

# Improving Reproducibility in Machine Learning Research

**Benjamin Haibe-Kains**

✉ [bhaibeka@uhnresearch.ca](mailto:bhaibeka@uhnresearch.ca)

🐦 [@bhaibeka](https://twitter.com/bhaibeka)

Senior Scientist, Princess Margaret Cancer Centre

Associate Professor, University of Toronto

Affiliate, OICR and Vector Institute

# Abstract

While biostatistics and machine learning are essential to analyze biomedical data, researchers are facing multiple challenges around research reproducibility and transparency. Given the proliferation of studies investigating the applications of biostatistics and machine learning in research and clinical studies, it is essential for independent researchers to be able to scrutinize and reproduce the results of a study using its materials, and build upon them in future studies.

Computational reproducibility is achievable when the data can easily be shared and the required computational resources are relatively common. However, the complexity of current algorithms and their implementation, the need for specific computer hardware and the use of sensitive biomedical data represent major obstacles in healthy-related research. In this talk, I will describe the various aspects of a typical biomedical study that are necessary for reproducibility and the platforms that exist for sharing these materials with the scientific community.

# Disclosures

- I am a co-Founder of the MAQC (Massive Analysis and Quality Control) Society
- I am part of the Scientific Advisory Board of:
  - Consortium de recherche biopharmaceutique (CQDM), Quebec, Canada
  - Break Through Cancer, Commonwealth Cancer Consortium, United States
  - Canadian Institute of Health Research - Institute of Genetics, Canada
- I am part of the Executive Committee of the Terry Fox Digital Health and Discovery Platform, Canada
- I am part of the Board of Directors of AACR International - Canada, The American Association for Cancer Research, United States
- I am a consultant for Code Ocean Inc, United States

# The Radiomics Heroes @ BHKLAB



Caryn Geady



Mattea Welch



Sejin Kim



Jun Won Kim



Joseph Marsilla



Michal Kazmierski



Reza Reiazi

# What is reproducibility?

## Replication

Scientific claims are confirmed by completely independent investigations

## Reproducibility

Ability of independent investigators to re-create the results claimed by the original investigators using the original samples, assays, data and analysis techniques

## *Computational Reproducibility*

Running the same computer code on the same data to obtain the same results

arXiv.org

**Reproducible Research: A Retrospective**

Roger D. Peng, Stephanie C. Hicks [Submitted on 23 Jul 2020]

# Why does reproducibility matter?

- As data analysis pipelines become increasingly complex, a high level of **transparency** is required to allow the scientific community to properly scrutinize a study
- An analysis may contain many (hyper)parameters, understanding their impact on results is necessary to assess the **robustness** of the claims
- Sharing all study materials maximize the **reusability** of the research outputs (higher impact and return on investment)

# The dark side of AI

nature

## International evaluation of an AI system for breast cancer screening

<https://doi.org/10.1038/s41586-019-1799-6>

Received: 27 July 2019

Accepted: 5 November 2019

Published online: 1 January 2020

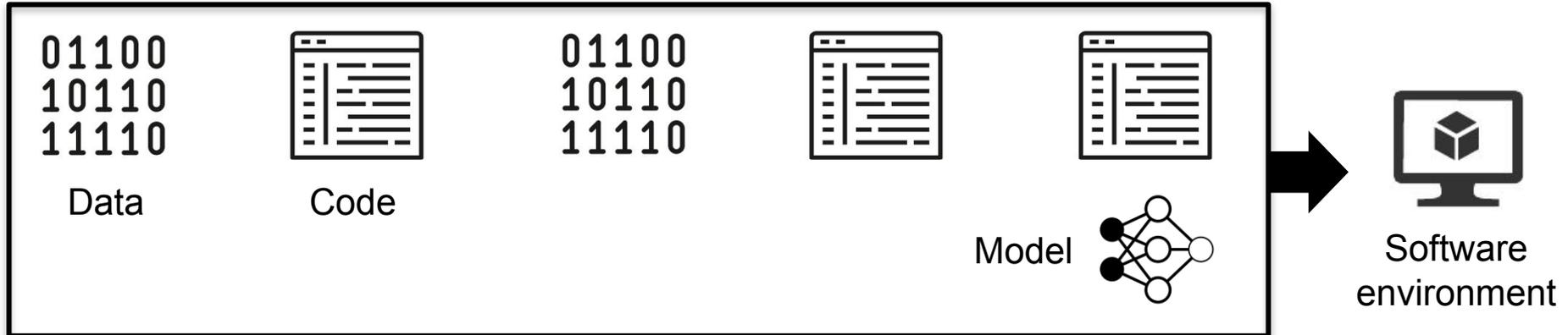
Scott Mayer McKinney<sup>1,3,4\*</sup>, Marcin Sieniek<sup>1,4</sup>, Varun Godbole<sup>1,4</sup>, Jonathan Godwin<sup>2,14</sup>,  
Natalia Antropova<sup>2</sup>, Hutan Ashrafian<sup>3,4</sup>, Trevor Back<sup>2</sup>, Mary Chesus<sup>2</sup>, Greg C. Corrado<sup>1</sup>,  
Ara Darzi<sup>3,4,5</sup>, Mozziyar Etemadi<sup>6</sup>, Florencia Garcia-Vicente<sup>6</sup>, Fiona J. Gilbert<sup>7</sup>,  
Mark Halling-Brown<sup>8</sup>, Demis Hassabis<sup>2</sup>, Sunny Jansen<sup>9</sup>, Alan Karthikesalingam<sup>10</sup>,  
Christopher J. Kelly<sup>10</sup>, Dominic King<sup>10</sup>, Joseph R. Ledsam<sup>2</sup>, David Melnick<sup>8</sup>, Hormuz Mostofi<sup>1</sup>,  
Lily Peng<sup>1</sup>, Joshua Jay Reicher<sup>11</sup>, Bernardino Romera-Paredes<sup>2</sup>, Richard Sidebottom<sup>12,13</sup>,  
Mustafa Suleyman<sup>2</sup>, Daniel Tse<sup>14</sup>, Kenneth C. Young<sup>1</sup>, Jeffrey De Fauw<sup>2,15</sup> & Shravya Shetty<sup>1,15\*</sup>

<sup>1</sup>Google Health, Palo Alto, CA, USA. <sup>2</sup>DeepMind, London, UK. <sup>3</sup>Department of Surgery and Cancer, Imperial College London, London, UK. <sup>4</sup>Institute of Global Health Innovation, Imperial College London, London, UK. <sup>5</sup>Cancer Research UK Imperial Centre, Imperial College London, London, UK. <sup>6</sup>Northwestern Medicine, Chicago, IL, USA. <sup>7</sup>Department of Radiology, Cambridge Biomedical Research Centre, University of Cambridge, Cambridge, UK. <sup>8</sup>Royal Surrey County Hospital, Guildford, UK. <sup>9</sup>Verily Life Sciences, South San Francisco, CA, USA. <sup>10</sup>Google Health, London, UK. <sup>11</sup>Stanford Health Care and Palo Alto Veterans Affairs, Palo Alto, CA, USA. <sup>12</sup>The Royal Marsden Hospital, London, UK. <sup>13</sup>Thirlestaine Breast Centre, Cheltenham, UK. <sup>14</sup>These authors contributed equally: Scott Mayer McKinney, Marcin T. Sieniek, Varun Godbole, Jonathan Godwin. <sup>15</sup>These authors jointly supervised this work: Jeffrey De Fauw, Shravya Shetty. \*e-mail: scottmayer@google.com; tsed@google.com; sshetty@google.com

