

# CSC 2431

# Machine Learning in Computational Biology

## *Pharmacogenomics*

### Benjamin Haibe-Kains

**Principal Investigator**, Bioinformatics and  
Computational Genomics Laboratory



The Princess Margaret  
Cancer Centre  
University Health Network

**Assistant Professor**, Medical Biophysics,  
University of Toronto



UNIVERSITY OF  
**TORONTO**

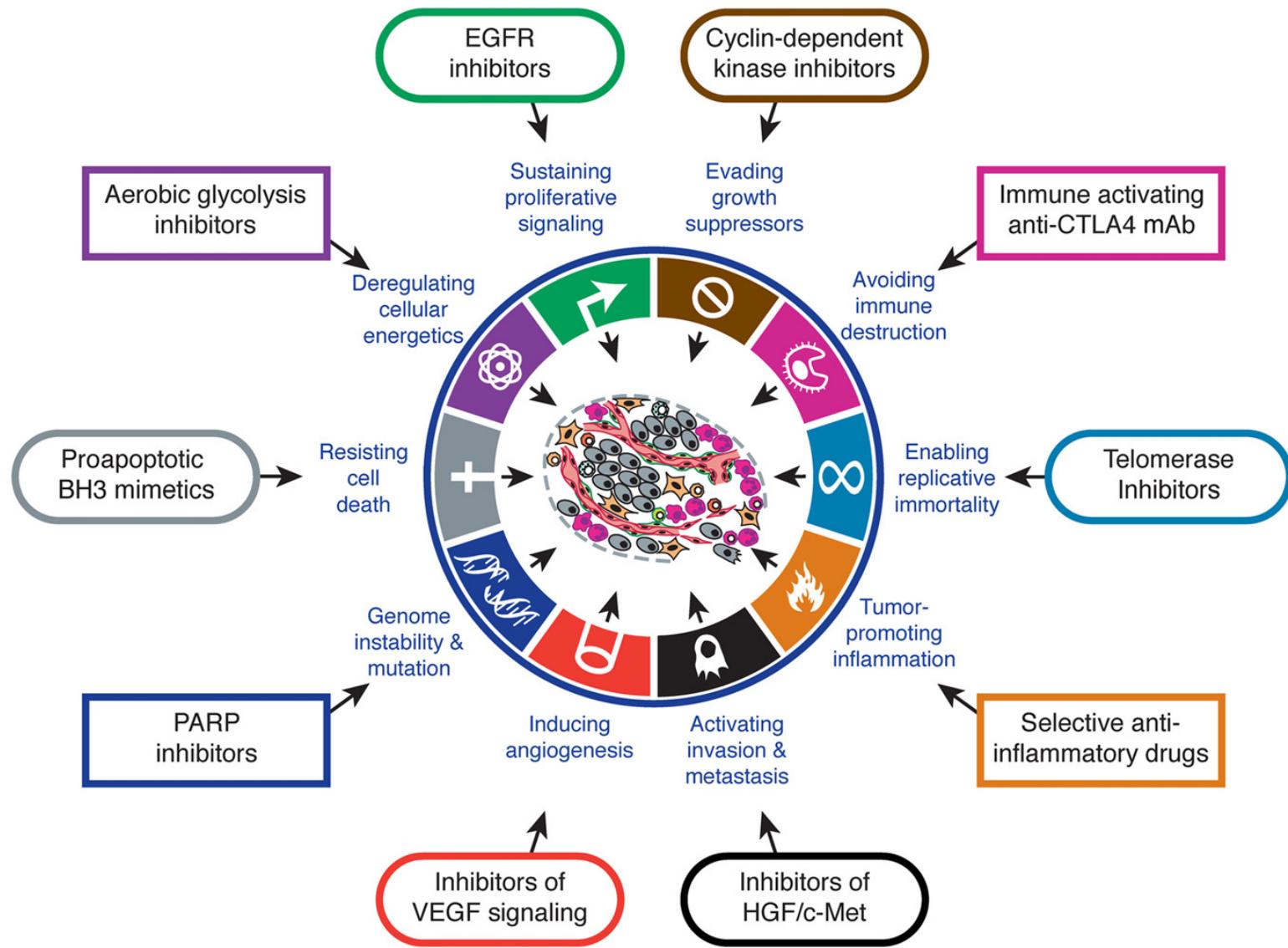
# Research interests

- Website: <http://www.pmgenomics.ca/bhklab/>
- Molecular subtyping
- Prognostic and predictive biomarkers (“signatures”)
- Pharmacogenomics:
  - Drug response prediction
  - Drug repurposing
  - Drug combination/synergy
- Tumor epi-stroma crosstalk

# Contents

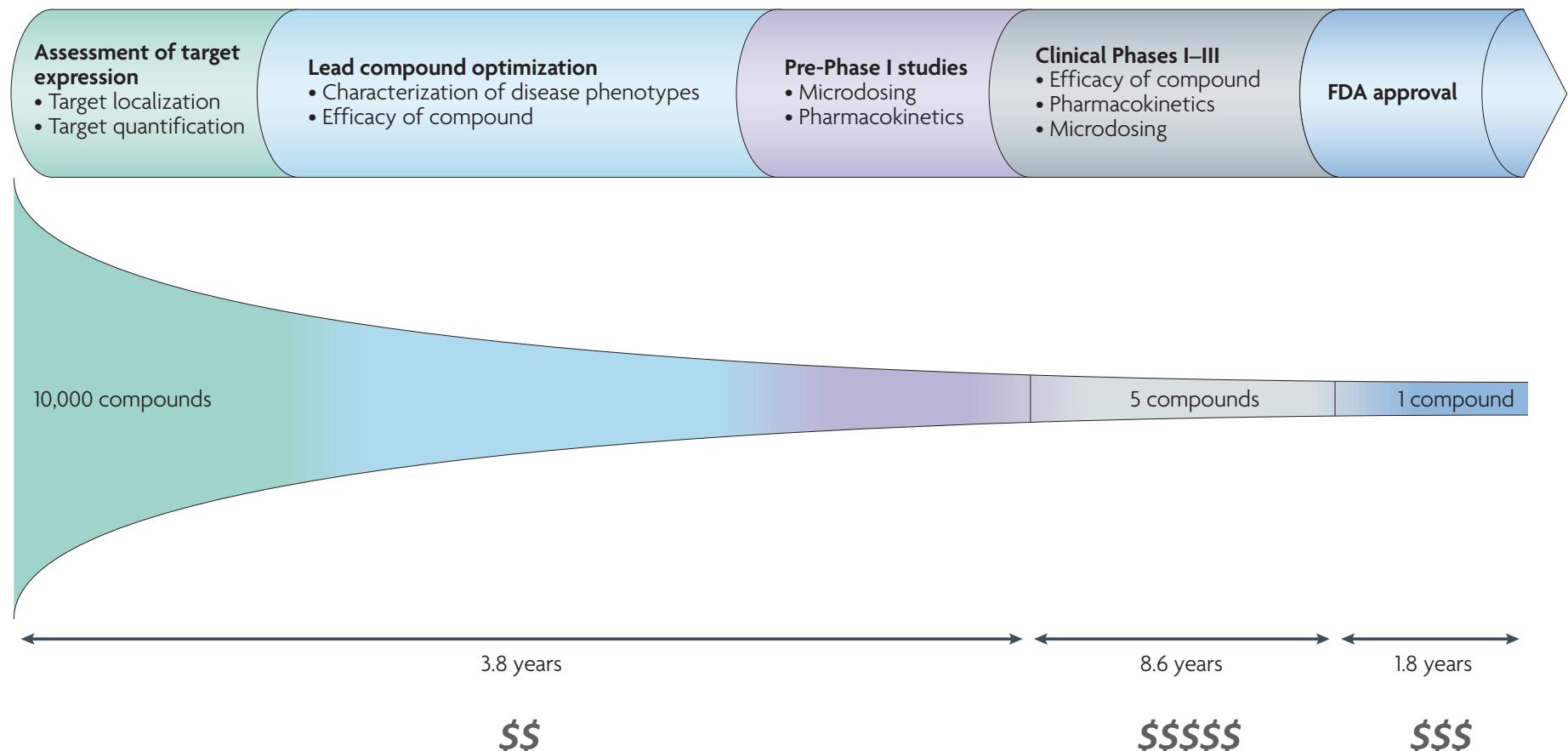
- Therapeutic strategies in cancer
- Biomarkers
- Basics of pharmacogenomic studies
- Machine learning in biomarker discovery
  - Previous efforts
  - Challenges
- Conclusion

# Therapeutic strategies in cancer



Hanahan and Weinberg, Cell, 2011

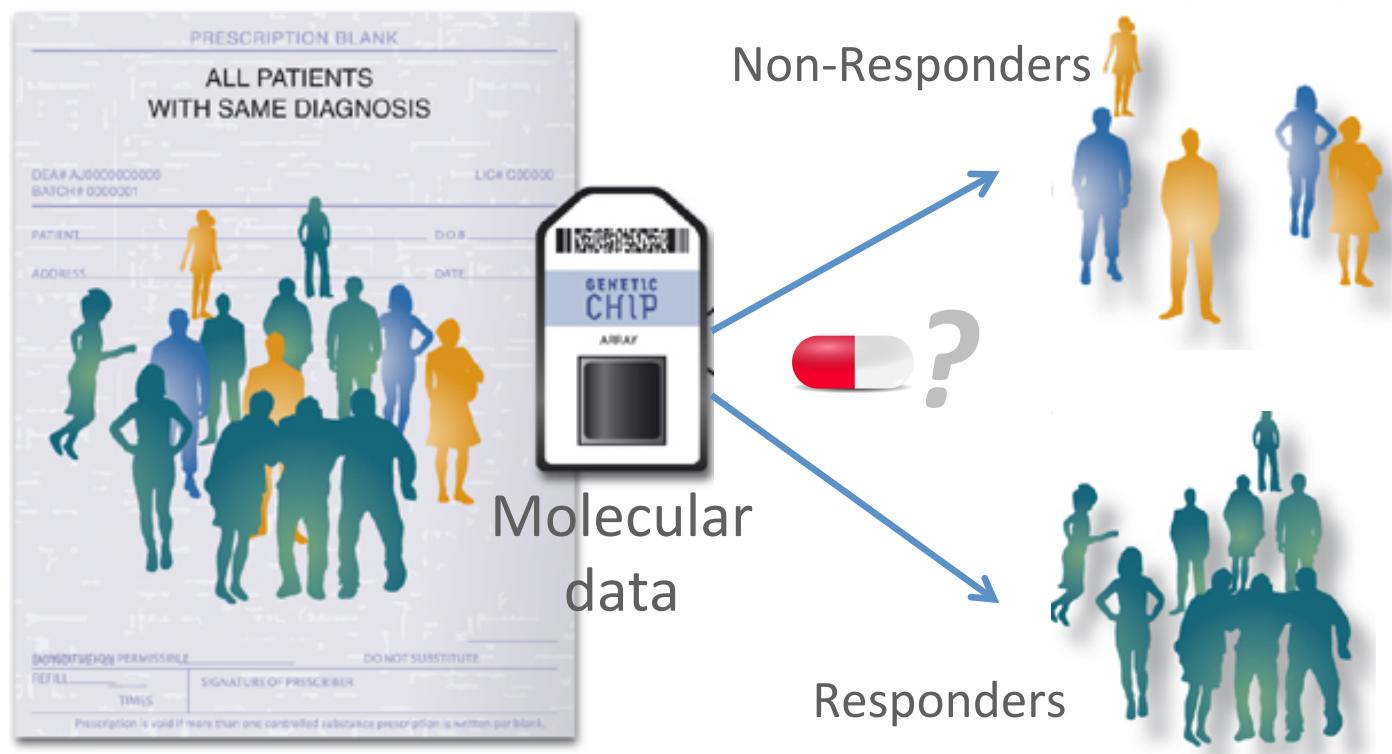
# Drug development



Willman et al., *Nature Reviews Drug Discovery* 2008

# Model systems used for biomarker discovery

- Cost for developing new (targeted) drugs is high and number of drugs being approved is dramatically slowing down  
→ ***Need for companion tests to identify patients who are likely to respond to targeted therapies***



# Biomarkers

- A **biomarker** is anything that can be used as a measurable indicator of
  - a particular disease state → diagnostic or prognostic biomarker
  - a particular response to a drug → predictive biomarker
- Biomarkers can be measured using different technologies
  - DNA sequencing, gene/protein expression, DNA methylation, etc.
  - Low-throughput: usually univariate biomarkers
  - High-throughput: multivariate biomarkers, often referred to as *signatures* (10-100 features)

