



feature

A network-based approach to quantifying the impact of biologically active substances

Julia Hoeng^{1,3}, julia.hoeng@pmintl.com, Renée Deehan^{2,3}, Dexter Pratt², Florian Martin¹, Alain Sewer¹, Ty M. Thomson², David A. Drubin², Christina A. Waters¹, David de Graaf² and Manuel C. Peitsch¹

Society increasingly demands close scrutiny of the potential health risks of long-term exposure to biologically active substances, such as therapeutic drugs or environmental toxins. Such risks are typically assessed *a posteriori* through clinical epidemiology studies. However, disease might take decades to manifest, at a point where changes in therapeutic regime, life style or exposure would not prevent disease onset. Moreover, disease risk as assessed correlatively in epidemiological studies is not intended to elucidate the mechanisms that link perturbations in molecular signaling to disease and, thus, provides fewer options for intervention. Here, we propose that network-based approaches to pharmacology are a valuable way to not only quantify biological network perturbations caused by active substances, but also identify mechanisms and biomarkers modulated in response to exposure and related to disease onset. We also discuss progress towards a generalizable approach for a mechanistic biological impact assessment.

Novel computational methods that derive the quantitative biological impact [defined as a biological impact factor (BIF)] from underlying system-wide data using defined causal biological (i.e. molecular) network models as the substrate for data analysis are currently under

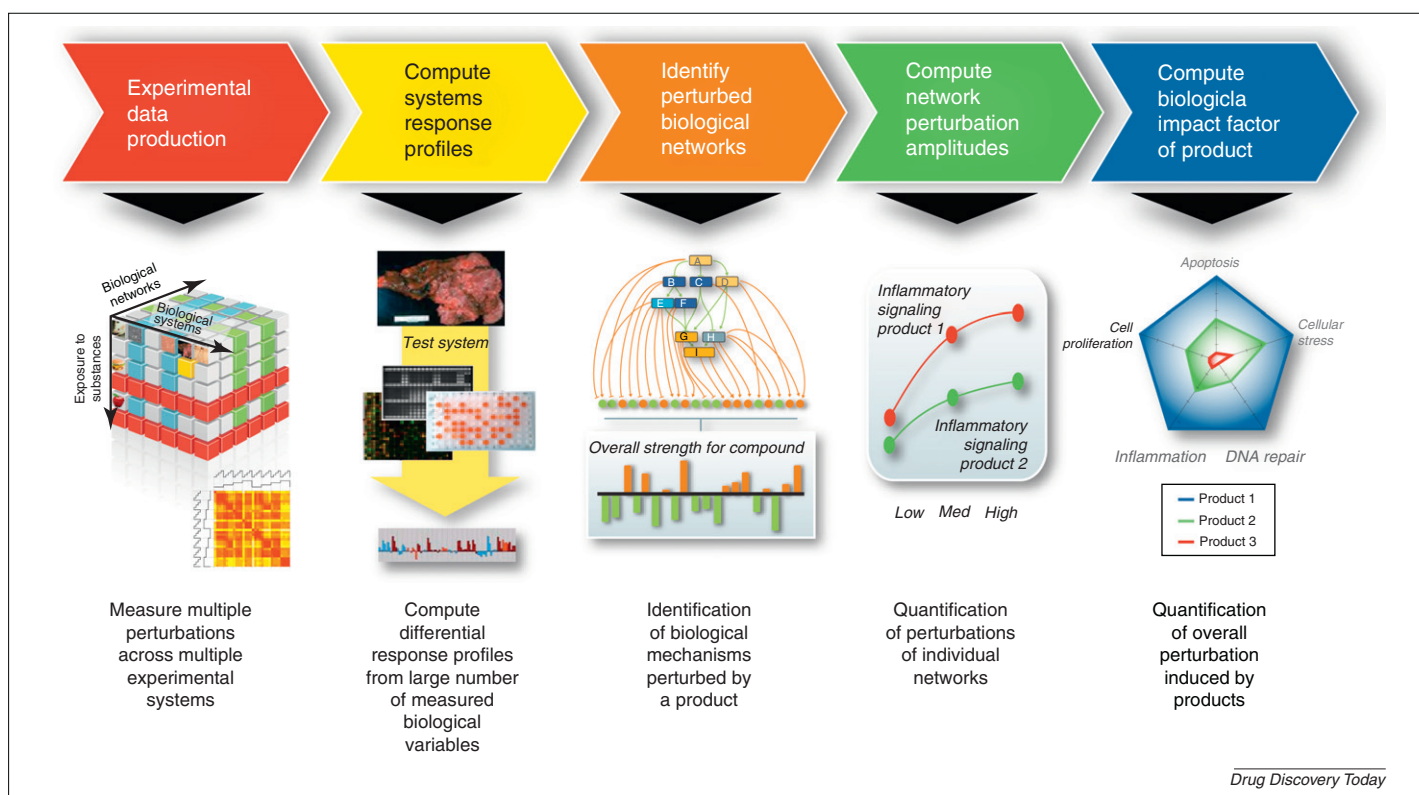
development. The approach described here enables biological impact assessment of active substances at the pharmacological level *a priori*, and can identify mechanisms of action through the application of causal biological network models. The impact of a specific biological network perturbation caused by a single, or a mixture of, biologically active substance(s) is determined for every described molecular entity in the network, thereby identifying causal mechanistic effects induced by each substance or mixture. As the approach is based on system-wide experimental data, this quantitative method takes into account entire biological systems and thereby the many biological networks that can be perturbed by the active substance(s). This enables a quantitative and objective assessment of each molecular entity (or node) in the described biological network(s) to serve, alone or as part of a signature, as a molecular biomarker closely expressing the overall state of perturbation (activation or inhibition compared with control) of every biological network in the system and its correlation with events such as disease onset or progression. Furthermore, this approach facilitates the quantitative comparison of biological impact across individuals and species at the mechanistic level, although gene-level comparisons are

