

COSMOfrag: A Novel Tool for High-Throughput ADME Property Prediction and Similarity Screening Based on Quantum Chemistry

Martin Hornig* and Andreas Klamt

COSMOlogic GmbH and Co. KG, Burscheider Str. 515, 51381 Leverkusen, Germany

Received May 12, 2005

The COSMO–RS (Continuum Solvation Model for Real Solvents) method has proven its broad applicability for the accurate prediction of thermodynamic, environmental, or physiological properties. On the basis of quantum chemical calculations with COSMO, COSMO–RS calculations were unavoidably restricted to small- to medium-sized compound sets, because of the time demand of the COSMO calculations. The COSMOfrag method, presented here, overcomes this restriction by replacing the costly quantum chemistry step with a selection of suitable fragments from a database of, presently, 40 000 DFT/COSMO precalculated molecules. Since, in the COSMO–RS picture, any molecular information is gathered in the so-called σ profiles, COSMOfrag replaces the single σ profile with a composition of partial σ profiles, selected by the use of extensive similarity searching algorithms. On five representative datasets, the accuracy loss of COSMOfrag versus full COSMO–RS calculations has been shown to be only in the range of 0.05 log units. From the performance point of view, it is now possible to carry out COSMO–RS property calculations for more than 100 000 compounds a day per standard PC CPU.

INTRODUCTION

The virtual screening of compound libraries is well-established in modern drug discovery and design. Calculations of physicochemical properties of the drug candidates are essential for the estimation of their pharmacokinetics. To evaluate the so-called ADME (absorption, distribution, metabolism, and elimination) parameters, basically, aqueous solubility and lipophilicity are in demand,¹ though for the consideration of acidic or basic compounds, partitioning and solubility become pH-dependent and pK_a , or rather $\log D$, is needed additionally.² Lipophilicity is most commonly assessed in the form of partition coefficients, usually *n*-octanol/water (P_{OW}). However, cyclohexane/water or 1,2-dichloroethane/water³ are partly considered as more appropriate measures for lipophilicity with regard to membrane permeability. Because of the smaller water fraction in cyclohexane or 1,2-dichloroethane, these partition coefficients better account for hydrogen-bond desolvation.

The importance of molecular electrostatics and the related hydrogen-bonding and hydrophobic interactions are widely accepted in different areas of drug design. For example, models such as CoMFA⁴ in 3D-quantitative structure–activity relationships (QSAR) or molecular polar surface area descriptors (PSA)^{5–7} for quantitative structure–property relationships (QSPR) demonstrate the broad acceptance of models describing the electrostatics of molecules.

The COSMO–RS method, a combination of the quantum chemical continuum solvation model (COSMO) and a statistical thermodynamics treatment for real solvents (RS) simulations, is a novel, widely applicable tool for accurate predictions of many kinds of thermodynamic as well as physiological properties.^{8–10} In this approach, all information

about solutes and solvents is gathered from initial density functional (DFT) COSMO calculations. On the basis of this very fundamental and broad knowledge of structure and electrostatics of the molecule in solution, a large set of physicochemical properties is accessible by means of the polarization (or screening) charge density σ on the molecular surface. The COSMO–RS theory has introduced this surface polarization charge density as a novel and highly significant description of the surface electrostatics. It turned out to be more local and better transferable than the electrostatic potential (ESP) itself. Since, in COSMO–RS, properties such as $\log P_{OW}$ are calculated as surface integrals of σ functions, this theory and its σ perspective provide a qualitative and quantitative understanding of the widely recognized relation between such properties and surface electrostatics.

A straightforward and logical extension of the COSMO–RS methodology is the calculation of similarity coefficients based on σ .¹¹ This approach allows for the comparison of similarities of molecular surfaces and their electrostatics independent of the structure of the molecules, enabling scaffold hopping in a natural way.

COSMO–RS property calculations using the COSMOtherm program¹² only require fractions of a second per compound. The overall speed of the COSMO–RS method is mainly limited by the time demand of the underlying quantum chemical calculations for the molecules. With the high quality level BP-TZVP (geometry optimization with the BP functional^{13–15} and TZVP basis set¹⁶), such calculations take about 4 h, on average, for molecules with up to 40 heavy atoms on a 3 GHz CPU, using the TURBOMOLE program package (University of Karlsruhe, Karlsruhe, Germany).^{17,18} This is acceptable for chemical engineering applications where normally only a few new molecules are considered, besides many common compounds that can be taken from a database. A database of carefully prepared BP-TZVP–

* Corresponding author phone: +49-2171-731683, e-mail: hornig@cosmologic.de.

COSMO files for 3000 common compounds and solvents is available (COSMObase).¹⁹

This differs strongly in the area of drug design. Here, often up to hundreds of thousands or even millions of potential drug candidates have to be prescreened regarding their physicochemical properties or biological activities, each of them being typically in the molecular weight range of 300–500, that is, having about 25–40 heavy atoms. Therefore, we have introduced a slightly more approximated “drug calculation level”, which uses BP–SVP^{20,21} single-point DFT/COSMO calculations on semiempirical MOPAC²² AM1/COSMO geometries. This level reduces the computation time of typical drug molecules by approximately a factor of 30, that is, roughly, to 8 min per drug. Still, on this level, a prescreening of compound numbers as large as that is unfeasible even on large parallel computer clusters. For these applications, a very fast bypass for the demanding DFT/COSMO calculations called COSMOfrag²³ has been developed. This is described in the present paper. The basic idea is to avoid the time-consuming DFT/COSMO calculation of the screening charge densities (σ profiles) for each individual molecule and to replace it with a composition of partial σ profiles taken from locally most-similar fragments of molecules whose DFT/COSMO files are stored in a database. It should be noted that the database does not consist of molecule fragments but of entire molecules, and the fragmentation is individually composed from these molecules. These fragment-based σ profiles can then be used as a starting point for any COSMO–RS calculation, that is, physicochemical, physiological, or environmental²⁴ property calculations; similarity searching; or even receptor binding approaches.

