*Environ Health Perspect*

**Note to readers with disabilities:** *EHP* strives to ensure that all journal content is accessible to all readers. However, some figures and Supplemental Material published in *EHP* articles may not conform to 508 standards due to the complexity of the information being presented. If you need assistance accessing journal content, please contact ehp508@niehs.nih.gov. Our staff will work with you to assess and meet your accessibility needs within 3 working days.

## Supplemental Material

**The Carcinogenome Project: *In Vitro* Gene Expression Profiling of Chemical Perturbations to Predict Long-Term Carcinogenicity**

Amy Li, Xiaodong Lu, Ted Natoli, Joshua Bittker, Nisha S. Sipes, Aravind Subramanian, Scott Auerbach, David H. Sherr, and Stefano Monti

## Table of Contents

**Table S11A.** GSEA analysis of signature DN_GTX_LIVER enrichment in Drugmatrix Cell Culture Low Dose with respect to phenotype Genotoxicity (GSEA results summary).

**Table S11B.** GSEA analysis of signature DN_GTX_LIVER enrichment in Drugmatrix Cell Culture Low Dose with respect to phenotype Genotoxicity (Gene ranking).

**Table S12.** Cmap Perturbagen Classes (PCL) with differential connectivity between carcinogens and non-carcinogens.

**Table S13.** Cmap Perturbagen Classes (PCL) with differential connectivity between genotoxicants and non-genotoxicants.

**Table S14.** Cmap Perturbagens with differential connectivity between carcinogens and non-carcinogens.

**Table S15.** Cmap Perturbagens with differential connectivity between genotoxicants and non-genotoxicants.

**Figure S1.** Overview of Experimental Design and Analysis Aims: (A) Data generation and annotation: Chemicals with long-term in vivo chemical annotation, as annotated by the Carcinogenic Potency Project, were procured. HepG2 cells are exposed to each chemical and followed by gene expression profiling. The number of unique chemicals and unique profiles by category (carcinogen, non-carcinogen, others) were catalogued. (B) Data analysis: analysis of the data consists of 1) analysis of transcriptional bioactivity using the Transcriptional Activity Scores (TAS), 2) prediction of carcinogenicity and genotoxicity, 3) mechanisms of action analysis using differential pathway enrichment analysis, and 4) comparison to other signatures such as signatures of carcinogenicity (Drugmatrix), small molecule perturbations (Connectivity Map) and Aryl hydrocarbon receptor (AhR) Receptor activity (Tox21).

**Figure S2.** Distribution of Transcriptional Activity Scores (TAS) grouped by chemical genotoxicity within each dose level. P-values indicate the significance of unpaired one-sided two-group TAS comparison between TAS of genotoxic chemicals and TAS of non-genotoxic chemicals within each dose group (* = $p < 0.05$) (see methods). The lower, middle, upper hinges correspond to the 25th, 50th (median), and 75th percentile. The upper and lower whiskers extend to the smaller and largest value at most 1.5 * IQR (inter-quartile range) from the hinge. Data points beyond the whiskers are represented as dots. Following multiple hypothesis testing, the FDR values are reported as follows: Dose rank 1: FDR = 0.12, Dose rank 2: FDR = 0.88, Dose rank 3: FDR = 0.12, Dose rank 4: FDR = 0.24, Dose rank 5: FDR = 0.55, Dose rank 6: FDR = 0.12.

**Figure S3.** Distribution of Transcriptional Activity Scores (TAS) grouped by chemical genotoxicity within each dose level, separated by different chemical procurement sources: (A) Sigma Aldrich chemicals with max dose of 20uM and (B) NTP chemicals with max dose of 40uM. P-values indicate the significance of unpaired one-sided two-group TAS comparison between TAS of genotoxic chemicals and TAS of non-genotoxic chemicals within each dose group (* = $p < 0.05$) (see methods). The lower, middle, upper hinges correspond to the 25th, 50th (median), and 75th percentile. The upper and lower whiskers extend to the smaller and largest value at most 1.5 * IQR (inter-quartile range) from the hinge. Data points beyond the whiskers are represented as dots.

