

## Supplementary Materials and Legends

### Identification of Early Liver Toxicity Gene Biomarkers Using Comparative

#### Supervised Machine Learning

Brandi Patrice Smith<sup>1,2</sup>(brandis2@illinois.edu), Loretta Sue Auvil<sup>3</sup>(lauvil@illinois.edu), Michael Welge<sup>3,4</sup>(mwelge@illinois.edu), Colleen Bushell<sup>3,4,5</sup>(cbushell@illinois.edu), Rohit Bhargava<sup>6,8,9</sup>(rxb@illinois.edu), Navin Elango<sup>7</sup>(navin.elango@corveva.com), Kamin Johnson<sup>7</sup>(kamin.johnson@corveva.com), Zeynep Madak-Erdogan<sup>\*1,3,4,8,9</sup>(zmadake2@illinois.edu)

<sup>1</sup> Department of Food Science and Human Nutrition, University of Illinois, Urbana-Champaign, Urbana, IL, USA

<sup>2</sup> Illinois Informatics Institute, University of Illinois, Urbana-Champaign, Urbana, IL USA

<sup>3</sup> National Center for Supercomputing Applications, University of Illinois, Urbana-Champaign, Urbana, IL, USA

<sup>4</sup> Carl R. Woese Institute for Genomic Biology, University of Illinois, Urbana-Champaign, Urbana, IL, USA

<sup>5</sup> Carle Illinois College of Medicine, University of Illinois, Urbana-Champaign, Urbana, IL, USA

<sup>6</sup>Department of Bioengineering, University of Illinois, Urbana-Champaign, Urbana, IL, USA

<sup>7</sup> Corteva Agrisciences, The Agriculture Division of DowDupont, Indianapolis, IN, USA

<sup>8</sup> Cancer Center at Illinois, University of Illinois, Urbana-Champaign, Urbana, IL, USA

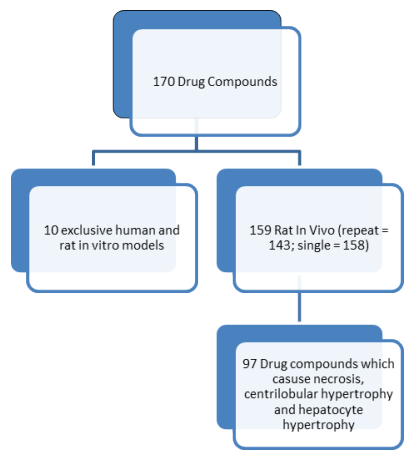
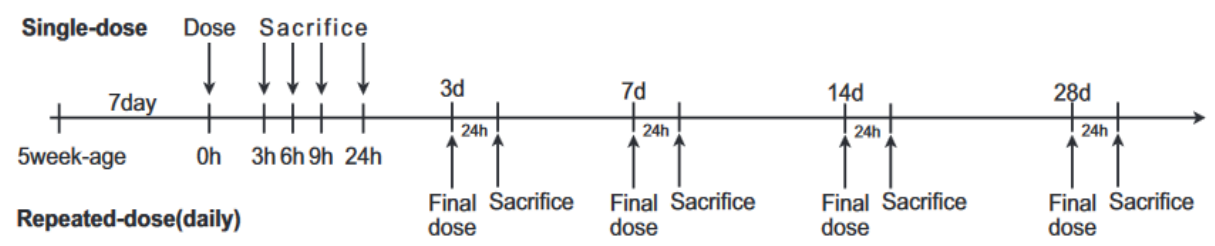
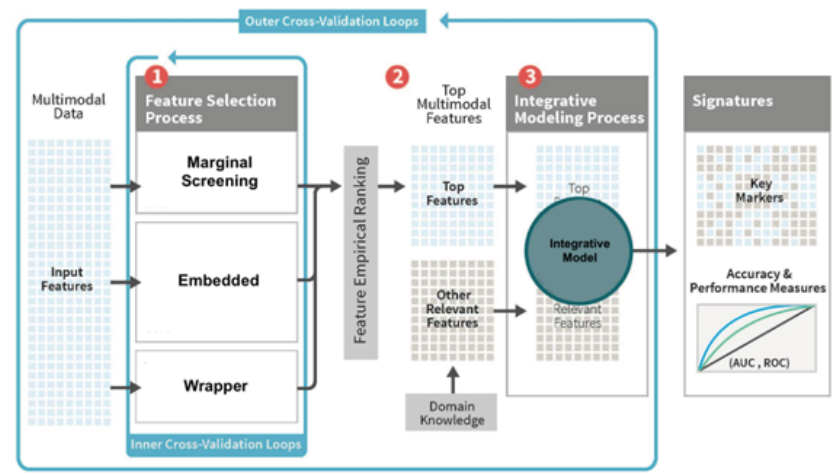
<sup>9</sup> Beckman Institute for Advanced Science and Technology, University of Illinois at Urbana-Champaign, Urbana, IL, USA

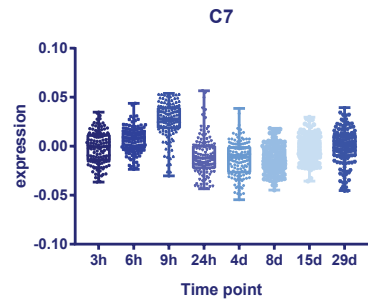
**Supplementary Figure 1: A.** Overview of the compounds in TG-GATES data. **B.** Experimental design. Timeline and doses used for experiments in TG-GATES database **C.** Flow chart for development of the ML pipeline.