GENERAL COSMO–RS THEORY

COSMO–RS is a model combining quantum theory, dielectric continuum models, surface interactions, and statistical thermodynamics. The theory of COSMO–RS has been described in detail in several articles.^{25–27} Therefore, we will only give a short survey of the basic concept here and refer the interested reader to these articles for details.

COSMO–RS considers a liquid system as an ensemble of molecules of different kinds, thus, as a solvent or solvent mixture and solutes. A precondition is a DFT/COSMO²⁸ calculation for each kind of molecule X, to get the total energy E_{COSMO}^X and the polarization (or screening) charge density (SCD) σ on its molecular surface. The COSMO calculation has to be carried out only once per compound, and thus, COSMO files can be stored for future use. The σ value is a good local descriptor of molecular surface polarity.²⁹

For the purpose of an efficient statistical thermodynamics calculation, the liquid ensemble of molecules is now considered as an ensemble of pairwise interacting molecular surfaces. The most important parts of the specific interaction between molecular surfaces, that is, electrostatics (es) and hydrogen bonding (hb), are expressed by the SCDs σ and σ' of the contacting surface pieces:

$$E_{\text{es}}(\sigma, \sigma') = \frac{\alpha'}{2}(\sigma + \sigma')^2 \quad (1)$$

and

$$E_{\text{hb}}(\sigma, \sigma') = c_{\text{hb}} \min\{0, \sigma\sigma' + \sigma_{\text{hb}}^2\} \quad (2)$$

The three parameters α' , c_{hb} , and σ_{hb} have been adjusted to a large number of thermodynamic data. Since all relevant interactions depend on σ , the distribution functions (histograms) $p^X(\sigma)$ are required for the statistical thermodynamics.

$p^X(\sigma)$, in the following σ profile, displays the composition of the ensemble of surface pieces with respect to σ . The σ profile of a special molecule has a characteristic shape and provides a vivid picture of the molecular polarity (see Figure 1 and Klamt et al.^{26,28}). Furthermore, we need the σ profile $p^S(\sigma)$ of the ensemble S, which is simply calculated as a sum of the molecular σ profiles weighted by molar fractions.

Next, the chemical potentials of the compounds in the solvent are calculated by a novel, exact, and very efficient statistical thermodynamics procedure. The first step is the iterative solution of the equation

$$\mu_S(\sigma) = -\frac{RT}{a_{\text{eff}}} \ln \left[\int d\sigma' p_S(\sigma') \exp \left\{ \frac{a_{\text{eff}}}{RT} [\mu_S(\sigma') - E(\sigma, \sigma')] \right\} \right] \quad (3)$$

This implicit equation, in which a_{eff} denotes an effectively independent piece of molecular area and $E(\sigma, \sigma')$ denotes the sum of the energy contributions of eqs 1 and 2, can be solved by iteration within milliseconds on a PC.

The resulting function $\mu_S(\sigma)$, the σ potential, describes the solvent behavior regarding electrostatics, H-bond affinity, and hydrophobicity (see Figure 2). It should be pointed out that the σ potential does not only express enthalpic aspects of the solvent–solute interactions but, as well, solvation entropy, as was explicitly demonstrated and explained by Klamt in the example of hydrocarbon solubility in water.³⁰

In a second step, the σ potential is integrated over the surface of each compound X, yielding the chemical potential of X in S:

$$\mu_S^X = \int p^X(\sigma) \mu_S(\sigma) d\sigma + \mu_{\text{combS}}^X \quad (4)$$

In this equation, the surface integral is evaluated as a σ integral, making use of the σ profile of solute X. The combinatorial contribution μ_{combS}^X to μ takes into account size and shape effects of the solute and solvent.²⁷ Usually, it is small compared to the first term in eq 4, which results from the surface interactions. It is sufficient to consider it as a solvent-specific constant, here.

Starting from a quantum chemical calculation for each compound, we found, as a result of a few statistical thermodynamical steps, an expression for the pseudochemical potential of an almost arbitrary chemical compound X in an almost arbitrary solvent S, which may be a pure compound or a mixture. This allows for the calculation of any partition coefficient as well as solubility. The few adjustable parameters required in COSMO–RS have been fitted to a large set of experimental data.²⁵

COSMOFRAG METHODOLOGY

The polarization charge density σ is a rather local feature of the molecular surface. Therefore, it reasonably can be assumed that structurally similar regions of molecules give