**Figure S4.** Sensitivity and specificity rates of classifiers at threshold of 0.3 in predictive models of carcinogenicity and genotoxicity. Boxplots have the following specifications: the lower, middle, upper hinges corresponding to the 25th, 50th (median), and 75th percentile, the upper and lower whiskers extend to the smaller and largest value at most 1.5 * IQR (inter-quartile range) from the hinge, and data points beyond the whiskers represented as dots.

**Figure S5.** Prediction probabilities on unlabeled chemicals for prediction of carcinogenicity in (A) all unlabeled profiles (B) profiles with Trascriptional Activity Scores (TAS) > 0.4, and for prediction of genotoxicity in (C) all unlabeled profiles (D) profiles with TAS > 0.4.

**Figure S6.** Heatmap of pathway enrichment scores (GSVA) for top 40 upregulated and downregulated differential pathways of carcinogenicity (A) and genotoxicity (B) for profiles with TAS > 0.2. Columns are clustered using the ward method with Euclidean distances. Rows are ordered by the frequency of the pathway categories among the top 40 (direction sensitive).

**Figure S7.** Pathway enrichment (pathways in Reactome 2016) of directionally inconsistent signatures between Drugmatrix and L1000 using Enrichr (Chen et al. 2013; Kuleshov et al. 2016).

**References Cited in Supplemental Material**

**Additional File**- Excel Document

# 1. Supplemental Tables (see Supplemental Excel File)

# 2. Supplemental Figures

Figure S1



## A. Data generation and annotation

**Long-term in vivo chemical annotation**

Liver carcinogenicity

The Carcinogenic Potency Project

carcinogenic | noncarcinogenic

Compound₁ Compound₂ ... Compoundₙ

**Short term in vitro exposures and gene expression profiling**

HEPG2 profiling — 6 doses, triplicates per dose

| Chemical Type | # Chemicals | # Profiles |
|---|---|---|
| Carcinogens | 131 | 2358 |
| Non-carcinogens | 172 | 3096 |
| Others | 33 | 594 |
| Total | 336 | 6048 |

## B. Analysis of transcriptional bioactivity

Genes — high transcriptional impact — Samples

Transcriptional Activity Scores

**Prediction of carcinogenicity/genotoxicity**

predictive model → Carcinogen / Non-carcinogen

**Mechanisms of Action analysis**

Gene expression matrix (Gene 1, Gene 2, ... Gene n)

Gene membership in pathways — MSigDB Molecular Signatures Database

Gene-set projection
- GSVA
Differential expression
- limma

Differentially expressed pathways (Pathway 1, Pathway 2, ... Pathway n)

**Comparison to other signatures**

L1000 gene expression data

Signatures of carcinogenicity (Drugmatrix) or small molecule perturbations (Cmap) or AhR receptor activity (Tox21)

