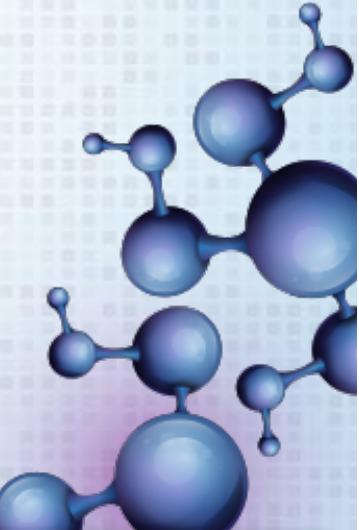




Broad-DREAM Gene Essentiality Prediction Challenge



Integrative Cancer
Biology Program



IBM Research



Unil
UNIL | Université de Lausanne

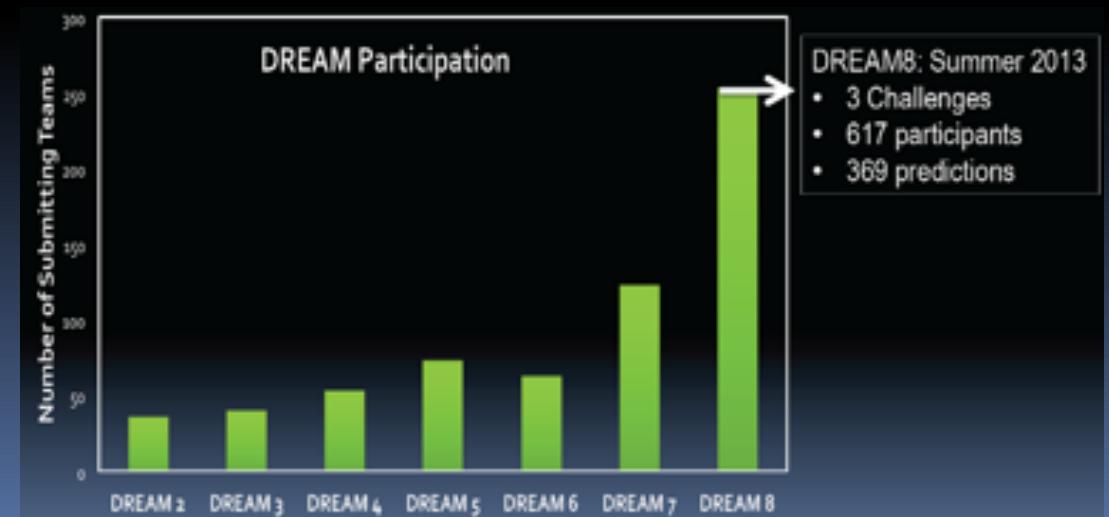
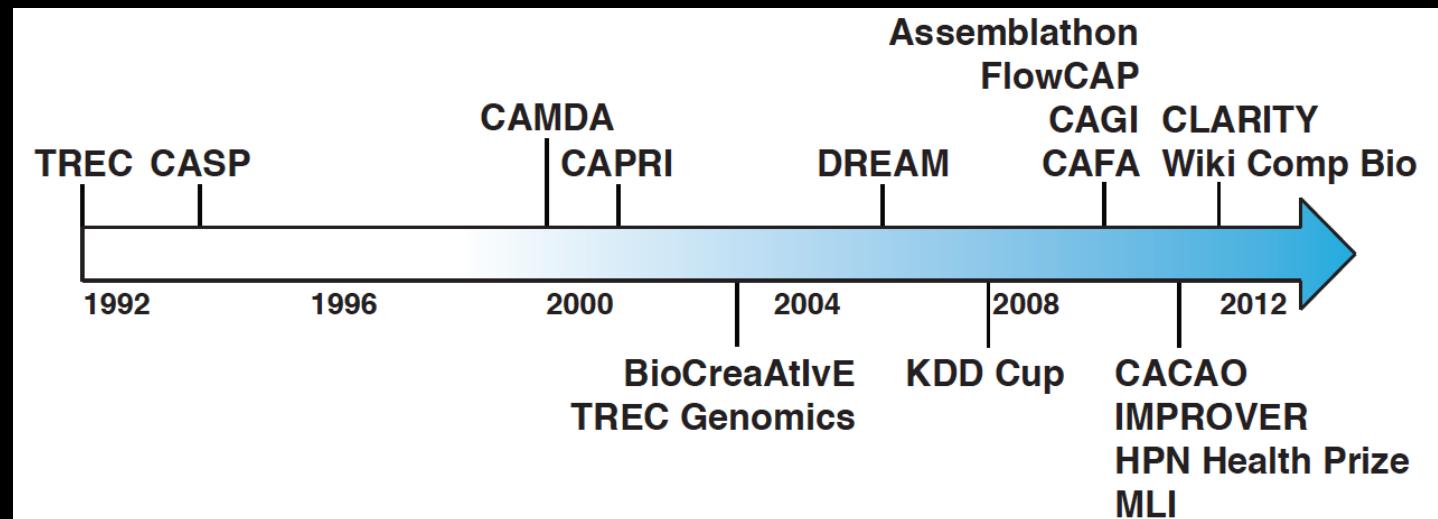




Outline

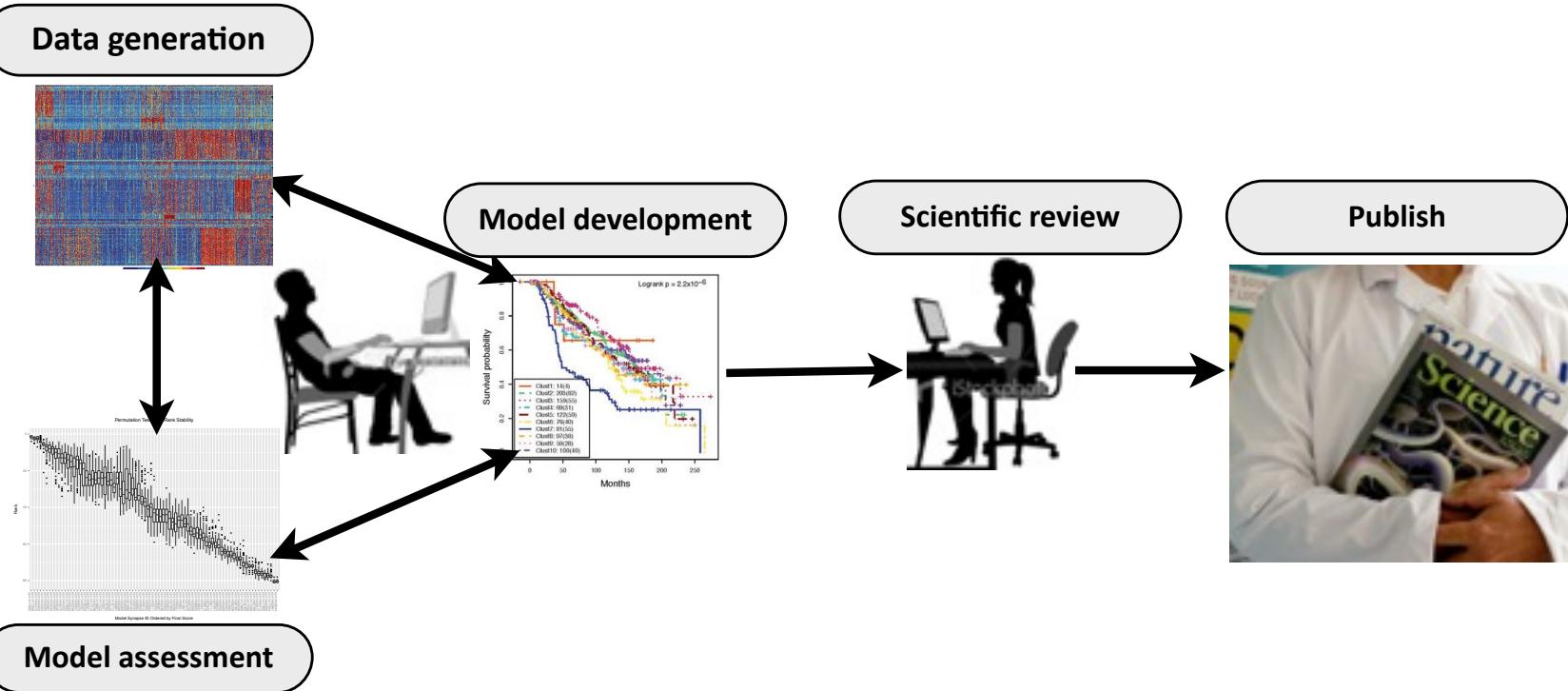
- DREAM overview (Adam Margolin, OHSU)
- Scientific overview (Barbara Weir, Aviad Tsherniak, Broad Institute)
- Challenge details (Mehmet Gonen, Sage Bionetworks)
- Synapse and forum registration (Bruce Hoff, Sage Bionetworks)
- Challenge demo (Mehmet Gonen, Sage Bionetworks)
- Compute resources (Venkat Balagurusamy, IBM)
- Questions and answers (all)

Crowd-sourcing in computational biology





The self-assessment trap



Molecular Systems Biology 7; Article number 537; doi:10.1038/msb.2011.70
Citation: Molecular Systems Biology 7: 537
© 2011 EMBO and Macmillan Publishers Limited. All rights reserved 1744-4292/11
www.molecularsystemsbiology.com

CORRESPONDENCE

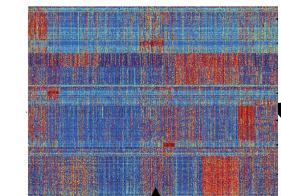
The self-assessment trap: can we all be better than average?

molecular
systems
biology

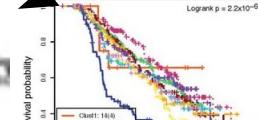


The self-assessment trap

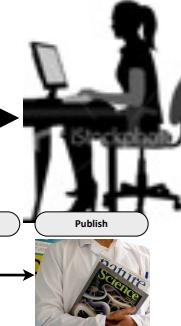
Data generation



Model development



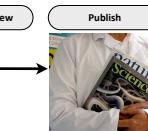
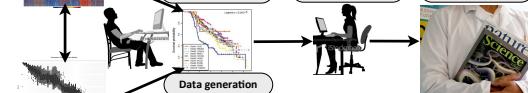
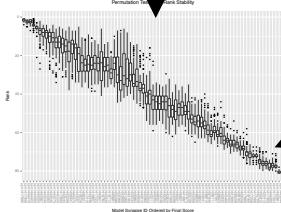
Scientific review



Publish



Model assessment



Molecular Systems Biology 7; Article number 537; doi:10.1038/msb.2011.70
Citation: Molecular Systems Biology 7:537
© 2011 EMBO and Macmillan Publishers Limited. All rights reserved 1744-4292/11
www.molecularsystemsbiology.com

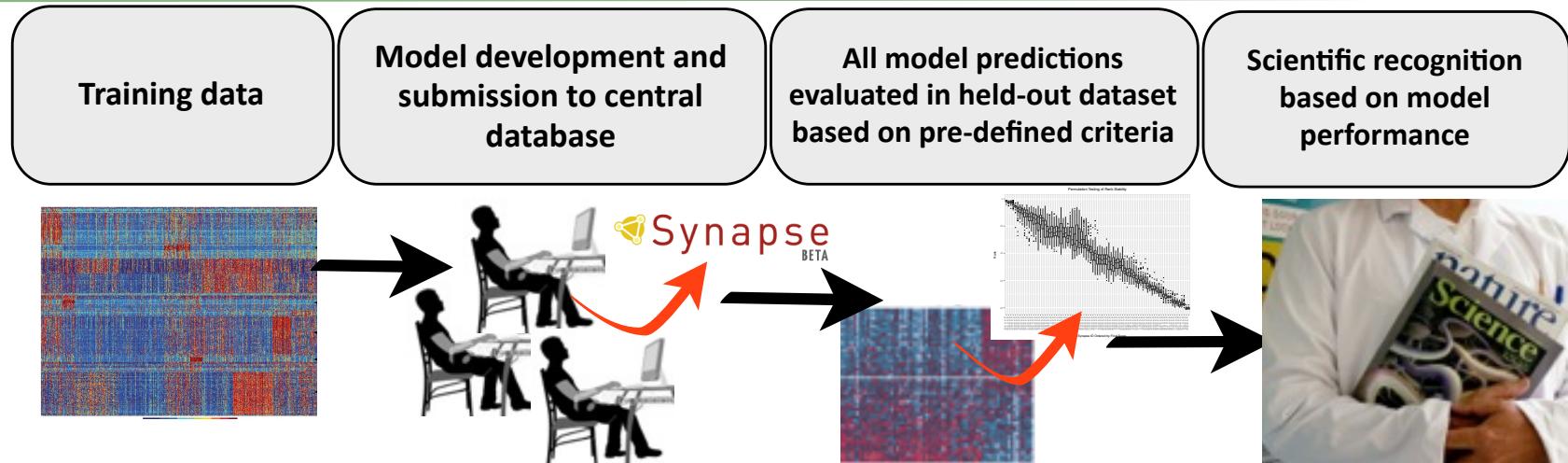
CORRESPONDENCE

The self-assessment trap: can we all be better than average?

