

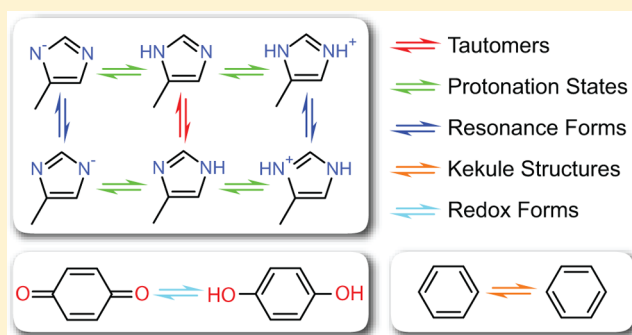
The Valence State Combination Model: A Generic Framework for Handling Tautomers and Protonation States

Sascha Urbaczek,[†] Adrian Kolodzik,[†] and Matthias Rarey*

University of Hamburg, Center for Bioinformatics (ZBH), Bundesstrasse 43, 20146 Hamburg, Germany

Supporting Information

ABSTRACT: The consistent handling of molecules is probably the most basic and important requirement in the field of cheminformatics. Reliable results can only be obtained if the underlying calculations are independent of the specific way molecules are represented in the input data. However, ensuring consistency is a complex task with many pitfalls, an important one being the fact that the same molecule can be represented by different valence bond structures. In order to achieve reliability, a cheminformatics system needs to solve two fundamental problems. First, different choices of valence bond structures must be identified as the same molecule. Second, for each molecule all valence bond structures relevant to the context must be taken into consideration. The latter is especially important with regard to tautomers and protonation states, as these have considerable influence on physicochemical properties of molecules. We present a comprehensive method for the rapid and consistent generation of reasonable tautomers and protonation states for molecules relevant in the context of drug design. This method is based on a generic scheme, the Valence State Combination Model, which has been designed for the enumeration and scoring of valence bond structures in large data sets. In order to ensure our method's consistency, we have developed procedures which can serve as a general validation scheme for similar approaches. The analysis of both the average number of generated structures and the associated runtimes shows that our method is perfectly suited for typical cheminformatics applications. By comparison with frequently used and curated public data sets, we can demonstrate that the tautomers and protonation state produced by our method are chemically reasonable.



INTRODUCTION

One of the most fundamental requirements in cheminformatics is the consistent handling of molecules from different sources. There is always the implicit assumption that the results of cheminformatics software applications are only dependent on the actual compounds and not on the way these are provided in the input data. Yet, apart from problems arising from the interpretation of data from chemical file formats, there are certain ambiguities in the way molecules are represented which considerably complicate this task. Virtually all modern cheminformatics systems are based on a description of molecules by valence bond structures (Lewis structures). The inherent limitations of this molecular representation and their implications on tautomer generation have been recently discussed in detail by Sayle.¹ In the following, we will largely follow the nomenclature used in his publication and refer back to particular aspects mentioned therein.

The main problem with respect to consistency is the fact that different valence bond structures can represent the same molecule. Some of these correspond to distinct chemical entities, e.g., tautomers and protonation states, whereas others are artifacts of valence theory, i.e., resonance forms and Kekule structures. In some contexts even oxidation states may be

interpreted as alternative forms of the same molecule (see Figure 1 for examples).

From a formal point of view, each of these valence bond structures could be chosen as a representation for a particular compound. In practice, not all members of this set of alternatives are equally likely to be encountered due to automated normalization procedures and manual curation. However, despite all these efforts, a certain degree of ambiguity cannot be entirely avoided. The resulting implications for cheminformatics systems in general² and large compound databases in particular³ have been thoroughly investigated in the literature. In his publication, Sayle¹ has identified five specific tasks associated with the ambiguities of molecular representations. With respect to consistency; these are comparison (#1) and, more importantly, canonicalization (#2). A cheminformatics system must be able to reliably identify and treat alternative valence bond structures as the same molecule. This is usually done by conversion to a canonical form which serves as input for subsequent methods.

Received: December 6, 2013

Published: February 18, 2014

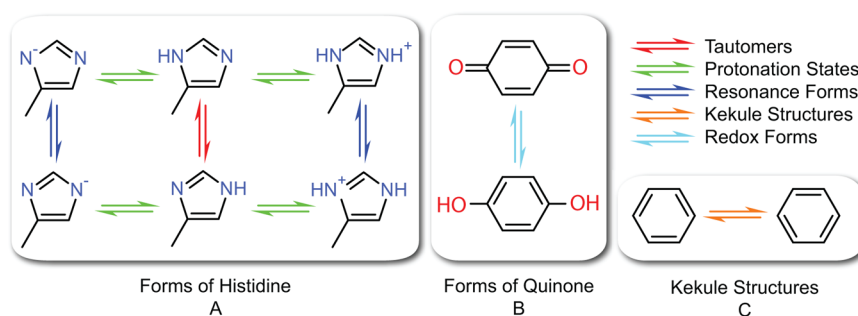


Figure 1. (A) Different valence bond structures of the imidazole ring of histidine including prototropic tautomers, protonation states, and resonance forms. (B) Two oxidation forms (quinone and hydroquinone) may in some context be considered as the same molecule. (C) Kekule structures are valence bond structures of aromatic rings with alternating single and double bonds.

The generation of unique identifiers, e.g., InChI,^{4,5} is a typical application scenario for canonicalization procedures.

Another quite opposite problem arises with regard to the general reliability of cheminformatics calculations. In many cases, it is necessary to consider multiple valence bond structures to sufficiently represent a molecule. The most prominent examples are certainly tautomers and protonation states, which will be summarized under the term protomers in the following. Since these correspond to actual physical entities, their respective ratios can have significant influence on a compound's observed physicochemical properties.^{6–11} The problem is, however, not exclusive to this scenario, as different resonance forms also play a role during the calculation of partial charges.¹² The respective tasks identified by Sayle¹ are (complete) enumeration (#3) and selection (#4). Both refer to the generation of valence bond structures, the difference being that selection (#4) restricts the results to a subset containing only relevant, e.g., energetically stable, solutions. Virtual screening techniques such as molecular docking are applications in which selection (#4) plays an important role. Relying on only one valence bond structure can lead to false-negative results as particular protomers may interact differently with target proteins. On the other hand, a large number of (possibly energetically unfavorable) alternatives can result in an increased false-positive rate and unnecessarily high runtimes. The general implications on structure-based and ligand-based screening methods have been investigated in several publications.^{13–15} The final task mentioned by Sayle¹ is prediction (#5), which extends selection by additionally ranking the relevant solutions by their respective energy.

The basic problem associated with the interconversion of valence bond structures is to transform groups of atoms according to specific rules with respect to bond orders and atomic properties (formal charges, bound hydrogens). As has been proposed by Sayle,¹ the methods developed for that purpose can be roughly divided into two categories: (1) Local approaches rely on pattern matching to identify relevant groups of atoms. These patterns are associated with rules describing the respective changes in the molecule. Pattern-based methods thus only use transformations that were anticipated in advance, thereby reducing the risk of generating unexpected and probably unwanted results. On the other hand, there is always the possibility of omitting relevant structures due to missing patterns. This can occur even if rules of a similar type are already included in the pattern library. Transformations covering long bond paths are a typical example for that problem. There are multiple publications describing local methodologies in the literature.^{13,16–18} (2) Global approaches

predefine substructures in a molecule, identify atoms with variable states within, and subsequently enumerate valid valence bond structures. This is usually done in a more generic manner than matching specific patterns, so that the results can easily contain completely artificial, i.e. chemically unreasonable, results. These either have to be omitted directly during or removed after the enumeration procedure. The omission of transformations in more complex structures, however, is generally not a problem. Global approaches have also been described in the literature^{19–21} and other sources.²² It must be noted that the previous differentiation between the two types of methods has been introduced mainly for classification purposes. Local approaches, for instance, often include a number of long-range patterns which, in combination with the underlying transformation engine, makes them suitable for the handling of the vast majority of molecules relevant in the field of drug design.

