

# *PharmacoGx: Data Sharing and Research Reproducibility in Pharmacogenomics*

**Benjamin Haibe-Kains**

Princess Margaret Cancer Centre  
University Health Network  
University of Toronto  
Ontario Institute of Cancer Research

***2 open postdoc positions:  
Re radiomics and single-cell RNA-seq***

June 25, 2016

# Reproducibility crisis

- ▷ Reproducibility in biomedical sciences has attracted a lot of attention in the last 10 years

**Reproducibility of published microarray gene expression**  
**Signal in Noise: Believe it or not: how much can we**  
**Reported Repro**  
**Serum Proteomi**  
**Ovarian Cancer** rely on published data on potential drug targets?

Jo  
M  
M

*Keith A. Baggerly, Florian Prinz, Thomas Schlange and Khusru Asadullah*

*S. Morris, Sarah R. Edmonson,*

*Kevin R. Coombes*

## Research Findings

# Raise standards for preclinical cancer research

**ACTIVITY FROM CELL LINES:  
STATISTICS AND REPRODUCIBLE  
HIGH-THROUGHPUT BIOLOGY**

C. Glenn Begley and Lee M. Ellis propose how methods, publications and incentives must change if patients are to benefit.

Y<sup>1</sup> AND KEVIN R. COOMBES<sup>2</sup>

# Why data and code sharing?

- ▷ Data are precious due to limited
  - Amount of samples
  - Resources
  - Budget

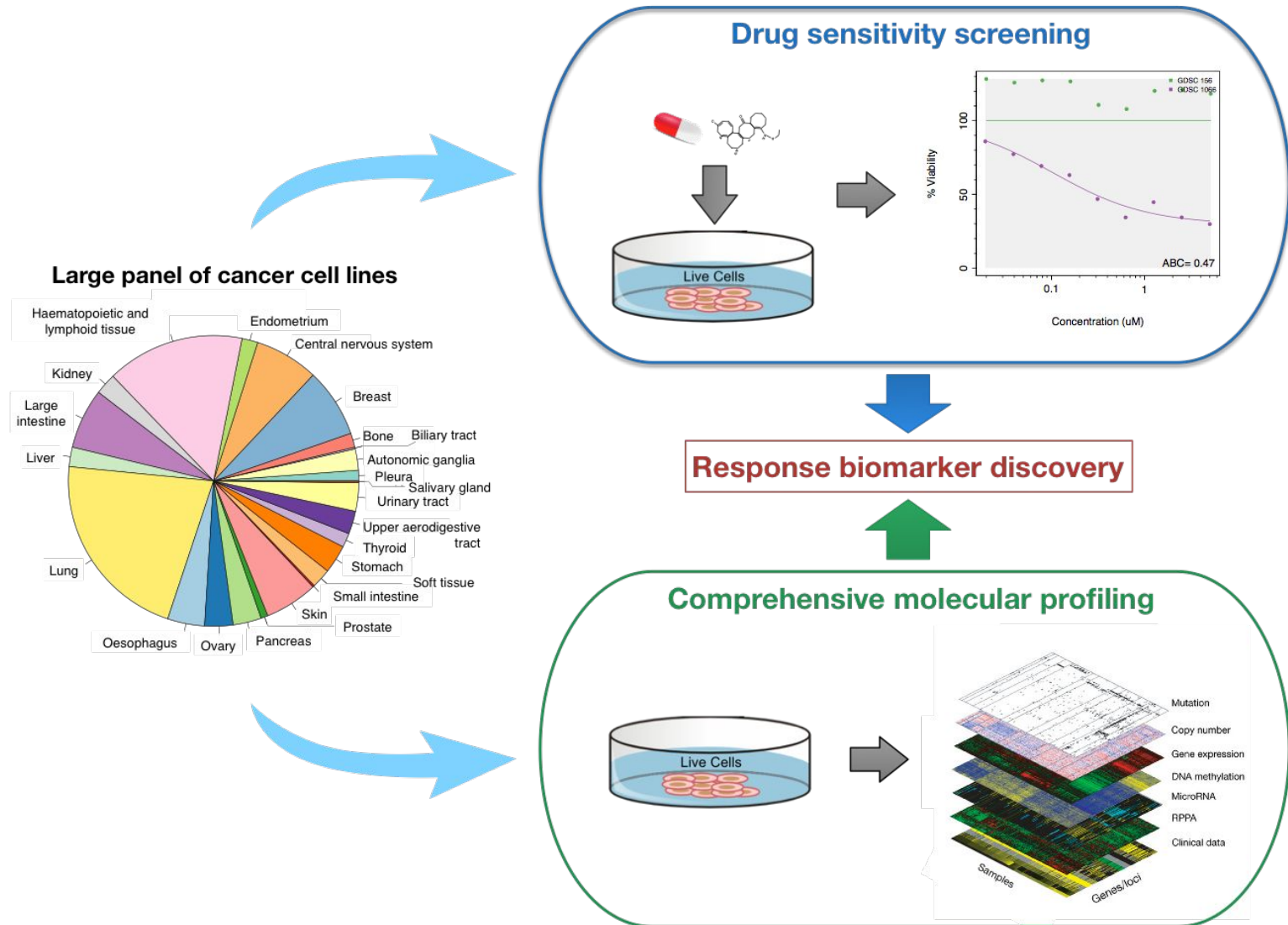


*“Anyone who believes in indefinite growth in anything physical, on a physically finite planet, is either mad or an economist.”*

