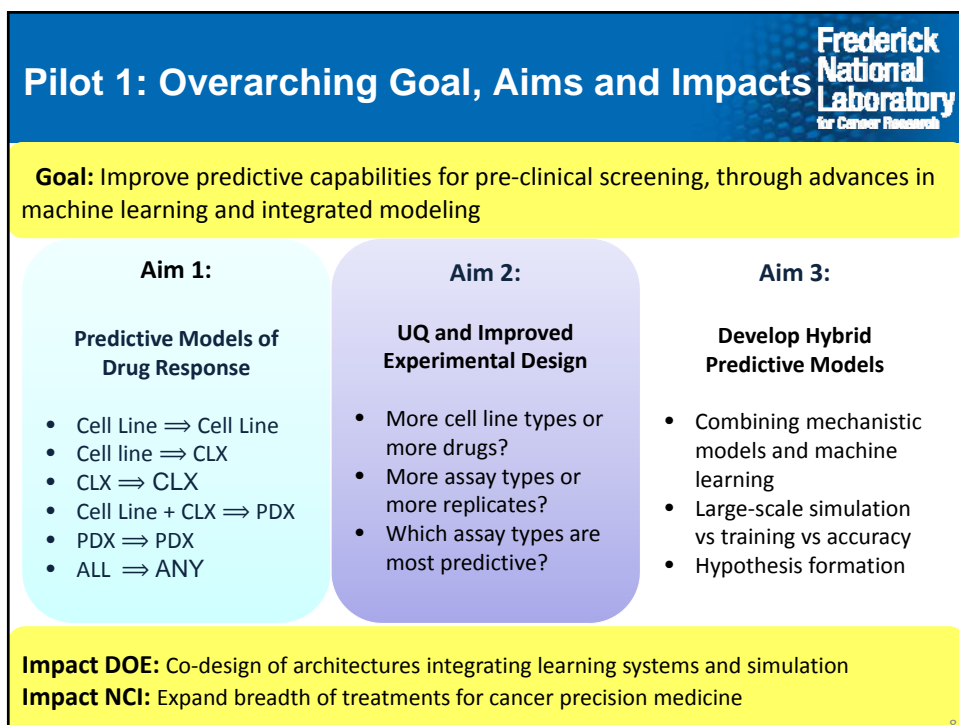
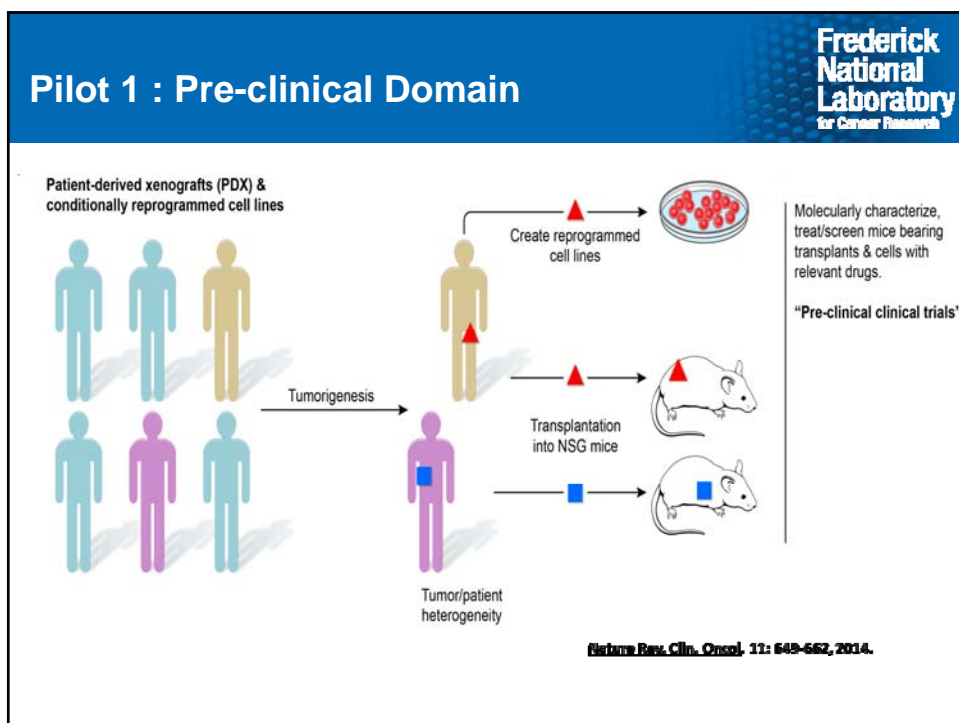