Calculate similarities

Enriched Sets
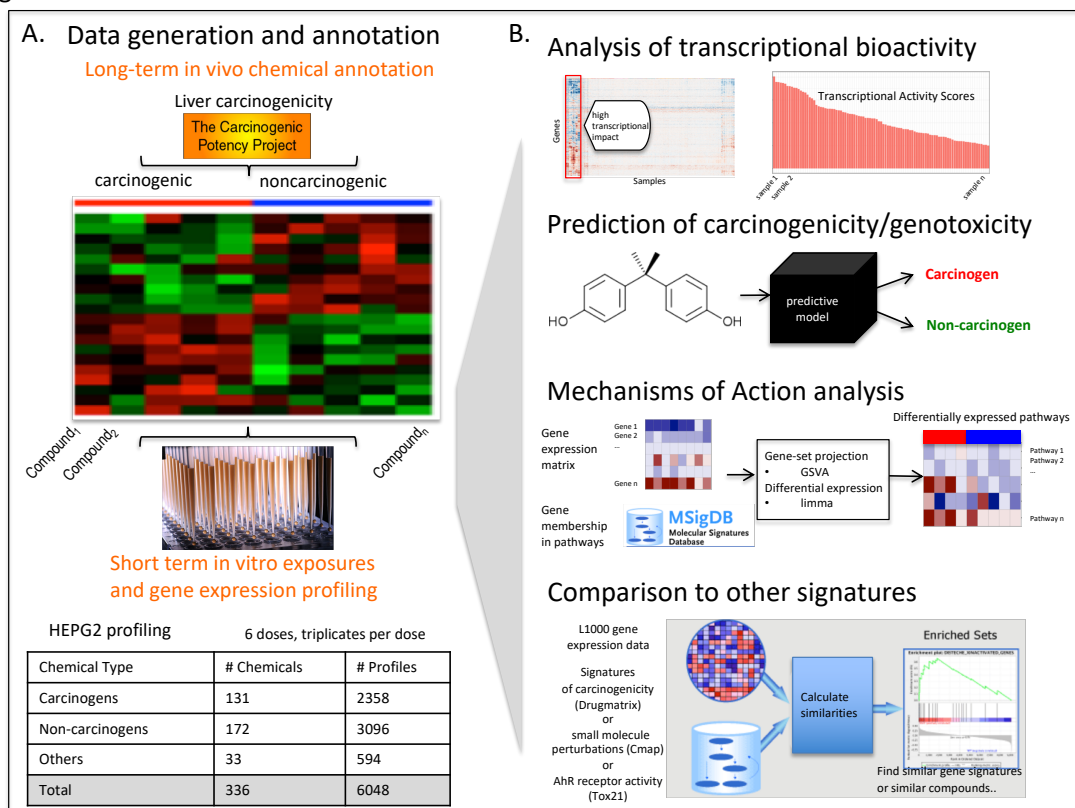
Find similar gene signatures or similar compounds..

*Figure S1:* Overview of Experimental Design and Analysis Aims: (A) Data generation and annotation: Chemicals with long-term in vivo chemical annotation, as annotated by the Carcinogenic Potency Project, were procured. HepG2 cells are exposed to each chemical and followed by gene expression profiling. The number of unique chemicals and unique profiles by category (carcinogen, non-carcinogen, others) were catalogued. (B) Data analysis: analysis of the data consists of 1) analysis of transcriptional bioactivity using the Transcriptional Activity Scores (TAS), 2) prediction of carcinogenicity and genotoxicity, 3) mechanisms of action analysis using differential pathway enrichment

analysis, and 4) comparison to other signatures such as signatures of carcinogenicity (Drugmatrix), small molecule perturbations (Connectivity Map) and Aryl hydrocarbon receptor (AhR) Receptor activity (Tox21).
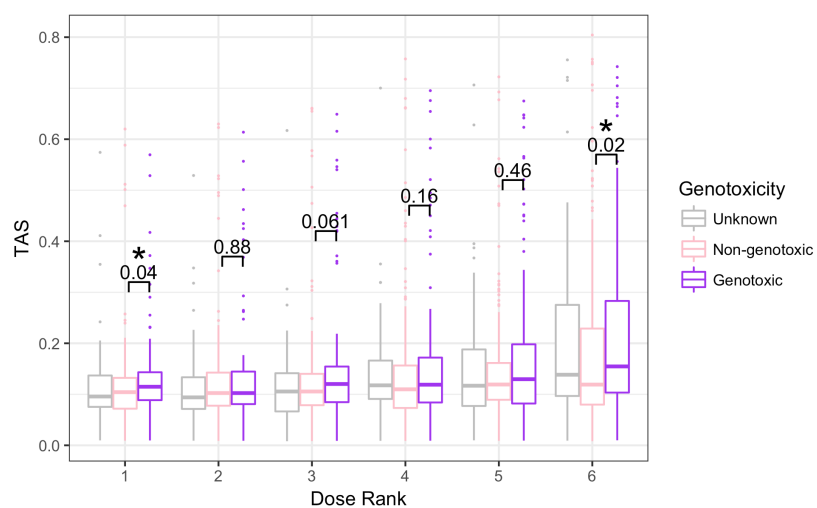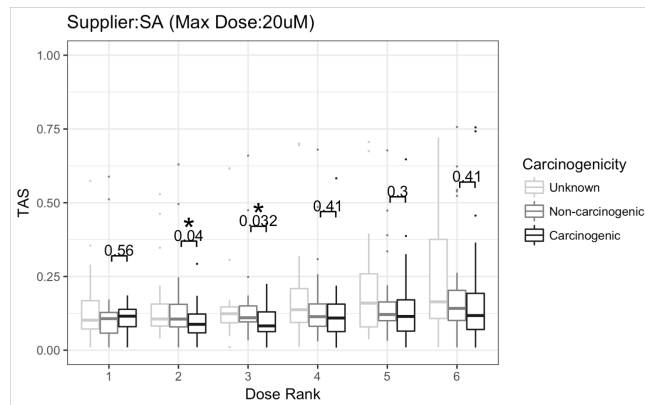
Figure S2