molecular
systems
biology

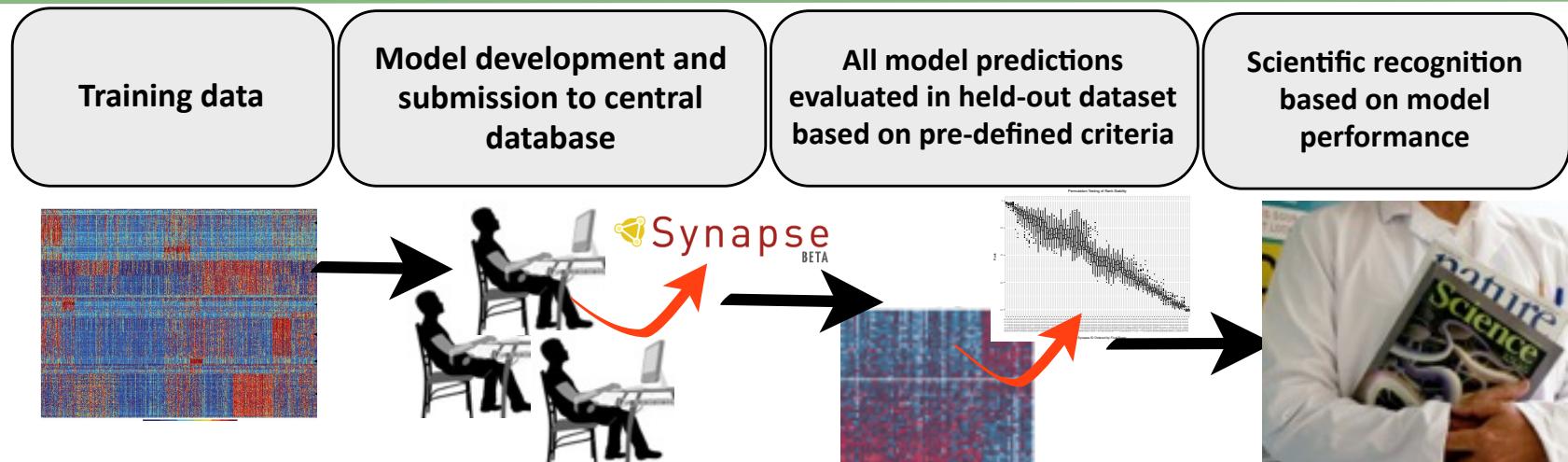


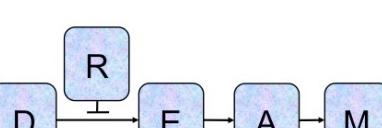
The Sage / DREAM breast cancer prognosis challenge





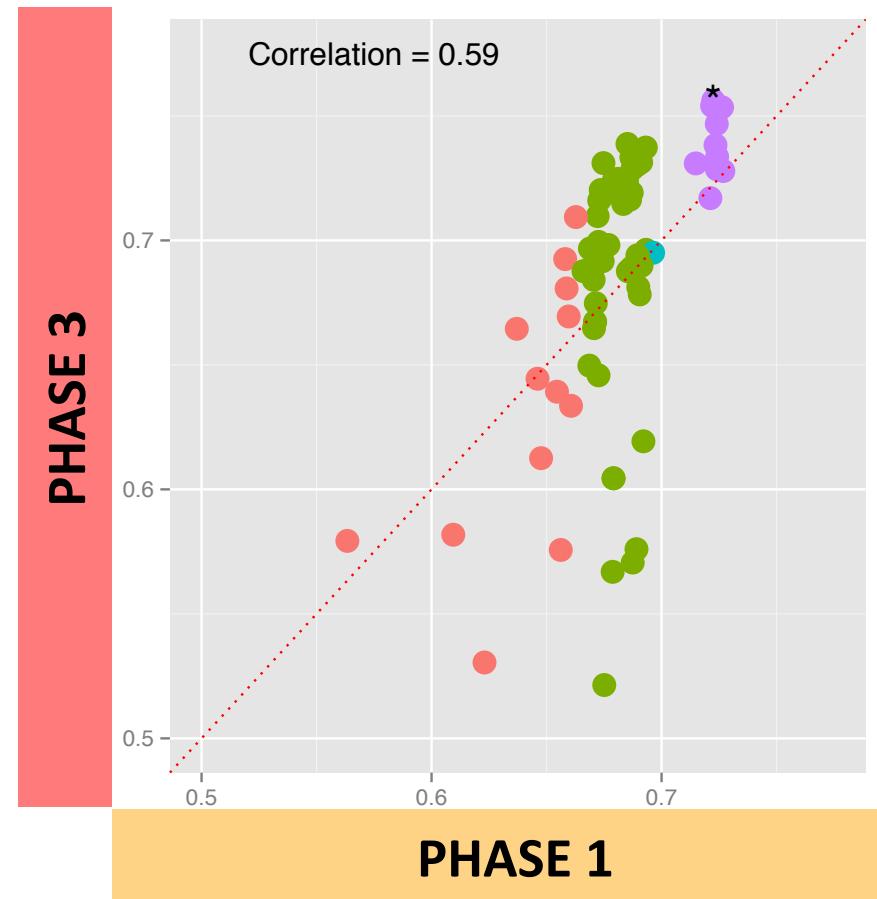
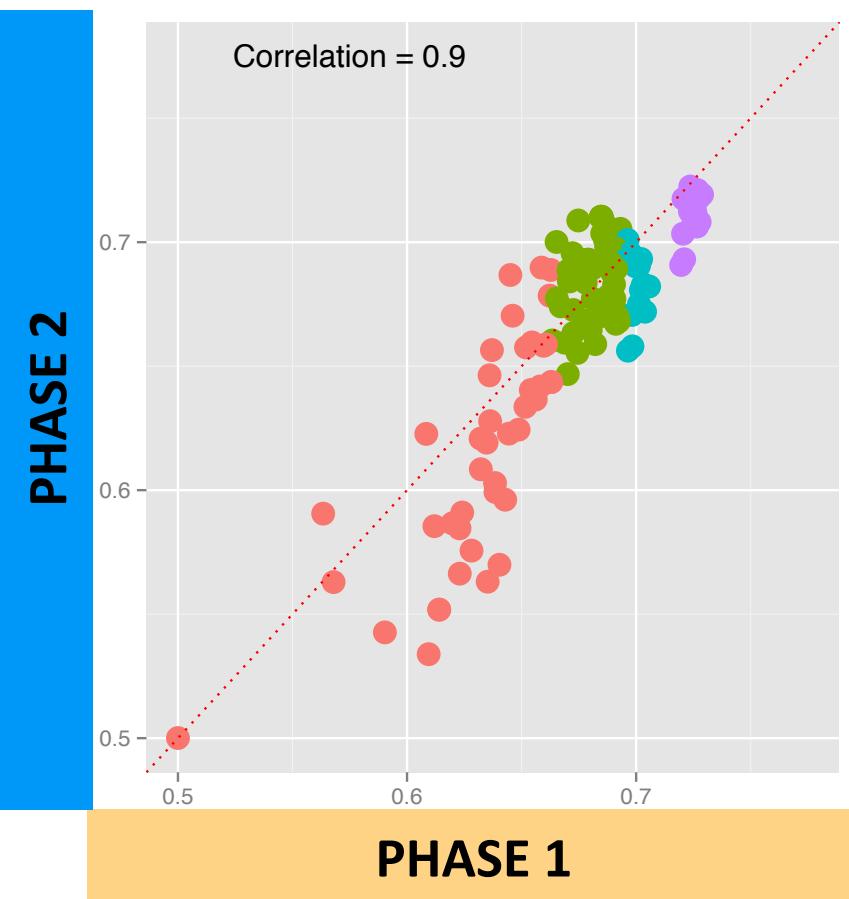
The Sage / DREAM breast cancer prognosis challenge



| PHASE 1 | 1,000 METABRIC samples | Unlimited model submission | 500 held-out METABRIC samples | Real-time leaderboard |
|---------|----------------------------|----------------------------|--|---|
| PHASE 2 | 1,000 METABRIC samples | 5 models per team | Additional 481 held-out METABRIC samples |  <pre> graph LR D[D] --> R[R] D --> E[E] R --> E E --> A[A] A --> M[M] M --> A </pre> |
| PHASE 3 | All 1,981 METABRIC samples | 5 models per team | 184 “OsloVal” samples |  Science Translational Medicine  |



Model concordance index scores are consistent across challenge phases



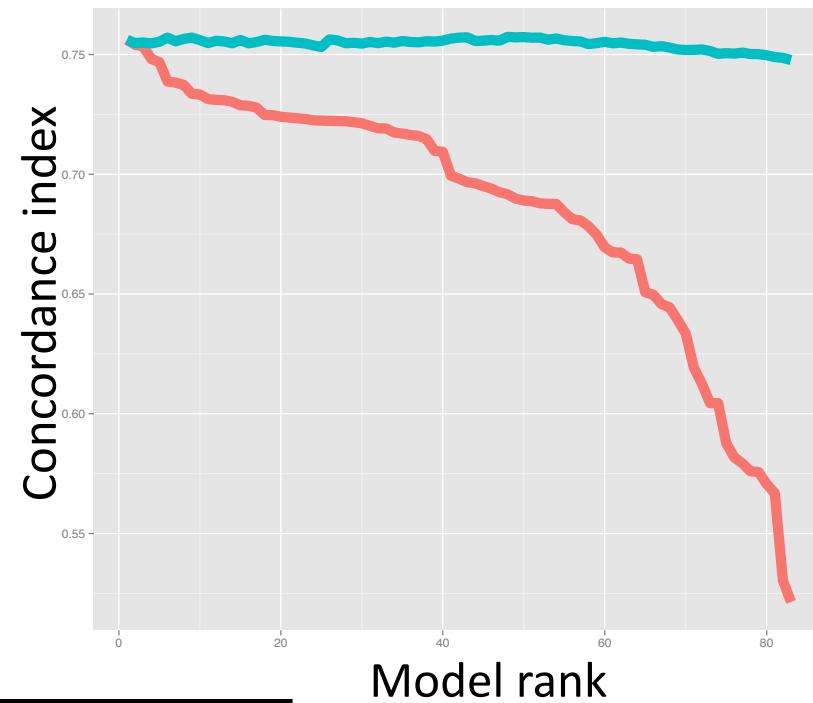


Leveraging the wisdom of crowds yields robust model performance

PHASE 2



PHASE 3



Individual model score

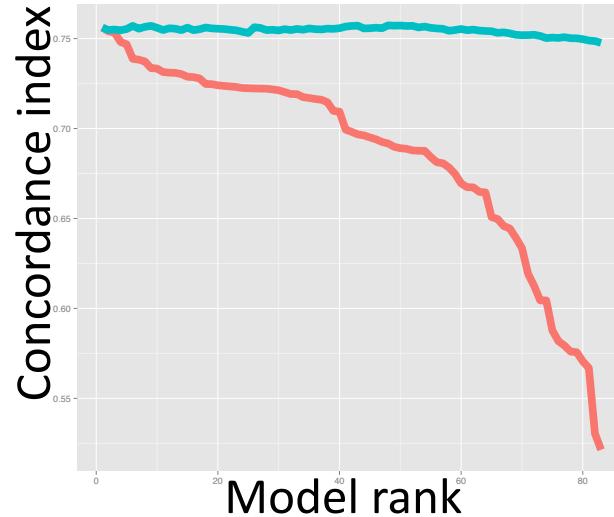
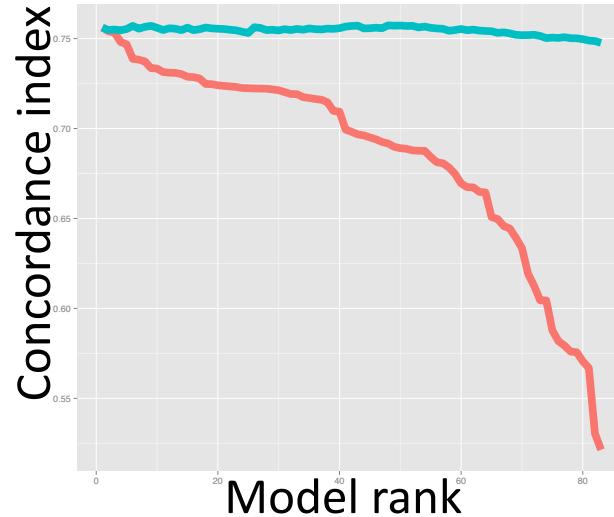
Consensus model score