## JDACS4C Pilot 1 Highlights

---

Predictive Modeling for Pre-Clinical Screening





Frederick  
National  
Laboratory  
for Cancer Research

## Pilot 1: Cross Laboratory Team

**ANL:** Rick Stevens, Frank Alexander, Jillian Aurisano, Prasanna Balaprakash, Tom Brettin, Jim Davis, Emily Dietrich, Nicoli Dryden, Hal Finkel, Ian Foster, Monisha Ghosh, Ushma Kriplani, Ravi Madduri, Sergei Maslov, Bob Olson, Dan Olson, Mike Papka, Lorenzo Pesce, John Santerre, Maulik Shukla, Venkat Vishwanath, Fangfang Xia

**LANL:** Marian Anghel, Tanmoy Bhattacharya, Judith Cohn, Paul Dotson, Will Fischer, Kumkum Ganguly, Jason Gans, Cristina Garcia-Cardona, Nick Hengartner, William Hlavacek, John Hogden, Patrick Kelly, Miranda Lynch, Ben McMahon

**LLNL:** Jonathan Allen, Ya Ju Fan, Adam Zemla

**ORNL:** Mike Lueze, Arvind Ramanathan

**NCI:** James Doroshov, Yvonne Evrard, Susan Holbeck, Eric Stalhberg, George Zaki



P1 Hackathon 2



CANDLE Hackathon 1



P1 Hackathon 3

Frederick  
National  
Laboratory  
for Cancer Research

## Pilot 1: Recent Key Results

- **Data environment in place**<sub>ANL,NCI,LANL</sub>
- **PDX model production, NCI-60 RNAseq**<sub>NCI</sub>
- **Aim1 : Predictive Models of Drug Response**<sub>ANL, UC, LLNL,LANL</sub>
  - Shallow models (feature selection and drug/tumor specific signatures)
  - Deep models to utilize larger training datasets (single drug, drug pairs)
  - Demonstration of convolution success\* in P1B3 (CANDLE benchmark)
  - Initial transfer learning models and autoencoders (expression, drugs)
- **Aim2 : UQ and Improved Experimental Design**<sub>LANL,ANL</sub>
  - QA, clustering and mapping studies of NCI60, CCLE, GDSC, GDC
  - Initial tradeoff studies (sample size vs error)
  - Model transfer problem formulation and initial development
  - Feature analysis {transcriptome > proteome > kinome}, {exp >> SNPs}, {RNAseq ≥ microarray}
- **Aim3 : Develop Hybrid Predictive Models**<sub>ANL, UIUC</sub>
  - Initial molecular network based convolution experiments\*
- **PILOT1 Benchmarks for CANDLE (three model problems)**
  - NCI workshop in April 18-19<sup>th</sup>
  - NVIDIA workshop May 9<sup>th</sup>

## Pilot 1: Data Sources and Integration

- **NCI-60 Cell Lines**
  - 61 cell lines, 92,691 compounds, 40 assay types
- **NCI Sarcoma Project**
  - 74 cell lines, 445 compounds
- **NCI Small Cell Lung Cancer Project (SCLC)**
  - 76 cell lines, 525 compounds
- **NCI Patient-derived Models and Xenografts (PDM)**
  - 274 PDM/PDX Models
- **Broad-Novartis Cancer Cell Line Encyclopedia (C)**
  - 504 cell lines, 24 compounds
- **Genomics of Drug Sensitivity in Cancer (GDSC)**
  - 1,074 cell lines, 265 compounds
- **Genomic Data Commons (GDC)**
  - 14,531 samples, 29 primary sites, 38 cancer types