Here, we present the valence state combination model, a new concept for the description and classification of valence bond structures based on the NAOMI²³ framework. Using this model, we have developed, based on similar ideas as the ones presented by Sayle et al.,²² an extended and significantly improved method for the generation of valence bond structures which falls into the general category of global approaches. By application of a generic scoring scheme, this method combines the inherent consistency of the global strategy with the high reliability generally attained by local approaches. In contrast to previously published global methods, our approach consistently deals with all aspects relevant for the generation of protomers, including resonance forms and ionization states. Our method has been used to solve three common cheminformatics tasks, namely the generation of a canonical form (canonicalization), the generation of a preferential representation (normalization), and the generation of a set of reasonable protomers (generation). We have tested each application with respect to consistency using a general and comprehensible validation scheme. Furthermore, we have assessed the general suitability of our approach for common cheminformatics applications on the basis of these three operations. The criteria for the evaluation comprise runtime, the average number of generated structures, and the quality of the resulting protomers.

METHODOLOGY

Valence State Combination Model. Valence bond structures of molecules are generally represented as graphs in which nodes correspond to atoms and edges correspond to bonds. Each atom is associated with an element and a formal charge and each bond with a localized bond order (single,

double, or triple). In the NAOMI model,²³ this description is extended by an atom-based valence state descriptor. A valence state is a chemically valid combination of bond orders and formal charge for a particular element (see Figure 2). This

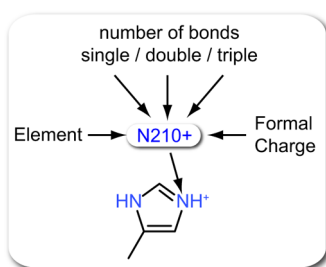


Figure 2. Example of a valence state descriptor for a nitrogen atom. The descriptor comprises the atom's element, bond order distribution, and formal charge.

additional descriptor is used to ensure the chemical validity of a molecule. A valence bond structure is valid if a valence state with the given bond orders and formal charge exists for each atom. Furthermore, valence states provide the means to systematically classify and generate different valence bond structures of molecules as explained below.

A set of valence states for all atoms of the molecule is called a valence state combination (VSC). A VSC is valid if a distribution of bond orders compatible with these valence states exists. Valid VSCs thus correspond to valence bond structures associated with a particular heavy atom skeleton. Note that bond orders are not part of the VSC representation; they are used for validation purposes only. Relations between valence bond structures can be determined by comparison of their corresponding VSCs (see Figure 3).

The description of these relations is based on atoms with different valence states, considering both their number and their types. Depending on the changed properties, substitutions of valence states for atoms are classified as protonation type, tautomer type, and resonance type as shown in Table 1. The involved states are called donors (higher number of single bonds) or acceptors (lower number of single bonds). The respective numbers of substitutions in VSCs are denoted as $\Delta_{\text{type}}(D \rightarrow A)$ and $\Delta_{\text{type}}(A \rightarrow D)$.

Table 2 lists the six basic relation types together with their conditions. Distinct valence bond structures with identical VSCs correspond to Kekule forms. They differ only in their respective bond order distribution. If all substitutions between two VSCs are of the protonation type, two cases need to be distinguished. When changing a donor to an acceptor or vice versa, the formal charge of the respective atom changes due to the addition or removal of hydrogen atoms. If the number of substitutions of donors and acceptors is not equal, the total charge of the molecule is altered, resulting in a different ionization state. Otherwise, the net charge of the molecule is identical, meaning that protons are merely occupying different locations. Tautomers and mesomers contain only changes of the tautomer-type and the resonance-type, respectively. Additionally, the number of donors and acceptor substitutions must be equal. Otherwise, the VSCs represent different redox forms of the molecule.

Substitution types can also occur in mixed constellations, and the resulting relations are best described as combinations of the just presented basic types. The 1-hydroxy-2-pyridone men-

tioned by Sayle¹ is an interesting example. The valence bond structures shown in Figure 4 can be best characterized as different resonance forms, a zwitterionic and a neutral one, with different proton positions.

The algorithms presented in the following chapters are based on the VSC representation of molecules. One of its major advantages is the fact that all of the potentially relevant molecule states can be consistently generated by considering different types of valence state substitutions. By explicitly handling all the different cases described in this section, a high degree of generality can be achieved.

Overview. The complete workflow for the generation of valence bond structures is shown in Figure 5. In the first step, the molecule is subdivided into multiple nonoverlapping substructures which are then treated independently. This partitioning reduces the computational costs for both the generation and the subsequent scoring of VSCs. A partition is considered valid if the independent enumeration of VSCs of each part and a subsequent combination of these lead to the same VSCs as if the enumeration would have been performed on the whole molecule. A partition is optimal if it is valid and has the smallest possible substructures. In the following sections, two partitioning schemes (generic and heuristic) are presented. Both are applied for the solution of different cheminformatics tasks described in later sections.

After partitioning, the atoms of each substructure are checked for alternative valence state assignments. Which valence states are included strongly depends on the context and will be explained in more detail later. As well as partitioning, valence state selection has a strong influence on the computational costs of the subsequent steps. The more alternatives are selected, the more VSCs must be generated and potentially scored. An optimal selection scheme thus only selects valence states for atoms that actually need to be modified. Again, two selection schemes (generic and heuristic) for different applications will be presented.

In the next step, VSCs are generated for each substructure using the alternative states selected in the previous step. Each of these VSCs is checked for validity by attempting to calculate a bond order distribution. VSCs for which this is not possible are invalid and therefore rejected. During the calculation, additional boundary conditions, e.g., the oxidation state of the initial molecule, are preserved.

The resulting VSCs are all chemically valid but may still contain undesired valence bond structures. These include unstable tautomers, unlikely protonation states, unreasonable resonance forms, or unusual representations of functional groups. In order to identify and eventually remove these VSCs, a pattern-based scoring scheme is applied. The resulting score expresses how well a particular substructure of the molecule is represented by the respective VSC. It must be stressed that the scoring scheme has not been designed to accurately predict the ratios between different molecular species. Its two main purposes are the elimination of completely artificial representations, i.e., energetically inaccessible states, and the coarse categorization of the remaining VSCs into stability classes. After eliminating all undesired VSCs, the final valence bond structures are completely enumerated by combining the VSC of the different substructures.