*Figure S2:* Distribution of Transcriptional Activity Scores (TAS) grouped by chemical genotoxicity within each dose level. P-values indicate the significance of unpaired one-sided two-group TAS comparison between TAS of genotoxic chemicals and TAS of non-genotoxic chemicals within each dose group (* = $p < 0.05$) (see methods). The lower, middle, upper hinges correspond to the 25th, 50th (median), and 75th percentile. The upper and lower whiskers extend to the smaller and largest value at most 1.5 * IQR

(inter-quartile range) from the hinge. Data points beyond the whiskers are represented as

dots. Following multiple hypothesis testing, the FDR values are reported as follows: Dose

rank 1: FDR = 0.12, Dose rank 2: FDR = 0.88, Dose rank 3: FDR = 0.12, Dose rank 4:

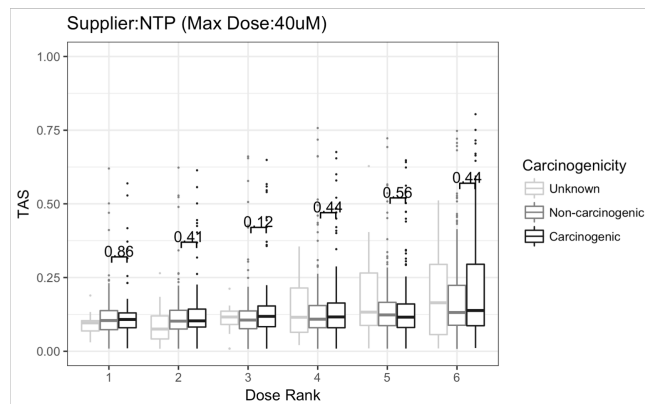FDR = 0.24, Dose rank 5: FDR = 0.55, Dose rank 6: FDR = 0.12.



Figure S3   A.   Supplier:SA (Max Dose:20uM)

B.   Supplier:NTP (Max Dose:40uM)

*Figure S3:* Distribution of Transcriptional Activity Scores (TAS) grouped by chemical

genotoxicity within each dose level, separated by different chemical procurement

sources: (A) Sigma Aldrich chemicals with max dose of 20uM and (B) NTP chemicals

with max dose of 40uM. P-values indicate the significance of unpaired one-sided two-

group TAS comparison between TAS of genotoxic chemicals and TAS of non-genotoxic

chemicals within each dose group (* = p< 0.05) (see methods). The lower, middle, upper

hinges correspond to the 25th, 50th (median), and 75th percentile. The upper and lower

