



IMPROVER

SYSTEMS BIOLOGY VERIFICATION

www.sbvimprover.com

The Systems Toxicology Computational Challenge:

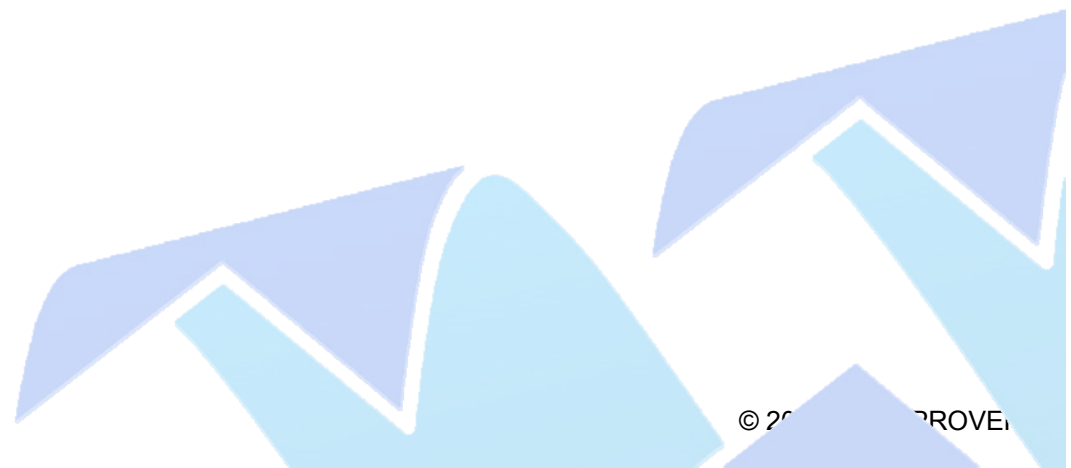
*Marker of Exposure Response
Identification*

Overview

Carine Poussin, PhD

July 11th 2016, ISMB, Orlando

sbv IMPROVER in a Nutshell



IMPROVER: Industrial Methodology for Process Verification in Research



Project initiated 6 years ago and funded by Philip Morris International R&D

Aims to verify methods & data in systems biology / toxicology

Verifies individual methods using double blind performance assessment

BIOINFORMATICS

REVIEW

Vol. 28 no. 9 2012, pages 1193–1201
doi:10.1093/bioinformatics/bts116

Systems biology

Advance Access publication March 14, 2012

Industrial methodology for process verification in research (IMPROVER): toward systems biology verification

Pablo Meyer^{1,†}, Julia Hoeng^{2,†}, J. Jeremy Rice^{1,†}, Raquel Norel¹, Jörg Sprengel³, Katrin Stolle², Thomas Bonk², Stephanie Corthesy³, Ajay Royyuru^{1,*}, Manuel C. Peitsch^{2,*} and Gustavo Stolovitzky^{1,*}

¹IBM Computational Biology Center, Yorktown Heights, 10598 NY, USA, ²Philip Morris Products SA, Research and Development, 2000, Neuchâtel, Switzerland and ³IBM Life Sciences Division, 8802, Zurich, Switzerland

Bioinformatics 2012 28(9):1193-1201

_computational
BIOLOGY

COMMENTARY

Verification of systems biology research in the age of collaborative competition

Pablo Meyer¹, Leonidas G Alexopoulos², Thomas Bonk³, Andrea Califano⁴, Carolyn R Cho⁵, Alberto de la Fuente⁶, David de Graaf⁷, Alexander J Hartemink⁸, Julia Hoeng³, Nikolai V Ivanov³, Heinz Koeppel⁹, Rune Linding¹⁰, Daniel Marbach¹¹, Raquel Norel¹, Manuel C Peitsch³, J Jeremy Rice¹, Ajay Royyuru¹, Frank Schacherer¹², Joerg Sprengel¹³, Katrin Stolle³, Dennis Vitkup⁴ & Gustavo Stolovitzky¹

Collaborative competitions in which communities of researchers compete to solve challenges may facilitate more rigorous scrutiny of scientific results.