Partitioning of Molecules. The partitioning algorithm is based on the exclusion of atoms and bonds from the molecular graph and the subsequent identification of the remaining connected components. These will be referred to as Multi State

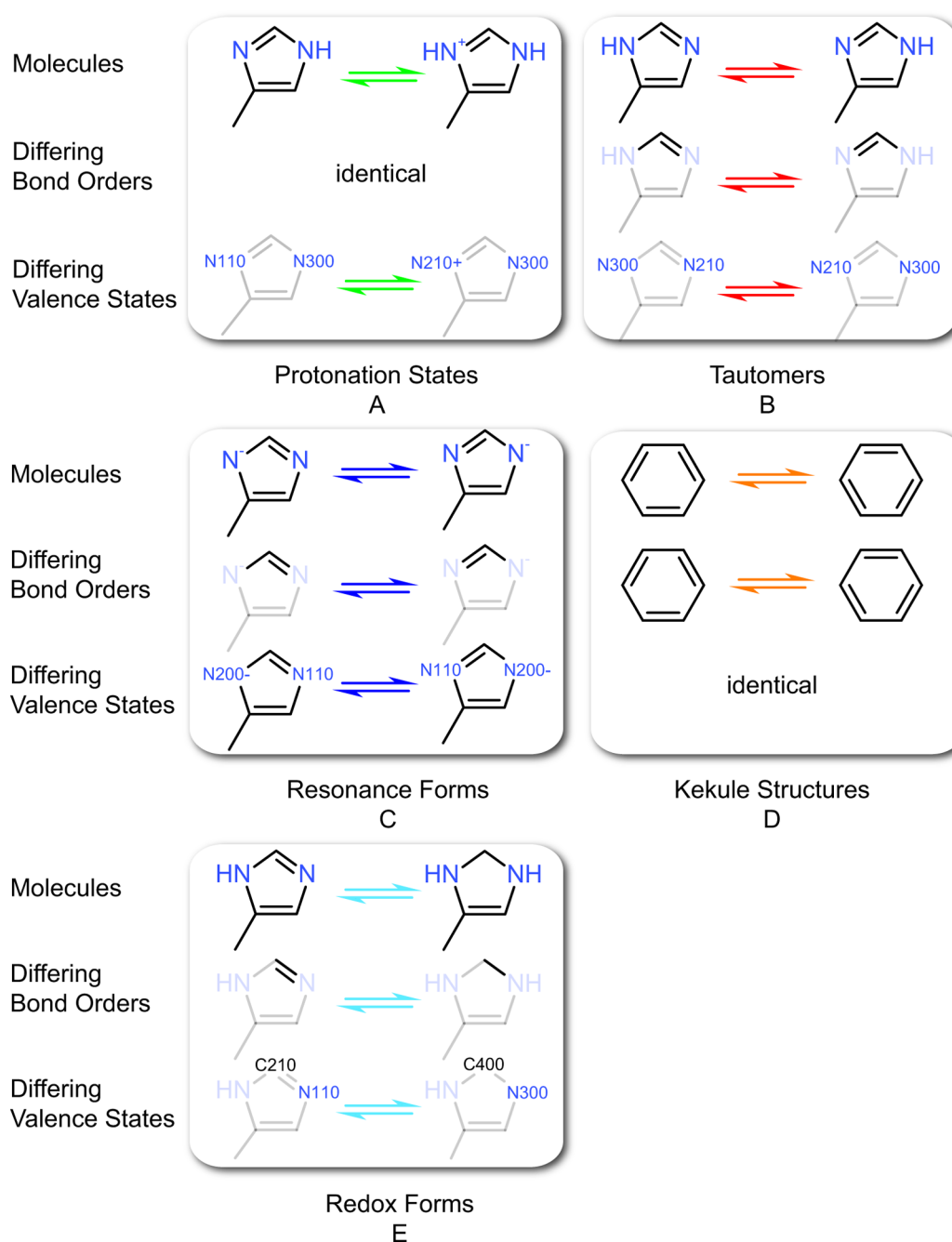


Figure 3. Differences between protonation states (A), tautomers (B), resonance forms (C), Kekule structures (D), and redox forms (E).

Table 1. Substitution Types for Valence States Including the Affected Properties^a

type	double bonds	# bonds	charge	examples	
				donor	acceptor
protonation	0	±	±	O200	O100-
resonance	±	0	±	O100-	O010
tautomer	±	±	0	O200	O010

^aChanged properties are marked with a ± and unchanged properties with 0. The pairs of valence states on the right side of the table represent common substitutions for oxygen atoms.

Table 2. Relations between Valence Bond Structures on the Basis of Valence State Substitution

relation	substitution type	condition
kekule	none	
ionization	protonation	$\Delta(D \rightarrow A) \neq \Delta(A \rightarrow D)$
protonation	protonation	$\Delta(D \rightarrow A) = \Delta(A \rightarrow D)$
mesomer	resonance	$\Delta(D \rightarrow A) = \Delta(A \rightarrow D)$
tautomer	tautomer	$\Delta(D \rightarrow A) = \Delta(A \rightarrow D)$
redox	resonance	$\Delta(D \rightarrow A) \neq \Delta(A \rightarrow D)$
	tautomer	$\Delta(D \rightarrow A) \neq \Delta(A \rightarrow D)$

Partitions (MSP) in the following discussion. The **generic** partitioning scheme only involves the exclusion of sp³-hybridized carbon atoms (corresponds to valence state

C400). There are only two particular cases in which atoms with valence state C400 are included in MSPs: first, if the atom is bound to an atom with valence state C210, which in turn has

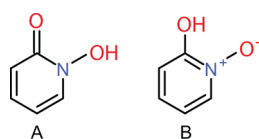


Figure 4. Example for the combination of valence state substitutions. The relation between the pyridone form (A) and the pyridine form (B) cannot be described by one of the basic types from Table 2.

at least one neighbor with the element nitrogen, oxygen or sulfur, and second, if the atom is part of a ring and is the only atom with valence state C400 in this ring. Bonds are excluded if one of the connected atoms is excluded.

The MSPs resulting from the generic partitioning scheme are usually large, and it is often possible to further reduce their size. This is achieved by removing bonds within the MSPs with the goal to effectively split them into smaller substructures. The exclusion of a bond is only valid if its bond order in the current structure is identical in all relevant VSCs. Since the final bond orders are not known at this point, the decision that a bond will keep its current type must be in accordance with the subsequent scoring procedure. This means that VSCs with a different bond order would be rejected in the following steps in any case.

The heuristic partitioning scheme builds on the results from the generic scheme and uses a set of rules to identify additional bonds for exclusion. These rules are based on the classification of each MSP into conjugated rings, conjugated chains, and functional groups. Rings are considered conjugated if all of their atoms are part of the respective MSP. Conjugated chains consist only of carbon atoms which have a multiple bond and are bound only to other carbon atoms. The remaining connected components represent functional groups. In a first step, bonds connecting functional groups with conjugated rings or conjugated chains are investigated. A bond is excluded if it is a single bond and the atom from the functional groups does not fulfill one of the following two criteria: (1) It has a valence state of type N300. (2) It has a valence state of type O200 or S200 and only one non-hydrogen bond. In these cases, a change in bond order is not unlikely, as is shown for two examples in Figure 6.

Since conjugated chains consist of only carbon atoms, they are merely bridges between the other two types of

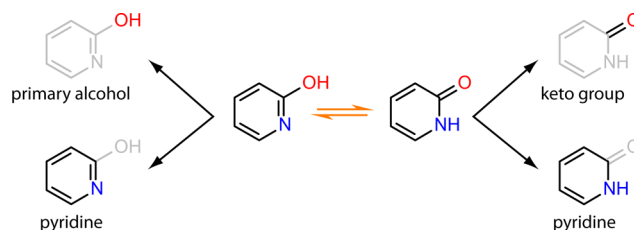


Figure 6. Two examples for which functional groups and conjugated rings have to be treated as a union to avoid missing VSCs.

substructures. Therefore, if a conjugated chain has only one bond to another structure (ring or functional group), this bond can be safely excluded. This is also done if the chain has multiple bonds which were previously excluded by the functional group rule. The complete partitioning of the NAD⁺ molecule is shown as an example in Figure 7.