### Code availability

The code used for training the models has a large number of dependencies on internal tooling, infrastructure and hardware, and its release is therefore not feasible. However, all experiments and implementation details are described in sufficient detail in the Supplementary Methods section to support replication with non-proprietary libraries. Several major components of our work are available in open source repositories: Tensorflow (<https://www.tensorflow.org>); Tensorflow Object Detection API ([https://github.com/tensorflow/models/tree/master/research/object\\_detection](https://github.com/tensorflow/models/tree/master/research/object_detection)).

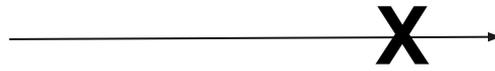
# Key elements of a typical AI pipeline



# Sharing raw and processed data

- For data that are not sensitive and can be released publicly
  - Large files: established repositories specific to each field/domain, such as the European Genotype-phenotype Archive (*EGA*) for genomics
  - Small files: immutable data repositories assigning unique identifiers such as *Zenodo*, *Harvard Dataverse*, *Dryad*, or *figshare*
- For data that are private and/or sensitive (eg, patient health information)
  - Set up a Data Access Committee (DAC) and a transparent process to review requests (with timeline)
  - Use established repositories that can forward requests to DAC and grant access upon approval, such as *EGA* or *dbGaP*
  - Ensure that the data transfer is secure and automated (manual transfer of data on hard disk and shipping by mail is not ideal...)

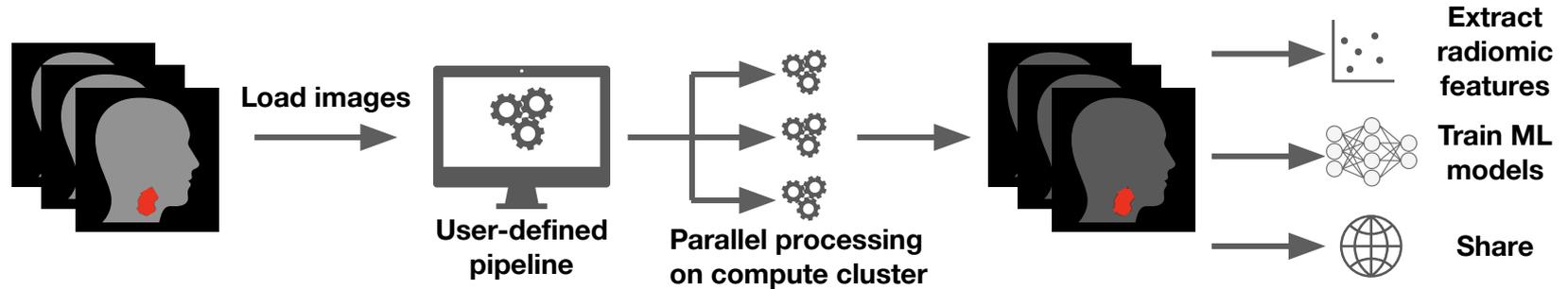
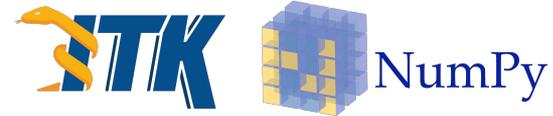
# Raw data ≠ analysis-ready



# In the lab: Imgtools for processing large imaging datasets

[github.com/bhklab/imgtools](https://github.com/bhklab/imgtools)

- Built in Python using tested open-source components
- Rich set of standardised image processing operations
- Easy sharing of pipeline specifications
- Support for DICOM (including RT Structure Set)



# In the lab: Automating data processing and versioning

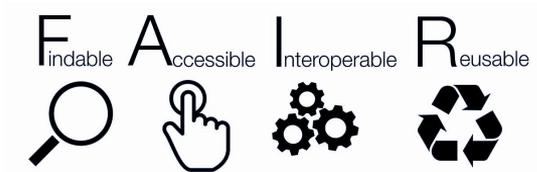
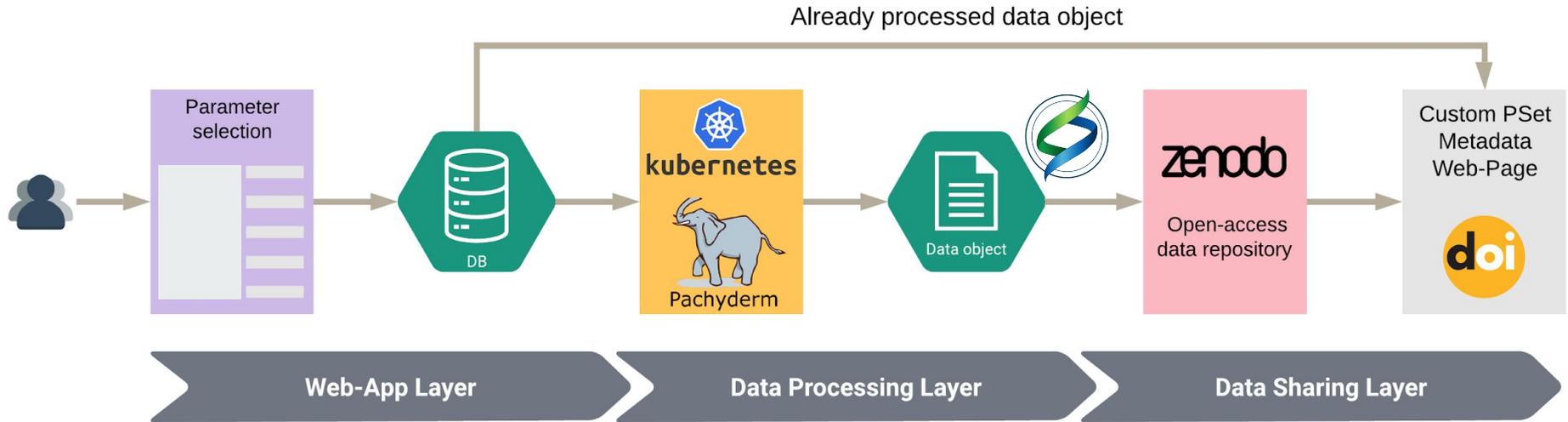
[orcestra.ca](http://orcestra.ca)

