1 **Title:** Cross-organism toxicogenomics with group factor analysis

2

3 **Authors:** Tommi Suvitaival[1], Juuso A. Parkkinen[1], Seppo Virtanen[1], Samuel Kaski[1,2]

4

5 **Affiliations:** [1]Helsinki Institute for Information Technology HIIT, Department of Information and

6 Computer Science, Aalto University; [2]Helsinki Institute for Information Technology HIIT,

7 Department of Computer Science, University of Helsinki

8   **Keywords:** Bayesian modeling, factor modeling, information retrieval, multi-view modeling,

9   toxicogenomics

10

11   **List of abbreviations and acronyms:**

| Abbrebiation | Meaning |
|---|---|
| ATC | Anatomical therapeutic chemical |
| DILI | Drug-induced liver injury |
| FDA | Food and Drug Administration |
| GFA | Group factor analysis |
| GSEA | Gene set enrichment analysis |
| QSAR | Quantitative structure-activity relationship |
| TGP | Japanese Toxicogenomics Project |

12

## Abstract

We investigate the problem of detecting toxicogenomic associations that generalize across organisms, that is, statistical dependencies between transcriptional responses of multiple organisms and toxicological outcomes. We apply an interpretable probabilistic model to detect cross-organism toxicogenomic associations and propose an approach for drug toxicity analysis based on the interactive retrieval of drugs with similar toxicogenomic properties. We show that our approach can give relevant information about the properties of a drug even when direct prediction of toxicity is not feasible. Moreover, we show that a search from a cross-organism database can improve accuracy in the analysis.

## Introduction

Evaluation of potential toxicity of new drugs and other chemical compounds is highly important for safety reasons. The toxic effects of new drugs cannot be tested directly on humans due to the obvious ethical issues, and new drugs thus go through a series of *in silico* and *in vitro* analyses, and then an animal experimentation phase. Organisms from yeast[1] to the worm *C. elegans*[2], zebrafish[3] and murine animals[4] are used in the drug development process, starting with simple organisms and moving towards organisms more similar to humans. Since all toxic effects do not generalize across the model organisms and setups, after the animal studies and even after the drug has entered the market, new toxic side effects are often discovered among the large population of consumers.

The earlier the toxic responses can be detected, the more potential harm can be avoided and resources saved. Computational tools for predictive toxicity have been developed and applied at each stage of the drug development cycle[5,6]. Quantitative structure-activity relationship (QSAR) assessment has traditionally been the most prominent *in silico* toxicity prediction procedure, where toxicological profiles, such as lethal concentrations, are predicted based on structural descriptors of the compounds[7]. Recently, the focus has shifted to identification of critical perturbations in biological pathways that lead to adverse outcomes, based on high-throughput screening methods[8].

### *Toxicogenomics*

Toxicogenomics has emerged in the cross-section of toxicology and bioinformatics, with the aim of finding predictive associations between transcriptomic and toxicological responses[9,10]. The rationale is that drug-treatment transcriptional data consist of various response patterns, some of which are related to drug toxicity. The identification of these toxicity-associated transcriptional response patterns is essential for understanding the molecular mechanisms behind toxicity and for enabling the prediction of toxicity[11]. However, distinguishing toxic adverse effects from intended therapeutic

46   effects and from various types of noise factors, such as batch effects, is highly non-trivial. Moreover,

47   transcriptomic response patterns vary over tissues and cell types, making this more complicated. As

48   toxicogenomic studies are typically performed *in vitro*, it would be important to identify those

49   toxicogenomic associations that generalize to humans as well.

50

51   The ToxCast project[12] is an example of large-scale high-throughput *in vitro* screening for predicting

52   *in vivo* toxicity. The TG-GATEs database from the Japanese toxicogenomics project[13] is another

53   interesting toxicogenomic resource with transcriptional drug-treatment data available from

54   organisms both *in vitro* and *in vivo*. Additionally, the database includes toxic outcome observations

55   such as blood level measurements and observed liver injuries from rats *in vivo*.

56

57   Liver toxicity is among the most common types of drug toxicity in humans[5]. The drug-induced liver

58   injury (DILI) labelings[14] have been designed to describe the risk of hepatotoxicity in humans: The

59   labels are continuously updated as the Food and Drug Administration (FDA) acquires more

60   information about the potential side effects of the drugs on the market. The DILI labels are

61   available for most of the drug compounds with experimental data at the TG-GATEs database.

62   ### *Data translation with machine learning*

63   The next step that follows detection of responses to drug compounds in a model organism is

64   translation of these responses to humans. In this work, we build on the hypothesis that responses

65   shared across organisms are more likely to generalize to humans as well. This is analogous to

66   searching for conserved genomic regions or responses, but on the more abstract level of statistical

67   relationships in the response profiles.

68

69 To detect "conserved responses," we need to examine databases of drug-response experiments from

70 multiple model organisms, or *domains*. The conserved response patterns can then be utilized to

71 make predictions about the human response based on experimental data from model organisms, that

72 is, to carry out *data translation* from one domain to another.

73

74 We define *data translation* as an analogue of language translation: of finding how a phenomenon in

75 one domain or organism is expressed in another, assuming it generalizes across domains, and then

76 predicting it. Data translation is a key part of *translational medicine*, which involves many

77 additional aspects.

78

79 In summary, our goal is to develop machine learning methods for discovering responses conserved

80 across organisms and for generalizing the responses to humans. The generalization of the responses

81 has so far been an unsolved problem. For discovering conserved responses, Le & Bar-Joseph[15] have

82 presented an approach for clustering genes across organisms based on their response patterns.

83 Suitaival *et al.*[16] focused on quantifying the responses to external covariates, such as the drug

84 treatment, that are conserved across organisms. Both of these approaches assume that a group of

85 genes responds to the covariate in a coherent fashion.

86

87 In this article, we assume that drug responses can be modeled as factors, each of which describes a

88 biological process that is disturbed by the treatment. Individual genes may be members of many of

89 these processes and the genes may be different across organisms. Also the level and direction of

90 responses may vary across genes and organisms while still following the abstract conserved pattern.

