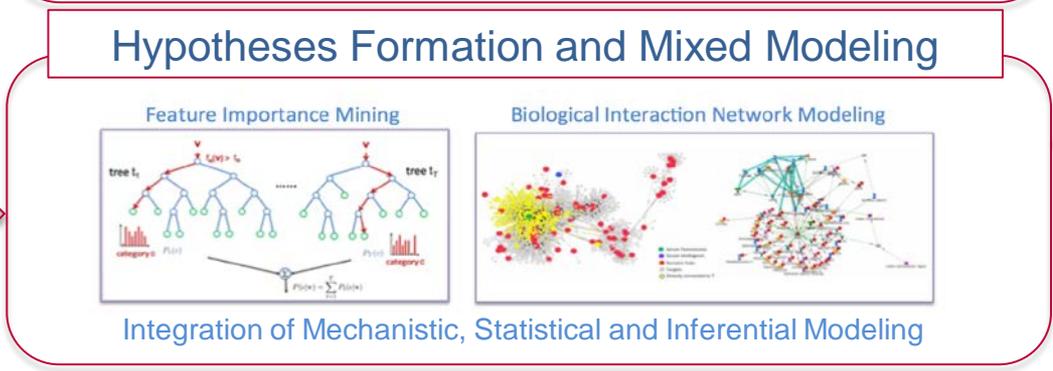
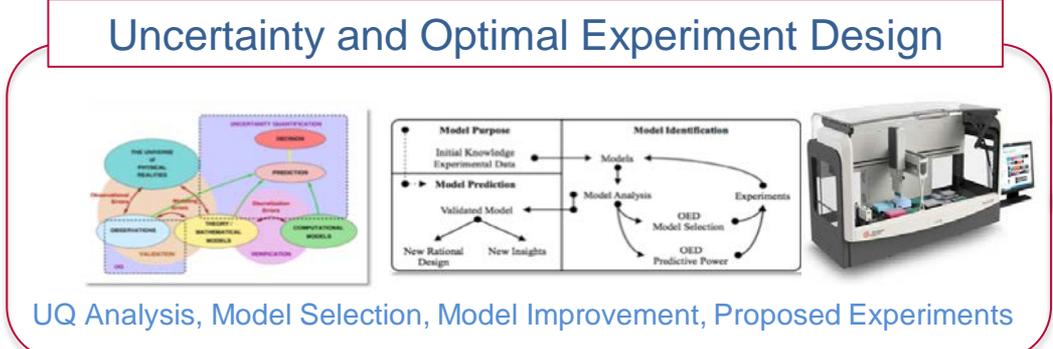
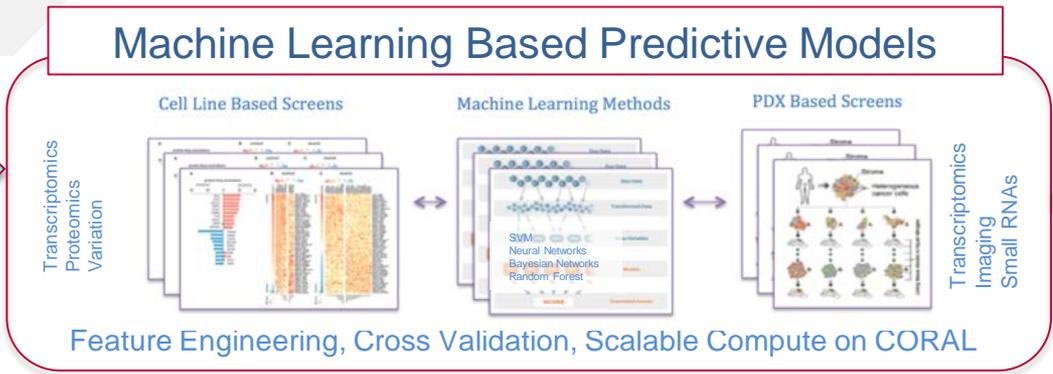
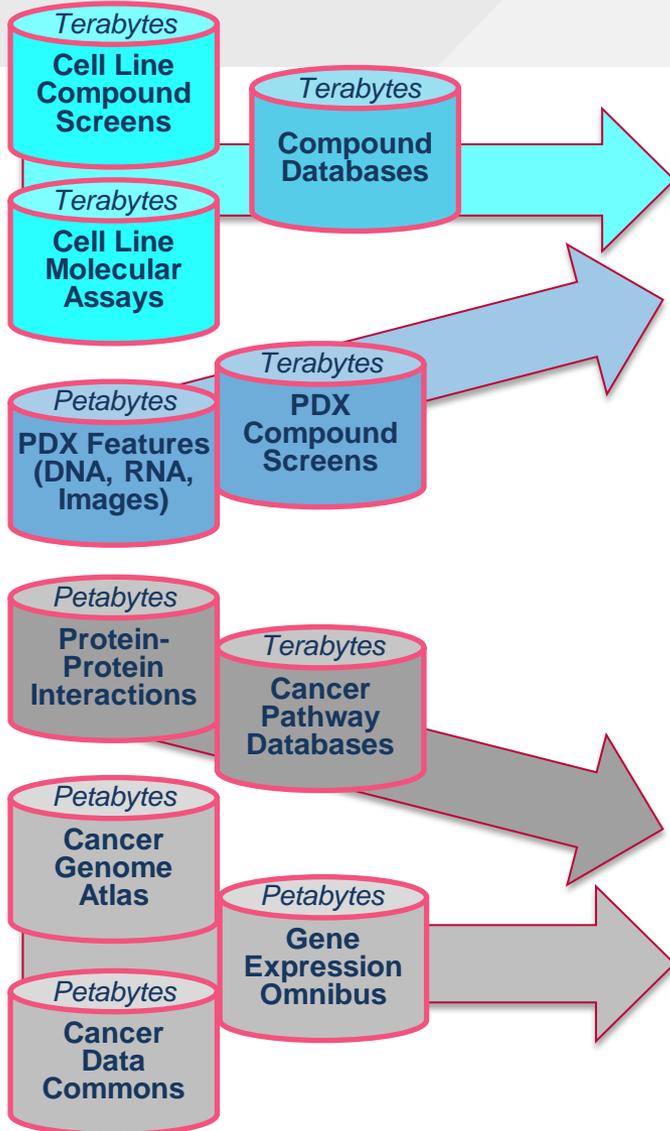


DOE and NIH Partnerships In Predictive Oncology

*James H. Doroshow, M.D.
Deputy Director for Clinical and Translational Research
National Cancer Institute, NIH*

*Pilot 1 Lead with:
Rick Stevens (Argonne National Laboratory, University of Chicago)
and Frank Alexander (Los Alamos National Laboratory)*

