Cite This: J. Chem. Inf. Model. XXXX, XXX, XXX-XXX

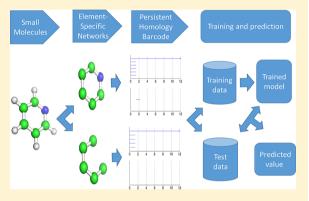
# Quantitative Toxicity Prediction Using Topology Based Multitask **Deep Neural Networks**

Kedi Wu<sup>†</sup> and Guo-Wei Wei\*,†,‡,¶

<sup>†</sup>Department of Mathematics, <sup>‡</sup>Department of Electrical and Computer Engineering, and <sup>¶</sup>Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, Michigan 48824, United States

Supporting Information

**ABSTRACT:** The understanding of toxicity is of paramount importance to human health and environmental protection. Quantitative toxicity analysis has become a new standard in the field. This work introduces element specific persistent homology (ESPH), an algebraic topology approach, for quantitative toxicity prediction. ESPH retains crucial chemical information during the topological abstraction of geometric complexity and provides a representation of small molecules that cannot be obtained by any other method. To investigate the representability and predictive power of ESPH for small molecules, ancillary descriptors have also been developed based on physical models. Topological and physical descriptors are paired with advanced machine learning algorithms, such as the deep neural network (DNN), random forest (RF), and gradient boosting decision



tree (GBDT), to facilitate their applications to quantitative toxicity predictions. A topology based multitask strategy is proposed to take the advantage of the availability of large data sets while dealing with small data sets. Four benchmark toxicity data sets that involve quantitative measurements are used to validate the proposed approaches. Extensive numerical studies indicate that the proposed topological learning methods are able to outperform the state-of-the-art methods in the literature for quantitative toxicity analysis. Our online server for computing element-specific topological descriptors (ESTDs) is available at http://weilab. math.msu.edu/TopTox/.

## 1. INTRODUCTION

Toxicity is a measure of the degree to which a chemical can adversely affect an organism. These adverse effects, which are called toxicity end points, can be either quantitatively or qualitatively measured by their effects on given targets. Qualitative toxicity classifies chemicals into toxic and nontoxic categories, while quantitative toxicity data set records the minimal amount of chemicals that can reach certain lethal effects. Most toxicity tests aim to protect human from harmful effects caused by chemical substances and are traditionally conducted in in vivo or in vitro manner. Nevertheless, such experiments are usually very timeconsuming and cost intensive, and even give rise to ethical concerns when it comes to animal tests. Therefore, computer-aided methods, or in silico methods, have been developed to improve prediction efficiency without sacrificing too much of accuracy. The quantitative structure—activity relationship (QSAR) approach is one of the most popular and commonly used approaches. The basic QASR assumption is that similar molecules have similar activities. Therefore, by studying the relationship between chemical structures and biological activities, it is possible to predict the activities of new molecules without actually conducting lab experiments.

There are several types of algorithms to generate QSAR models: linear models based on linear regression and linear discriminant analysis; nonlinear models including nearest neighbor, 2,3

support vector machine, 1,4,5 and random forest. These methods have advantages and disadvantages due to their statistical natures. For instance, linear models overlook the relatedness between different features, while the nearest neighbor method largely depends on the choice of descriptors. To overcome these difficulties, more refined and advanced machine learning methods have been introduced. Multitask (MT) learning<sup>8</sup> was proposed partially to deal with data sparsity problems, which are commonly encountered in QSAR applications. The idea of MT learning is to learn the so-called "inductive bias" from related tasks to improve accuracy using the same representation. In other words, MT learning aims at learning a shared and generalized feature representation from multiple tasks. Indeed, MT learning strategies have brought new insights to bioinformatics since compounds from related assays may share features at various feature levels, which is extremely helpful if the data set is small. Successful applications include splice-site and MHC-I binding prediction<sup>9</sup> in sequence biology, gene expression analysis, and system biology. 10

Recently, deep learning (DL), 11,12 particularly convolutional neural network (CNN), has emerged as a powerful paradigm to render a wide range of the-state-of-the-art results in signal and

Received: September 15, 2017 Published: January 9, 2018

information processing fields, such as speech recognition 13,14 and natural language processing. 15,16 Deep learning architecture is essentially based on artificial neural networks. The major difference between deep neural network (DNN) models and non-DNN models is that DNN models consist of a large number of layers and neurons, making it possible to construct abstract features.

Geometric representation of molecules often contains too much structural detail and thus is prohibitively expensive for most realistic large molecular systems. However, traditional topological methods often reduce too much of the original geometric information. Persistent homology, a relatively new branch of algebraic topology, offers an interplay between geometry and topology. 17,18 It creates a variety of topologies of a given object by varying a filtration parameter. As a result, persistent homology can capture topological structures continuously over a range of spatial scales. Unlike commonly used computational homology which results in truly metric free representations, persistent homology embeds geometric information in topological invariants, e.g., Betti numbers, so that "birth" and "death" of isolated components, rings, and cavities can be monitored at all geometric scales by topological measurements.

Recently, we have introduced persistent homology for the modeling and characterization of nanoparticles, proteins, and other biomolecules. 19-24 We proposed molecular topological fingerprint (TF) to reveal topology-function relationships in protein folding and protein flexibility. 19 This approach was integrated machine-learning algorithms for protein classification.<sup>22</sup> However, it was found that primitive persistent homology has a limited power in protein classification due to its oversimplification of biological information.<sup>25</sup> Most recently, element specific persistent homology (ESPH) has been introduced to retain crucial biological information during the topological simplifica-tion of geometric complexity. <sup>26–28</sup> The integration of ESPH and machine learning gives rise to some of the most accurate predictions of protein—ligand binding affinities<sup>27,28</sup> and mutation induced protein stability changes.<sup>26,28</sup> However, ESPH has not been validated for its potential utility in small molecular characterization, analysis, and modeling. In fact, unlike proteins, small molecules involve a large number of element types and are more diversified in their chemical compositions. They are also rich in structural variability in structures, including cis-trans distinctions and chiral and achiral stereoisomers. Small molecular properties are very sensitive to their structural and compositional differences. Therefore, it is important to understand the representability and predictive power of ESPH in dealing with small molecular diversity, variability, and sensitivity.

The objective for this work is to introduce element specific topological descriptors (ESTDs) constructed via ESPH for quantitative toxicity analysis and prediction of small molecules. We explore the representational and predictive powers of ESTDs for small molecules. Physical descriptors constructed from microscopic models are also developed both as ancillary descriptors and as competitive descriptors to further investigate the proposed topological methods. These new descriptors are paired with advanced machine learning algorithms, including MT-DNN, single-task DNN (ST-DNN), random forest (RF), and gradient boosting decision tree (GBDT), to construct topological learning strategies for illustrating their predictive power in quantitative toxicity analysis. We demonstrate that the proposed topological learning provides a very competitive description of relatively small drug-like molecules. Additionally, the inherent correlation among different quantitative toxicity end

points makes our topology based multitask strategy a viable approach to quantitative toxicity predictions.

## 2. METHODS AND ALGORITHMS

In this section, we provide a detail discussion about molecular descriptors used in this study, including element-specific topological descriptors and auxiliary descriptors calculated from physical models. Moreover, an overview of machine learning algorithms, including ensemble methods (random forest and gradient boosting decision tree), deep neural networks, singletask learning and multitask learning, is provided. Emphasis is given to advantages of multitask deep convolutional neural network for quantitative toxicity end point predictions and how to select appropriate parameters for network architectures. Finally, we provide a detailed description of our learning architecture, training procedure, and evaluation criteria.

2.1. Element Specific Topological Descriptor (ESTD). In this subsection, we give a brief introduction to persistent homology and ESTD construction. An example is also given to illustrate the construction.

2.1.1. Persistent Homology. For atomic coordinates in a molecule, algebraic groups can be defined via simplicial complexes, which are constructed from simplices, i.e., generalizations of the geometric notion of nodes, edges, triangles, tetrahedrons, etc. Homology associates a sequence of algebraic objects, such as abelian groups, to topological spaces and characterizes the topological connectivity of geometric objects in terms of topological invariants, i.e., Betti numbers, which are used to distinguish topological spaces. Betti-0, Betti-1, and Betti-2, respectively, represent independent components, rings, and cavities in a physical sense. A filtration parameter, such as the radius of a ball, is used to continuously vary over an interval so as to generate a family of structures. Loosely speaking, the corresponding family of homology groups induced by the filtration is a persistent homology. The variation of the topological invariants, i.e., Betti numbers, over the filtration gives rise to a unique characterization of physical objects, such as molecules.

Simplex. Let  $u_0$ ,  $u_1$ , ...,  $u_k$  be a set of points in  $\mathbb{R}^d$ . A point  $x = \sum_{i=0}^k \lambda_i u_i$  is called an *affine combination* of the  $u_i$  if  $\sum_{i=0}^k \lambda_i = 1$ . The k + 1 points are said to be *affinely independent*, if and only if  $u_i - u_{0i}$ ,  $1 \le i \le k$  are linearly independent. We can find at most d linearly independent vectors and at most d + 1 affinely independent points in  $\mathbb{R}^d$ .