# Examples of univariate biomarkers

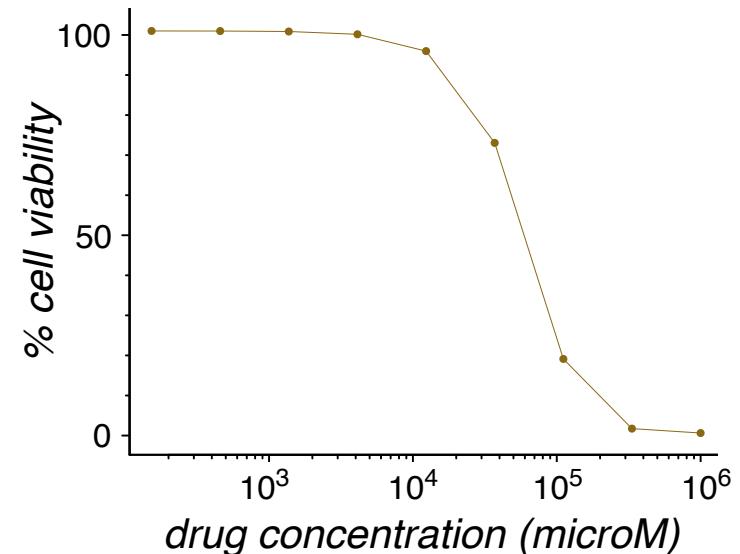
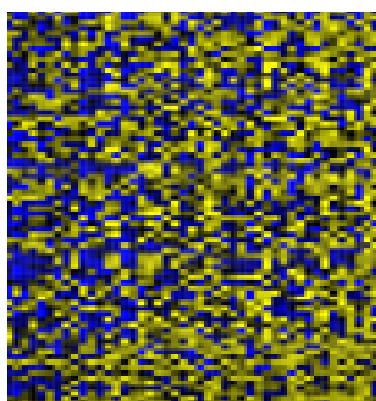
- Genetic aberrations:
  - ✓ BRAF V600E mutation predicts response to vemurafenib and dabrafenib in melanoma
  - ✓ BCR-ABL1 gene fusion predicts response to nilotinib in chronic myelogenous leukemia
- Protein expressions:
  - ✓ ER expression predicts response to tamoxifen in breast cancer
- Gene/transcript expressions:
  - ✓ NQO1 expression predicts response to tanespimycin in multiple cancer types
  - ✓ ERBB2 expression predicts response to lapatinib in breast cancer

# Model systems used for biomarker discovery

- *In situ*: Patient tumors
- *In vivo*:
  - Patient-derived xenografts (PDX)
  - Genetically-engineered mouse models (GEMM)
- *In vitro*: Cancer cell lines
  - **Cheap and high-throughput**
  - **Sensitivity data available for many approved and experimental drugs**
  - **But no cell lines are like real tumors**

# Cell-based drug screening studies

## cell lines

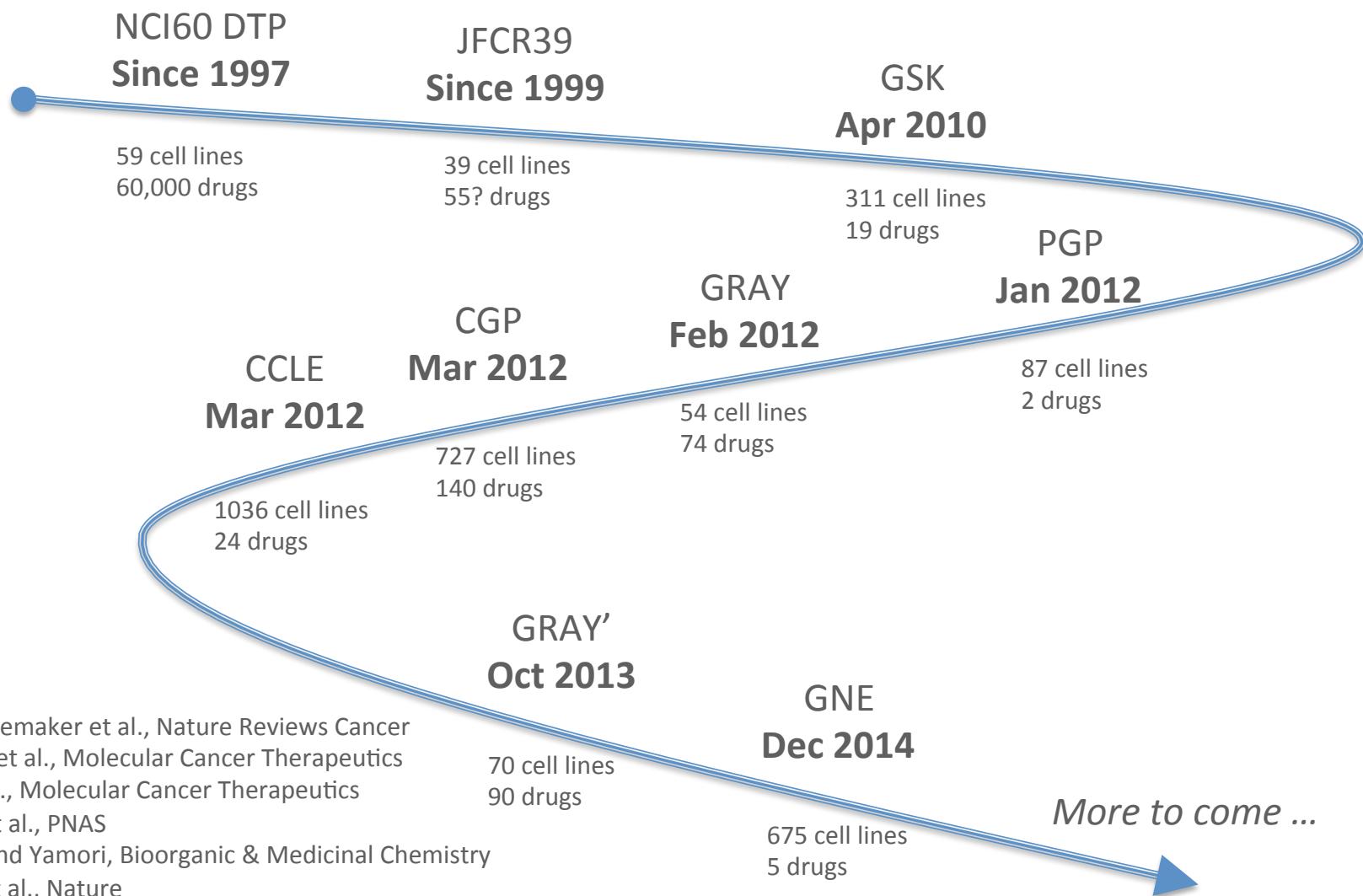


*molecular  
data*

*sensitivity  
(growth inhibition)*

→ resistant vs. sensitive cell lines

# Pharmacogenomic datasets



NCI60 DTP: Shoemaker et al., Nature Reviews Cancer

GSK: Greshock et al., Molecular Cancer Therapeutics

PGP: Hook et al., Molecular Cancer Therapeutics

GRAY: Heiser et al., PNAS

JFCR39: Kong and Yamori, Bioorganic & Medicinal Chemistry

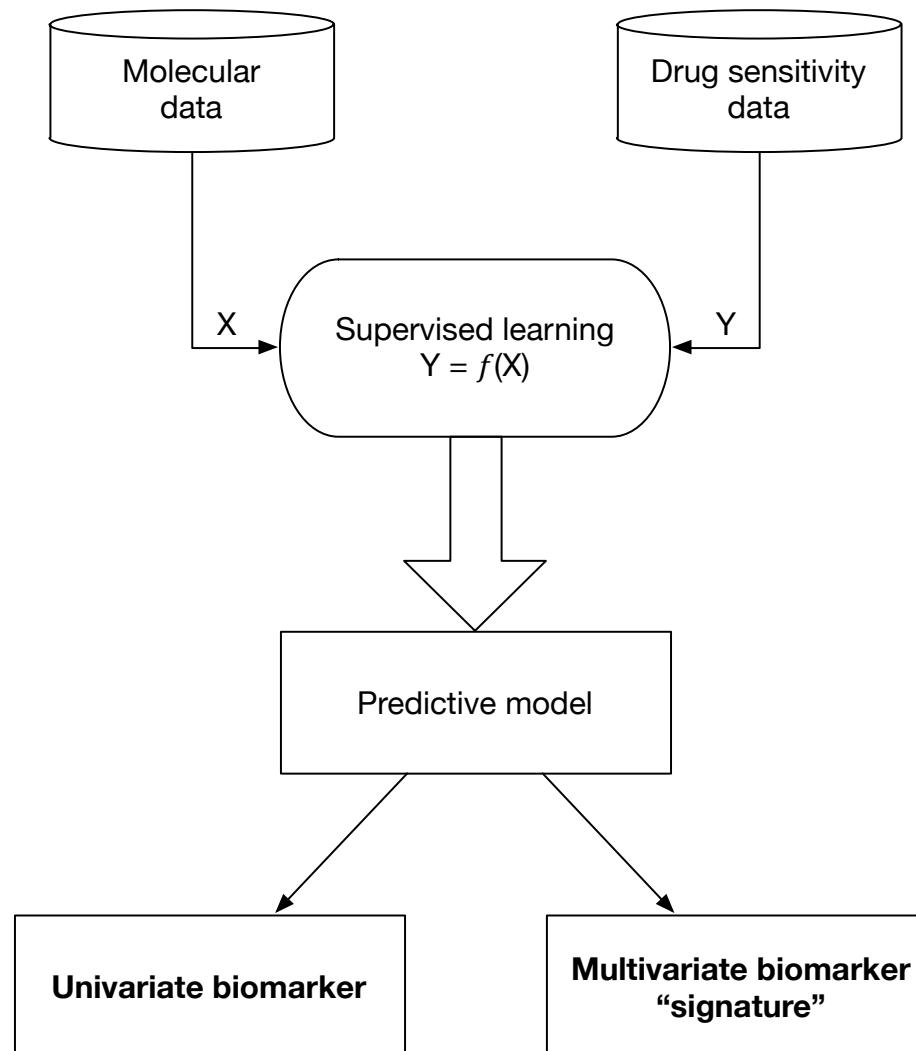
CGP: Garnett et al., Nature

CCLE: Barretina et al., Nature

GRAY updated: Daemen et al., Genome Biology

GNE: Klijn et al., Nature Biotechnology

# Machine learning in biomarker discovery



# First data mining approaches on NCI60

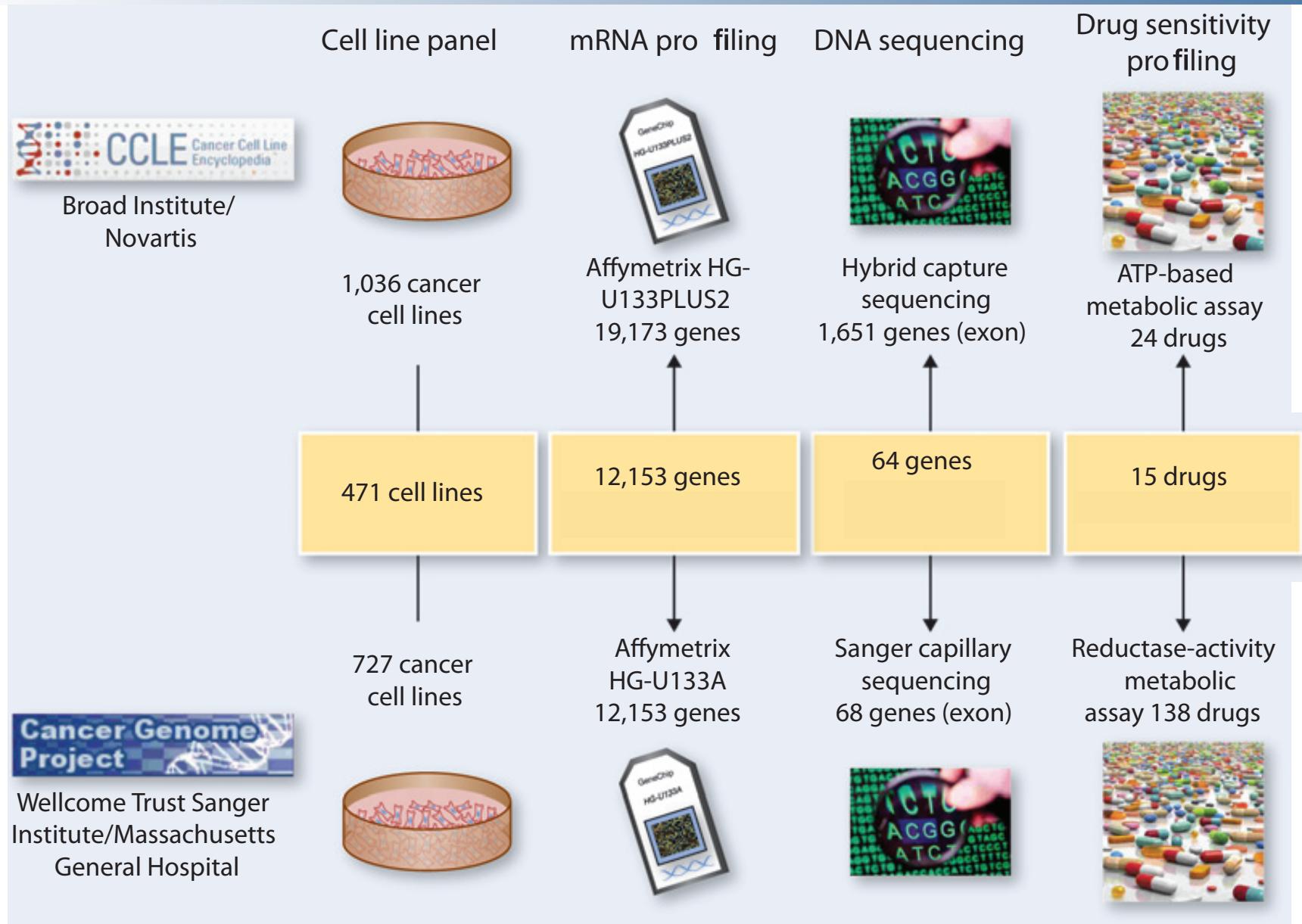
- **COMPARE algorithm** (Paull et al, *JNCI* 1998): unsupervised learning method to cluster drugs with similar IC<sub>50</sub> measures in the NCI60 panel, implemented in **CellMiner** (Shankavaram et a, *BMC Genomics* 2009)
- **Simple supervised learning** (Staunton, *PNAS* 2001):
  1. Define drug phenotypes (resistant vs. sensitive)
  2. Select most relevant features (high correlation between gene expressions and drug phenotypes)
  3. Build ensemble of univariate models (weighted voting classification)
  4. Assess accuracy in cross-validation

