### Frederick National Laboratory for Cancer Research



### **Big Data Big Theory**

Eric Stahlberg, George Zaki

New York Scientific Data Summit 2018 - Brookhaven National Laboratory

#### August 6, 2018

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



### **FNLCR and Data Science**

## **Overview of Frederick National** Laboratory for Cancer Research (FNLCR)

- FNLCR is the only Federally Funded Research and Development Center (FFRDC) dedicated exclusively to biomedical research
  - Operated in the public interest by Leidos Biomedical Research, Inc (formerly SAIC-Frederick) on behalf of the National Cancer Institute
- Main campus located on 70 acres at Ft. Detrick, MD
  - Leidos Biomed employees co-located with NCI researchers and other contractors on the NCI Campus at Frederick
  - Additional Leidos Biomed scientists at Bethesda and Rockville sites





### Mission

Provide a unique national resource for the development of new technologies and the translation of basic science discoveries into novel agents for the prevention, diagnosis and treatment of cancer and AIDS.

**Frederick** 

Laboratory for Cancer Research

4

## **Research & Development at FNLCR**

### Research & Development

- **Basic Research**: New knowledge about AIDS and cancer
- **Applied R&D**: New diagnostics and therapeutics
- **Clinical Research:** Clinical trials and laboratory analysis
- **cGMP manufacturing:** Biologicals and vaccine production

### Specialties

- Genomics, proteomics, and metabolomics
- Bioinformatics and imaging
- Nanotechnology
- Animal models
- Tumor cell biology and virology
- Immunology and inflammation
- Data Science and Information Technology underpin and support all R&D activities and specialties



Frederick

Laboratory for Cancer Research



### **Data Science and Biomedical Informatics**

Leverage leading edge data science and information technology skills, tools, and capabilities to accelerate translation of biomedical data to scientific discoveries, medical treatments, diagnostic and prevention tools for cancer and AIDS patients.



Enabling physicians, scientists, and patients to make critical decisions based on knowledge gained from <u>all</u> and not only a fraction of data and information available to them.



Frederick

Laboratory

for Cancer Research

Frederick National Laboratory for Cancer Research

### **Predictive Oncology**



# Predictive Oncology – Leveraging and Learning

**Frederick** 

Laboratory for Cancer Research





### **Predictive Oncology Learning System**



### Predictive Oncology General Adoption Challenges

- Potential criteria to achieve broad adoption of predictive oncology models
  - Credible scientifically substantiated
  - Reliable performs as expected for specified conditions
  - Affordable generally accessible in a costeffective implementation
  - In patient's best interest able to aid patient diagnosis or treatment in a timely manner





9



### NCI Precision Oncology Extending the Frontiers

Frederick National Laboratory for Cancer Research

- Identify promising new treatment options through the use of advanced computation to rapidly develop, test and validate predictive pre-clinical models for precision oncology.
- Deepen understanding of cancer biology and identify new drugs through the integrated development and use of new simulations, predictive models and next-generation experimental data.
- Transform cancer care by applying advanced computational capabilities to population-based cancer data to understand the impact of new diagnostics, treatments and patient factors in real world patients.

### **Challenge Areas for Predictive Oncology**

### Challenges for cancer

- Insufficient data for describing all possibilities
  - Over 250,000 unique cancer characterizations
  - Observation gaps absence of specific confirming data
  - Bridging molecular with preclinical and preclinical to clinical domains
- Data fusion and scientific credibility
  - Achieving coherence across scales
  - Achieving coherence and quality across organizations
- Achieving reliability
  - Consistency of response for characterized conditions
  - Accounting for uncertainty of unknown factors
  - Similarity of behavior across similar models







Frederick National Laboratory

### Joint Design of Advanced Computing Solutions for National **Cancer (JDACS4C) Developing Key Capabilities**



Frederick

Laboratory for Cancer Research

# Joint Design of Advanced Computing Solutions for Cancer





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



- Timely, accurate, adoptable and affordable precision oncology predictions to inform clinical decisions
  - Spanning domain of available information from molecular insights to population factors
- Requires integrating new capabilities across multiple observational domains
- Integrate JDACS4C pilots across the domain of scales
- How to develop predictive models that are efficient while maintaining experimental, scientific and computational credibility?



- Pathfinders to survey and find ways in each area
- Deployed seven engineering regiments concurrently with unique assignments – each about 350 miles to complete
- Common goals to connect with the other by a set time
- Within each regiment, two teams leapfrogged each other
- Technical impediments (like permafrost) required novel engineering, shared solutions across the effort



Frederick

Laboratory

for Cancer Research

# Multi-domain computing for integrated prescriptive cancer predictions

Frederick National Laboratory for Cancer Research



### **Complex System Modeling Drives Demand** for Computing





### JDACS4C Pilots Pioneering Computational Capabilities in New Areas

### Frederick National Laboratory



### **Complex System Modeling Drives Demand for Computing**





### **Complex System Modeling Drives Demand for Computing**

#### Ensembles of models –

- Promises for affordability and customizability
- Potential challenges in credibility and ease of explanation



Frederick

Laboratory

for Cancer Research

### From ensemble to predictive model



- Can ensemble data lead to a reduced representation using a predictive model? What is the cost?
  - Deep learning can be an avenue to get to that point without having to do the computational investment in the ensemble.
  - How to build a predictive model for that corpus of data?
    - Data size
    - Model complexity
    - Hyperparameters
    - Available compute
- CANDLE is a tool to scale the training of these DL models on next generation of Exascale systems.

#### **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

Collaborate with DOE computing centers, HPC vendors and ECP co-design and software technology projects

needs for deep dearning



Frederick

Laboratory

for Cancer Research

- Frederick National Laboratory for Cancer Research
- CANDLE is DOE Funded contribution to JDACS4C
- Four year project.
- 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
- FNL brings NCI connection to CANDLE
  - Translating computational environment to broader cancer research community
  - Portability and standardization of model representations
  - Conventions and methods for model validation and evaluation in cancer

### **CANDLE Software Stack**

Hyperparameter Sweeps, Data Management (e.g. DIGITS, Swift, etc.)

Workflow

Network description, Execution scripting API (e.g. Keras, Mocha)

Scripting

Tensor/Graph Execution Engine (e.g. Theano, TensorFlow, LBANN-LL, etc.)

Engíne

Architecture Specific Optimization Layer (e.g. cuDNN, MKL, etc.)

Optimization

Frederick National Laboratory for Cancer Research



### **CANDLE - Multi-level Parallelism**

# Parallelism Targets in CANDLE

10,000 x 10-1000 x 10-100 = 1M - 1000M "cores"



Data Parallel 10x-1000x Data Parallel 10x-1000x Model Model Model Model Model Parallel Parallel Parallel Parallel Parallel 10x-100x 10x-100x 10x-100x 10x-100x 10x-100x ... ...

### Summary

- Predictive oncology faces big data and big theory challenges to achieve wide adoption.
- Learning is an important aspect of improving models.
- DOE and NCI have embarked on pilot efforts to explore challenges and opportunities in different domains.
- Achieving coherence from the molecular domain to the clinical domain presents computational, engineering and scientific challenges.
- Deep learning and CANDLE point to potential paths forward.





# Thank you! Questions?

## George Zaki (george.zaki@nih.gov)