An affine combination,  $x = \sum_{i=0}^{k} \lambda_i u_i$  is a *convex combination* if  $\lambda_i$  are nonnegative. A k-simplex, which is defined to be the convex hull (the set of convex combinations) of k + 1 affinely independent points, can be formally represented as

$$\sigma = \left\{ \sum_{i=0}^{k} \lambda_{i} u_{i} \middle| \sum_{i=0}^{k}$$

where  $\{u_0, u_1, ..., u_k\} \subset \mathbb{R}^d$  is a set of affinely independent points. Examples of k-simplex for the first few dimensions are shown in Figure 1. Essentially, a 0-simplex is a vertex, a 1-simplex is an edge, a 2-simplex is a triangle, and a 3-simplex is a tetrahedron. A *face*  $\tau$  of  $\sigma$  is the convex hull of a nonempty subset of  $u_i$  and is proper if the subset does not contain all k + 1 points. Equivalently, we can write as  $\tau \leq \sigma$  if  $\tau$  is a face or  $\sigma$ , or  $\tau < \sigma$  if  $\tau$  is proper. The *boundary* of  $\sigma$  is defined to be the union of all

Simplicial Complex. A simplicial complex is a finite collection of simplices K such that  $\sigma \in K$  and  $\tau \leq \sigma$  implies  $\tau \in K$ , and  $\sigma$ ,

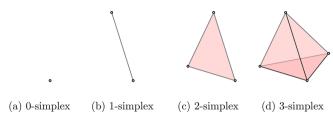


Figure 1. Examples of simplex of different dimensions.

 $\sigma_0 \in K$  implies  $\sigma \cap \sigma_0$  is either empty or a face of both. The *dimension* of K is defined to be the maximum dimension of its simplices.

Chain Complex. Given a simplicial complex K and a constant p as dimension, a p-chain is a formal sum of p-simplices in K, denoted as  $c = a_i \sigma_i$ . Here  $\sigma_i$  are the p-simplices and the  $a_i$  are the coefficients, mostly defined as 0 or 1 (module 2 coefficients) for computational considerations. Specifically, p-chains can be added as polynomials. If  $c_0 = \sum a_i \sigma_i$  and  $c_1 = \sum b_i \sigma_i$ , then  $c_0 + c_1 = \sum (a_i + b_i)\sigma_i$ , where the coefficients follow  $\mathbb{Z}_2$  addition rules. The p-chains with the previously defined addition form an abelian group and can be written as  $(C_p, +)$ . A boundary operator of a p-simplex  $\sigma$  is defined as

$$\partial_p \sigma = \sum_{j=0}^p (-1)^j [u_0, u_1, ..., \widehat{u}_j, ..., u_p]$$
(2)

where  $[u_0, u_1, ..., \widehat{u}_j, ..., u_p]$  means that vertex  $u_j$  is excluded in computation. Given a p-chain  $c = a_i \sigma_i$ , we have  $\partial_p c = \sum a_i \partial_p \sigma_i$ . Notice that  $\partial_p$  maps p-chain to the  $\{p-1\}$ -chain and that boundary operation commutes with addition, a boundary homomorphism  $\partial_p : \sigma_p \to \sigma_{p-1}$  can be defined. The chain complex can be further defined using such boundary homomorphism as follows:

$$\cdots \to C_{p+1} \xrightarrow{\partial_{p+1}} C_p \xrightarrow{\partial_p} C_{p-1} \xrightarrow{\partial_{p-1}} \cdots \xrightarrow{\partial_1} C_0 \xrightarrow{\partial_0} 0 \tag{3}$$

Cycles and Boundaries. A p-cycle is defined to be a p-chain c with empty boundary ( $\partial_p c = 0$ ), and the group of p-cycles of K is denoted as  $Z_p = Z_p(K)$ . In other words,  $Z_p$  in the kernel of the p-th boundary homomorphism,  $Z_p = \ker \partial_p$ . A p-boundary is a p-chain, say c, such that there exists  $d \in C_{p+1}$  and  $\partial_p d = c$ , and the group of p-boundaries is written as  $B_p = B_p(K)$ . Similarly, we can rewrite  $B_p$  as  $B_p = \operatorname{im} \partial_{p+1}$  since the group of p-boundaries is the image of the (p+1)th boundary homomorphism.

Homology Groups. The fundamental lemma of homology says that the composition operator  $\partial_p \bigcirc \partial_{p+1}$  is a zero map. With this lemma, we conclude that im  $\partial p + 1$  is a subgroup of ker  $\partial_p$ . Then the *pth homology group* of simplicial complex is defined as the *pth* cycle group modulo the *pth* boundary group,

$$H_p = Z_p/B_p \tag{4}$$

and the *pth Betti number* is the rank of this group,  $\beta_p = \text{rank } H_p$ . Geometrically, Betti numbers can be used to describe the connectivity of given simplicial complexes. Intuitively,  $\beta_0$ ,  $\beta_1$ , and  $\beta_2$  are numbers of connected components, tunnels, and cavities, respectively, for the first few Betti numbers.

Filtration and Persistence. A filtration of a simplicial complex *K* is a nested sequence of subcomplexes of *K*.

$$\emptyset = K_0 \subseteq K_1 \subseteq \dots \subseteq K_n = K \tag{5}$$

For each  $i \leq j$ , there exists an inclusion map from  $K_i$  to  $K_j$  and therefore an induced homomorphism  $f_p^{i,j} \colon H_p(K_i) \to H_p(K_j)$  for each dimension p. The filtration defined in eq 5 thus

corresponds to a sequence of homology groups connected by homomorphisms.

$$0 = H_n(K_0) \to H_n(K_1) \to \dots \to H_n(K_n) = H_n(K)$$
(6)

for each dimension *p*. The *p*th persistent homology groups are defined as the images of the homomorphisms induced by inclusion,

$$H_p^{i,j} = \operatorname{im} f_p^{i,j} \tag{7}$$

where  $0 \le i \le j \le n$ . In other words,  $H_p^{i,j}$  contains the homology classes of  $K_i$  that are still alive at  $K_j$  for given dimension p and each pair i, j. We can reformulate the pth persistent homology group as

$$H_v^{i,j} = Z_v(K_i) / (B_v(K_i) \cap Z_v(K_i))$$
(8)

The corresponding pth persistent Betti numbers are the ranks of these groups,  $\beta_p^{i,j} = \text{rank } H_p^{i,j}$ . The birth, death, and persistence of a Betti number carry important chemical and/or biological information, which is the basis of the present method.

2.1.2. Persistent Homology for Characterizing Molecules. As introduced before, persistent homology indeed reveals long lasting properties of a given object and offers a practical method for computing topological features of a space. In the context of toxicity prediction, persistent homology captures the underlying invariants and features of small molecules directly from discrete point cloud data. A intuitive way to construct simplicial complex from point cloud data is to utilize Euclidean distance, or essentially to use a so-called "Vietoris—Rips complex". A Vietoris—Rips complex is defined to be a simplicial complex whose k-simplices correspond to unordered (k+1)tuples of points which are pairwise within radius  $\epsilon$ .

However, a particular radius  $\epsilon$  is not sufficient since it is difficult to see if a hole is essential. Therefore it is necessary to increase radius  $\epsilon$  systematically and see how the homology groups and Betti numbers evolve. The persistence 18,29 of each Betti number over the filtration can be recorded in barcodes. 30,31 The persistence of topological invariants observed from barcodes offers an important characterization of molecular structures. For instance, given the 3D coordinates of a small molecule, a short-lived Betti-0 bar may be the consequence of a strong covalent bond while a long-lived Betti-0 bar can indicate a weak covalent bond. Similarly, a long-lived Betti-1 bar may represent a chemical ring. Such observations motivate us to design persistent homology based topological descriptors. However, it is important to note that the filtration radius is not a chemical bond and topological connectivity is not a physical relationship. In other words, persistent homology offers a representation of molecules that is entirely different from classical theories of chemical and/or physical bonds. Nevertheless, such a representation is systematical and comprehensive and thus is able to unveil structure-function relationships when it is coupled with advanced machine learning algorithms.

Example. Figure 2 is an detailed example of how our ESTDs are calculated and how they can reveal the structural information on pyridine. An all-elements representation of pyridine is given in Figure 2a, where carbon atoms are in green, nitrogen atoms are in blue, and hydrogen atoms are in white. Without considering covalent bonds, there are 11 isolated vertices (atoms) in Figure 2a. Keep in mind that if the distance between two vertices is less than the filtration value then these two vertices do not connect. Thus, at filtration value 0, we should have 11 independent components and no loops, which are respectively

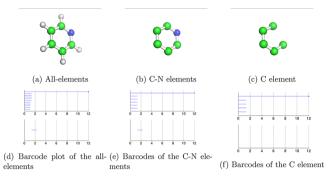


Figure 2. Illustration of pyridine and its persistent homology barcode plots. (a-c) Only carbon atoms are used for persistent homology computation, respectively. (d-f) From top to bottom, the results are computed for 0-dimension (Betti-0) and 1-dimension (Betti-1), respectively.

reflected by the 11 Betti-0 bar and 0 Betti-1 bar in Figure 2d. As the filtration value increases to 1.08 Å, every carbon atom starts to connect with its nearest hydrogen atoms, and consequently, the number of independent components (also the number of Betti-0 bar) reduces to 6. When filtration value reaches 1.32 Å, we are left with 1 Betti-0 bar and 1 Betti-1 bar. This indicates that there is only one independent component and the hexagonal carbon-nitrogen ring appears since the filtration value has exceeded the length of both the carbon-carbon and carbonnitrogen bonds. As the filtration value becomes sufficiently large, the hexagonal ring is eventually filled and there is only one totally connected component left.