Pilot 1: Predictive Models for Preclinical Screening



Data Sharing for NCI-DOE Collaboration

In vitro Drug Response
and Molecular
Characterization Data

Developmental Therapeutics Program (DTP): NCI-60 Cancer Cell Line Panel Data Sets

- Use the publicly-available database of NCI-60 panel molecular characteristics. Includes 60 cell lines from 9 different tissue types
- ~350,000 small molecules and natural product extracts have been screened for growth inhibition/cell kill

Large scale Data Sets

- mRNA expression
 - multiple microarray platforms performed by multiple groups
- Exome sequence
- SNP array
- MicroRNA expression (2 datasets)
- Proteomics
 - Reverse phase protein arrays (2 datasets)
 - Mass spec
- Metabolomics
- DNA methylation and copy number

Smaller scale Data Sets

- Activity measurements
 - Drug efflux pumps
 - p53 function
- Viral infectivity
- Individual proteins or phospho-proteins
- Karyotypes
- mRNA by RT-PCR

Developmental Therapeutics Program (DTP): Sarcoma and SCLC Cell Line Data Sets

Sarcoma Data Sets

63 sarcoma lines screened with 100
FDA-approved and 345
investigational agents

- Exon Array Data (Affymetrix)
- MicroRNA expression
- IC50

<http://sarcoma.cancer.gov>

SCLC Data Sets

62 SCLC lines screened with 103
FDA-approved and 420
investigational agents

- Exon Array Data (Affymetrix)
- MicroRNA expression
- IC50

<http://sclccelllines.cancer.gov>

Sarcoma Cell Line Screen of Oncology Drugs and Investigational Agents Identifies Patterns Associated with Gene and microRNA Expression.

Teicher BA, Polley E, Kunkel M, Evans D, Silvers T, Delosh R, Laudeman J, Ogle C, Reinhart R, Selby M, Connelly J, Harris E, Monks A4, Morris J.

Mol Cancer Ther. 2015 Nov;14(11):2452-62.

Small Cell Lung Cancer Screen of Oncology Drugs, Investigational Agents, Gene and microRNA Expression.

Polley E, Kunkel M, Evans D, Silvers T, Delosh R, Laudeman J, Ogle C, Reinhart R, Selby M, Connelly J, Harris E, Fer N, Sonkin D, Kaur G, Monks A, Malik S, Morris J, Teicher BA.

JNCI 2016,

In Press

The NCI ALMANAC: Testing All Pairwise Combinations of Approved Cancer Drugs

- The NCI ALMANAC (A Large Matrix of AntiNeoplastic Agent Combinations)
- Currently just over 100 small molecule oncology drugs are FDA-approved.
- Test all possible pairwise combinations: ~5000 drug pairs
- Test each drug pair in each of the cell lines in the NCI-60 panel:
 - ~300,000 experiments
 - ~4.3 million wells
- Screen run at Frederick National Labs & 2 contract locations

Data Sharing for NCI-DOE Collaboration

In vivo Patient-Derived
Xenograft (PDX) Models
Drug Response and
Molecular Characterization
Data

NCI Patient-Derived Models (PDM) Repository at FNLCR

A national repository of clinically-annotated patient-derived xenografts (PDXs) and cell cultures (PDCs) to serve as a resource for academic discovery efforts and public-private partnerships for drug discovery.

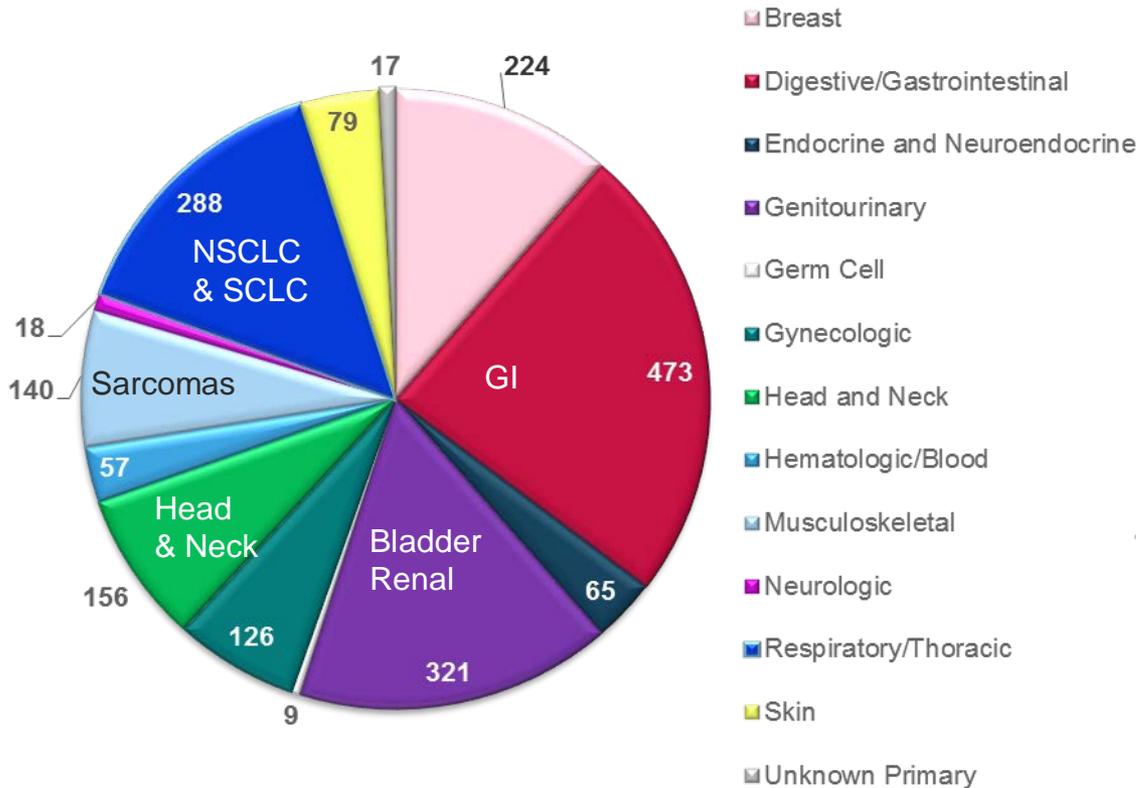
NCI plans to provide a long-term home for >1000 PDX and PDC models each produced from tissues and blood supplied by NCI-designated Cancer Centers, NCORP, and ETCTN participants.