# Building and testing predictors with CGP/CCLE

- The Cancer Genome Project (**CGP**) initiated by the Sanger Institute
  - **138** drugs
  - **727** cancer cell lines
- The Cancer Cell Line Encyclopedia (**CCLE**) initiated by Novartis/Broad Institute
  - **24** drugs
  - **1036** cancer cell lines
- Both used regularized regression (*elasticnet*) with mutations, CNVs and gene expressions used as input
- Several summary measures of the dose-response curves were used as output



# CCLE ∩ CGP



# CCLE ∩ CGP: Shared drugs

<b>Paclitaxel</b>	Microtubules depolymerization inhibitor
<b>PD-0325901, AZD6244</b>	Mitogen-activated protein kinase kinase (MEK) inhibitor
<b>AZD0530 (Saracatinib)</b>	Proto-oncogene tyrosine-protein Src inhibitor
<b>Nutlin-3</b>	Ubiquitin-protein ligase MDM2 inhibitor
<b>Nilotinib</b>	BCR-ABL fusion protein inhibitor
<b>17-AAG (Tanespamycin)</b>	Heat shock protein (Hsp90) inhibitor
<b>PD-0332991</b>	CDK4/6-Cyclin D inhibitor
<b>PLX4720, Sorafenib</b>	RAF kinase inhibitors
<b>Crizotinib, TAE684</b>	ALK kinase inhibitors
<b>Erlotinib, Lapatinib</b>	EGFR/HER2 kinase inhibitors
<b>PHA-665752</b>	Proto-oncogene c-MET kinase inhibitor

# Studies leveraging CGP/CCLE

JAMIA

## Comparison and validation of genomic predictors for anticancer drug sensitivity

Simon Papillon-Cavanagh,<sup>1</sup> Nicolas De Jay,<sup>1</sup> Nehme Hachem,<sup>1</sup> Catharina Olsen,<sup>2</sup> Gianluca Bontempi,<sup>2</sup> Hugo J W L Aerts,<sup>3,4</sup> John Quackenbush,<sup>4</sup> Benjamin Haibe-Kains<sup>1</sup>

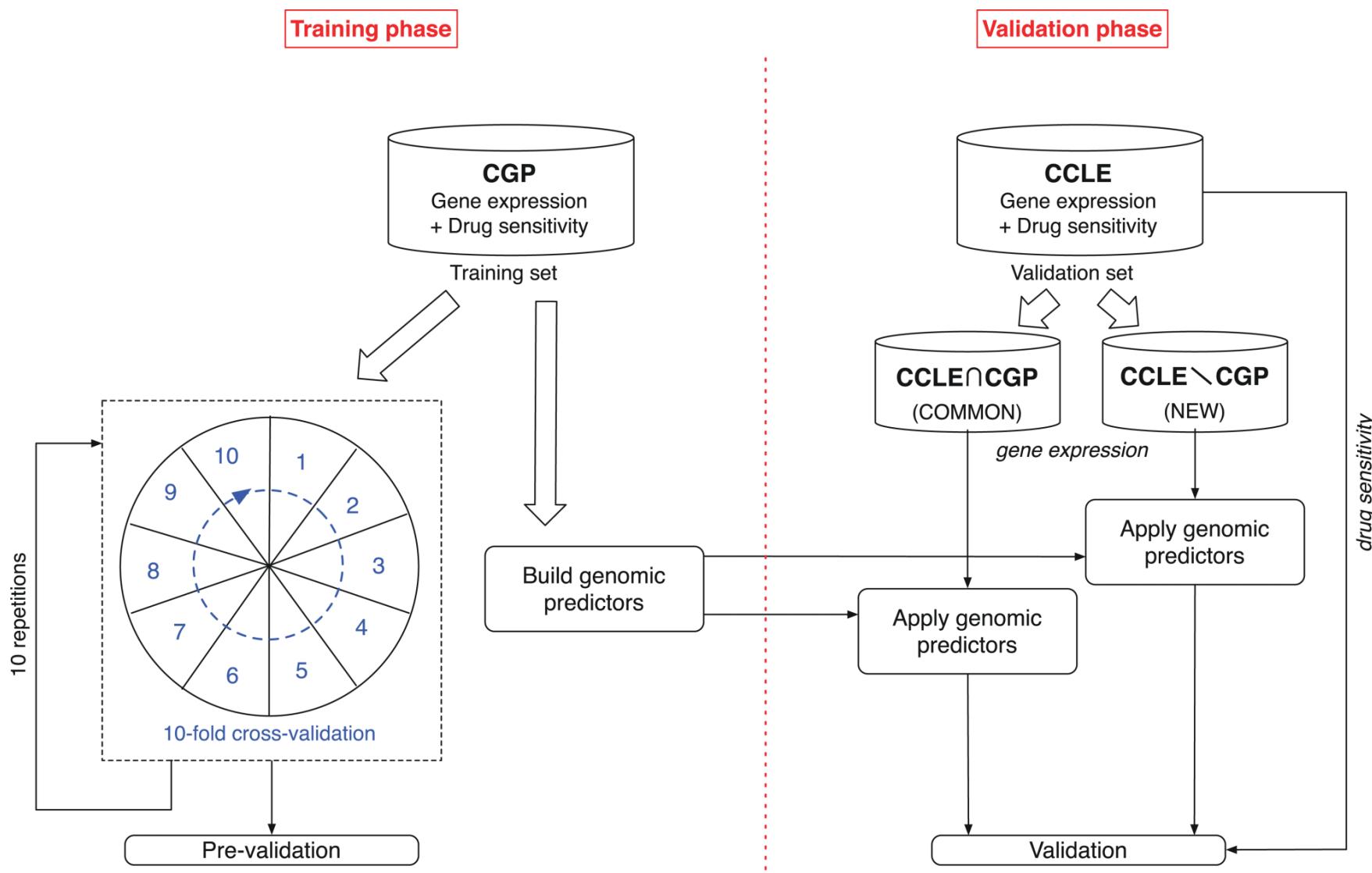
Papillon-Cavanagh S, et al. *J Am Med Inform Assoc* 2013;0:1–6. doi:10.1136/amiajnl-2012-001442

*Pacific Symposium on Biocomputing 2014*

## SYSTEMATIC ASSESSMENT OF ANALYTICAL METHODS FOR DRUG SENSITIVITY PREDICTION FROM CANCER CELL LINE DATA\*

IN SOCK JANG<sup>1</sup>, ELIAS CHAIBUB NETO, JUSTIN GUINNEY, STEPHEN H. FRIEND, ADAM A. MARGOLIN<sup>1</sup>

# Papillon-Cavanagh et al, experimental design

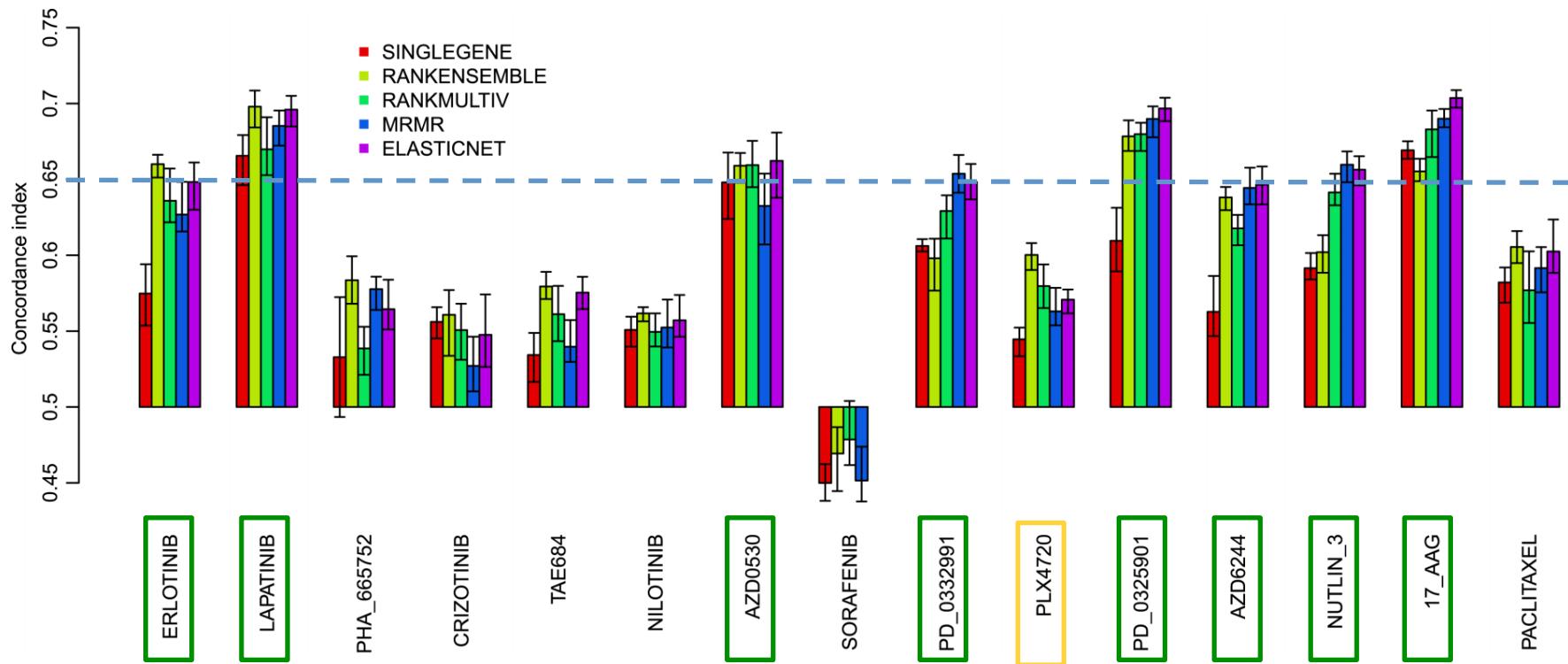


# Papillon-Cavanagh et al, modeling

- **SINGLEGENE**: Univariate linear regression model with the gene the most correlated to drug sensitivity  $[-\log_{10}(\text{IC}_{50})]$
- **RANKENSEMBLE**: Average of the predictions of the top 30 models
- **RANKMULTIV**: Multivariate model with the top 30 genes
- **MRMR**: Multivariate model with the 30 most correlated and less redundant genes
- **ELASTICNET**: Regularized multivariate model (L1/L2 penalization)