confounded by genomic and/or genetic variations. This capability provides a means to translate between *in vivo* and *in vitro* model system biology and human biology.

This approach provides both potential predictive capabilities and explicit listing of all assumptions through deterministic scoring algorithms. The algorithms, metrics and biological network models presented here will be further developed, published and then made available to the public. We believe that this approach enables application of network pharmacology and systems biology beyond toxicological assessment [1–4] (<http://www.fda.gov/oc/initiatives/criticalpath/whitepaper.html>), and can be applied in areas such as drug development, consumer product testing and environmental impact analysis. Here, we outline the five steps of the strategy (Fig. 1) and the progress made to date.

Step 1: design experiments for data production

For research to translate to human systems, data collected from clinical studies are the most applicable. However, owing to the challenges in obtaining large human data sets, it is necessary to consider non-human models *in vivo* as well as models based on cellular and organotypical (3D)

**FIGURE 1**

The five-step approach to assessing biological impact factors, and the progress made to date.

cultures *in vitro* that represent key aspects of human disease. Data derived from these systems provide at least some insights into the biological network perturbations caused by substances, enabling researchers to identify mechanism-specific biomarkers for use in human studies, and eventually link these mechanisms to the onset of disease for impact assessments.

Although experimental systems *in vitro* and *in vivo* are known to have many shortcomings, we propose that taking a systematic approach to their use will minimize these issues (Fig. 2). Such a systematic approach requires consideration of several constraints:

- **Exposure.** The exposure regimen for a substance or complex stimulus should reflect the range and circumstances of exposure in everyday settings. It is therefore imperative to define a set of standard exposure regimens to be applied systematically to equally well-defined experimental systems. Furthermore, each assay should be designed to collect time- and dose-dependent data to capture both early and late events and to ensure that a representative dose range is covered.
- **Experimental systems.** Experimental systems, if possible, should cover two complementary purposes: (i) animal models that reproduce defined features of the human disease and

are adequate for the exposure and (ii) cellular and organotypical systems should be selected to reflect the cell types and tissues involved in the disease etiology, and priority should be given to primary cells or organ cultures that recapitulate as much as possible the human biology *in vivo*. It is also crucial to match each human culture *in vitro* with the most equivalent culture derived from animal models *in vivo*. This enables creation of a 'translational continuum' from animal model to human biology *in vivo* using the matched systems *in vitro* as 'hubs'.

- **Measurements.** High-throughput system-wide measurements for gene expression, protein expression and post-translational modifications, such as phosphorylation and metabolite profiles, will be generated and correlated with functional outcomes of system exposure. Functional outcome measurements are crucial to the strategy as they serve as anchors for the assessment and represent clear steps in the disease etiology. Although animal models and cellular systems do not always completely translate to human disease, some of the key steps can be reproduced and these represent a major asset in understanding how biological network perturbations can lead to disease.

