# T1000: A reduced toxicogenomics gene set for improved decision making

**Othman Soufan** [1], **Jessica Ewald** [2], **Charles Viau** [1], **Doug Crump** [3], **Markus Hecker** [4], **Niladri Basu** [Corresp., 2], **Jianguo Xia** [Corresp. 1]

[1] Institute of Parasitology, McGill University, Montreal, Canada

[2] Faculty of Agricultural and Environmental Sciences, McGill University, Montreal, Canada

[3] National Wildlife Research Centre, Carleton University, Ottawa, Canada

[4] School of the Environment & Sustainability and Toxicology Centre, University of Saskatchewan, Saskatoon, Canada

Corresponding Authors: Niladri Basu, Jianguo Xia
Email address: niladri.basu@mcgill.ca, jeff.xia@mcgill.ca

There is growing interest within regulatory agencies and toxicological research communities to develop, test, and apply new approaches, such as toxicogenomics, to more efficiently evaluate chemical hazards. Given the complexity of analyzing thousands of genes simultaneously, there is a need to identify reduced gene sets.Though several gene sets have been defined for toxicological applications, few of these were purposefully derived using toxicogenomics data. Here, we developed and applied a systematic approach to identify 1000 genes (called Toxicogenomics-1000 or T1000) highly responsive to chemical exposures. First, a co-expression network of 11,210genes was built by leveraging microarray data from the Open TG-GATEs program. This network was then re-weighted based on prior knowledge of their biological (KEGG, MSigDB) and toxicological (CTD) relevance. Finally, weighted correlation network analysis was applied to identify 258 gene clusters. T1000 was defined by selecting genes from each cluster that were most associated with outcome measures. For model evaluation, we compared the performance of T1000 to that of other gene sets (L1000, S1500, Genes selected by Limma, and random set) using two external datasets. Additionally, a smaller (T384) and a larger version (T1500) of T1000 were used for dose-response modeling to test the effect of gene set size. Our findings demonstrated that the T1000 gene set is predictive of apical outcomes across a range of conditions (e.g.,*in vitro*and *in vivo*, dose-response, multiple species, tissues, and chemicals), and generally performs as well, or better than other gene sets available.

1  **T1000: A reduced toxicogenomics gene set for improved decision making**

2  [1]Othman Soufan, [2]Jessica Ewald, [1]Charles Viau, [3]Doug Crump, [4]Markus Hecker, [2,*]Niladri Basu

3  and [1,5,*]Jianguo Xia

4  [1]Institute of Parasitology, McGill University, Montreal, Quebec, Canada; [2]Faculty of

5  Agricultural and Environmental Sciences, McGill University, Montreal, Quebec, Canada;

6  [3]Ecotoxicology and Wildlife Health Division, Environment and Climate Change Canada,

7  National Wildlife Research Centre, Carleton University, Ottawa, Canada; [4]School of the

8  Environment & Sustainability and Toxicology Centre, University of Saskatchewan, Saskatoon,

9  Canada; [5]Department of Animal Science, McGill University, Montreal, Quebec, Canada.


10

11  *Correspondance: Niladri Basu
12  Email:         niladri.basu@mcgill.ca
13  Tel:           (514) 398-8642
14  Address:       CINE Building,
15                 21111 Lakeshore Road
16                 Ste. Anne de Bellevue
17                 QC, Canada
18                 H9X 3V9.
19
20  *Correspondence: Jianguo Xia
21  Email:         jeff.xia@mcgill.ca
22  Tel:           (514) 398-8668
23  Address:       Office P107, Parasitology Building,
24                 21111 Lakeshore Road
25                 Ste. Anne de Bellevue
26                 QC, Canada
27                 H9X 3V9.
28

29

30

## Abstract

There is growing interest within regulatory agencies and toxicological research communities to

develop, test, and apply new approaches, such as toxicogenomics, to more efficiently evaluate

chemical hazards. Given the complexity of analyzing thousands of genes simultaneously, there is

a need to identify reduced gene sets. Though several gene sets have been defined for

toxicological applications, few of these were purposefully derived using toxicogenomics data.

Here, we developed and applied a systematic approach to identify 1000 genes (called

Toxicogenomics-1000 or T1000) highly responsive to chemical exposures.  First, a co-

expression network of 11,210 genes was built by leveraging microarray data from the Open TG-

GATEs program. This network was then re-weighted based on prior knowledge of their

biological (KEGG, MSigDB) and toxicological (CTD) relevance. Finally, weighted correlation

network analysis was applied to identify 258 gene clusters. T1000 was defined by selecting

genes from each cluster that were most associated with outcome measures. For model evaluation,

we compared the performance of T1000 to that of other gene sets (L1000, S1500, Genes selected

by Limma, and random set) using two external datasets. Additionally, a smaller (T384) and a

larger version (T1500) of T1000 were used for dose-response modeling to test the effect of gene

set size. Our findings demonstrated that the T1000 gene set is predictive of apical outcomes

across a range of conditions (e.g., *in vitro* and *in vivo*, dose-response, multiple species, tissues,

and chemicals), and generally performs as well, or better than other gene sets available.

## Introduction

Over the past decade there have been profound steps taken across the toxicological sciences and regulatory communities to help transform conventional toxicity testing largely based on animal models and apical outcome measurements to an approach that is founded on systems biology and predictive science (Kavlock et al. 2018; Knudsen et al. 2015; Villeneuve & Garcia-Reyero 2011). On the scientific side, efforts are being exemplified by emergent notions such as the Adverse Outcome Pathway framework (AOP; Ankley et al., 2010) and New Approach Methods (ECHA 2016). On the regulatory side, these are exemplified by changes to, for example, chemical management plans in Canada, the United States and REACH (ECHA 2007) across the European Union.

A core tenet underlying the aforementioned transformations, as catalyzed by the 2007 U.S. National Research Council report "Toxicity Testing in the 21st Century" (Andersen & Krewski 2009), is that perturbations at the molecular-level can be predictive of those at the whole organism-level. Though whole transcriptome profiling is increasingly popular, it still remains costly for routine research and regulatory applications. Additionally, building predictive models with thousands of features introduces problems due to the high dimensionality of the data and so considering a smaller number of genes has the potential to increase classification performance (Alshahrani et al. 2017; Soufan et al. 2015b). Identifying smaller panels of key genes that can be measured, analyzed and interpreted conveniently remain an appealing option for toxicological studies and decision making

74   In recent years, several initiatives across the life sciences have started to identify reduced gene

75   sets from whole transcriptomic studies. For example, the Library of Integrated Network-Based

76   Cellular Signatures (LINCS) project derived L1000, which is a gene set of 976 'Landmark'

77   genes chosen to infer the expression of 12,031 other highly connected genes in the human

78   transcriptome (Subramanian et al. 2017). In the toxicological sciences, the U.S. Tox21 Program

79   recently published S1500+, which is a set of 2,753 genes designed to be both representative of

80   the whole-transcriptome, while maintaining a minimum coverage of all biological pathways in

81   Kyoto Encyclopedia of Genes and Genomes (KEGG) (Kanehisa et al. 2007) and Molecular

82   Signatures Database (MSigDB) (Liberzon et al. 2015a). The first 1,500 genes were selected by

83   analyzing microarray data from 3,339 different studies, and the rest were nominated by members

84   of the scientific community (Mav et al. 2018).  L1000 and S1500 gene sets were originally

85   proposed to serve a different purpose. The 978 landmark genes of L1000 are chosen to infer

86   expression of other genes more accurately, while genes of S1500 are selected to achieve more

87   biological pathway coverage. Compared to L1000, the S1500 gene set attains more toxicological

88   relevance through the gene nomination phase, though its data-driven approach relies upon

89   microarray data primarily derived from non-toxicological studies. It worth nothing that about

90   33.7% of genes are shared between both signatures. Even though some differences can be

91   realized between L1000 and S1500, they are both strong candidates of gene expression modeling

92   and prediction (Haider et al. 2018).

93

94   The objectives of the current study were to develop and apply a systematic approach to identify

95   highly-responsive genes from toxicogenomic studies, and from these to nominate a set of 1000

96   genes to form the basis for the T1000 (Toxicogenomics-1000) reference gene set.  Co-expression

97  network analysis is an established approach using pairwise correlation between genes and

98  clustering methods to group genes with similar expression patterns (van Dam et al. 2018). First, a

99  co-expression network was derived using *in vitro* and *in vivo* data from human and rat studies

100  from the Toxicogenomics Project-Genomics Assisted Toxicity Evaluation System (Open TG-

101  GATEs) database. Next, the connections within the co-expression network were adjusted to

102  increase the focus on genes in KEGG pathways, the MSigDB, or the Comparative

103  Toxicogenomics Database (CTD) (Davis et al. 2017). This incorporation of prior biological and

104  toxicological knowledge was motivated by loose Bayesian inference to refine the

105  computationally-prioritized transcriptomic space. Clusters of highly connected genes were

106  identified from the resulting co-expression network, and machine learning models were applied

107  to prioritize clusters based on their association with apical endpoints. Clustering genes based on

108  expression data has been shown to be instrumental in functional annotation and sample

109  classification (Necsulea et al. 2014), with the rationale that genes with similar expression

110  patterns are likely to participate in the same biological pathways (Budinska et al. 2013). From

111  each cluster key genes were identified for inclusion in T1000. Testing and validation of T1000

112  was realized through two separate datasets (one from Open TG-GATEs and one from the U.S.