Goals:

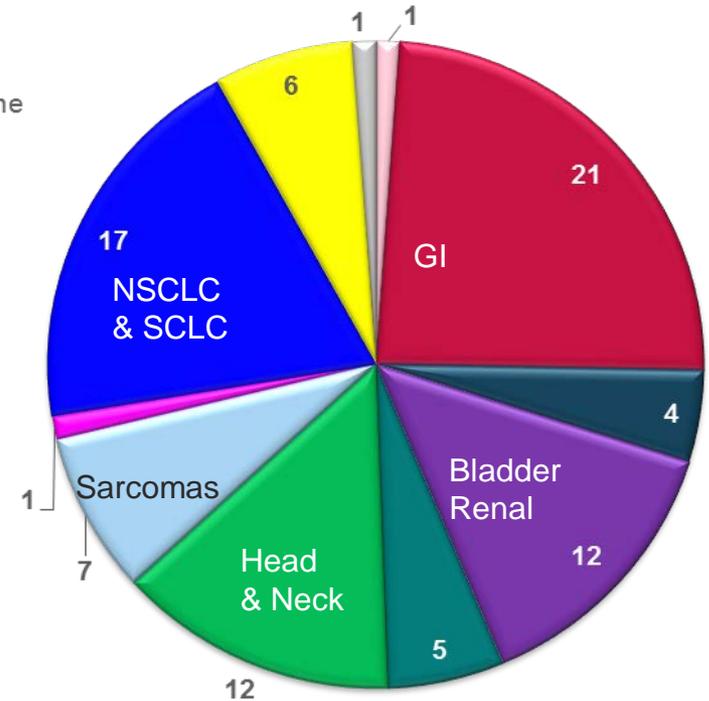
- ~50 unique patient models (solid & derived tumor line) per disease with sufficient size of each molecularly-characterized subgroup to power validation and/or efficacy studies
- Comprehensive molecular characterization of samples and early-passage PDXs: targeted mutation panel, WES, RNASeq, histology, growth curves, and proteomics/phospho-proteomics (pilot study), whole-mouse imaging (pilot study)
- All models and associated data made available through a publicly available website

NCI Patient-Derived Models Repository

Tumor Tissue Specimens Organized by Disease Body Location



Models with Confirmed Histology and Genetic Characterization Completed or In Progress



As of Apr 29, 2016

Total Number of Specimens (tumor resections and biopsies) Received: 711

As of Apr 29, 2016

Total: 87. These numbers increase monthly as new models complete "distribution lot" growth

Preclinical MPACT Drug Studies

- PDX models selected based on clinical MPACT actionable mutation of interest (aMOI) criteria. Models are tested in both the “clinically assigned” cohort as well as the other 3 cohorts. Selection is independent of histology.
 - Additional studies have been added to test single-agent cohorts to determine if models with responses were due to a single agent or the combination
- Target is 8-10 models assigned/cohort to be able to analyze data with sufficient statistical strength

	No aMOI Identified	Clinically-Assigned Cohort (all models tested in all cohorts)			
		Everolimus	Trametinib	ABT-888 + Temozolomide	MK1775 + Carboplatin
Overall Number of Models Sequenced for Cohort Assignment	19	2	15	2	24
Number of Models Selected for Preclinical MPACT (as of 4/29/2016)	3	2	12 (3 models from same patient)	2	8

Drug Studies	Preclinical MPACT (4 cohorts)	Single-agent: ABT-888 and Temozolomide	Single-agent: MK1775 and Carboplatin
Completed	17	16	6
In Progress	0	4	2
In Queue	10	10	9

NCI PDM Repository: Model Characterization Data Sets

- Patient diagnosis, subtype, limited medical history
- Histopathology images and standardized report: for patient tissue (when available) and every PDX
- Sequencing files for patient tissue (when available) and 4-6 PDXs/model
 - Whole exome sequence (VCF and FASTQ files; BAM files also available)
 - RNASeq (FASTQ and RSEM files)
 - NCI Cancer Gene Panel variant results (MPACT assay; variant status for 62 cancer-associated genes)
- Data from models developed under NCI tissue procurement protocols, as well as those acquired from Jackson laboratories and models currently being acquired from other academic and commercial sources will be included with this collaboration.

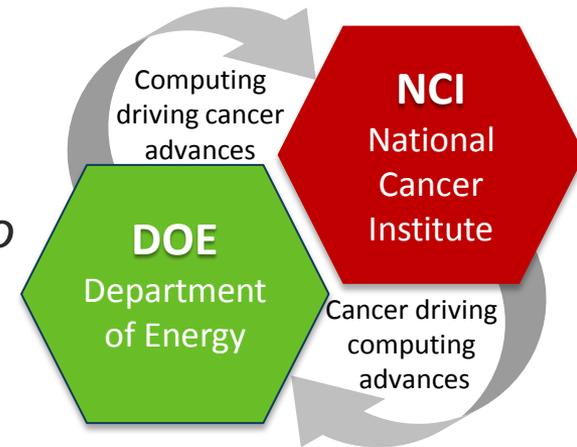
These data are stored in the NCI Patient-Derived Models Repository, hosted at FNLCR, and are available for download by the DOE for this project

NCI PDM Repository: PDX Drug-Response Data Sets

- Ongoing preclinical MPACT drug response studies
 - Currently have 17 models tested on all 4 MPACT trial arms; plan to add at least 15 more models based on clinical trial mutation profile selection criteria
- Administrative Supplement announcement released 5/2/2016:
 - Test 30 PDX models, including models from the PDM Repository, across a diverse range of histologies using NCI-IND agents, including FDA-approved agents for non-approved indications, to facilitate the development of improved early-phase drug development plans for these compounds.

Predictive Models for Preclinical Screening in Oncology

1. Use molecular characterization data from both *in vitro* cell lines and *in vivo* PDXs as well as drug response data in machine-learning to build computational models of observed drug response or resistance
2. Using drug response data from single agent and combination studies to propose a set of compact molecular signatures that are predictive of sensitivity or resistance to specific agents or classes of agents
 - Signatures will not be restricted to one data type (i.e., not only genetic variants); these signatures can incorporate several levels of molecular/phenotypic characteristics
3. Test molecular signatures and predicted effective agents on PDX models, or other *in vitro* models, and feed results back into machine-learning modeling
 - Prioritize molecular assays and sample characteristics to improve prediction of response
4. Explore "deep learning" methodologies applied to the DTP and PDX datasets



Acknowledgements

Melinda Hollingshead
Yvonne A. Evrard
Alice Chen
Michelle Ahalt-Gottholm
Michelle Eugeni
Sergio Alcoser
P. Mickey Williams
Anand Datta
Jason Lih
Bishu Das
Han Si
Dianne Newton
Carrie Bonomi
Kelly Dougherty
Cheryl Davis
John Carter