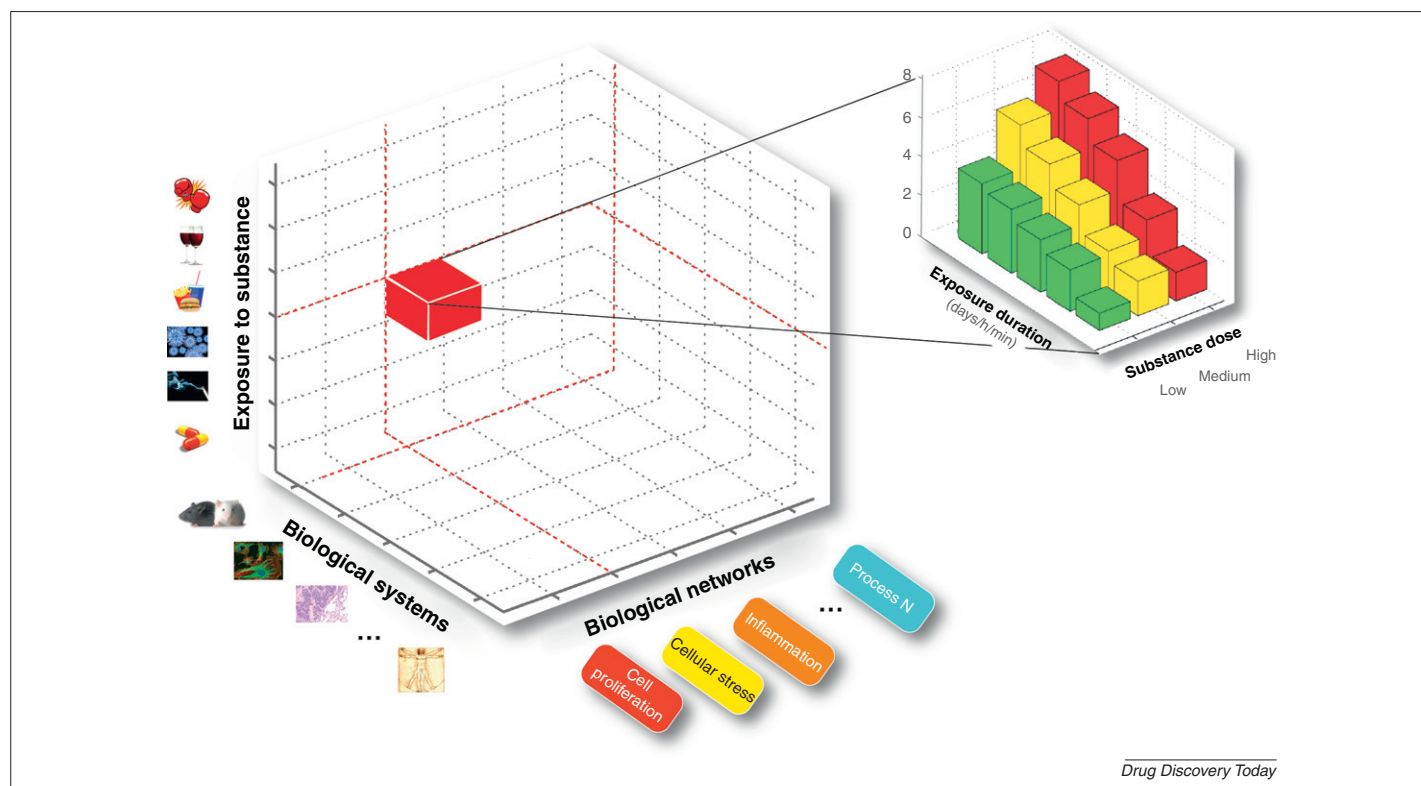
Step 2: compute systems response profiles

The quality-controlled measurements generated in the first step constitute a systems response profile (SRP) for each given exposure in a given experimental system. The SRP therefore expresses the degree to which each individual molecular entity is changed as a consequence of the exposure of the system and is the result of rigorous quality controls and statistical analysis. In this way, different measurements and data types can be integrated and co-analyzed to provide a more accurate quantitative representation of the biology.

Next, measurable elements (e.g. mRNA expression) will be causally integrated into biological network models through the use of prior knowledge. This, coupled with the computational methods in development, enables mechanistic assessment and understanding of biological network perturbations caused by active substances (see below).