113  National Toxicology Program) that were not used for gene selection.  The current study is part of

114  the larger EcoToxChip project (Basu et al. 2019).

115

116  **Materials & Methods**

117  **Overview**

118  The work was conducted in four discrete phases as follows:  I) data preparation and gene co-

119  expression network generation; II) network clustering to group relevant genes; III) gene selection

120    and prioritization; and IV) external testing and performance evaluation.  Within these four study

121    phases there were eight activities or steps (**Figure 1**): 1) data preparation, 2) constructing co-

122    expression networks; 3) computing prior scores from toxicogenomics resources; 4) re-weighting

123    co-expression scores and applying graph clustering; 6) building local prediction models for each

124    cluster; 6) building a global prediction model using representative genes from each cluster; 7)

125    dose-response analysis and apical outcome correlation using an external dataset; and 8)

126    prediction accuracy analysis using an external dataset.

127

128    **Phase I: Data preparation and gene co-expression network generation**

129    The goal of phase I was to construct two network representations of the interactions between

130    toxicologically-relevant genes, with one based on TG-GATES microarray data (step 1) and the

131    other based on the KEGG, MSigDB, and CTD databases (step 3).

132

133    Step 1: data preparation

134    The derivation of T1000 was based on five public microarray datasets of toxicological relevance

135    (**Table 1**): four datasets from Open TG-GATEs (Igarashi et al. 2014b), and one dataset generated

136    by Thomas *et al* (referred to as the dose-response dataset in this manuscript) (Thomas et al.

137    2013).  **Table 1** provides a summary of all microarray datasets used in this study. For building

138    the initial T1000 gene set, we used three of the four Open TG-GATEs datasets (see datasets 1-3

139    in **Table 1**). For the performance evaluation and testing phase, we leveraged the fourth dataset

140    from Open TG-GATEs (see dataset 4 in **Table 1**), which was not used for gene ranking or

141    selection so that it could serve as an external validation dataset. The dose-response dataset was

142    used for an additional external validation (see dataset 5 in **Table 1**).

143

144    Open TG-GATEs is one of the largest publicly accessible toxicogenomics resources (Igarashi et

145    al. 2014b). This database comprises data from 170 compounds (mostly drugs) with the aim of

146    improving and enhancing drug safety assessment. It contains gene expression profiles and

147    traditional toxicological data derived from *in vivo* (rat) and *in vitro* (primary rat hepatocytes and

148    primary human hepatocytes) studies. To process the raw gene expression data files of Open TG-

149    GATEs, the Affy package (Gautier et al. 2004) was used to produce Robust Multi-array Average

150    (RMA) probe set intensities (Irizarry et al. 2003b). Gene annotation for human and rat was

151    performed using Affymetrix Human Genome U133 Plus 2.0 Array annotation data and

152    Affymetrix Rat Genome 230 2.0 Array annotation data, respectively. Genes without annotation

153    were excluded. When the same gene was mapped multiple times, the average value was used.

154    Finally, all profiles for each type of experiment were joined into a single matrix for downstream

155    analysis.

156

157    The dose-response dataset was used to externally evaluate the ability of T1000 genes to predict

158    apical endpoints (Thomas et al. 2013). Briefly, this dataset contains Affymetrix HT Rat230 PM

159    microarray data following *in vivo* exposure of rats to six chemicals (TRBZ: 1,2,4-

160    tribromobenzene, BRBZ: bromobenzene, TTCP: 2,3,4,6-tetrachlorophenol, MDMB: 4,4'-

161    methylenebis(*N,N'*-dimethyl)aniline, NDPA: N-nitrosodiphenylamine, and HZBZ:

162    hydrazobenzene). In exposed animals, both gene expression and apical outcomes (liver: absolute

163    liver weight, vacuolation, hypertrophy, microvesiculation, necrosis; thyroid: absolute thyroid

164    weight, follicular cell hypertrophy, follicular cell hyperplasia; bladder: absolute bladder weight,

165    increased mitosis, diffuse transitional epithelial hyperplasia, increased necrosis epithelial cell)

166    were measured, permitting the comparison of transcriptionally-derived benchmark doses ($BMD_t$)

167    with traditional benchmark doses derived from apical outcomes (Yang et al. 2007). The apical

168    outcome-derived benchmark dose ($BMD_a$) for each treatment group was defined as the

169    benchmark dose from the most sensitive apical outcome for the given chemical-duration group.

170

171    <u>Step 2: constructing a co-expression network</u>

172    In a co-expression network, nodes represent genes and edges represent the Pearson's correlation

173    of expression values of pairs of genes.  In the current study, we constructed three co-expression

174    networks using gene expression profiles from Open TG-GATEs datasets (human *in vitro*, rat *in*

175    *vitro*, and rat *in vivo*) (**Table 1**). If an interaction with a correlation coefficient of 60% or higher

176    was present in all three networks, that gene-gene interaction was then accepted and mapped into

177    one integrated co-expression network by averaging the absolute values of the pairwise

178    correlation coefficients between individual genes. The final integrated co-expression network

179    had 11,210 genes from a total of 20,502 genes.

180

181    <u>Step 3: computing prior scores from toxicogenomics resources</u>

182    The CTD, KEGG, and Hallmark databases were mined to integrate existing toxicogenomics and

183    broader biological knowledge into one network that represents the prior knowledge space. CTD

184    is manually curated from the literature to serve as a public source for toxicogenomics

185    information, currently including over 30.5 million chemical-gene, chemical-disease, and gene-

186    disease interactions (Davis et al. 2017). Following the recommendations of Hu et al. (2015), only

187    "mechanistic/marker" associations were extracted from the CTD database, thus excluding

188    "therapeutic" associations that are presumably less relevant to toxicology. The extracted

189   subgraph contained 2,889 chemicals, 950 diseases annotated as toxic endpoints (e.g.

190   neurotoxicity, cardiotoxicity, hepatotoxicity and nephrotoxicity), and 22,336 genes. KEGG

191   pathways are a popular bioinformatics resource that help to link, organize, and interpret genomic

192   information through the use of manually drawn networks describing the relationships between

193   genes in specific biological processes (Kanehisa et al. 2007). The MSigDB Hallmark gene sets

194   have been developed using a combination of automated approaches and expert curation to

195   represent known biological pathways and processes while limiting redundancy (Liberzon et al.

196   2015b). To build the prior knowledge space, we first encoded information from the three

197   databases into feature vectors describing each gene. Then, we applied dimensionality reduction

198   and K-means clustering to detect those genes that contributed most to the prior knowledge space.

199

200   Each feature vector consisted of 239 dimensions, representing information encoded from

201   Hallmark, KEGG and CTD. For the Hallmark and KEGG features, we used "1" or "0" to

202   indicate if a gene was present or absent for each of the 50 Hallmark gene sets (Liberzon et al.

203   2015b) and 186 KEGG pathways (Kanehisa & Goto 2000). These features were transformed into

204   z-scores. For the CTD features, we computed the degree, betweenness centrality, and closeness

205   centrality of each gene, based on the topology of the extracted CTD subgraph. The topology

206   measures were log-scaled for each gene in the network. The resulting prior knowledge space

207   consisted of a 239-dimension vector for each of the 22,336 genes, with each vector containing 50

208   z-score normalized Hallmark features, 186 z-score normalized KEGG features, and three log-

209   scaled CTD network features.

210

211    The 239-dimensional prior knowledge space was then projected onto a two-dimensional space

212    using principal component analysis (PCA) and clustered using K-means (K=3). Genes that were

213    furthest from the centroids (i.e., highest contributing ones) of the K-means clusters were more

214    enriched with pathways and gene-chemical-disease interactions (see **Supplemental Information**

215    **1**). Thus, we used the Euclidean distance of genes from the cluster centroids to rank genes based

216    on the prior knowledge space. The ranked list was used to generate prior scores such that the first

217    ranked gene would have a prior score of 100% and the last ranked gene would have a prior score

218    close to 0%. The computational steps for computing the prior score are shown in **Supplemental**

219    **Information 1**. Although the focus was on prioritizing 1000 genes, at this stage of building the

220    prior knowledge, it was necessary to collect information for all potentially relevant genes. Thus,

221    this was done for 22,336 genes.

222

223    **Phase II: Network clustering for relevant grouping of genes**

224    In this phase, we re-weighted the interactions in the co-expression network based on the prior

225    knowledge space and then detected clusters of highly connected genes in the updated network.