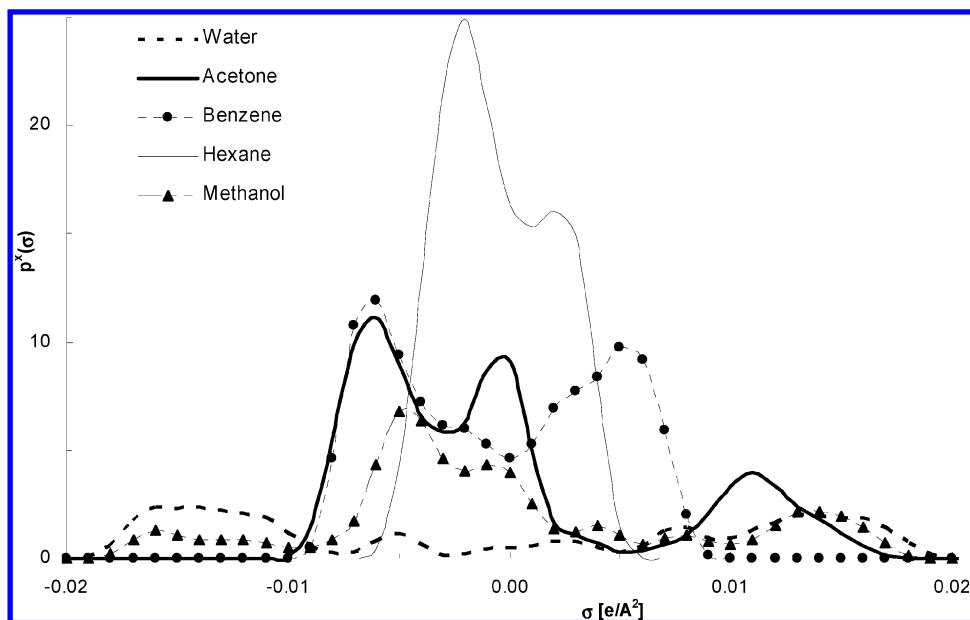


Figure 1. Solvent σ profiles. These profiles show the amount of molecular surface in a given interval of polarization charge density σ .

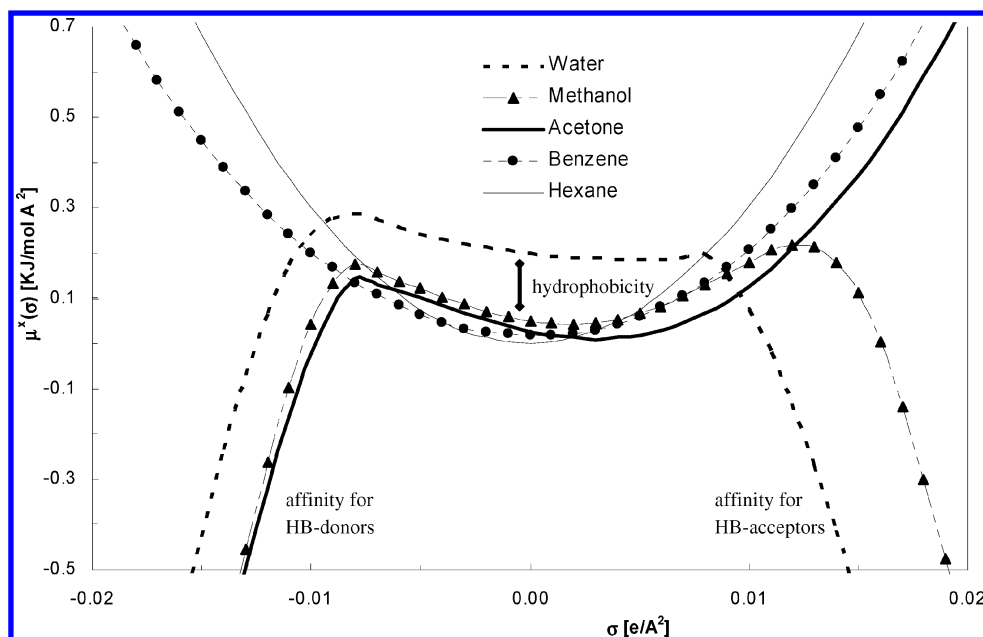


Figure 2. σ potentials of solvents. These curves show the chemical potential of surface pieces of polarization charge density σ in a solvent. Thus, they quantify the affinity of a solvent to a surface of polarity σ .

similar contributions to the σ profile. As a simple example, the contribution of the sp^3 oxygens in water and methanol can be considered, which exhibit an almost identical contribution to the σ profiles, as can be seen in Figure 1. Thus, it is plausible to assume that the σ profiles of larger new molecules can well be approximated by contributions taken from other, locally most similar molecules. Since for most COSMO-RS applications only the σ profile and some information about the area and volume of the molecule is needed, the basic idea of COSMOfrag is to compose the σ profile of new molecules from existing σ profiles of molecules that have already been precalculated and stored in a database (Figure 3). For this purpose, a database of presently 40 000 COSMO files of highly diverse, smaller basic and larger drug-like compounds has been prepared.

Conformers. A full conformational analysis of such large numbers of compounds is hardly feasible. Despite that, in

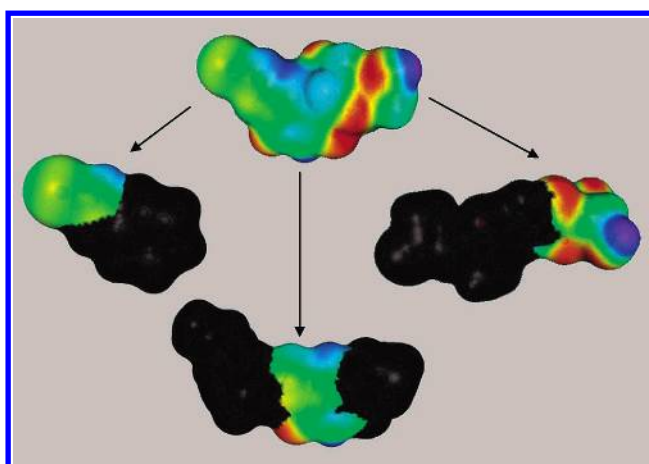
some cases, a single conformation of a molecule is insufficient for property predictions of the highest accuracy, and even, depending on the solvent, the favorable conformer may differ. Non-high-throughput calculations with COSMOtherm allow for the utilization of different conformer COSMO files by weighting them using the COSMO energies and their chemical potentials. Proceeding as differentiated like that is impracticable within a high-throughput application. Therefore, the COSMOfrag database (CFDB) consists of one single conformation for each database compound.