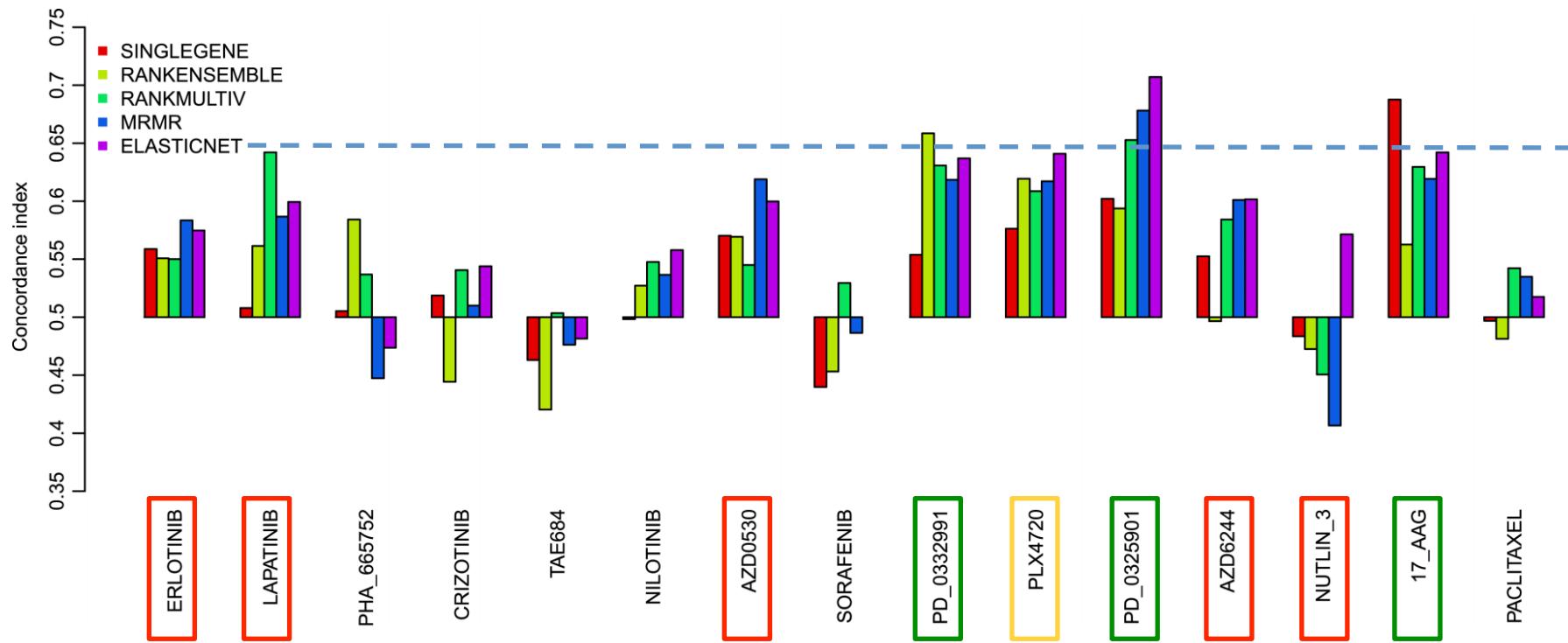
# Papillon-Cavanagh et al, results in cross-val

## CGP in 10-fold cross-validations



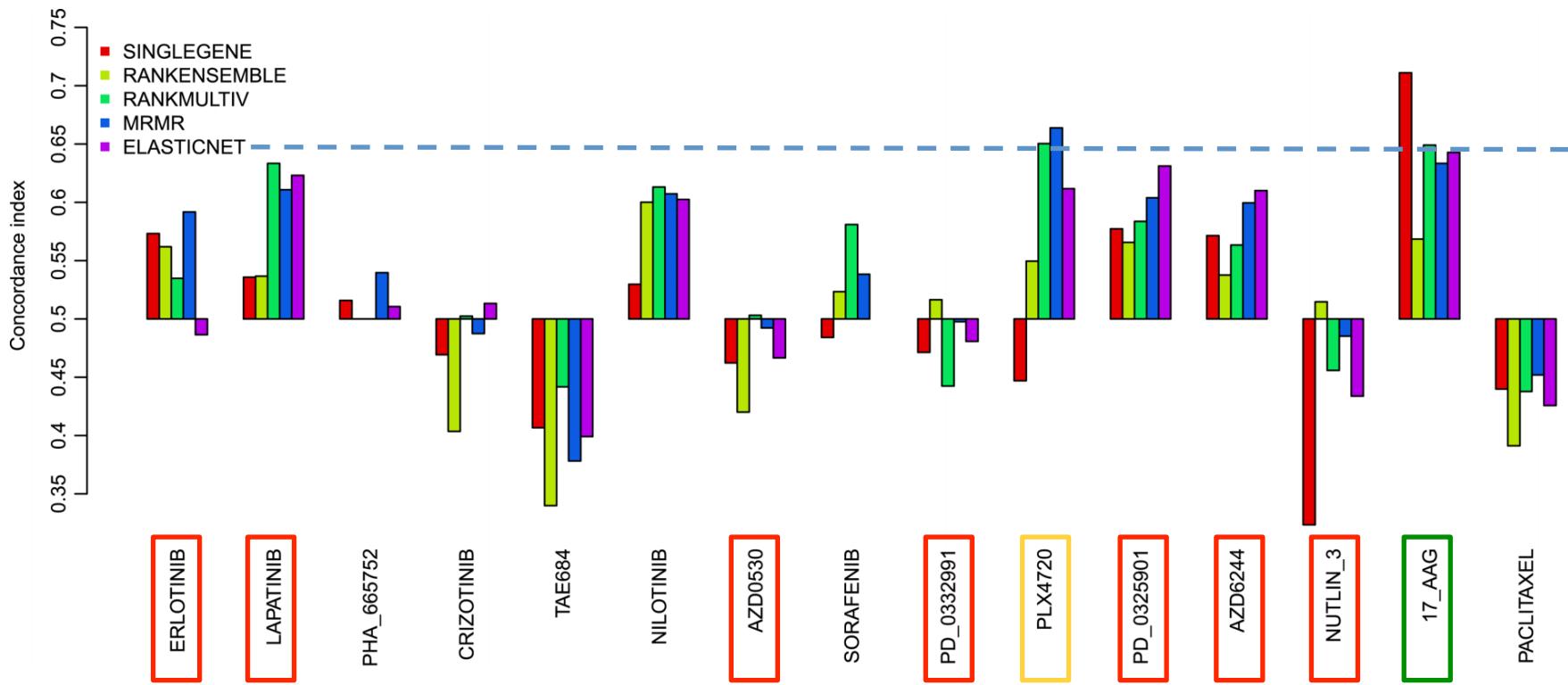
# Papillon-Cavanagh et al, test set #1

Trained on CGP, tested on CCLE  
**Common cell lines**



# Papillon-Cavanagh et al, test set #2

Trained on CGP, tested on CCLE  
**New cell lines**



# Jang et al, experimental design

- Evaluation of >110,000 predictive modeling strategies to establish best practices in biomarker discovery
- 5 experimental factors:
  - *Molecular features*
  - *Compound*
  - *Response summary*
  - *Continuous vs. categorical models*
  - *Algorithms*
- All combinations are tested in 5-fold cross-validation

# Jang et al, results

- Factors explaining most of the variance:
  - the type of *molecular features* used to build the model
  - the *compound* being predicted by the model
  - then *algorithm*
- Strong interaction between compound and response summary factors in CCLE but not in CGP
- Gene expression data are the most predictive
- Elasticnet and ridge regression are the best performers

# NCI-DREAM Drug Sensitivity Prediction Challenge

- Training set: multiassay molecular profiling of **35** breast cancer cell lines (mutations, CNV, DNA methylation, gene and protein expressions) treated with **28 drugs**
- Test: A community effort to assess and improve drug sensitivity prediction algorithms
- Go: James C Costello<sup>1,2,13,14</sup>, Laura M Heiser<sup>3,14</sup>, Elisabeth Georgii<sup>4,14</sup>, Mehmet Gönen<sup>4</sup>, Michael P Menden<sup>5</sup>, Nicholas J Wang<sup>3</sup>, Mukesh Bansal<sup>6</sup>, Muhammad Ammad-ud-din<sup>4</sup>, Petteri Hintsanen<sup>7</sup>, Suleiman A Khan<sup>4</sup>, John-Patrick Mpindi<sup>7</sup>, Olli Kallioniemi<sup>7</sup>, Antti Honkela<sup>8</sup>, Tero Aittokallio<sup>7</sup>, Krister Wennerberg<sup>7</sup>, NCI DREAM Community<sup>9</sup>, James J Collins<sup>1,2,10</sup>, Dan Gallahan<sup>11</sup>, Dinah Singer<sup>11</sup>, Julio Saez-Rodriguez<sup>5</sup>, Samuel Kaski<sup>4,8</sup>, Joe W Gray<sup>3</sup> & Gustavo Stolovitzky<sup>12</sup>
- Predict compounds
- Performance estimator: weighted average of probabilistic concordance indices across drugs

## ANALYSIS

computational  
BIOLOGY

VOLUME 32 NUMBER 12 DECEMBER 2014 NATURE BIOTECHNOLOGY

# NCI-DREAM, the top predictor is quite fancy

- Bayesian multitask multiple kernel learning method that leveraged four machine-learning principles:
  - **kernelized regression** computes outputs from similarities between cell lines
  - **Bayesian inference** to learn drug-specific parameters of the kernelized regression
  - **multiview learning** to combine different “views” of the data (data discretization, pathway-based summarization, data combination, ...)
  - **multitask learning** to simultaneously model kernel weights based on drug sensitivities across all the drugs

# NCI-DREAM, the second best predictor is less fancy

- Ensemble combination of the random forests trained on each data type
- Data type-specific random forest predictions weighted by their  $R^2$
- As usual, ensemble classifiers combining all predictions performed better (wisdom of crowds)
- Gene expression data were the most informative followed by methylation
- Other teams found the proteomic (RPPA) data being very predictive too

# Cell lines to patients

- Geeleher et al. used CGP to build predictors of drug response
  - *Ridge* and *elasticnet* regression models
  - *Lasso* logistic regression model using extreme drug phenotypes
- They tested the predictive value of their models in patient tumors from clinical trials
  - Decent predictive ability for docetaxel in breast cancer but poor predictions for the other drugs

Paul Geeleher<sup>1</sup>, Nancy J Cox<sup>2</sup> and R Stephanie Huang<sup>1\*</sup>

*Genome Biology* 2014, **15**:R47

# Challenges

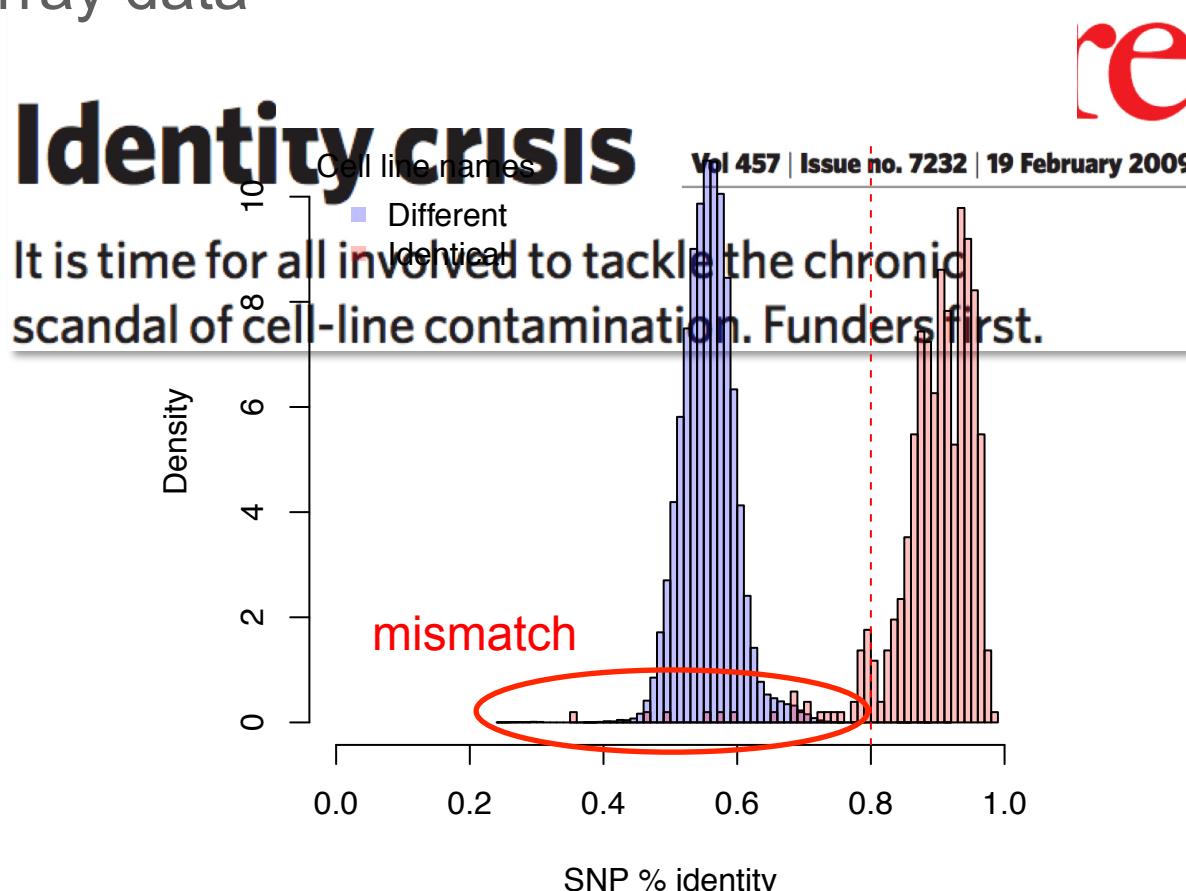
- Over 1000 PubMed articles just for gene expression-based biomarkers in cancer cell lines
- But few biomarkers translated into clinical settings ...
- Challenges:
  - Analysis reproducibility and validation
  - Cell line identity crisis
  - (In)Consistency in drug phenotypes
  - Experimental protocols
  - Summary measures for drug phenotypes