226

227    <u>Step 4: Re-weighting co-expression scores (Bayesian) and applying graph clustering</u>

228    In a Bayesian fashion, the pairwise connections between genes in the co-expression network

229    were re-weighted by multiplying the correlation with the mean prior score. For example, given $P$

230    $(A)$ and $P(B)$ as prior scores of genes A and B, the correlation score $S(A, B)$ is re-weighted as

231    follows (Eq. 1):

232

233    $S(A, B)_{new} = S(A, B) * ((P(A) + P(B))/2)$ (1)

234

235    After re-weighting the connections, we detected clusters of highly connected genes using the

236    Markov Cluster Algorithm (MCL) (**Figure 1**, part c) (Van Dongen & Abreu-Goodger 2012). The

237    MCL approach groups together nodes with strong edge weights and then simulates a random

238    flow through a network to find more related groups of genes based on the flow's intensity of

239    movement. It does not require the number of clusters to be pre-specified. An inflation parameter

240    controls the granularity of the output clustering and several values within a recommended range

241    (1.2-5.0) were tried (Van Dongen & Abreu-Goodger 2012). After running several experiments

242    and optimizing for the granularity of the clustering, the inflation parameter was set to 3.3, which

243    generated 258 clusters that consisted of 11,210 genes.  The average number of genes in each

244    cluster was 43.4 with the min-max ranging from 1 to 8,423.

245

246    **Phase III: Gene selection and prioritization**

247    The goal of phase III was to select the top genes from each cluster to form T1000 (step 5), and

248    then produce a final ranking of the 1000 selected genes (step 6).

249

250    Step 5: building local prediction models for each cluster

251  From the training microarray datasets, specific samples were labelled binary as "dysregulated" or

252  "non-dysregulated". Dysregulated refers to exposure cases with potential toxic outcomes and

253  non-dysregulated included controls and exposures with non-toxic outcomes. For the *in vitro*

254  datasets, gene expression changes were associated with lactate dehydrogenase (LDH) activity

255  (%). The activity of LDH, which serves as a proxy for cellular injury or dysregulation, was

256  binarized such that values above 105% and below 95% were considered "dysregulated". While

257  conservative, we note that these cut-off values were situated around the 5% and 95% marks of

258  the LDH distribution curve (see **Supplemental Figure S1** and **Supplemental Information S2**

259  for more details).

260
261  For the *in vivo* datasets (kidney and liver datasets from Open TG-GATEs), gene expression

262  changes were associated with histopathological measures. The magnitude of pathologies was

263  previously annotated into an ordinal scale: present, minimal, slight, moderate and severe

264  (Igarashi et al. 2014a). This scale was further reduced into a binary classification with the first

265  three levels considered "non-dysregulated" while the latter two were considered "dysregulated".

266

267  For each of the 258 gene clusters, random forest (RF) classifiers were used to rank genes based

268  on their ability to separate changes in gene expression labelled as "dysregulated" from those

269  labelled "non-dysregulated", using the Gini impurity index of classification (Nguyen et al. 2013;

270  Qi 2012; Tolosi & Lengauer 2011). RF is one of the most widely used solutions for feature

271  ranking, and as an ensemble model, it is known for its stability (Chan & Paelinckx 2008). In

272  order to cover more biological space and ensure selected genes represent the whole

273  transcriptome, a different RF classifier is built for each cluster and used to select representative

274  genes (Sahu & Mishra 2012).

275

276    We selected the top genes from each cluster based on the performance of the RF classifier. For

277    example, when selecting the 1,000 top genes from two clusters (A and B), if the cross-validation

278    prediction accuracy estimated for models A and B were 60% and 55%, respectively, then 522

279    ((60%/(60%+55%))*1000) and 478 ((55%/(60%+55%))*1000) genes would be selected from

280    clusters A and B. However, if cluster A contained only 520 genes, the remaining two genes

281    would be taken from group B, if possible. So, the cluster size is only used if it contains

282    insufficient genes. We repeated this process until 1000 genes were selected.

283

284    <u>Step 6: building a global prediction model using representative genes from each cluster</u>

285    After choosing top $k$ genes from each cluster, we aggregated them into a single list of 1000 genes

286    and built a final RF model to get a global ranking of the genes. We refer to this final ranked list

287    as T1000 (see **Supplemental Table S1** for a full list of selected genes and summary annotation;

288    see **Supplemental Information S3** for the cluster assignment of the genes).

289

290    **Phase IV: External testing and performance evaluation**

291    The goal of phase IV was to test the performance of the T1000 gene set using external datasets,

292    and thus transition from gene selection activities to ones that focus on the evaluation of T1000.

293

294    <u>Step 7: Dose-response analysis with an external dataset</u>

295    Overall, the aim of the evaluation was to assess the ability of T1000 gene sets to predict apical

296    outcomes according to previously published methods (Farmahin et al. 2017). Additionally, we

297    repeated step 4 of the T1000 approach to select the top 384 (T384; i.e., a number conducive to

298    study in a QCPR microplate format as per the EcoToxChip project; (Basu et al. 2019)) and 1,500

299    (T1500 see **Supplemental Information S4**; i.e., a number pursued in other endeavours like

300    S1500) genes to investigate the effect of gene set size on apical outcome prediction.

301

302    Raw gene expression data (CEL files) for the dose-response dataset were downloaded from GEO

303    (Accession No. GSE45892), organized into chemical-exposure-duration treatment groups, and

304    normalized using the RMA method (Irizarry et al. 2003a). Only expression measurements

305    corresponding to genes in the T1000 gene (or T384 and T1500) set were retained, resulting in

306    reduced gene expression matrices for each treatment group ($t = 24$). The reduced gene

307    expression matrices were analyzed using BMDExpress 2.0 to calculate a toxicogenomic

308    benchmark dose ($BMD_t$) for each treatment group (Yang et al. 2007). Here, the $BMD_t$ was

309    calculated as the dose that corresponded to a 10% increase in gene expression compared to the

310    control (Farmahin et al. 2017). Within BMDExpress 2.0, genes were filtered using one-way

311    ANOVA (FDR adjusted p-value cut-off = 0.05). A $BMD_t$ was calculated for each differentially

312    expressed gene by curve fitting with exponential (degree 2-5), polynomial (degree 2-3), linear,

313    power, and Hill models. For each gene, the model with the lowest Akaike information criterion

314    (AIC) was used to derive the $BMD_t$.

315

316    The $BMD_t$s from individual genes were used to determine a treatment group-level $BMD_t$ using

317    functional enrichment analysis with Reactome pathways (Farmahin et al. 2017).  Note, we chose

318    here to functionally enrich with Reactome since we utilized KEGG previously to derive the

319    T1000 list. After functional enrichment analysis, significantly enriched pathways (p-value <

320    0.05) were filtered such that only pathways with > 3 genes and > 5% of genes in the pathway

321  were retained. The treatment group-level $BMD_t$ was calculated by considering the mean gene-

322  level $BMD_t$ for each significantly enriched pathway and selecting the lowest value. If there were

323  no significantly enriched pathways that passed all filters, no $BMD_t$ could be determined for that

324  treatment group. The similarity of the $BMD_t$ to the benchmark dose derived from apical

325  outcomes ($BMD_a$) was assessed by calculating the $BMD_t/BMD_a$ ratio and the correlation

326  between $BMD_t$ and $BMD_a$ for all treatment groups (Farmahin et al. 2017). Following the same

327  procedures, $BMD_t/BMD_a$ ratio and correlation statistics were determined from genes belonging

328  to L1000, S1500, and Linear Models for Microarray Data (Limma) (Smyth 2005) to provide a

329  reference for the performance of T1000 genes.

330

331  <u>Step 8: Prediction accuracy analysis with an external dataset</u>

332  In this step, we applied five supervised machine learning methods to the TG-GATES rat kidney

333  *in vivo* dataset, with the objective to predict which exposures caused significant "dysregulation",

334  according to the criteria defined in step 4. This dataset was purposefully not used earlier when

335  deriving T1000 so that it could serve later as a validation and testing dataset. The five machine

336  learning models used were K-nearest neighbors (KNN; K = 3) (Cover & Hart 1967), Decision

337  Trees (DT), Naïve Bayes Classifier (NBC), Quadratic Discriminant Analysis (QDA) and

338  Random Forests (RF).

339

340  The performance of each method was evaluated with five-fold cross-validation and measured

341  using six different metrics (Equations 2 – 7). TP represents the number of true positives, FP the

342  number of false positives, TN the number of true negatives and FN the number of false

343  negatives. The $F_1$ score (also called the balanced F-score)  is a performance evaluation measure

344     that computes the weighted average of sensitivity and precision (He & Garcia 2009), and is well-

345     suited for binary classification models. The $F_{0.5}$ score (Davis & Goadrich 2006; Maitin-Shepard

346     et al. 2010; Santoni et al. 2010) is another summary metric that gives twice as much weight to

347     precision than sensitivity. The evaluation was performed on a Linux based workstation with 16

348     cores and 64 GB RAM for processing the data and running the experiments.

349

350     $sensitivity = TP/(TP + FN)$         (2)

351     $specificity = TN/(TN + FP)$         (3)

352     $precision = TP/(TP + FP)$         (4)

353     $GMean = \sqrt{sensitivity \times specificity}$        (5)

354     $F_1 Score = 2 \times \frac{precision \times sensitivity}{precision + sensitivity}$       (6)

355     $F_{0.5} Score = 1.25 \times \frac{precision \times sensitivity}{0.25 \times precision + sensitivity}$      (7)

356

357     The performance of the T1000 gene list was evaluated by comparing it to the performance of

358     randomly selected genes, the top differentially expressed genes, and other notable gene sets. For

359     the random gene set, we generated a list of 1000 random genes, out of 22,336 genes, three times

360     and reported the best. For the differentially expressed gene set, we selected the 1000 top-ranked

361     genes based on analyzing the rat kidney dataset with Limma (Smyth 2005). Finally, to

362     benchmark the performance of T1000 against other notable gene sets, we considered S1500

363     (Merrick et al. 2015) and L1000 (Subramanian et al. 2017).