Attempts have been made to compose COSMOfrag databases from sets of conformers generated by two different quantum chemical procedures. One CFDB was built up from the lowest-energy conformers optimized on the MOPAC AM1 gas-phase level, the other with MOPAC AM1 COSMO optimized structures, representing the opposite cases of polar and nonpolar surroundings. The results of any calculation

Table 1. Results Statistics for Calculations of Different Physicochemical, Physiological, and Environmental Properties with COSMOfrag and on Full COSMO Files

dataset	N^a	property (log units)	COSMOfrag		(full) COSMO	
			RMS	MUE	RMS	MUE
pesticides ^b	107/105	water solubility	0.62	0.50	0.60	0.44
pesticides ^c	53/50	soil sorption	0.75	0.60	0.72	0.62
BOSS ^b	150/147	water solubility	0.70	0.56	0.66	0.52
PHYSPROP ^d	2570	pow	0.62	0.50	0.59	0.47
Abraham ^e	170/166	intestinal absorption	15.24	9.78	14.86	10.22

^a Number of compounds; for COSMOfrag calculations, they are mostly smaller, because of missing appropriate or selected inappropriate fragment molecules in the database. ^b Published in ref 9. ^c Published in ref 24. ^d Selected from PHYSPROP.³³ ^e Data from ref 34.

**Figure 3.** COSMOfrag decomposition of the drug sorivudine into three fragments. The parts of the COSMO surface not used for the decomposition are colored in black.

listed in Table 1 on both of these special CFDBs, without exception, were slightly worse (1–10% of root mean square) when compared with the standard CFDB. It, therefore, can be concluded that the influence of conformational selection on the prediction results, on average, is rather small. Even a better suitability of the polar versus the nonpolar CFDB, for example, for water solubility predictions, could not be found.

Molecules entering the CFDB are subject to a standard procedure regarding their geometrical optimization. This heuristic procedure has been developed to gain reliable quantum chemically optimized geometries from 2D structures. As a first step, a 3D starting structure is generated using CORINA,³¹ followed by an AM1/COSMO geometry optimization using a modified MOPAC⁷³² version, which is customized to produce better geometries, especially for amines, sulfonamides, and phenols. Finally, a single-point DFT/COSMO calculation (BP–SVP) with TURBOMOLE is performed. In particular cases, the MOPAC geometry optimizations, likewise MOPAC7 or an up-to-date version, result in wrong structures in comparison to the given 3D structure. This is typically the case for compounds containing multiple sulfur or phosphorus atoms. The optimization for these molecules is alternatively carried out with a DFT geometry optimization (BP–SVP) with TURBOMOLE.

COSMOfrag Similarity Algorithms. As core functionality, COSMOfrag possesses a molecular perception routine, which analyzes a molecule with respect to the hybridization states of atoms, bond orders, rings, ring properties such as aromaticity, and even their stereochemical classification. This perception can be started from most of the different common electronic file formats, such as SDF files, SMILES code,

XYZ files, and others. Some of them include bond tables with bond orders; others include elements and geometry only; some have explicit hydrogens, and other implicit hydrogen atoms have to be analyzed. It is most important that such a perception routine is able to end with a unique internal representation of the molecule, independent of all the different ways the chemical structure can be represented in the original input. Furthermore, care must be taken that equivalent atoms in a molecule also have an equivalent description. As an example, the nitro group may be considered. Normally represented by a four-valent nitrogen atom with a formal charge of +1, one single-bonded oxygen with a formal charge of −1, and a double-bonded oxygen atom, it is, as well, frequently represented by a neutral five-valent nitrogen with two neutral double-bonded oxygens. Although the former description may be closer to chemical correctness, in COSMOfrag, the second convention is applied, as it describes the two oxygen atoms as chemically equivalent. In the same way, all cases of partially ionic bond descriptions had to be reduced to a neutral multiple-bond representation. A unique representation of aromatic bonds is also of crucial importance; however, the usual Kekule description of alternating single and double bonds leads to nonunique representations.