Nature Biotechnology 2011 Sep 8;29(9):811-5

Previous sbv IMPROVER crowd-sourcing challenges

1. Diagnostic Signature Challenge (2012)

The identification of gene expression signatures and computational methods for diagnostic classification

Sub-Challenge 1: Psoriasis
Identify normal vs. psoriatic skin based on the transcriptome of a skin biopsy.

Sub-Challenge 2: Multiple sclerosis
Identify control vs. affected or remitting vs. relapsing patients based on the transcriptome of peripheral blood mononuclear cells (PBMCs).

Sub-Challenge 3: Chronic Obstructive Pulmonary Disease (COPD)
Identify affected vs. non-affected subjects based on the transcriptome of bronchial brushings.

Sub-Challenge 4: Lung cancer
Identify stages 1 and 2 of squamous cell carcinoma (SCC) vs. adenocarcinoma (AC) based on the transcriptome of the tumor.

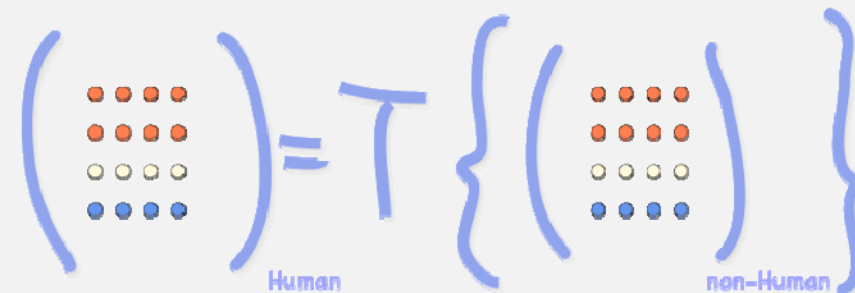
Diagnostic Signature Benchmarking
Use our Diagnostic Signature Benchmarking tool to see how you compare with your peers

1 DOWNLOAD 2 SUBMIT 3 COMPARE

<https://sbvimprover.com/challenge-1>

2. Species Translation Challenge (2013)

The translatability of biological system perturbations across species



3. Network Verification Challenge (2014-2015)

The verification and enhancement of biological causal networks representative of various biological processes

Cell fate
Cell stress
Cell proliferation
Inflammation
Tissue repair/angiogenesis

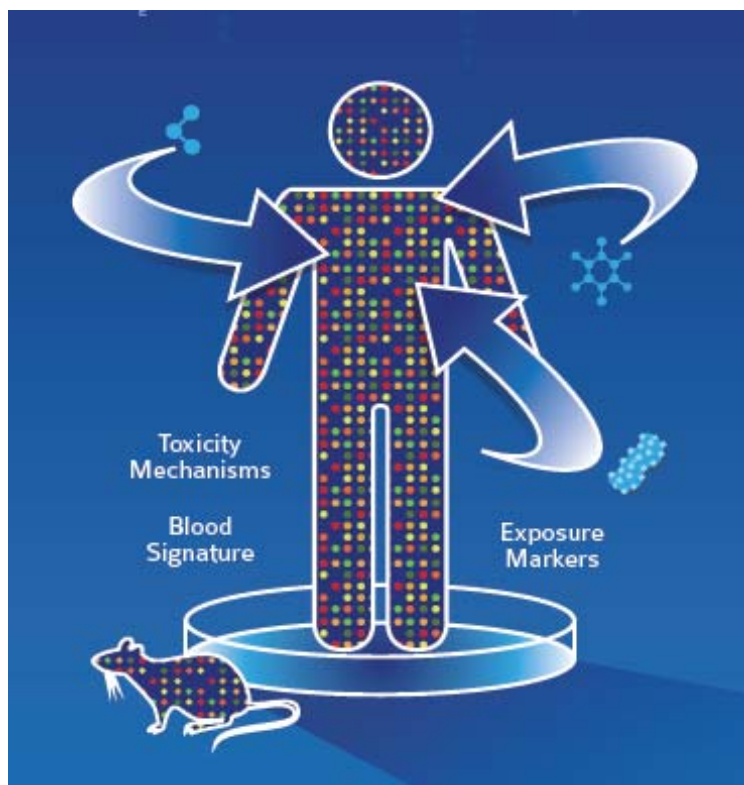
bionet.sbvimprover.com

Challenge outcome including best performing methods and lessons learned were:

- Presented in symposium
- Published in peer-reviewed journals

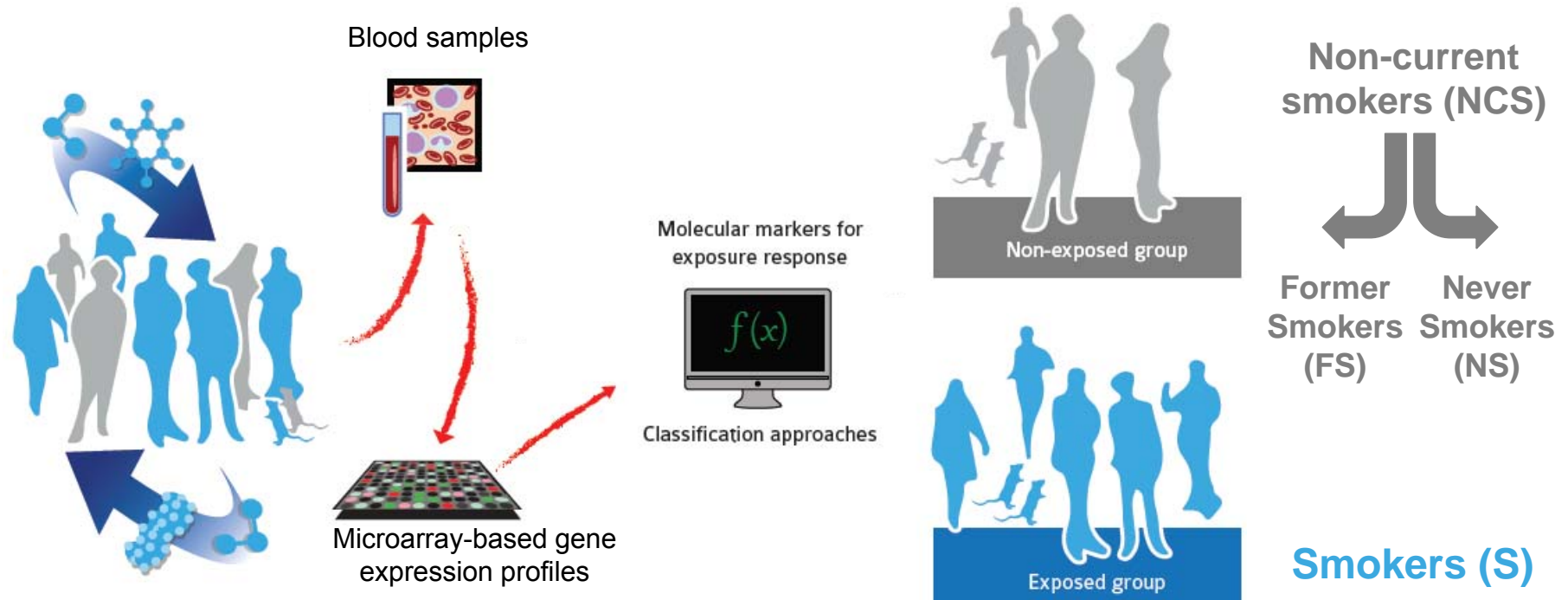
<https://sbvimprover.com/sbv-improver/publications/>

2016 - The Systems Toxicology Computational Challenge



<https://sbvimprover.com/challenge-4>

Overview



To develop a classification approach that identifies a blood gene signature capable of associating subjects to the correct exposure group

The Systems Toxicology Computational Challenge

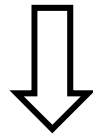
Marker of Exposure Response Identification

Goal: to develop **inductive** blood gene signature-based classification models to predict **smoking exposure** or **cessation status**

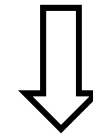
Sub-challenge 1 (SC1)
Human blood signature as
exposure response marker



Sub-challenge 2 (SC2)
Species translatable blood gene
signature as exposure response marker