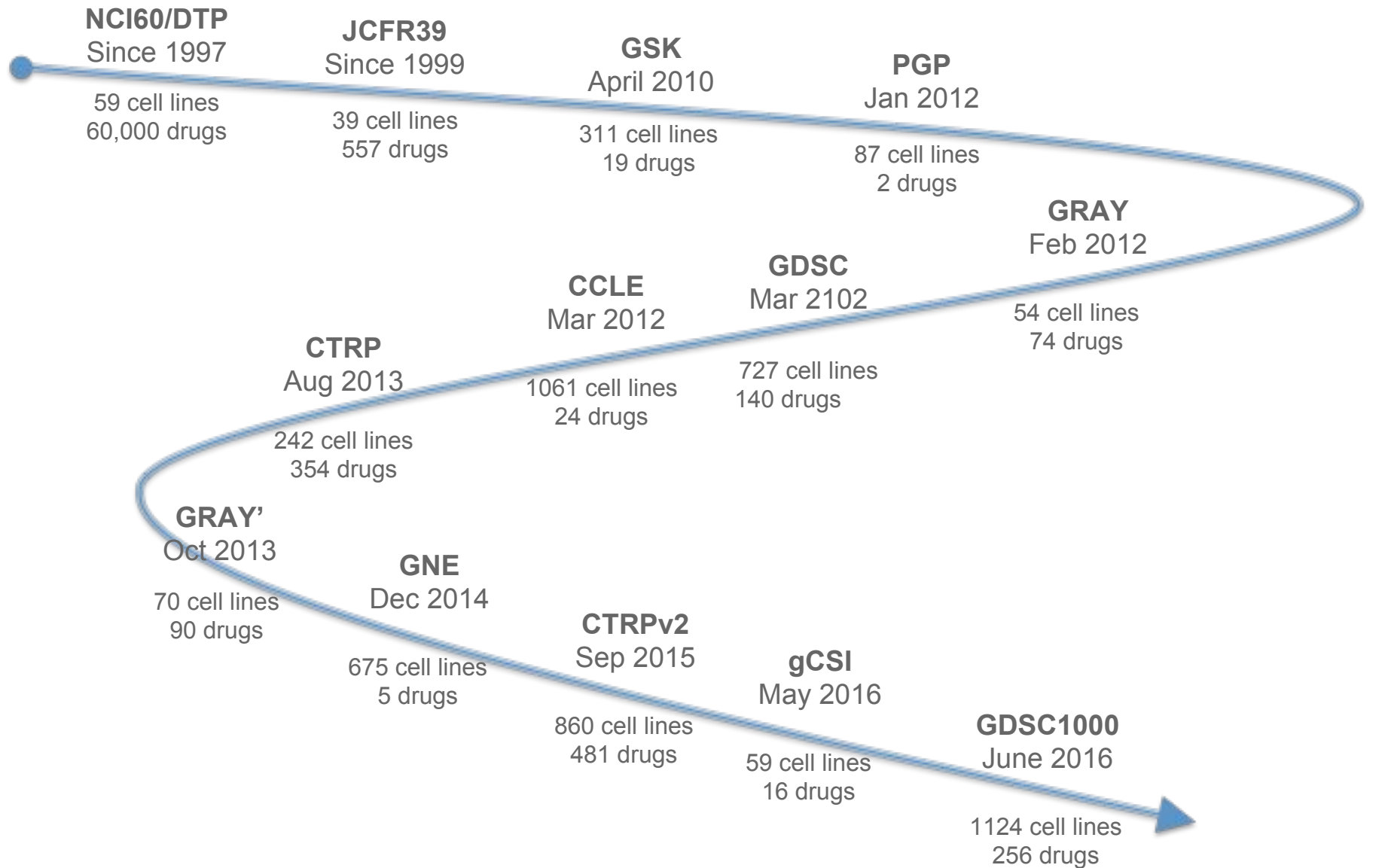
— Kenneth E. Boulding

- ▷ Benefits of sharing data and code
  - Replicability
  - Reproducibility
  - Reusability
  - Post-publication peer review

# High-throughput *in vitro* drug screening



# Long history of data sharing in pharmacogenomics



*More to come...*

# Predictors trained on one dataset hardly validate on an independent set

JAMIA

Comparison and validation of genomic predictors for anticancer drug sensitivity

Simon Papillon-Cavanagh,<sup>1</sup> Nicolas De Jay, Gianluca Bontempi,<sup>2</sup> Hugo J W L Aerts,<sup>3,4</sup>

Papillon-Cavanagh S, et al. *J A*

*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>



Dong et al. *BMC Cancer* (2015) 15:489  
DOI 10.1186/s12885-015-1492-6

RESEARCH ARTICLE

Open Access

Anticancer drug sensitivity prediction in cell lines from baseline gene expression through recursive feature selection



Zuoli Dong<sup>1†</sup>, Naiqian Zhang<sup>1†</sup>, Chun Li<sup>2</sup>, Haiyun Wang<sup>3</sup>, Yun Fang<sup>1</sup>, Jun Wang<sup>1\*</sup> and Xiaoqi Zheng<sup>1\*</sup>

Bioinformatics

OXFORD JOURNALS

**Improved large-scale prediction of growth inhibition patterns using the NCI60 cancer cell line panel**

Isidro Cortés-Ciriano<sup>1</sup>, Gerard J.P. van Westen<sup>2</sup>, Guillaume Bouvier<sup>1</sup>, Michael Nilges<sup>1</sup>, John P. Overington<sup>2</sup>, Andreas Bender<sup>3\*</sup> and Thérèse E. Malliavin<sup>1\*</sup>



# Comparative studies

2013

**nature**

ANALYSIS RESEARCH

## Inconsistency in large pharmacogenomic studies

Benjamin Haibe-Kains<sup>1,2</sup>, Nehme El-Hachem<sup>1</sup>, Nicolai Juul Birkbak<sup>3</sup>, Andrew C. Jin<sup>4</sup>, Andrew H. Beck<sup>4\*</sup>, Hugo J. W. L. Aerts<sup>5,6,7\*</sup> & John Quackenbush<sup>5,8\*</sup>

2015

**nature**

## Revisiting inconsistency in large pharmacogenomic studies

**Pharm**  
**two c**

The Cancer Cell L

Zhaleh Safikhani, Mark Freeman, Petr Smirnov, Nehme El-Hachem, Adrian She, Rene Quevedo, Anna Goldenberg, Nicolai Juul Birkbak, Christos Hatzis, Leming Shi, Andrew H Beck, Hugo JW L Aerts, John Quackenbush, Benjamin Haibe-Kains  
doi: <http://dx.doi.org/10.1101/026153>