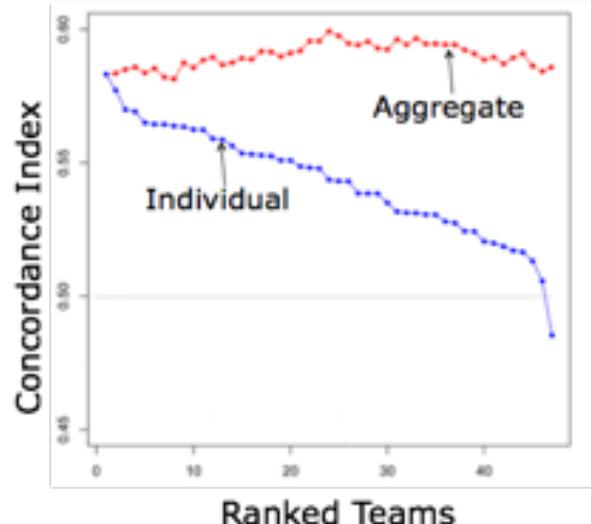
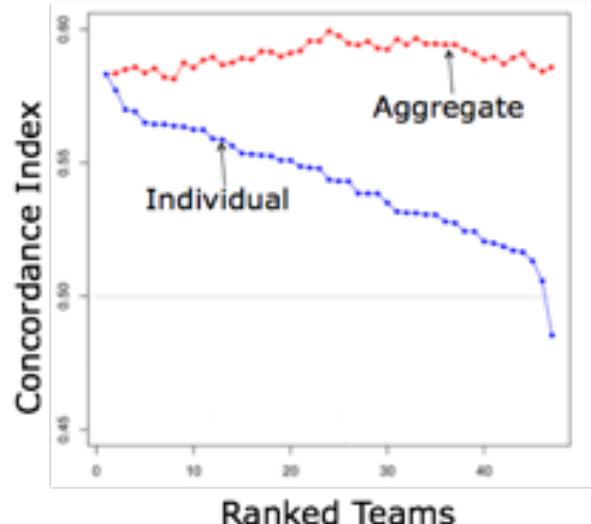


Vox populi across multiple challenges

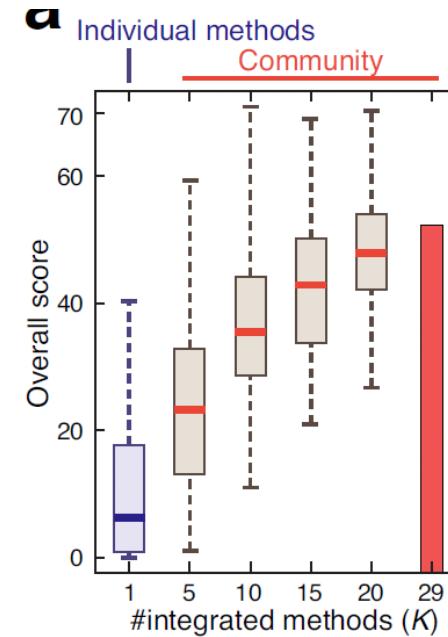
Breast cancer prognosis



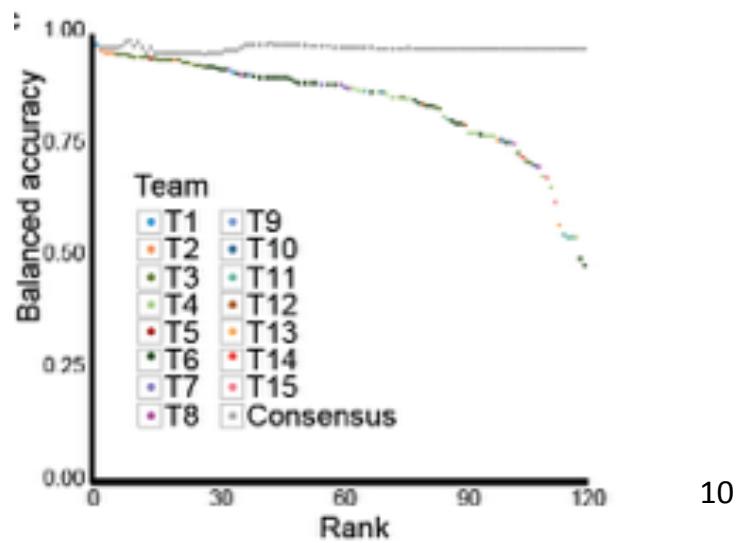
Drug sensitivity prediction

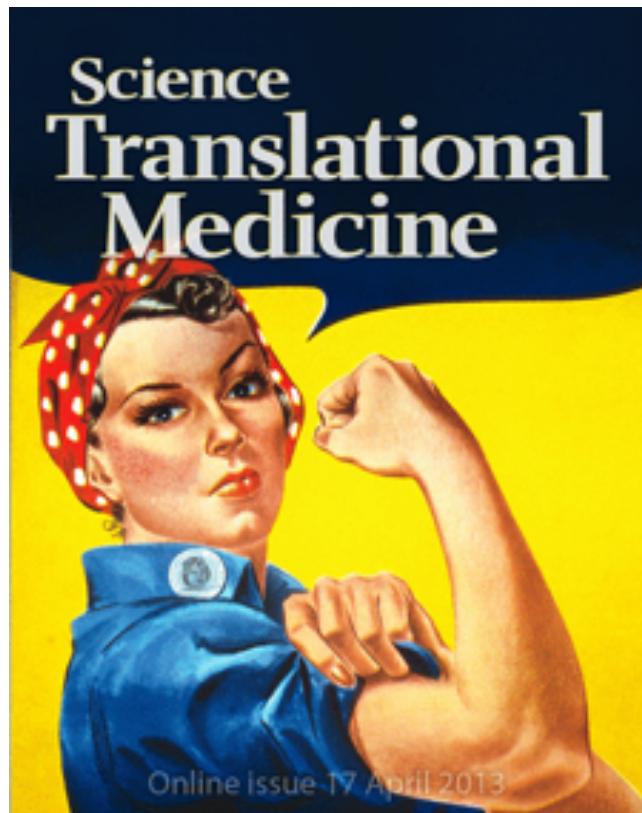


Network inference



Sequence analysis (in silico #1)





17 APRIL 2013, VOL. 5, #181

COVER STORY | RESEARCH ARTICLE AND REPORT

DREAMing of Biomedicine's Future

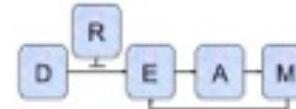
An Open Challenge yields fresh insights (Margolin et al.), a new prognostic model for breast cancer (Cheng et al.), and a modern approach to peer review (Editor's Summary).

FREE

All articles free and open access



CROWDSOURCING



Systematic Analysis of Challenge-Driven Improvements in Molecular Prognostic Models for Breast Cancer

Adam A. Margolin,^{1,†} Erhan Bilal,^{2†} Erich Huang,^{1,3,4†} Thea C. Norman,¹ Lars Ottestad,⁵ Brigham H. Mecham,^{1,6} Ben Sauerwine,⁷ Michael R. Kellen,¹ Lara M. Mangravite,¹ Matthew D. Furia,^{1,8} Hans Kristian Moen Vollan,^{9,10,11} Oscar M. Rueda,¹¹ Justin Guinney,¹ Nicole A. Deflaux,¹ Bruce Hoff,¹ Xavier Schildwachter,¹ Hege G. Russnes,^{9,10,12} Daehoon Park,¹³ Veronica O. Vang,^{9,10} Tyler Pirtle,⁷ Lamia Youseff,⁷ Craig Citro,⁷ Christina Curtis,¹⁴ Vessela N. Kristensen,^{9,10,15} Joseph Hellerstein,⁷ Stephen H. Friend,^{1,*} Gustavo Stolovitzky,² Samuel Aparicio,^{16,17,18†} Carlos Caldas,^{11,19,20†} Anne-Lise Barresen-Dale^{9,10†}



RESEARCH ARTICLE

COMPUTATIONAL MODELING

Development of a Prognostic Model for Breast Cancer Survival in an Open Challenge Environment

Wei-Yi Cheng, Tai-Hsien Ou Yang, Dimitris Anastassiou*





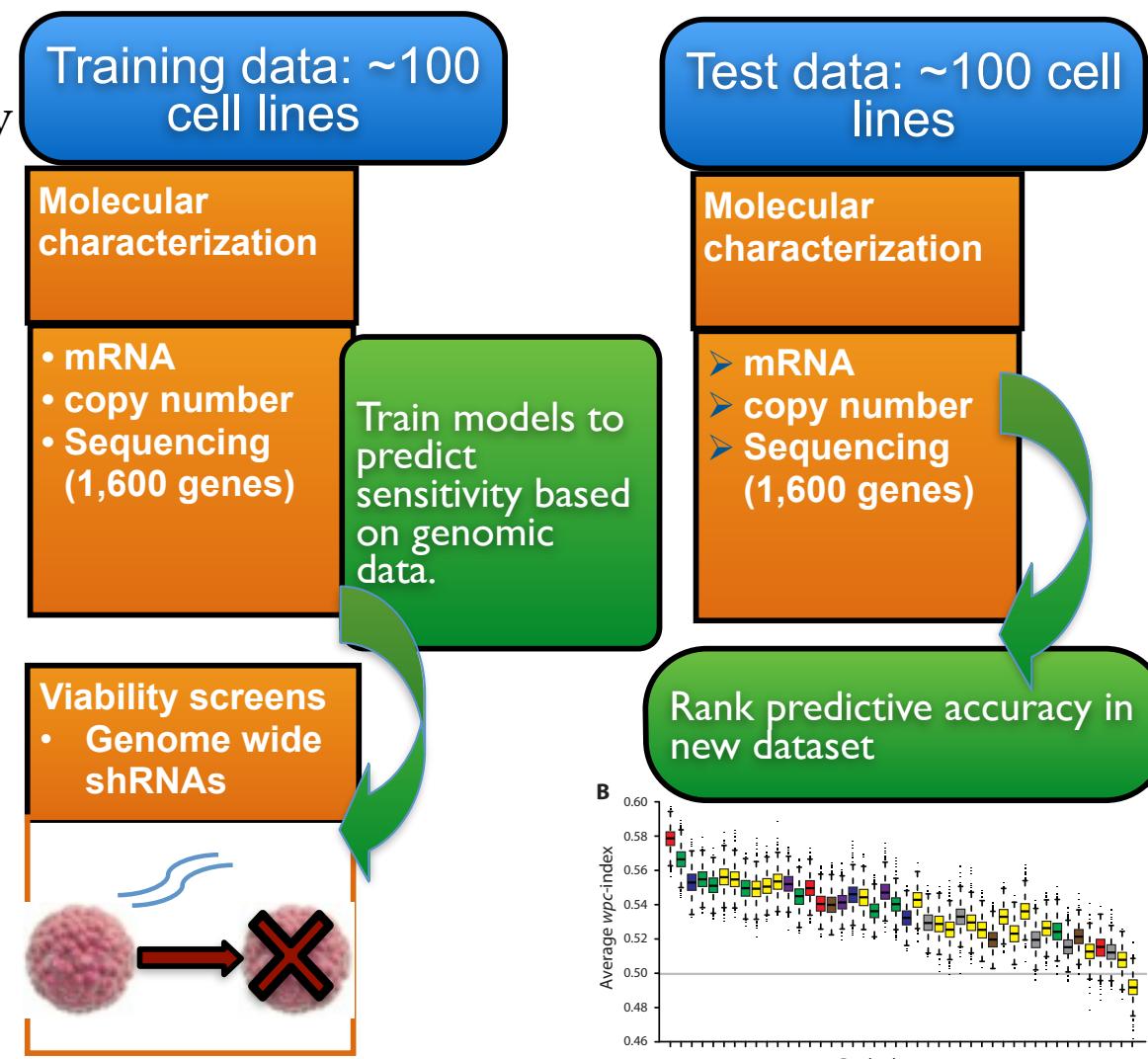
Broad gene essentiality prediction challenge

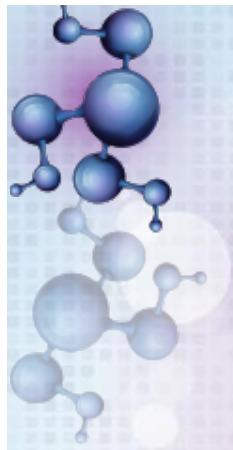
- Project Achilles dataset of 200 cell lines.