# Analysis reproducibility and validation

- Irreproducible research, see Anil Potti's scandal at Duke University  
(Potti et al, *Nat Med* 2011) **\*retracted**
- Unsuccessful validation of MPI predictors built on the NCI60 cell line panel to predict therapy response  
(Wang et al., *JNCI* 2013)

# Cell line identity crisis

- Mislabeling or contamination of the cancer cell lines
- One can check the identity of the cell lines by leveraging the SNP array data



# (In)Consistency of pharmacogenomic data

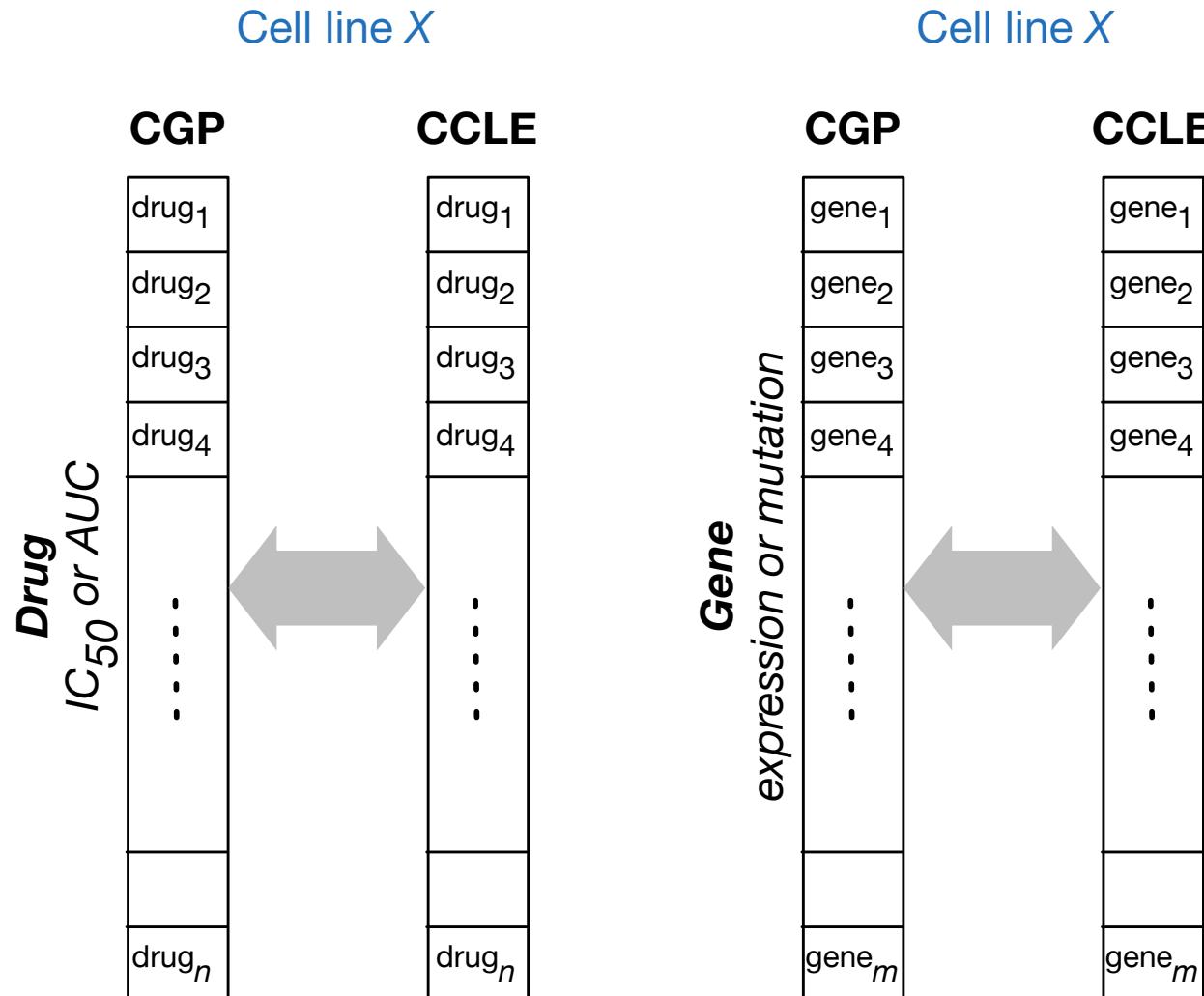
- **Input data** → Molecular profiles
  - For each [cell line, gene], are *molecular measurements* consistent across studies?
- **Output data** → Drug phenotypes
  - For each [cell line, drug], are *sensitivity measurements* consistent?



Nehme El-Hachem

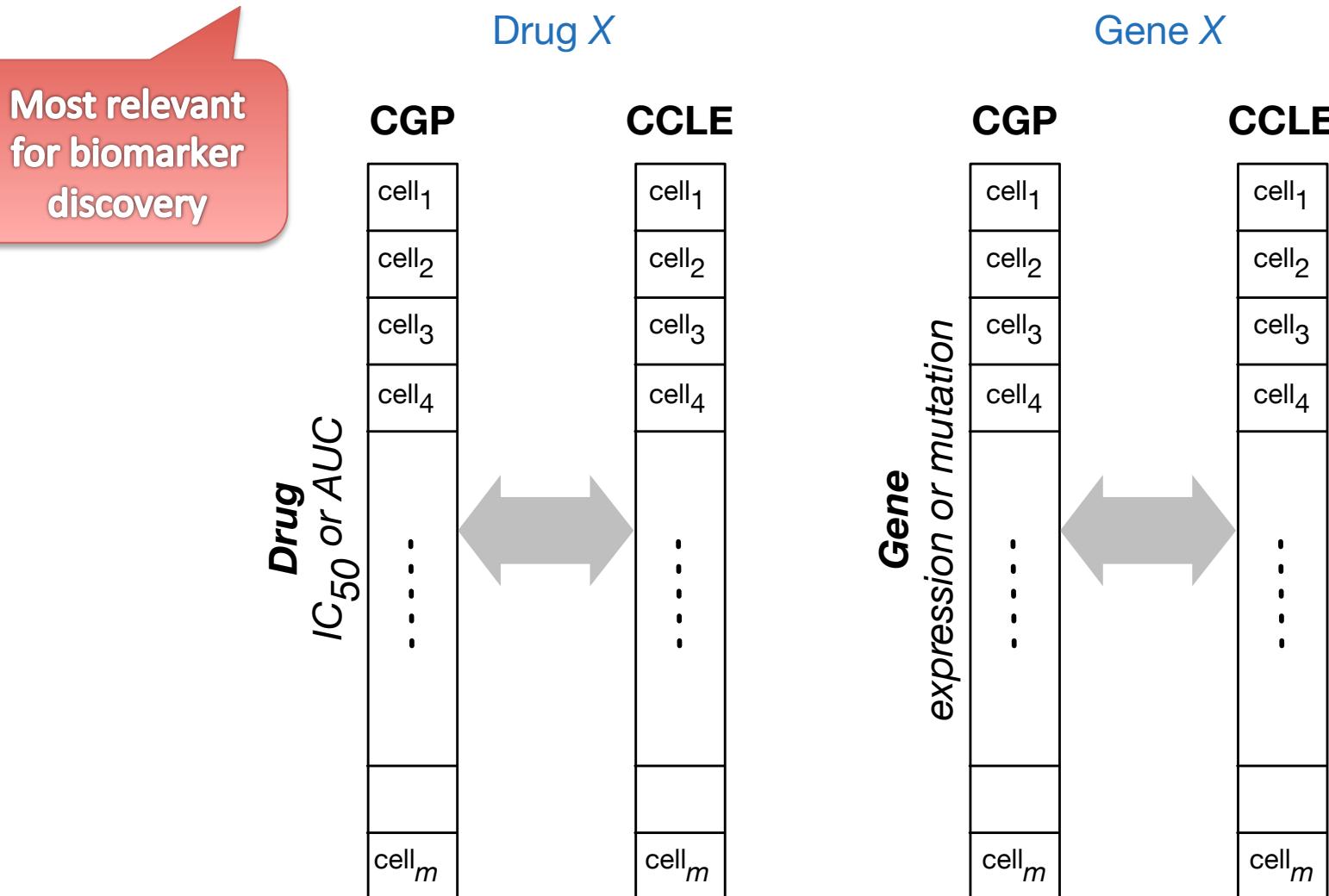
# Concordance “between cell lines”

- Identical cell lines in CGP and CCLE should exhibit similar profiles



# Concordance “across cell lines”

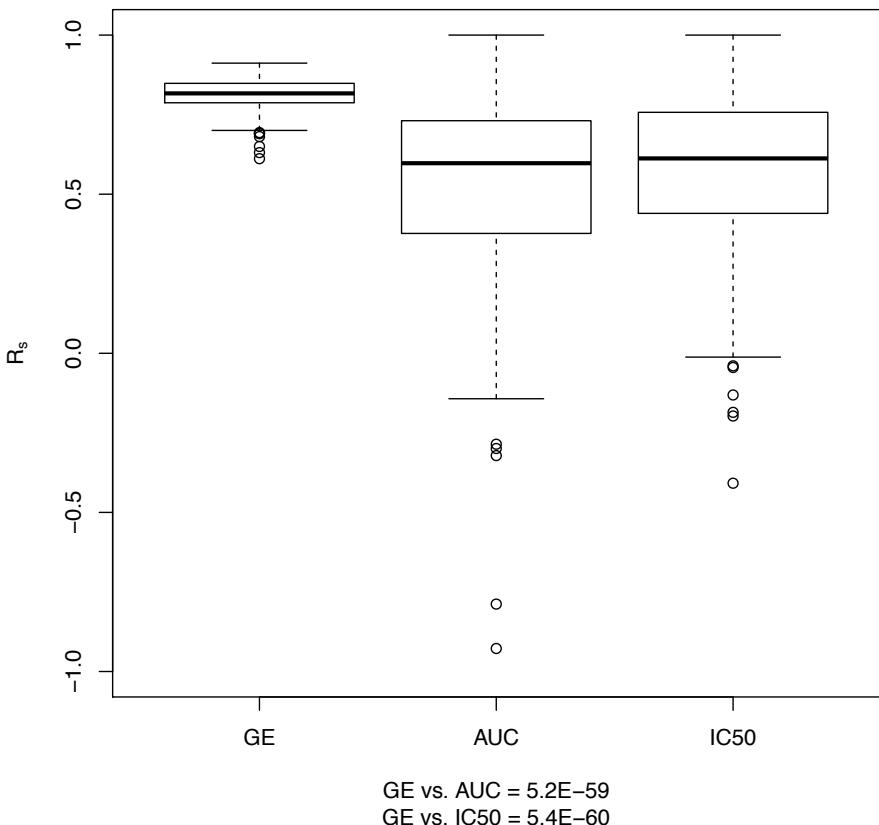
- Identical drugs or genes in CGP and CCLE should exhibit similar measurements in cell lines