Susan Holbeck
Jerry Collins
Marie Hose
Dane Liston
Karen Schweikart
Penny Svetlik
Larry Rubinstein
Eric Polley (Mayo)
Bev Teicher

**Investigators and Patients at NIH
Clinical Center, NCI Cancer Centers,
ETCTN, and NCORP Sites Supplying
Tissues and Blood**

DOE and NIH Partnerships in Predictive Oncology

Frederick National Laboratory Advisory
Committee (FNLAC) Meeting
May 11, 2016

Rick Stevens[†] and Frank Alexander[‡]

[†]Argonne National Laboratory

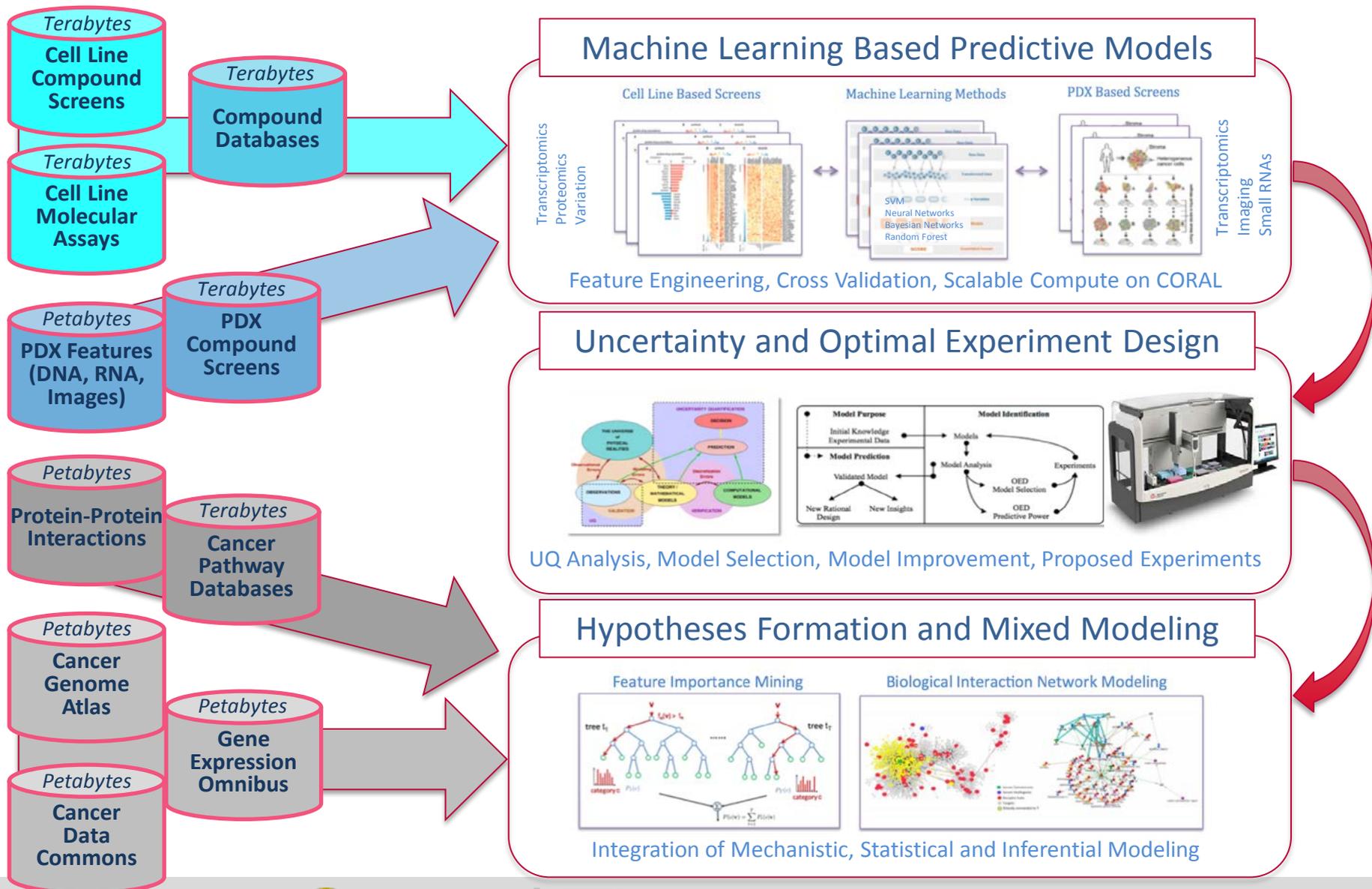
[†]University of Chicago

[‡]Los Alamos National Laboratory

Aims for Preclinical Screening Pilot

- Reliable machine learning based predictive models of drug response that enable the projection of screening results from and between cell-lines and PDX models
- Uncertainty quantification and optimal experimental design to assert quantitative limits on predictions and to recommend experiments that will improve predictions
- Improved modeling paradigms that support the graded introduction of mechanistic models into the machine learning framework and to rigorously assess the potential modeling improvements obtained thereof

Pilot 1: Predictive Models for Preclinical Screening



Preliminaries and Organization

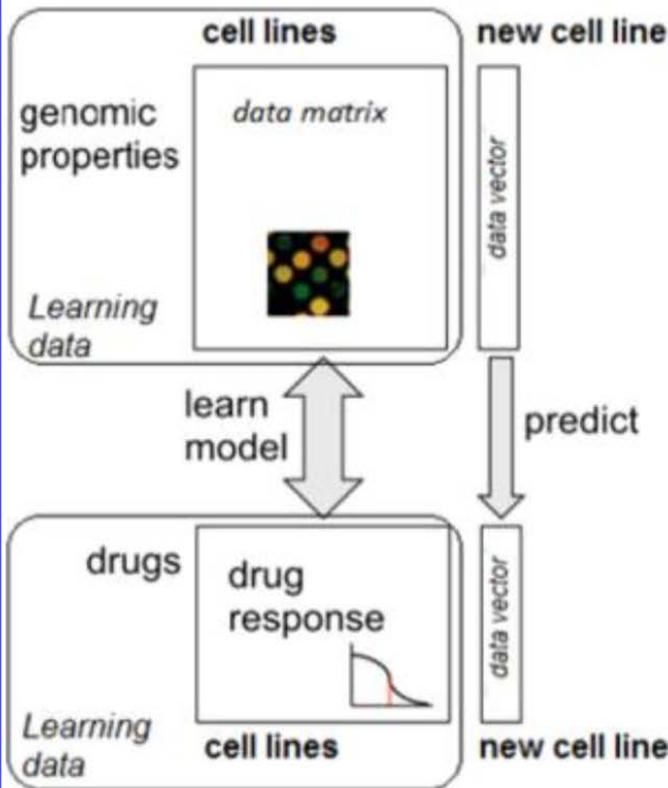
- **Established collaborative team with NIH and DOE laboratory members**
 - ~20 members: computer scientists, mathematicians, physicists, biologists, etc.
- **Established data sharing agreement, CDA and credentials for Lab members for NCI networks**
- **Established core working groups and initial development goals**
 - Data, Machine Learning, Experimental Design
- **Established communication plan, collaboration infrastructure and regular meeting schedule**
 - Bi-Weekly conference, Monthly three pilot conference
 - Quarterly face-to-face on Methods, Tools and Science
- **Developing three year pilot technical roadmap**
 - Technical goals and project work breakdown structure
 - Milestones and schedule
- **Demonstrating a few key early progress items**
 - Points of capability demonstration providing useful test cases