**bioRxiv**  
beta  
THE PREPRINT SERVER FOR BIOLOGY

2016

**nature**

**Reprod**

**proc**

Peter M. J  
Jeff Settle

## Assessment of pharmacogenomic agreement [version 1; referees: 1 approved]

Zhaleh Safikhani<sup>1,2</sup>, Nehme El-Hachem<sup>3</sup>, Rene Quevedo<sup>1,2</sup>, Petr Smirnov<sup>1</sup>, Anna Goldenberg<sup>4,5</sup>, Nicolai Juul Birkbak<sup>6</sup>, Christopher Mason<sup>7-9</sup>, Christos Hatzis<sup>10,11</sup>, Leming Shi<sup>12,13</sup>, Hugo JW L Aerts<sup>14,15</sup>, John Quackenbush<sup>14,16</sup>,  Benjamin Haibe-Kains<sup>1,2,5</sup>

## Integrating heterogeneous drug sensitivity data from cancer pharmacogenomic studies

Nikita Pozdeyev<sup>1</sup>, Minjae Yoo<sup>1</sup>, Ryan Mackie<sup>1</sup>, Rebecca E. Schweppe<sup>1</sup>, Aik Choon Tan<sup>1,\*</sup>, Bryan R. Haugen<sup>1,\*</sup>

# Challenges in pharmacogenomic analyses

- ▶ Cell line and drug identifiers are not standardized
  - Difficult to assess overlap between datasets
- ▶ Heterogeneous experimental designs
- ▶ No consensus on data formats



**PharmacGx: an R package for analysis of large pharmacogenomic datasets**  
doi: 10.1093/bioinformatics/btv723

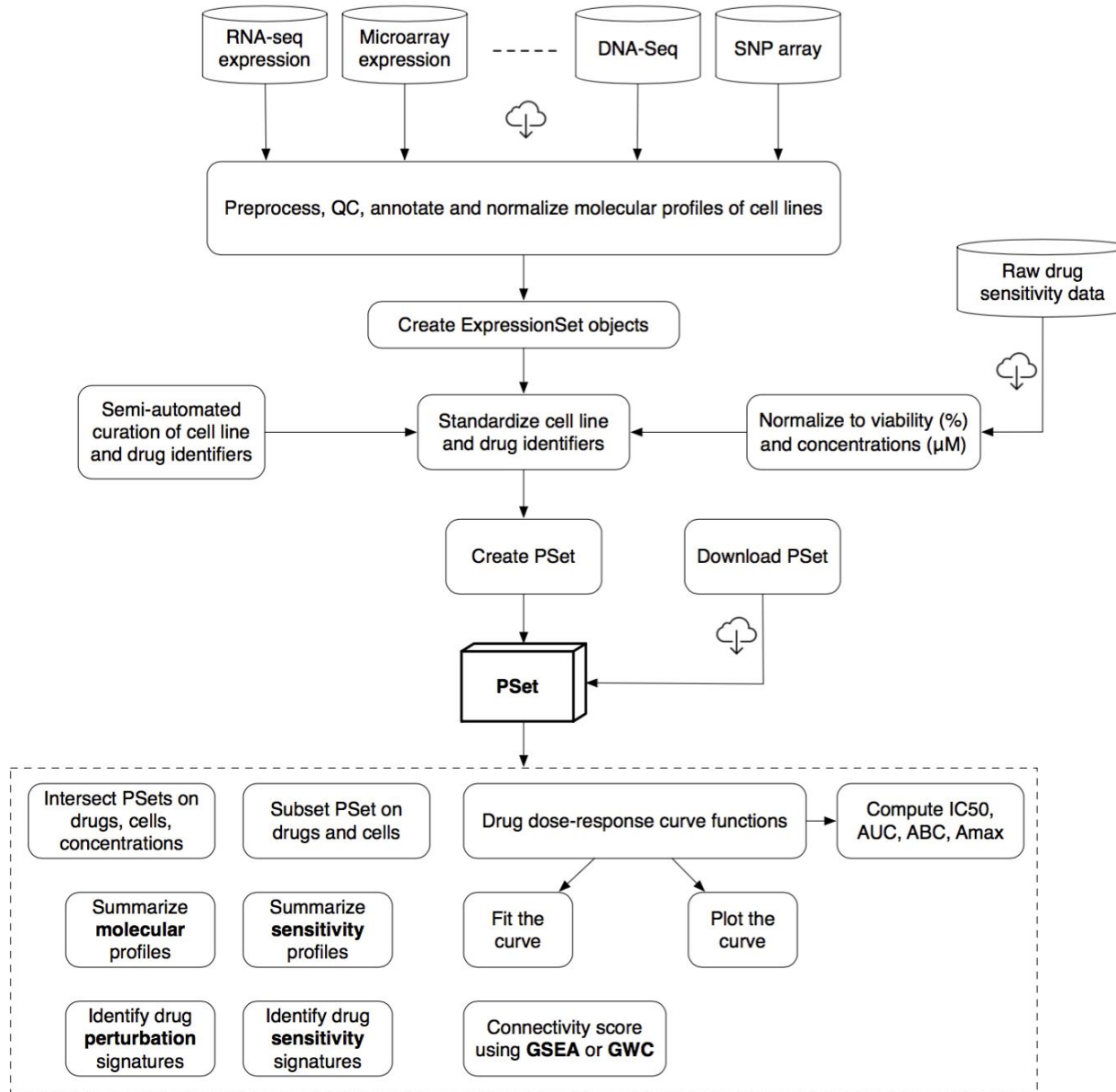
Petr Smirnov<sup>1,†</sup>, Zhaleh Safikhani<sup>1,2,†</sup>, Nehme El-Hachem<sup>3</sup>, Dong Wang<sup>1</sup>,  
Adrian She<sup>1</sup>, Catharina Olsen<sup>1,4,5</sup>, Mark Freeman<sup>1</sup>, Heather Selby<sup>6,7</sup>,  
Deena MA Gendoo<sup>1,2</sup>, Patrick Grossmann<sup>6</sup>, Andrew H. Beck<sup>8</sup>,  
Hugo JWL Aerts<sup>6</sup>, Mathieu Lupien<sup>1,2,9</sup>, Anna Goldenberg<sup>10,11</sup> and  
Benjamin Haibe-Kains<sup>1,2,\*</sup>

<https://github.com/haibe-kains/PharmacGx>

(available on CRAN, under review for BioC)



# PharmacoGx in a nutshell



# PharmacoSet S4 class

## @ annotation:

- \$ name: Acronym of the pharmacogenomic dataset.
- \$ dateCreated: When the object was created.
- \$ sessionInfo: Software environment used to create the object.
- \$ call: Set of parameters used to create the object.