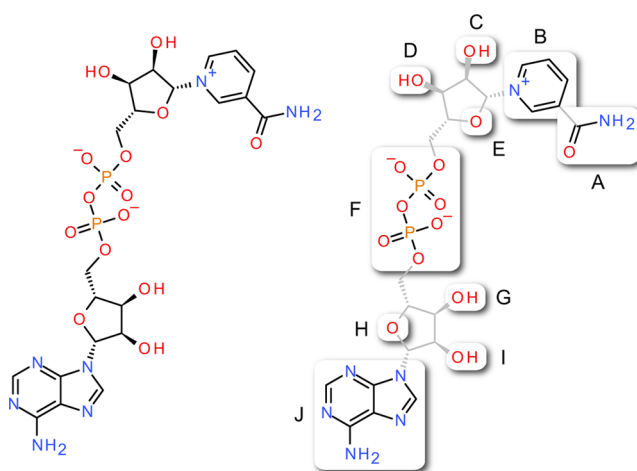


Figure 7. Partitioning of NAD⁺ into functional groups and conjugated rings. The amide group and the pyridine ring have been separated, whereas the amino group remains connected to the purine.

Selection of Valence States. The selection of valence states is based on the substitution types introduced above (see Table 1). Each substitution corresponds to a pair of valence states which are known in advance and can be retrieved starting

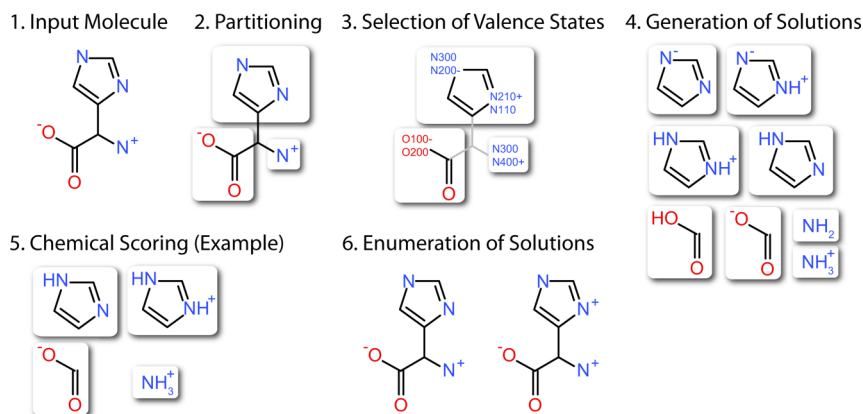


Figure 5. Overview of the generation of protonation states for an input molecule (1). In a first step, the molecule is partitioned into substructures (2) which are handled separately. In the next step, alternative valence states are selected (3). Afterward, valid VSCs are generated (4). These are scored (5), and only the best solutions for each zone are retained. The final list of valence bond structures results from the combination of all remaining VSCs.

from any valence state. A list of alternatives for an atom can thus be easily obtained by consecutively and uniquely adding the respective members of the pairs for each of the relevant substitution types. In order to select an alternative assignment, the compatibility with the atom's topology must be ensured. This means that the number of bonds of the valence state must be larger than or equal to the atom's number of non-hydrogen bonds. Otherwise, the assignment would correspond to the removal of non-hydrogen bonds. Although this may be interesting with respect to transformations such as ring-chain tautomerism, it will not be further considered here.

For the sake of generality, the **generic** selection procedure includes all possible valence states for each atom in a MSP. This usually results in potentially many more alternatives than are actually needed. The **heuristic** selection scheme aims at reducing this number by explicitly excluding valence states for particular atoms. The problem at this point is similar to the one discussed in the previous section. The final VSCs are not yet known, and the decisions must be in accordance with the subsequent scoring procedure in order to avoid missing VSCs.

The exclusion of particular valence states in the **heuristic** selection scheme is based solely on an analysis of the atom's environment. For atoms in functional groups, this includes their direct neighbors from the same functional group. These are transformed into a SMILES-like identifier which reflects the valence bond structure of the input molecule. This identifier is looked up in a list of predefined structures. If the identifier is present, information concerning the exclusion of particular substitution types is retrieved. In this way, groups that already have a preferred representation in the initial valence bond structure need not be modified. The information provided from the patterns is described in Figure 8.

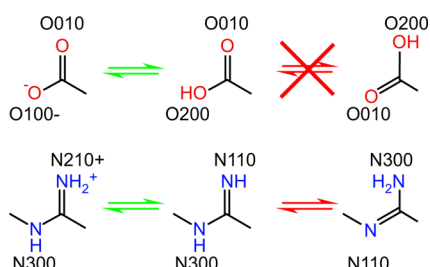


Figure 8. Selection of valence states for carboxylic acid and amidine groups. Due to the group's symmetry, different tautomers of carboxylic acids are not considered.

For the generation of tautomers, carboxylic acids are irrelevant. Due to the symmetry of the group the transfer of the hydrogen from one oxygen to the other would only result in a different rotamer. In this case both oxygen atoms are excluded from tautomer substitution. With respect to protonation, both the charged and the neutral form need to be included. This means that both oxygens are not excluded from protonation substitution. The same procedure is applied to atoms in conjugated rings with the ring constituting the atom's environment. If the identifier is not included, the **generic** scheme is used to identify alternative states for the atom.

Generation of Valid VSCs. Prior to the generation of VSCs, each MSP is analyzed to ensure that the generation of additional states is at all possible. MSPs can be ignored if no atom with alternative valence states could be found. For

tautomers and mesomers, i.e. if new bond order distributions are to be generated, MSPs can also be omitted if only either donors or acceptors are present. In this case, no substitution of valence states is possible (see Figure 9 for examples). Changing the number of donors and acceptors corresponds to changing the oxidation state of the molecule, which is not desired in most contexts.

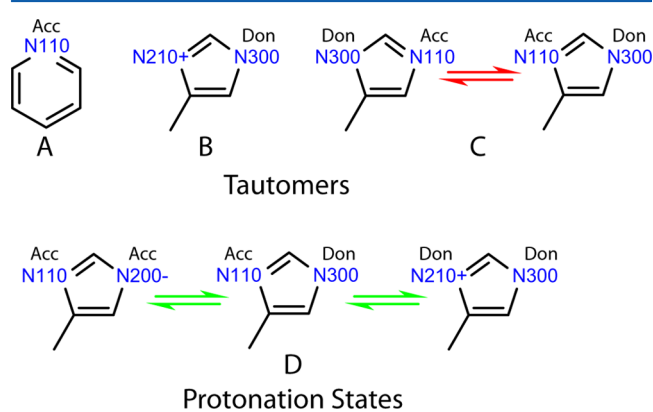


Figure 9. Criteria for the generation of additional states. The generation of tautomers requires at least one tautomer acceptor and one tautomer donor in a zone. Pyridine (A) has only a single tautomer acceptor, and the imidazolium ion (B) only has one tautomer donor. No tautomers can be generated in such cases. Imidazole (C) contains a tautomer acceptor as well as a tautomer donor and can tautomerize. Protonation states (D) can also be generated if a molecule only contains either protonation-type donors or acceptors.