Goal 1: ML Models of Drug Response

- Pilot1 group are building models with NCI-60 data (extending to PDX) datasets using multiple ML methods to predict drug response
 - Binary encoding of response as well as regression models
- **Methods include: SVM, Lasso, Ridge, Bayesian, Random Forest, Extreme Gradient Boosting, Ada Boost, Neural Networks, etc.**
 - Models are constructed using NCI drug response screen data (e.g. IC50) and one or more molecular characterization data sets (mRNA expression, SNP array, EXM sequence, μ RNA, proteomics, metabolomics, methylation, etc.)
 - n-fold cross validation is done with many random splits to training and testing sets
 - ROC AUCs are computed and used to compare methods and parameters
- **PROGRESS: we have built and cross validated > 50K models, by end of year we will have a comprehensive collection of models for each agent and agent class that shows differential activity in NCI-60 and PDX collections**
 - Models have focused on gene expression but we are expanding the “feature space” for the models to include additional data types in the near future
 - Priority is to building models for recent combination screens

Two Modeling Approaches: By Cell Lines, By Drugs

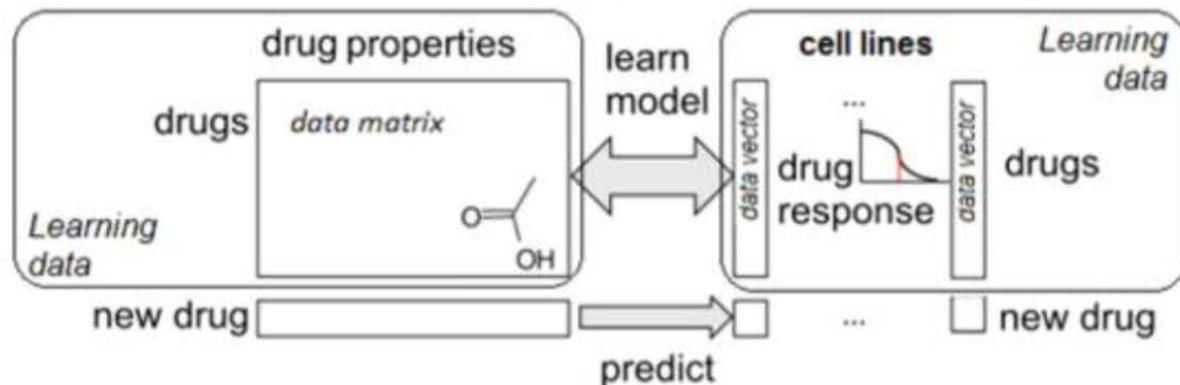
A. Generalization across cell lines



Steps in Modeling

- Feature Selection
- Outlier Removal
- Predictor Generation
- Cross Validation
- Independent Validation

B. Generalization across drugs



Goal 2: Catalog of Molecular Signatures

- For each agent or class of agents we will apply feature selection methods to the models to generate where possible a compact molecular signature that retains prediction performance
 - Typical reduced signatures include $O(10)$ - $O(100)$ features from $\gg 50,000$ starting features
 - Features may be genes, SNPs, μ RNA etc.
- Analysis of molecular signatures to provide insight to potential mechanisms
 - Enrichment analysis will be applied to the signatures to identify associated pathways
 - Pathways will be identified that associate with both sensitive and resistant response phenotypes
- **PROGRESS:** We have constructed compact signatures for many NCI-60 based models and started maps to pathways
 - Feature selection methods based on iterating xgboost and random forest that are successful at identifying compact (< 50) gene signatures that retain 90%-99% predictive capacity of much large feature sets (> 300)