# ORCESTR

Orchestration platform for reproducing multimodal data

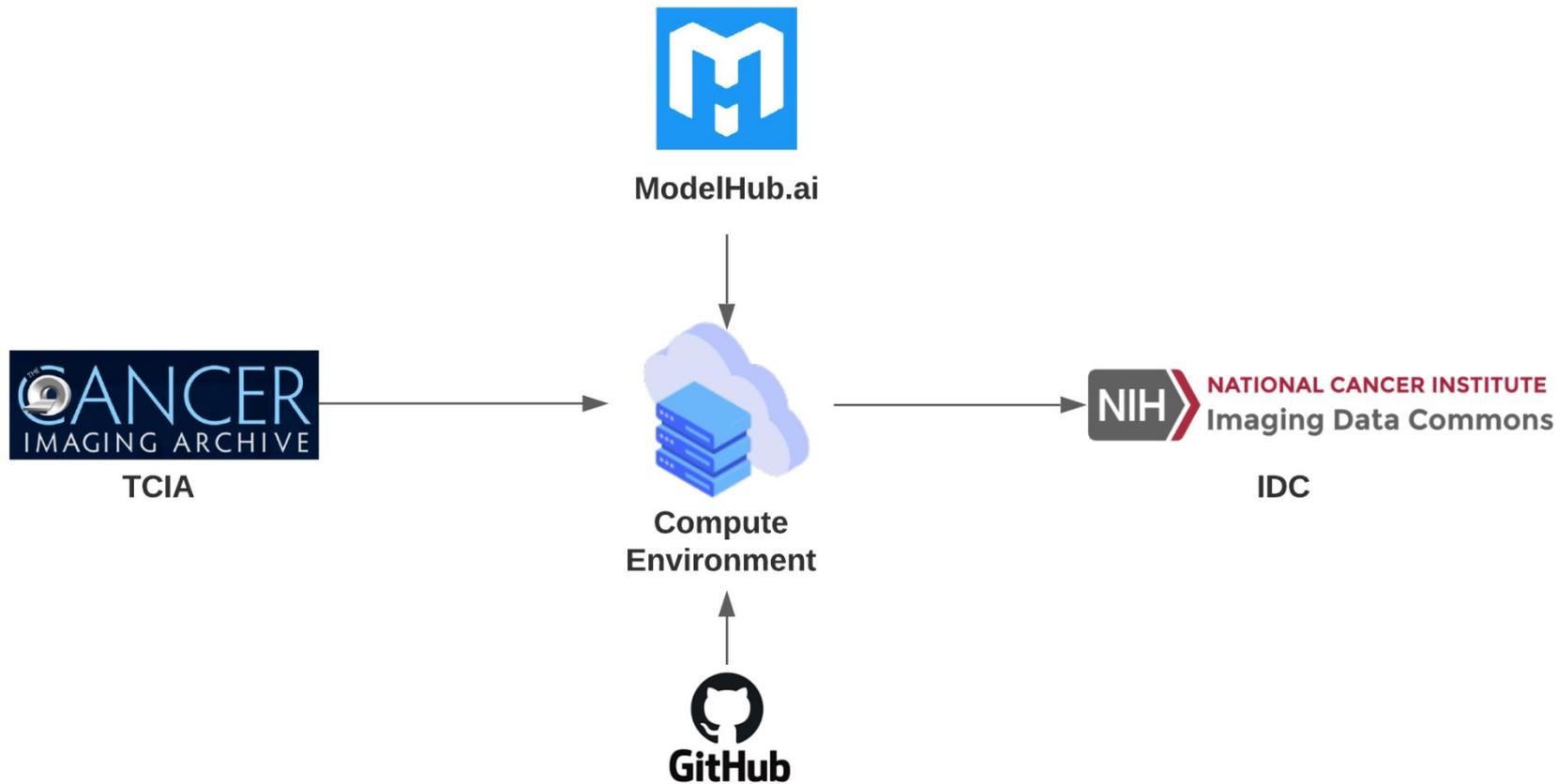
Pharmacogenomics Data 	Toxicogenomics Data 	Xenographic Pharmacogenomics Data 
Radiogenomics Data 	Clinical Genomics Data 	Radiomics Data 

# In the lab: Automating data processing and versioning



→ Maximize cross-dataset operability by standardizing processing pipelines and annotations (meta and clinical data)

# For the community



# Sharing raw and processed data

For data that are private and/or sensitive (eg, patient health information)

- Set up a Data Access Committee (DAC) and a transparent process to review requests (with timeline)
- Use established repositories that can forward requests to DAC and grant access upon approval, such as *EGA* or *dbGaP* for genomics
- Ensure that the data transfer is secure and automated (manual transfer of data on hard disk and shipping by mail is not ideal...)

**TCIA** is using a different approach

- Most studies are open, few are protected (eg, head and neck due to potential face reconstruction)
- Requesters signed a Data Usage Agreement (DUA) with TCIA
- TCIA takes on the liability
- Identity of requesters can be communicated to the data generators

# Sharing computer code

- Three levels of computer code:
  - Preprocessing of the raw data, including curation of metadata, quality control and normalization
  - Training of the predictive model (fitting of the model, tuning, testing)
  - Validation of the final (fully specified) predictive model on independent data
- Code can easily be shared on established versioning systems, like *GitHub* or *BitBucket*
- Multiple frameworks exists for deep learning: *PyTorch*, *TensorFlow*
- Higher level frameworks like *Keras* and *PyTorch Lightning* allow for better organization of the code, making it easier to read
- Documentation is a must!
- Think about your license: open-source? Non-commercial?

# Sharing predictive models

- Final models can also be shared to facilitate application on new data
- Platforms exist to easily share models, their architecture and full set of parameters
  - *ModelHub, ModelDepot, Model Zoo*
- An example of code goes a long way to showcase how the final model can be used
  - Describe how to format/process the input data, how to interpret the outputs (predictions), and the limitations of the model
- A model may reveal key aspects of your dataset, and be used to reveal the identity of some participants
  - Differential privacy can sometimes be ensured
  - Otherwise share your model via a data access committee

# The power of virtualization

- Having access to data and code is not sufficient to reproduce the results of a study
- You also need to set up the software environment
  - Operating system, software libraries (versions?)
  - Even the type of hardware matters (eg, numerical precision)
- Go **virtual!** (I mean *virtual machines* and *containers*...)
  - Specify all the required libraries and their versions
  - Once the container is built, add your code and data
  - Complete a successful run, check your results and... voilà!
- Users can now easily **reproduce** all your results, **scrutinize** and **challenge** them
- Many platforms exist to help build and share containers in the cloud:  
*Code Ocean, Gigantum, Docker Hub, ...*

# Levels of reproducibility - Level 1



01100  
10110  
11110

Raw  
data



Code for  
QC and  
preprocessing

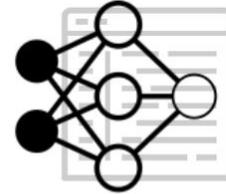


01100  
10110  
11110

Processed  
data



Code for  
model  
training



Final model  
and code

**BONUS**



Software  
environment

- All elements are shared, the study is fully reproducible
- Providing a container will **really** help the community build upon your study