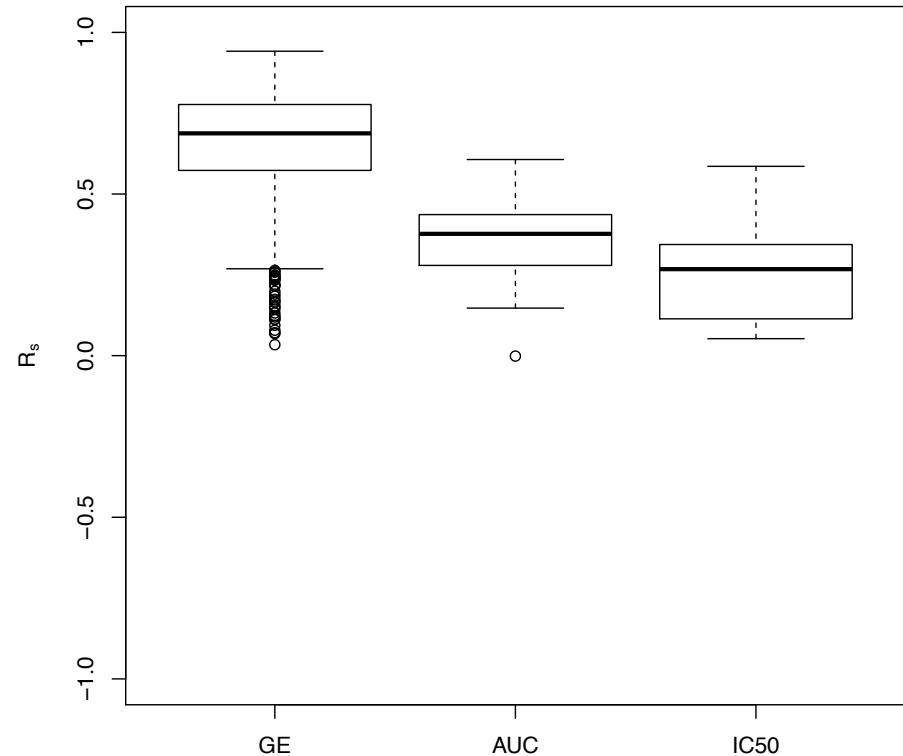
# Concordance of transcriptomic and pharmacologic data

## Spearman correlation

*Between cell lines*

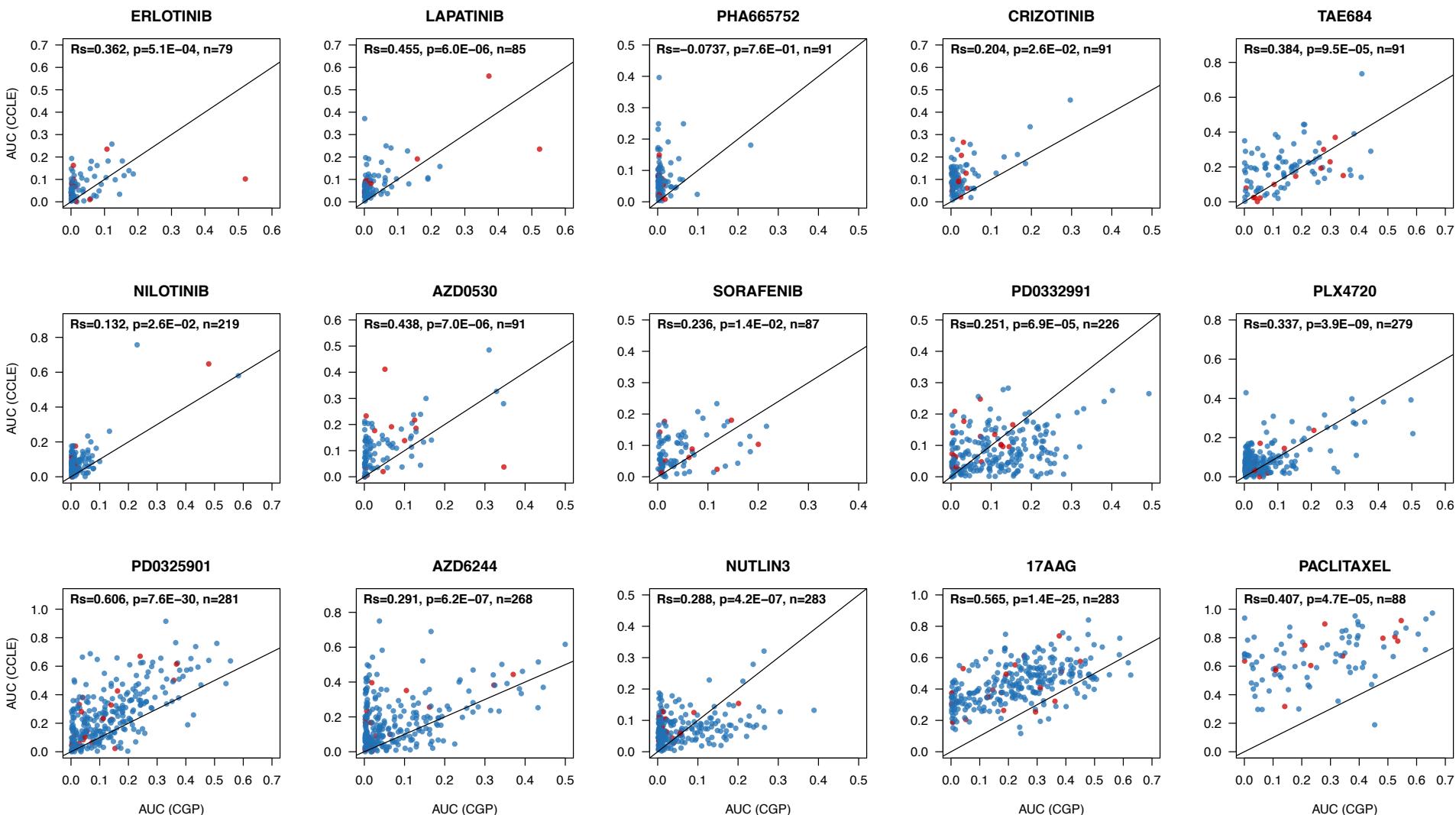


*Across cell lines*



→ Gene expression profiles are significantly more correlated than drug sensitivity measures

# Across cell lines → published AUC



# Consistency of gene-drug associations

- Univariate biomarkers
- Model for gene-drug association:

$$Y = \beta_0 + \beta_i G_i + \beta_t T$$

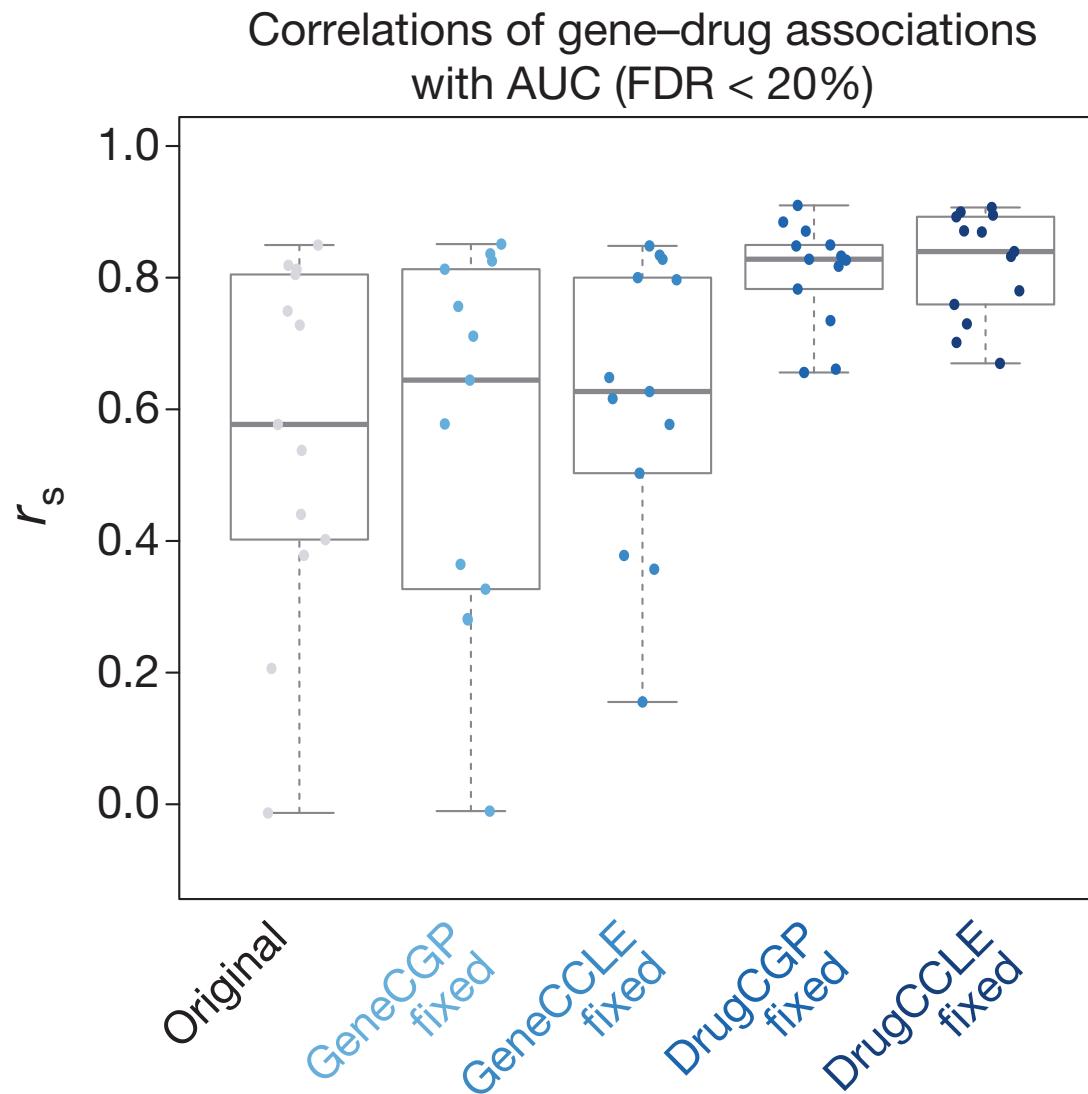
where  $Y$  = drug sensitivity (AUC)

$G_i$  = gene expression of gene  $i$

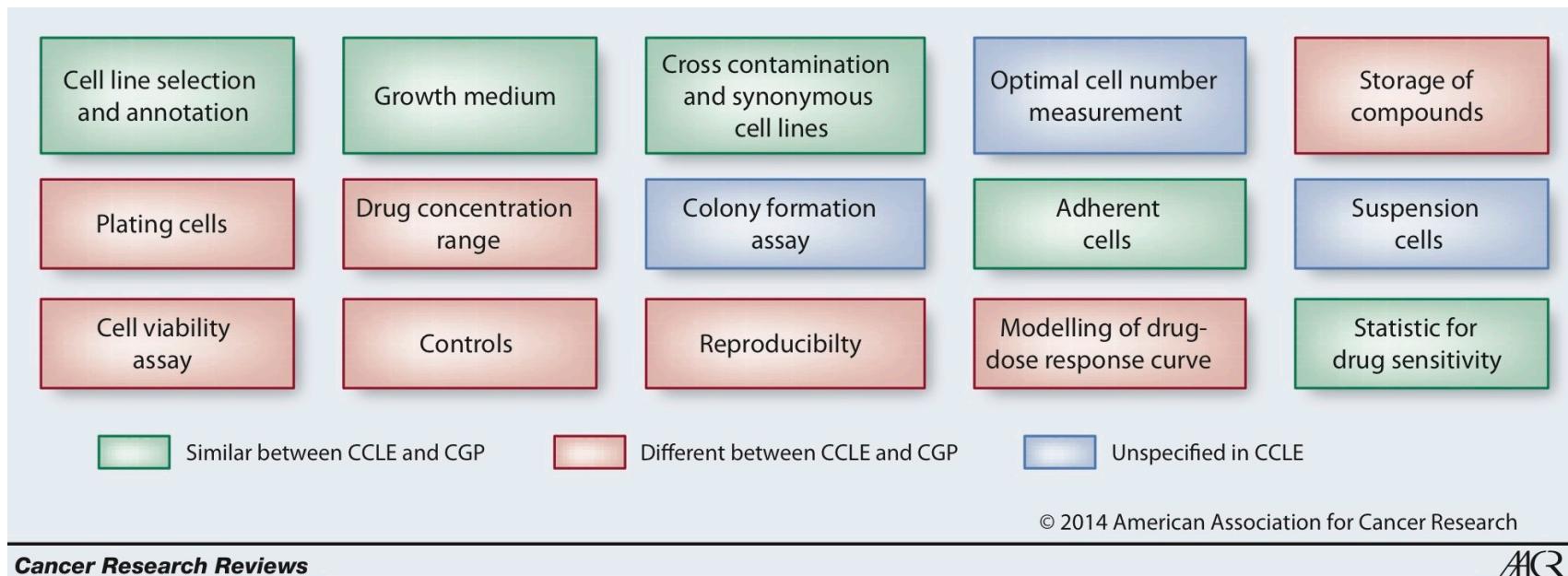
$T$  = tissue type

- Strength and significance of association are given by  $\beta_i$  and its corresponding p-value