Enrichment Analysis and Pathway Mappings

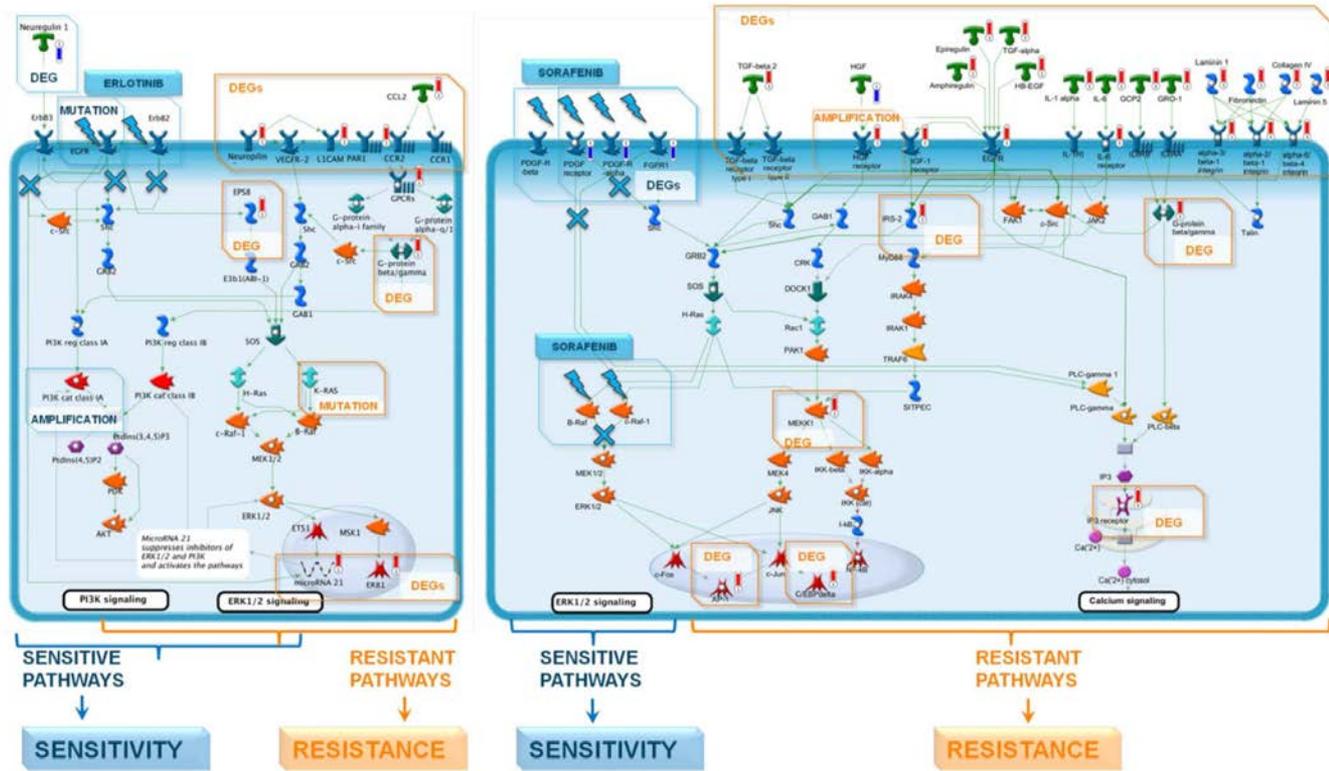


Fig 2. Causal network to depict functional relations between sensitivity-specific and resistance-specific signature genes. The network was reconstructed from canonical signaling pathways regulated by signature genes and a signature specific direct interaction network. Sensitivity-specific signature genes are highlighted with blue thermometers, resistance-specific genes with red thermometers.

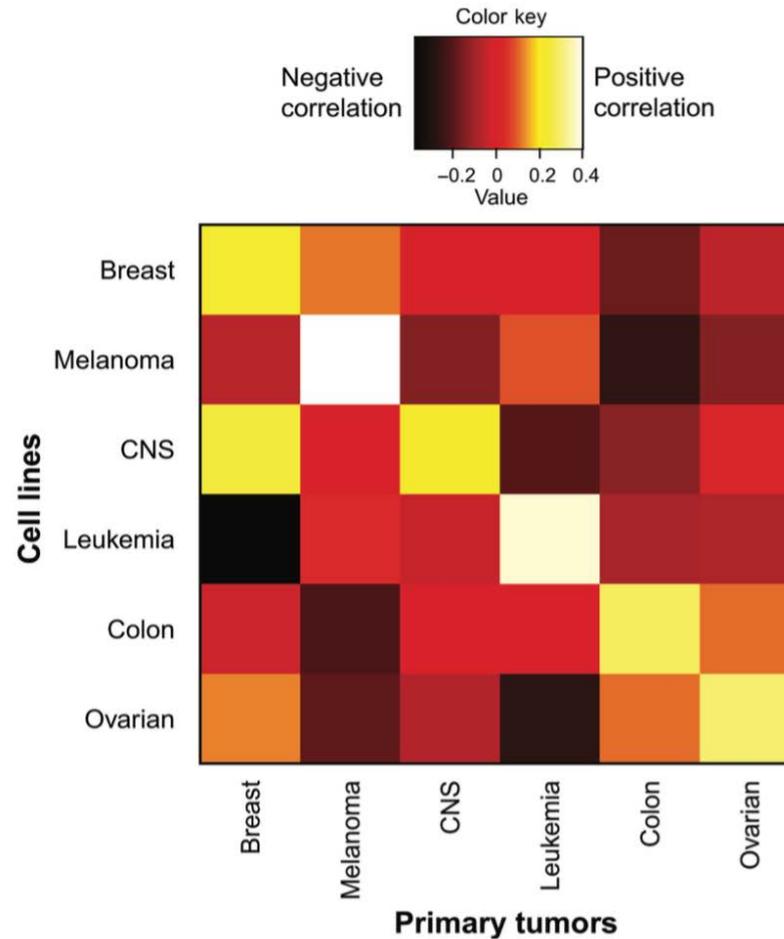
doi:10.1371/journal.pone.0130700.g002

Li B, Shin H, Gulbekyan G, Pustovalova O, Nikolsky Y, Hope A, et al. (2015) Development of a Drug-Response Modeling Framework to Identify Cell Line Derived Translational Biomarkers That Can Predict Treatment Outcome to Erlotinib or Sorafenib. PLoS ONE 10(6): e0130700. doi:10.1371/journal.pone.0130700

Goal 3: Test Signatures and Agents on PDX Models and other *in vitro* models

- Using NCI-60 derived signature/agent computational models to predict responses in PDX and other *in vitro* models
 - Provide input to multi-ARM PDX trials
 - Design test cases using optimal experimental design
 - Feed results into databases to revise the ML models
- Develop framework for high-throughput predictions with uncertainty quantification
 - Rank ordering of agents for specific tumor profiles
 - Qualitative (S/R) and quantitative (IC50) response predictions
- **PROGRESS:** we are evaluating several methods for extrapolating models from cell lines, PDX and primary tumors
 - Community has mixed results in projecting models based on expression
 - Including molecular data beyond gene expression may improve results