## Results

The genes comprising T1000 cover a wide biological space of toxicological relevance. For
illustration, co-expression networks, before and after applying Steps 2 and 3 (i.e., networks built
on the Open TG-GATEs data that are subsequently updated with prior information from KEGG,
MSigDb, and CTD), are shown in **Figure 2**. In part (a) of **Figure 2**, a sample co-expression
network composed of 150 genes (i.e., 150 for visualization purposes only; of the 11,210 genes
identified) has, in general, similar color and size of all the nodes of the network. While this
covers a broad toxicological space, it does not necessarily identify or prioritize the most
important genes. After subjecting the data to steps 2 and 3, two clusters of genes with different
node sizes and colors were identified (**Figure 2b**). Through this refined network, we then
applied a prediction model to each cluster to identify the most representative genes resulting in
the final co-expression network of the T1000 genes (**Figure 2c**). The complete list of T1000
genes with their gene symbols and descriptions, as well as their regulation state (up- or down-
regulated) is provided in **Supplemental Table S1**.

379    To understand the biological space covered by T1000, we analyzed T1000's top enriched

380    Reactome pathways (as KEGG was used to develop T1000). Reactome is a manually curated

381    knowledgebase of human reactions and pathways with annotations of 7,088 protein-coding genes

382    (Croft et al. 2014). Visual examination of the Reactome enrichment map (**Figure 3**) reveals that

383    'biological oxidations' (largest circle in **Figure 3**) contained the most enriched pathways

384    followed by 'fatty acid metabolism'. This is logical given that xenobiotic and fatty acid

385    metabolism, mediated by cytochrome P450 (CYP450) enzymes, feature prominently across the

386    toxicological literature (Guengerich 2007) (Hardwick 2008).

387

388    **Evaluation of T1000 to predict apical outcomes**

389    $BMD_t$ analysis of the dose-response dataset was performed with the T1000 gene list and the

390    BMDExpress software program (Yang et al. 2007). The maximum number of BMDs calculated

391    was 21 because for three of the experimental groups a $BMD_a$ (benchmark dose, apical outcome)

392    did not exist due to a lack of observed toxicity (**Table 3**). The T384 gene set performed similarly

393    with Limma; however, increasing the size of this gene set to T1000 resulted in performance

394    evaluation metrics that rivaled that of all other gene sets of the same size or larger (L1000,

395    Limma, and S1500). Further increasing the size of T1000 to T1500 did not increase the

396    performance as the correlation slightly decreased while the average ratio of $BMD_t/BMD_a$ got

397    slightly closer to one (**Figure 4**).

398

399    In a second validation study, we applied T1000 to study the Rat Genome 230 2.0 Array for

400    Kidney dataset from the Open TG-GATEs program.  This dataset was not included in any model

401    training or parameter tuning steps. This helped to establish another external validation of T1000

402     in terms of its generalized ability to predict apical outcomes for datasets derived from different

403     tissues. When compared to baseline gene sets mapped using Limma and L1000, T1000 achieved

404     a relative improvement of the $F_1$Score by 6% and 17%, respectively, thus outperforming

405     comparison gene sets (**Table 4, Figure 5)**. When considering the absolute difference of $F_1$Score

406     between T1000 and the second best (i.e., Limma), T1000 achieved an improvement of 1.2%. The

407     improvement was 1.5% for $F_{0.5}$Score confirming that T1000 led to fewer false positive

408     predictions. In the context of high throughput screening, such small improvements in $F_1$Score or

409     $F_{0.5}$Score may represent large cost savings (Soufan et al. 2015a) as false positives may lead to

410     added experiments that would otherwise be unnecessary. Detailed performance scores of each

411     individual machine learning model are provided in **Supplemental Table S2**. Please refer to

412     **Supplemental Information S5** for more comparisons including expression space visualization

413     using PCA and gene set coverage evaluation.

414

## Discussion & Conclusions

416     There is great interest across the toxicological and regulatory communities in harnessing

417     transcriptomics data to guide and inform decision-making (Basu et al. 2019; Council 2007;

418     ECHA 2016; Mav et al. 2018; Thomas et al. 2019).  In particular, transcriptomic signatures hold

419     great promise to identify chemical-specific response patterns, prioritize chemicals of concern,

420     and predict quantitatively adverse outcomes of regulatory concern, in a cost-effective manner.

421     However, the inclusion of full transcriptomic studies into standard research studies faces

422     logistical barriers and bioinformatics challenges, and thus, there is interest in the derivation and

423     use of reduced but equally meaningful gene sets.

424

425    Here we outlined a systematic, data-driven approach to identify highly-responsive genes from

426    toxicogenomics studies.  From this, we prioritized a list of 1,000 genes termed the T1000 gene

427    set. We demonstrated the applicability of T1000 to 7,172 expression profiles, showing great

428    promise in future applications of this gene set to toxicological evaluations.  Our approach to

429    select T1000 followed the same rationale of how the LINCS program derived the L1000 dataset

430    (Liu et al. 2015), though here we purposefully included additional steps to bolster the

431    toxicological relevance of the resulting gene set.  Generating a list of ranked genes based on

432    toxicologically relevant input data and prior knowledge is another key feature of T1000.

433    There are some limitations associated with our current study. For instance, the co-expression

434    network was based on data from the Open TG-GATEs program. While this is arguably the

435    largest toxicogenomics resource available freely, the program is founded on one *in vivo* model

436    (rat), two *in vitro* models (primary rat and human hepatocytes), 170 chemicals that are largely

437    drugs, and microarray platforms. Thus, there remain questions about within- and cross- species

438    and cell type differences, the environmental relevance of the tested chemicals, and the biological

439    space captured by the microarray. The multi-pronged and -tiered bioinformatics approach was

440    designed to yield a toxicologically robust gene set, and the approach can be ported to other

441    efforts that are starting to realize large toxicogenomics databases such as our own EcoToxChip

442    project (Basu et al. 2019). In addition, our approach in selecting T1000 genes was purely data-

443    driven without considering input from scientific experts as was done by the NTP to derive the

444    S1500 gene set (Mav et al. 2018).  It is unclear how such gene sets (e.g., T1000, S1500) will be

445    used by the community and under which domains of applicability, and thus there is a need to

446    perform case studies in which new approach methods are compared to traditional methods

447    (Kavlock et al. 2018).

448

449  The toxicology community still favors using generic bioinformatics resources, such as KEGG

450  and Reactome, to help organize genes though these are not necessarily applicable to most

451  toxicological use cases. Here we externally validated T1000 against two *in vivo* datasets of

452  toxicological prominence (a kidney dataset of 308 experiments on 41 chemicals from Open TG-

453  GATEs and a dose-response study of 30 experiments on six chemicals (Thomas et al. 2013). We

454  compared the performance of T1000 against existing gene sets (Limma, L1000 and S1500) as

455  well as panels of randomly selected genes.  In doing so, we demonstrate T1000's versatility as it

456  is predictive of apical outcomes across a range of conditions (e.g., *in vitro* and *in vivo*, dose-

457  response, multiple species, tissues, and chemicals), and generally performs as well, or better than

458  other gene sets available. Our approach represents a promising start to yield a toxicologically-

459  relevant gene set.  We hope that future efforts will start to use and apply T1000 in a diverse

460  range of settings, and from these we can then start to make updates to the composition of the

461  T1000 gene set based on improved understanding of its performance characteristics and user

462  experiences.

464 **Supplemental data**

465 Supplemental data are available at PeerJ online.

466

480

481 **References**

482 Alshahrani M, Soufan O, Magana-Mora A, and Bajic VB. 2017. DANNP: an efficient artificial
483         neural network pruning tool. *PeerJ Computer Science* 3:e137.
484 Andersen ME, and Krewski D. 2009. Toxicity testing in the 21st century: bringing the vision to
485         life. *Toxicol Sci* 107:324-330. 10.1093/toxsci/kfn255
486 Basu N, Crump D, Head J, Hickey G, Hogan N, Maguire S, Xia J, and Hecker M. 2019.
487         EcoToxChip: A next-generation toxicogenomics tool for chemical prioritization and
488         environmental management. *Environ Toxicol Chem* 38:279-288. 10.1002/etc.4309
489 Budinska E, Popovici V, Tejpar S, D'Ario G, Lapique N, Sikora KO, Di Narzo AF, Yan P,
490         Hodgson JG, Weinrich S, Bosman F, Roth A, and Delorenzi M. 2013. Gene expression