# Levels of reproducibility - Level 2



01100  
10110  
11110

Raw  
data



Code for  
QC and  
preprocessing

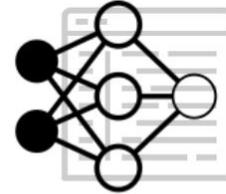


01100  
10110  
11110

Processed  
data



Code for  
model  
training



Final model  
and code

**BONUS**



Software  
environment

- Raw data could not be shared but processed data available
- If the raw data and processing are well documented, it is *OK*
- But new processing approaches or meta-analyses might not be feasible

# Levels of reproducibility - Level 3



01100  
10110  
11110

Raw  
data



Code for  
QC and  
preprocessing

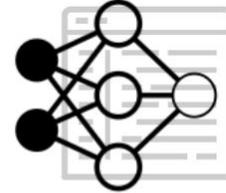


01100  
10110  
11110

Processed  
data



Code for  
model  
training



Final model  
and code



Software  
environment

- The study is NOT reproducible
- But the application of the model in new datasets can still be explored

# Levels of reproducibility - Level UNDEF



Raw  
data



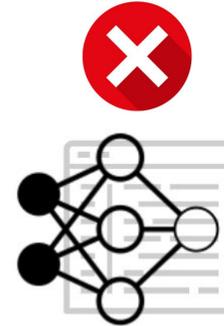
Code for  
QC and  
preprocessing



Processed  
data



Code for  
model  
training



Final model  
and code



Software  
environment

- Advertisement for a cool but closed technology

# Initial model vs derivatives

- As machine learning models become more complex (eg, deep neural networks), more computational resources are required
  - Large amounts of CPUs, GPUs or memory
  - Specific hardware that are expensive and/or hard to find (even proprietary technologies)
- It is only a matter of time until the technology becomes cheaper and more accessible
- Meantime, further research can be done in **simplifying** the initial model
  - Decrease the dependency on specific hardware
  - Clever updates in how the model is trained or its architecture can go a long way
  - Broaden the user base, increase impact



# Research Reproducibility in Practice

# Tools to improve reproducibility - Code Ocean

- Sharing data
  - Public data: Zenodo, Harvard Dataverse, Dryad, figshare, ...
  - Private/sensitive data: dbGaP, NCBI SRA, EBI EGA, ...
- Sharing computer code
  - GitHub, GitLab, Bitbucket, ...
- Sharing software environment
  - **Code Ocean**, Gigantum, Docker Hub, ...
- Sharing computational models
  - ModelHub, ModelDepot, Model Zoo, ...
- Sharing analysis results
  - Zenodo, Harvard Dataverse, Dryad, figshare, ...



# Code Ocean interface



Published Integrative Pharmacogenomics Analysis of Patient Derived Xenografts (Arvind Singh Mer & Benjamin Haibe-Kains)

Capsule File Help

Edit Original



Files  
Tabs

Core Files	
▶ metadata	2.69 KB ✓
▶ environment	1.09 KB ✓
▼ code	77.28 KB ✓
R figure_1.R	759 B ✓
R figure_2.R	4.25 KB ✓
R figure_3.R	6.32 KB ✓
R figure_4.R	3.08 KB ✓
R figure_5.R	1.62 KB ✓
R figure_6.R	6.67 KB ✓
R figure_7.R	6.74 KB ✓
R Gene_Drug_association.R	4.1 KB ✓
R imp_functions.R	7.38 KB ✓
R install_xeva.R	756 B ✓
LICENSE	34.32 KB ✓
README.md	808 B ✓
run.sh	499 B ✓
▼ data Manage Datasets	84.59 MB
c2.cp.reactome.v6.0.symb...	322.98 KB
drug_pathway_assco.Rda	131.72 KB
ICC_for_genes.Rda	3.38 MB
icc_gsea.Rda	25.31 KB
LICENSE	18.88 KB
passage_data.Rda	80.63 MB
PDX_passage_correlation.Rda	5 KB
pdxe_BRCA_gene_drug_ass...	81.81 KB
PDXE_Drug_targetGenes_List...	1.4 KB
.gitignore	7 B ✓

Results  
Your files will appear in the timeline.  
[View latest results](#)

## Integrative Pharmacogenomics Analysis of Patient Derived Xenografts

[View Raw](#)

This capsule contains all of the analysis conducted for the manuscript *Integrative Pharmacogenomics Analysis of Patient Derived Xenografts*.

The Biorxiv preprint of the manuscript can be found [here](#).

The link to the GitHub repository for the Xeva tool described in the manuscript can be found [here](#).

### Code, Data, Results

- All associated code is provided in the code directory.
- The data directory contains all necessary data.
- Output plots will be saved in the results directory.

### How to replicate analysis

- All scripts can be executed using `run.sh` or directly using `Rscript`.
- Necessary packages can be installed using `install_xeva.R`.

### Reproducible Run

or launch a cloud workstation

lab Studio jupyter ↩️ Shiny

Timeline

- Dec 31, 2018  
Published Version 1.0  
Currently viewing
- Benjamin Haibe-Kains [committed](#)  
Dec 31, 2018  
Version 1.0
- Author ran Dec 31, 2018 00:01:16
- ▼ Published Result
  - figure-1.pdf 137.29 KB
  - figure-2.pdf 24.54 KB
  - figure-3.pdf 68.62 KB
  - figure-4.pdf 5.18 KB
  - figure-5.pdf 6.89 KB
  - figure-6.pdf 45.11 KB
  - figure-7.pdf 67.56 KB
  - Output 674 B
- Jun 4, 2018  
Created capsule

Reproducibility



# How do we use Code Ocean in my lab?

- **Over the course of the analysis**
  - Implementation of a capsule to run all the analyses from the beginning
  - Good tool for collaboration as anyone is using the same software environment
- **Internal review prior to the manuscript submission**
  - Independent review of the code and rerun of all the analyses from the raw or low-level processed data
- **At the time of the manuscript submission**
  - Strengthen the study, show the reviewers that the research is transparent and increase its impact (reusability)
- **After publication**
  - For older publications or collaborative papers with time-sensitive submission
  - A great way to update the results when new data are released, improved methods are available or bugs are found

# Continuous development

Using **Code Ocean** during the implementation of an analysis pipelines has multiple benefits:

- Quickly install software packages
- A well-controlled software environment facilitates the debugging process
  - You can literally share the capsule with a colleague to get help
- Optimize collaboration as all parties use the exact same data and software environment
  - Co-editing of code can be challenging
- Running code many times for large analyses may require large computational resources
  - **Co-Exist** (beta)
- It ensures full reproducibility by construction

# Pre-publication review

The analysis is done, results looking good!