- Molecular characterization: expression, sequencing, copy number
- Genome wide RNAi screens
- 100 training samples.
- 100 test samples released with challenge results.

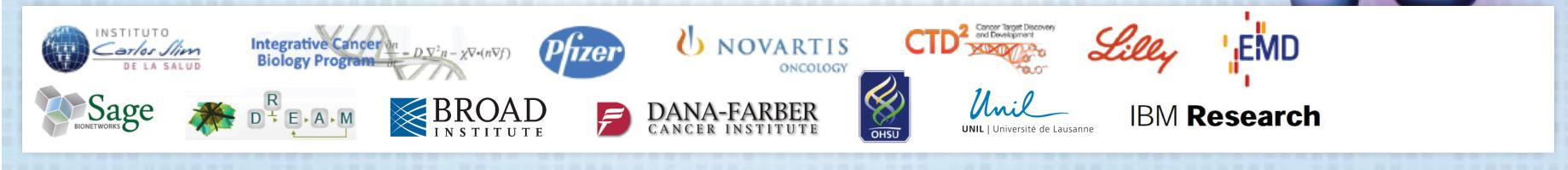
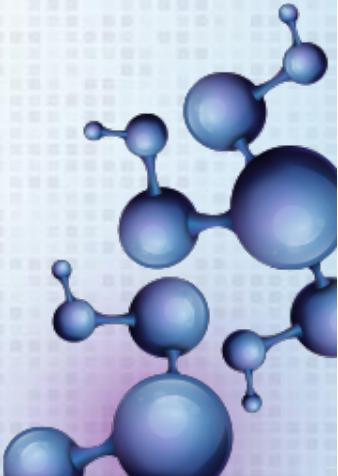
- <https://www.synapse.org/#!/Challenges:DREAM>

- Or navigate from www.synapse.org





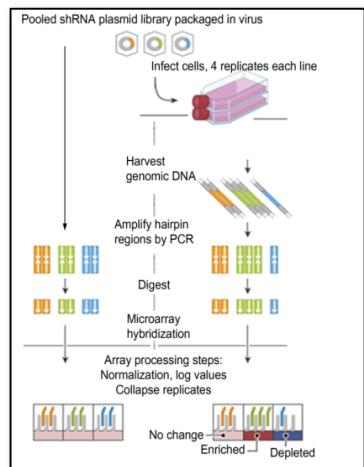
Broad-DREAM Gene Essentiality Prediction Challenge



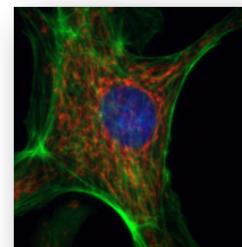
Outline

1. DREAM overview (Adam Margolin, OHSU)
2. Scientific overview (Barbara Weir, Aviad Tsherniak, Broad Institute)
3. Challenge details (Mehmet Gonen, Sage Bionetworks)
4. Synapse and forum registration (Bruce Hoff, Sage Bionetworks)
5. Challenge demo (Mehmet Gonen, Sage Bionetworks)
6. Compute resources (Venkat Balagurusamy, IBM)
7. Questions and answers (all)

Project Achilles

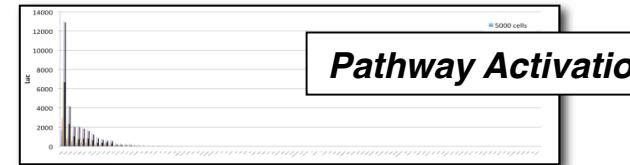


Dependencies

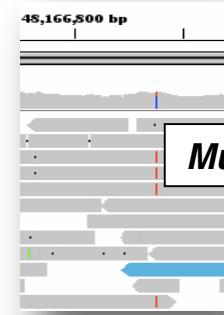


Genome-wide RNAi screens in a large number of cell lines

Extensive molecular characterization of cell lines



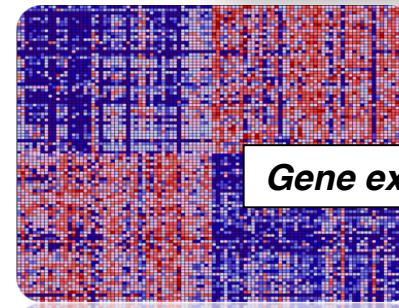
Pathway Activation



Mutational analysis



Copy number



Gene expression

Project Achilles

- Project Achilles is a large collaboration between the Broad Institute and the Dana-Farber Cancer Institute to ***identify tumor vulnerabilities that are linked with predictive biomarkers.***
 - **The Broad-DREAM Gene Essentiality Challenge** is designed to help address this goal, specifically to look for models that can predict these vulnerabilities from biomarkers (i.e. molecular features of the tumor)
- Current **public** dataset is Achilles v2.0
 - 'Systematic investigation of genetic vulnerabilities across cancer cell lines reveals lineage-specific dependencies in ovarian cancer', PNAS 2011 Jul 26;108(30):12372-7.
 - 102 cell lines, array deconvolution of screening data
 - Available, along with additional information about Project Achilles, on the Project Achilles portal (next slide)
- Challenge dataset is **unpublished** and new for you to use!
 - 89 cell lines (half held back for scoring) that also have copy number and expression array data
 - Additional cell lines coming (*data is being generated now*)
 - Challenge will be re-launched in late July / early August when the new data is ready

Project Achilles

Project Achilles Portal

- Genes can be queried for essentiality values across all cell lines for published data
- Data at the shRNA level and gene-level collapsed by ATARiS are shown and can be:
 - easily downloaded
 - opened with GENE-E
 - launched into our data analysis GenePattern module PARIS.
- More functionality, including data mining tools, is currently being developed.

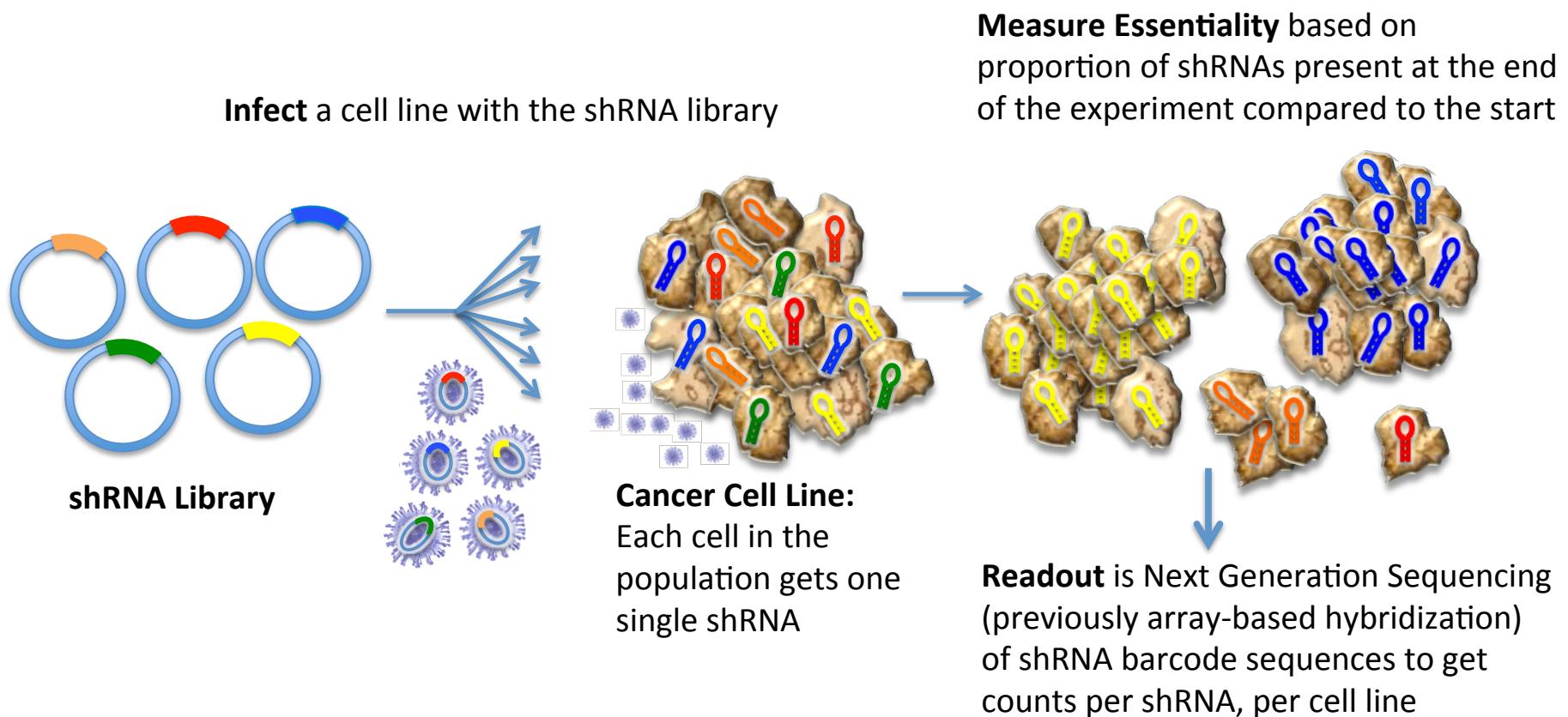


The screenshot shows the Project Achilles website interface. At the top, there is a purple header bar with the text "Project Achilles | BROAD INSTITUTE" and navigation links for "Data", "Tools", "About", "Contact", and "Admin". On the far right of the header, it says "Logout: bweir". Below the header is a large, stylized graphic of a DNA double helix composed of colored squares (blue, red, yellow) against a purple background with text. To the right of the graphic is a search bar with the placeholder "Search Genes" and a "Search" button, along with a dropdown menu showing "Dataset: Achilles v2.4.2 (02/20/14)". Below the search bar is a detailed description of the project's goal: "Project Achilles is a systematic effort aimed at identifying and cataloging genetic vulnerabilities across hundreds of genetically characterized cancer cell lines. The project uses a genome-wide shRNA library to silence individual genes and identify those genes that affect cell survival. Large-scale functional screening of cancer cell lines provides a complementary approach to those studies that aim to characterize the molecular alterations (mutations, copy number alterations, etc.) of primary tumors, such as The Cancer Genome Atlas. The overall goal of the project is to link cancer genetic dependencies to their molecular characteristics in order to identify molecular targets and guide therapeutic development."