The algorithm for the generation of valid VSCs is based on a backtracking procedure with pruning. The atoms of the MSP are processed in a specific order which is established prior to the actual assignment procedure. The algorithm starts with terminal atoms, i.e. atoms with only one bond in the MSP, followed by internal atoms with at least one terminal neighbor. The remaining atoms are processed last. The order of the atoms inside the three classes is arbitrary and does not affect the result. As a combinatorial problem, the procedure can be represented by a tree, where each node corresponds to the assignment of a valence state to an atom. Inner nodes thus represent partial VSCs while the tree's leaves correspond to complete VSCs. For each node, the chemical validity of the corresponding VSC is verified. In most cases, this can be performed without actually generating bond orders for the bonds of the MSP. The checks are based on the compatibility between valence states of different atoms with respect to the expected bond types as well as their oxidation states: (1) For atoms with only one bond in the MSP, the assignment of a valence state is equivalent to the assignment of a bond order to the corresponding bond. The compatibility with the atom's neighbors can be easily checked by ensuring that the count of this particular bond type is not exceeded. This check is always performed when an atom with terminal neighbors is encountered. (2) When reaching a leaf, the valence states with an uneven number of multiple bonds are counted. If this number is uneven, no valid bond order distribution exists, and the VSC can be further ignored. (3) The number of donors in the initial valence bond structures is counted in order to retain the molecule's oxidation state. VSCs differing in the number of donors compared to the initial valence bond structures can be discarded. Note that since information about being a donor or

acceptor is also stored in the valence states, VSCs not fulfilling this boundary condition can be easily identified. (4) Eventually, for each VSC passing all previous checks, a recursive bond localization routine is used which assigns bond orders to all bonds in the MSP. If this routine is successful, the solution represents a valid valence bond structure and is stored.

Scoring of VSCs. Scores for each VSC are calculated under consideration of the bond order distribution generated in the previous step. The scoring procedure is mainly based on the recognition of predefined structural fragments contained within particular substructures, i.e., conjugated rings and functional groups, of the molecule. The final score of the VSC (S_{VSC}) is calculated as the sum of the individual scores obtained for each of these substructures (see eq 1). Please note that due to changes in bond orders and valence states, the scores have to be recalculated for each VSC.

$$S_{\text{VSC}} = \sum S_{\text{ring}} + \sum S_{\text{group}} \quad (1)$$

$$S_{\text{ring}} = \sum \text{cycle} + \sum S_{\text{sub}} \quad (2)$$

$$S_{\text{group}} = \sum S_{\text{subgroup}} \quad (3)$$

The structural fragments in the substructures are identified using canonical SMILES-like identifiers. These are generated on the basis of the bond types and valence states of the respective VSC. The predefined data are stored in multiple databases which can be queried with the identifiers in order to retrieve the score associated with a fragment.

In case of conjugated rings, the score S_{ring} comprises two types of contributions, one from the ring itself, S_{cycle} , and one from its substituents, S_{sub} (see eq 2). The reference point for S_{cycle} is the isolated aromatic system without exocyclic double bonds, e.g., pyrrole for a five-membered ring with one nitrogen atom. In case there are multiple structures fulfilling this requirement, e.g., the 1H and 2H tautomers of 1,2,3-triazole, one is arbitrarily selected. The score of the reference system is set to an arbitrary value of 100. If a ring with an identical heavy atom connectivity does contain a structural deviation from the reference, e.g., an sp^3 hybridized carbon atom, the associated fragment has an individual score. This can be higher or lower depending on the stability assigned to this particular arrangement. The substructures representing ring substituents comprise the ring atom, the exocyclic atom, and the exocyclic atom's direct neighbors. The associated scores have fixed values and are independent from the concrete ring system they are connected to. Again, one particular representation of the substituent, the one with an exocyclic single bond and without charges, receives an arbitrary reference score of 100. Functional groups are first treated as a whole; i.e., an identifier for the complete group is generated. If the pattern was present, the associated score is directly set as the score of the substructure. Otherwise, the group is partitioned into smaller pieces which serve as starting points for further queries. In this case, the score for the group is composed of the scores of the smaller fragments (see eq 3). The reference system for a subgroup is preferably neutral and corresponds to the most stable tautomeric form where possible.

If no predefined data are available in any of the three cases, a generic score is calculated according to eq 4:

$$S_{\text{generic}} = \max(0, 80 - \sum P) \quad (4)$$

This is done by subtracting various penalties (P) which are summarized together with the respective conditions in Table 3.

Table 3. Classification and Conditions for the Penalties used during the Calculation of Generic Scores

substructure	type	penalty	condition
ring	aromaticity	20	nonaromatic ring (Hueckel's rule)
ring	charge	20	single charge in ring
ring	charge	80	multiple charges in ring and substituents
ring	stability	80	three consecutive donors ^a in the ring
substituent	bond order	20	substituent has exocyclic double bond
substituent	charge	20	single charge in substituent
substituent	charge	80	multiple charges in substituent
group	charge	80	multiple positive charges in group

^aDonors are atoms with the following valence states: O200, N300, S200.

Since S_{generic} is used only as a fallback, the respective maximal score is deliberately set lower than that of the reference system. If the sum of the penalties (P) exceeds 80, the score of the substructure is set to zero.

The relative differences between the scores of rings, substituents, and functional groups have been derived from multiple pairs of tautomers and ionization states for which the major form was known from either experiments or theoretical calculations.²⁴ The databases currently contain 252 entries in total (113 in cycles, 121 in subgroups, 18 in substituents). Examples for ring and functional groups patterns are shown in Figure 10.

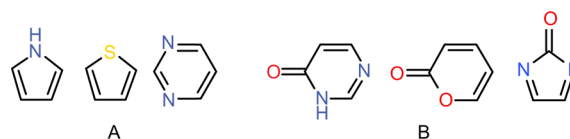


Figure 10. Examples for ring and functional groups patterns. (A) A reference score of 100 is assigned to isolated aromatic rings without exocyclic double bonds. (B) The score for rings with exocyclic double bonds comprises one contribution from the ring and another from the carbonyl substituent.

VSC MODEL APPLICATIONS

In the following applications, we will consider resonance forms, prototropic tautomers, and protonation states as instances of the same molecule, whereas oxidation forms are interpreted as distinct chemical species. The method is, however, not restricted to this assumption in general and can be easily modified so that different types of valence bond structures are perceived as identical.

Canonicalization. The generation of a canonical representation is the first workflow in which our method is applied. Canonical representations are mainly used to determine whether two valence bond structures represent the same molecule. In this context, it does not matter if the result corresponds to the most stable form or even a chemically reasonable one.

The workflow starts with the partitioning of the molecule into MSPs and the selection of alternative valence states as described above. The atoms of each MSP are then sorted in a

canonical way using a variant of the Morgan extended sums algorithm.²⁵ The backtracking algorithm in the generation step processes the atoms in this exact order until the first valid VSC has been found. This VSC serves as the canonical form of the respective substructure. Since no additional scoring is needed, the canonical VSCs of each substructure can be directly combined to yield the canonical representation for the complete molecule.