Human signatures



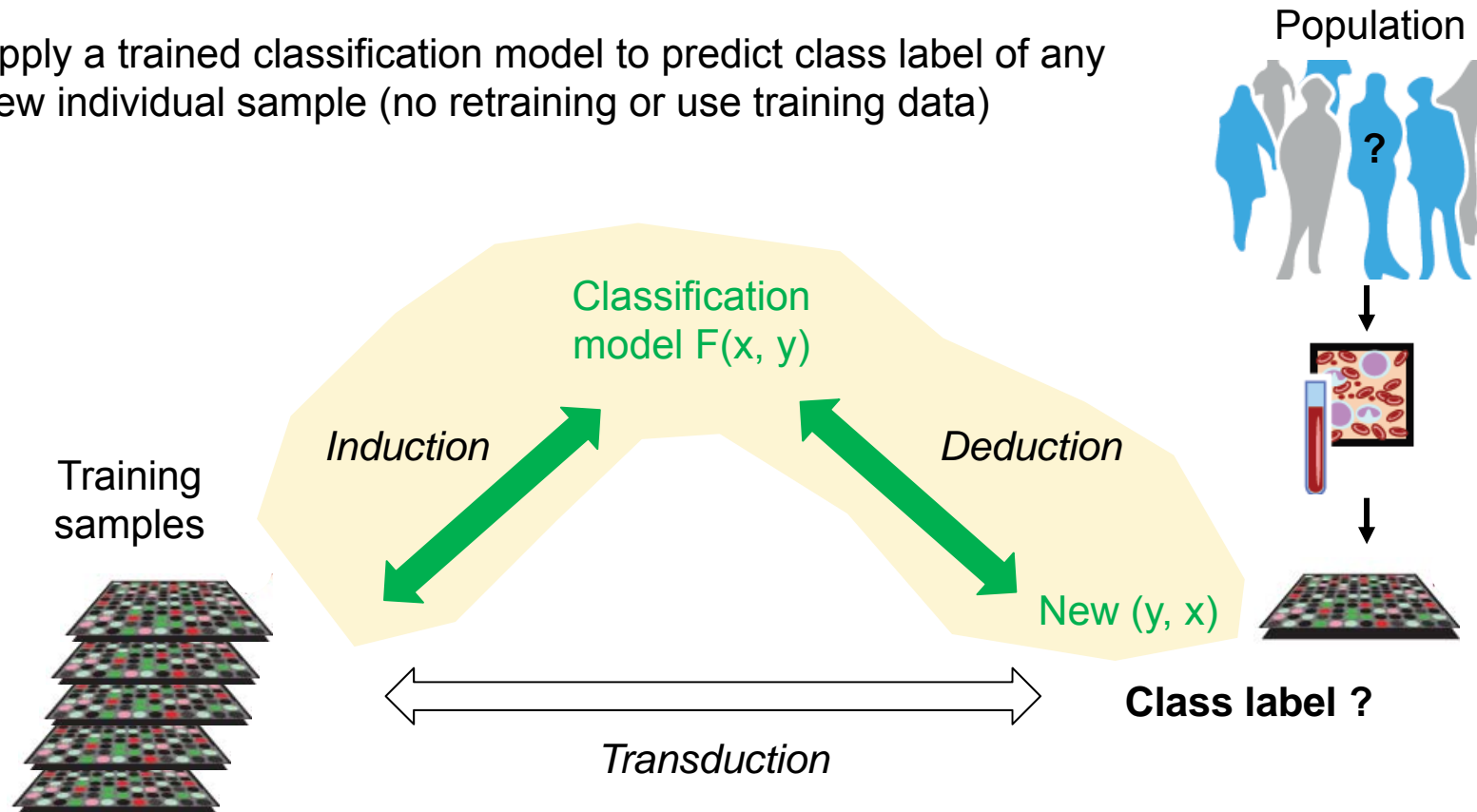
Species-independent signatures

Robust and **sparse** gene signatures → do not exceed 40 genes

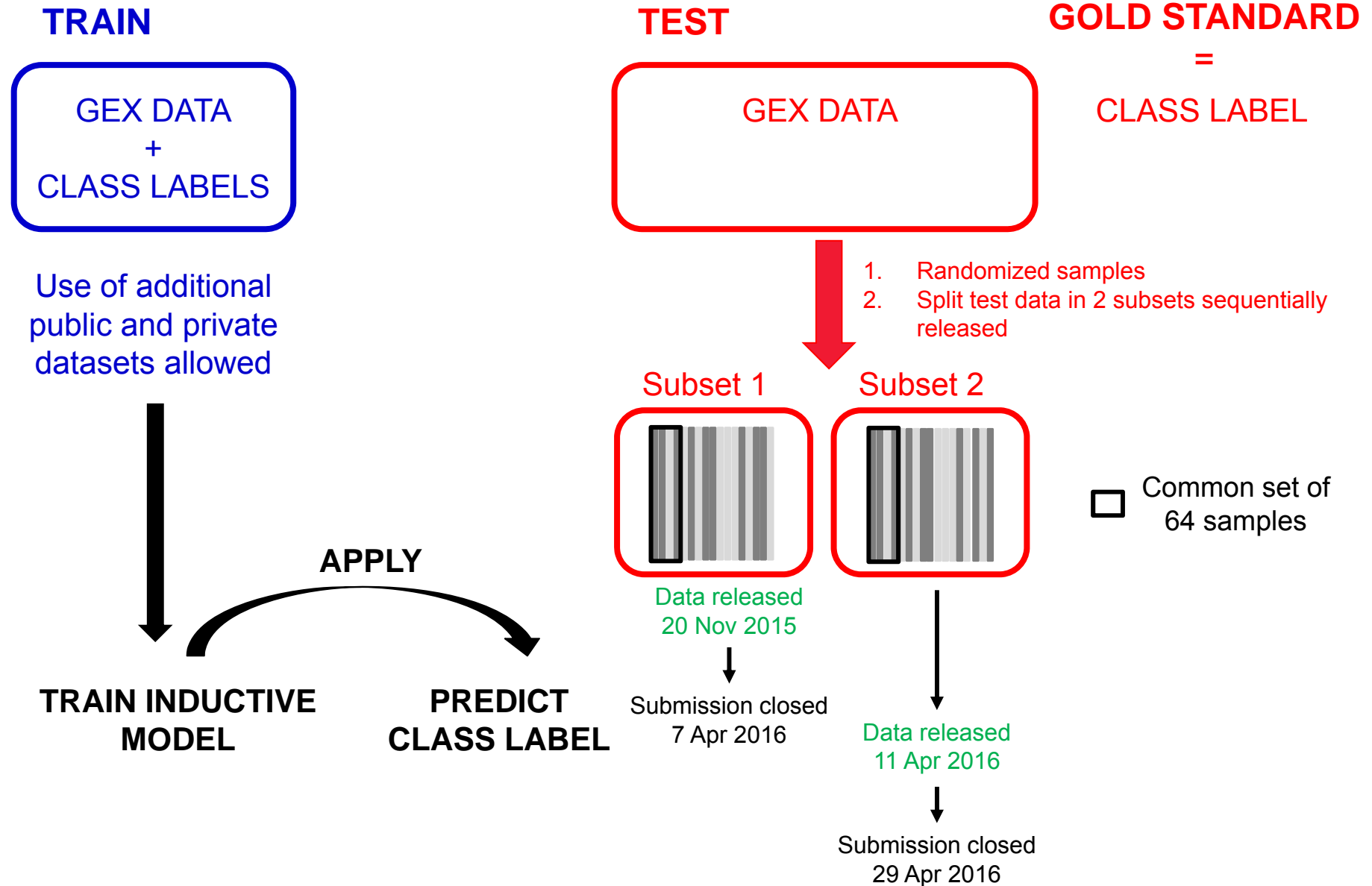
Classification models

Inductive classification model as opposed to transductive

Apply a trained classification model to predict class label of any new individual sample (no retraining or use training data)