- <http://www.broadinstitute.org/achilles>

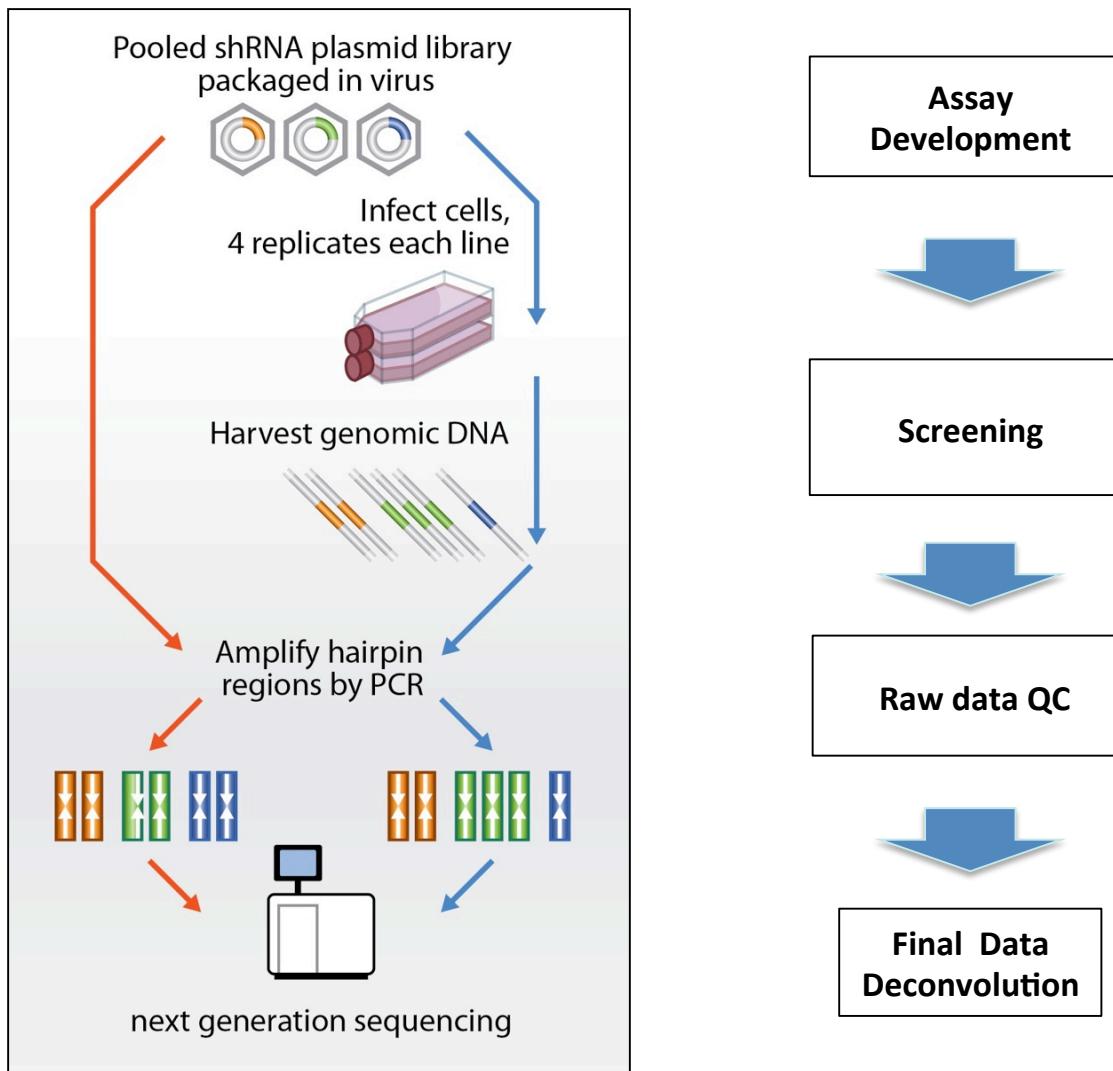
shRNA Loss of Function Screens

Overview

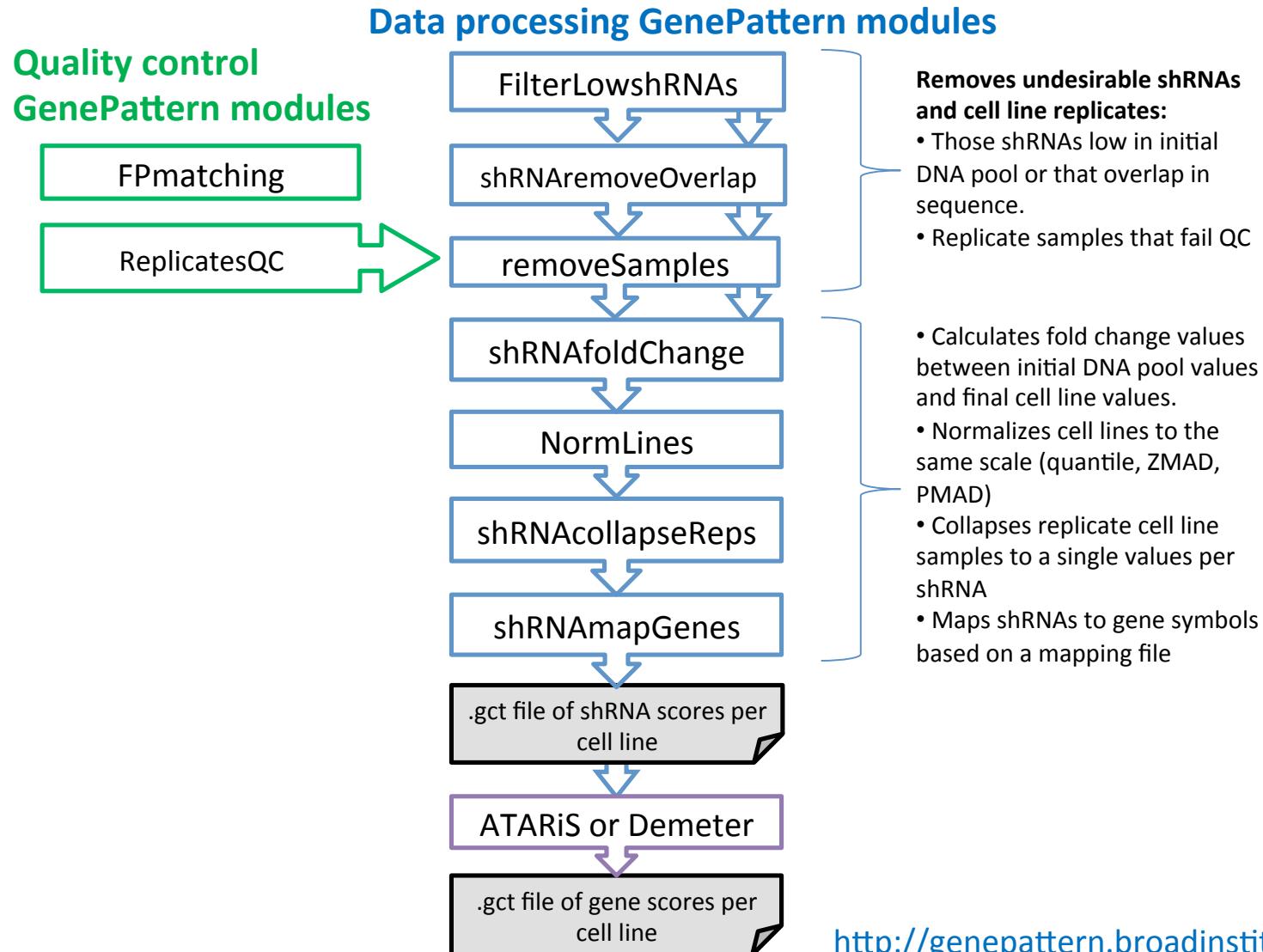


Red and Green in this example are depleted during the experiment and therefore have a negative effect on viability (are 'essential' to viability). Orange is neutral to viability. Blue and Yellow are enriched during the experiment and therefore have a positive effect on viability.

Project Achilles



Project Achilles

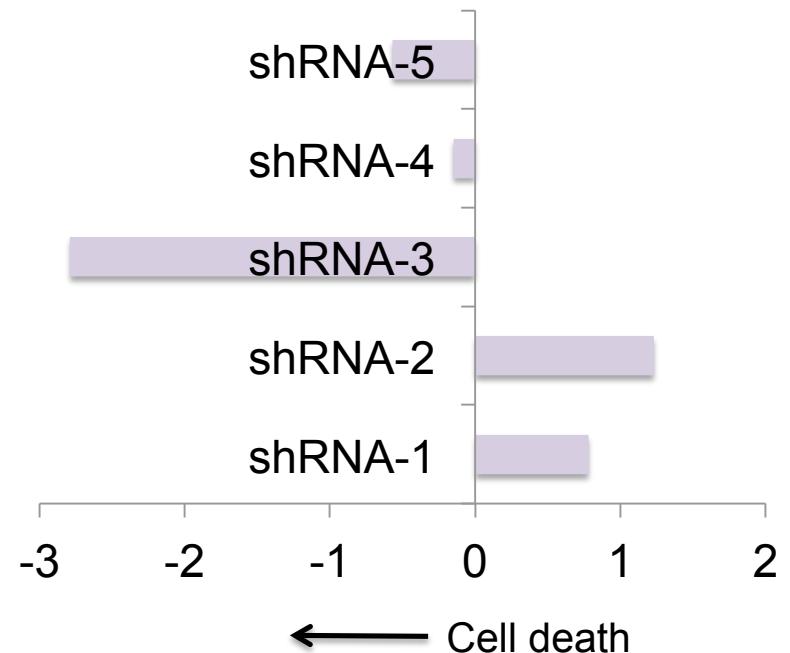


<http://genepattern.broadinstitute.org/gp/>

RNAi Data is Challenging

- shRNAs targeting the same gene produce different effects
- Sources of variability
 - Varying suppression levels
 - Off-target effects
 - Noise
- No experimental way to fully assess shRNA performance
 - Especially off-target suppression
- Naïve averaging of the values will hurt the signal

Viability data for one gene in one sample

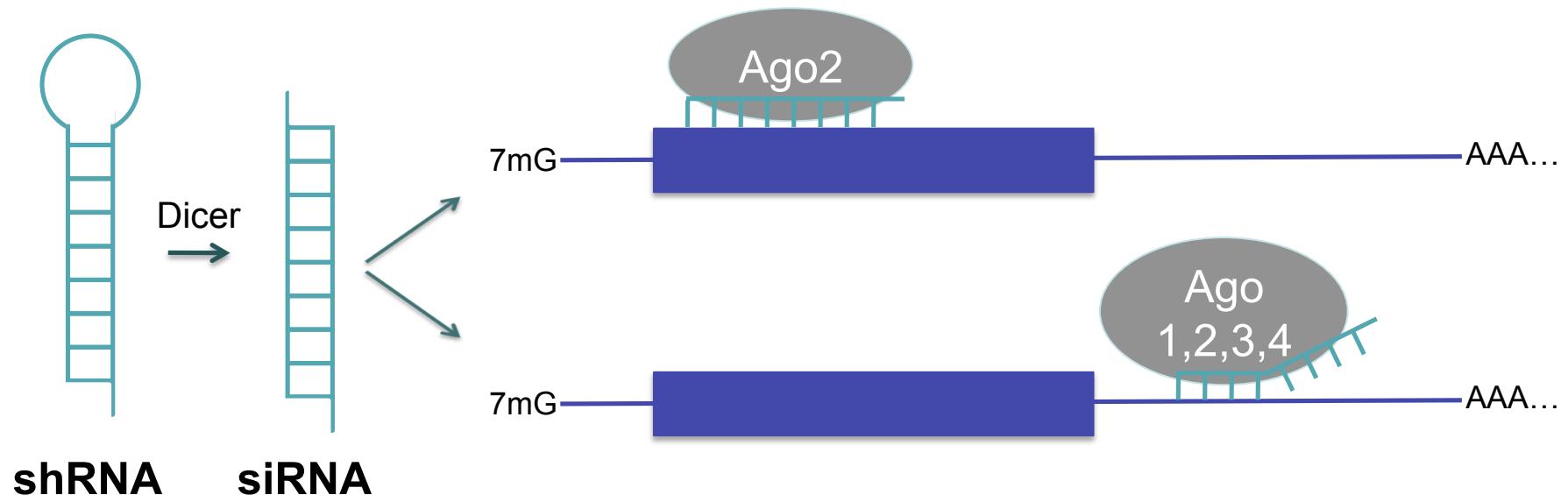


Two pathways for shRNAs



RNAi

Perfect (or near-perfect) match to mRNA causes mRNA cleavage and degradation: **18 – 22 nt**