- Time to clean and document the code → **Build a Code Ocean capsule!**
- Paper ready to submit but ... let's have an independent review first
- Another lab member, not involved or only peripherally involved in the study checks the code
- And rerun everything from the raw data or low-level processed data
  
- **Surprising how often the results change!**

# Post-publication review

- Your analysis yields unexpected results (expect a fair share of scrutiny)
- You work hard on making your study fully reproducible and prepare the data and computer code for public release
- Your paper is finally accepted! Let the fun begin ...



# Adoption in the lab

- 16 publications with a Code Ocean capsule since 2017
- A capsule is mandatory for **all** new studies led by the lab

## Challenges:

- Private data?
- Proprietary software?
- Need for expensive computational resources?

The screenshot shows the Code Ocean website interface. At the top, there is a navigation bar with the Code Ocean logo, 'EXPLORE', 'HELP', a 'New Capsule' button, and a 'Sign Up' link. Below the navigation bar is a search bar containing the text 'haibe-kains'. The search results are displayed in a grid format. The first result is titled 'MEDICAL SCIENCES April 2017' and 'PharmacGx' by Benjamin Haibe-Kains. It includes a scatter plot showing a positive correlation between two variables, with a regression line and statistical values:  $PCQ=0.97$ ,  $p=4.51451548901003e-23$ , and  $R^2=0.9522$ . The second result is titled 'BIOINFORMATICS November 2018' and 'MetaGxData: Clinically Annotated Breast, Ovar...' by Benjamin Haibe-Kains & Michael Zon. It includes a heatmap and a line graph. The third result is titled 'BIOINFORMATICS July 2020' and 'SYNERGxDB NAR Web Server' by Heewon Seo & Benjamin Haibe-Kains. It includes a logo for SYNERGxDB and a line graph. Each result has a 'Research' icon and a tag indicating the field of study.

# Platforms we used in the lab

- **R** or **Python** for the analyses
  - R packages on **CRAN / Bioconductor**
  - Python packages on **Pypi / Anaconda**
- **PyTorch Lightning** for deep learning
- **GitHub** for code versioning
- **Code Ocean** to share code and software environment
- **Zenodo** for sharing specific research outputs (processed data, results,...)
- European Genotype-phenotype Archive (**EGA**) and The Cancer Imaging Archive (**TCIA**) for sharing large genomics and imaging data, respectively
- **ORCESTRA** to process, version and share the data \*
- **ModelHub** to share deep neural networks

# Lessons learned

- **Making your research fully reproducible takes time and expertise**
  - Include these costs in the budget of your grant proposals
- **Put research reproducibility at the center of your lab's modus operandi and mission**
  - Discuss research reproducibility during interviews
  - Set expectations for new recruits
  - Make it mandatory, no leeway (exceptions must be strongly justified)
- **Resist the pressure!**
  - We are all in a hurry to publish, competition is tough ...
  - Editors and Reviewers appreciate transparency
  - It is an investment, addressing the reviewer's comments and building upon previous results will be greatly facilitated

# Lessons learned

- Knowing that data/code and software environment will be shared is a **strong motivation to improve the quality of the research itself**
  - Containers facilitate the pre- and post-publication review
- Containers make it easier to **collaborate**
  - All participants use exactly the same software environment
  - Easier to check whether the code runs properly after an update
- Sharing research outputs **increases impact**, opportunities for new collaborations
- Open Science and research reproducibility **do not prevent commercialization**



PNAS

February 10, 2015 112 (6) 1645-1646

**Opinion: Reproducible research can still be wrong: Adopting a prevention approach**

Jeffrey T. Leek and Roger D. Peng

# Transparency and reproducibility in artificial intelligence

Benjamin Haibe-Kains<sup>1,2,3,4,5</sup>✉, George Alexandru Adam<sup>3,5</sup>, Ahmed Hosny<sup>6,7</sup>, Farnoosh Khodakarami<sup>1,2</sup>, Massive Analysis Quality Control (MAQC) Society Board of Directors\*, Levi Waldron<sup>8</sup>, Bo Wang<sup>2,3,5,9,10</sup>, Chris McIntosh<sup>2,5,9</sup>, Anna Goldenberg<sup>3,5,11,12</sup>, Anshul Kundaje<sup>13,14</sup>, Casey S. Greene<sup>15,16</sup>, Tamara Broderick<sup>17</sup>, Michael M. Hoffman<sup>1,2,3,5</sup>, Jeffrey T. Leek<sup>18</sup>, Keegan Korthauer<sup>19,20</sup>, Wolfgang Huber<sup>21</sup>, Alvis Brazma<sup>22</sup>, Joelle Pineau<sup>23,24</sup>, Robert Tibshirani<sup>25,26</sup>, Trevor Hastie<sup>25,26</sup>, John P. A. Ioannidis<sup>25,26,27,28,29</sup>, John Quackenbush<sup>30,31,32</sup> & Hugo J. W. L. Aerts<sup>6,7,33,34</sup>

Published online: 14 October 2020

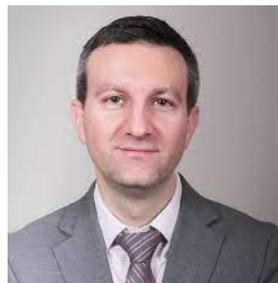
# Acknowledgements



Andrew Hope



Scott Bratman



Tom Purdie



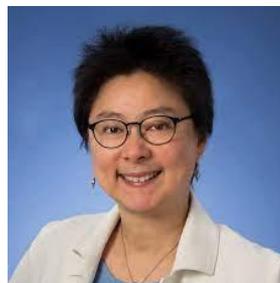
Chris McIntosh



Alejandro Berlin



Michael Milosevic



Fei-Fei Liu



Hugo Aerts



Benjamin Kann



**Thank you for your attention**

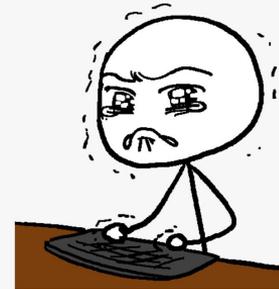
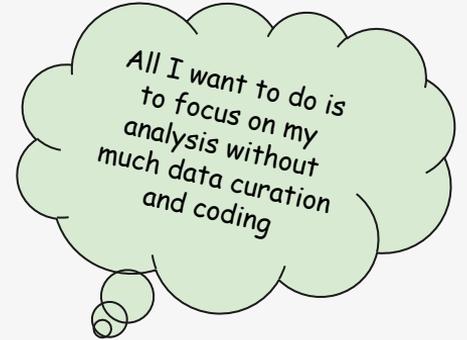
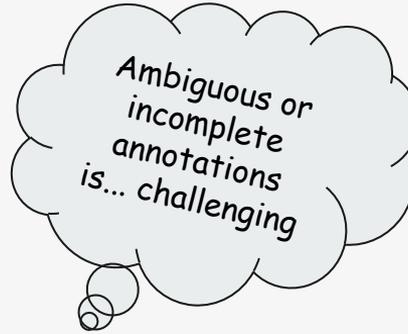
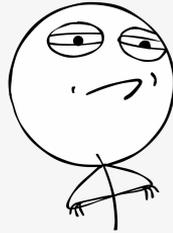
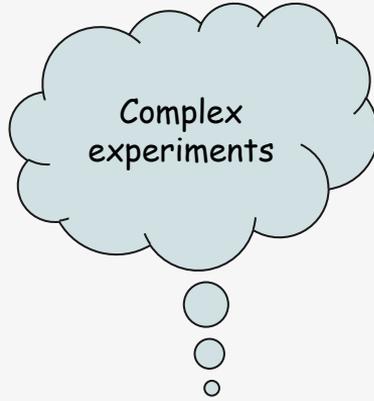
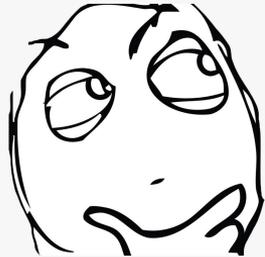
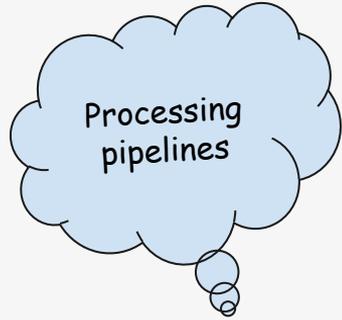