Step 3: build biological network models

Whereas the SRPs derived in the previous step represent the experimental data from which biological impact will be determined, it is the causal biological network models that are the substrate for SRP analysis. Application of this strategy requires development of detailed causal

**FIGURE 2**

A systematic approach to experimental design for biological impact factor assessment. Several well-chosen biological systems are exposed to substances in a time- and dose-dependent manner to generate systems-wide data that will be interpreted in the context of each biological network that is relevant to disease onset.

network models of mechanistic biological processes relevant to risk assessment. Such a framework provides a layer of mechanistic understanding beyond examination of gene lists that have been used in more classical toxicogenomics [5]. The strategy to build such models uses biological expression language (BEL), the computable framework for biological network representation produced by Selventa (<http://www.selventa.com>), enabling its application to the evaluation of the biological process of interest based on high-throughput data. The framework will be made available for public use in 2012.

Construction of such a network is an iterative process. Selection of biological boundaries of the network is guided by literature investigation of signaling pathways relevant to the process of interest (e.g. cell proliferation in the lung). Causal relationships describing these pathways are extracted from the Selventa Knowledgebase to nucleate the network with those relationships derived from relevant cell types. The literature-based network can be verified using high-throughput data sets with available phenotypic endpoints.

An example is the microarray analysis of human bronchial epithelial cells perturbed with

an inhibitor of the key cell cycle regulator cyclin-dependent kinase 1 (CDK1) in conjunction with proliferation assays. These data sets are analyzed using reverse causal reasoning (RCR), a method for identifying predictions of the activity states of biological entities (nodes in the network) that are statistically significant and consistent with the measurements taken for a given high-throughput data set [6–8] (Pratt, D.P. *et al.*, unpublished data).

RCR prediction of literature network nodes consistent with the observations of cell proliferation in experiments used to generate the high-throughput data verify whether the network is competent to capture mechanisms regulating the biological process being represented. Additionally, network-relevant nodes predicted by RCR, which were not already represented in the literature network, are integrated. This approach generates a comprehensive biological network with nodes and edges (directional connections between nodes) derived from literature as well as nodes derived from relevant high-throughput data sets.

These networks contain key features that enable process scoring. Topology is maintained and networks of causal relationships (signaling pathways) can be traced from any point in the

network to a measurable entity. Furthermore, the models are dynamic and the assumptions used to build them can be modified or restated and enable adaptability to different tissue contexts and species. This allows for iterative testing and improvement as new knowledge becomes available. To date, we have completed and published two such networks (cell proliferation [9] and cell stress [10]) and are preparing another for publication (cellular stress; Martin, F. *et al.*, unpublished data). The networks are provided in XGMML format and can be viewed using freely available software, such as Cytoscape. The third and fourth networks, describing inflammatory processes and DNA damage, necroptosis, apoptosis, senescence and autophagy, are also in preparation for submission and any additional networks that will be built and the knowledge contained therein will be made freely available to the scientific community.

Step 4: compute NPA scores for biological networks from SRPs

To enable a quantitative comparison of the perturbation of biological networks, we developed a computational approach that translates SRPs into network perturbation amplitudes (NPAs) scores. NPA is an algorithm applied to

experimental data within the context of a causal model of a biological network. Specifically, measurements that are causally mapped as downstream effects of perturbation to individual elements in the model are aggregated via an NPA algorithm into a biological network-specific score. By providing a measure of biological network perturbation, NPA scores allow correlation of molecular events with phenotypes that characterize the network at the cell, tissue or organ level.