### *Generative model for cross-organism toxicogenomics*

Inspired by the CAMDA challenge[17], we address the following research questions: (1) Can we associate drug-induced toxicological responses observed in humans or rats to changes observed at the molecular level, and are these associations predictive? (2) Can we find toxicogenomic associations that are conserved across organisms? Could these associations be utilized to replace animal studies with *in vitro* assays?

In other words, we seek simultaneous associations between transcriptional data and toxicological outcome data, and between transcriptional data from multiple organisms. Associations that generalize both across organisms and across levels of biological complexity have the potential of enabling data translation between the molecular level and the organ or population level.

The biological properties and their resemblance to the human vary across the cells extracted from animals grown *in vivo* and cell lines grown *in vitro*. Even though this resemblance to the human is still largely unkown, they all are grown with the purpose of experimenting chemical compounds intended for human use. By taking a data-driven approach to identifying conserved responses, we do not make prior assumptions about the organisms' similarity to the human. To stress these points, we refer to each of the types of biological sample as a model organism, even though a cell line is not an entire representation of the animal from which it is originally extracted from. Moreover, we view a cell line grown *in vitro* as a different model organism than what a cell extract from an animal of the same species grown *in vivo* is.

We propose a generative model-based approach to answer the two research questions. To do this, we make the following modeling assumptions: (1) The data consist of drug-induced transcriptional responses patterns, that is, consistent gene expression changes for a subset of the drugs and genes,

116 and noise from various sources. (2) Drugs may activate multiple response patterns, and the patterns

117 may be partially overlapping in terms of affected genes. (3) We are especially interested in response

118 patterns that are associated with observed toxic outcomes and are conserved across organisms.

119

120 It turns out that a recently introduced model family, group factor analysis[18] (GFA), when applied to

121 toxicogenomic data, matches these assumptions. It is a multi-view model that in an unsupervised

122 fashion detects statistical dependencies between multiple data sets having co-occurring samples. In

123 this context, samples correspond to drug treatments, which are the same in all the data sets. We call

124 the data sets *views*, because they are matched by their samples.

125

126 The associations found by the model are represented by factors that are interpretable in terms of

127 factor loadings of the data variables, in this case genes. This interpretability allows the user to

128 formulate testable hypotheses, for instance about the mechanisms of action of a drug and about their

129 association to toxicological outcomes. The associations can also be used for predicting one data

130 view based on another, for example, predicting toxic outcomes based on transcriptomic responses.

131

132 For cross-organism toxicogenomic analysis, group sparsity is an especially useful feature of GFA.

133 The model can distinguish patterns that are shared across all the data sources from patterns that are

134 specific to a single source or shared by a subset of the sources. In this paper, we will apply GFA to

135 studying biological responses that are conserved across organisms.

## Results

137 We demonstrate the potential of the model to detect responses that generalize across organisms in

138 two practical use cases with the TG-GATEs data[13], consisting of three sets of transcriptional drug-

139 treatment measurements: human *in vitro*, rat *in vitro* and rat *in vivo*. In Case 1, the task is to find

140  associations between transcriptional changes and pathological findings from *in vivo* rat livers. In

141  Case 2, the task is to search for drugs having a similar risk of drug-induced liver injury (DILI) in

142  humans at the population level, based on data about transcriptional changes in model organisms.

### *Case 1: Finding associations between transcriptomic responses and*

144  *pathological findings*

145  In the first case, we are interested in two types of associations to start with: First, associations

146  between the molecular level and the organ-level, and second, molecular-level associations between

147  the different organisms. In order to detect responses that are most likely to generalize to humans, we

148  require both of these constraints to hold for the associations that we focus on. Focusing on these

149  maximally conserved associations will also be beneficial for filtering out structured noise that arises

150  from the laboratory effects and from the properties of the model organisms.

151

152  Applying GFA to the combination of three transcriptomic data sets and pathological findings for rat

153  *in vivo*, we obtain a set of factors that capture the required kind of associations. Each factor is

154  interpretable as a biological process associated with specific pathological findings at the organ-level

155  and is generalized across a subset of the organisms at the molecular level (Figure 1). This result

156  indicates that the model learns biologically meaningful response structure in the transcriptomic data.

157  For example, Factor *B* associates changes in metabolic processes to degeneration in the liver tissue,

158  while Factor *C* associates changes in the cell-cycle to increased mitosis in the liver.

159

160  Although the associations are biologically meaningful, given the small amount of available data,

161  their predictive power is not significant (results not shown; the low power was not due to the

162  method, which was tested additionally using a standard L1-regularized regression model). As more

163  toxicogenomic data accumulates, the predictive power of the associations needs to be revisited.

### *Case 2: Modeling-based data retrieval for human drug toxicity analysis*

Direct prediction of toxicity for a new drug is not a trivial task, but we have demonstrated that the detected conserved associations are biologically meaningful. Predicting the toxicity of a drug on humans is even more difficult due to the lack of direct experimental data. Analyzing drug toxicity in humans is possible indirectly, using available drug toxicity classifications of approved drugs. These data are not perfect, however, as the toxic potential of many drugs has been over-estimated for increased safety[14]. Some drugs have been categorized as risky based on only indirect evidence of other drugs, with similar therapeutic potential or chemical properties, having shown toxic outcomes.

## Interactive toxicity analysis framework

We propose an alternative approach for the risk-analysis of a novel drug by formulating the prediction task as an information retrieval problem. We assume that transcriptomic response data in existing databases of model organism experiments carries relevant information on drug toxicity in humans. The level of relevance may, however, vary across different experimental practices and model organisms. For instance, *in vivo* experiments are likely to be more informative than *in vitro* experiments.