Once a unique representation of all atoms, bonds, and rings in the molecule is achieved, the second major step is the definition of the most useful measure for the local similarity of atoms and the atomic environment. For COSMOfrag, atoms should be considered as most similar if their partial molecular surfaces and surface polarities, that is, SCDs σ , are most similar. But since the latter is not known, at least for the new molecule under consideration, we have to ensure that the local geometries and the electronic effects of the surrounding atoms are most similar. Obviously, two similar atoms should at least be identical with respect to their elements and hybridizations. By the usage of hashing algorithms, this information is turned into a unique real number for each atom, a similarity index of the lowest order (zeroth order). Since hydrogen atoms are not considered explicitly, the number of implicit hydrogen atoms is also included in the zeroth-order similarity index. In the next step, a similarity index of the first order can be defined by the propagation of the zeroth-order similarity indices of the neighbor atoms to the central atoms and the addition of this new information to the zeroth-order similarity index of the central atom. Bond orders of the bonds used for the propagation are explicitly taken into account. By doing so, the first-order identity of two atoms ensures the zeroth-order identity of all neighboring atoms. In the same way, we can

now generate the next higher similarity indices out of the similarity indices of all neighbors and continue this up to any level we require. Identity of the *i*th-order similarity index will ensure chemical identity up to the *i*th-order neighbor spheres of the atoms. Additionally, detailed information on ring sizes, cis and trans isomers, and hydrogen-bond donor or acceptor atoms is incorporated into the similarity indices. Since, for example, a carbon atom of cyclohexane must be distinguished from a central atom of a long-chain alkane, the information about the minimum ring size to which an atom belongs is integrated into the similarity indices, starting at zeroth order for three- and four-membered rings and continuing in this way for the higher similarity indices with the next higher ring sizes. Information about the cis- or trans-position with respect to double bonds is taken into account at the second order, and the information concerning whether a typical hydrogen-bond donor or acceptor atom can make a favorable intramolecular hydrogen bond, by forming a five or six ring, is included in the second-order similarity indices as well. Especially, the ability to form intramolecular H bonds may strongly change the properties of atoms in a molecule and are, thus, important for finding similar atoms in other molecules in the sense of COSMO-RS.

Indeed, the hashing algorithm used in COSMOfrag for the generation of local similarity indices is even slightly more refined in order to take care of potential problems arising from aromatic and resonance effects. To ensure the highest similarity along conjugated bond pathways, the propagation of the atomic similarity indices is additionally carried out two times along conjugated bonds, before propagation along all bonds is started. In this way, information about adjoining atoms in conjugated systems is contained in the index, which later on is denoted as the similarity index of the zeroth order. This ensures that structural similarity along conjugated bond paths is given much more attention during the fragmentation process than that along nonconjugated pathways. Thus, a high similarity of conjugated or aromatic systems and their neighbors must be given to be selected as a suitable fragment for such regions in a new target molecule.

COSMOfrag makes use of these similarity indices in two ways. First, it calculates the sum of the similarity indices of order seven for all atoms of a molecule. The resulting molecular similarity index is essentially an identity index, because, to the best of our knowledge, identical indices imply identity of the molecular structures. In COSMOfrag, this index is called a “unique name” and is used to detect the identity of molecular structures in the database. However, more important is the use of the atomic similarity indices for database screenings regarding the highest similarities of atoms. For that purpose, each of the eight similarity indices of an atom (zeroth to seventh order) is converted into a five-digit ASCII word and, afterward, combined to a 40-digit string (Figure 4). Then, the identity of the atomic ASCII similarity strings up to the fifth digit ensures zeroth-order similarity, identity up to 10th ensures first-order similarity, and so forth.

The atom-similarity strings of all atoms of the COSMOfrag database are stored as a sorted ASCII list. For the 40 000 molecules with an average of about 17 heavy atoms, this results in a list of more than 700 000 entries. For an atom out of a new molecule, the most similar atom in the database can now be found by a very efficient binary search.

```
*@_Ej )PLEU . *MDP 2°-8 8e'f'b ?,*h) F±[÷ MÉy-Ä )ÖöSy -Y$7 EJZYAKPII
* @_Ej )PLEU . *MDP 2°-8 8e'f'b ?,*h) F±@Ü± MÉ's )ÖöSy -Y$7 QNNAWWEYD
* @_Ej )PLEU . *MDP 2°-8 8e'f'b ?,*h) F±kY MÉ-Ç )ÖöSy -Y$7 EMGRWWEYS
* @_Ej )PLEU . *MDP 2°-8 8e'f'b ?,-i, F±AdÄ MÉ,8C )ÖöSy -Y$7 BZSIUAFEI
* @_Ej )PLEU . *MDP 2°-8 8e'f'b ?,-'-% F±ÆIF MÉjÜ; )ÖöSy -Y$7 CEJISAFHI
```

Figure 4. Random section of the COSMOfrag database file: five-digit atom words are separated by blanks. The first eight words are the atom similarity codes of zeroth to seventh order, followed by two bycodes containing additional information. At the end of each line, the molecule of each atom is marked, using a nine-letter unique name constructed from the molecule identity index, followed by information on the neighbor atoms and implicit hydrogens. In this section, all atoms are identical up to the fourth order, while atoms 1–3 are the most similar (up to fifth order).

```
f= COSMO/Z/XOEIAKNC.ccf CFDB w={1111111111110001111111100000000000}
f= COSMO/L/LWUCLIXML.ccf CFDB w={000011011111}
f= COSMO/J/KLBMBRKL.ccf CFDB w={10010000011000000}
```

Figure 5. COSMO metafile coding the parts of database molecules to be used as pictures for the construction of a new molecule. Database molecules are named with their unique name.

Obviously, in many cases, there will be more than one atom in the database having the same maximum similarity level. All these atoms are considered as candidates for fragment formation for the new molecule under consideration. In a final step, the fragments are built from all the candidate atoms in such a way that a small number of fragments is ensured. The result of this extensive selection process is written in a COSMO metafile (Figure 5), displaying the chosen fragment molecules and their selected atoms.