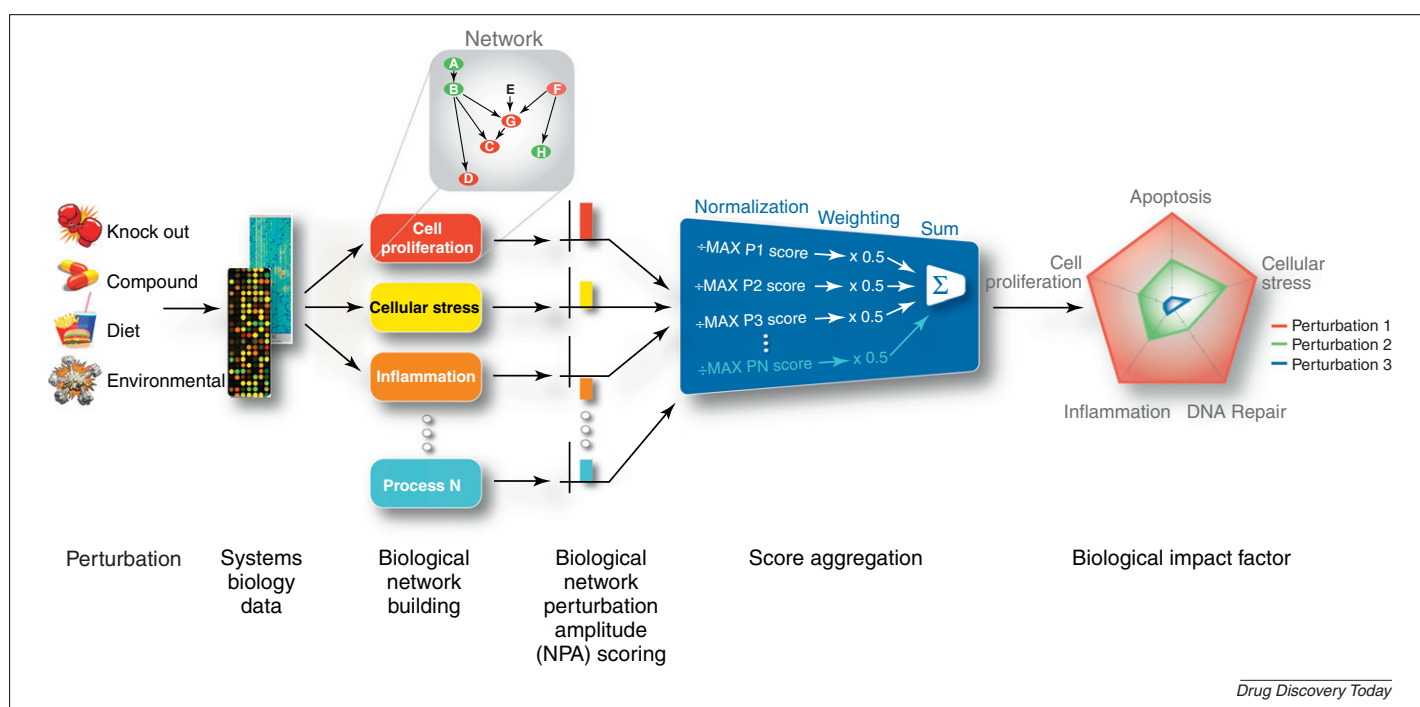
Different methods for integrating measurements that are optimized for different data modalities and to quantify different aspects of system response have been developed, along with statistics to evaluate the quality of an NPA score and to assess the comparability of scores across different experimental systems or treatments (Martin, F et al., unpublished data). These statistics are vital for the interpretation of an NPA score across different exposures and experimental systems. Although an NPA score, with its accompanying statistics, provides a quantifiable measure of the degree of perturbation of a biological network, it does so at a high level of abstraction. For example, an NPA score might reveal that a network is significantly perturbed by a substance, but not indicate exactly where in the network these perturbations take effect. Thus, to complement the network-level NPA methods, several topological NPA methods to

quantify and compare how a perturbation propagates through the network topology were also developed. These topological NPA methods measure the patterns of activation within a network, which further support biological interpretation of the perturbation. Ultimately, topological NPA methods facilitate transparency in the estimate of biological impact, where a perturbation can be traced back from the final estimate of impact, through the NPA scores of contributing biological networks, to the precise pattern of effect within each network.

Step 5: compute biological impact factors for biological systems

Ultimately, a holistic score can be computed that represents the system-wide and pan-mechanistic biological impact of a given substance in a mixture. The final step in estimating the biological impact of a perturbing agent is to aggregate the NPA scores (Fig. 3), which express the impact on each individual biological networks, into one holistic value that expresses the overall impact on the entire biological system. NPA scores for each contributing network are aggregated to produce an estimate of biological impact in a process that requires both normalizing the scores between networks and weighting the contribution of each network. The design of the aggregation algorithm must thus address the issue of defining the relative