The interactive toxicity analysis takes place through a table-lookup procedure: Given a query compound and a measure of similarity, the expert receives a ranked list of database compounds in the order of the similarity of transcriptomic response. To the extent there are associations between the molecular level and the organ-level, the properties of the top-ranked database compounds are likely to be similar to the query compound. Based on the list, an expert user can then construct a hypothesis about the expected properties of the drug and about the uncertainty around these properties. In an illustrative example of the retrieval result for a query (Table 1), many of the top-ranked drug compounds retrieved from the database are shown to share toxic and therapeutic properties with the query.

189

190     The idea of searching for similar drugs has earlier been introduced as "connectivity mapping"[19] and

191     applied to drug discovery and drug repositioning[20,21]. It has also been applied to drug toxicity

192     analysis[22,23]. Recently, Xing *et al.*[24] introduced an online resource for making queries to the TG-

193     GATEs database. We use the retrieval method behind that tool as one of the two baseline

194     approaches in the experiments that follow. In the connectivity mapping approaches the similarity

195     measure for the retrieval relevance is based on the gene set enrichment[25] computed on the list of the

196     most differentially expressed genes for the query drug. These approaches have either focused on a

197     single cell type or simply averaged over multiple cell types, neglecting the likely differences

198     between organisms.

199

200     We propose to carry out toxicity analysis by modeling-based retrieval that takes into account the

201     translatability of data between different organisms. In particular, we use the GFA to detect shared

202     transcriptomic responses between the three model organisms in the database: human *in vitro*, rat *in*

203     *vitro* and rat *in vivo*. Now, we can examine the similarity in the responses in the lower-dimensional

204     latent space of the model. More importantly, we can focus our examination into the part of the

205     latent space that is shared between the model organisms (details in the section *Material and*

206     *Methods*). The shared latent factors describe the drug-responses that are conserved across the model

207     organisms, and thus are likely to have potential for the generalization to humans as well.

208

209     We evaluate the retrieval using as ground truth the drug-induced liver injury (DILI) label and

210     concern classes[14], as well as more detailed information about the drugs' mechanism of action based

211     on the anatomical therapeutic chemical[26] (ATC) classes. We compare with rank-based connectivity

212     mapping[19] and simple correlation between the differential expression profiles. As a measure of

213     performance, we use mean average precision.

## Retrieval from single-organism database

Transcriptomic drug response data are informative about both the toxicity and mechanisms of action (Figure 2), resulting from off-target and on-target effects of the drug, respectively. For all organisms, types of validation classes and used similarity measures, retrieval based on the transcriptomic database lead to a higher performance than expected by chance. This indicates that the transcriptomic response data on model organisms is informative of the toxicity of the drugs on humans at the population level. However, the results are not conclusive of the relative performance of the individual organisms. Retrieval performance is observed to be almost as sensitive to the choice of the similarity measure as it is to the choice of the organism.

## Retrieval from cross-organism database

We study the potential of cumulating biological information from existing model organism experiments to increase the amount of knowledge that can be extracted from human *in vitro* experiments. We focus on human *in vitro* experiments, because they are more ethical and less expensive than *in vivo* experiments and could potentially replace *in vivo* animal studies in the future.

We examine model-based retrieval performance from a cross-organism database of transcriptional measurements, given a human *in vitro* sample as a query. The results show that retrieval performance is improved by using the cross-organism database of experiments compared to single-organism retrieval, when the retrieval is based on responses conserved across the model organisms (Figure 3). The outcome is consistent on all the three validation classes. This is indirect evidence for the hypothesis that compared to organism-specific responses, conserved responses of model organisms are more likely to generalize to humans at the population level.

## Discussion

We have analyzed drug toxicity using a new machine learning approach that identifies cross-organism toxicogenomic associations. This is a key step towards developing methods for predictive toxicology. The identification of associations that generalize reliably across multiple organisms, especially from in vitro to in vivo, is essential for toxicity analysis. This approach has potential for predicting drug toxicity in humans based on in vitro experiments, thus reducing the need for animal studies *in vivo.*

The TG-GATEs data set with experiments on three model organisms has given us the opportunity to take a data-driven approach for cross-organism toxicogenomics. The group factor analysis model for toxicogenomic responses is flexible about the type of responses: neither genes nor biological pathways are restricted to be the same between the organisms. Minimum two model organisms are needed for identifying conserved responses. A new experiment in one organism can then be generalized via retrieval. The model can operate in the "small $n$, large $p$" regime thanks to the probabilistic approach and the sparsity assumptions.

We have shown how our probabilistic model finds biologically relevant associations between transcriptomic drug responses and pathological findings from rats, and that many of these associations generalize across in vivo and in vitro organisms. However, the predictive performance of these linear associations is very limited, probably due to limited amount of data, as the pathological findings have been observed only for a few rat samples.

Since quantitative linear prediction of toxicological outcomes is limited in performance, we propose an alternative toxicity analysis scheme. It is based on information retrieval, where the task is to search for the most relevant drugs from the database of existing experiments, given a new query

261 drug. Based on the most relevant drugs retrieved, the user can then construct a hypothesis of the

262 toxicity and other properties of the query drug. This can support expert decision making.

263

264 We first studied the retrieval performance using the differential gene expression data only, and

265 confirmed earlier findings[22,23] about the suitability of the retrieval approach to the task of

266 identification of toxic drug compounds. We then showed that when we do retrieval based on cross-

267 organism associations, we were able to improve the retrieval performance, as compared to single-

268 organism retrieval. This indicates that the cross-organism associations detected by the model are

269 relevant for human toxicity and give hope that the *in vivo* animal studies could be replaced with *in*

270 *vitro* studies in the future.

## Materials and Methods

272 We report the pre-processing done for the data before modeling, the model description, and the

273 technical details of the two experiments (Cases 1 & 2). The details of Cases 1 and 2 are described in

274 the subsections *Model-based exploratory analysis* and *Retrieval of relevant experiments*,

275 respectively.