491          patterns unveil a new level of molecular heterogeneity in colorectal cancer. *J Pathol*
492          231:63-76. 10.1002/path.4212
493  Chan JC-W, and Paelinckx D. 2008. Evaluation of Random Forest and Adaboost tree-based
494          ensemble classification and spectral band selection for ecotope mapping using airborne
495          hyperspectral imagery. *Remote Sensing of Environment* 112:2999-3011.
496  Council NR. 2007. *Toxicity testing in the 21st century: a vision and a strategy*: National
497          Academies Press.
498  Cover TM, and Hart PE. 1967. Nearest neighbor pattern classification. *Information Theory,*
499          *IEEE Transactions on* 13:21-27.
500  Croft D, Mundo AF, Haw R, Milacic M, Weiser J, Wu G, Caudy M, Garapati P, Gillespie M,
501          Kamdar MR, Jassal B, Jupe S, Matthews L, May B, Palatnik S, Rothfels K, Shamovsky
502          V, Song H, Williams M, Birney E, Hermjakob H, Stein L, and D'Eustachio P. 2014. The
503          Reactome pathway knowledgebase. *Nucleic Acids Res* 42:D472-477.
504          10.1093/nar/gkt1102
505  Davis AP, Grondin CJ, Johnson RJ, Sciaky D, King BL, McMorran R, Wiegers J, Wiegers TC,
506          and Mattingly CJ. 2017. The comparative toxicogenomics database: update 2017. *Nucleic*
507          *acids research* 45:D972-D978.
508  Davis J, and Goadrich M. 2006. The relationship between Precision-Recall and ROC curves.
509          Proceedings of the 23rd international conference on Machine learning: ACM. p 233-240.
510  ECHA. 2007. Understanding REACH. *Available at*
511          *https://echa.europa.eu/regulations/reach/understanding-reach* (accessed April 11, 2019.
512  ECHA. 2016. New Approach Methodologies in Regulatory Science. In: Agency EC, editor.
513          European Chemicals Agency (ECHA). Helsinki.
514  Farmahin R, Williams A, Kuo B, Chepelev NL, Thomas RS, Barton-Maclaren TS, Curran IH,
515          Nong A, Wade MG, and Yauk CL. 2017. Recommended approaches in the application of
516          toxicogenomics to derive points of departure for chemical risk assessment. *Archives of*
517          *toxicology* 91:2045-2065.
518  Gautier L, Cope L, Bolstad BM, and Irizarry RA. 2004. affy--analysis of Affymetrix GeneChip
519          data at the probe level. *Bioinformatics* 20:307-315. 10.1093/bioinformatics/btg405
520  Guengerich FP. 2007. Mechanisms of cytochrome P450 substrate oxidation: MiniReview. *J*
521          *Biochem Mol Toxicol* 21:163-168.
522  Haider S, Black MB, Parks BB, Foley B, Wetmore BA, Andersen ME, Clewell RA, Mansouri K,
523          and McMullen PD. 2018. A Qualitative Modeling Approach for Whole Genome
524          Prediction Using High-Throughput Toxicogenomics Data and Pathway-Based Validation.
525          *Front Pharmacol* 9:1072. 10.3389/fphar.2018.01072
526  Hardwick JP. 2008. Cytochrome P450 omega hydroxylase (CYP4) function in fatty acid
527          metabolism and metabolic diseases. *Biochem Pharmacol* 75:2263-2275.
528          10.1016/j.bcp.2008.03.004
529  He H, and Garcia EA. 2009. Learning from imbalanced data. *Knowledge and Data Engineering,*
530          *IEEE Transactions on* 21:1263-1284.
531  Igarashi Y, Nakatsu N, Yamashita T, Ono A, Ohno Y, Urushidani T, and Yamada H. 2014a.
532          Open TG-GATEs - Pathological items. *Available at*
533          *https://dbarchive.biosciencedbc.jp/en/open-tggates/data-12.html*2017).
534  Igarashi Y, Nakatsu N, Yamashita T, Ono A, Ohno Y, Urushidani T, and Yamada H. 2014b.
535          Open TG-GATEs: a large-scale toxicogenomics database. *Nucleic acids research*
536          43:D921-D927.

537 Irizarry RA, Bolstad BM, Collin F, Cope LM, Hobbs B, and Speed TP. 2003a. Summaries of
538         Affymetrix GeneChip probe level data. *Nucleic acids research* 31:e15-e15.
539 Irizarry RA, Hobbs B, Collin F, Beazer-Barclay YD, Antonellis KJ, Scherf U, and Speed TP.
540         2003b. Exploration, normalization, and summaries of high density oligonucleotide array
541         probe level data. *Biostatistics* 4:249-264. 10.1093/biostatistics/4.2.249
542 Kanehisa M, Araki M, Goto S, Hattori M, Hirakawa M, Itoh M, Katayama T, Kawashima S,
543         Okuda S, and Tokimatsu T. 2007. KEGG for linking genomes to life and the
544         environment. *Nucleic acids research* 36:D480-D484.
545 Kanehisa M, and Goto S. 2000. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic acids*
546         *research* 28:27-30.
547 Kavlock RJ, Bahadori T, Barton-Maclaren TS, Gwinn MR, Rasenberg M, and Thomas RS. 2018.
548         Accelerating the Pace of Chemical Risk Assessment. *Chem Res Toxicol* 31:287-290.
549         10.1021/acs.chemrestox.7b00339
550 Knudsen TB, Keller DA, Sander M, Carney EW, Doerrer NG, Eaton DL, Fitzpatrick SC,
551         Hastings KL, Mendrick DL, Tice RR, Watkins PB, and Whelan M. 2015. FutureTox II:
552         in vitro data and in silico models for predictive toxicology. *Toxicol Sci* 143:256-267.
553         10.1093/toxsci/kfu234
554 Liberzon A, Birger C, Thorvaldsdottir H, Ghandi M, Mesirov JP, and Tamayo P. 2015a. The
555         Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell Syst* 1:417-
556         425. 10.1016/j.cels.2015.12.004
557 Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, and Tamayo P. 2015b. The
558         molecular signatures database hallmark gene set collection. *Cell systems* 1:417-425.
559 Liu C, Su J, Yang F, Wei K, Ma J, and Zhou X. 2015. Compound signature detection on LINCS
560         L1000 big data. *Mol Biosyst* 11:714-722. 10.1039/c4mb00677a
561 Maitin-Shepard J, Cusumano-Towner M, Lei J, and Abbeel P. 2010. Cloth grasp point detection
562         based on multiple-view geometric cues with application to robotic towel folding.
563         Robotics and Automation (ICRA), 2010 IEEE International Conference on: IEEE. p
564         2308-2315.
565 Mav D, Shah RR, Howard BE, Auerbach SS, Bushel PR, Collins JB, Gerhold DL, Judson RS,
566         Karmaus AL, and Maull EA. 2018. A hybrid gene selection approach to create the
567         S1500+ targeted gene sets for use in high-throughput transcriptomics. *PLoS One*
568         13:e0191105.
569 Merrick BA, Paules RS, and Tice RR. 2015. Intersection of toxicogenomics and high throughput
570         screening in the Tox21 program: an NIEHS perspective. *International journal of*
571         *biotechnology* 14:7-27.
572 Necsulea A, Soumillon M, Warnefors M, Liechti A, Daish T, Zeller U, Baker JC, Grutzner F,
573         and Kaessmann H. 2014. The evolution of lncRNA repertoires and expression patterns in
574         tetrapods. *Nature* 505:635-640. 10.1038/nature12943
575 Nguyen C, Wang Y, and Nguyen HN. 2013. Random forest classifier combined with feature
576         selection for breast cancer diagnosis and prognostic. *Journal of Biomedical Science and*
577         *Engineering* 6:551.
578 NTP. 2018. High-Throughput Transcriptomics and the S1500+ Gene Set Strategy. *Available at*
579         *https://ntp.niehs.nih.gov/results/tox21/s1500-gene-set-consensus-strategy-*
580         *index.html*2018).
581 Qi Y. 2012. Random forest for bioinformatics. *Ensemble machine learning*: Springer, 307-323.