# Gene expression or drug sensitivity at fault?



# Heterogeneous experimental protocols



Cancer Research Reviews



Review

Cancer  
Research

## Enhancing Reproducibility in Cancer Drug Screening: How Do We Move Forward? ☰

Christos Hatzis<sup>1,3</sup>, Philippe L. Bedard<sup>10,11</sup>, Nicolai Juul Birkbak<sup>13</sup>, Andrew H. Beck<sup>4</sup>, Hugo J.W.L. Aerts<sup>5,7</sup>, David F. Stern<sup>2,3</sup>, Leming Shi<sup>8,14,15</sup>, Robert Clarke<sup>9</sup>, John Quackenbush<sup>5,6</sup>, and Benjamin Haibe-Kains<sup>10,12</sup>

# Pharmacological assay

- In 2010, GlaxoSmithKline tested
  - 19 compounds
  - on 311 cancer cell lines
- 194 cell lines in common with CGP and CCLE

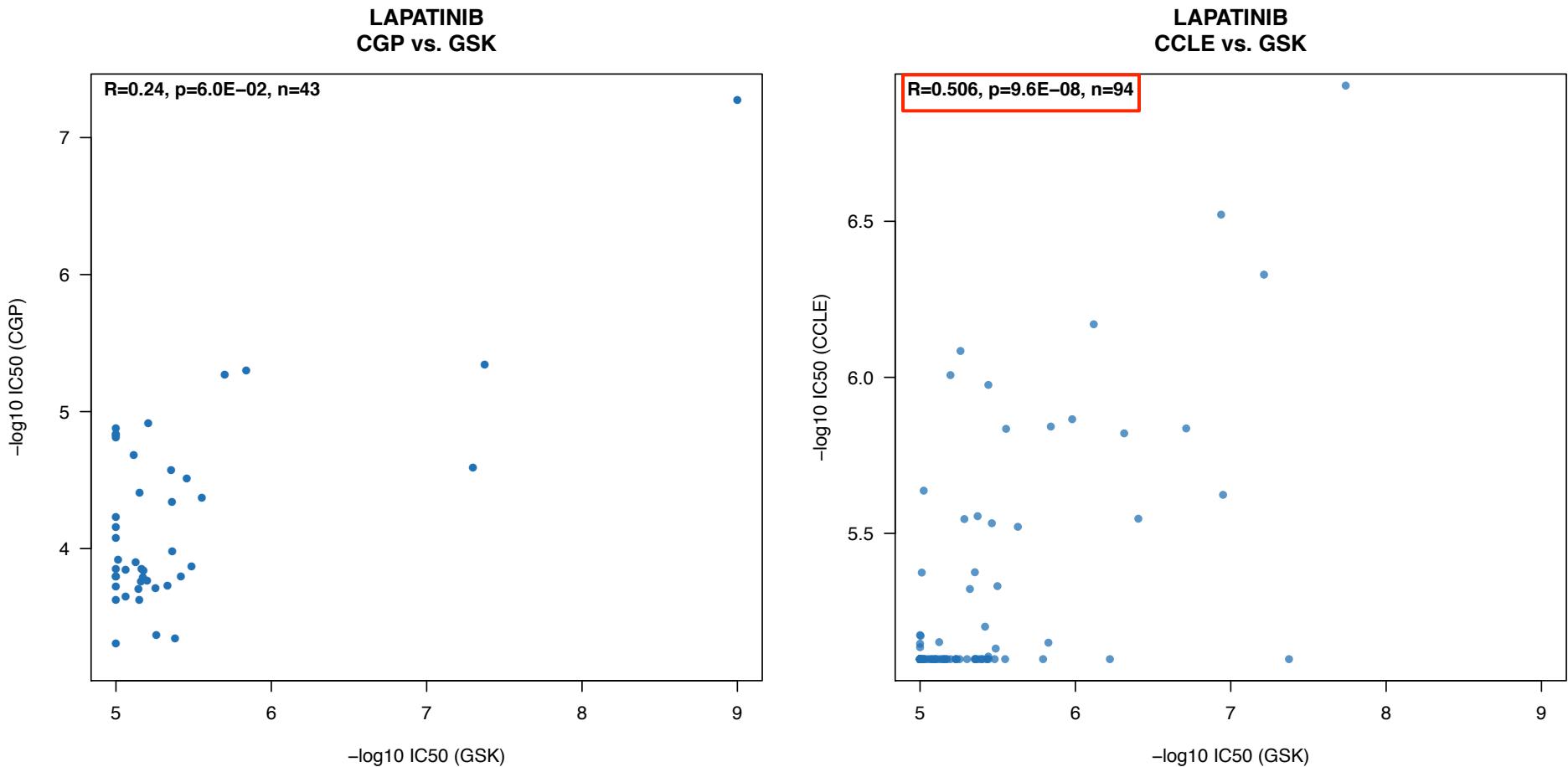


- 2 drugs in common, lapatinib and paclitaxel
- **CCLE and GSK used the same pharmacological assay**  
(Cell Titer Glo luminescence assay, Promega)