## @ datasetType: Either 'sensitivity', 'perturbation', or 'both'

## @ cell: data frame annotating all cell lines investigated in the study.

## @ drug: data frame annotating all the drugs investigated in the study.

## @ sensitivity:

- \$ n: Number of experiments for each cell line treated with a given drug
- \$ info: Metadata for each pharmacological experiment.
- \$ raw: All cell viability measurements at each drug concentration from the drug dose-response curves.
- \$ phenotype: Drug sensitivity values summarizing each dose-response curve (IC<sub>50</sub>, AUC, etc.)

## @ perturbation:

- \$ n: Number of experiments for each cell line perturbed by a given drug, for each molecular data type
- \$ info: 'The metadata for the perturbation experiments is available for each molecular type by calling the appropriate info function'

@ molecularProfiles: List of ExpressionSet objects containing the molecular profiles of the cell lines, such as mutations, gene expressions, or copy number variations.

→ MultiAssayExperiment

# PharmacoGx enables meta-analysis

- ▶ Cellosaurus to uniquely identify

**Datasets available today:**

[web.expasy.org/cellosaurus/](http://web.expasy.org/cellosaurus/)

**CMAP, GDSC, CCLE and gCSI**

- ▶ Drugs annotated with PubChem ID, InChiKey and SMILES

- Exact and fuzzy matching based on structure similarity

**In the oven:**

- ▶ Ensembl annotations for omics profiles

- ▶ Functional drug words, GSK, GNE,

summarize pharmacogenomic studies

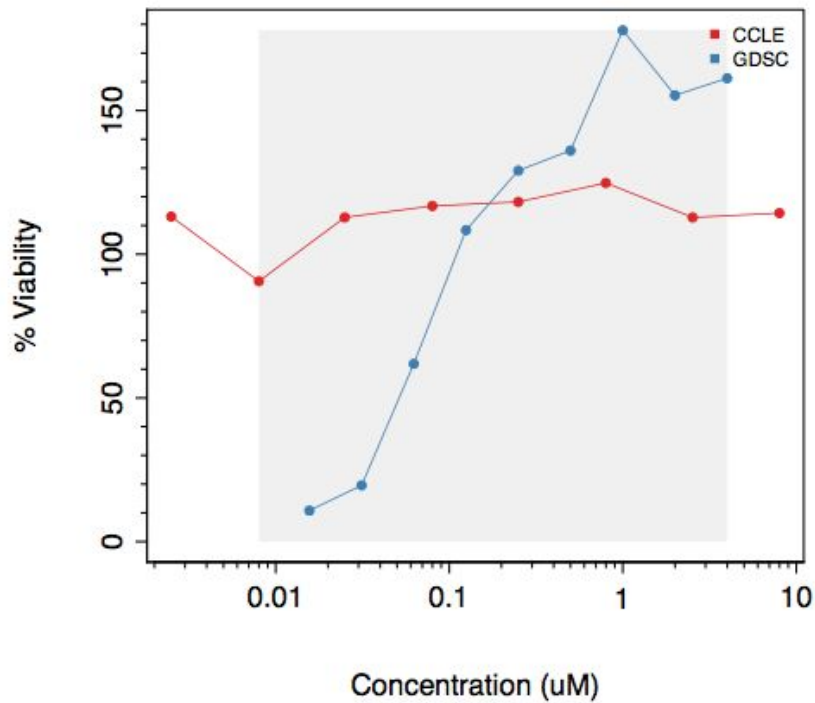
**CTRPv2, GRAY**

- *DownloadPSet()*
- *IntersectPSets()*
- *SubsetTo()*
- *summarize\*()*

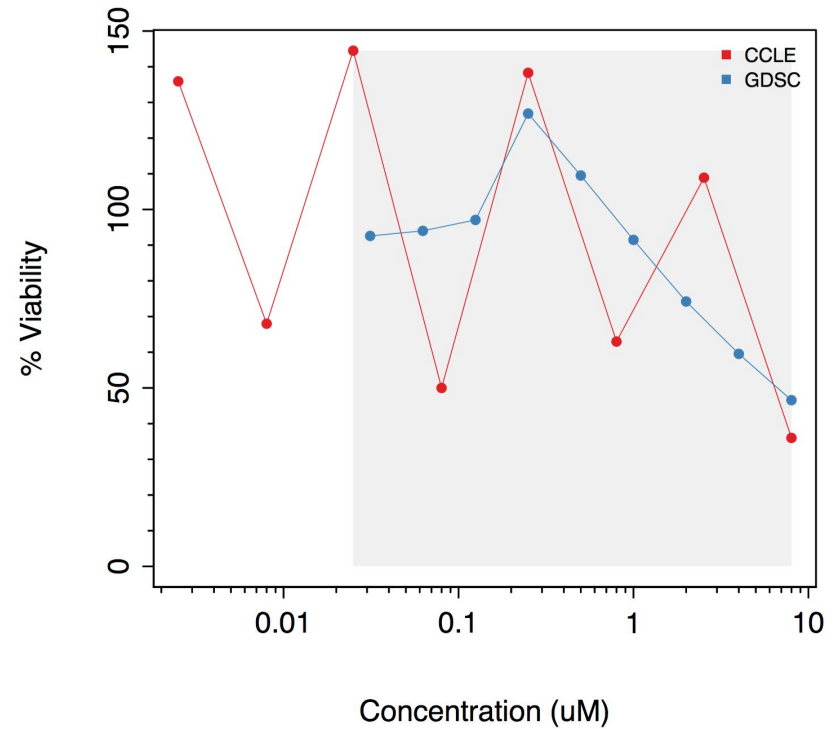


# Filtering of noisy dose-response curves

PD-0332991:HCC70

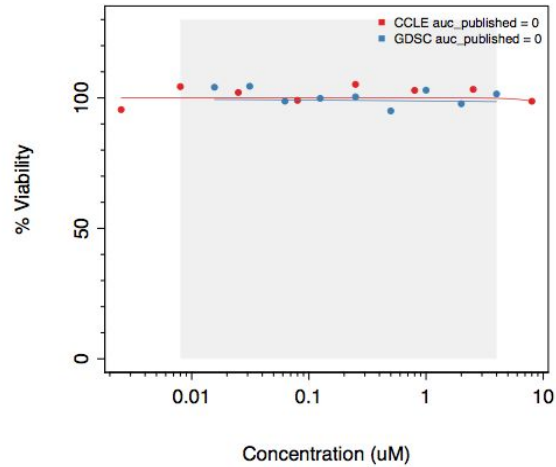


Nutlin-3:LS-513

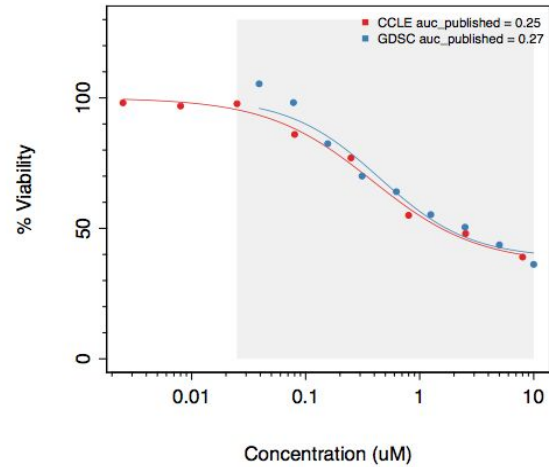


# Fitting of drug dose-response curves

AZD6244:COLO-320-HSR



PLX4720:HT-29



*Highly consistent*

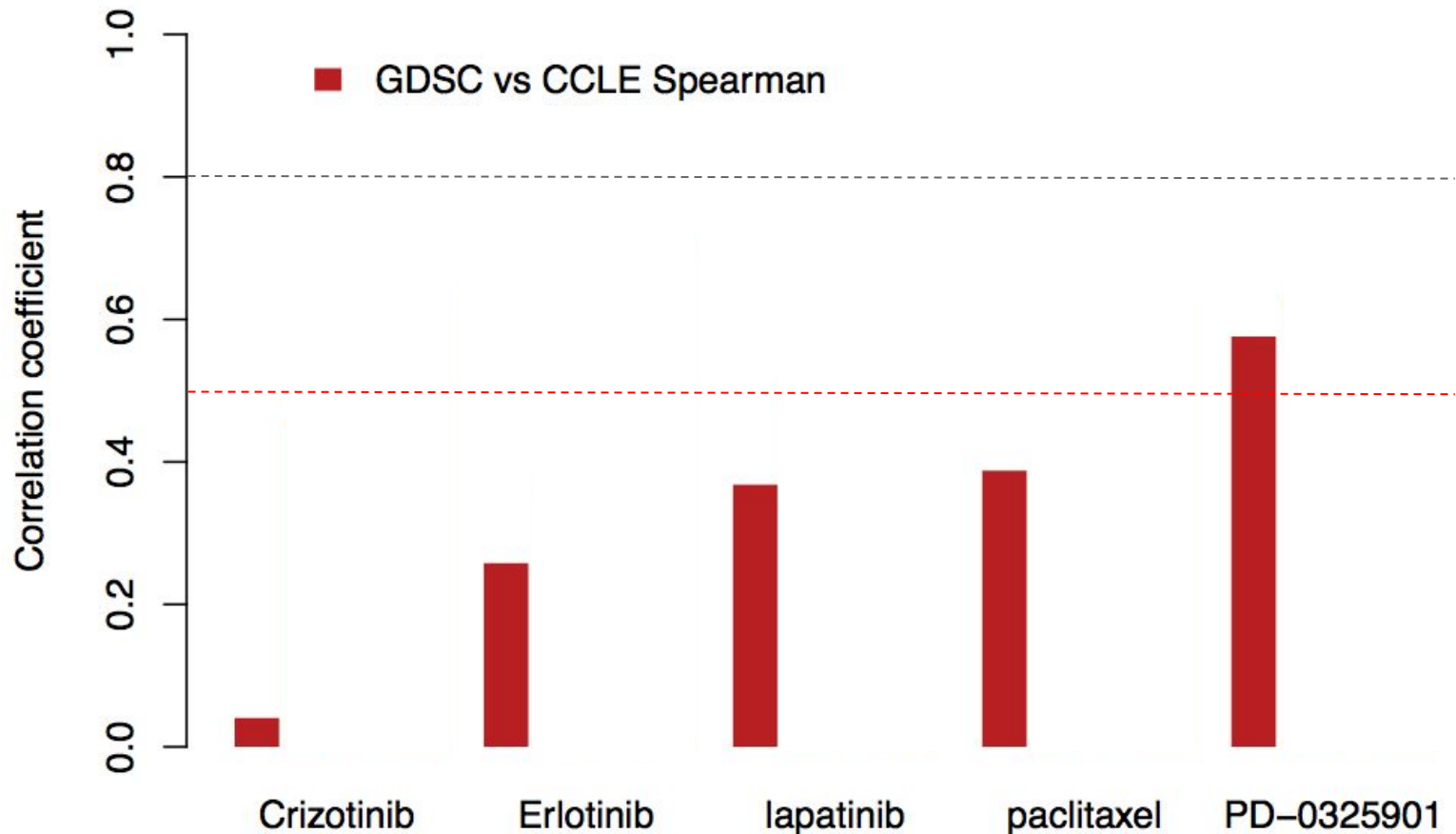


# Correlations of drug sensitivity data

2013 Inconsistency in large pharmacogenomics studies

2015 Revisiting inconsistency in large pharmacogenomic studies  
Pharmacogenomic agreement between two cancer cell line data sets