582    Sahu B, and Mishra D. 2012. A novel feature selection algorithm using particle swarm
583         optimization for cancer microarray data. *Procedia Engineering* 38:27-31.
584    Santoni FA, Hartley O, and Luban J. 2010. Deciphering the code for retroviral integration target
585         site selection. *PLoS computational biology* 6:e1001008.
586    Smyth GK. 2005. Limma: linear models for microarray data. *Bioinformatics and computational
587         biology solutions using R and Bioconductor*: Springer, 397-420.
588    Soufan O, Ba-alawi W, Afeef M, Essack M, Rodionov V, Kalnis P, and Bajic VB. 2015a.
589         Mining Chemical Activity Status from High-Throughput Screening Assays. *PLoS One*
590         10:e0144426. 10.1371/journal.pone.0144426
591    Soufan O, Kleftogiannis D, Kalnis P, and Bajic VB. 2015b. DWFS: a wrapper feature selection
592         tool based on a parallel genetic algorithm. *PLoS One* 10:e0117988.
593         10.1371/journal.pone.0117988
594    Subramanian A, Narayan R, Corsello SM, Peck DD, Natoli TE, Lu X, Gould J, Davis JF, Tubelli
595         AA, Asiedu JK, Lahr DL, Hirschman JE, Liu Z, Donahue M, Julian B, Khan M, Wadden
596         D, Smith IC, Lam D, Liberzon A, Toder C, Bagul M, Orzechowski M, Enache OM,
597         Piccioni F, Johnson SA, Lyons NJ, Berger AH, Shamji AF, Brooks AN, Vrcic A, Flynn
598         C, Rosains J, Takeda DY, Hu R, Davison D, Lamb J, Ardlie K, Hogstrom L, Greenside
599         P, Gray NS, Clemons PA, Silver S, Wu X, Zhao WN, Read-Button W, Wu X, Haggarty
600         SJ, Ronco LV, Boehm JS, Schreiber SL, Doench JG, Bittker JA, Root DE, Wong B, and
601         Golub TR. 2017. A Next Generation Connectivity Map: L1000 Platform and the First
602         1,000,000 Profiles. *Cell* 171:1437-1452 e1417. 10.1016/j.cell.2017.10.049
603    Thomas RS, Bahadori T, Buckley TJ, Cowden J, Deisenroth C, Dionisio KL, Frithsen JB, Grulke
604         CM, Gwinn MR, Singh A, Richard AM, Williams AJ, Deisenroth C, Grulke CM,
605         Patlewicz G, Shah I, Cowden J, Wambaugh JF, Harrill JA, Paul-Friedman K, Houck KA,
606         Gwinn MR, Linnenbrink M, Setzer RW, Sams R, Judson RS, Simmons SO, Knudsen TB,
607         Thomas RS, Lambert JC, Bahadori T, Swank A, Wetmore BA, Ulrich EM, Sobus JR,
608         Phillips KA, Dionisio KL, Isaacs KK, Strynar M, Tornero-Valez R, Newton SR, Buckley
609         TJ, Frithsen JB, Villeneuve DL, Hunter ES, III, Simmons JE, Higuchi M, Hughes MF,
610         Padilla S, Shafer TJ, and Martin TM. 2019. The next generation blueprint of
611         computational toxicology at the U.S. Environmental Protection Agency.
612         10.1093/toxsci/kfz058
613    Thomas RS, Wesselkamper SC, Wang NCY, Zhao QJ, Petersen DD, Lambert JC, Cote I, Yang
614         L, Healy E, and Black MB. 2013. Temporal concordance between apical and
615         transcriptional points of departure for chemical risk assessment. *toxicological sciences*
616         134:180-194.
617    Tolosi L, and Lengauer T. 2011. Classification with correlated features: unreliability of feature
618         ranking and solutions. *Bioinformatics* 27:1986-1994. 10.1093/bioinformatics/btr300
619    van Dam S, Vosa U, van der Graaf A, Franke L, and de Magalhaes JP. 2018. Gene co-expression
620         analysis for functional classification and gene-disease predictions. *Brief Bioinform*
621         19:575-592. 10.1093/bib/bbw139
622    Van Dongen S, and Abreu-Goodger C. 2012. Using MCL to extract clusters from networks.
623         *Bacterial Molecular Networks*: Springer, 281-295.
624    Villeneuve DL, and Garcia-Reyero N. 2011. Vision & strategy: Predictive ecotoxicology in the
625         21st century. *Environ Toxicol Chem* 30:1-8. 10.1002/etc.396
626    Yang L, Allen BC, and Thomas RS. 2007. BMDExpress: a software tool for the benchmark dose
627         analyses of genomic data. *BMC genomics* 8:387.

# Figure 1

Framework of the T1000 approach for gene selection and prioritization.

Phase I includes generating co-expression networks (a) and gene-chemical-toxicity endpoint graphs. Phase II involves re-weighting of the co-expression scores (b) to identify genes in Phase III that contribute more to the clustering (c). Phase IV is an external evaluation of the prioritized gene list.
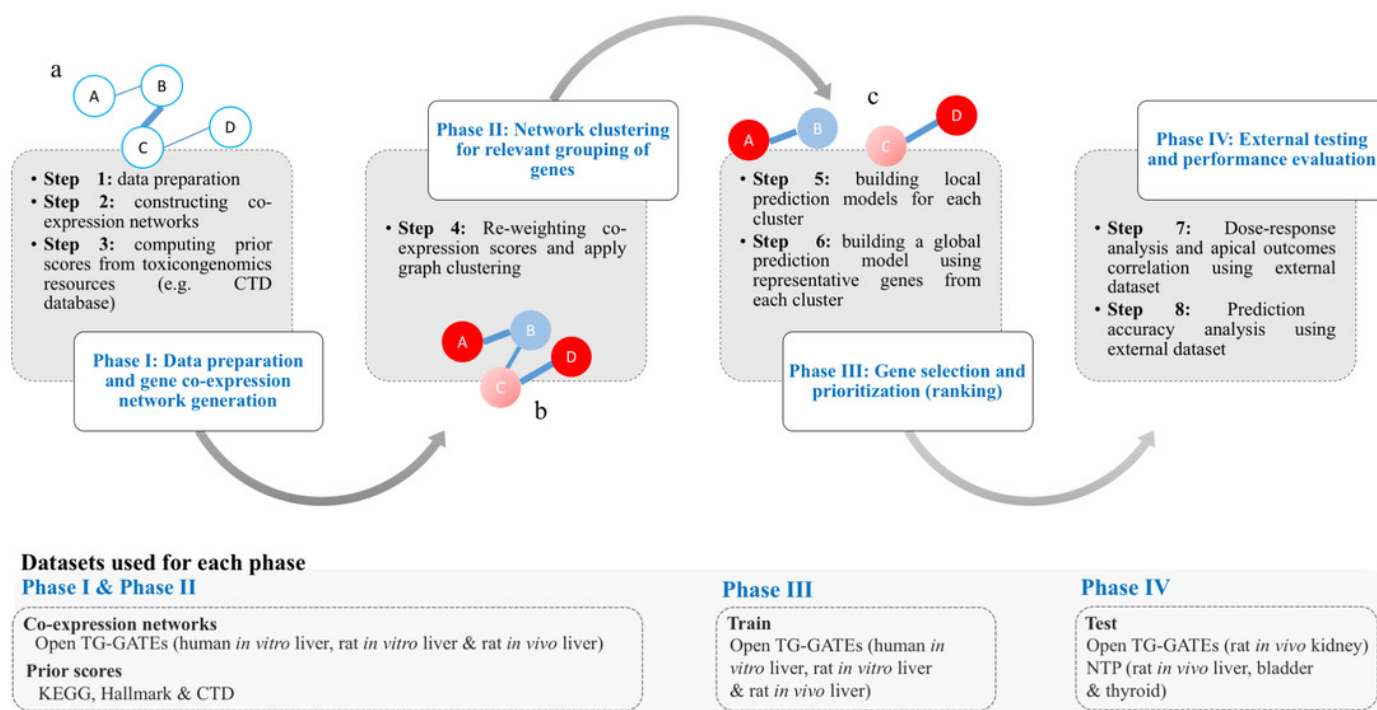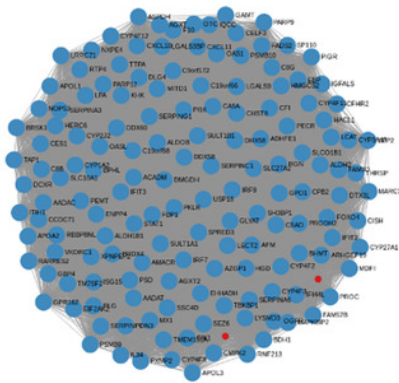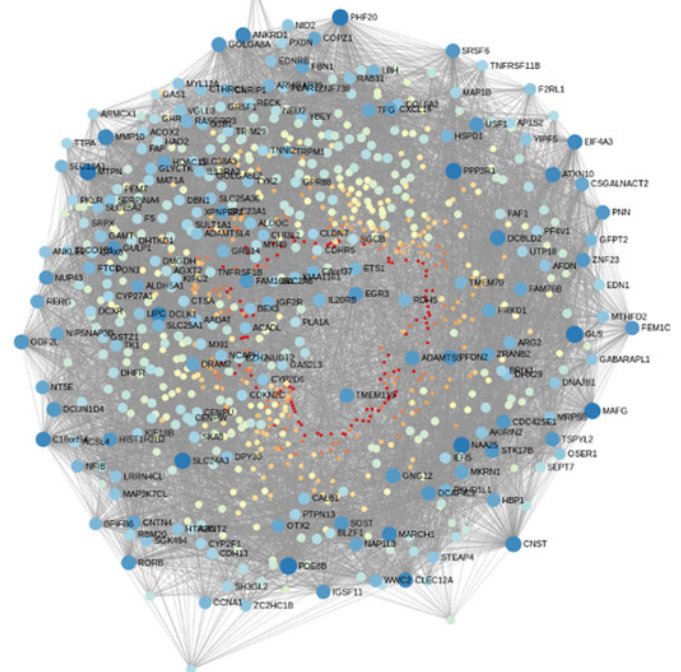
# Figure 2

Visual representation of co-expression networks before and after performing Steps 2 and 3 of the T1000 selection process.

a) Co-expression network of 150 genes after step 1. b) Re-weighted co-expression network of 150 genes after step 3 (same genes as part a). c) Re-weighted co-expression network of T1000 genes after step 5. The color indicates the intensity of betweenness centrality (or amount of influence a gene has along the shortest path of bridged pairs of genes) and size of the node indicates degree (or the number of edges incident to a gene), which are two common metrics to describe the topological structure of a network. A darker blue color reflects higher intensity and a darker red a lower intensity. Yellow indicates a median intensity. A larger size of the node indicates a greater number of connections.

a) Co-expression network of a group of genes before clustering



b) Co-expression network after applying prior weights and clustering



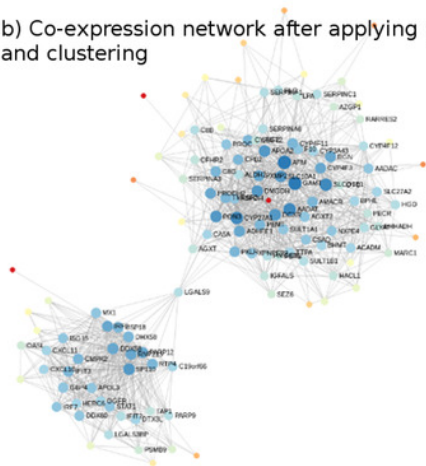c) T1000 generated co-expression network after clustering

# Figure 3

Reactome enrichment map of the T1000 gene set.