RESULTS AND DISCUSSION

Results. A selection of five datasets has been chosen to evaluate the accuracy of the prediction of different physicochemical, physiological, or environmental properties. Table 1 displays the statistics of the COSMOfrag calculations versus the experimental results, on one hand, and the calculations on the full COSMO files, on the other. Owing to the individual fragmentation based on the described concept of maximum similar substructures, the accuracy loss of COSMOfrag is always below 0.05 log units compared to direct DFT/COSMO calculations. These results demonstrate the ability of the COSMOfrag algorithms to compose sets of reasonable fragments as substitutes for a molecule under consideration and support the approach of partial σ profiles.

Since the CFDB has been constructed and extended according to maximum structural diversity in the optimal representation of typical basic and life science compounds, it meanwhile has achieved a status that ensures a good representation of most compounds appearing in life science

Table 2. Performance Statistics for the Calculations Listed in Table 1

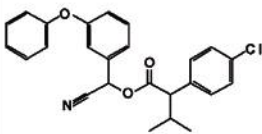
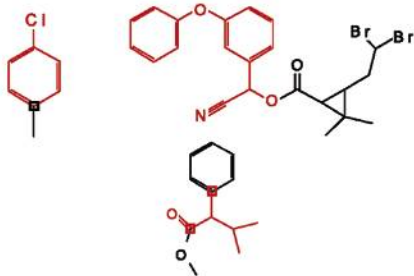
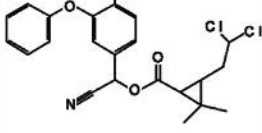
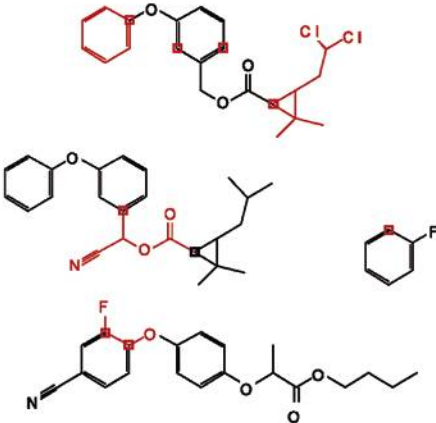
dataset	N	property	CPU ^a time	avg [s/comp.]
pesticides	107	water solubility	70 s	0.63
pesticides	53/50	soil sorption	40 s	0.75
BOSS	150/147	water solubility	150 s	1.0
PHYSPROP	2570	log Pow	1750 s	0.59
Abraham	170/166	intestinal absorption	90 s	0.53

^a Calculations carried out on a 3 GHz standard PC.

Table 3. Example Fragmentations and Solubilities of a Handpicked Number of Compounds from the Pesticides Data Set⁹

	Compound name/ CAS No.	Compound structure	Fragment structures	Log(XH2O) COSMO	Log(XH2O) Meta	Log(XH2O) Exp.
1	Indole-3- acetic acid [87-51-4]			-3.62	-3.53	-3.81
2	Cycloate [1134-23-2]			-5.10	-5.50	-5.20
3	Dinoterb [1420-07-1]			-4.43	-5.12	-6.47
4	Carbofuran [1563-66-2]			-4.50	-4.74	-4.58
5	Trifluralin [1582-09-8]			-7.70	-7.72	-8.00
6	Desmedipham [13684-56-5]			-6.29	-6.32	-6.38
7	Procymidone [32809-16-8]			-6.14	-	-6.54

Table 3 (Continued)

	Compound name/ CAS No.	Compound structure	Fragment structures	Log(XH2O) COSMO	Log(XH2O) Meta	Log(XH2O) Exp.
8	Fenvalerate [51630-58-1]			-9.49	-8.96	-9.37
9	Cyfluthrin [68359-37-5]			-9.05	-9.40	-10.04

* denotes multiple weighted fragment atoms

or drug design projects. The CFDB will be further extended in the future by parsing additional datasets for less well-represented compounds and adding these wherever possible. Presently, for a very small portion of typical datasets, a reasonable fragmentation is not possible because of missing fragments (see Table 1). Beyond such rare fragmentation failures, bad fragmentations may rarely occur in other cases. Because of wrong or incompatible conformations of the fragment molecules in the database or the weakness of the similarity algorithms, unreasonable fragmentations may occur. Such unfavorable metafiles often can be detected in the COSMOfrag results by their noticeable total COSMO charge. Altogether, 2–4% of the molecules of a dataset, on average, cannot be satisfactorily processed by the present COSMOfrag release.

Nevertheless, it must be pointed out that COSMOfrag is only applicable to neutral compounds, because ionic compounds can be much less well-fragmentized and are not represented in the CFDB for this reason. As a consequence, pK_a prediction that involves ionic species is not feasible with COSMOfrag. Also, some other COSMOtherm features that involve total energy differences of molecular species or conformations are out of the scope of COSMOfrag.

Performance Aspects. Apart from accuracy, computing time is the most important aspect when evaluating COSMOfrag. In relation to the time demand of the quantum chemical calculations, the property calculation with COSMOtherm is extremely fast. However, in connection with COSMOfrag, the COSMOtherm percentage of the overall run time lies between 30 and 80%, depending on the property

to be computed and the performance of the metafile generation on the special dataset (see Table 2). Basically, the computation of QSPR properties such as intestinal absorption or soil sorption with COSMOtherm is significantly faster than the calculation of partition coefficients or solubilities. The time demand of the COSMOfrag metafile generation, on the other hand, strongly depends on the representation of the given structures within the CFDB. If no long-ranging similarity can be found in the database, the number of fragments increases and, similarly, the computing time does as well. It can be stated that the overall performance typically lies below 1 s per compound, therefore, allowing for the calculation of 100–150 000 compounds a day on a 3 GHz standard PC.