It is worth to mentioning that Figure 2d does not inform the existence of the nitrogen atom in this molecule. Much chemical information is missing during the topological simplification. This problem becomes more serious as the molecular becomes larger and its composition becomes more complex. A solution to this problem is the element specific persistent homology or multicomponent persistent homology.<sup>26</sup> In this approach, a molecule is decomposed into multiple components according to the selections of element types and persistent homology analysis is carried out on each component. The all-atom persistent homology shown in Figure 2d is a special case in the multicomponent persistent homology. Additionally, barcodes in Figure 2e are for all carbon and nitrogen elements, while barcodes in Figure 2f are for carbon only. By a comparison of these two barcodes, one can conclude that there is a nitrogen atom in the molecule and it must be on the ring. In this study, all persistent homology computations are carried out by Dionysus (http:// mrzv.org/software/dionysus/) with Python bindings.

Element Specific Networks. The key to accurate prediction is to engineer ESTDs from corresponding element specific networks (ESNs) on which persistent homology is computed. As the example above shows, it is necessary to choose different element combinations in order to capture the properties of a given molecule. Carbon (C), nitrogen (N), and oxygen (O) are commonly occurring elements in small molecules. Unlike proteins where hydrogen atoms are usually excluded due to their absence in the database, for small molecules it is beneficial to include hydrogen atoms in our ESTD calculations. Therefore, ESNs of single-element types include four type elements  $\mathcal{A} = \{H, C, N, O\}$ . Additionally, we also consider element combinations that involve two or more element types in an element specific network. In particular, the barcode of the network consisting of N and O elements in molecule might reveal hydrogen bond interaction strength.

Networks with a wide variety of element combinations were tested and a good selection of such combinations is shown in Table 1. Specifically, two types of networks are used in the present work, namely, single- and two-element networks.

Table 1. Element Specific Networks Used to Characterize Molecules

network type	element specific networks
single-element	$\{a_i\}$ , where $a_i \in \mathcal{A}$ , $\mathcal{A} = \{H, C, N, O\}$
two-element	$\{b_i, c_j\}$ , where $b_i \in \mathcal{B}, c_j \in \mathcal{C}, i \in \{1,, 3\}, j \in \{1,, 9\}$ ,
	and $i < j$
	here $\mathcal{B} = \{C, N, O\}$ and $C = \{C, N, O, F, P, S, Cl, Br, I\}$

Denote  $a_i$  as the *i*th atom of element type a and  $\{a_i\}$  the set of all atoms of element type a in a molecule. Then  $\{a_i\}$  with  $a \in \mathcal{A}$  includes four different single-element type networks. Similarly, Table 1 lists 21 different two-element networks. Therefore, a total of 25 element-specific networks is used in the present work.

Filtration Matrix. Another importance aspect is the filtration matrix that defines the distance in persistent homology analysis.  $^{19,32}$  We denote the Euclidean distance between atom i at  $(x_i, y_i, z_i)$  and atom j at  $(x_j, y_j, z_j)$  to be

$$d_{i,j} = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2}$$
(9)

By a direct filtration based on the Euclidean distance, one can capture the information on covalent bonds easily as shown in Figure 2d. However, intramolecular interactions such as hydrogen bonds and van der Waals interactions cannot be revealed. In other words, the Betti-0 bar of two atoms with certain hydrogen bonding effect cannot be captured since there already exist shorter Betti-0 bar (covalent bonds). To circumvent such deficiencies we use filtration matrix to redefine the

$$M_{i,j} = \begin{cases} d_{i,j}, & \text{if } d_{i,j} \ge r_i + r_j + |\Delta d| \\ d_{\infty}, & \text{otherwise} \end{cases}$$
 (10)

where  $r_i$  and  $r_i$  are the atomic radius of atoms i and j, respectively. Here  $\Delta d$  is the bond length deviation in the data set and  $d_{\infty}$  is a large number which is set to be greater than the maximal filtration value. By setting the distance between two atoms that have a covalent bond to a sufficiently large number, we are able to use topology to capture important intramolecular interactions, such as hydrogen bonds, electrostatic interactions and van der Waals interactions.

Topological Dimension. Finally we need to consider the dimensions of topological invariants. For large molecules such as proteins, it is important to compute the persistent homology of first three dimensions, which will result in Betti-0, Betti-1, and Betti-2 bar. The underlying reason is that proteins generally consists of thousands of atoms, and Betti-1 and Betti-2 bar usually contain very rich geometric information such as internal loops and cavities. However, small molecules are geometrically relatively simple and their barcodes of high dimensions are usually very sparse. Additionally, small molecules are chemically complex due to their involvement of many element types and oxidation states. As such, high dimensional barcodes of element specific networks carry little information. Therefore, we only consider Betti-0 bar for small molecule modeling.

2.1.3. ESTDs for Small Molecules. A general process for our ESTD calculation can be summarized as follows.

1. 3D coordinates of atoms of selected atom types are selected, and their Vietoris-Rips complexes are constructed. Note that distance defined in eq 10 is used for persistent homology barcodes generation.

- 2. The maximum filtration size is set to 10 Å considering the size of small molecules. After barcodes are obtained, the first 10 small intervals of length 0.5 Å are considered. In other words, ESTDs will be calculated based on the barcodes of each subinterval  $Int_i = [0.5i, 0.5(i + 1)], i = 0, ..., 9$ .
  - Within each Int<sub>i</sub>, search Betti-0 bar whose birth time falls within this interval and Betti-0 bar that dies within Int<sub>i</sub>, respectively, and denote these two sets of Betti-0 bar as S<sub>birthi</sub> and S<sub>deathi</sub>.
  - Count the number of Betti-0 bars within S<sub>birth,</sub> and S<sub>death,</sub> and these two counts yield two ESTDs for the interval Int<sub>i</sub>.
- 3. In addition to intervalwise descriptors, we also consider global ESTDs for the entire barcodes. All Betti-0 bar birth and death times are collected and added into  $S_{\text{birth}}$  and  $S_{\text{death}}$ , respectively. The maximum, minimum, mean, and sum of each set of values are then computed as ESTDs. This step gives eight more ESTDs.

Therefore, for each element specific network, we have a total of  $28~(2\times10~\text{intervals}+8)$  ESTDs. Since we consider a total 25 single- and two-element networks, we have a total 700  $(25\times28)$  ESTDs.

Finally, we would like to emphasize the essential ideas of our choice of ESTDs. In step 2 of the ESTD generation process, we collect all birth and death time of Betti-0 bar in order to capture the hydrogen bonding and van der Waals interactions. These intramolecular interactions are captured by eliminating the topological connectivity of covalent bonds. The birth position can signal the formation of hydrogen bonding, and the death position represents the disappearance of such effects, which in turn reflects the strength of these effects. In step 3 of the above process, we consider all potential element-specific intramolecular effects together and use statistics of these effects as global descriptors for a given molecule. This would help us to better characterize small molecules.

The topological feature vector that consists of ESTDs for the ith molecule in the tth prediction task (one task for each toxicity prediction), denoted as  $\mathbf{x}_i^t$ , can be used to approximate of the topological functional  $f^t$  of MT-DNN. This optimization process will be carefully discussed in section 2.4.3.

- 2.2. Auxiliary Molecular Descriptors. In addition to ESTDs, we are also interested in constructing a set of microscopic features based on physical models to describe molecular toxicity. This set of features should be convenient for being used in different machine learning approaches, including deep learning and non deep learning, and single-task and multitask ones. To make our feature generation feasible and robust to all compounds, we consider three types of basic physical information, i.e., atomic charges computed from quantum mechanics or molecular force fields, atomic surface areas calculated for solvent excluded surface definition, and atomic electrostatic solvation free energies estimated from the Poisson model. To obtain this information, we first construct optimized 3D structure of for each molecule. Then the aforementioned atomic properties are computed. Our feature generation process can be divided into several steps:
  - 1. **Structure**. Optimized 3D structures were prepared by LigPrep in Schrödinger suites (2014-2) from the original 2D structures, using options {-i 0 -nt -s 10 -bff 10}.

- Charge. Optimized 3D structures were then fed in antechamber,<sup>33</sup> using parametrization: AM1-BCC charge, Amber mbondi2 radii, and general Amber force field (GAFF).<sup>34</sup> This step leads to pqr files with corresponding charge assignments.
- 3. Surface. ESES online server<sup>35</sup> was used to compute atomic surface area of each molecule, using pqr files from the previous step. This step also results in molecular solvent excluded surface information.
- 4. Energy. MIBPB online server<sup>36</sup> was used to calculate the atomic electrostatic solvation free energy of each molecule, using surface and pqr files from previous steps.

Auxiliary molecular descriptors were obtained according to the above procedure. Specifically, these molecular descriptors come from steps 2, 3, and 4. To make our method scalable and applicable to all kinds of molecules, we manually construct element-specific molecular descriptors so that it does not depend on atomic positions or the number of atoms. The essential idea of such construction is to derive atomic properties of the each element type, which is very similar to the idea of ESPH.

We consider 10 different commonly occurring element types, i.e., H, C, N, O, F, P, S, Cl, Br, and I and three different types of descriptors—charge, surface area, and electrostatic solvation free energies. Given an element type and a descriptor type, we compute the statistics of the quantities obtained from the aforementioned physical model calculation, i.e., summation, maximum, minimum, mean, and variance, giving rise to five physical descriptors. To capture absolute strengths of each element descriptor, we further generate five more physical descriptors after taking absolute values of the same quantities. Consequently, we have a total of 10 physical descriptors for each given element type and descriptor type. Thus, 300 (10 descriptor ×10 element types × 3 descriptor types) molecular descriptors can be generated at element type level.