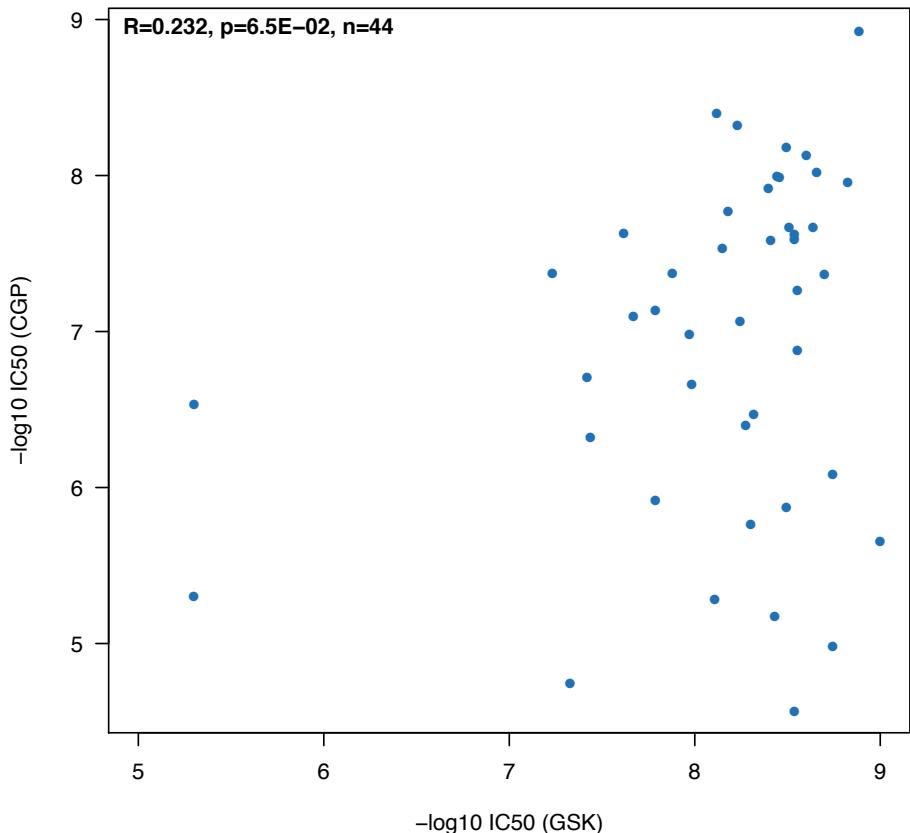
Nehme El-Hachem  
Donald Wang

# Comparison with GSK for lapatinib

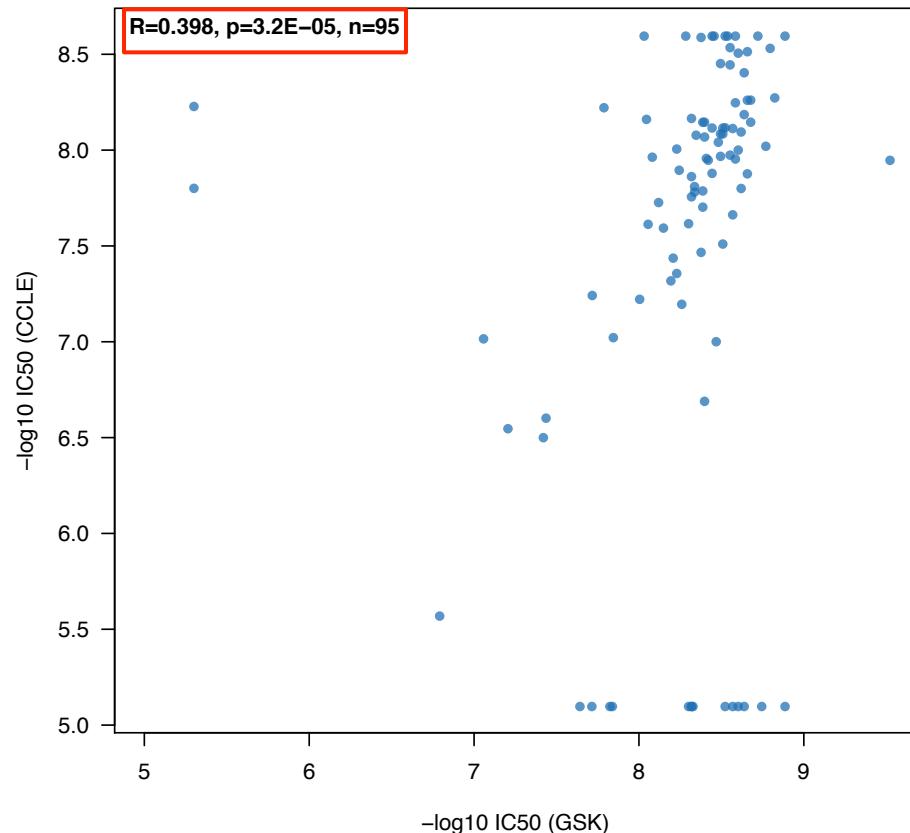


# Comparison with GSK for paclitaxel

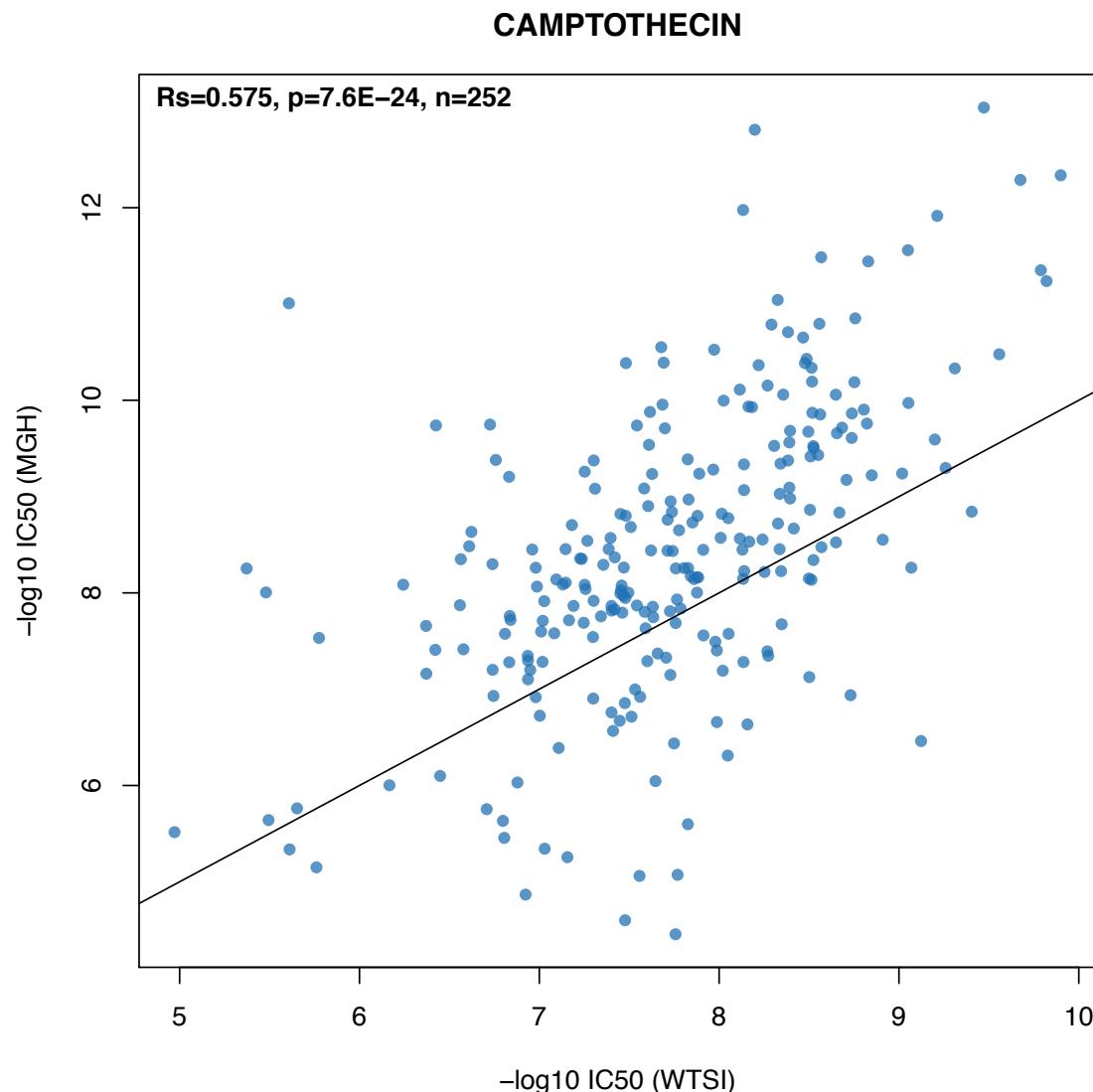
PACLITAXEL  
CGP vs. GSK



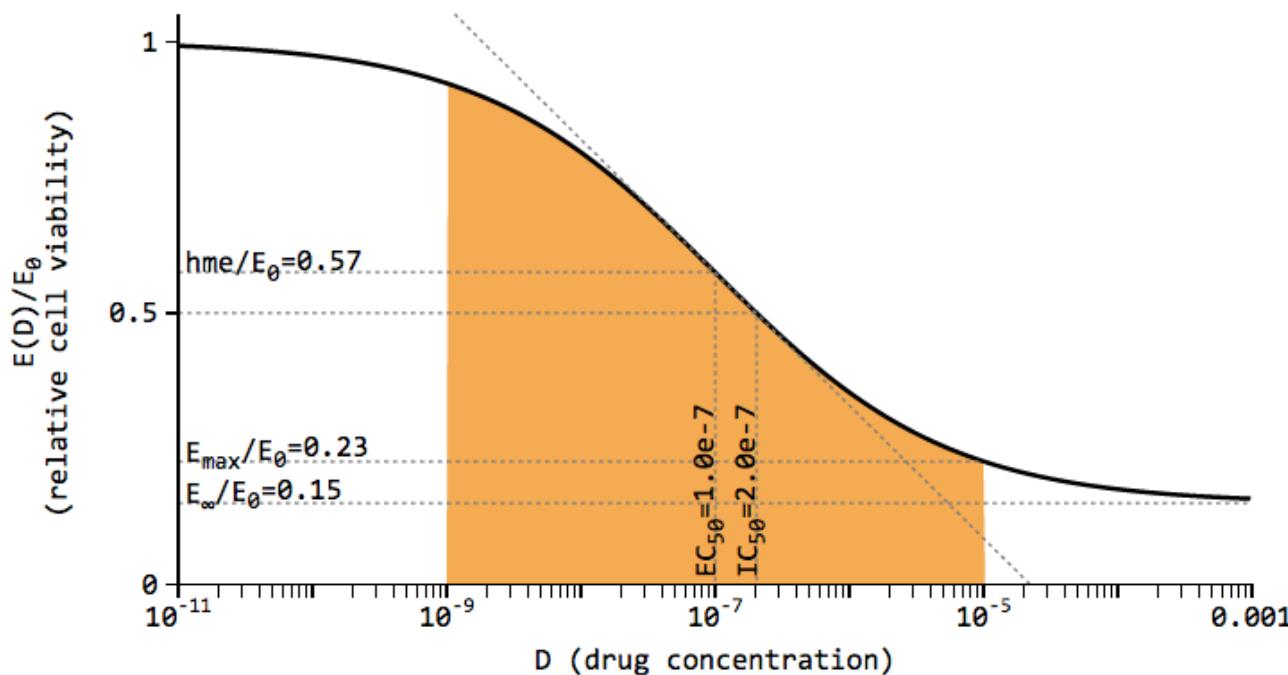
PACLITAXEL  
CCLE vs. GSK



# Identical protocols in different laboratories?



# Complexity of drug dose-response curves



<http://lincs.hms.harvard.edu/explore/10.1038-nchembio.1337/fallahi-sichani-2013/>

# Complexity of drug dose-response curves

- The choice of drug sensitivity summary measure will strongly affect the biomarker discovery process



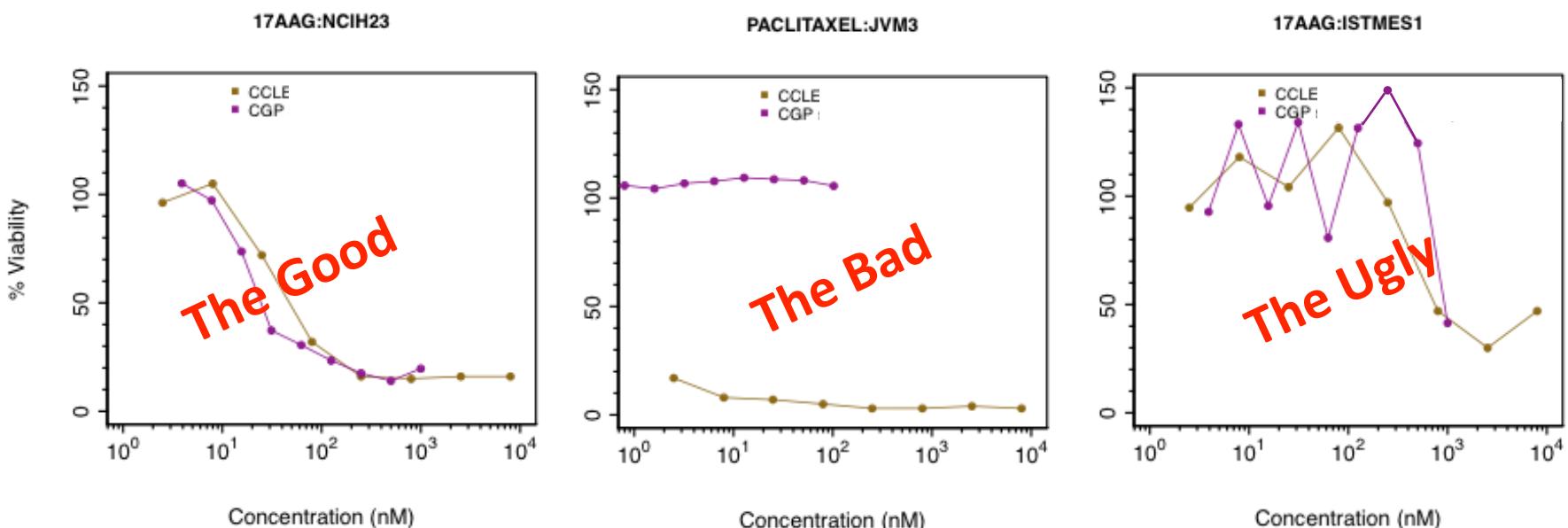
*Pacific Symposium on Biocomputing 2014*

## SYSTEMATIC ASSESSMENT OF ANALYTICAL METHODS FOR DRUG SENSITIVITY PREDICTION FROM CANCER CELL LINE DATA\*

IN SOCK JANG<sup>1</sup>, ELIAS CHAIBUB NETO, JUSTIN GUINNEY, STEPHEN H. FRIEND, ADAM A. MARGOLIN<sup>1</sup>

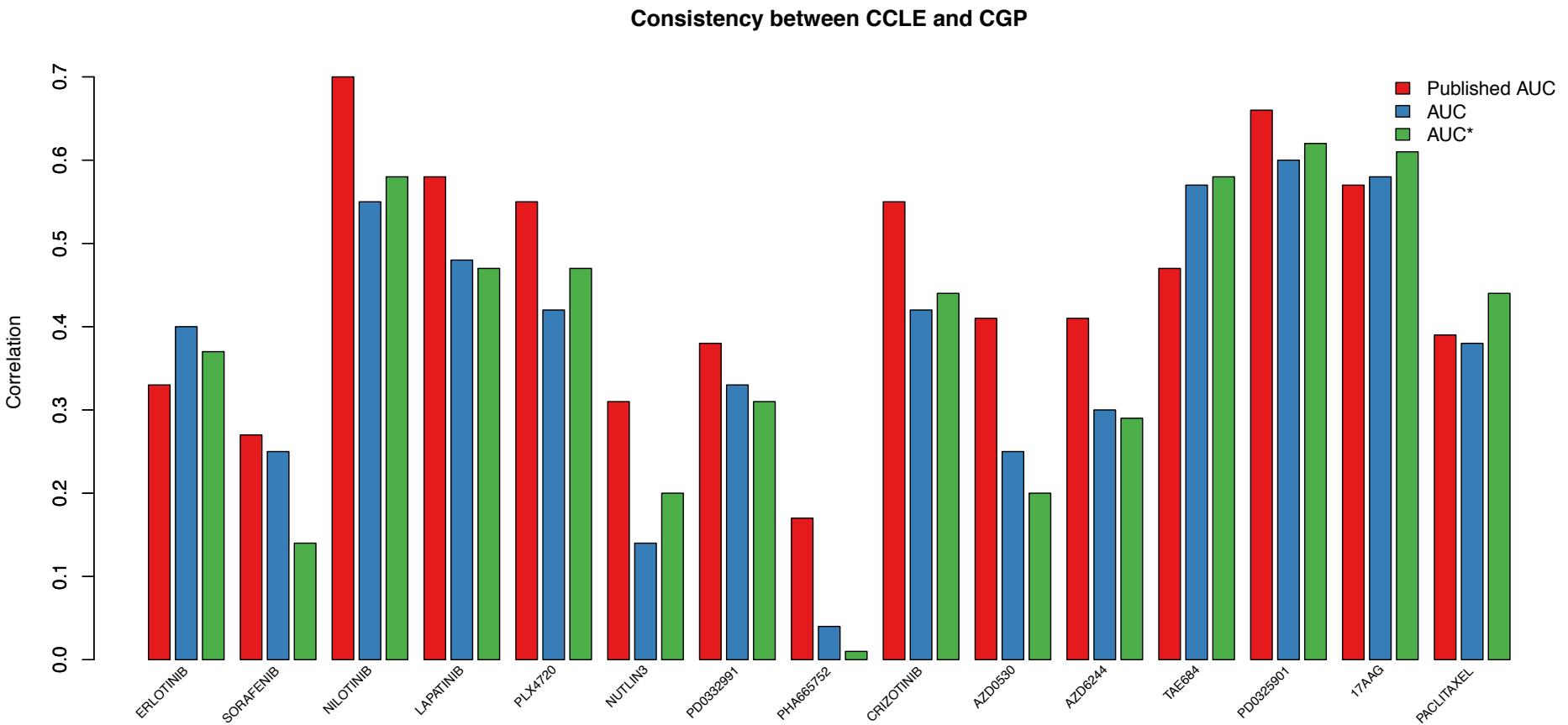
# Summarization of drug dose-response curves

- CGP and CCLE used different statistics for drug sensitivity
- We went back to the raw cell viability data to re-estimate AUC using the standard trapezoidal rule



Zhaleh Safikhani

# Homogeneous summarization does not help



Zhaleh Safikhani

# Conclusions

# Pharmacogenomics

- Pharmacogenomics holds great potential for
  - Discovering new targeted therapies
  - Defining drug mechanisms of action
  - Identifying robust biomarkers of drug response
- In this course we focused on “drug sensitivity” datasets as opposed to other pharmacogenomic datasets such as
  - Drug perturbations, aka Connectivity Map (Lamb et al. Science 2006)
  - Pharmacokinetic (PK) data to determine the fate of substances administered externally to a living organism
  - Pharmacodynamic (PD) data to study effects of drugs on the body and the mechanisms of drug action and the relationship between drug concentration and effect.

# Take home messages

- Prediction of drug response is challenging
- Not much success to date
  - Even in simple models such as cancer cell lines
  - It will only be more difficult *in vivo*!
- Potential solutions:
  - Improve reproducibility of cell-based HTS studies
  - Perform meta-analysis to extract the consistent signal from the data (*PharmacoGx* coming soon; Petr Smirnov)
- This is just the beginning...
  - How to deal with tumor heterogeneity?
  - How to design effective drug combinations?

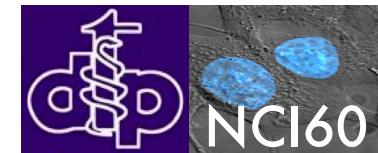
# Mining existing pharmacogenomic data



Genentech



S<sup>↑</sup>2C™



Courtesy of Mathieu Lupien

# Acknowledgements

## From the lab

- Zhaleh Safikhani
- Petr Smirnov
- Rene Quevedo
- Nehme El-Hachem
- Donald Wang
- Simon Papillon-Cavanagh
- Nicolas de Jay
- Adrian She
- Deena Gendoo



Yale University  
School of Medicine

- Christos Hatzis



- Nicolai Juul Birkbak



- Jacques Archambault



- Hugo Aerts
- John Quackenbush



- Andrew Beck



- Leming Shi

- Carl Virtanen

The Princess Margaret  
Cancer Foundation UHN

*Thank you for your attention!*