For the canonicalization to work correctly, the results must be identical for each possible valence bond structure of the molecule provided as input. This can only be achieved if the substructures generated in the partitioning step and the lists of valence states identified in the selection step are both identical in each case. The heuristic algorithms for partitioning and selection are therefore inappropriate, and the generic variants are applied. Since only a single valid VSC must be generated in the end, the size of the substructures and the number of alternative states are of less importance for the resulting compute time. Nevertheless, in order to further accelerate the process, all charged valence states are transformed into their neutral states where possible (considering the number of hydrogens) using the protonation-type substitution. Consequently, only tautomer-type and resonance-type substitutions need to be considered in the next steps.

The canonicalization procedure applied to the atoms of each substructure differs only in one aspect from the CANON algorithm used for the generation of USMILES.²⁶ In the CANON algorithm, the atomic invariants correspond to the atom's valence state in combination with the number of attached hydrogens. This means that the initial ranks of atoms can normally be deduced by comparing valence states and hydrogens. In case of a yet unknown valence bond structure, the final valence state of an atom is, however, not defined. Instead, a list of valence states is used to describe the topology of each atom and provide the initial ranks. Furthermore, the number of non-hydrogen bonds serves as a replacement for the number of hydrogens.

Normalization. The aim of normalization is the generation of a canonical valence bond structure which additionally adheres to common conventions for the representation of molecules. This task seems, at least at first glance, quite similar to the previously described canonicalization. The main difference results from the necessity of a scoring step in order to determine the best suited choice for the molecule. This implies that multiple VSCs have to be generated and compared with each other. Here, we have chosen a neutral form as normalized representation, meaning that all atoms are neutralized when possible (considering bound hydrogens). The only exception to this rule is functional groups which are represented in a zwitterionic form by convention, e.g., nitro groups and n-oxides. The method is, however, not restricted to this preference and can be easily modified so that, for instance, the preferred ionization state is generated.

Again, the workflow starts with the partitioning of the molecule into substructures and the selection of alternative valence states. Due to the enumeration of VSCs in the later steps, the size of the substructures and the number of states are relevant factors. Therefore, the heuristic strategies for both partitioning and state selection are used. In contrast to canonicalization, the initial substructures and alternative valence states do not have to be identical for each starting structure. The additional scoring step ensures that the results are consistent.

In the next two steps, valid VSCs are generated and scores are assigned as explained in the sections above. For each substructure, only those solutions with the highest score are retained. If there is only one VSC left for a substructure, it can be directly assigned, and no further steps are necessary. Otherwise, a canonical solution has to be picked from the VSCs with the highest score. This is done using the canonicalization method described in the previous section. However, since this method only works correctly in case of identical MSPs and lists of valence states, a preprocessing step is required. The respective MSP is repartitioned by exclusion of bonds having the same bond type in all VSCs. Additionally, all valence states which could not be found in one of the remaining VSCs are removed from the lists of alternatives. This eventually creates the necessary conditions for the canonicalization procedure.

Generation. The last application of our method is the generation of a set of reasonable tautomers and protonation states of a molecule. The resulting molecules can be used as input for methods that rely on the positions of hydrogen atoms such as docking. They can also serve as a starting point for the determination of the energetically most stable form of a molecule under consideration of the molecule's local environment, e.g., bound to a protein. The inclusion of multiple resonance structures, although possible with our method, is not considered useful in this context.

The initial steps of the workflow are identical to those described for normalization. But instead of canonically selecting one of the remaining VSCs of each zone, the combinations are enumerated in order to generate a set of molecules. One major difference from the previously presented approaches is the possibility of generating duplicates due to molecular symmetry. This is avoided by removing VSCs from each zone that would lead to identical valence bond structures in the resulting molecules. For this purpose, automorphism classes for atoms are calculated using the Morgan algorithm, which is also used for the canonicalization. In combination with the respective valence state of an atom in a VSC, these classes can be used to generate a string representation of each VSC in a zone, which are used to identify and remove duplicates.

For molecules containing more than one ionizable group, it is usually not desirable to enumerate all combinations of VSCs from the respective zones. To avoid chemically unreasonable species with a high number of charges, the maximum number of charges in the complete molecule is restricted by three simple rules: (1) The number of charged groups must be smaller than four, (2) the number of pairs of oppositely charged groups is smaller than two, and (3) the maximum number of positive charges in a ring system is restricted to one.

■ RESULTS AND DISCUSSION

The three applications presented in the previous sections are the basis for the evaluation of our method in terms of consistency, quality, and performance. Throughout these studies, the following commonly used public data sets served as input: (1) ZINC clean leads^{27,28} (ZINC-CL), (2) LigandExpo component dictionary^{29,30} (LEXPO-CD), (3) Drugbank^{31,32} (DRUGBANK), and (4) ChEMBL.^{33,34} All calculations and runtime measurements were performed on a PC with an Intel Core i5-3570 CPU (3.40 GHz) and 8 GB of main memory.

Consistency. Independence from the initial valence bond structure of a molecule is a fundamental requirement of the presented method and has been thoroughly investigated for

Table 4. Runtimes for the Three Workflows with Different Data Sets

	ZINC-CL	LEXPO-CD	DRUGBANK	ChEMBL
# total molecules	5735035	17310	6583	1318187
runtime canonicalization [ms/cmpd]	0.28	0.41	0.45	0.71
runtime normalization [ms/cmpd]	0.31	0.50	0.6	0.73
runtime generation [ms/cmpd]	0.45	0.75	0.75	1.37

each of the three applications. Consistency can be verified by a simple and straightforward procedure. The starting point is a set containing different representations of the same molecule, e.g., as different molecule entries in a file. After applying the respective workflow to each representation, the resulting molecules are converted to USMILES for comparison. If the method is consistent, all resulting USMILES are identical. In case of enumeration, lists of USMILES must be compared.

The best way to ensure consistency would be to test all possible valence bond structures of the molecule with the procedure described above. This is, however, not feasible in many cases due to the prohibitively large number of resulting molecular states. We therefore decided to reduce the set by exclusion of protonation and ionization states (see Table 2 for our definition), since the main complexity of the task results from valence bond structures with different bond order distributions.

The input structures needed for the assessment of our method's consistency were generated using a workflow corresponding to the one described for the canonicalization of molecules. But instead of selecting a canonical form, we generated valid VSCs without any scoring step and enumerated all possible combinations. Identical results could be achieved for all three workflows, canonicalization, normalization, and generation, with all four data sets mentioned above.

Runtimes. Table 4 lists the runtimes for the three workflows with the above-mentioned data sets. The results for canonicalization and normalization are comparable in both cases, whereas the time needed for the generation of a set of states is higher. This is not surprising since the workflow involves the enumeration of multiple molecule states and the built-in elimination of duplicates based on automorphism classes. In all cases, an average runtime lower than 1.5 ms per molecule is measured, thus showing that our method is suitable for processing large data sets. The similarity of results for canonicalization and normalization are most probably a consequence of the normalization procedures used during the curation of the used data sets. As has been explained above, the runtimes for normalization are highly dependent on the input form of the molecule, and the process is accelerated by reasonable initial representations.

Normalization. The main purpose of normalization is to transform different input forms of the same molecule into an identical and at the same time chemically reasonable representation. We have already shown that our normalization workflow is consistent for the four data sets used in this study. Here, we will focus on the second aspect. We believe that the best way to investigate if results are chemically reasonable is to compare the resulting valence bond structure with those found in frequently used and curated public data sets.