Additionally when all atoms are included for computation, 10 more physical descriptors can be constructed in a similar way (5 statistical quantities of original values and another 5 for absolute values) for each element descriptor type (charge, surface area, and electrostatic solvation free energies). This step yields another 30 molecular descriptors. As a result, we organize all of the above information into a 1D feature vector with 330 components, which is readily suitable for ensemble methods and DNN.

These auxiliary molecular descriptors result in an independent descriptor set. When adding these molecular descriptors to the previously mentioned ESTDs, we have a full descriptor set.

- **2.3. Descriptor Selection.** The aforementioned descriptor construction process results in a large amount of descriptors, which naturally leads to the concern of descriptor ranking and overfitting. Therefore, we rank all descriptors according to their feature importance and use various feature importance thresholds as a selection protocol. Here the feature importance is defined to be Gini importance<sup>37</sup> weighted by the number of trees in a forest calculated by our baseline method GBDT with scikit-learn, and train separate models to examine their predictive performances on test sets. Four different values are chosen  $(2.5 \times 10^{-4}, 5 \times 10^{-4}, 7.5 \times 10^{-4}, and 1 \times 10^{-3})$  and detailed analysis of their performances are also presented in a later section.
- **2.4. Topological learning algorithms.** In this subsection, we integrate topology and machine learning to construct topological learning algorithms. Two types of machine learning

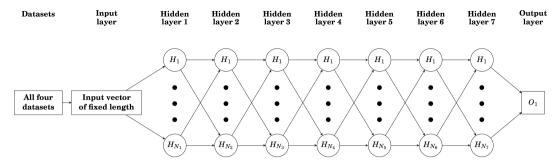


Figure 3. Illustration of the ST-DNN architecture.

algorithms, i.e., ensemble methods and DNN algorithms are used in this study. Training details are also provided.

2.4.1. Ensemble Methods. To explore strengths and weaknesses of different machine learning methods, we consider two popular ensemble methods, namely, random forest (RF) and gradient boosting decision tree (GBDT). These approaches have been widely used in solving QSAR prediction problems, as well as solvation and protein—ligand binding free energy predictions. <sup>27,39,40</sup> They naturally handle correlation between descriptors and usually do not require a sophisticated feature selection procedure. Most importantly, both RF and GBDT are essentially insensitive to parameters and robust to redundant features. Therefore, we choose these two machine learning methods as baselines in our comparison.

We have implemented RF and GBDT using the scikit-learn package (version 0.13.1).<sup>38</sup> The number of estimators is set to 2000 and the learning rate is optimized for GBDT method. For each set, 50 runs (with different random states) were done and the average result is reported in this work. Various descriptors groups discussed in section 2.2 are used as input data for RF and GBDT. More specifically, the maximum feature number is set to the square-root of the given descriptor length for both RF and GBDT models to facilitate training process given the large number of features, and it is shown that the performance of the average of sufficient runs is very decent.

2.4.2. Single-Task Deep Learning Algorithms. A neural network acts as a transformation that maps an input feature vector to an output vector. It essentially models the way a biological brain solves problems with numerous neuron units connected by axons. A typical shallow neural network consists of a few layers with neurons and uses back propagation to update weights on each layer. However, it is not able to construct hierarchical features and thus falls short in revealing more abstract properties, which makes it difficult to model complex nonlinear relationships.

A single-task deep learning algorithm, compared to shallow networks, has a wider and deeper architecture—it consists of more layers and more neurons in each layer and reveals the facets of input features at different levels. A single-task deep learning algorithm is defined for each individual prediction task and only learns data from the specific task. A representation of such single task deep neural network (ST-DNN) can be found in Figure 3, where  $N_i$  (i=1,...,7) represents the number of neurons on the ith hidden layer.

2.4.3. Multitask Learning. Multitask learning is a machine learning technique which has shown success in qualitative Merck and Tox21 prediction challenges. The main advantage of MT learning is to learn multiple tasks simultaneously and exploit commonalities as well as differences across different tasks. Another advantage of MT learning is that a small data set

with incomplete statistical distribution to establish an accurate predictive model can often be significantly benefited from relatively large data sets with more complete statistical distributions.

Suppose we have a total of T tasks and the training data for the tth task are denoted as  $(\mathbf{x}_i^t, y_i^t)_{i=1}^{N_t}$ , where t = 1, ..., T,  $i = 1, ..., N_t$ ,  $N_t$  is the number of samples of the tth tasks, with  $\mathbf{x}_i^t$  and  $y_i^t$  being the topological feature vector that consists of ESTDs and target toxicity end point of the ith molecule in the tth task, respectively. The goal of MTL and topological learning is to minimize the following loss function for all tasks simultaneously:

$$\operatorname{argmin} \sum_{i=1}^{N_t} L(y_i^t, f^t(\mathbf{x}_i^t; \{\mathbf{W}^t, \mathbf{b}^t\}))$$
(11)

where  $f^t$  is a functional of the topological feature vector  $\mathbf{x}_i^t$  parametrized by a weight vector  $\mathbf{W}^t$  and bias term  $\mathbf{b}^t$ , and L is the loss function. A typical cost function for regression is the mean squared error, thus the loss of the tth task can be defined as

loss of task 
$$t = \frac{1}{2} \sum_{i=1}^{N_t} L(\mathbf{x}_i^t, y_i^t)$$
  

$$= \frac{1}{2} \sum_{i=1}^{N_t} (y_i^t - f^t(\mathbf{x}_i^t; \{\mathbf{W}^t, \mathbf{b}^t\})^2$$
(12)

To avoid the overfitting problem, it is usually beneficial to customize the above loss function (12) by adding a regularization term on weight vectors, giving us an improved loss function for the *t*th task:

loss of task 
$$t = \frac{1}{2} \sum_{i=1}^{N_t} (y_i^t - f^t(\mathbf{x}_i^t; \{\mathbf{W}^t, \mathbf{b}^t\})^2 + \beta \left\| \mathbf{W}^t \right\|_2^2$$

$$(13)$$

where  $\|\cdot\|$  denotes the  $L_2$  norm and  $\beta$  represents a penalty constant.

In this study, the goal of topology based MTL is to learn different toxicity end points jointly and potentially improve the overall performance of multiple toxicity end points prediction models. More concretely, it is reasonable to assume that different small molecules with different measured toxicity end points comprise distinct physical or chemical features, while descriptors such as the occurrence of certain chemical structures can result in similar toxicity properties. A simple representation of multitask deep neural network (MT-DNN) for our study is shown in Figure 4, where  $N_i$  (i = 1, ..., 7) represents the number of neurons on the ith hidden layer and  $O_1$ – $O_4$  represent four predictor outputs.

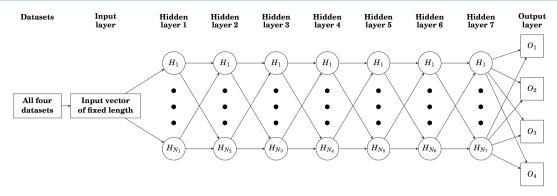


Figure 4. Illustration of the MT-DNN architecture.

2.4.4. Network Parameters and Training. Due to the large number of adjustable parameters, it is very time-consuming to optimize all possible parameter combinations. Therefore, we tune parameters within a reasonable range and subsequently evaluate their performances. The network parameters we use to train all models are four deep layers with each layer having 1000 neurons and an ADAM optimizer with 0.0001 as learning rate. It turns out that adding dropout or  $L^2$  decay does not necessarily increase the accuracy, and as a consequence, we omit these two techniques. The underlying reason may be that the ensemble results of different DNN models is essentially capable of reducing bias from individual predictions. A list of hyperparameters used to train all models can be found in Table 2.

Table 2. Proposed Hyperparameters for MT-DNN

number of epochs	1000
number of hidden layers	7
number of neurons on each layer	1000 for first 3 layers, and 100 for the next 4 layers
optimizer	ADAM
learning rate	0.001

The hyperparameter selection of DNN is known to be very complicated. In order to come up with a reasonable set of hyperparameters, we perform a grid search of each hyperparameter within a wide range. Hyperparameters in Table 2 are chosen so that we can have a reasonable training speed and accuracy. In each training epoch, molecules in each training set are randomly shuffled and then divided into mini-batches of size 200, which are then used to update parameters. When all mini-batches are traversed, an training "epoch" is done. All the training processes were done using Keras wrapper<sup>41</sup> with Theano (v0.8.2)<sup>42</sup> as the backend. All trainings were run on an Nvidia Tesla K80 GPU, and the approximate training time for a total of 1000 epochs is about 80 min.

**2.5. Evaluation Criteria.** Golbraikh et al.<sup>43</sup> proposed a protocol to determine if a QSAR model has a predictive power.

$$q^2 > 0.5 \tag{14}$$

$$R^2 > 0.6 \tag{15}$$

$$\frac{R^2 - R_0^2}{R^2} < 0.1 \tag{16}$$

$$0.85 \le k \le 1.15 \tag{17}$$

where  $q^2$  is the squared leave one out correlation coefficient for the training set,  $R^2$  is the squared Pearson correlation coefficient between the experimental and predicted toxicities for the

test set,  $R_0^2$  is the squared correlation coefficient between the experimental and predicted toxicities for the test set with the *y*-intercept being set to zero so that the regression is given by Y = kX. In addition to (15)-(17), the prediction performance will also be evaluated in terms of root-mean-square error (RMSE) and mean absolute error (MAE). The prediction coverage, or fraction of chemical predicted, of corresponding methods is also taken into account since the prediction accuracy can be increased by reducing the prediction coverage.