2016 Reproducible pharmacogenomic profiling of cancer cell line panels

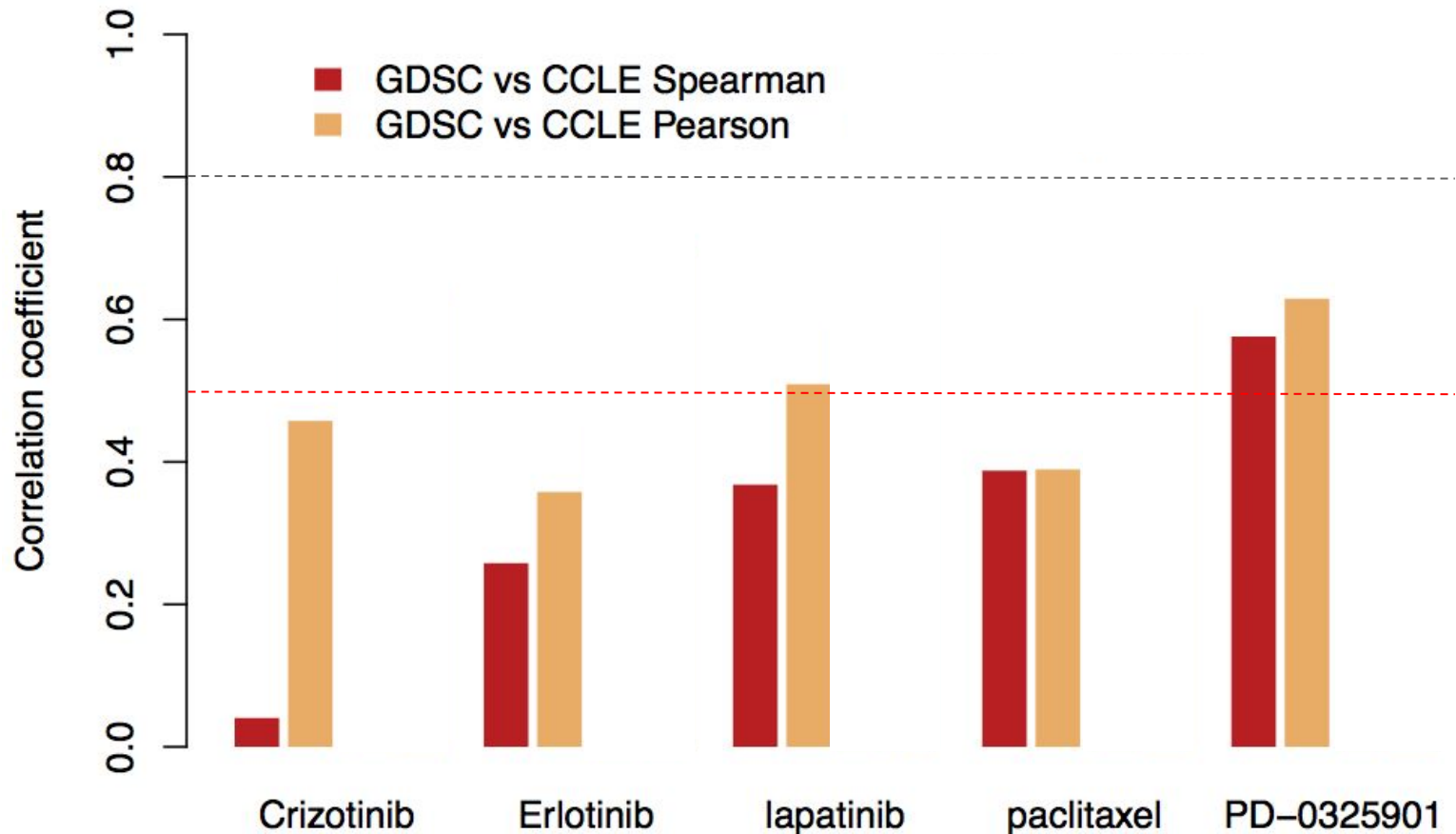


# Correlations of drug sensitivity data

2013 Inconsistency in large pharmacogenomics studies

**2015 Revisiting inconsistency in large pharmacogenomic studies**  
**Pharmacogenomic agreement between two cancer cell line data sets**