**Thank you for your attention**



# APPENDIX

# Challenges



# The reality of many biomedical labs

## Today

- No coding expertise, use of tools like Excel and PRISM to analyze the data
  - Prone to errors, untraceable
  - **Not reproducible**
- Lack of data coordination
- No software development or difficult to maintain
- Journal guidelines typically not properly enforced
  - Easy to get away from sharing data or code

## Tomorrow?

- Funding agencies support data coordination and research reproducibility
- Journals require and assess research reproducibility
- Use of open programming languages to develop analysis pipelines
  - Documentation
  - Automated/reproducible
- Data are annotated and shared internally and externally
- Data and code grow/improve over time

# How do we use Code Ocean in my lab?

- **Over the course of the analysis**
  - Implementation of a capsule to run all the analyses from the beginning
  - Good tool for collaboration as anyone is using the same software environment
- **Internal review prior to the manuscript submission**
  - Independent review of the code and rerun of all the analyses from the raw or low-level processed data
- **At the time of the manuscript submission**
  - Strengthen the study, show the reviewers that the research is transparent and increase its impact (reusability)
- **After publication**
  - For older publications or collaborative papers with time-sensitive submission
  - A great way to update the results when new data are released, improved methods are available or bugs are found

# Continuous development

Using **Code Ocean** during the implementation of an analysis pipelines has multiple benefits:

- Quickly install software packages
- A well-controlled software environment facilitates the debugging process
  - You can literally share the capsule with a colleague to get help
- Optimize collaboration as all parties use the exact same data and software environment
  - Co-editing of code can be challenging
- Running code many times for large analyses may require large computational resources
  - **Co-Exist** (beta)
- It ensures full reproducibility by construction

# Pre-publication review

The analysis is done, results looking good!

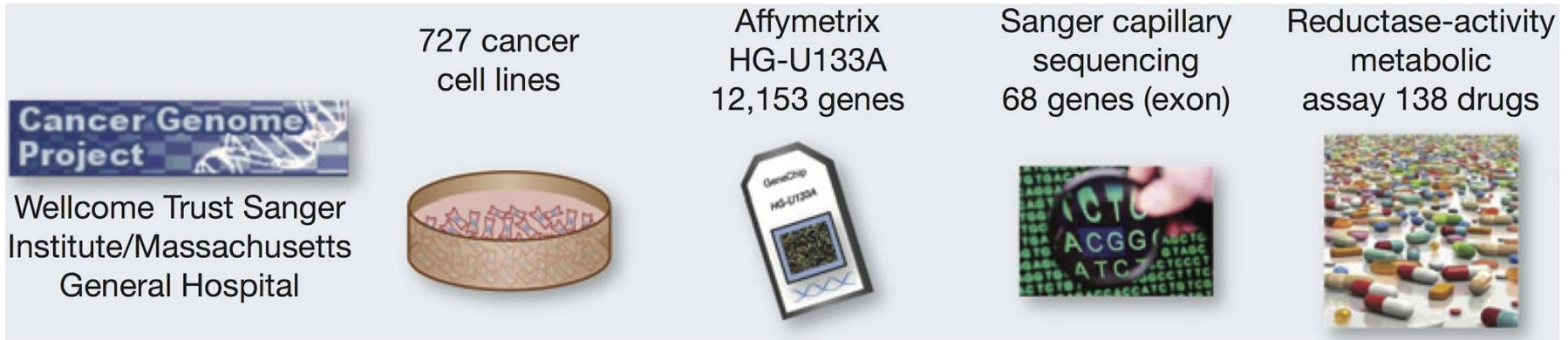
- Time to clean and document the code → **Build a Code Ocean capsule!**
- Paper ready to submit but ... let's have an independent review first
- Another lab member, not involved or only peripherally involved in the study checks the code
- And rerun everything from the raw data or low-level processed data
- **Surprising how often the results change!**

# Post-publication review

- Your analysis yields unexpected results (expect a fair share of scrutiny)
- You work hard on making your study fully reproducible and prepare the data and computer code for public release
- Your paper is finally accepted! Let the fun begin ...



# Genomics of Drug Sensitivity in Cancer (GDSC)



*“By linking drug activity to the functional complexity of cancer genomes, systematic pharmacogenomic profiling in cancer cell lines provides a powerful biomarker discovery platform to guide rational cancer therapeutic strategies”*

# Cancer Cell Line Encyclopedia (CCLE)



Broad Institute/  
Novartis



1,036 cancer  
cell lines



Affymetrix HG-  
U133PLUS2  
19,173 genes



Hybrid capture  
sequencing  
1,651 genes (exon)



ATP-based  
metabolic assay  
24 drugs

*“The generation of genetic predictions of drug response in the preclinical setting and their incorporation into cancer clinical trial design could speed the emergence of ‘personalized’ therapeutic regimens.”*

# Comparing large pharmacogenomic studies

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>

Published: 27 November 2013

RESEARCH ARTICLE

## Revisiting inconsistency in large pharmacogenomic studies

Zhaleh Safikhani<sup>1,2</sup>, Petr Smirnov<sup>1</sup>, Mark Freeman<sup>1</sup>, Nehme El-Hachem<sup>3</sup>, Adrian She<sup>1</sup>, Quevedo Rene<sup>1,2</sup>, Anna Goldenberg<sup>4,5</sup>, Nicolai J. Birkbak<sup>6</sup>, Christos Hatzis<sup>7,8</sup>, Leming Shi<sup>9,10</sup>, Andrew H. Beck<sup>11</sup>, Hugo J.W.L. Aerts<sup>13,14</sup>, John Quackenbush<sup>12,14</sup>,  Benjamin Haibe-Kains<sup>1,2,5,15</sup>