#### 3. RESULTS

In this section, we first give a brief description of the data sets used in this work. We then carry out our predictions by using topological and physical features in conjugation with ST-DNN and MT-DNN, and two ensemble methods, namely, RF and GBDT. The performances of these methods are compared with those of QSAR approaches used in the development of TEST software. He for the quantitative toxicity end points that we are particularly interested in, a variety of methodologies were tested and evaluated, hincluding hierarchical method, FDA method, single model method, group contribution method, and nearest neighbor method.

As for ensemble models (RF and GBDT), the training procedure follows the traditional QSAR pipeline. <sup>47</sup> A particular model is then trained to predict the corresponding toxicity end point. Note that except for specifically mentioned, all our results shown in following tables are the average outputs of 50 numerical experiments. Similarly, to eliminate randomness in neural network training, we build 50 models for each set of parameters and then use their average output as our final prediction.

Additionally, consensus of GBDT and MT-DNN is also calculated (the average of these two predictions) and its performance is also listed in tables for every data set. Finally, the best results across all descriptor combinations are presented.

**3.1. Overview of Data Sets.** This work concerns quantitative toxicity data sets. Four different quantitative toxicity data sets, namely, 96 h fathead minnow LC<sub>50</sub> data set (LC<sub>50</sub> set), 48 h *Daphnia magna* LC<sub>50</sub> data set (LC<sub>50</sub>-DM set), 40 h Tetrahymena pyriformis IGC<sub>50</sub> data set (IGC<sub>50</sub> set), and oral rat LD<sub>50</sub> data set (LD<sub>50</sub> set), are studied in this work. Among them, the LC<sub>50</sub> set reports at the concentration of test chemicals in water in milligrams per liter that cause 50% of fathead minnows to die after 96 h. Similarly, the LC<sub>50</sub>-DM set records the concentration of test chemicals in water in milligrams per liter that cause 50% *Daphnia maga* to die after 48 h. Both sets were originally downloadable from the ECOTOX aquatic toxicity database via the web site <a href="http://cfpub.epa.gov/ecotox/">http://cfpub.epa.gov/ecotox/</a> and were preprocessed using filter criterion including media type,

test location, etc.  $^{44}$  The third set, IGC $_{50}$  set, measures the 50% growth inhibitory concentration of Tetrahymena pyriformis organism after 40 h. It was obtained from Schultz and co-workers. The end point LD $_{50}$  represents the amount of chemicals that can kill half of rats when orally ingested. The LD $_{50}$  was constructed from ChemIDplus databse (http://chem.sis.nlm.nih. gov/chemidplus/chemidheavy.jsp) and then filtered according to several criteria.  $^{44}$ 

The final sets used in this work are identical to those that were preprocessed and used to develop the Toxicity Estimation Software Tool (TEST).<sup>44</sup> TEST was developed to estimate chemical toxicity using various QSAR methodologies and is very convenient to use as it does not require any external programs. It follows the general QSAR workflow—it first calculates 797 2D molecular descriptors and then predicts the toxicity of a given target by utilizing these precalculated molecular descriptors.

All molecules are in either 2D sdf format or SMILE string, and their corresponding toxicity end points are available on the TEST web site. It should be noted that we are particularly interested in predicting quantitative toxicity end points so other data sets that contain qualitative end points or physical properties were not used. Moreover, different toxicity end points have different units. The units of LC50, LC50-DM, IGC50 end points are  $-\log_{10}(T \text{ mol/L})$ , where T represents corresponding end point. For LD<sub>50</sub> set, the units are  $-\log_{10}(LD_{50} \text{ mol/kg})$ . Although the units are not exactly the same, it should be pointed out that no additional attempt was made to rescale the values since end points are of the same magnitude order. These four data sets also differ in their sizes, ranging from hundreds to thousands, which essentially challenges the robustness of our methods. A detailed statistics table of four data sets is presented in Table 3.

The number inside the parentheses indicates the actual number of molecules that we use for developing models in this work. Note that for the first three data sets (i.e.,  $LC_{50}$ ,  $LC_{50}$ -DM, and  $IGC_{50}$  set), all molecules were properly included. However, for  $LD_{50}$  set, some molecules involved element As were dropped out due to force field failure. Apparently, the TEST tool encounters a similar problem since results from two TEST models are unavailable, and the coverage (fraction of molecules predicted) from various TEST models is always smaller than one. The overall coverage of our models is always higher than that of TEST models, which indicates a wider applicable domain of our models.

**3.2. Feathead Minnow LC**<sub>50</sub> **Test Set.** The feathead minnow LC<sub>50</sub> set was randomly divided into a training set (80% of the entire set) and a test set (20% of the entire set),<sup>44</sup> based on which a variety of TEST models were built. Table 4 shows the performances of five TEST models, the TEST consensus obtained by the average of all independent TEST predictions, four proposed methods, and two consensus results obtained from averaging over present RF, GBDT, ST-DNN, and MT-DNN results. TEST consensus gives the best prediction among TEST results, reporting a correlation coefficient of 0.728 and RMSE of

Table 4. Comparison of Prediction Results for the Fathead Minnow  $LC_{50}$  Test Set

method	$R^2$	$\frac{R^2 - R_0^2}{R^2}$	k	RMSE	MAE	coverage
hierarchical <sup>44</sup>	0.710	0.075	0.966	0.801	0.574	0.951
single model <sup>44</sup>	0.704	0.134	0.960	0.803	0.605	0.945
FDA <sup>44</sup>	0.626	0.113	0.985	0.915	0.656	0.945
group contribution <sup>44</sup>	0.686	0.123	0.949	0.810	0.578	0.872
nearest neighbor <sup>44</sup>	0.667	0.080	1.001	0.876	0.649	0.939
TEST consensus <sup>44</sup>	0.728	0.121	0.969	0.768	0.545	0.951
	Re	sults with	ESTDs			
RF	0.661	0.364	0.946	0.858	0.638	1.000
GBDT	0.672	0.103	0.958	0.857	0.612	1.000
ST-DNN	0.675	0.031	0.995	0.862	0.601	1.000
MT-DNN	0.738	0.012	1.015	0.763	0.514	1.000
consensus	0.740	0.087	0.956	0.755	0.518	1.000
Results w	ith Only	Auxiliary	Molecu	lar Descri	ptors	
RF	0.744	0.467	0.947	0.784	0.560	1.000
GBDT	0.750	0.148	0.962	0.736	0.511	1.000
ST-DNN	0.598	0.044	0.982	0.959	0.648	1.000
MT-DNN	0.771	0.003	1.010	0.705	0.472	1.000
consensus	0.787	0.105	0.963	0.679	0.464	1.000
	Results	with All	Descript	tors		
RF	0.727	0.322	0.948	0.782	0.564	1.000
GBDT	0.761	0.102	0.959	0.719	0.496	1.000
ST-DNN	0.692	0.010	0.997	0.822	0.568	1.000
MT-DNN	0.769	0.009	1.014	0.716	0.466	1.000
Consensus	0.789	0.076	0.959	0.677	0.446	1.000

0.768 log(mol/L). As Table 4 indicates, our MT-DNN model outperforms TEST consensus both in terms of  $R^2$  and RMSE with only ESTDs as input. When physical descriptors are independently used or combined with ESTDs, the prediction accuracy can be further improved to a higher level, with  $R^2$  of 0.771 and RMSE of 0.705 log(mol/L). The best result is generated by consensus method using all descriptors, with  $R^2$  of 0.789 and RMSE of 0.677 log(mol/L).

**3.3.** Daphnia magna LC<sub>50</sub> Test Set. The Daphnia magna LC50 set is the smallest in terms of set size, with 283 training molecules and 70 test molecules, respectively. However, it brings difficulties to building robust QSAR models given the relatively large number of descriptors. Indeed, five independent models in TEST software give significantly different predictions, as indicated by RMSEs shown in Table 5 ranging from 0.810 to 1.190 log units. Though the RMSE of group contribution is the smallest, its coverage is only 0.657% which largely restricts this method's applicability. Additionally, its R² value is inconsistent with its RMSE and MAE. Since ref 44 states that "The consensus method achieved the best results in terms of both prediction accuracy and coverage", these usually low RMSE and MAE values might be typos.