The gradient of colors represents p-adjusted of enrichment, where a high intensity red color corresponds to more significance for the enriched term. The different sized circles reflect the number of matched genes between T1000 and the enriched reference gene set. The thickness of the edges indicates the ratio of common genes between the enriched gene sets on both sides of the edge.
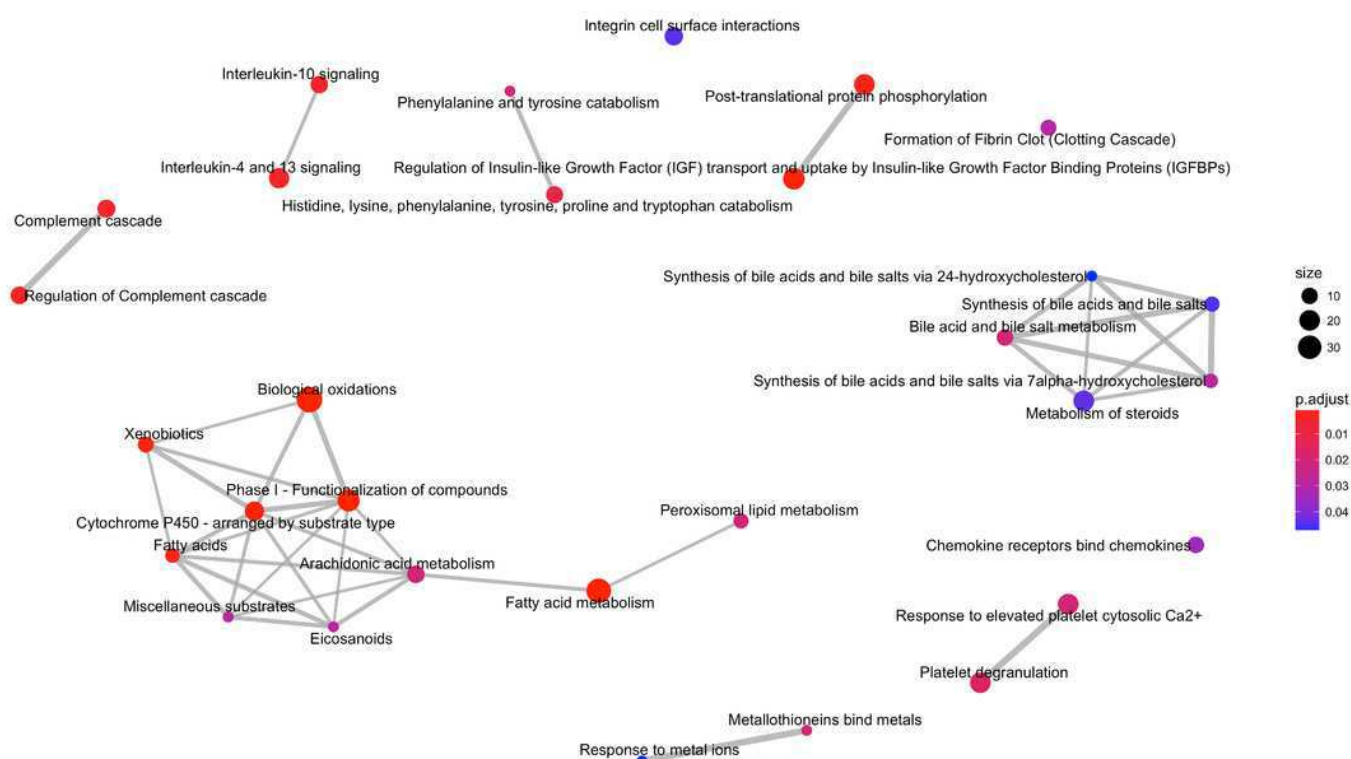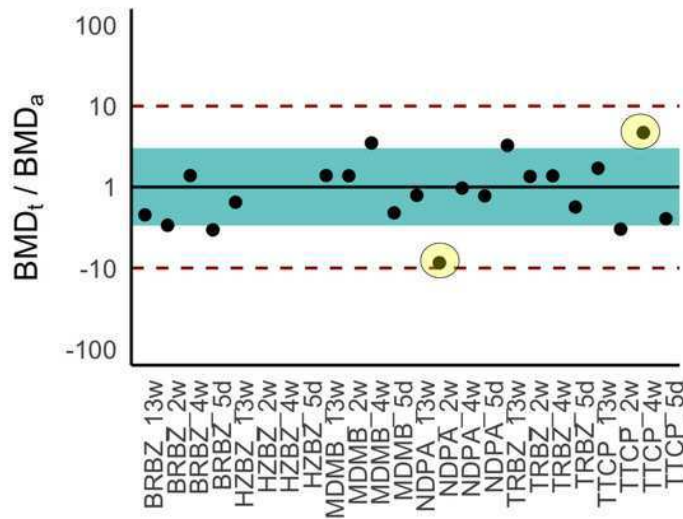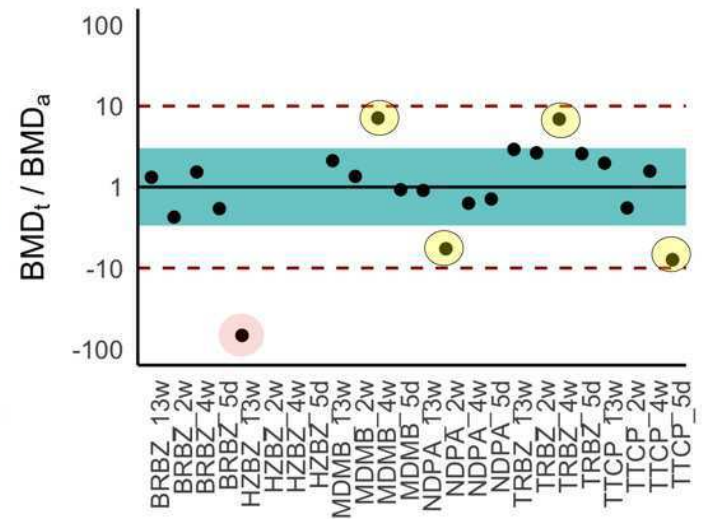
# Figure 4

Ratios of BMDt/BMDa for each experimental group determined with various gene sets as indicated atop the plots.

The limits of the blue rectangular band and dotted lines represent 3-fold and 10-fold of unity, respectively. Ratios could not be calculated for three experimental groups (HZBZ 5 day, 2 week, 4 week) due to a lack of apical outcomes. Red circles represent mean ratios greater than 10-fold, while the yellow ones represent ratios greater than 3-fold. The fewer circles, the more the gene set is indicative of potential relevance to the examined apical endpoints (see Supplementary Figure 2 and 3 for T384 and T1500 plots, respectively).
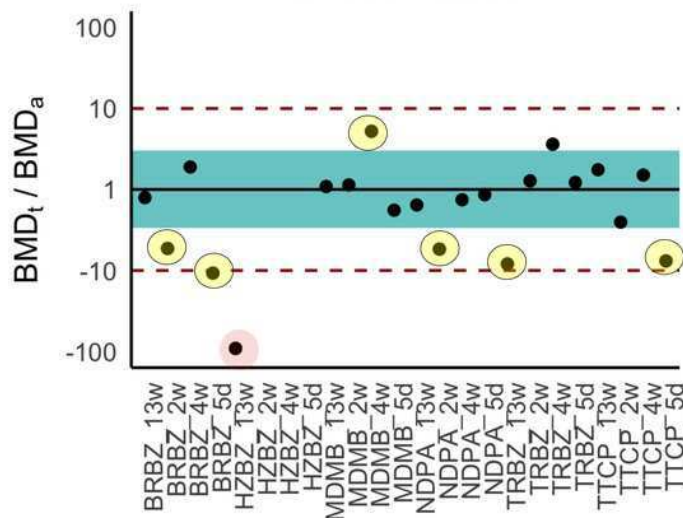
## T1000 Ratios



## L1000 Ratios
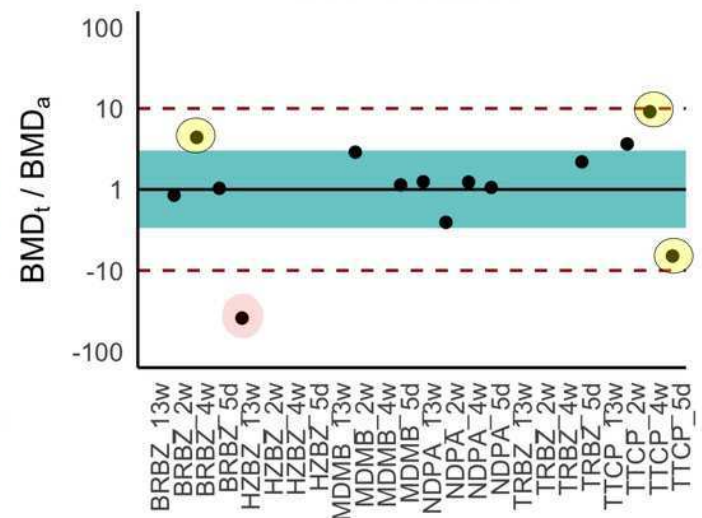


## S1500 Ratios



## Limma Ratios

# Figure 5

Average classification performance over different classification models.