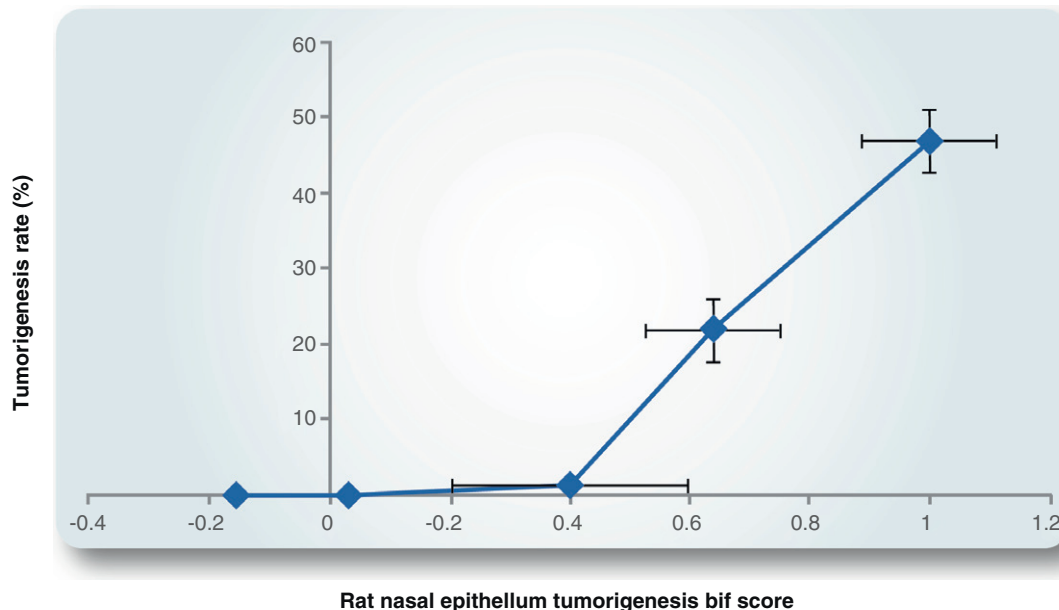
contribution of each biological network to the overall state of the system. Finally, for a BIF to be used as a predictor for medium- and long-term disease outcome, it must be calibrated using a combination of experimental and, if available, epidemiological data. One example of the BIF concept is to estimate the respiratory nasal epithelium tumorigenesis in rats in response to inhaled chemical products. For a simple BIF, the proliferation and inflammation networks were determined to be most relevant for tumorigenesis. Proliferation and inflammation were assumed to contribute equally to tumorigenesis, and thus weighted scores from the two networks equally; however, further investigation into methods for weighting the contributions of different biological networks is required and is currently under development. NPA scores were normalized for each network to the highest score across the different doses. We evaluated nasal epithelium tumorigenesis BIF in rats using transcriptomic data obtained following exposure to multiple doses of formaldehyde for 13 weeks (Fig. 4) [11], and found excellent correlation between the BIF and tumorigenesis rates for rats exposed to the same doses of formaldehyde for two years [12]. In this study, formaldehyde dose was also correlated with tumorigenesis rate; however, the BIF network pharmacology approach elucidates the biological mechanisms that potentially link exposure to disease risk. This



Drug Discovery Today

FIGURE 3

Computational process to derive a biological impact factor for a given biologically active substance using systems-wide experimental data analyzed in the context of biological networks linked to disease onset.



Drug Discovery Today

FIGURE 4

Increases in the tumorigenesis rate in rats after 2 years of formaldehyde exposure correlated with increases in biological impact factor (BIF) score.

shows that even a simple BIF, derived from systems-wide data obtained in short-term experiments, can be a good predictor of long-term disease onset as it represents the impact of an exposure on two disease-linked mechanisms.

The nasal epithelium BIF illustrates these key points. First, for a BIF to be used quantitatively (i.e. to predict the quantitative degree of the disease outcome), it must be calibrated with known disease outcomes. For example, in Fig. 4, there is a threshold effect with tumorigenesis only becoming significant above a BIF of 0.4. Even if a BIF is not calibrated (i.e. long-term disease onset and outcome data are not available), it can be used to rank biological network perturbations based on their expected biological outcomes. Second, component networks can overlap. For example, both the inflammation and proliferation networks contain the transcription factor nuclear factor (NF)- κ B. In general, because component networks contribute to a common outcome (in this case tumorigenesis), we expect such an overlap between component networks. The method used to compute a BIF from the NPA scores should thus take this overlap into account to avoid double-counting particular effects. This might affect how NPA scores are computed for individual networks, meaning that the manner in which the component networks are scored should depend on how they will be used to construct a BIF.