2016 Reproducible pharmacogenomic profiling of cancer cell line panels

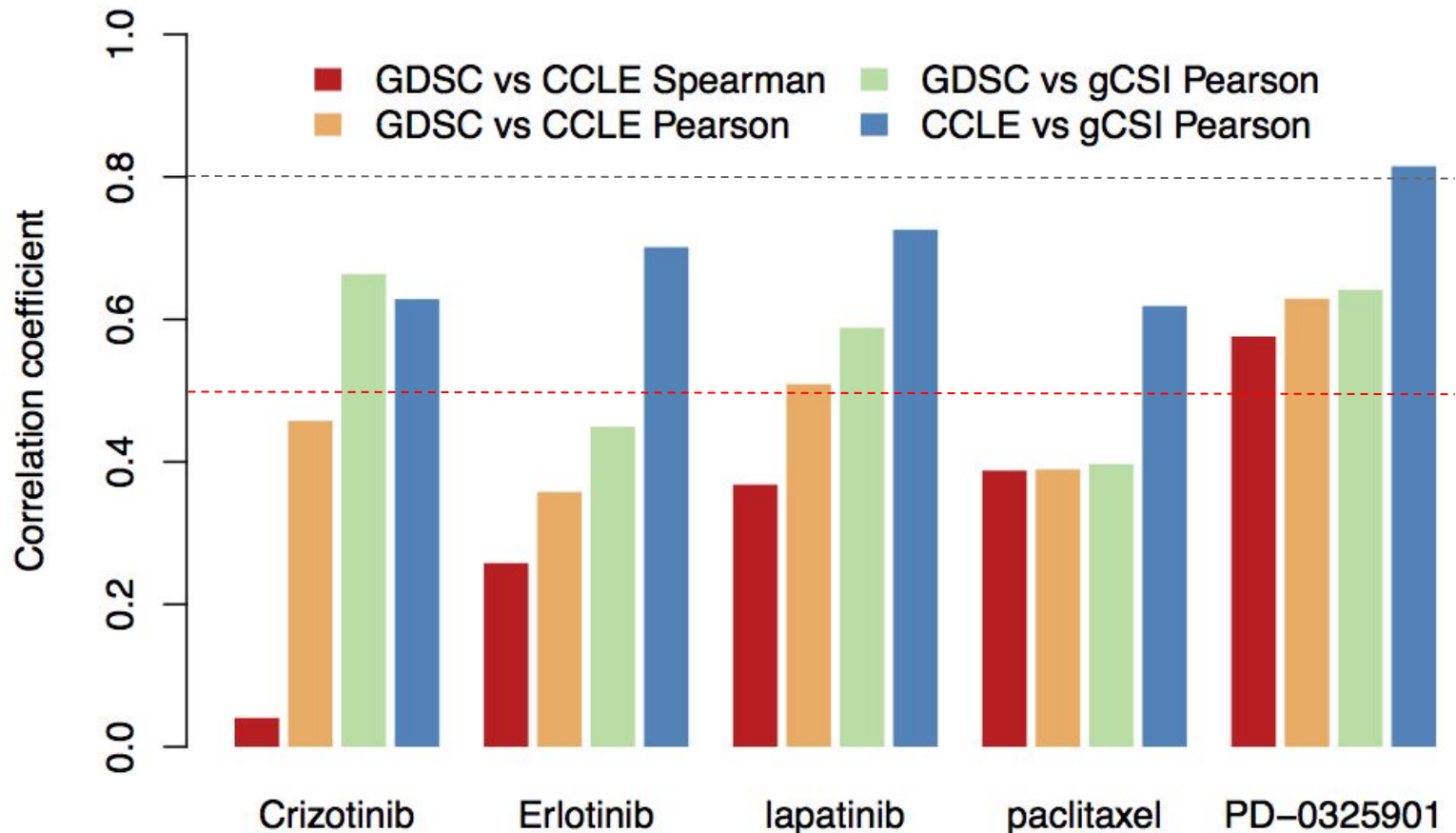


# Correlations of drug sensitivity data

2013 Inconsistency in large pharmacogenomics studies

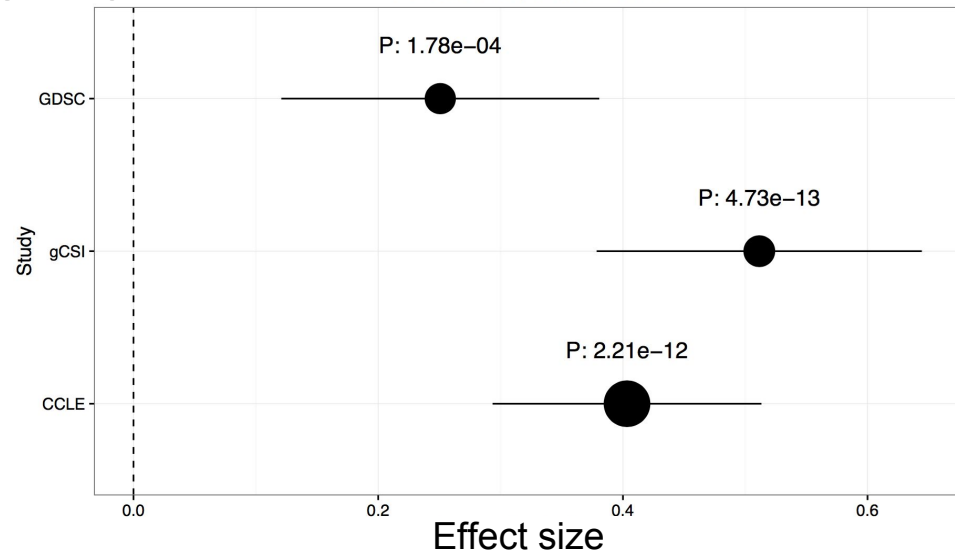
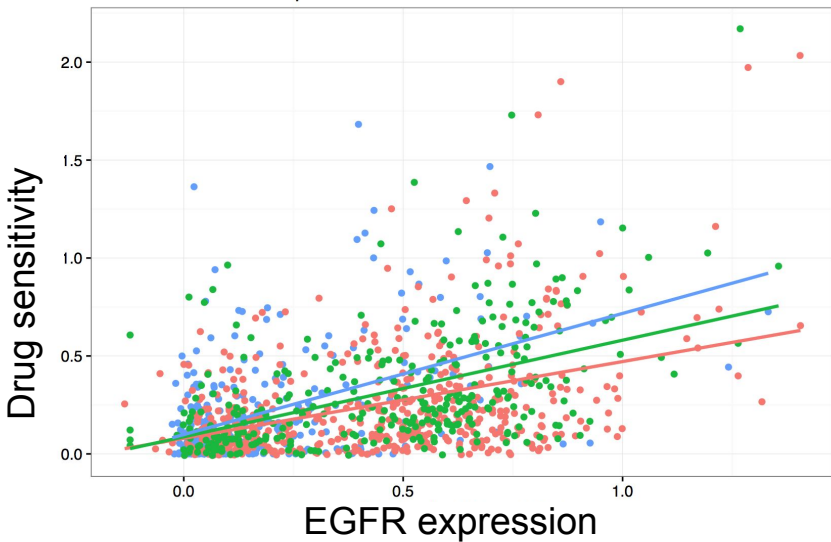
2015 Revisiting inconsistency in large pharmacogenomic studies  
Pharmacogenomic agreement between two cancer cell line data sets