The procedure applied for this purpose is again based on the comparison of USMILES. Directly using the input molecule and the normalized version is, unfortunately, not suitable in many cases. As has been explained above, a canonicalization step at the end of the workflow is used to arbitrarily select one

of multiple equally acceptable solutions. This makes the comparison to a reference structure, which has most likely been normalized by a different procedure, pointless. We therefore decided to enumerate all combinations of VSCs with the highest score and to check if the input structure is contained within the obtained set (best). A negative result does, however, not necessarily mean that our method generated an unreasonable result. The representation in the data set could simply correspond to a VSC which received a lower score based on our scoring scheme. For that reason, we additionally enumerated all VSCs with a score of at least 75% of the best score and also searched in this extended set (extended). The results of this validation procedure are summarized in Table 5.

Table 5. Classification of the Input Structures from Three Data Sets into Mutually Exclusive Categories for the Generation of Tautomers

	LEXPO-CD	DRUGBANK	ChEMBL
# total molecules	17310	6583	1318187
# molecules (best)	16837	6431	1252408
# molecules (extended)	364	118	52491
# molecules (not found)	135	48	11433

The differences encountered during the process can be subdivided into two classes. First, there are input structures which are not found in the best set, but in the extended set. These correspond in many cases to keto and enol tautomers of aromatic heterocycles, which are ranked differently by our scoring method (see 138 in Figure 11). Second, there are input structures which are not present in either of both sets. After visual inspection, we think that the results generated by our method are in general at least equally acceptable and in some cases even better than the representations found in the data set. The latter especially applies to charged structures for which a reasonable neutral form can be formulated (see 3MC in Figure 11). The normalized molecules generated by our method are provided as Supporting Information for all entries of LEXPO-CD and DRUGBANK which were not included in either of the two sets.

Finding the input structure in a set of equally scored alternatives is, however, only one aspect of the method's performance. Additionally, one has to make sure that the success is not simply based on the enumeration of an unreasonably large number of representations. For that reason, the sizes of the respective sets are also an important performance indicator and are shown in Table 6.

The average number of generated states is considerably lower than the result of an exhaustive enumeration. Only for a small percentage of molecules (less than 0.5%) does the number of equivalent structures actually exceed a size of five. This is in all cases caused by the combination of states from independent zones, e.g., molecules having multiple imidazole rings.

Generation. The aim of our generation workflow is to generate a set of chemically reasonable protomers of a molecule

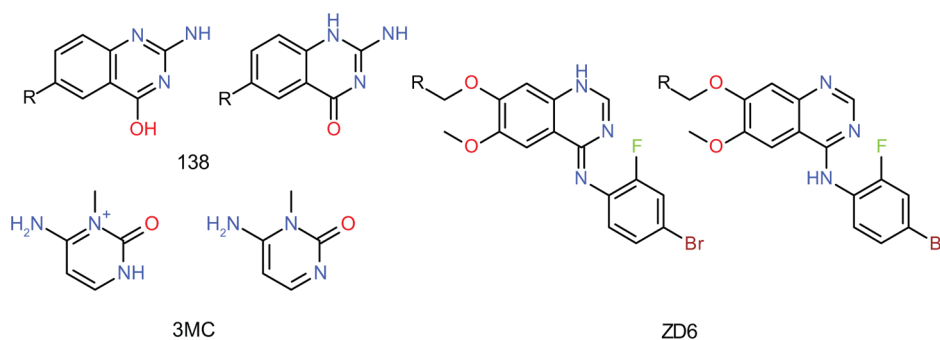


Figure 11. Examples of differences between the normalized forms generated by our method (right side) and those found in the Ligand Expo data set (left side).

Table 6. Number of Molecules with More than One and More than Five Tautomers in the Best Set^a

	ZINC-CL	LEXPO-CD	DRUGBANK	ChEMBL
# total molecules	5735035	17310	6583	1318187
# tautomers >1	207207	1483	520	93430
# tautomers >5	671	5	5	4699
# average	2.27	2.21	2.37	8.3

^aThe provided average refers only to cases with more than one tautomer.

for typical cheminformatics applications, e.g., docking calculations. Considering this context, the resulting set should only contain states which are realistically expected to be stable in a protein–ligand complex. In order to assess the quality of our results, we used ZINC-CL as a reference set since it was generated for the exact same scenario. The procedure is identical to the one described for the evaluation of the normalization workflow. The input structure is searched in two sets, one containing the states with the highest score (best) and one containing states with a score of at least 75% of the highest score (extended). The results of the procedure are summarized in Table 7.

Table 7. Classification of the Input Structures from the ZINC-CL Data Set into Mutually Exclusive Categories for the Generation of Protomers

	ZINC-CL
# total molecules	5735035
# molecules (best)	4764463
# molecules (not best)	914921
# molecules (not found)	55651

As has already been discussed above, one important parameter for the evaluation of the method's performance certainly is the number of generated states. The results for all four data sets are summarized in Table 8.

CONCLUSION

The simple fact that the same molecule can be represented by different valence bond structures constitutes a complex challenge for cheminformatics applications. It complicates the determination of molecular identity and makes the results of cheminformatics calculations prone to inconsistencies. Furthermore, it imposes the task of selecting the best suited structure or structures for the respective context of application. The identification, description, and consistent handling of these

Table 8. Number of Molecules with More than One and More than Five Protomers in the Best Set^a

	ZINC-CL	LEXPO-CD	DRUGBANK	ChEMBL
# total molecules	5735035	17310	6583	1318187
# protomers >1	1007976	2221	770	159663
# protomers >5	9240	183	78	13231
# average	2.54	3.14	3.20	4.40

^aThe provided average refers only to cases with more than one protomer.

different molecular representations is thus a fundamental requirement in the field of cheminformatics.

To cope with these problems, we have introduced a formalism which describes different valence bond structures of a molecule on the basis of the recently published NAOMI model. Using this description, we developed a general method for their fast and consistent enumeration and presented three exemplary applications. In our validation, we have shown that the devised methodology can be successfully applied to relevant tasks in cheminformatics in a consistent manner. We have also demonstrated the low runtime of our approach which makes it suitable for processing large data sets.

ASSOCIATED CONTENT

Supporting Information

The Supporting Information includes the normalized structures of all entries from the Ligand Expo Component Dictionary and Drugbank whose input form was not included in the results of our method. These are provided as separate SMILES files. This material is available free of charge via the Internet at <http://pubs.acs.org/>.

AUTHOR INFORMATION

Corresponding Author

*E-mail: rarey@zbh.uni-hamburg.de.

Author Contributions

[†]Equal contribution.

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

The authors thank one of the reviewers for bringing the existence of the freely available source code of the method developed by Sayle and Delaney³⁵ to their attention.