### *Data pre-processing*

277 The data set of the Japanese Toxicogenomics Project (TGP) includes transcriptional data from three

278 model organisms: primary hepatocyte cells from humans and rats grown *in vitro*, and similar cells

279 extracted from rats *in vivo*. The conditions of the experiment can be summarized as three

280 experimental factors: administered drug compound, its dosage and time from the administration of

281 the compound. For the analysis in this work, we selected the subset of experimental factor levels

282 that are observed in all three organisms. This set includes 119 drug compounds administered at two

283 dosage levels (middle & high) and measurements made at two time points after the

284 treatment (8/9 h & 24 h). Histopathology of the liver had been examined from the extracted livers in

285    the rat *in vivo* experiments at the same time points and dosage levels, providing a pathological

286    finding class and severity grading for each sample. The data were downloaded from the website of

287    the CAMDA challenge[27], where the transcriptional observations were provided in a FARMS-

288    summarized[28] format.

289

290    For the modeling task, we considered each treatment – a combination of compound, dose and time –

291    as a single sample in the model. We selected transcriptomic probes, which have non-zero variance

292    across the samples and which appear in all the three transcriptomic microarray data sets. This was

293    done to make the data sets from different organisms balanced in their size in order to allow a fair

294    comparison between the relevant information content in them. However, the model itself does not

295    require the variables of the data sets to be matched and the analysis could alternatively be done on

296    all probes as well.

297

298    We computed the average differential expression of the treated samples against the corresponding

299    control samples. We represented the pathological finding classes for each sample as a grade-

300    weighted count. As the four data matrices (differential gene expression $\mathbf{X}^{\binom{\text{human}}{\text{in vitro}}}$, $\mathbf{X}^{\binom{\text{rat}}{\text{in vitro}}}$ and

301    $\mathbf{X}^{\binom{\text{rat}}{\text{in vivo}}}$, as well as pathological findings $\mathbf{Y}$) are now matched by their samples, we call the

302    matrices different *views* of the data.

### *Model*

304    We have $N$ observation vectors $\mathbf{x}_n^{(m)}$, corresponding to measured transcriptional and toxicological

305    responses to drug treatments indexed as $n = 1, \dots, N$. Observations from one measurement type $m$

306    are concatenated as columns of a data set $\mathbf{X}^{(m)}$. All data sets are matched by co-occurring

307    observations, that is, they can be regarded as *views*. We assume the transcriptomic data contain

308    complex drug-induced response patterns embedded in measurement noise. We are interested in

309 finding these patterns and, more importantly, in associating them to toxic outcomes. Response

310 patterns that are present in multiple views provide valuable information for interpretation and data

311 translation. The task suits well to the problem formulation of group factor analysis[18] (GFA), which

312 learns associations between matched data sets.

313

314 GFA is formulated as a Bayesian latent factor model, where the data are explained by factors. Each

315 observation $\mathbf{x}_n^{(m)}$ from the $m$th view is generated from a multivariate normal distribution

$$\mathbf{x}_n^{(m)} \sim N\big(\mathbf{W}^{(m)}\mathbf{z}_n, \mathbf{\Sigma}^{(m)}\big), \tag{1}$$

316 where $z_n$ are the latent factors for the $n$th observation, $\mathbf{W}^{(m)}$ are the factor loadings for the $m$th

317 view, and the noise covariance matrix is assumed to be diagonal, $\mathbf{\Sigma}^{(m)} = \tau_m^{-1}\mathbf{I}$, with a view-specific

318 precision $\tau_m$. The main task is to learn how factors are associated with the views: each factor

319 describes associations between any combination of the views. Thus, some factors are shared across

320 all the views, some are shared by a subset of the views, and the rest are specific to a single view.

321 For a view $m$ that is not associated with factor $k$, the $k$th column of $\mathbf{W}^{(m)}$ is automatically set to

322 zero by the model. With variables from each view seen as groups, this is equivalent to group-sparse

323 factor loadings.

324

325 GFA learns the associations by employing a group-sparse prior distribution for the factor loadings.

326 That is, each column of $\mathbf{W}^{(m)}$ is generated from a normal distribution

$$\mathbf{W}_{:,k}^{(m)} \sim N\left(\mathbf{0}, \left(\alpha_k^{(m)}\right)^{-1}\mathbf{I}\right), \tag{2}$$

327 where precision $\alpha_k^{(m)}$ is drawn from a gamma prior distribution,

$$\alpha_k^{(m)} \sim Gamma\big(\alpha_0, \beta_0\big), \tag{3}$$

328 with small values for the shape parameters $\alpha_0$ and $\beta_0$. Gamma distribution is conjugate to normal

329 distribution with a known mean. When the prior and the likelihood are conjugate, posterior

330    inference through Gibbs sampling is possible, as the posterior is of the same form as the likelihood

331    and the parameters of the posterior distribution can be directly calculated based on the parameters

332    of the prior and the likelihood. The model learns the sought-for associations for factor $k$ by setting

333    the $\left(\alpha_k^{(m)}\right)^{-1}$ of non-associated views $m$ close to zero, thus pushing all the elements in the factor

334    loadings  for those views jointly to zero. To complete the model description, a conjugate gamma

335    prior,

$$\tau_m \sim Gamma(\alpha_0^\tau, \beta_0^\tau) , \qquad (4)$$

336     is set for the noise precisions, and the latent variables are generated from a normal distribution

$$\mathbf{z}_n \sim N(\mathbf{0}, \mathbf{I}). \qquad (5)$$

337

338    Factors capture response patterns in the observed data, for instance, sets of genes in the

339    transcriptomic views that respond to sets of drug-treatments in a coherent fashion. Some of these

340    patterns are shared across views. Each factor and the corresponding loadings are assumed to

341    represent a biological process and we are interested in interpreting them. Thus, each factor is

342    assumed to be related to a sparse set of drugs and each loading to a sparse set of variables, for

343    example genes. Further, we assume that each drug induces a sparse set of response patterns

344    corresponding to sparsity of $\mathbf{z}_n$. Motivated by these assumptions, we modify the priors for GFA in a

345    way that leads to a more easily interpretable model.