2016 **Reproducible pharmacogenomic profiling of cancer cell line panels**

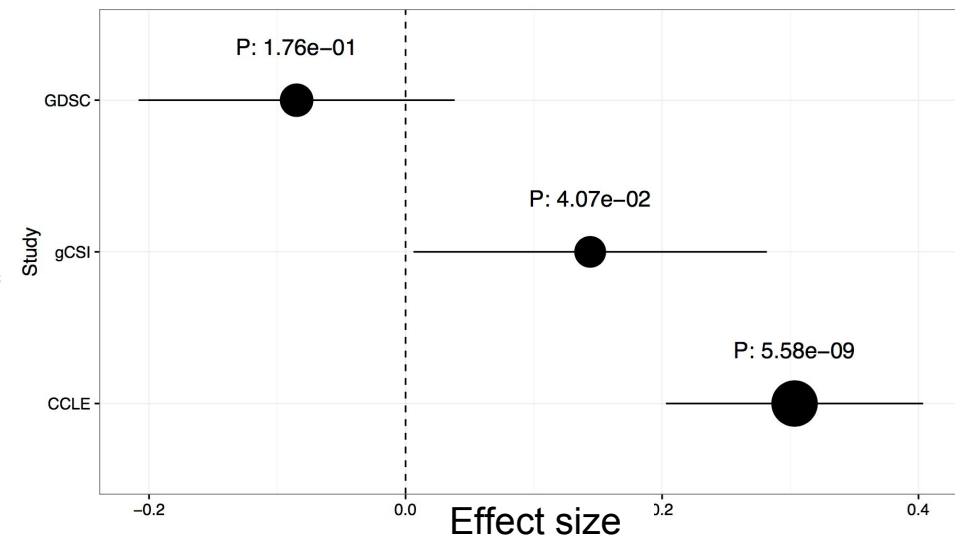
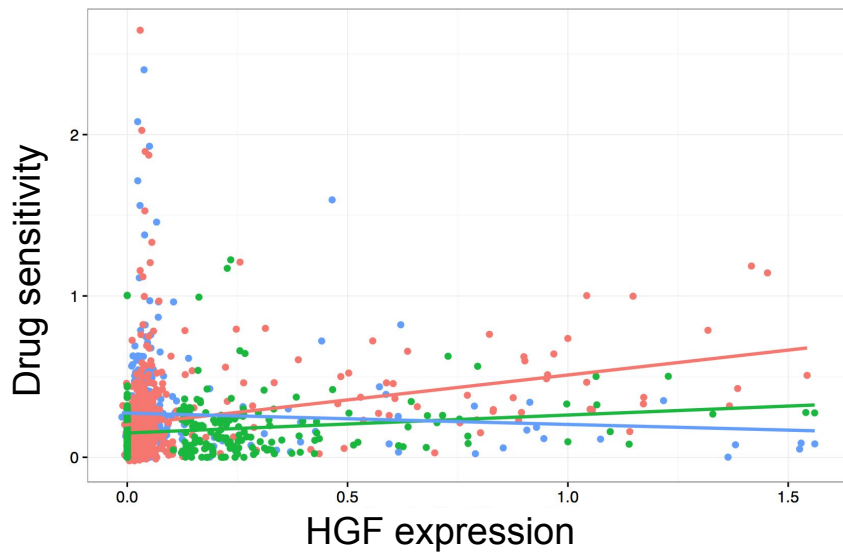


# Robust biomarker discovery

## Erlotinib



## Crizotinib



# Conclusions

- ▷ Pharmacogenomics is a hot field, new datasets and new players everyday
  - You can even stay in the game after pissing off the major league :-)
- ▷ Great need for standardization
  - Experimental protocols
  - Data processing
  - Annotations
- ▷ **PharmacoGx** provides a unified platform for meta-analysis of pharmacogenomic studies

*Our curation is far from perfect, we need your feedback to make it better!*



# Future directions

- ▷ **MultiAssayExperiment (MAE)** to replace the list of ExpressionSet objects and better integrate diverse molecular profiles -- *Workshop session 3*
- ▷ **PharmacoDb**: Companion web-application to facilitate exploration of the large compendium of published pharmacogenomics datasets
- ▷ Development of statistical/machine learning methods to jointly analyze heterogeneous pharmacogenomics datasets
- ▷ Extension to drug combinations (AstraZeneca-Sanger DREAM Challenge)

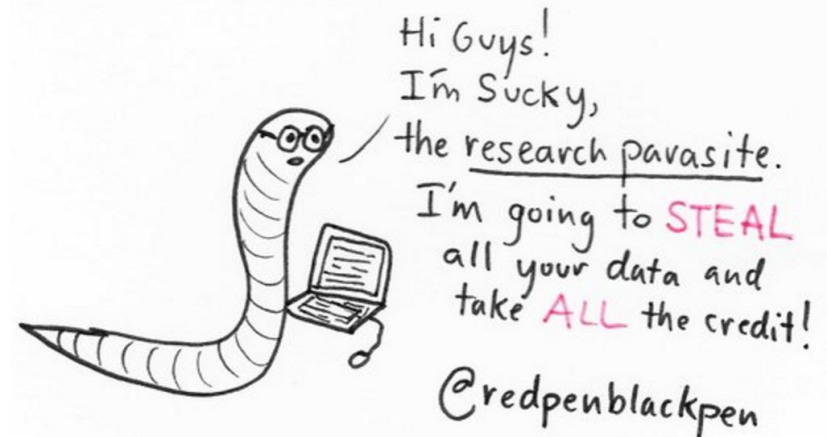
# PharmacoGx can be safely used by

Data vultures



And research parasites

Data vampires



**#IAmAResearchParasite**

# Research parasites



## Data Sharing

Dan L. Longo, M.D., and Jeffrey M. Drazen, M.D.

January 2016

Scientists?

[...] concern held by some is that a new class of person will emerge — people who had nothing to do with the design and execution of the study but who use another group's data for their own ends, possibly to reduce the research productivity planned by the original investigators, or even use the data to try to disprove what the original investigators had posited. There is concern among some front-line researchers that the system will be taken over by what some researchers have characterized as “research parasites.”

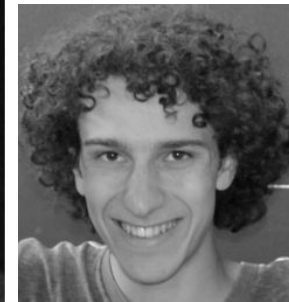
Doing Science?

# Acknowledgements

## BHK lab

*Princess Margaret Cancer Centre*

- ▷ **Zhaleh Safikhani**
- ▷ **Petr Smirnov**
- ▷ **Nehme El-Hachem**
- ▷ **Mark Freeman**
- ▷ **Ali Madani**



## Collaborators

- ▷ John Quackenbush
- ▷ Christos Hatzis
- ▷ Christopher Mason
- ▷ Leming Shi
- ▷ Anna Goldenberg
- ▷ Nicolai Juul-Birkbak
- ▷ Andrew Beck
- ▷ Hugo Aerts



Canadian  
Cancer  
Society



**Thank you  
for your attention!**

**Questions?**