Cell Line Expression vs. Primary Tumor

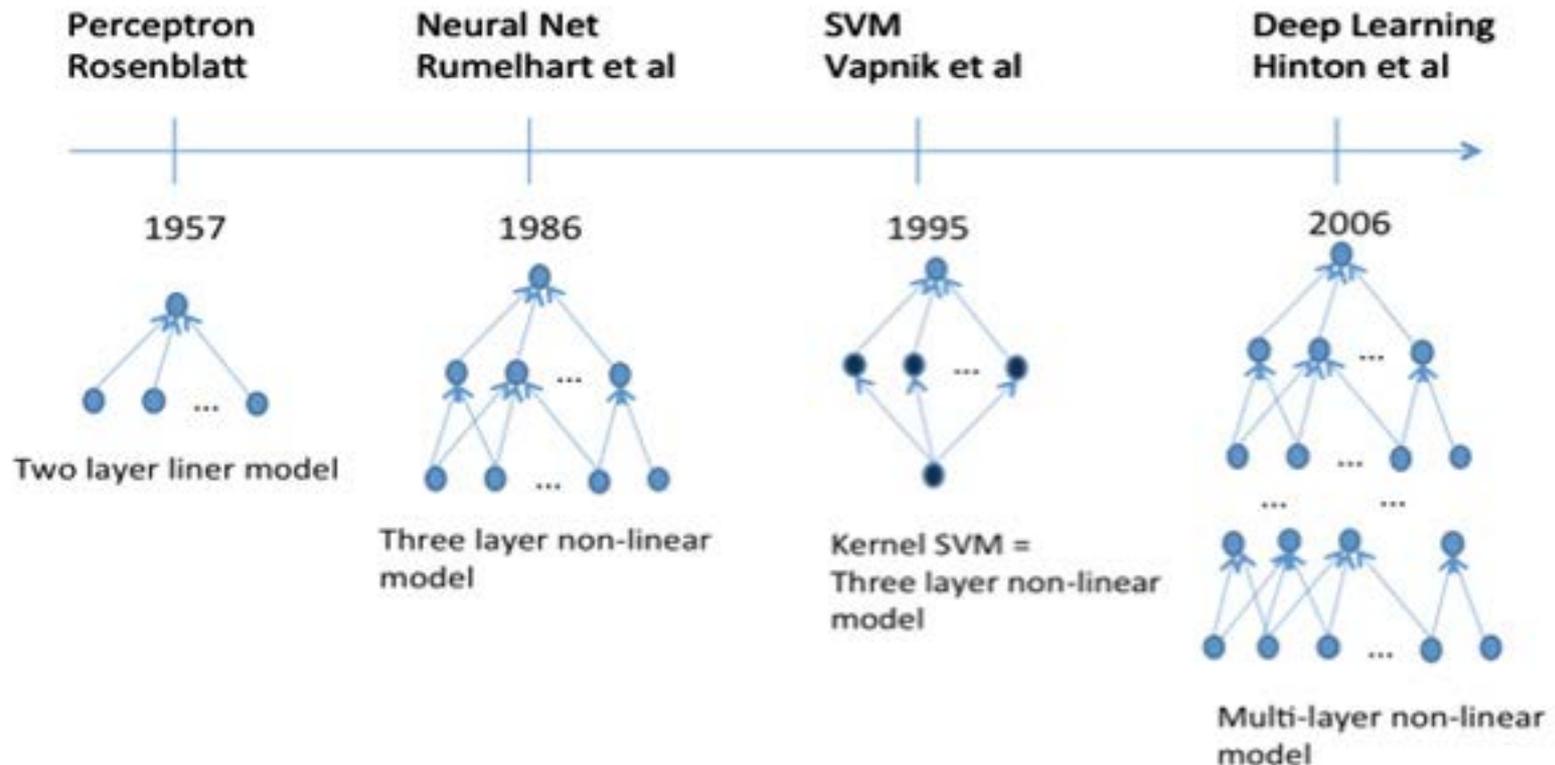


JNCI_J_Natl_Cancer_Inst_2013_Gillet.pdf

Goal 4: Deep Learning Based Models

- The pilot 1 team is developing a “deep learning” based model formulation that combines information about the drug, the drug target and the tumor/PDX or cell lines
- Goal is to improve drug response prediction by utilizing all available information sources (features) and combining data from many sources
- Deep learning has recently been successfully applied to problems in computer vision, game playing (Go), speech, handwriting, and many other areas
- Computer vendors are optimizing hardware to support deep learning
- There are positive early signs from the drug development community that deep learning methods can out perform traditional machine learning

Neural Network Complexity



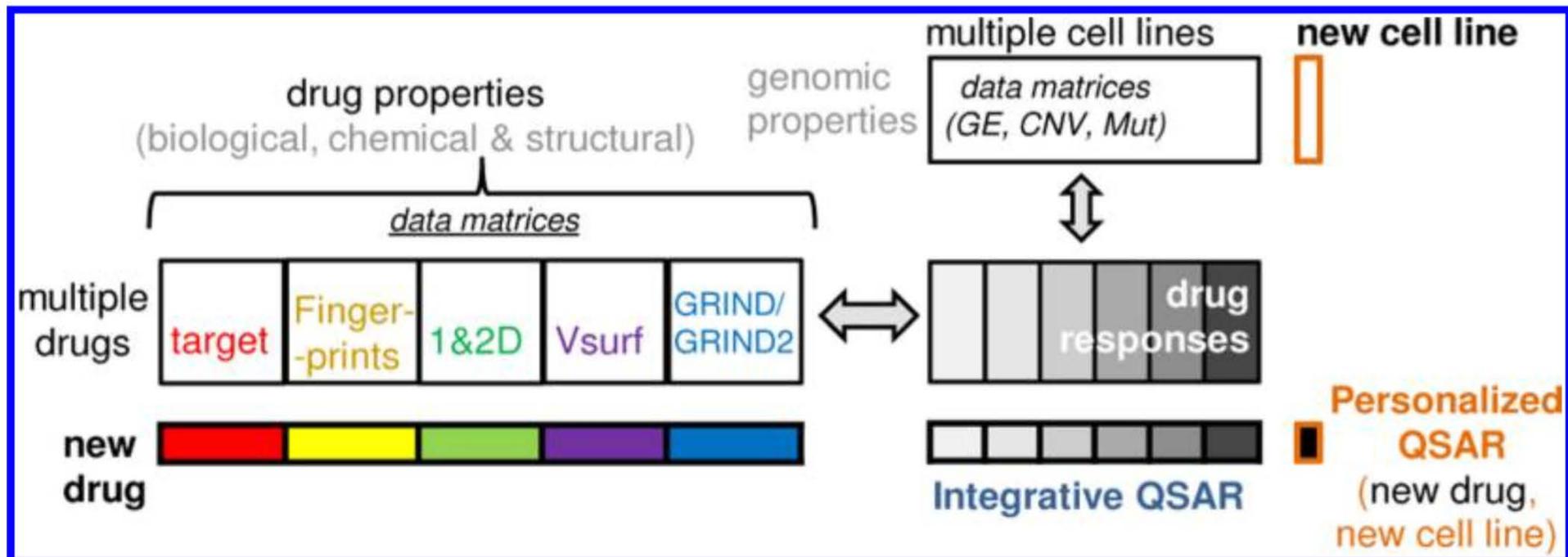
Deep Learning and Drug Screening @Johnson and Johnson

Jörg K. Wegner, Hugo Ceulemans, et. al. (NIPS2014)

“Deep learning outperformed all other methods with respect to the area under ROC (auc 0.83) curve and was significantly better than all commercial products. Deep learning surpassed the threshold to make virtual compound screening possible and has the potential to become a standard tool in industrial drug design.”