346

347    We extend the plain GFA by assuming that, in addition to the group sparsity, both the factors and

348    the factor loadings are element-wise sparse. With this extension, the GFA model becomes a multi-

349    view biclustering model, generalizing the factor analysis-based multiplicative biclustering model

350    (FABIA)[29] to multiple views of the data. Further, FABIA and GFA with the element-wise sparsity

351    structure extend the Bayesian plaid model[30] from additive responses to multiplicative responses.

352

353  We modify the priors of the GFA model to achieve the element-wise sparsity for the factors and the

354  factor loadings by drawing them both from a two-component mixture distribution. In the mixture,

355  the first component corresponds to a delta distribution $\delta_0$ with a peak at zero, and the second to a

356  normal distribution with a zero mean and an unknown precision. This construction corresponds to a

357  spike-and-slab prior[31,32], where the spike is a delta distribution and the slab is a normal distribution.

358

359  Mathematically, the spike-and-slab prior for the factors is written as

$$z_{n,k} \sim h_{n,k}^{(z)} N\left(\mathbf{0}, \left(\alpha_{n,k}^{(z)}\right)^{-1}\right) + \left(\mathbf{1} - h_{n,k}^{(z)}\right)\delta_0, \tag{6}$$

360  and for the factor loadings as

$$W_{d,k}^{(m)} \sim h_{d,k}^{(m)} N\left(\mathbf{0}, \left(\alpha_{d,k}^{(m)}\right)^{-1}\right) + \left(\mathbf{1} - h_{d,k}^{(m)}\right)\delta_0. \tag{7}$$

361  Binary variables $h_{n,k}^{(z)}$ and $h_{d,k}^{(m)}$ indicate whether $z_{n,k}$ and $W_{d,k}^{(m)}$, respectively, are set to zero or

362  drawn from a normal distribution. The $h_{d,k}^{(m)}$ are drawn from a Bernoulli distribution,

$$h_{d,k}^{(m)} \sim Bernoulli\left(\pi_k^{(m)}\right), \tag{8}$$

363  where the expectation $\pi_k^{(m)}$ is specific to each factor $k$ and view $m$. The $\pi_k^{(m)}$ is drawn from a beta

364  distribution

$$\pi_k^{(m)} \sim Beta(a_0, b_0) \tag{9}$$

365  with shape parameters $a_0$ and $b_0$. The beta prior distribution is conjugate to the Bernoulli

366  distribution, leading to a posterior, which is Bernoulli-distributed. A similar construction is used for

367  the $h_{n,k}^{(z)}$ but now the expectation is shared across observations. When $\pi_k^{(m)}$ is close to zero, the $k$th

368  column of $\mathbf{W}^{(m)}$ is suppressed to zero jointly, implementing group sparsity. We also find shared

369  noise for each view too limiting and instead allow variable-wise independent noise by assuming a

370  non-isotropic diagonal $\mathbf{\Sigma}^{(m)}$ whose elements are drawn independently from a gamma distribution.

371

372 Since all the priors are conjugate, we implement inference using Gibbs sampling. The sampler

373 learns the model for the TG-GATEs data set overnight on a standard desktop computer. A

374 variational Bayesian approximation, presented for the vanilla GFA model earlier[18], may be useful

375 for larger data sets.

### *Model-based exploratory analysis*

377 We study the biological interpretability of the learned associations which are represented by factors

378 of the model. More specifically, we focus on factors that are shared across all the views. In order to

379 do that, we need to define a threshold for a factor to be considered shared by the views. We

380 consider the $k$th factor as shared, if in each of the $m$ views there exists at least one non-zero value in

381 the loadings vector $\mathbf{W}_{:,k}^{(m)}$ of the $k$th factor. In Case 1, we study associations that generalize across

382 the transcriptomic views $\mathbf{X}^{\left(\substack{\text{human} \\ \text{in vitro}}\right)}$, $\mathbf{X}^{\left(\substack{\text{rat} \\ \text{in vitro}}\right)}$ and $\mathbf{X}^{\left(\substack{\text{rat} \\ \text{in vivo}}\right)}$, and the pathology view $\mathbf{Y}$.

383

384 For the interpretation of the model, we want to study the importance of individual variables of the

385 observed data to the detected association. For the $k$th factor representing an association between the

386 views, we do this by examining its loadings $\mathbf{W}_{:,k}^{(m)}$ across the $m$ views.

387

388 For biological interpretation, we rank variables of the observed data for each factor-view pair ($k$,$m$).

389 The ranking is done by sorting the loadings $\mathbf{W}_{:,k}^{(m)}$ by their magnitude. For the transcriptomic data

390 views, this procedure leads to a ranked list of transcriptomic microarray probes. The drug-response

391 behavior of the top-ranked probes can be seen as being explained by the factor based on which the

392 ranking was done.

393

394 To detect biological processes, whose changes in the $m$th transcriptomic view are explained by

395 the $k$th factor, we computed the hyper-geometric enrichment test[25] for gene ontology (GO) terms of

396 the transcriptomic probes for the factor-transcriptomic view pair. The $p$-values of the test were

397 controlled for false discovery with the Benjamini-Hochberg correction[33] at the level 0.05.

398 Associations between the enriched pathways and pathological findings were reported in Figure 1

399 based on factor loadings of the pathology view.

## *Retrieval of relevant items*

401 Retrieval means the search of relevant items given a query item. Given the query, the relevance of

402 the items in the database is computed based on a similarity measure, and the items are retrieved in

403 the ranked order of similarity.

404

405 In Case 2, the items are drug-treatments. We retrieved drug-treatments relevant to the query

406 treatment from the database based on their similarity in transcriptomic responses, either using a

407 single-view database $\mathbf{X}^{\binom{\text{human}}{\text{in vitro}}}$, $\mathbf{X}^{\binom{\text{rat}}{\text{in vitro}}}$ or $\mathbf{X}^{\binom{\text{rat}}{\text{in vivo}}}$, or using a multi-view database consisting of

408 all the three transcriptomic views.