Version 1 16 Sep 16

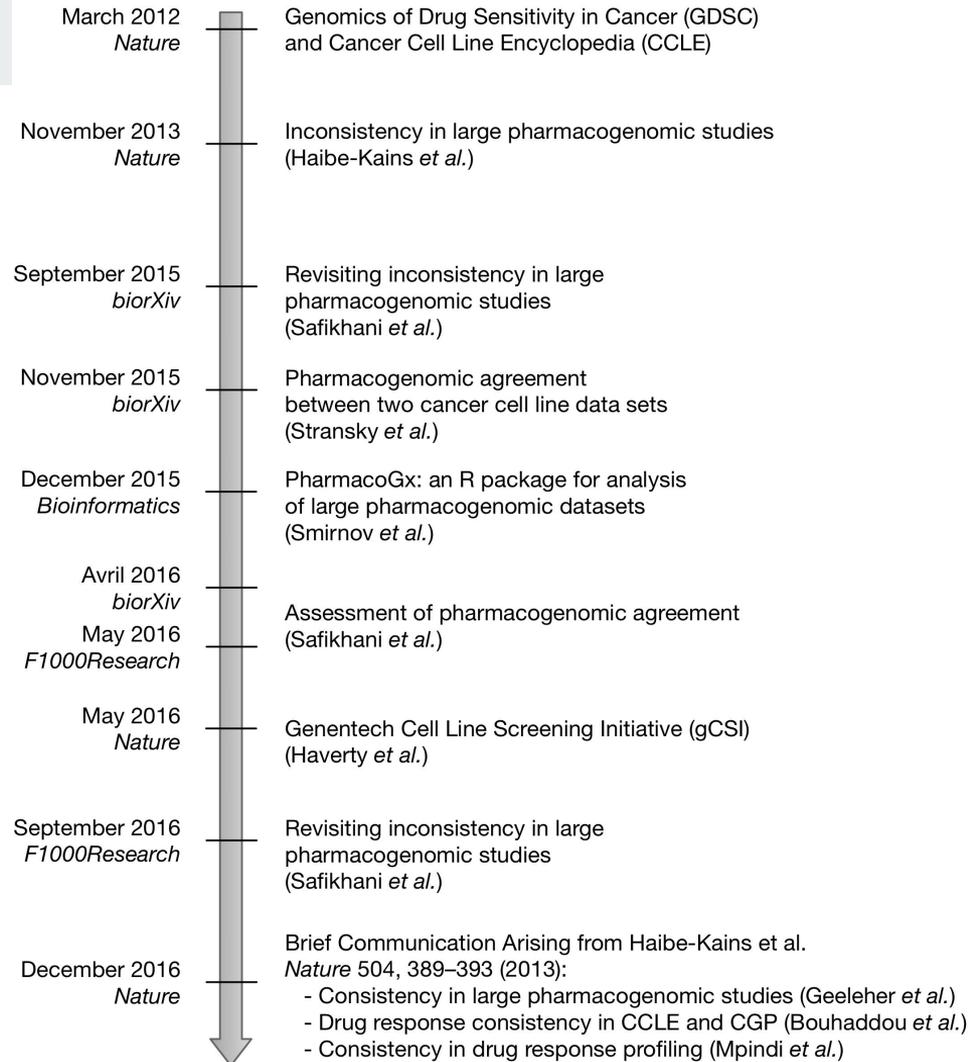
F1000Research  
Open for Science

 Open in Code Ocean



# An epic controversy

- Transparency (sharing of curated data, computer code and all results) helped us address critics
- Reproducibility and high-quality software helped us update, adapt and improve the analysis quickly
- The shared data and code triggered new research in the field



# Adoption in the lab

- 16 publications with a Code Ocean capsule since 2017
- A capsule is mandatory for **all** new studies led by the lab

## Challenges:

- Private data?
- Proprietary software?
- Need for expensive computational resources?

The screenshot displays the Code Ocean interface with a search bar containing 'haibe-kains'. Three search results are visible:

- MEDICAL SCIENCES April 2017**  
**PharmacGx**  
Benjamin Haibe-Kains  
This script uses the PharmacGx package to download the 2013 version of the Genomics of Drug Sensitivity in...  
*Associated article published in F1000Research on April 2016*  
Tags: f1000
- BIOINFORMATICS November 2018**  
**MetaGxData: Clinically Annotated Breast, Ovar...**  
Benjamin Haibe-Kains & Michael Zon  
A wealth of transcriptomic and clinical data on solid tumours are under-utilized due to unharmonized data storage and format. We have developed the MetaGxData package compendium, which includes manually-curated and standardized clinical, pathologic...  
Tags: gene-expression
- BIOINFORMATICS July 2020**  
**SYNERGxDB NAR Web Server**  
Heewon Seo & Benjamin Haibe-Kains  
Drug-combination data portals have recently been introduced to mine huge amounts of pharmacological ...  
*Associated article published in Nucleic Acids Research on May 2020*  
Tags: cancer, pharmacogenomics, drug

The SYNERGxDB logo is also visible at the bottom left of the search results area.

# Back to your own lab

1. **Restrict the use of Excel** and similar tools to *data visualization*
  - A computer program (code) is a more efficient way to transform and analyze data
2. **Nominate a lab member** to look into the various platforms to increase transparency and reproducibility. For instance:
  - R or Python for the analyses
  - GitHub for code versioning
  - Code Ocean to share code and software environment
  - Zenodo for sharing specific research outputs (processed data, results, ...)
3. **Implement a test run** with a published study
  - Can you reproduce the results?
  - Share the data, code, capsule, and results, enjoy the increased visibility/impact
4. **Build the expertise**, try with a new study to be submitted