Fragmentations. Table 3 shows a number of exemplary fragmentations generated with COSMOfrag and calculated water solubilities. Basically, it can be stated that the highest similarity of local polarity is crucial for a good property prediction. For most of the example cases, the algorithms were able to identify database molecules to be applied as fragments whose local surroundings were similar enough to assume similar electronic conditions. Compound 5 demonstrates the case of an aromatic system with strongly pulling substituents. Naturally, a fragmentation of such conjugated systems requires the selection of atoms from molecules with either the identical substituent pattern or an electronically comparable one. Subdividing of such an aromatic ring in multiple fragments is justifiable if each single fragment exhibits the special electronic conditions. The generation of many thousands of metafiles has shown that, in almost any

case, a suitable fragment molecule could be found in the database that meets the substituent pattern in demand. However, in particular cases, exactly fitting aromatic fragments are missing (e.g., fluorinated benzene ring of Compound 9). Even in such cases, an acceptable fragmentation can be achieved by superposition of aromatic fragments that only partly meet the substituent pattern. The completion of the COSMOfrag database concerning substituted aromatic or heteroaromatic fragments is an important goal for the future. Besides that, a few fragmentations (see Compound 7) provide indications of weaknesses in the current fragmentation algorithms. For the joining atoms of cyclopropane and the succinimide ring in this example, the algorithm does not demand comparable fragments also possessing a similar condensed ring system. Therefore, the single cyclopropane fragment is chosen that exhibits a completely different polarity, especially because of the four nitrile substituents. The indication of such poor fragmentations is given by COSMOfrag by means of the total COSMO charge, which is then above 0.5 in such cases.

Conformers. Conformational aspects also may have a strong influence on the prediction quality for the single molecule. This is, of course, not only the case for the fragment molecules. Compound 3 displays a case where a suboptimal conformer has been chosen on the side of the full COSMO calculations. The dinoterb geometry, in this case generated by a different 3D builder, exhibits no internal hydrogen bond that would be favorable here. This may be the reason for the large deviation between the COSMOtherm calculated and the experimental values here. The database fragment, on the other hand, shows the mentioned hydrogen bond, and the solubility prediction, therefore, is much closer to the experimental value. For the sake of completeness, the prediction for the dinoterb, as optimized by the heuristic standard optimization procedure, yields -4.85 , close to the COSMOfrag result.

SUMMARY AND OUTLOOK

The COSMOfrag method has been introduced as a high-quality shortcut for almost any kind of COSMO-RS calculation. It, therefore, makes the COSMO-RS method applicable for high-throughput tasks, especially in life science. Properties of different application areas, for example, water solubility or intestinal absorption, mostly published earlier, have been calculated with an almost negligible loss of accuracy. In the same way, COSMOfrag enables similarity screenings based on σ profiles for large numbers of compounds.

Though the COSMOfrag database of presently 40 000 molecules allows for the reliable calculation of properties for almost any class of compound in life science or drug design, it nevertheless will be extended further, especially where the electronically complicated substituted aromatics and heteroaromatics are concerned. Similarly, the refinement of the COSMOfrag hashing methodology and the systematic optimization of the CFDB molecule geometries will be continued.