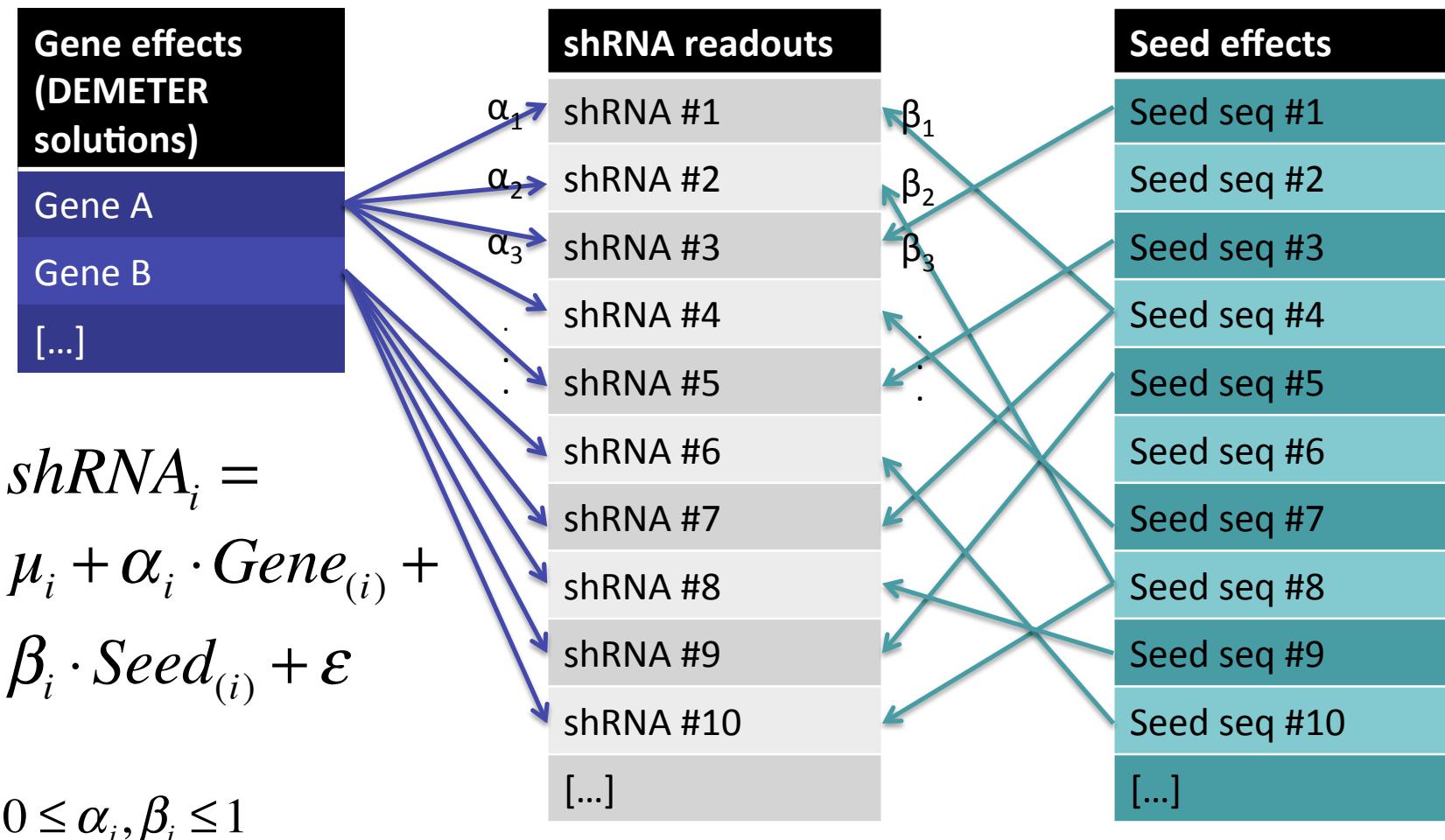
shRNA

siRNA

microRNA

'Seed' region binds to 3'UTR to represses translation and/or destabilize mRNA: **6 – 8 nt**

DEMETER infers gene dependencies by explicitly modeling both gene- and seed-effects



Challenge Questions

Pages

Broad-DREAM Gene Essentiality Prediction Challenge

1. Challenge Overview

1.1 Challenge Questions & Scoring

1.3 Timeline & Incentives

2. Data Description

2.1 Gene Essentiality Details

2.2 Data File Formats

2.3 Downloading Data Files

3. Participation

3.1 Submission File Formats

3.2 Submitting Results

4. Leaderboards

5. Compute Resources

6. Challenge Forum

7. Challenge Organizers

8. Literature

<<

Broad-DREAM Gene Essential... » 1.1 Challenge Questions ...

1.1 Challenge Questions & Scoring

[Edit Wiki](#) [+ Add a new Page](#)

The goal of this project is to use a crowd-based competition to develop predictive models that can infer gene dependencies (genes that are essential to cell viability when suppressed) in cancer cells using features of the cell line.

Participants will be asked to solve three sub-challenges:

1. Build a model that best predicts the gene essentialities of **14557** genes, using the molecular characteristics/features of the cancer cell lines.
In this sub-challenge any data can be used for prediction.

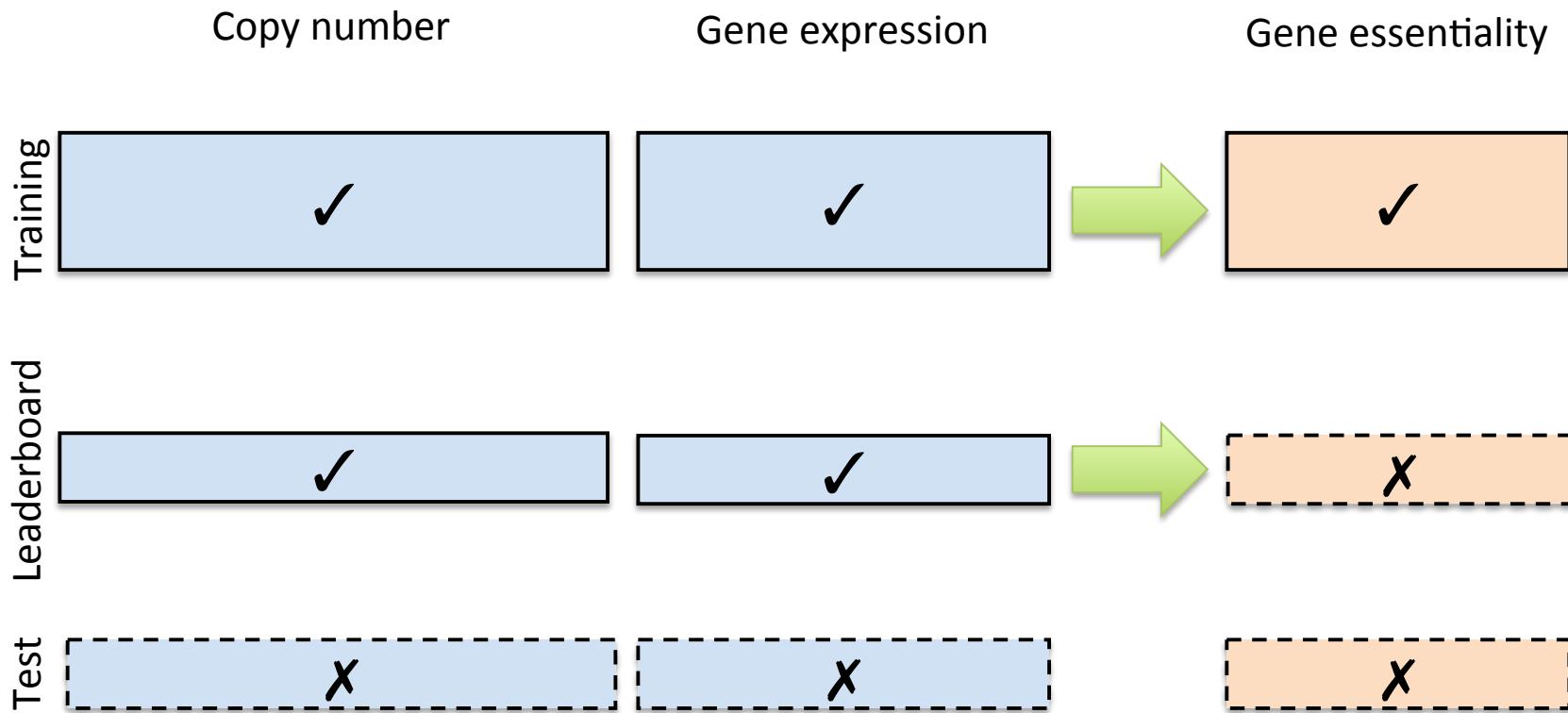
2. Identify the most predictive features for **each** gene essentiality of a prioritized list of **2637** genes. For each prioritized gene the aim is to: (1) select a small set of at most **10** predictive features (gene expression and copy-number) and (2) predict gene essentiality using **only** these features.
In this sub-challenge any data can be used for choosing the gene expression and copy-number features, but only those selected features can be used for prediction.

3. Identify the most predictive features for **all** gene essentialities of a prioritized list of **2637** genes. For the set of all prioritized genes, the aim is to: (1) identify a single list of at most **100** general predictive features (gene expression and copy-number) and (2) predict essentiality using **only** these features for **all** prioritized genes. In machine learning, this problem is known as multilabel or multitask feature selection.
In this sub-challenge any data can be used for choosing the gene expression and copy-number features, but only those selected features can be used for prediction.

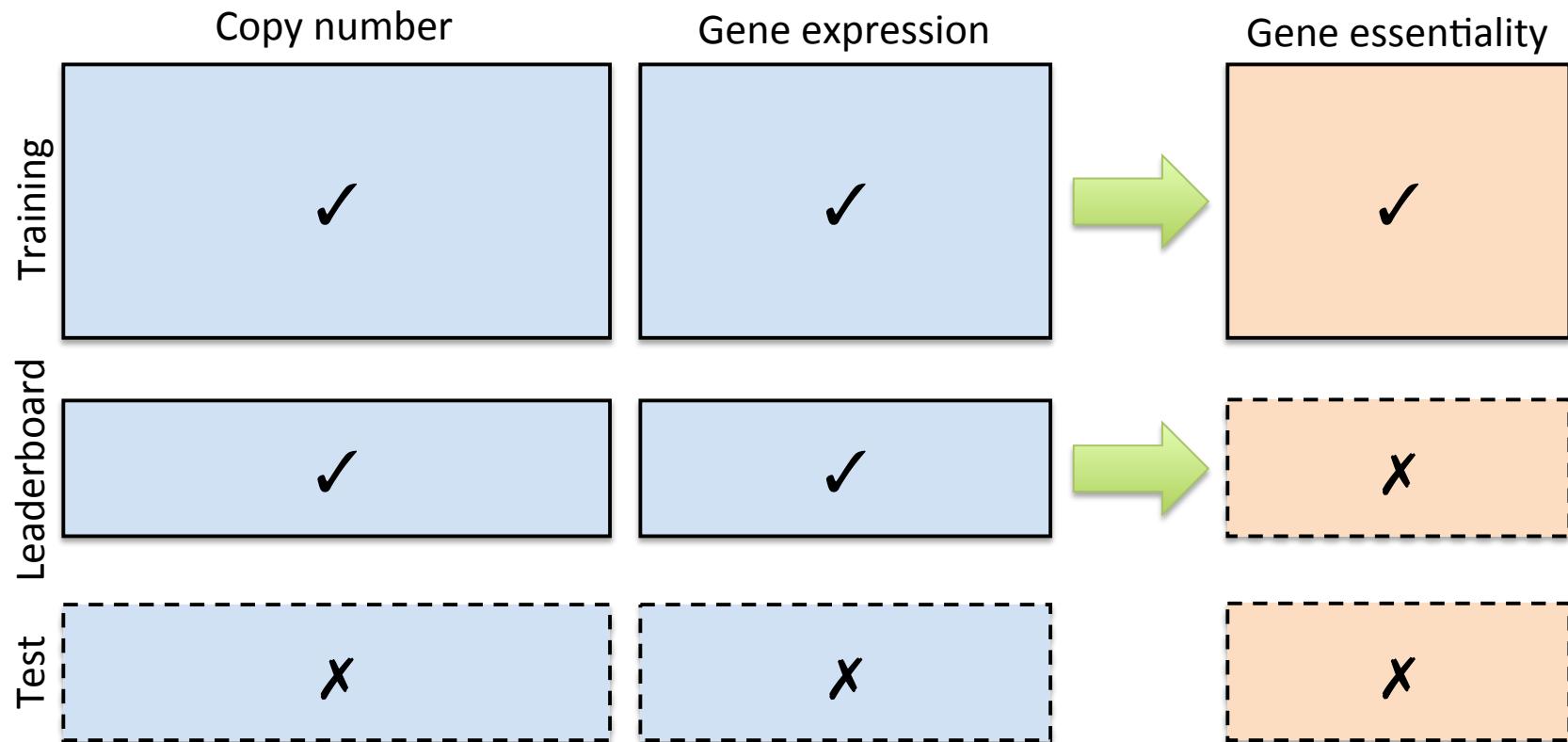
Challenge Timeline

- Phase I (89 cell lines): till end July
 - 45 in training set (fully available)
 - 22 in leaderboard set (hidden gene essentiality)
 - 22 in test set (fully hidden)
- Phase II (\approx 200 cell lines): till end August
 - \approx 100 in training set (fully available)
 - \approx 50 in leaderboard set (hidden gene essentiality)
 - \approx 50 in test set (fully hidden)
- Phase III (\approx 200 cell lines): till mid September
 - \approx 150 in training set (fully available)
 - \approx 50 in test set (hidden gene essentiality)