We also notice that our nonmultitask models that contain ESTDs result in very large deviation from experimental values. Indeed, the overfitting issue challenges traditional machine learning approaches especially when the number of samples is less

Table 3. Statistics of Quantitative Toxicity Data Sets

	total no. of mols	train set size	test set size	max value	min value
LC <sub>50</sub> set	823	659	164	9.261	0.037
LC <sub>50</sub> -DM set	353	283	70	10.064	0.117
IGC <sub>50</sub> set	1792	1434	358	6.36	0.334
LD <sub>50</sub> set	7413 (7403)	5931 (5924)	1482 (1479)	7.201	0.291

Table 5. Comparison of Prediction Results for the Daphnia magna LC<sub>50</sub> Test Set

method	$R^2$	$\frac{R^2 - R_0^2}{R^2}$	k	RMSE	MAE	coverag
hierarchical <sup>44</sup>	0.695	0.151	0.981	0.979	0.757	0.886
single model <sup>44</sup>	0.697	0.152	1.002	0.993	0.772	0.871
FDA <sup>44</sup>	0.565	0.257	0.987	1.190	0.909	0.900
group contribution <sup>44</sup>	0.671	0.049	0.999	0.803 <sup>a</sup>	0.620 <sup>a</sup>	0.657
nearest neighbor <sup>44</sup>	0.733	0.014	1.015	0.975	0.745	0.871
TEST consensus <sup>44</sup>	0.739	0.118	1.001	0.911	0.727	0.900
		Resul	ts with ESTDs			
RF	0.441	1.177	0.957	1.300	0.995	1.000
GBDT	0.467	0.440	0.972	1.311	0.957	1.000
ST-DNN	0.446	0.315	0.927	1.434	0.939	1.000
MT-DNN	0.788	0.008	1.002	0.805	0.592	1.000
consensus	0.681	0.266	0.970	0.977	0.724	1.000
		Results with Only A	uxiliary Molecular De	escriptors		
RF	0.479	1.568	0.963	1.261	0.946	1.000
GBDT	0.495	0.613	0.959	1.238	0.926	1.000
ST-DNN	0.430	0.404	0.921	1.484	1.034	1.000
MT-DNN	0.705	0.009	1.031	0.944	0.610	1.000
consensus	0.665	0.359	0.945	1.000	0.732	1.000
		Results w	rith All Descriptors			
RF	0.460	1.244	0.955	1.274	0.958	1.000
GBDT	0.505	0.448	0.961	1.235	0.905	1.000
ST-DNN	0.459	0.278	0.933	1.407	1.004	1.000
MT-DNN	0.726	0.003	1.017	0.905	0.590	1.000
consensus	0.678	0.282	0.953	0.978	0.714	1.000

These values are inconsistent with  $R^2 = 0.671$ .

than the number of descriptors. The advantage of MT-DNN model is to extract information from related tasks and our numerical results show that the predictions do benefit from MTL architecture. For models using ESTDs, physical descriptors, and all descriptors, the  $R^2$  has been improved from around 0.5 to 0.788, 0.705, and 0.726, respectively. It is worth mentioning that our ESTDs yield the best results, which proves the power of persistent homology. This result suggests that by learning related problems jointly and extracting shared information from different data sets, MT-DNN architecture can simultaneously perform multiple prediction tasks and enhances performances especially on small data sets.

3.4. Tetraphymena pyriformis IGC<sub>50</sub> Test Set. IGC<sub>50</sub> set is the second largest QSAR toxicity set that we want to study. The diversity of molecules of in IGC<sub>50</sub> set is low and the coverage of TEST methods is relatively high compared to previous  $LC_{50}$  sets. As shown in Table 6, the  $R^2$  of different TEST methods fluctuates from 0.600 to 0.764, and test consensus prediction again yields the best result for TEST software with  $R^2$  of 0.764. As for our models, the  $R^2$  of MT-DNN with different descriptors spans a range of 0.038 (0.732 to 0.770), which indicates that our MT-DNN not only takes care of overfitting problem but also is insensitive to data sets. Although ESTDs slightly underperform compared to physical descriptors, its MT-DNN results are able to defeat most TEST methods except for the FDA method. When all descriptors are used, predictions by GBDT and MT-DNN outperform TEST consensus, with  $R^2$ of 0.787 and RMSE of 0.455 log(mol/L). The best result is again given by consensus method using all descriptors, with  $R^2$ of 0.802 and RMSE of 0.438 log(mol/L).

**3.5. Oral Rat LD**<sub>50</sub> **Test Set.** The oral rat LD<sub>50</sub> set contains the largest molecule pool with 7413 compounds. However, none of methods is able to provide a 100% coverage of this data set.

Table 6. Comparison of Prediction Results for the Tetraphymena pyriformis IGC<sub>50</sub> Test Set

a 1	$R^2$	$R^2 - R_0^2$	1	DMCE	MAE	
method		$\frac{R^2 - R_0^2}{R^2}$	k	RMSE	MAE	coverage
hierarchical <sup>44</sup>	0.719	0.023	0.978	0.539	0.358	0.933
FDA <sup>44</sup>	0.747	0.056	0.988	0.489	0.337	0.978
group contribution <sup>44</sup>	0.682	0.065	0.994	0.575	0.411	0.955
nearest neighbor <sup>44</sup>	0.600	0.170	0.976	0.638	0.451	0.986
TEST consensus <sup>44</sup>	0.764	0.065	0.983	0.475	0.332	0.983
	Re	sults with	i ESTDs			
RF	0.625	0.469	0.966	0.603	0.428	1.000
GBDT	0.705	0.099	0.984	0.538	0.374	1.000
ST-DNN	0.708	0.011	1.000	0.537	0.374	1.000
MT-DNN	0.723	0.000	1.002	0.517	0.378	1.000
consensus	0.745	0.121	0.980	0.496	0.356	1.000
Results w	ith Only	Auxiliary	y Molecu	lar Descri	ptors	
RF	0.738	0.301	0.978	0.514	0.375	1.000
GBDT	0.780	0.065	0.992	0.462	0.323	1.000
ST-DNN	0.678	0.052	0.972	0.587	0.357	1.000
MT-DNN	0.745	0.002	0.995	0.498	0.348	1.000
consensus	0.789	0.073	0.989	0.451	0.317	1.000
	Results	s with All	Descrip	tors		
RF	0.736	0.235	0.981	0.510	0.368	1.000
GBDT	0.787	0.054	0.993	0.455	0.316	1.000
ST-DNN	0.749	0.019	0.982	0.506	0.339	1.000
MT-DNN	0.770	0.000	1.001	0.472	0.331	1.000
consensus	0.802	0.066	0.987	0.438	0.305	1.000

The results of single model method or group contribution method were not properly built for the entire set. 44 It was noted that LD<sub>50</sub> values of this data set are relatively difficult to predict as they have a higher experimental uncertainty.<sup>50</sup> As shown in Table 7, results of two TEST approaches, i.e., single model and group contribution, were not reported for this problem. The

Table 7. Comparison of Prediction Results for the Oral Rat  $LD_{50}$  Test Set

.1 1	n2	$R^2 - R_0^2$	1	DIAGE	3.645	
method	$R^2$	$\frac{R^2 - R_0^2}{R^2}$	k	RMSE	MAE	coverage
hierarchical <sup>44</sup>	0.578	0.184	0.969	0.650	0.460	0.876
FDA <sup>44</sup>	0.557	0.238	0.953	0.657	0.474	0.984
nearest neighbor <sup>44</sup>	0.557	0.243	0.961	0.656	0.477	0.993
TEST consensus <sup>44</sup>	0.626	0.235	0.959	0.594	0.431	0.984
	]	Results wi	th ESTD	s		
RF	0.586	0.823	0.949	0.626	0.469	0.999
GBDT	0.598	0.407	0.960	0.613	0.455	0.999
ST-DNN	0.601	0.006	0.991	0.612	0.446	0.999
MT-DNN	0.613	0.000	1.000	0.601	0.442	0.999
consensus	0.631	0.384	0.956	0.586	0.432	0.999
Results	with O	nly Auxilia	ry Molec	ular Descr	riptors	
RF	0.597	0.825	0.946	0.619	0.463	0.997
GBDT	0.605	0.385	0.958	0.606	0.455	0.997
ST-DNN	0.593	0.008	0.992	0.618	0.447	0.997
MT-DNN	0.604	0.003	0.995	0.609	0.445	0.997
consensus	0.637	0.350	0.957	0.581	0.433	0.997
	Resu	ılts with A	dl Descrij	ptors		
RF	0.619	0.728	0.949	0.603	0.452	0.997
GBDT	0.630	0.328	0.960	0.586	0.441	0.997
ST-DNN	0.614	0.006	0.991	0.601	0.436	0.997
MT-DNN	0.626	0.002	0.995	0.590	0.430	0.997
consensus	0.653	0.306	0.959	0.568	0.421	0.997

TEST consensus result improves overall prediction accuracy of other TEST methods by about 10%, however, other non-consensus methods all yield low  $R^2$  and high RMSE.

For our models, all results outperform those of nonconsensus methods of TEST. In particular, GBDT and MT-DNN with all descriptors yield the best (similar) results, giving slightly better results compared to TEST consensus. Meanwhile, our predictions are also relatively stable for this particular set as  $R^2$ s do not essentially fluctuate. It should also be noted that our ESTDs have slightly higher coverage than physical descriptors (all combined descriptors) since two molecules in the test set that contains the As element cannot be properly optimized for energy computation. However, this is not an issue with our persistent homology computation. The consensus method using all descriptors again yields the best results for all combinations, with optimal  $R^2$  of 0.653 and RMSE of 0.568  $\log(\text{mol/kg})$ .

# 4. DISCUSSION

In this section, we will discuss how ESTDs bring new insights to quantitative toxicity end points and how ensemble based topological learning can improve overall performances.

**4.1.** Impact of Descriptor Selection and Potential Overfitting. A major concern for the proposed models is descriptor redundancy and potential overfitting. To address this issue, four different sets of high-importance descriptors are selected by a threshold to perform prediction tasks as described

in section 2.3. Table 8 below shows the results of MT-DNN using these four different descriptor sets for LC50 set. Results for the other three remaining sets are provided in the Supporting Information.

Table 8 shows performance with respect to different numbers of descriptors. When the number of descriptors is increased from 222, 254, 308, 411 to 1030, RMSE does not increase and  $R^2$  does not change much. This behavior suggests that our models are essentially insensitive to the number of descriptors and thus there is little overfitting. MT-DNN architecture takes care of overfitting issues by successive feature abstraction, which naturally mitigates noise generated by less important descriptors. MT-DNN architecture can also potentially take advantage over related tasks, which in turn reduces the potential overfitting on single data set by the alternative training procedure.