Although the end-point calibrated BIF is presented as a means to evaluate disease onset,

such a BIF could be used to predict effects, including disease progression. In essence, the BIF offers the ability to describe quantitatively the long-term impact of network perturbations, and can be used as a scale for comparison or for threshold establishment based on an associated, BIF-calibrated, outcome. Whereas it is currently necessary to correlate defined exposure modalities (time and dose) of a specified substance (e.g. formaldehyde) or mixture with the rate of disease onset (e.g. tumorigenesis rate) (Fig. 4), such a mechanism-based BIF allows the correlation of biological network perturbations with disease onset as a function of exposure regimen. This would enable the mechanism-based estimation of the risk of long-term disease onset caused by substances for which no long-term epidemiology data are available. Additionally, the process of computing a BIF from system-wide measurements mapped to contributing biological networks enables the simultaneous identification of mechanistic biomarkers that can be used as assessment tools for product testing.

This approach to the prediction of disease onset caused by exposure to biologically active substances is based on the idea that disease risk can be estimated from the quantification of the (long-term or chronic) perturbation of biological networks that contribute to disease onset. Thus, experiments performed over hours, days, or weeks can be used to measure the degree of perturbation of individual networks, and these

can be aggregated into an estimate of risk for disease onset or prognosis for disease progression. Furthermore, time- and exposure-dependent changes of this risk estimate can be readily derived from appropriate experimental data and further provide an indication of risk modification as a function of time and changes in exposure. Applications of this framework include the evaluation of the degree of unwanted biological impact caused by either different manufactured products for safety comparisons or therapeutics (especially those for chronic use) and environmentally active substances to predict safety of long-term exposure and the relationship to adverse effect and onset of disease.

References

- Waters, M.D. and Fostel, J.M. (2004) Toxicogenomics and systems toxicology: aims and prospects. *Nat. Rev. Genet.* 5, 936–948
- Krewski, D. et al. (2011) New directions in toxicity testing. *Annu. Rev. Public Health* 32, 19.1–19.18
- Pleil, J.D. and Sheldon, L.S. (2011) Adapting concepts from systems biology to develop systems exposure event networks for exposure science research. *Biomarkers* 16, 99–105
- Ekins, S. et al. (2005) Techniques: application of systems biology to absorption, distribution, metabolism, and toxicity. *Trends Pharmacol. Sci.* 26, 202–209
- Kyosawa, N. et al. (2010) Practical application of toxicogenomics for profiling toxicant-induced biological perturbations. *Int. J. Mol. Sci.* 11, 3397–3412
- Smith, J.J. et al. (2009) Small molecule activators of SIRT1 replicate signaling pathways triggered by calorie restriction *in vivo*. *BMC Syst. Biol.* 3, 31

- 7 Laifenfeld, D. *et al.* (2010) The role of hypoxia in 2-butoxyethanol-induced hemangiosarcoma. *Toxicol. Sci.* 113, 254–266
- 8 Kumar, R. *et al.* (2010) Reverse causal reasoning: applying qualitative causal knowledge to the interpretation of high-throughput data. *BMC Genomics* 6, 419
- 9 Westra, J.W. *et al.* (2011) Construction of a computable cell proliferation network focused on non-diseased lung cells. *BMC Syst. Biol.* 2, 105
- 10 Schlage, W.K. *et al.* (2011) A computable cellular stress network model for non-diseased pulmonary and cardiovascular tissue. *BMC Syst Biol* 5, 168
- 11 Andersen, M.E. *et al.* (2010) Formaldehyde: integrating dosimetry, cytotoxicity, and genomics to understand dose-dependent transitions for an endogenous compound. *Toxicol. Sci.* 118, 716–731
- 12 Monticello, T.M. *et al.* (1996) Correlation of regional and nonlinear formaldehyde-induced nasal cancer with proliferating populations of cells. *Cancer Res.* 56, 1012–1022
- Julia Hoeng^{1,3}**
Renée Deehan^{2,3}
Dexter Pratt²
Florian Martin¹

Alain Sewer¹
Ty M. Thomson²
David A. Drubin²
Christina A. Waters¹
David de Graaf²
Manuel C. Peitsch¹

¹Philip Morris International R&D,
Philip Morris Products S.A., Neuchâtel, Switzerland

²Selventa, Cambridge, MA, USA

³These authors contributed equally to this work.