Cancer Data Processing, Storage and Machine Learning Workflow

## Pilot 1: Early Insights

- Effort is not limited by model representations – many model types have predictive capability on response problems
- Efforts are constrained by available data that capture the distributions expected in clinically relevant cases
- Convolutional neural networks have shown promise on expression and drug descriptors
- Shallow models can learn responses for single drugs/tumor types and be used for feature selection
- Transfer learning can be used to learn features in one problem and reapplied in another problem

**Frederick National Laboratory**  
for Cancer Research

## JDACS4C Pilot 3 Highlights

Population Information Integration, Analysis, and Modeling

Exascale technologies driving advances  
 DOE Department of Energy  
 NCI National Cancer Institute  
 Cancer driving computing advances

**Frederick National Laboratory**  
for Cancer Research

## Pilot 3: Population Domain

Surveillance data captured on each cancer patient for the entire population

Demographics → Pathology → Molecular Characterization → Initial Treatment → Subsequent Treatment → Progression Recurrence → Survival Cause of Death

Understand treatment and improve outcomes in the "real world"  
 SEER Cancer Information Resource  
 Exposome Genome  
 Prospectively support development of new diagnostics and treatments

## Pilot 3: Community Engagement

Frederick  
National  
Laboratory  
for Cancer Research

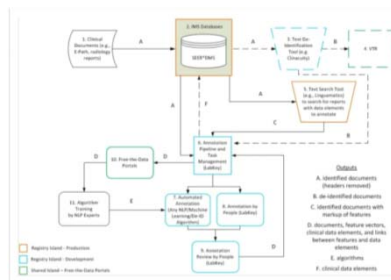
- NCI Team
- DOE Labs- Oak Ridge, Los Alamos, Lawrence Livermore, Argonne
- 4 SEER registries- Kentucky, Louisiana, Georgia, Seattle
- Contracts/Sub-Contracts
  - IMS as honest broker and agent for registries, hosting and sub-contracts
  - Software and support from LabKey, Linguamatics for Clinical Document Annotation & Processing (CDAP) Pipeline
  - Annotation Services
- Data acquisition and linkages (e.g. claims, pharmacy, radiation oncology)
- Clinical experts to lead use cases for breast, colorectal, lung, and prostate cancer
- NLP Workshop with NCI, DOE, CDC, FDA and academic partners

## Pilot 3: AIM 1 - Progress Update

Frederick  
National  
Laboratory  
for Cancer Research

### Annotation Framework:

- Utilizing clinical document annotation pipeline to annotate ALK, EGFR biomarkers in e-path reports to send to DOE for algorithm development
- Developing schema for annotation of recurrence, progression data elements in path reports for breast & colorectal cancer
- Plans to scale up pipeline to annotate up to 10,000 documents per month
  - Add biomarkers from breast, colorectal, lung, prostate and pathology elements from CAP protocols





Frederick National Laboratory for Cancer Research

## Pilot 3: AIM 1 - Progress Update

- **Text Comprehension:**
  - Developed and benchmarked rule-based, conventional ML-based, and DL-based NLP tools for e-paths for three information extraction tasks: (i) primary cancer site, (ii) histological grade, (iii) behavior
    - DL tools: CNN, HAN, MT-DNN
    - DL interpretability
    - Cross-registry robustness validation
  - Established reproducible experimental design pipeline for future studies with new data and different tasks
  
- **Text Synthesis:**
  - Generative models for e-path text synthesis under development
  - Character and word-based LSTM-RNNs e-path text synthesis under development
  - Initial validation in progress