The Rat Kidney dataset was used as an external validation dataset. Refer to step 8 in Phase IV: External testing and performance evaluation for information on F0.5Score, F1Score and GMean.

External Evaluation of Prediction Performance using Rat *(in vivo)* Kidney Dataset
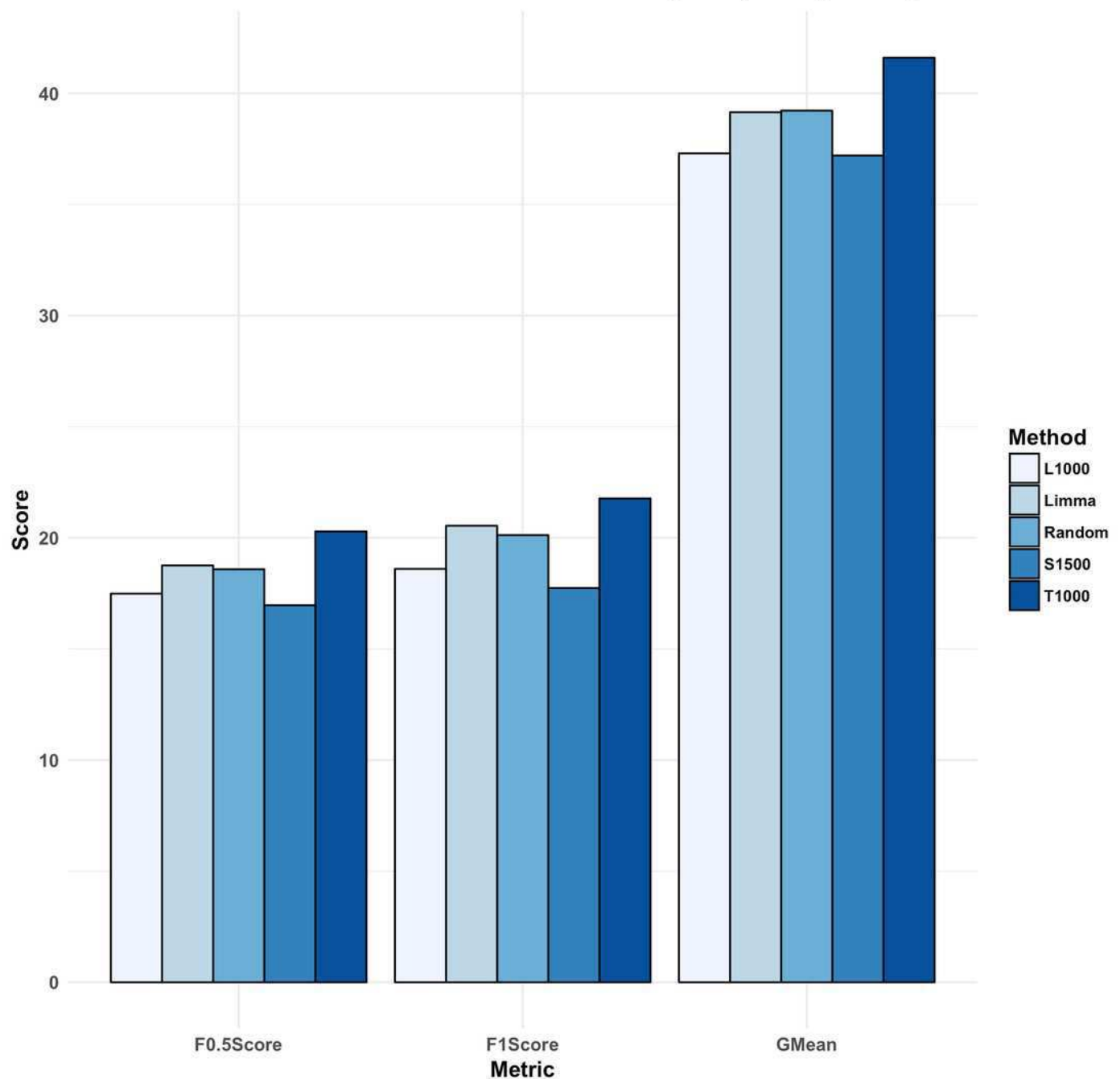
**Table 1**(on next page)

Summary of datasets used in the current study.

Datasets 1-3 were used to develop T1000 (see Phase I, II & III in Methods Section) and datasets 4 and 5 (see Phase IV in Methods Section) were used to evaluate the performance of the gene sets.

1

| Dataset # | Dataset | Organism | Organ | Exposure Type | Number of chemicals | Matrix size (% missing values) | Purpose in Current Study |
|---|---|---|---|---|---|---|---|
| 1 | Open TG-GATEs | Human | Liver | *in vitro* | 158 chemicals | 2,606 experiments x 20,502 genes (8.9%) | Training |
| 2 | Open TG-GATEs | Rat | Liver | *in vitro* | 145 chemicals | 3,371 experiments x 14,468 genes (11.6%) | Training |
| 3 | Open TG-GATEs | Rat | Liver | *in vivo* (single dose) | 158 chemicals | 857 experiments x 14,400 genes (11.5%) | Training |
| 4 | Open TG-GATEs | Rat | Kidney | *in vivo* (single dose) | 41 chemicals | 308 experiments x 14,400 genes (12.2%) | Testing |
| 5 | Dose-response | Rat | Liver, Bladder, Thyroid | *in vivo* (repeated dose) | 6 chemicals | 30 experiments x 14,400 genes (0%) | Testing (external validation) |
| Total | | | | | | 7,172 experiments | |

2

**Table 2**(on next page)

Descriptive comparison of T1000 against existing gene sets.

For the 'selection criteria' column, expression space coverage refers to the goal of finding a subset of genes that would achieve high correlation with the original full set of genes. Pathway coverage refers to finding a subset of genes that cover more pathways in a reference library.

Preprints

1

| Gene set | Selection criteria | Ranked gene list | Species | Data | Approach | Number of genes |
|---|---|---|---|---|---|---|
| L1000 | Expression space coverage | No | Human | L1000 data | PCA and clustering (Data mining) | 978 |
| S1500 (NTP 2018) | Pathway coverage that combines data-driven and knowledge-driven activities | No | Human | Public GEO expression datasets (mainly GEO 3339 gene expression series) | PCA, clustering, and other data-driven steps (Data mining) | 2861 (includes L1000 genes) |
| T1000 | Toxicological relevance using endpoint prediction | Yes | Human and Rat | Open TG-GATEs that is founded on co-expression networks from CTD, KEGG and Hallmark | Co-expression network and prior knowledge (Graph mining). PCA and clustering are used only for the prior knowledge. | 1000 |

2
3

# Table 3 **(on next page)**

Summary of correlation of apical endpoints to 24 experimental groups (6 chemicals x 4 exposure durations).

1

| | T384 (n = 384) | T1000 (n = 1000) | T1500 (n = 1500) | L1000 (n = 976) | S1500 (n = 2861) | Limma (n = 1000) |
|---|---|---|---|---|---|---|
| # of $BMD_t$s | 18 | 21 | 21 | 21 | 21 | 14 |
| Mean ratio ($BMD_t/BMD_a$) | 2.2 | 1.2 | 1.1 | 1.8 | 1.1 | 2.1 |
| Correlation ($BMD_t$, $BMD_a$) | 0.83 ($p < 0.001$) | **0.89** ($p < 0.001$) | 0.83 ($p < 0.001$) | 0.76 ($p < 0.001$) | 0.78 ($p < 0.001$) | 0.73 ($p < 0.01$) |

2

## Table 4<span>(on next page)</span>

Summary comparison of average classification performance using the testing Rat Kidney dataset.*

* Sensitivity would refer to the proportion of expression profiles that are correctly predicted to be dysregulated (or toxic) among all actual dysregulated profiles. (see (Equations 2 – 7) for definition of other performance evaluation metrics).

1

|  | Sensitivity | Specificity | Precision | GMean | F1 Score | F0.5 Score |
|---|---|---|---|---|---|---|
| **T1000** | **25.4%** | 72.1% | **19.5%** | **41.6%** | **21.8%** | **20.3%** |
| **Limma** | 24.6% | 70.8% | 17.8% | 39.2% | 20.5% | 18.8% |
| **Random** | 23.9% | 71.9% | 17.8% | 39.2% | 20.1% | 18.6% |
| **L1000** | 20.9% | 73.0% | 16.8% | 37.3% | 18.6% | 17.5% |
| **S1500** | 19.4% | **73.1%** | 16.5% | 37.2% | 17.7% | 17% |

2