REFERENCES

- (1) Sayle, R. So you think you understand tautomerism? *J. Comput.-Aided Mol. Des.* **2010**, *24*, 485–496.
- (2) Warr, W. A. Tautomerism in chemical information management systems. *J. Comput.-Aided Mol. Des.* **2010**, *24*, 497–520.
- (3) Sitzmann, M.; Ihlenfeldt, W.-D.; Nicklaus, M. Tautomerism in large databases. *J. Comput.-Aided Mol. Des.* **2010**, *24*, 521–551.
- (4) Heller, S.; McNaught, A.; Stein, S.; Tchekhovskoi, D.; Pletnev, I. InChI - the worldwide chemical structure identifier standard. *J. Cheminform.* **2013**, *5*, 7.
- (5) InChI version 1, software version 1.04 (2011)-Technical Manual. http://www.inchi-trust.org/fileadmin/user_upload/software/inchi-v1.04/INChI_TechMan.pdf (last accessed Dec 06, 2013).
- (6) Milletti, F.; Storch, L.; Sforza, G.; Cruciani, G. New and Original pK_a Prediction Method Using Grid Molecular Interaction Fields. *J. Chem. Inf. Model.* **2007**, *47*, 2172–2181.
- (7) Shelley, J.; Chollet, A.; Frye, L.; Greenwood, J.; Timlin, M.; Uchimaya, M. Epik: a software program for pK_a prediction and protonation state generation for drug-like molecules. *J. Comput.-Aided Mol. Des.* **2007**, *21*, 681–691.
- (8) Martin, Y. Let's not forget tautomers. *J. Comput.-Aided Mol. Des.* **2009**, *23*, 693–704.
- (9) Clark, T. Tautomers and reference 3D-structures: the orphans of in silico drug design. *J. Comput.-Aided Mol. Des.* **2010**, *24*, 605–611.
- (10) Cramer, R. Tautomers and topomers: challenging the uncertainties of direct physicochemical modeling. *J. Comput.-Aided Mol. Des.* **2010**, *24*, 617–620.
- (11) Greenwood, J.; Calkins, D.; Sullivan, A.; Shelley, J. Towards the comprehensive, rapid, and accurate prediction of the favorable tautomeric states of drug-like molecules in aqueous solution. *J. Comput.-Aided Mol. Des.* **2010**, *24*, 591–604.
- (12) Gilson, M.; Gilson, H.; Potter, M. Fast assignment of accurate partial atomic charges: an electronegativity equalization method that accounts for alternate resonance forms. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 1982–1997.
- (13) Oellien, F.; Cramer, J.; Beyer, C.; Ihlenfeldt, W.-D.; Selzer, P. M. The Impact of Tautomer Forms on Pharmacophore-Based Virtual Screening. *J. Chem. Inf. Model.* **2006**, *46*, 2342–2354.
- (14) ten Brink, T.; Exner, T. Influence of Protonation, Tautomeric, and Stereoisomeric States on Protein-Ligand Docking Results. *J. Chem. Inf. Model.* **2009**, *49*, 1535–1546.
- (15) Kalliokoski, T.; Salo, H.; Lahtela-Kakkonen, M.; Poso, A. The effect of ligand-based tautomer and protomer prediction on structure-based virtual screening. *J. Chem. Inf. Model.* **2009**, *49*, 2742–2748.
- (16) Kenny, P.; Sadowski, J. In *Cheminformatics in Drug Discovery*; Oprea, T., Ed.; Wiley-VCH Verlag GmbH & Co. KGaA: Weinheim, Germany, 2005; pp 271–285.
- (17) Milletti, F.; Storch, L.; Sforza, G.; Cross, S.; Cruciani, G. Tautomer enumeration and stability prediction for virtual screening on large chemical databases. *J. Chem. Inf. Model.* **2009**, *49*, 68–75.
- (18) Kochev, N. T.; Paskaleva, V. H.; Jeliaskova, N. Ambit-Tautomer: An Open Source Tool for Tautomer Generation. *Mol. Inf.* **2013**, *32*, 481–504.
- (19) Haranczyk, M.; Gutowski, M. Quantum Mechanical Energy-Based Screening of Combinatorially Generated Library of Tautomers. TauTGen: A Tautomer Generator Program. *J. Chem. Inf. Model.* **2007**, *47*, 686–694.
- (20) Thalheim, T.; Vollmer, A.; Ebert, R.-U.; Kühne, R.; Schüürmann, G. Tautomer Identification and Tautomer Structure Generation Based on the InChI Code. *J. Chem. Inf. Model.* **2010**, *50*, 1223–1232.
- (21) Will, T.; Hutter, M. C.; Jauch, J.; Helms, V. Batch tautomer generation with MolTPC. *J. Comput. Chem.* **2013**, *34*, 2485–2492.
- (22) Sayle, R.; Delany, J. In *Innovative Computational Applications: the Interface of Library Design, Bioinformatics, Structure Based Drug Design and Virtual Screening*; IIRG publishers: San Francisco, CA, 1999. http://www.daylight.com/meetings/emug99/Delany/taut_html/sld001.htm (accessed Jan 30, 2014).
- (23) Urbaczek, S.; Kolodzik, A.; Fischer, J. R.; Lippert, T.; Heuser, S.; Groth, I.; Schulz-Gasch, T.; Rarey, M. NAOMI - On the almost trivial task of reading molecules from different file formats. *J. Chem. Inf. Model.* **2011**, *51*, 3199–3207.
- (24) Raczynska, E.; Kosinska, W.; Osmialowski, B.; Gawinecki, R. Tautomeric Equilibria in Relation to Pi-Electron Delocalization. *Chem. Rev.* **2005**, *105*, 3561–3612.
- (25) Morgan, H. L. The generation of a unique machine description for chemical structures - a technique developed at Chemical Abstracts Service. *J. Chem. Doc.* **1965**, *5*, 107–113.
- (26) Weininger, D.; Weininger, A.; Weininger, J. L. SMILES. 2. Algorithm for Generation of Unique SMILES Notation. *J. Chem. Inf. Comput. Sci.* **1989**, *29*, 97–101.
- (27) Irwin, J.; Shoichet, B. ZINC - A Free Database of Commercially Available Compounds for Virtual Screening. *J. Chem. Inf. Model.* **2005**, *45*, 177–182.
- (28) ZINC Database - Version 12. <https://zinc.docking.org/> (Clean Leads Reference as SMILES downloaded Dec 03, 2013).
- (29) Feng, Z.; Chen, L.; Maddula, H.; Akcan, O.; Oughtred, R.; Berman, H.; Westbrook, J. Ligand Depot: a data warehouse for ligands bound to macromolecules. *Bioinformatics* **2004**, *20*, 2153–2155.
- (30) Ligand Expo; RCSB PDB. <http://ligand-expo.rcsb.org/> (chemical component dictionary as SMILES (OpenEye with stereo) downloaded Jul 10, 2012).
- (31) Knox, C.; Law, V.; Jewison, T.; Liu, P.; Ly, S.; Frolkis, A.; Pon, A.; Banco, K.; Mak, C.; Neveu, V.; Djoumbou, Y.; Eisner, R.; Guo, A.; Wishart, D. DrugBank 3.0: a comprehensive resource for “Omics” research on drugs. *Nucleic Acids Res.* **2011**, *39*, D1035–D1041.
- (32) DrugBank 3.0. <http://www.drugbank.ca/> (all drugs as SDF downloaded Dec 03, 2013).
- (33) Bellis, L. J.; et al. Collation and data-mining of literature bioactivity data for drug discovery. *Biochem. Soc. Trans.* **2011**, *39*, 1365–1370.
- (34) ChEMBLdb - Version 17. <https://www.ebi.ac.uk/> (ChEMBLdb including an SDF file downloaded Dec 03, 2013).
- (35) Source code of the tautomer generation method by Sayle and Delany. <http://www.daylight.com/meetings/emug99/Delany/tautomers/> (accessed Jan 30, 2014).