All graphs are generated using Graph is generated using Powerpoint software (Microsoft Corporation, Seattle, WA, USA).

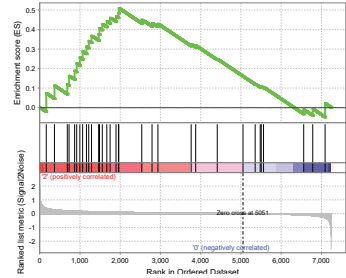
**Supplementary Figure 2:** GSEA analysis of clusters related to Figure 2. Upper panel graphs are generated by Graphpad© Prism8 software (GraphPad Software Inc., La Jolla, CA, [www.graphpad.com](http://www.graphpad.com)). Lower panels are generated using Gene Set Enrichment Analysis software (<https://www.gsea-msigdb.org/gsea/index.jsp>)<sup>31,32</sup>.

**Supplementary Figure 3:** Evaluation of average ROC for gene sets up to **A.** 50 and **B.** 100 . Figures are generated using Tableau software (Seattle, WA, USA, <https://www.tableau.com/>).

**A****B****C**

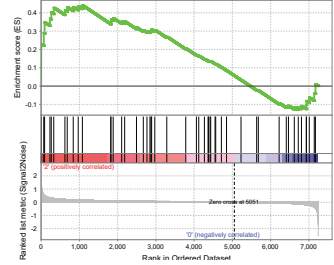
**A**

Enrichment plot:  
GO\_S\_ADENOSYLMETHIONINE\_DEPENDENT\_METHYLTRANSFERASE\_ACTIVITY

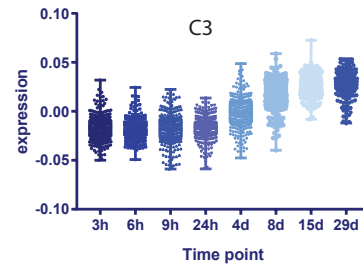


PRMT3  
SMYD2  
LCMT2  
PRMT1  
DNMT1  
FBL  
ICMT  
PCMT1  
DPH5  
TFB1M  
CXXC1  
FTSJ3  
WDR77  
TRMT1  
WDR5  
TFB2M  
RNMT  
EHMT2  
RG9MTD1

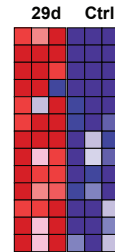
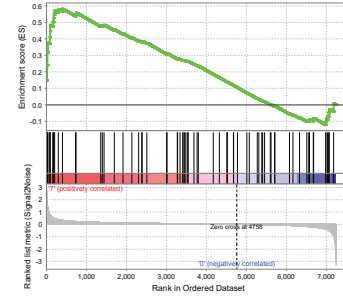
Enrichment plot:  
GO\_TRANSCRIPTIONAL\_ACTIVATOR\_ACTIVITY\_RNA\_P  
CORE\_PROMOTER\_PROXIMAL\_REGION\_SEQUENCE\_SPI



DBP  
ESRRA  
KLF15  
MAFB  
PLSCR1  
CEBPB  
KLF4  
SOX18  
NFIA  
DDIT3  
ATF4

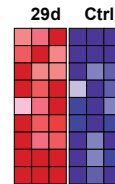
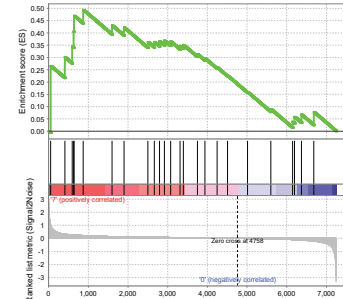
**B**

Enrichment plot:  
HALLMARK\_ESTROGEN\_RESPONSE\_LATE

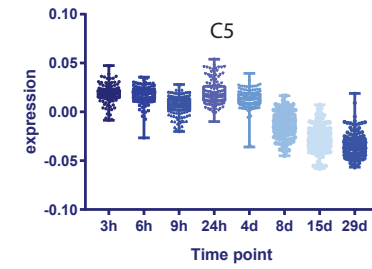


ASCL1  
PGR  
KLF4  
CD9  
CISH  
IMPA2  
BTG3  
ABCA3  
FGFR3  
PRLR  
SLC22A5  
PTGER3  
FRK

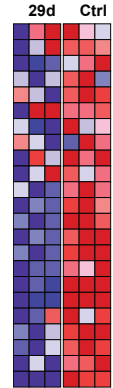
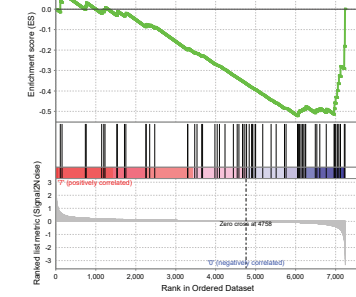
Enrichment plot: HALLMARK\_TGF\_BETA\_SIGNALING