REFERENCES AND NOTES

- (1) Bergstrom, C. A.; Strafford, M.; Lazorova, L.; Avdeef, A.; Luthman, K. Absorption Classification of Oral Drugs Based on Molecular Surface Properties. *J. Med. Chem.* **2003**, *46*, 558–570.
- (2) Tetko, I. V.; Poda, G. I. Application of ALOGPS 2.1 to Predict log D Distribution Coefficient for Pfizer Proprietary Compounds. *J. Med. Chem.* **2004**, *47*, 5601–5604.
- (3) Bouchard, G.; Carrupt, P.-A.; Testa, B.; Gobry, V.; Girault, H. H. Lipophilicity and Solvation of Anionic Drugs. *Chem. Eur. J.* **2002**, *8*, 3478–3484.
- (4) Cramer, R. D.; Patterson, D. E.; Bunce, J. D. Comparative Molecular Field Analysis (CoMFA). 1. Effect of Shape on Binding of Steroids to Carrier Proteins. *J. Am. Chem. Soc.* **1988**, *110*, 5959–5967.
- (5) Ertl, P.; Rohde, B.; Selzer, P. Fast calculation of molecular polar surface area as a sum of fragment based contributions and its application to the prediction of drug transport properties. *J. Med. Chem.* **2000**, *43*, 3714–3717.
- (6) van de Waterbeemd, H.; Camenisch, G.; Folkers, G.; Raevsky, O. A. Estimation of Caco-2 cell permeability using calculated molecular descriptors. *Quant. Struct.-Act. Relat.* **1996**, *15*, 480–490.
- (7) Palm, K.; Luthman, K.; Ungell, A.-L.; Strandlund, G. P. A. Correlation of drug absorption with molecular surface properties. *J. Pharm. Sci.* **1996**, *85*, 32–39.
- (8) Diedenhofen, M.; Eckert, F.; Klamt, A. Prediction of Infinite Dilution Activity Coefficients of Organic Compounds in Ionic Liquids Using COSMO-RS. *J. Chem. Eng. Data* **2003**, *48*, 475–479.
- (9) Klamt, A.; Eckert, F.; Hornig, M.; Beck, M. E.; Burger, T. Prediction of aqueous solubility of drugs and pesticides with COSMO-RS. *J. Comput. Chem.* **2002**, *23*, 275–281.
- (10) Klamt, A.; Eckert, F.; Diedenhofen, M.; Beck, M. E. First Principles Calculations of Aqueous pK_a Values for Organic and Inorganic Acids Using COSMO-RS Reveal an Inconsistency in the Slope of the pK_a Scale. *J. Phys. Chem. A* **2003**, *107*, 9380–9386.
- (11) Hornig, M.; Klamt, A. **2005**. Manuscript in preparation.
- (12) Eckert, F.; Klamt, A. *COSMOtherm*, version C2.1, revision 01.04; COSMOlogic KG: Leverkusen, Germany, 2004.
- (13) Vosko, S. H.; Wilk, L.; Nussair, M. Accurate spin-dependent electron liquid correlation energies for local spin density calculations: a critical analysis. *Can. J. Phys.* **1980**, *58*, 1200–1211.
- (14) Becke, A. D. Density-functional exchange-energy approximation with correct asymptotic behavior. *Phys. Rev. A: At., Mol., Opt. Phys.* **1988**, *38*, 3098–3100.
- (15) Perdew, J. P. Density-functional approximation for the correlation energy of the inhomogeneous electron gas. *Phys. Rev. B: Condens. Matter Mater. Phys.* **1986**, *33*.
- (16) Schaefer, A.; Huber, C.; Ahlrichs, R. Fully optimized contracted Gaussian basis sets of triple- ζ valence quality for atoms Li to Kr. *J. Chem. Phys.* **1994**, *100*, 5829–5835.
- (17) Ahlrichs, R.; Bär, M.; Häser, M.; Horn, H.; Kölmel, C. Electronic structure calculations on workstation computers: The program system TURBOMOLE. *Chem. Phys. Lett.* **1989**, *162*, 165–169.
- (18) Schaefer, A.; Klamt, A.; Sattel, D.; Lohrenz, J. C. W.; Eckert, F. COSMO Implementation in TURBOMOLE: Extension of an efficient quantum chemical code towards liquid systems. *Phys. Chem. Chem. Phys.* **2000**, *2*, 2187–2193.
- (19) Eckert, F.; Klamt, A. *COSMObase*; version C2.1, revision 01.04; COSMOlogic KG: Leverkusen, Germany, 2004.
- (20) Schaefer, A.; Horn, H.; Ahlrichs, R. Fully optimized contracted Gaussian basis sets for atoms lithium to krypton. *J. Chem. Phys.* **1992**, *97*, 2571–2577.
- (21) Eichkorn, K.; Treutler, O.; Öhm, H.; Häser, M.; Ahlrichs, R. *J. Chem. Phys.* **1994**, *100*, 5829.
- (22) *MOPAC2000*; Fujitsu: Tokyo, Japan, 2002.
- (23) Hornig, M.; Klamt, A. *COSMOfrag*, version 2.1; COSMOlogic KG: Leverkusen, Germany, 2005.
- (24) Klamt, A.; Eckert, F.; Diedenhofen, M. Prediction of soil sorption coefficients with a conductor-like screening model for real solvents. *Environ. Toxicol. Chem.* **2002**, *21*, 2562–2566.
- (25) Klamt, A. Conductor-like Screening Model for Real Solvents: A New Approach to the Quantitative Calculation of Solvation Phenomena. *J. Phys. Chem.* **1995**, *99*, 2224–2235.
- (26) Klamt, A.; Jonas, V.; Buerger, T.; Lohrenz, J. C. W. Refinement and Parametrization of COSMO-RS. *J. Phys. Chem. A* **1998**, *102*, 5074–5085.
- (27) Klamt, A.; Eckert, F. COSMO-RS: a novel and efficient method for the a priori prediction of thermophysical data of liquids. *Fluid Phase Equilib.* **2000**, *172*, 43–72.
- (28) Klamt, A.; Schueuermann, G. COSMO: a new approach to dielectric screening in solvents with explicit expressions for the screening energy and its gradient. *J. Chem. Soc., Perkin Trans. 2* **1993**, 799–805.
- (29) Klamt, A.; Eckert, F.; Hornig, M. COSMO-RS: a novel view to physiological solvation and partition questions. *J. Comput.-Aided Mol. Des.* **2001**, *15*, 355–365.
- (30) Klamt, A. Prediction of the mutual solubilities of hydrocarbons and water with COSMO-RS. *Fluid Phase Equilib.* **2003**, *206*, 223–235.

- (31) Gasteiger, J.; Rudolph, C.; Sadowski, J. CORINA Program: Molecular Networks GmbH, Erlangen, Germany. *Tetrahedron Comput. Methodol.* **1990**, *3*, 537–547.
- (32) Stewart, J. J. P. MOPAC program package (MOPAC7). *QCPE* **1993**, 455.
- (33) Howard, P.; Meylan, W. *PHYSPROP Database*; Syracuse Research Corp.: Syracuse, NY, **2000**.
- (34) Zhao, Y. H.; Le, J.; Abraham, M. H.; Hersey, A.; Eddershaw, P. J. Evaluation of human intestinal absorption data and subsequent derivation of a QSAR with the Abraham descriptors. *J. Pharm. Sci.* **2000**, *90*, 749.

CI0501948