Deep Learning Formulation

$O(10^7)$ instances x $O(10^7)$ features



Drugs + Cell Lines \Rightarrow DNN \Rightarrow IC50

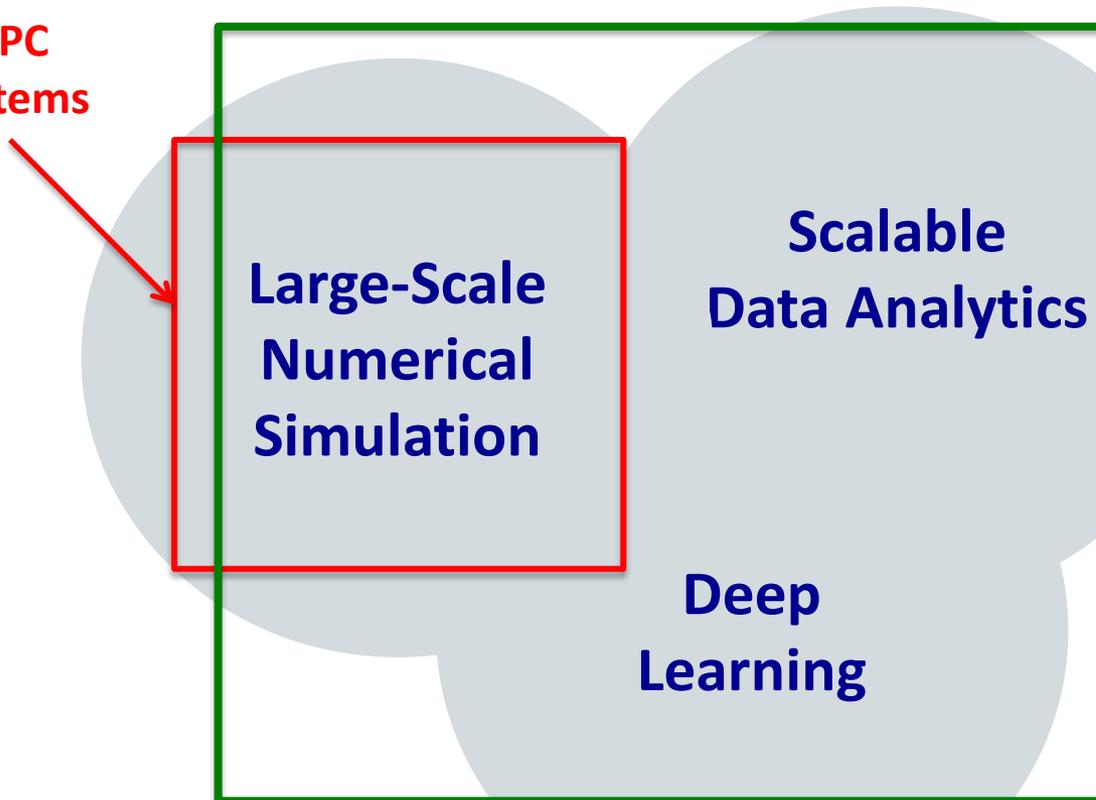
- Virtual screening new drugs on existing cell lines and PDX models
- Prediction of drug IC50 on new tumor/cell lines

Goal 5: Advanced Architectures

- A major goal for the DOE-NCI collaboration is to more deeply understand the architectural requirements for integration of simulation, scalable data analytics and deep learning
- The “use cases” from Pilots, 1, 2 and 3 will be used in evaluating the next generation supercomputers at DOE laboratories (Aurora, Summit and Sierra) and will provide feedback to designers for the planned exascale systems
- Key Objectives:
 - Architectural support for linking machine learning models to large scale simulations and ensembles of simulations
 - Community standard data analytics environments like “Apache Spark” on HPC platforms via container technologies
 - Scalable deep learning platforms that can fully utilize the interconnects and memory systems
 - Exploiting emerging non-volatile memory systems to accelerate applications

Integration of Simulation, Data Analytics and Machine Learning

Traditional
HPC
Systems



CORAL Supercomputers
And Exascale Systems

Acknowledgements

Yvonne A. Evrard

Fangfang Xia

Tom Brettin

Jim Davis

Maulik Shukla

Emily Dietrich

Monisha Gosh

Ian Foster

Venkat Vishwanath

John Santerre

Hal Finkel

Adam Zemla

Marisa Torres

Bob Olson

Susan Holbeck

Ravi Madduri

Marian Anghel

Cristina Garcia-Cardona

Judith Cohn

Amy Larson

Patrick Kelly

Mike Leuze

Intawat Nookaew

Arvind Ramanathan

Jonathan Allen

Ya Ju Fan

Chris Bun

Prasanna Balaprakash

BACKUP

Hybrid Models in Cancer

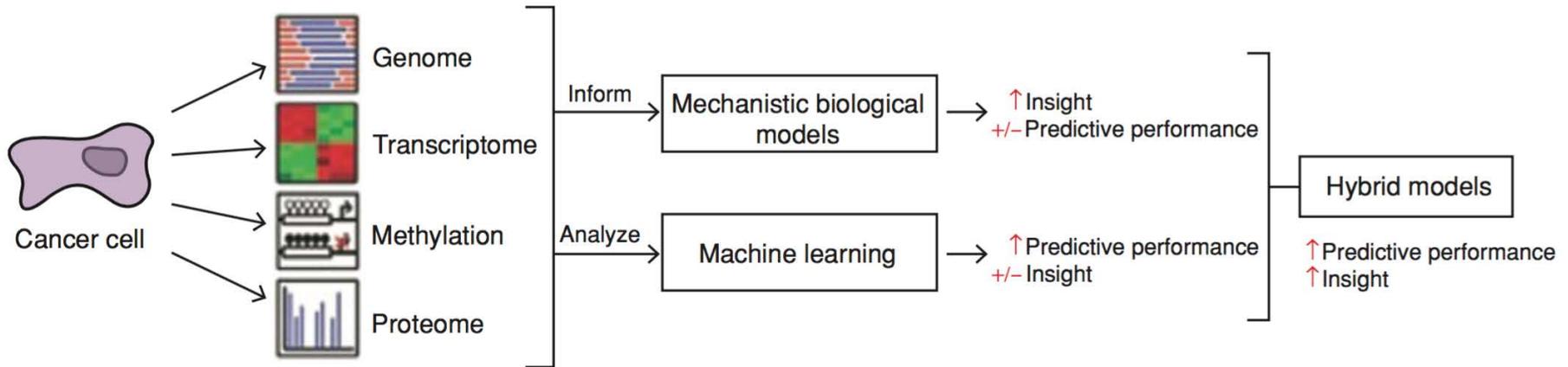


Figure 1. In two DREAM challenges, high throughput data characterizing cancer cells are used to build predictive models. Mechanistic models provide insight into the underlying biology, but do not take full advantage of the information within the data to achieve high performance. Machine learning methods are associative and extract maximum predictive value from the data, but do not always provide insight about mechanism. The future may bring hybrid models that combine the best of both approaches.

Predicting Cancer Drug Response: Advancing the DREAM

Russ B. Altman

Summary: The DREAM challenge is a community effort to assess current capabilities in systems biology. Two recent challenges focus on cancer cell drug sensitivity and drug synergism, and highlight strengths and weaknesses of current approaches. *Cancer Discov*; 5(3); 237-8. ©2015 AACR.