ID1  
SMAD6  
SMAD3  
SLC20A1  
ID3  
PPM1A  
FN1A  
FKBP1A

**C**

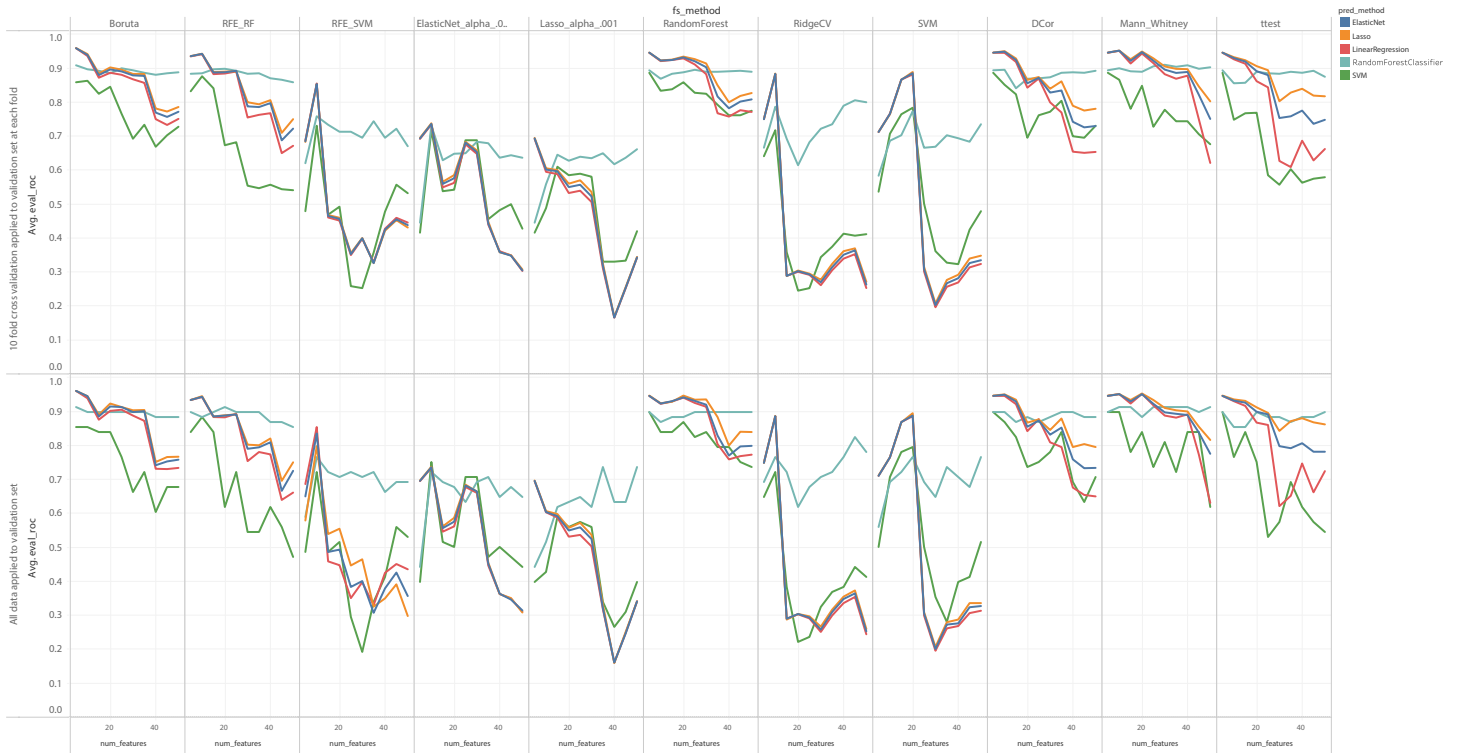
Enrichment plot: HALLMARK\_BILE\_ACID\_METABOLISM



PEX11A  
CROT  
GCLM  
ACSL5  
RBP1  
RETSAT  
GNMT  
NR0B2  
NR1H2  
ALDH1A1  
SLC23A1  
BBOX1  
IDH1  
AKR1D1  
FADS1  
HACL1  
AQP9  
AMACR  
HSD3B7  
FADS2  
CY8B1  
HSD3B1  
ABCG8

**A**

## Validation Results/Up tp 50 genes

**B**

## Validation Results/Up tp 100 genes