whiskers extend to the smaller and largest value at most 1.5 * IQR (inter-quartile range) from the hinge. Data points beyond the whiskers are represented as dots.
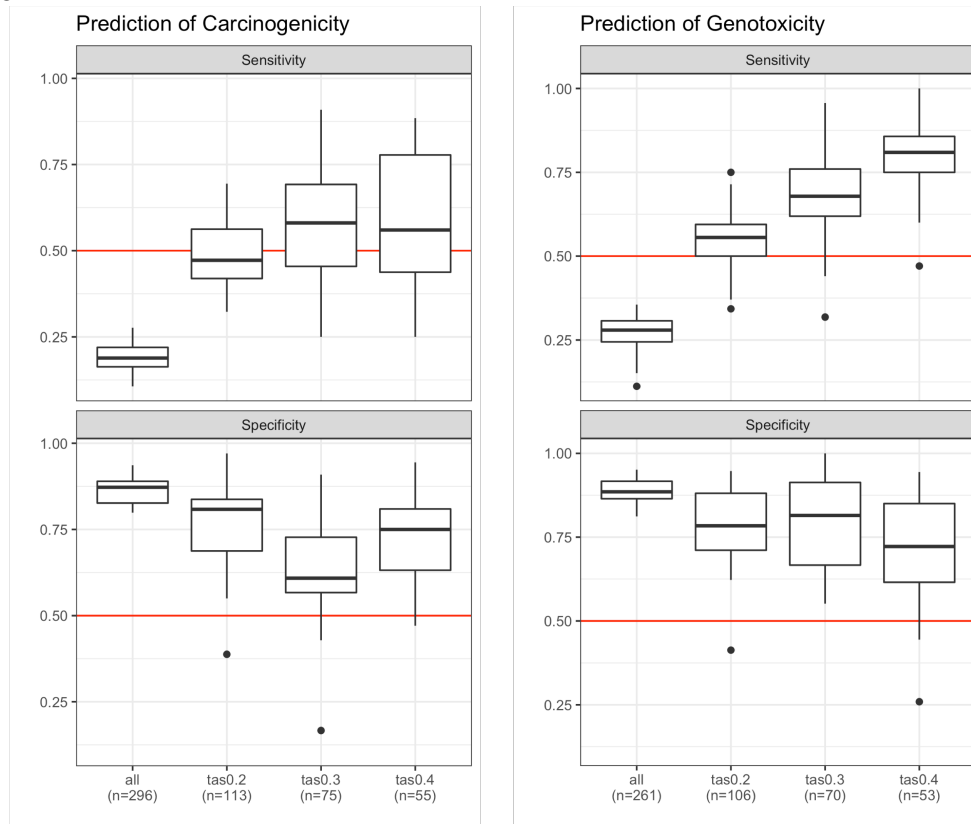
Figure S4