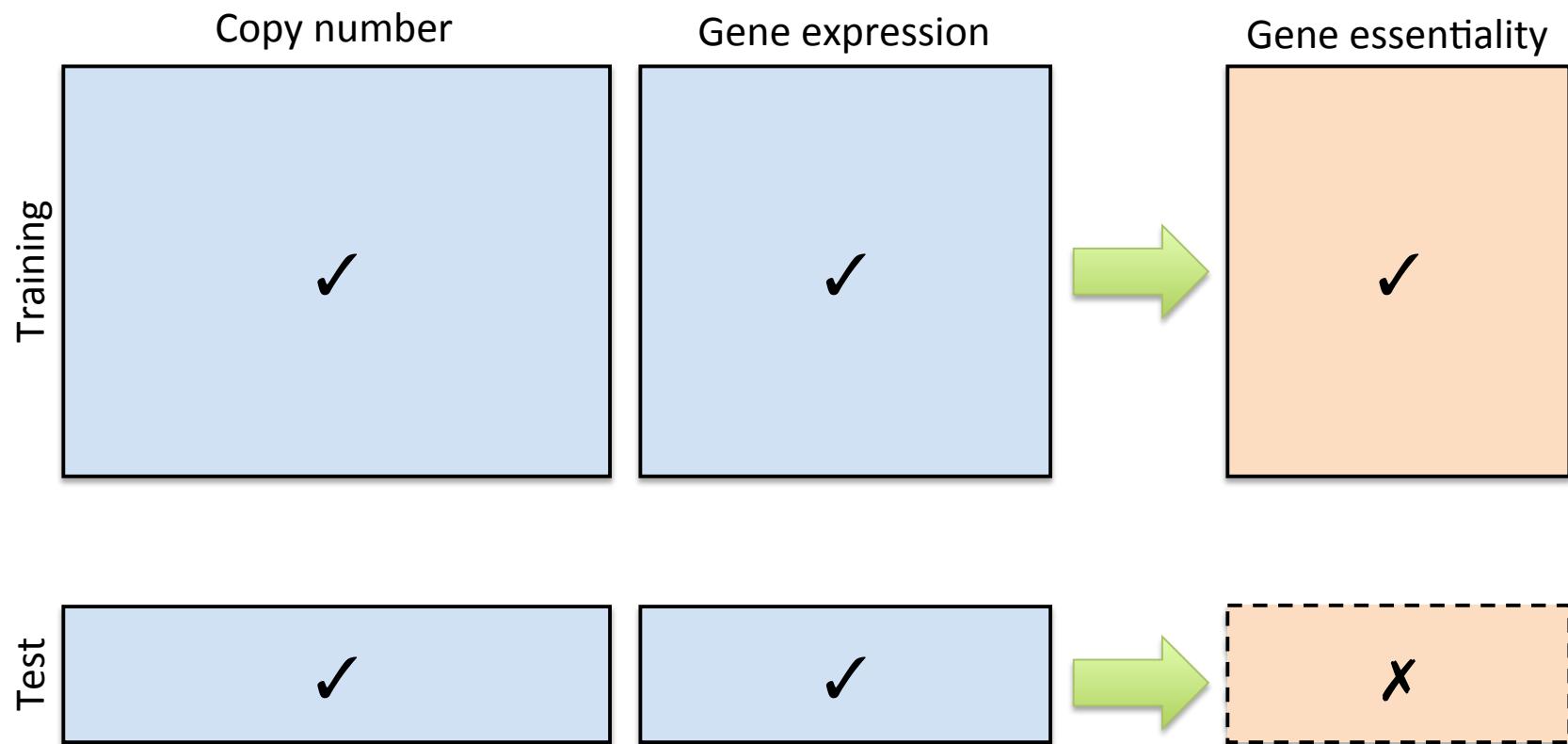
Overall Challenge Setup (Phase I)



Overall Challenge Setup (Phase II)



Overall Challenge Setup (Phase III)



Sub-challenge 1 Question

- Build a model that best predicts the gene essentialities of thousands of genes, using the molecular characteristics/features of the cancer cell lines

23288 copy number features 18960 expression features

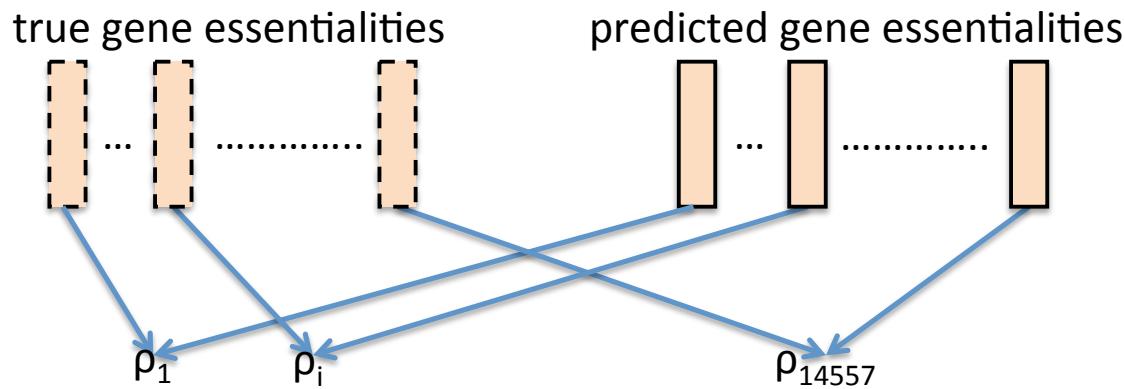


14557 gene essentialities



Sub-challenge 1 Scoring

- We use Spearman's rank correlation coefficient to evaluate the performance

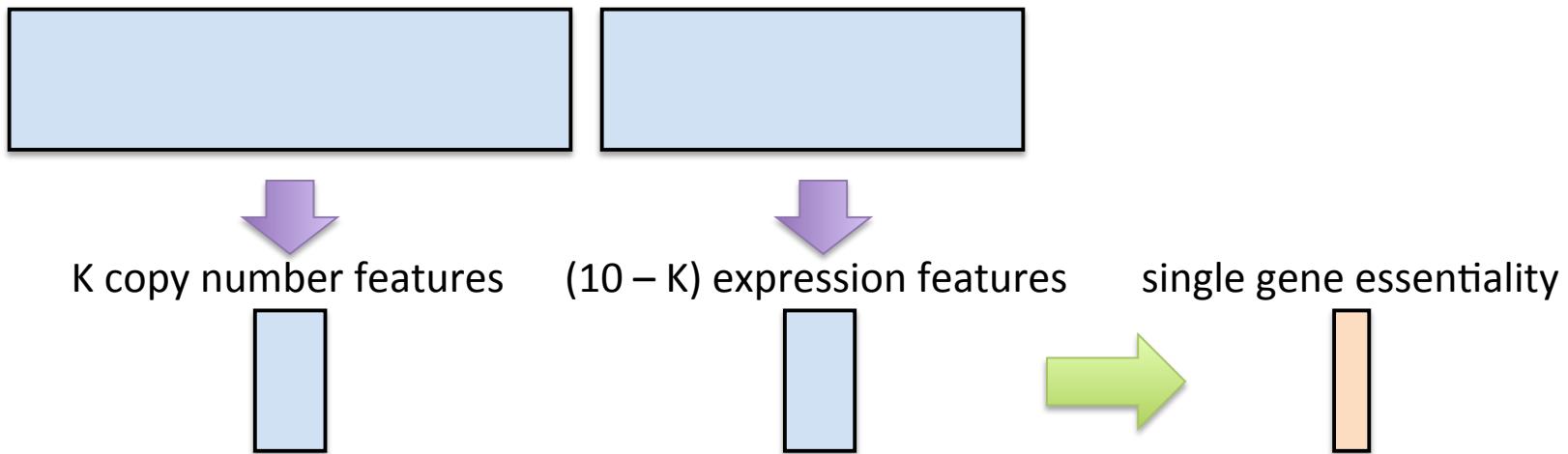


- Overall score is the mean of 14557 correlation values

Sub-challenge 2 Question

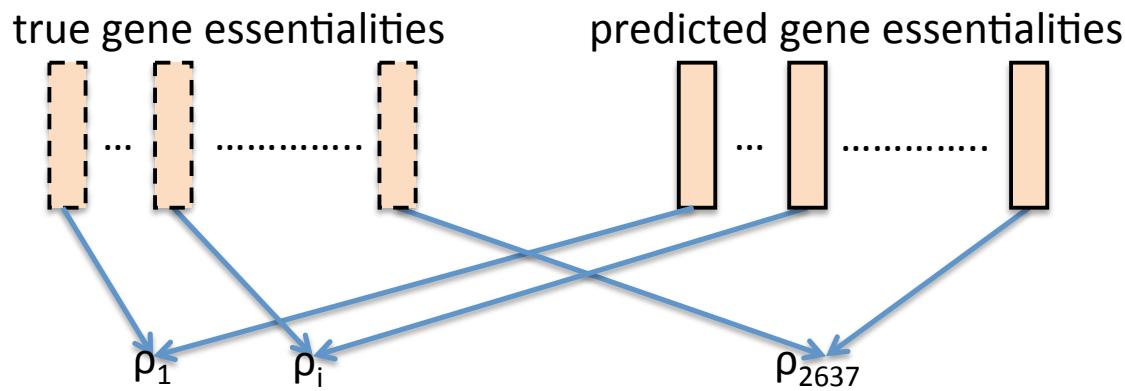
- Identify the most predictive features for each gene essentiality of a prioritized list of 2637 genes

23288 copy number features 18960 expression features



Sub-challenge 2 Scoring

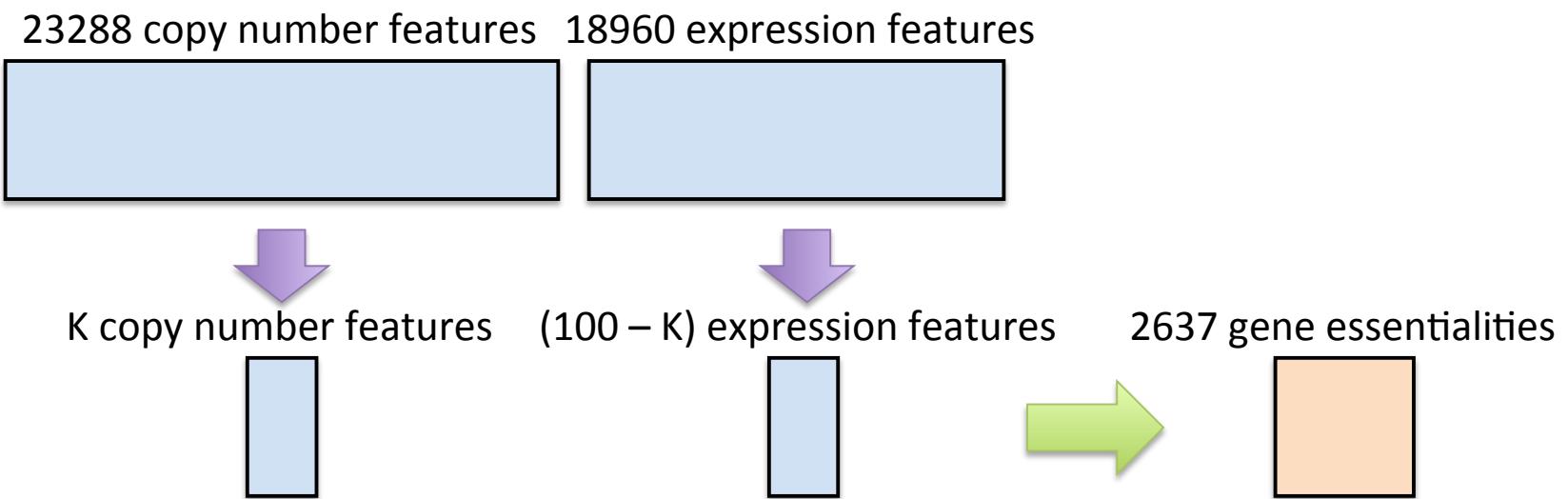
- We use Spearman's rank correlation coefficient to evaluate the performance