**ORCESTRA**

# Research reproducibility is also relevant for data processing



Login/Register

Pachyderm is  
offline

## ORCESTRA

Orchestration platform for reproducing multimodal data

Pharmacogenomics Data



Toxicogenomics Data



Xenographic  
Pharmacogenomics Data



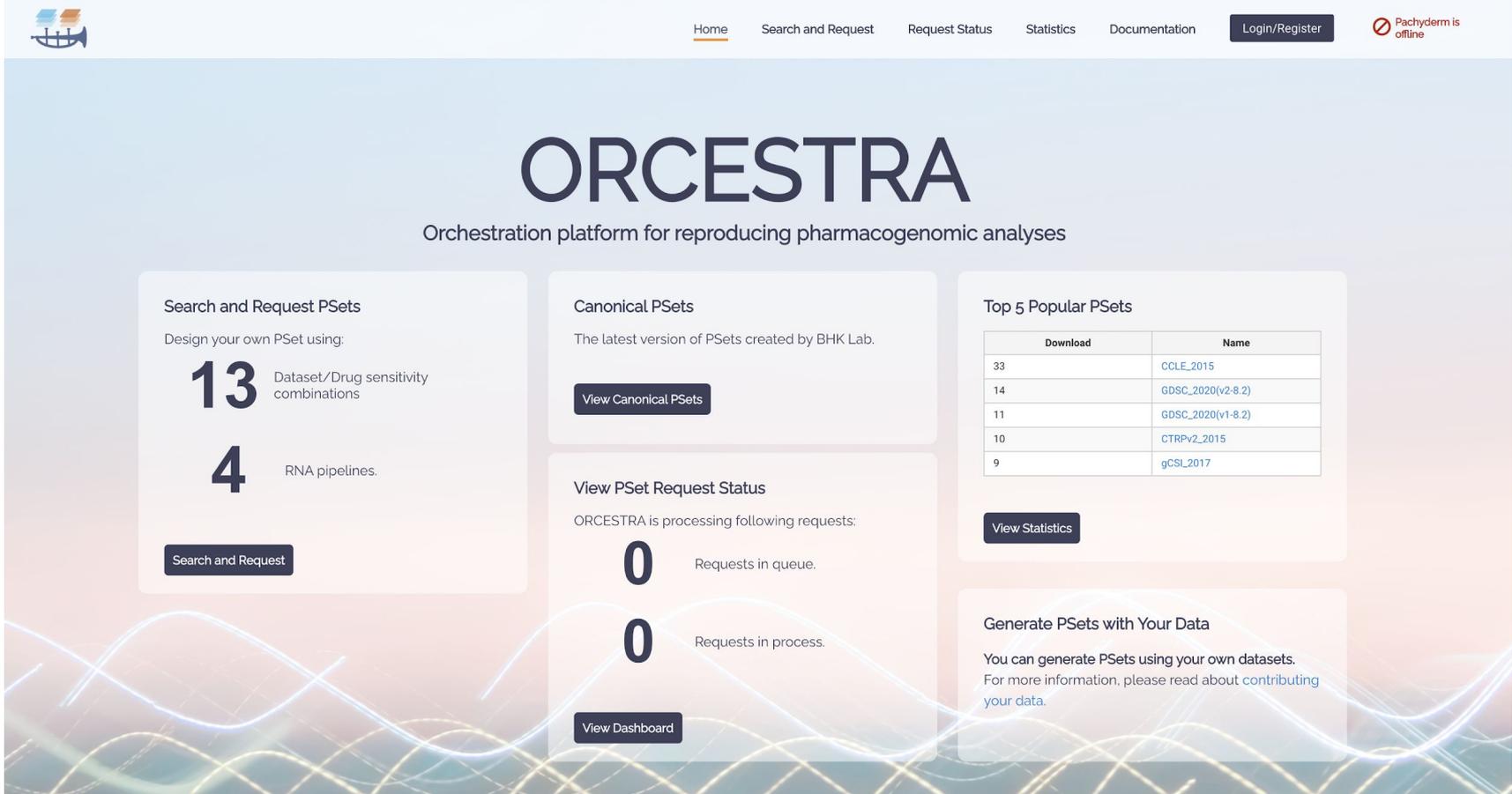
Clinical Genomics Data

Coming soon

Radiogenomics Data



# Research reproducibility is also relevant for data processing

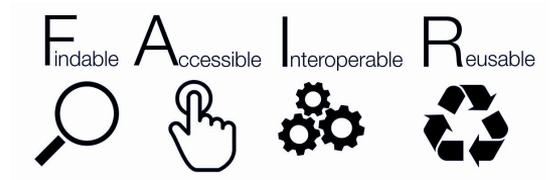
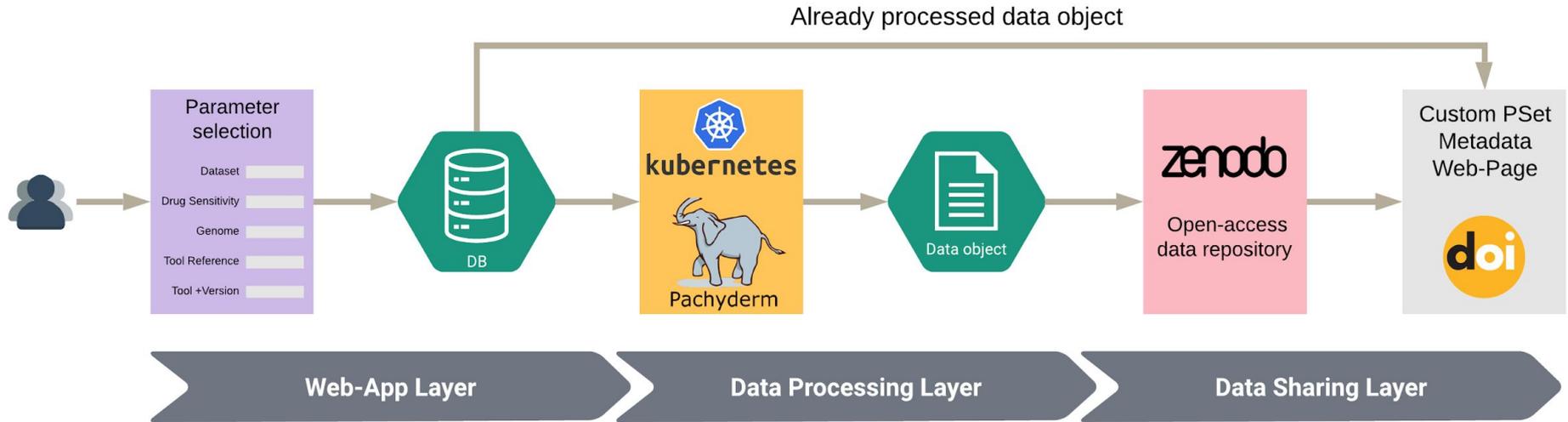


The image shows the ORCESTRA website interface. At the top, there is a navigation bar with links for Home, Search and Request, Request Status, Statistics, and Documentation. A Login/Register button and a status indicator for Pachyderm (offline) are also present. The main heading is 'ORCESTRA' with the subtitle 'Orchestration platform for reproducing pharmacogenomic analyses'. Below this, there are several key features and statistics:

- Search and Request PSets:** Design your own PSet using:
  - 13 Dataset/Drug sensitivity combinations
  - 4 RNA pipelines.A 'Search and Request' button is provided.
- Canonical PSets:** The latest version of PSets created by BHK Lab. A 'View Canonical PSets' button is available.
- View PSet Request Status:** ORCESTRA is processing following requests:
  - 0 Requests in queue.
  - 0 Requests in process.A 'View Dashboard' button is provided.
- Top 5 Popular PSets:** A table listing the most popular PSets.
- Generate PSets with Your Data:** A section explaining that users can generate PSets using their own datasets, with a link to learn more about contributing data.

Download	Name
33	CCLE_2015
14	GDSC_2020(v2-8.2)
11	GDSC_2020(v1-8.2)
10	CTRPv2_2015
9	gCSL2017

# Reproducible data processing: ORCESTRA workflow



→ Maximize cross-dataset operability by standardizing processing pipelines and compound annotations

# ORCESTRAs for pharmacogenomics