Similar behaviors have also been observed for the remaining three data sets, as presented in the Supporting Information. Therefore, our MT-DNN architecture is very robust against feature selection and can avoid overfitting.

**4.2. Predictive Power of ESTDs for Toxicity.** One of the main objectives of this study is to understand toxicity of small molecules from a topological point of view. It is important to see if ESTDs alone can match those methods proposed in TEST software. When all ESTDs (group 6) and MT-DNN architecture are used for toxicity prediction, we observe following results:

- LC<sub>50</sub> and LC<sub>50</sub>-DM sets. Models using only ESTDs achieve higher accuracy than the TEST consensus method.
- $LD_{50}$  set. Consensus result of ESTDs tops TEST software in terms of both  $R^2$  and RMSE and MT-DNN results outperform all nonconsensus TEST methods.
- IGC<sub>50</sub> set. ESTDs are perform slightly worse than TEST consensus. However, MT-DNN with ESTDs still yields better results than most nonconsensus TEST methods except FDA.

It is evident that our ESTDs along with MT-DNN architecture have a strong predictive power for all kinds of toxicity end points. The ability of MT-DNN to learn from related toxicity end points has resulted in a substantial improvement over ensemble methods such as GBDT. Along with physical descriptors calculated by our in-house MIBPB, we can obtain state-of-the-art results for all four quantitative toxicity end points.

**4.3.** Alternative Element Specific Networks for Generating ESTDs. Apart from the element specific networks proposed in Table 1, we also use alternative element specific networks listed below in Table 9 to perform the same prediction tasks. Instead of using two types of element-specific networks, we only consider two-element networks to generate ESTDs, which essentially puts more emphasis on intramolecular interaction aspect. Eventually, this new construction yields 30 different element specific networks (9 + 8 + 7 + 6), and a total of 840 ESTDs  $(30 \times 28)$  is calculated and used for prediction. On the LC<sub>50</sub>, IGC<sub>50</sub>, and LD<sub>50</sub> sets, overall performances of the

Table 8. Results of Selected Descriptor Groups for LC<sub>50</sub> Set

threshold	no. of descriptors	$R^2$	$\frac{R^2 - R_0^2}{R^2}$	k	RMSE	MAE	coverage
0.0	1030	0.769	0.009	1.014	0.716	0.466	1.000
$2.5 \times 10^{-4}$	411	0.784	0.051	0.971	0.685	0.459	1.000
$5 \times 10^{-4}$	308	0.764	0.062	0.962	0.719	0.470	1.000
$7.5 \times 10^{-4}$	254	0.772	0.064	0.958	0.708	0.468	1.000
$1 \times 10^{-3}$	222	0.764	0.063	0.963	0.717	0.467	1.000

Table 9. Alternative Element Specific Networks Used to Characterize Molecules

network type	element specific networks
two-element	$\begin{array}{l} \{b_{i},c_{j}\},\text{where}b_{i}\in\mathcal{B},c_{j}\in C,i\in\{1,,4\},j\in\{1,,10\},\text{and}i< j\\ \text{where}\mathcal{B}=\{\text{H, C, N, O}\}\text{and}C=\{\text{H, C, N, O, F, P, S, Cl, Br, I}\} \end{array}$

new ESTDs can be improved slightly. However, on the LC<sub>50</sub>-DM set, the accuracy is comparably lower (still higher than TEST consensus). Detailed performances of these ESTDs are presented in the Supporting Information. Thus, the predictive power of our ESTDs is not sensitive to the choice of element specific networks as long as reasonable element types are included.

4.4. Potential Improvement with Consensus Tools. In this work, we also propose consensus method as discussed in section 3. The idea of consensus is to train different models on the same set of descriptors and average across all predicted values. The underlying mechanism is to take advantage of system errors generated by different machine learning approaches with a possibility to reduce bias for the final prediction.

As we notice from section 3, consensus method offers a considerable boost in prediction accuracy. For reasonably large sets except LC50-DM set, consensus models turn out to give the best predictions. When it comes to a small set (LC<sub>50</sub>-DM set), consensus models perform worse than MT-DNN. It is likely due to the fact that a large number of descriptors may cause overfitting issues for most machine learning algorithms, and consequently generate large deviations, which eventually result in a large error of consensus methods. Thus, it should be a good idea to preform prediction tasks with both MT-DNN and consensus methods, depending on the size of data sets, to take advantage of both approaches.

# 5. CONCLUSION

Toxicity refers to the degree of damage a substance causes to an organism, such as an animal, bacterium, or plant, and can be qualitatively or quantitatively measured by experiments. Experimental measurement of quantitative toxicity is extremely valuable, but is typically expensive and time-consuming, in addition to potential ethic concerns. Theoretical prediction of quantitative toxicity has become a useful alternative in pharmacology and environmental science. A wide variety of methods has been developed for toxicity prediction in the past. The performances of these methods depend not only on the descriptors, but also on machine learning algorithms, which makes the model evaluation a difficult task.

In this work, we introduce a novel method, called element specific topological descriptor (ESTD), for the characterization and prediction of small molecular quantitative toxicity. Additionally physical descriptors based on established physical models are also developed to enhance the predictive power of ESTDs. These new descriptors are integrated with a variety of advanced machine learning algorithms, including two deep neural networks (DNNs) and two ensemble methods (i.e., random forest (RF) and gradient boosting decision tree (GBDT)) to construct topological learning strategies for quantitative toxicity analysis and prediction.

Four quantitative toxicity data sets, i.e., 96 h fathead minnow LC<sub>50</sub> data set (LC<sub>50</sub> set), 48 h Daphnia magna LC<sub>50</sub> data set (LC<sub>50</sub>-DM set), 40 h Tetrahymena pyriformis IGC<sub>50</sub> data set (IGC<sub>50</sub> set), and oral rat LD<sub>50</sub> data set (LD<sub>50</sub> set), are used in the present study. Comparison has also been made to the stateof-the-art approaches given in the literature Toxicity Estimation Software Tool (TEST)<sup>44</sup> listed by United States Environmental Protection Agency. Our numerical experiments indicate that the proposed ESTDs are as competitive as individual methods in TEST. Aided with physical descriptors and MT-DNN architecture, ESTDs are able to establish new state-of-the-art predictions for quantitative toxicity data sets. Additionally, MT deep learning algorithms are typically more accurate than ensemble methods such as RF and GBDT.

It is worthy to note that the proposed new descriptors are very easy to generate and thus have almost 100% coverage for all molecules, indicating their broader applicability to practical toxicity analysis and prediction. In fact, our topological descriptors are much easier to construct than physical descriptors, which depend on physical models and force fields. The present work indicates that ESTDs are a new class of powerful descriptors for small molecules.

## ASSOCIATED CONTENT

# S Supporting Information

The Supporting Information is available free of charge on the ACS Publications website at DOI: 10.1021/acs.jcim.7b00558.

Detailed performances of four groups of descriptors based on feature importance threshold with MT-DNN are presented in Tables S1-S4. Results with the ESTDs proposed in the Discussion using different algorithms are listed in Tables S5-S12 (PDF)

## AUTHOR INFORMATION

#### **Corresponding Author**

\*E-mail: wei@math.msu.edu.

ORCID ®

Kedi Wu: 0000-0003-3737-5415 Guo-Wei Wei: 0000-0002-5781-2937

#### **Funding**

This work was supported in part by NSF Grants IIS-1302285 and DMS-1721024 and the MSU Center for Mathematical Molecular Biosciences Initiative.

# Notes

The authors declare no competing financial interest. Software for computing ESTDs and auxiliary molecular descriptors is available as an online server at http://weilab. math.msu.edu/TopTox/ and http://weilab.math.msu.edu/ MIBPB/, respectively. The source code for computing ESTDs can be found in the Supporting Information.

### REFERENCES

- (1) Deeb, O.; Goodarzi, M. In Silico Quantitative Structure Toxicity Relationship of Chemical Compounds: Some Case Studies. Curr. Drug Saf. 2012, 7, 289-297.
- (2) Kauffman, G. W.; Jurs, P. C. QSAR and K-Nearest Neighbor Classification Analysis of Selective Cyclooxygenase-2 Inhibitors Using Topologically-Based Numerical Descriptors. J. Chem. Inf. Comput. Sci. **2001**, *41*, 1553–1560.
- (3) Ajmani, S.; Jadhav, K.; Kulkarni, S. A. Three-Dimensional QSAR Using the K-Nearest Neighbor Method and Its Interpretation. J. Chem. Inf. Model. 2006, 46, 24-31.
- (4) Si, H.; Wang, T.; Zhang, K.; Duan, Y.-B.; Yuan, S.; Fu, A.; Hu, Z. Quantitative Structure Activity Relationship Model for Predicting the Depletion Percentage of Skin Allergic Chemical Substances of Glutathione. Anal. Chim. Acta 2007, 591, 255-264.
- (5) Du, H.; Wang, J.; Hu, Z.; Yao, X.; Zhang, X. Prediction of Fungicidal Activities of Rice Blast Disease Based on Least-Squares