- Overall score is the mean of 2637 correlation values

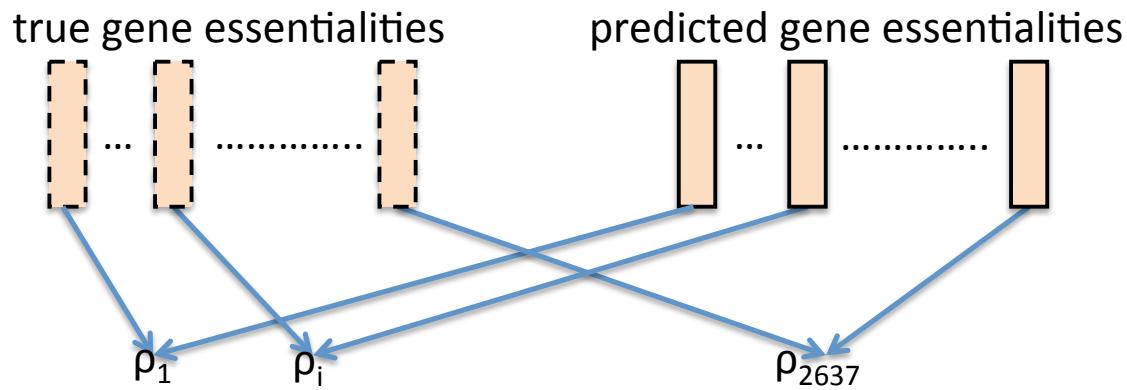
Sub-challenge 3 Question

- Identify the most predictive features for all gene essentialities of a prioritized list of 2637 genes



Sub-challenge 3 Scoring

- We use Spearman's rank correlation coefficient to evaluate the performance



- Overall score is the mean of 2637 correlation values

External Data Usage Policy

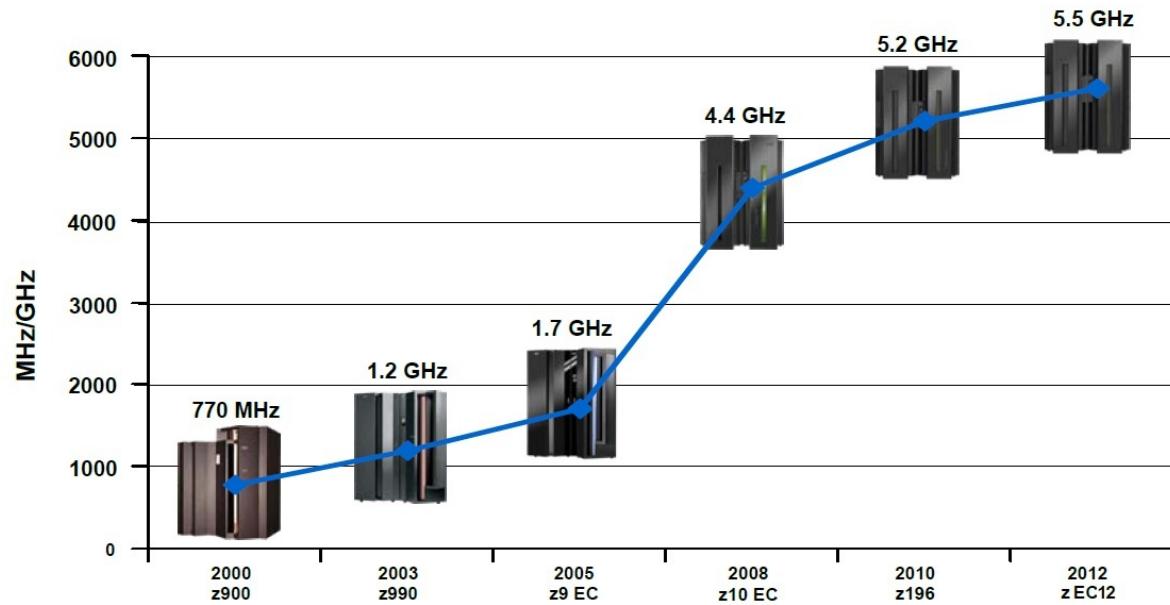
- Sub-challenge 1:
 - Other data sources can be used freely
- Sub-challenges 2 and 3:
 - Other data sources can be used for feature selection step only
 - Prediction step must be performed using selected copy number and expression features only

IBM Computing Cluster for Gene Essentiality Prediction Challenge

IBM System Z team

June 26, 2014

IBM System Z – fastest processor in the Industry



Available for the GEP Challenge – 20 processors, 297 GB memory, 1 TB storage (zpublic server)

http://en.wikipedia.org/wiki/IBM_zEC12_%28microprocessor%29

© 2014 IBM Corporation



How to get a user account in IBM System Z server

- Once a challenge user is approved by Sage to access the data, they become **eligible to get** a user account in IBM System Z Server
- IBM creates userid login and password and sends the participant an email message (at least once a day)
- **Welcome message and login details** will be sent to each approved user

Subject: The DREAM Gene Essentiality Challenge: Your IBM server account

Dear DREAM challenge participant,

The IBM System Z Team (Contact: Venkat Balagurusamy; Email: vkbala@us.ibm.com)

IP address: 63.90.228.10

Internet host name : zpublic.wsc.ihost.com



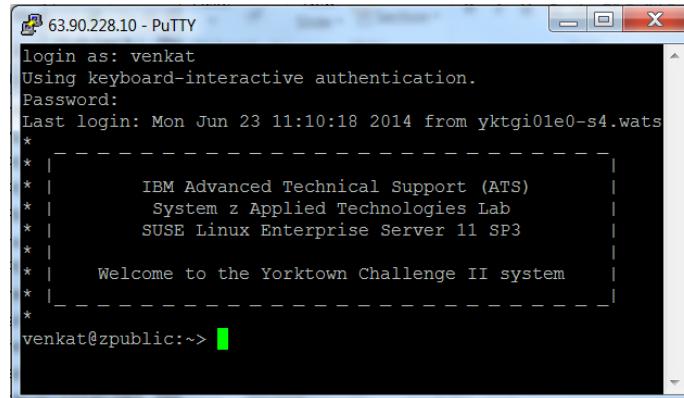
How to access the IBM System Z server

Once you have received the login details and the welcome message from Sagebase website (by email), you are ready to login to IBM System Z server

- **Windows users** – can use a secure shell client e.g., **PuTTY** to connect to the server with their login details
 - **File transfer:** can use e.g., **WinSCP**
- **Mac users** – use “ssh” command in a **terminal** window in their macbook
- **Command format:** “**ssh userid@zpublic.wsc.ihost.com**” OR “**userid@63.90.228.10**”
- **Password change after first login:** After their first successful login with the initial password (case-sensitive), they will be asked to change the password immediately. When they change it, the system will disconnect and you can relogin with the changed password.

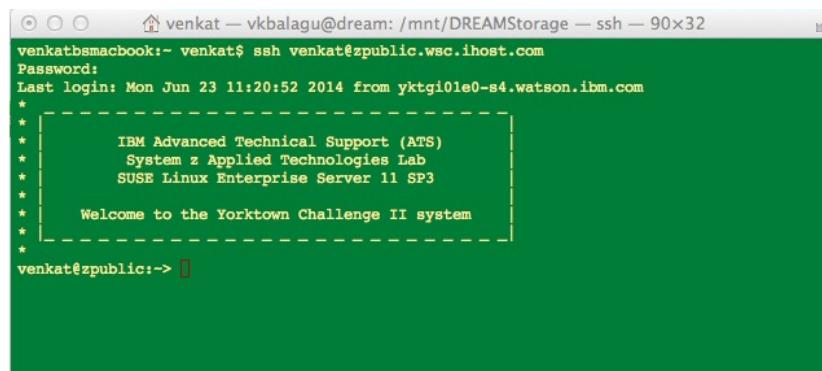
Screen shots of the server after changing password and login

Windows
screen (PuTTY)



```
63.90.228.10 - PuTTY
login as: venkat
Using keyboard-interactive authentication.
Password:
Last login: Mon Jun 23 11:10:18 2014 from yktgi01e0-s4.wats
* |-----*
* |       IBM Advanced Technical Support (ATS)
* |       System z Applied Technologies Lab
* |       SUSE Linux Enterprise Server 11 SP3
* |
* |       Welcome to the Yorktown Challenge II system
* |-----*
venkat@zpublic:~>
```

Macbook
screen
(terminal)



```
venkatbsmacbook:~ venkat$ ssh venkat@zpublic.wsc.ihost.com
Password:
Last login: Mon Jun 23 11:20:52 2014 from yktgi01e0-s4.watson.ibm.com
* |-----*
* |       IBM Advanced Technical Support (ATS)
* |       System z Applied Technologies Lab
* |       SUSE Linux Enterprise Server 11 SP3
* |
* |       Welcome to the Yorktown Challenge II system
* |-----*
venkat@zpublic:~>
```

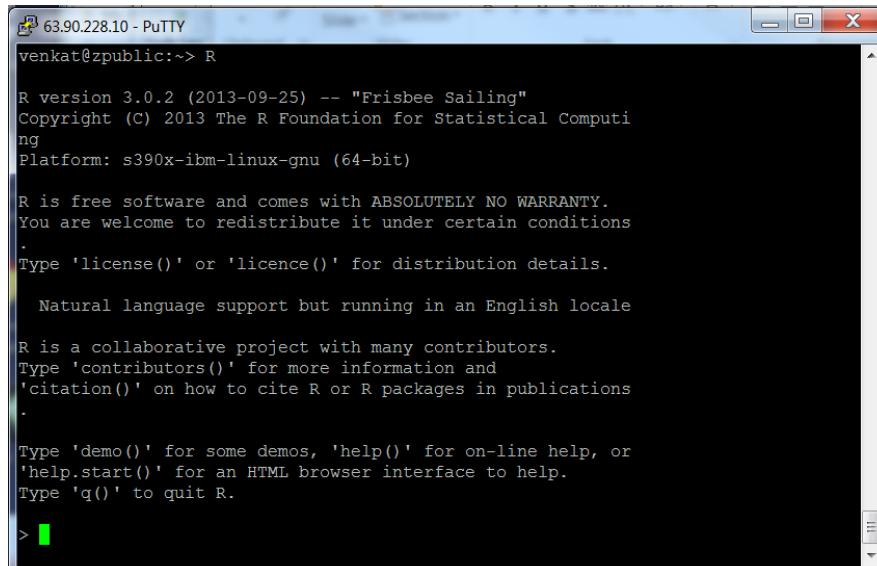
IBM System Z runs in
SUSE Linux OS

Can use ssh command
in other Linux/Unix
systems

Running the software suite R – statistical analysis and graphics tool

R with many of the commonly used packages are installed and running in the server

“R” can be started in the shell, by simply typing “R”



The screenshot shows a PuTTY terminal window titled "63.90.228.10 - PuTTY". The user has typed "R" at the prompt. The terminal displays the R startup message, which includes:

```
venkat@zpublic:~> R
R version 3.0.2 (2013-09-25) -- "Frisbee Sailing"
Copyright (C) 2013 The R Foundation for Statistical Computing
Platform: s390x-ibm-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions
.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications
.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> |
```

Software suite PLINK is also available.
Octave is also available

Other freely downloadable softwares can be installed on request

User account folder path in the server: /mnt/DREAMStorage

6

Challenge data: As the data size is only ~ 50 MB, users can download their own copy !

© 2014 IBM Corporation

Best practices for using the server

Efficient memory and cpu use:

- Jobs submitted by users will be sharing a total memory of 297 GB
- So, if users submit large jobs that take significant amount of memory, say 25 to 30 GB, they are requested to consider efficient use of memory. This will enable multiple users to run many jobs at the same time among the 20 available processors
- Coordinate with Venkat when you need to run a large memory job that will restrict other users. He will **send a note to other users** to review a particular time slot for a job that requires large memory.
- For jobs that take *very long* (possibly many hours) users can **parallelize them to run in more than one processor** wherever possible to have their results *fast*



Contact details for help with problems

For problems related to login, running R packages, installing new R packages, and any other questions related to using the server:

Contact: Venkat Balagurusamy, IBM System Z team

Email: vkbalagu@us.ibm.com

Acknowledgements: Gustavo Stolovitzky, Program Director, CBC (IBM)

Donna Dillenberger, Manager, System Z group (IBM)

DREAM and Sage team

THANK YOU and HAVE FUN SOLVING THE CHALLENGE !