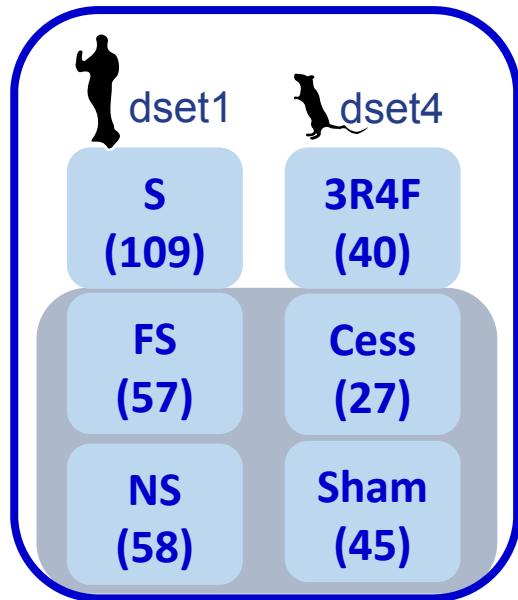
Training and Test data



Study datasets

Gene expression data generated from human and mouse blood samples

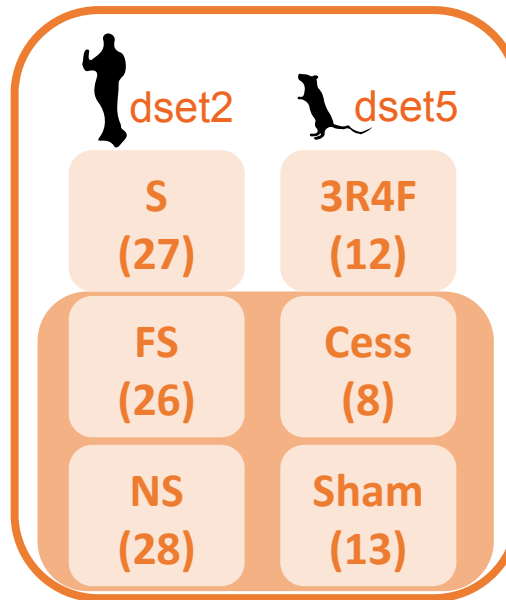
TRAIN



SC1 ✓

SC2 ✓ ✓

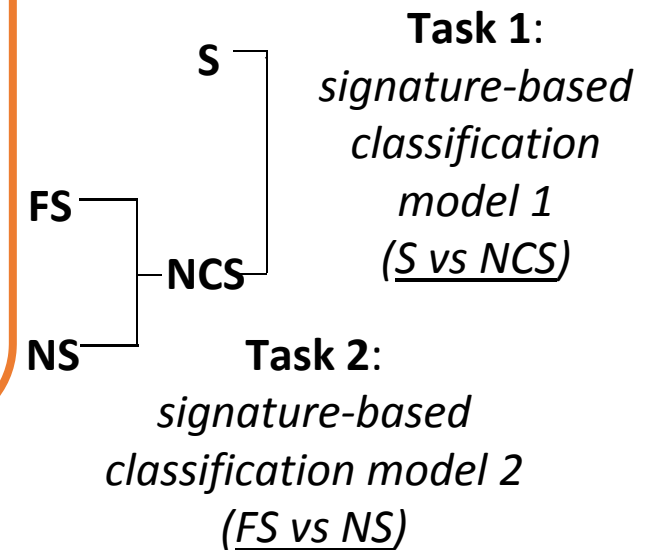
TEST



✓

✓ ✓

Step-wise class predictions



S/3R4F: Smokers / 3R4F (exposure to smoke from a reference cigarette)

FS/Cess: Former smokers / Cessation

NS/Sham: Never smokers / Sham

NCS: Non-current smoker

Freedom to use two separate models for 2-class prediction for each step, or directly a 3-class prediction model

Submissions and Timelines

Submissions for each sub-challenge

- 4 class label prediction files
- 2 gene lists
- 1 write up

Sample ID	Smoker	Non-current smoker
Sample 1	P1	P2
Sample 2	0.95	0.05
Sample 3	0.94	0.04
....		
Sample M	0.85	0.15

Sample ID	Former Smoker	Non-current smoker
Sample 1	P1	P2
Sample 2	0.95	0.05
Sample 3	0.94	0.04
....		
Sample M	0.85	0.15

- P_x = confidence value that a sample belongs to a class

11 • $P1+P2=1$; $P1 \neq P2$