***Figure S4:*** Sensitivity and specificity rates of classifiers at threshold of 0.3 in predictive models of carcinogenicity and genotoxicity. Boxplots have the following specifications: the lower, middle, upper hinges corresponding to the 25th, 50th (median), and 75th percentile, the upper and lower whiskers extend to the smaller and largest value at most 1.5 * IQR (inter-quartile range) from the hinge, and data points beyond the whiskers represented as dots.
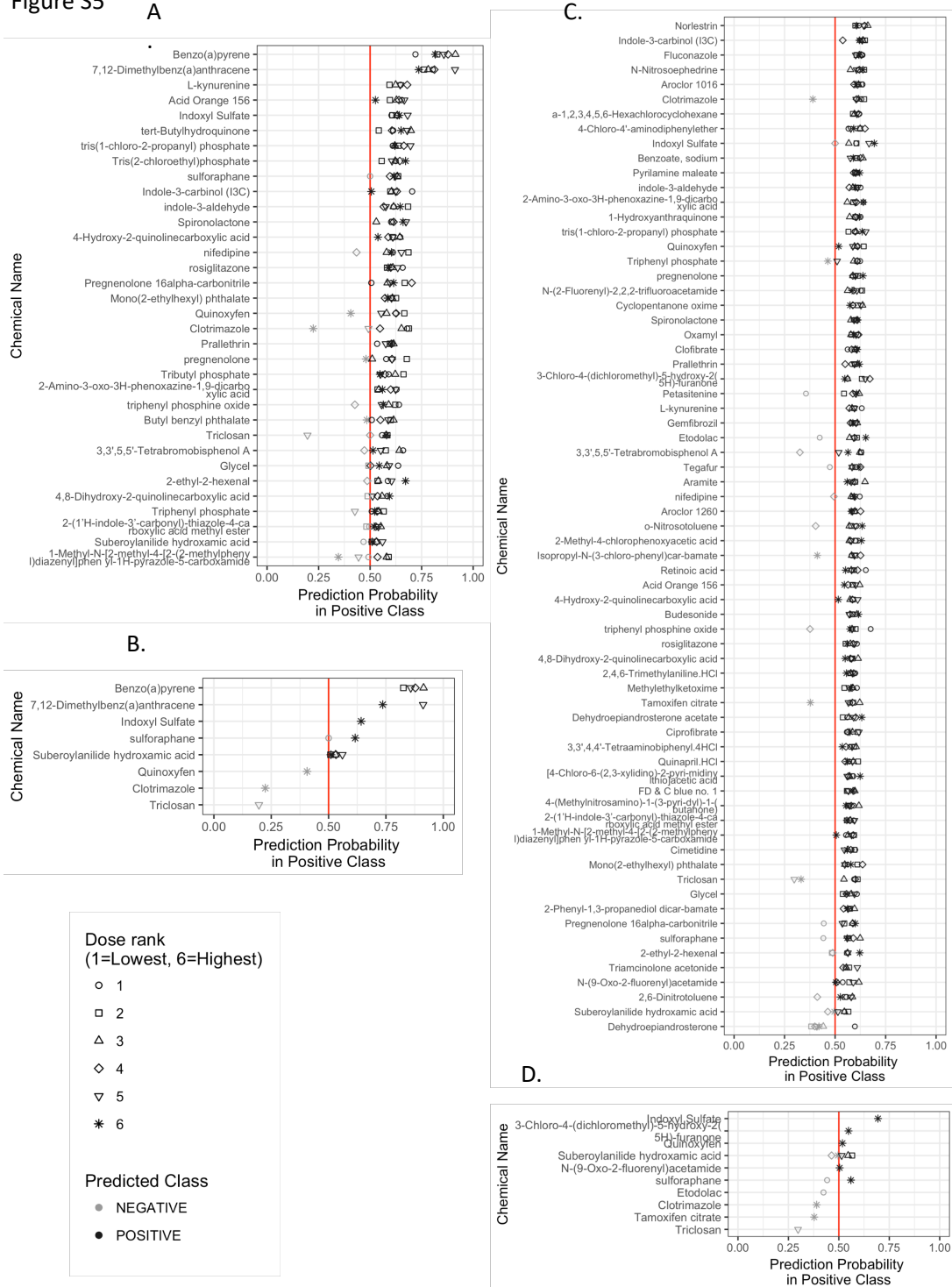
**Figure S5:** Prediction probabilities on unlabeled chemicals for prediction of carcinogenicity in (A) all unlabeled profiles (B) profiles with Trascriptional Activity