409

410 For single-view retrieval, we considered two similarity measures. In the first measure

411 ("correlation"), similarity is defined simply as the correlation between the transcriptomic profiles of

412 the query and the database from the organism in question. As the second measure ("rank-based"),

413 we used a ranked-based approach, also known as connectivity mapping[19]. To compute the similarity

414 of the items, we followed the procedure by Iorio *et al.*[20] In brief, we used a signature of the 250

415 most differentially expressed genes, and computed the average enrichment score similarity between

416 the query signature and the entire ranked list of genes of each of the database items.

## Multi-view database

The simple approach used to compare the query against a single-view database is not directly applicable, when the database and query come from different views or from a different set of views. In either of the cases, we can utilize GFA to detect cross-view associations that then enable the data translation between the query and the database domains and allow us to retrieve relevant items across views.

The database contains data matrices $\mathbf{X}^{(m)} \in \mathbb{R}^{N \times D_m}$ representing views $m = 1, \ldots, M$. In each view, items are organised as rows and variables as columns. Items are co-occurring between the views. The query item $\mathbf{x}^{(query)}$ may be observed in a subset of the database views. In the experiment of this article, the query item is an observation vector from the human *in vitro* transcriptomic view, while the database consist of all the three transcriptomic views.

Since the data domains of the query and the database now are different, similarity search cannot be done in the original data domain as it was done with a single-view database. Latent representation of GFA allows us to carry out the similarity search between items that are observed in different domains. First, we learn a GFA model for the database items. Then, using the learned factors, we learn a latent representation for the query item. Having a latent representation for both the query item and the database items, we can carry out the similarity search in the latent space of the model. Again, we use correlation as a similarity measure, but now in the latent space instead of the original data domain.

## Validation

We validate the retrieval outcome using external information for the items. First, we use the drug-induced liver injury (DILI) label and concern classes[14], which describe the toxic risks of the drugs

observed for the large population of consumers. Second, we use the anatomical therapeutic

chemical (ATC) codes[26] at level 4 to give more detailed information about the drugs' mechanisms

of action.

We measure the retrieval performance in terms of mean average precision at retrieving items with

the same class with the query. We compare the retrieval performance to the performance that

follows the randomization of the class information. For the randomization, we report the mean and

confidence intervals with the width of two standard deviations.

## Acknowledgements

## References

1      Hartwell LH, Szankasi P, Roberts CJ, Murray AW, Friend SH. Integrating genetic

approaches into the discovery of anticancer drugs. Science 1997; 278(5340):1064–1068.

2      Kaletta T, Hengartner MO. Finding function in novel targets: *C. elegans* as a model

organism. Nature Reviews Drug Discovery 2006; 5(5):387–399.

3      Zon LI, Peterson RT. *In vivo* drug discovery in the zebrafish. Nature Reviews Drug

Discovery 2005; 4(1):35–44.

4      Sharpless NE, DePinho RA. Model organisms: The mighty mouse: genetically engineered

mouse models in cancer drug development. Nature Reviews Drug Discovery 2006; 5(9):741–754.

5       Collins FS, Gray GM, Bucher, JR. Toxicology. Transforming environmental health protection. Science 2008; 319(5865): 906–907.

6       Hardy B, Apic G, Carthew P, Clark D, Cook D, Dix I, Escher S, Hastings J, Heard DJ, Jeliazkova N, Judson P, Matis-Mitchell S, Mitic D, Myatt G, Shah I, Spjuth O, Tcheremenskaia O, Toldo L, Watson D, White A, Yang C. Toxicology ontology perspectives. ALTEX 2012; 29(2): 139–156.

7       Willighagen EL, Wehrens R, Buydens LMC. Molecular chemometrics. Crit Rev Anal Chem 2006; 36(3-4): 189–198.

8       Krewski D, Westphal M, Al-Zoughool M, Croteau MC, Andersen ME. New directions in toxicity testing. Annu Rev Public Health 2011; 32: 161–178.

9       Chen M, Zhang M, Borlak J, Tong W. A decade of toxicogenomic research and its contribution to toxicological science. Toxicol Scis 2012; 130(2): 217–228.

10      Zhou T, Chou J, Watkins PB, Kaufmann WK. Toxicogenomics: transcription profiling for toxicology assessment. In: Luch A, editor. Vol. 1, Molecular, Clinical and Environmental Toxicology; Basel (Switzerland): Birkhäuser; 2009. p. 325–366. (Experientia Supplementum; vol. 99).

11      Hartung T, van Vliet E, Jaworska J, Bonilla L, Skinner N, Thomas R. Food for thought ... systems toxicology. ALTEX 2012; 29(2): 119–128.

12      Thomas RS, Black M, Li L, Healy E, Chu T-M, Bao W, Andersen M, Wolfinger R. A comprehensive statistical analysis of predicting *in vivo* hazard using high- throughput *in vitro* screening. Toxicol Sci 2012; 128(2): 398–417.

13      Uehara T, Ono A, Maruyama T, Kato I, Yamada H, Ohno Y, Urushidani T. The Japanese toxicogenomics project: application of toxicogenomics. Molecular Nutrition & Food Research 2010; 54(2): 218–227.

14    Chen M, Vijay V, Shi Q, Liu Z, Fang H, Tong W. FDA-approved drug labeling for the study of drug-induced liver injury. Drug Discov Today 2011; 16(15): 697–703.

15    Le H-S, Bar-Joseph Z. Cross species expression analysis using a Dirichlet process mixture model with latent matchings. In: Lafferty JD, Williams CKI, Shawe-Taylor J, Zemel RS, Culotta A, editors. Advances in Neural Information Processing Systems 23. 24th Annual Conference on Neural Information Processing Systems; 2010 Dec 6–9; Vancouver, BC, Canada. Red Hook, NY: Curran Associates; 2011. p. 1270–1278.

16    Suvitaival T, Huopaniemi I, Oresic M, Kaski S. (2011). Cross-species translation of multi-way biomarkers. Honkela T, Duch W, Girolami M, Kaski S, editors.  Artificial Neural Networks and Machine Learning – ICANN 2011. 21st International Conference on Artificial Neural Networks; 2011 June 14–17; Espoo, Finland. Berlin/Heidelberg (Germany): Springer; 2011. Part I: p. 209–216. (Lecture Notes in Computer Science; vol 6791).