17

Frederick National Laboratory for Cancer Research

## CANDLE Highlights

---

CANcer Distributed Learning Environment

**Frederick National Laboratory**  
for Cancer Research

## CANDLE: Deep Learning Across JDACS4C

### ECP-CANDLE Project : CANCER Deep Learning Environment

**CANDLE Goals**

- Develop an exscale deep learning environment for cancer
- Building on open source Deep learning frameworks
- Optimization for CORAL and exascale platforms
- Support all three pilot project needs for deep learning
- Collaborate with DOE computing centers, HPC vendors and ECP co-design and software technology projects

**Frederick National Laboratory**  
for Cancer Research

## CANDLE: CANCER Distributed Learning Environment

- **CANDLE is DOE supported contribution to JDACS4C**
- **Four year multi-lab project commencing in September 2016**
- **Focuses on creating scalable, open and portable Deep Learning framework**
- **Supports Deep Learning needs for all JDACS4C pilots**
  - DOE scientific leads bring pilot-specific deep learning challenges
- **Open Source software release**
- **Reference benchmarks released in January 2017**
- **Workshops and early community building**
  - CANDLE @ NIH – April 18-19, 2017 – More than 60 attendees involving more than 12 NIH institutes/centers
  - CANDLE @ GTC – May 9, 2017

## JDACS4C Summary

- **Building, motivating and expanding interdisciplinary collaborations and partnerships**
- **Piloting cross-disciplinary activities to enable integrated precision oncology**
  - Identifying key linkages among and gaps between domains
- **Delivering scientific insight together with new capabilities**
  - Scalable frameworks, environments, conventions and standards
  - Cutting-edge predictive models for cancer across multiple domains
  - Advanced machine learning accounting for uncertainty in data
- **Enabling open and team science to more rapidly achieve goals in precision oncology**

## Acknowledgements

- |  |   |
|--|---|
| <ul style="list-style-type: none"> <li>• <b>NCI CBIIT</b> <ul style="list-style-type: none"> <li>– Warren Kibbe, Carl McCabe, Betsy Hsu and the CBIIT team</li> </ul> </li> <li>• <b>NCI DCCPS</b> <ul style="list-style-type: none"> <li>– Lynne Pemberthy, Paul Fearn, Jessica Boten and DCCPS team</li> <li>– Louisiana, Seattle, Georgia and Kentucky cancer registries</li> </ul> </li> <li>• <b>NCI DCTD</b> <ul style="list-style-type: none"> <li>– Jim Doroshov, Susan Holbeck</li> </ul> </li> <li>• <b>FNLCR</b> <ul style="list-style-type: none"> <li>– Data Science and IT Program</li> <li>– RAS Initiative - Dwight Nissley, Frank McCormick and RAS team</li> <li>– PDX - Yvonne Evrard and PDX team</li> </ul> </li> </ul> | <ul style="list-style-type: none"> <li>• <b>Department of Energy</b> <ul style="list-style-type: none"> <li>– Dimitri Kusnesov, Steve Binkley, Doug Wade, Carolyn Lauzon</li> </ul> </li> <li>• <b>Argonne National Lab</b> <ul style="list-style-type: none"> <li>– Rick Stevens, Tom Brettin, Fangfang Xia, Emily Dietrich, and the ANL team</li> </ul> </li> <li>• <b>Lawrence Livermore</b> <ul style="list-style-type: none"> <li>– Jason Paragas, Amy Gryshuk, Fred Streit, Felice Lightstone, Brian Van Essen, Dave Rakestraw and the LLNL team</li> </ul> </li> <li>• <b>Los Alamos National Lab</b> <ul style="list-style-type: none"> <li>– Marian Anghel and LANL team</li> </ul> </li> <li>• <b>Oak Ridge National Lab</b> <ul style="list-style-type: none"> <li>– Gina Tourassi, Gil Weigand, Arvind Ramanathan, Joe Lake, and the ORNL team</li> </ul> </li> </ul> |
|--|---|

*And many, many more supporting and working with JDACS4C !*