- Support Vector Machines and Project Pursuit Regression. J. Agric. Food Chem. 2008, 56, 10785-10792.
- (6) Svetnik, V.; Liaw, A.; Tong, C.; Culberson, J. C.; Sheridan, R. P.; Feuston, B. P. Random Forest: A Classification and Regression Tool for Compound Classification and QSAR Modeling. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 1947–1958.
- (7) Liu, P.; Long, W. Current Mathematical Methods Used in QSAR/QSPR Studies. *Int. J. Mol. Sci.* **2009**, *10*, 1978–1998.
- (8) Caruana, R. Learning to Learn; Springer, 1998; pp 95-133.
- (9) Widmer, C.; Rätsch, G. Multitask Learning in Computational Biology. Workshop on Unsupervised and Transfer Learning; 2012; pp 207–216.
- (10) Xu, Q.; Yang, Q. A Survey of Transfer and Multitask Learning in Bioinformatics. *J. Comput. Sci. Eng.* **2011**, *5*, 257–268.
- (11) LeCun, Y.; Bengio, Y.; Hinton, G. Deep Learning. *Nature* **2015**, 521, 436–444.
- (12) Schmidhuber, J. Deep Learning in Neural Networks: An Overview. *Neural Networks* **2015**, *61*, 85–117.
- (13) Dahl, G. E.; Dong, Y.; Deng, L.; Acero, A. Context-Dependent Pre-Trained Deep Neural Networks for Large-Vocabulary Speech Recognition. *IEEE Trans. Audio, Speech, and Language Process* **2012**, 20, 30–42.
- (14) Deng, L.; Hinton, G.; Kingsbury, B. New Types of Deep Neural Network Learning for Speech Recognition and Related Applications: An Overview. Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference On. 2013; pp 8599—8603.10.1109/ICASSP.2013.6639344
- (15) Socher, R.; Bengio, Y.; Manning, C. D. Deep Learning for NLP (without Magic). Tutorial Abstracts of ACL 2012. 2012; pp 5–5.
- (16) Sutskever, I.; Vinyals, O.; Le, Q. V. Sequence to Sequence Learning with Neural Networks. *Adv. Neural Infor. Process. Syst.* **2014**; pp 3104–3112.
- (17) Edelsbrunner, H.; Letscher, D.; Zomorodian, A. Topological Persistence and Simplification. *Discrete Comput. Geom.* **2002**, 28, 511–533.
- (18) Zomorodian, A.; Carlsson, G. Computing Persistent Homology. *Discrete Comput. Geom.* **2005**, 33, 249–274.
- (19) Xia, K. L.; Wei, G. W. Persistent Homology Analysis of Protein Structure, Flexibility and Folding. *Int. J. Numer. Method Biomed. Eng.* **2014**, 30, 814–844.
- (20) Xia, K. L.; Feng, X.; Tong, Y. Y.; Wei, G. W. Persistent Homology for the Quantitative Prediction of Fullerene Stability. *J. Comput. Chem.* **2015**, *36*, 408–422.
- (21) Xia, K. L.; Zhao, Z. X.; Wei, G. W. Multiresolution Topological Simplification. *J. Comput. Biol.* **2015**, 22, 887.
- (22) Xia, K. L.; Zhao, Z. X.; Wei, G. W. Multiresolution Persistent Homology for Excessively Large Biomolecular Datasets. *J. Chem. Phys.* **2015**, *143*, 134103.
- (23) Xia, K. L.; Wei, G. W. Persistent Topology for cryo-EM Data Analysis. Int. J. Numer. Method Biomed. Eng. 2015, 31, e02719.
- (24) Wang, B.; Wei, G. W. Object-Oriented Persistent Homology. J. Comput. Phys. 2016, 305, 276–299.
- (25) Cang, Z. X.; Mu, L.; Wu, K.; Opron, K.; Xia, K.; Wei, G.-W. A Topological Approach to Protein Classification. *Mol. Based Math. Biol.* **2015**, 3, 140–162.
- (26) Cang, Z. X.; Wei, G. W. Analysis and Prediction of Protein Folding Energy Changes upon Mutation by Element Specific Persistent Homology. *Bioinformatics* **2017**, *33*, 3549–3557.
- (27) Cang, Z. X.; Wei, G. W. Integration of Element Specific Persistent Homology and Machine Learning for Protein-Ligand Binding Affinity Prediction. *Int. J. Numer. Method Biomed. Eng.* **2017**, e2914.
- (28) Cang, Z. X.; Wei, G. W. TopologyNet: Topology Based Deep Convolutional Neural Networks for Biomolecular Property Predictions. *PLoS Comput. Biol.* **2017**, *13* (7), e1005690.
- (29) Edelsbrunner, H.; Letscher, D.; Zomorodian, A. Topological Persistence and Simplification. *41st Ann. Symp. Found. Comput. Sci.* **2000**, 454–463.

- (30) Carlsson, G.; Zomorodian, A.; Collins, A.; Guibas, L. J. Persistence Barcodes for Shapes. *Int. J. Shape Model.* **2005**, *11*, 149–187.
- (31) Ghrist, R. Barcodes: The Persistent Topology of Data. *Bull. Am. Math. Soc.* **2008**, *45*, 61–75.
- (32) Cang, Z. X.; Mu, L.; Wei, G. W. Representability of Algebraic Topology for Biomolecules in Machine Learning Based Scoring and Virtual Screening. *PLoS Comput. Biol.* **2018**, *14* (1), e1005929.
- (33) Wang, J.; Wang, W.; Kollman, P. A.; Case, D. A. Automatic Atom Type and Bond Type Perception in Molecular Mechanical Calculations. *J. Mol. Graphics Modell.* **2006**, 25, 247–260.
- (34) Wang, J.; Wolf, R. M.; Caldwell, J. W.; Kollman, P. A.; Case, D. A. Development and Testing of a General Amber Force Field. *J. Comput. Chem.* **2004**, 25, 1157–1174.
- (35) Liu, B.; Wang, B.; Zhao, R.; Tong, Y.; Wei, G. W. ESES: Software for Eulerian Solvent Excluded Surface. *J. Comput. Chem.* **2017**, 38, 446–466.
- (36) Chen, D.; Chen, Z.; Chen, C.; Geng, W. H.; Wei, G. W. MIBPB: A Software Package for Electrostatic Analysis. *J. Comput. Chem.* **2011**, 32, 756.
- (37) Breiman, L. Random Forests. Mach. Learn. 2001, 45, 5-32.
- (38) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; Vanderplas, J.; Passos, A.; Cournapeau, D.; Brucher, M.; Perrot, M.; Duchesnay, E. Scikit-Learn: Machine Learning in Python. *J. Mach. Learn. Res.* **2011**, *12*, 2825–2830.
- (39) Wang, B.; Wang, C.; Wu, K.; Wei, G.-W. Breaking the Polar-Nonpolar Division in Solvation Free Energy Prediction. *J. Comput. Chem.* **2018**, *39*, 217–233.
- (40) Wang, B.; Zhao, Z.; Nguyen, D. D.; Wei, G. W. Feature Functional Theory Binding Predictor (FFT-BP) for the Blind Prediction of Binding Free Energy. *Theor. Chem. Acc.* **2017**, *136*, 55.
- (41) Chollet, F. Keras. https://github.com/fchollet/keras (accessed 2015).
- (42) Theano Development Team.. Theano: A Python Framework for Fast Computation of Mathematical Expressions. *arXiv.org* **2016**, No. abs/1605.02688.
- (43) Golbraikh, A.; Shen, M.; Xiao, Z.; Xiao, Y.-D.; Lee, K.-H.; Tropsha, A. Rational Selection of Training and Test Sets for the Development of Validated QSAR Models. *J. Comput.-Aided Mol. Des.* **2003**, *17*, 241–253.
- (44) Martin, T. User's Guide for T.E.S.T. (version 4.2) (Toxicity Estimation Software Tool): A Program to Estimate Toxicity from Molecular Structure; USEPA, 2016.
- (45) Martin, T. M.; Harten, P.; Venkatapathy, R.; Das, S.; Young, D. M. A Hierarchical Clustering Methodology for the Estimation of Toxicity. *Toxicol. Mech. Methods* **2008**, *18*, 251–266.
- (46) Martin, T. M.; Young, D. M. Prediction of the Acute Toxicity (96-H LC50) of Organic Compounds to the Fathead Minnow (Pimephales Promelas) Using a Group Contribution Method. *Chem. Res. Toxicol.* **2001**, *14*, 1378–1385.
- (47) Nantasenamat, C.; Isarankura-Na-Ayudhya, C.; Naenna, T.; Prachayasittikul, V. A Practical Overview of Quantitative Structure-Activity Relationship. *EXCLI J.* **2009**, *8*, 74–88.
- (48) Akers, K. S.; Sinks, G. D.; Schultz, T. W. Structure—Toxicity Relationships for Selected Halogenated Aliphatic Chemicals. *Environ. Toxicol. Pharmacol.* **1999**, *7*, 33–39.
- (49) Zhu, H.; Tropsha, A.; Fourches, D.; Varnek, A.; Papa, E.; Gramatica, P.; Oberg, T.; Dao, P.; Cherkasov, A.; Tetko, I. V. Combinatorial QSAR Modeling of Chemical Toxicants Tested Against Tetrahymena Pyriformis. *J. Chem. Inf. Model.* **2008**, *48*, 766–784.
- (50) Zhu, H.; Martin, T. M.; Ye, L.; Sedykh, A.; Young, D. M.; Tropsha, A. Quantitative Structure- Activity Relationship Modeling of Rat Acute Toxicity by Oral Exposure. *Chem. Res. Toxicol.* **2009**, 22, 1913–1921.