17    The CAMDA Organizing Committee. The CAMDA challenges [Internet]. [cited 2013 Sep 23]. Available from: http://dokuwiki.bioinf.jku.at/doku.php/contest_dataset.

18    Virtanen S, Klami A, Khan SA, Kaski S. Bayesian group factor analysis. In: Lawrence N, Girolami M, editors. JMLR W&CP 22. 15th International Conference on Artificial Intelligence and Statistics; 2012 Apr 21–23; La Palma, Canary Islands. JMLR; 2012. p. 1269–1277.

19    Lamb J, Crawford ED, Peck D, Modell JW, Blat IC, Wrobel MJ, Lerner J, Brunet, J-P, Subramanian A, Ross KN, Reich M, Hieronymus H, Wei G, Armstrong SA, Haggarty SJ, Clemons PA, Wei R, Carr SA, Lander ES, Golub TR. The Connectivity Map: Using gene-expression signatures to connect small molecules, genes, and disease. Science 2006; 313(5795): 1929–1935.

20    Iorio F, Rittman T, Ge H, Menden M, Saez-Rodriguez J. Transcriptional data: a new gateway to drug repositioning? Drug Discov Today 2013; 18(7-8): 350–357.

21    Qu XA, Rajpal DK. Applications of Connectivity Map in drug discovery and development. Drug Discov Today 2012; 17(23-24): 1289–1298.

22      Caiment F, Tsamou M, Jennen D, Kleinjans J. Assessing compound carcinogenicity *in vitro* using connectivity mapping. Carcinogenesis 2014; 35(1):201–207.

23      Smalley JL, Gant TW, Zhang S-DD. Application of connectivity mapping in predictive toxicology based on gene-expression similarity. Toxicology 2010; 268(3): 143–146.

24      Xing L, Wu L, Liu Y, Ai N, Lu X, Fan X. LTMap: a web server for assessing the potential liver toxicity by genome-wide transcriptional expression data. J Appl Toxicol. Forthcoming 2013.

25      Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, Mesirov JP. Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. Proc Natl Acad Sci U S A 2005; 102(43): 15545–15550.

26      ATC classification index with DDDs 2013. Oslo: WHO Collaborating Centre for Drug Statistics Methodology. 2012 [cited 2013 Sep 23]. Available from: http://www.whocc.no/atc_ddd_index/.

27      The CAMDA Organizing Committee. Preprocessed TGP data [Internet]. 2013 [cited 2013 Mar 1]. Available from: http://dokuwiki.bioinf.jku.at/doku.php/tgp_prepro /.

28      Hochreiter S, Clevert D-A, Obermayer K. A new summarization method for Affymetrix probe level data. Bioinformatics 2006; 22(8): 943–949.

29      Hochreiter S, Bodenhofer U, Heusel M, Mayr A, Mitterecker A, Kasim A, Khamiakova T, Van Sanden S, Lin D, Talloen W, Bijnens L, Göhlmann HWH, Shkedy Z, Clevert D-A. FABIA: factor analysis for bicluster acquisition. Bioinformatics 2010; 26 (12): 1520–1527.

30      Caldas J, Kaski S. Bayesian biclustering with the plaid model. In: Proceedings of the IEEE Workshop on Machine Learning for Signal Processing (MLSP); 2008 Oct 16–19; Cancun, Mexico. IEEE; 2008. p. 291-296.

31      Ishwaran H, Rao JS. Spike and slab variable selection: Frequentist and Bayesian strategies. Ann Stat 2005; 33(2): 730–773.

32      Mitchell T, Beauchamp JJ. Bayesian variable selection in linear regression. J Am Stat Assocn 1988; 83(404): 1023–1032.

33      Benjamini Y, Hochberg Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. J Roy Stat Soc B 1995; 57(1):289–300.

456

**Figure 1:** The model detects drug response patterns that generalize across organisms and are associated to organ-level changes driven by toxicity. Also the biological interpretation of the associations represented by a factor generalizes across organisms: changes at the molecular level are interpretable as a biological process. The "eye diagram" shows identified associations between pathological findings (left) and enriched gene ontology (GO) terms (right), represented by factors of the model (middle). Line widths between pathological findings and factors indicate the magnitude of factor loadings learned by the model. Line widths between factors and GO terms indicate the strength of the enrichment. Associations are shown individually for each organism and factor: organisms are indicated as small nodes attached to the nodes of the factors. Factors are named alphabetically from A to H; organisms are human *in vitro* (1), rat *in vitro* (2) and rat *in vivo* (3).

**Figure 2:** All model organisms are informative of the human population-level risk of toxicity. The figure shows how much information the retrieved similar drugs give about the DILI concern, DILI label and ATC level four class, of the query drug. The figure shows the top-10 mean average precision (y-axis) for each organism (x-axis) when used for the retrieval. Retrieval based on differential expression data gives above-random results for each organism using both the correlation and rank-based similarity measure. For the randomized results, shaded areas indicate the 95 % confidence intervals.

**Figure 3:** GFA-based cross-organism approach leads to a higher performance in the retrieval of similar compounds to a human *in vitro* query. The figure shows the top-*k* mean average precision as a function of the number *k* of retrieved highest-ranking samples. GFA utilizes the cross-organism

480  associations learned from the database while the other methods rely on the human *in vitro* data only.