Scores (TAS) > 0.4, and for prediction of genotoxicity in (C) all unlabeled profiles (D) profiles with TAS > 0.4
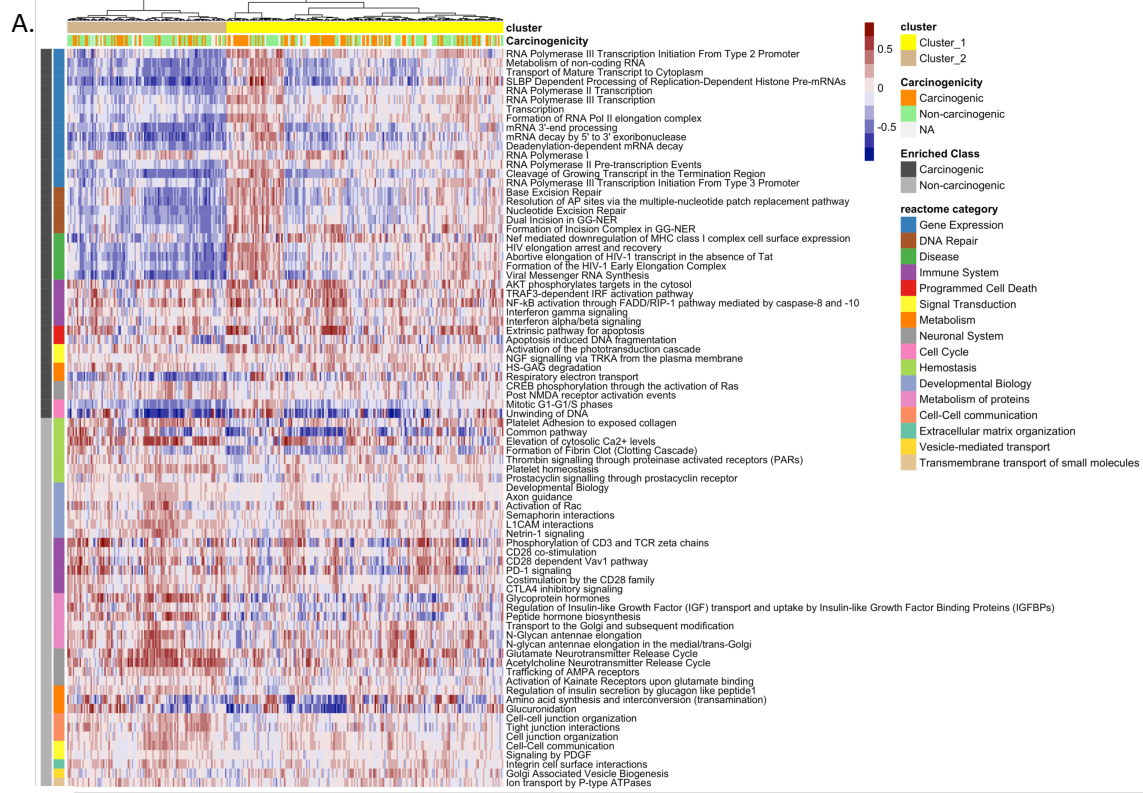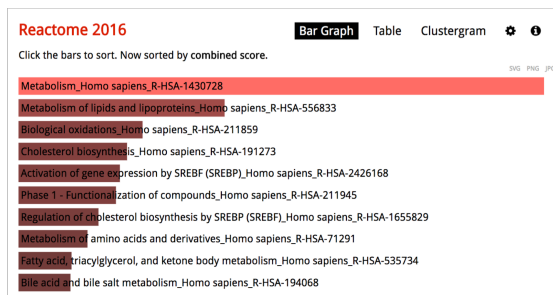
Figure S6

***Figure S6.*** Heatmap of pathway enrichment scores (GSVA) for top 40 upregulated and downregulated differential pathways of carcinogenicity (A) and genotoxicity (B) for profiles with TAS > 0.2. Columns are clustered using the ward method with Euclidean distances. Rows are ordered by the frequency of the pathway categories among the top 40 (direction sensitive).
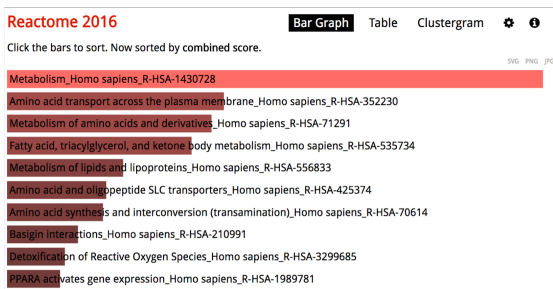
Figure S7

UP_GTX_LIVER



DN_GTX_LIVER



***Figure S7:*** Pathway enrichment (pathways in Reactome 2016) of directionally inconsistent signatures between Drugmatrix and L1000 using Enrichr (Chen et al. 2013; Kuleshov et al. 2016).

## 3. References Cited in Supplemental Material

Chen EY, Tan CM, Kou Y, Duan Q, Wang Z, Meirelles G V., et al. 2013. Enrichr: Interactive and collaborative HTML5 gene list enrichment analysis tool. BMC Bioinformatics 14; doi:10.1186/1471-2105-14-128.

Kuleshov M V., Jones MR, Rouillard AD, Fernandez NF, Duan Q, Wang Z, et al. 2016. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. Nucleic Acids Res 44:W90–W97; doi:10.1093/nar/gkw377.