481  For the randomized results, shaded areas indicate the 95 % confidence intervals.
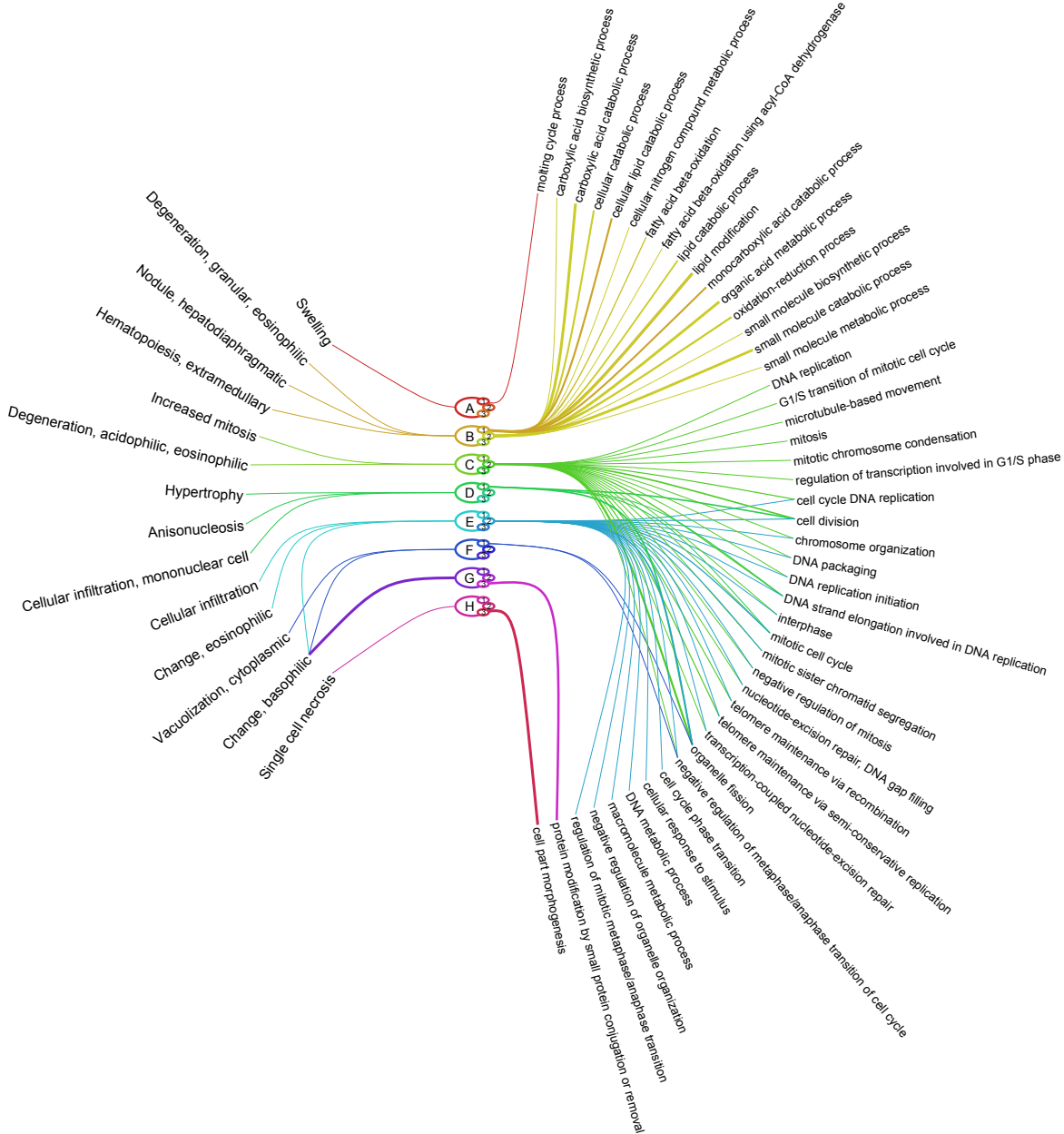
## Tables

482

483  **Table 1:** An example retrieval result shows notable similarity to the query both by toxic and

484  therapeutic properties. Using imipramine as a query, the five most similar compounds are retrieved
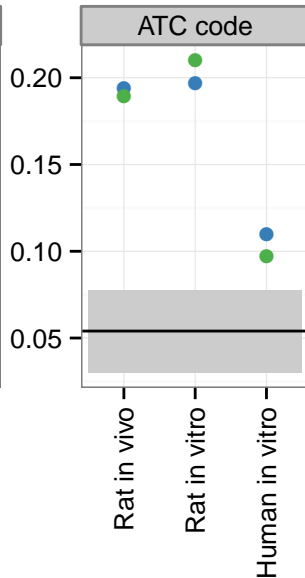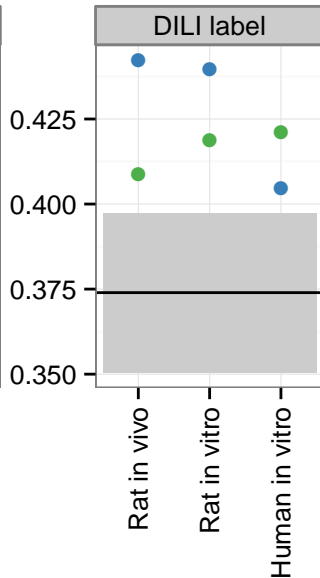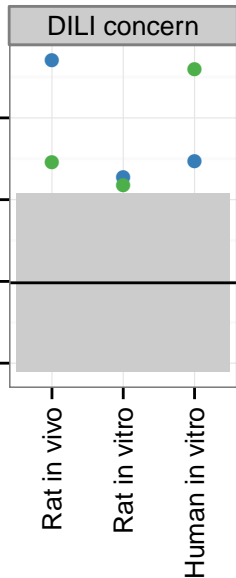
485  based on the GFA model. The table shows the class labels of the retrieved compounds.

| Rank | Compound | DILI concern | DILI label | ATC code |
|---|---|---|---|---|
| Query | Imipramine | Less | Adverse reaction | Non-selective monoamine reuptake inhibitors |
| 1 | Chlorpheniramine | No | No mentioned | |
| 2 | Amitriptyline | Less | Adverse reaction | Non-selective monoamine reuptake inhibitors |
| 3 | Ranitidine | Less | Adverse reaction | H2-receptor antagonists |
| 4 | Hydroxyzine | No | No mentioned | Diphenylmethane derivatives |
| 5 | Tacrine | Most | Warning and precaution | Anticholinesterases |

486

Organ-level
(Pathological findings)

Factors

Molecular level
(GO terms)