

**Development of QSAR models using C-QSAR program: a regression program that has dual databases of over 21,000 QSAR models**

Rajeshwar P. Verma & Corwin Hansch

Department of Chemistry, Pomona College, 645 North College Avenue,  
Claremont, California 91711, USA.

Correspondence should be addressed to R.P.V. (rverma@pomona.edu )

Telephone: (909)607-4249; Fax: (909)607-7726

**Supplier Information**

BioByte

**Keywords**

C-QSAR program, quantitative structure-activity relationship (QSAR)

## **Abstract**

The interest in the application of quantitative structure-activity relationships has steadily increased in recent decades because it has repeatedly proven itself to be a low-cost, high-return investment. Potential use of QSAR models for screening of chemical databases or virtual libraries before their synthesis appears equally attractive to chemical manufacturers, pharmaceutical companies and government agencies. In the present protocol, we describe the use of C-QSAR program, a regression program used in the development of QSAR models for drug designers (especially for those who do not have extensive experience in statistics), which handles linear, parabolic and bi-linear equations with various transformation of variables, and has 'jack-knifing' capability on all types of equations. The program has a very user-friendly method of data entry as well as of verification of structures and parameters. Auto-loading of preferred parameters is a time-saving feature of C-QSAR program. C-QSAR also produces a variety of 2-D graphs as output.

## **Introduction**

C-QSAR<sup>1</sup> is a regression program used in the development of QSAR models for drug designers (especially for those who do not have extensive experience in statistics), which handles linear, parabolic and bi-linear equations with various transformation of variables, and has 'jack-knifing' capability on all types of equations. The program has a very user-friendly method of data entry as well as of verification of structures and parameters. Auto-loading of preferred parameters is a time-saving feature of C-QSAR program. It

also produces a variety of 2-D graphs as output. The C-QSAR program has dual databases of over 21,000 QSAR equations relating bio- and physico-chemical activities to structural parameters. Presently, BIO contains 12,950+ and PHYS 8,900+ equations. This is valuable for the users to validate new QSAR equations as they are being developed that is to see if the emerging structure-activity relationship bears a resemblance to others with known mechanisms. There are a number of commercial and free software programs available, which may be used to calculate descriptors and/or develop a QSAR model and are listed in **Table 1**. In C-QSAR program, the information associated with each data set, which is either in biological (BIO) or physical (PHYS) database, has been shown in **Table 2**.

## **Materials**

### Equipment

- Modern personal computer (Operating system should be Windows or Macintosh)
- Reflection for ReGIS Graphics (It is a terminal emulation program for Windows that allows any PC user to access a Unix or OpenVMS system, emulating terminals ranging from VT52 to VT400, WYSE to UNISYS).
- Versa Term Pro (It is for text and color graphics program for Macintosh, emulating DEC VT220 and Tektronix terminals, with automatic switching between them.
- C-QSAR Program (see EQUIPMENT SETUP)

## EQUIPMENT SETUP

- C-QSAR Program: The details about the C-QSAR package are available at <http://www.biobyte.com/bb/prod/cqsar.html> and their installation instructions at <http://www.biobyte.com/bb/prod/qsarinstall.pdf>.
- For **Table 1-4**, the classification of C-QSAR program, and the details about the physicochemical parameters, please see the Supplementary Information section of the attached pdf.

## Procedure

*(By using Reflection for ReGIS Graphics from a personal computer with Microsoft Windows XP)*

1. Download the Reflection for ReGIS Graphics software and contact BioByte to provide the access of C-QSAR program and to create login and password.
2. To start C-QSAR, click Start → Programs → WRQ Reflection → Host-ReGIS Graphics → A window (WRQ Reflection for ReGIS Graphics) will appear.

Connection → Connect to Best Network → cqsar.com → click 'OK'

BioByte AlphaServer 500/333

**Username:** [write down the user name and then press ENTER]

**Password:** [write down the password and then press ENTER]

\$

\$ **qsar** [press ENTER]

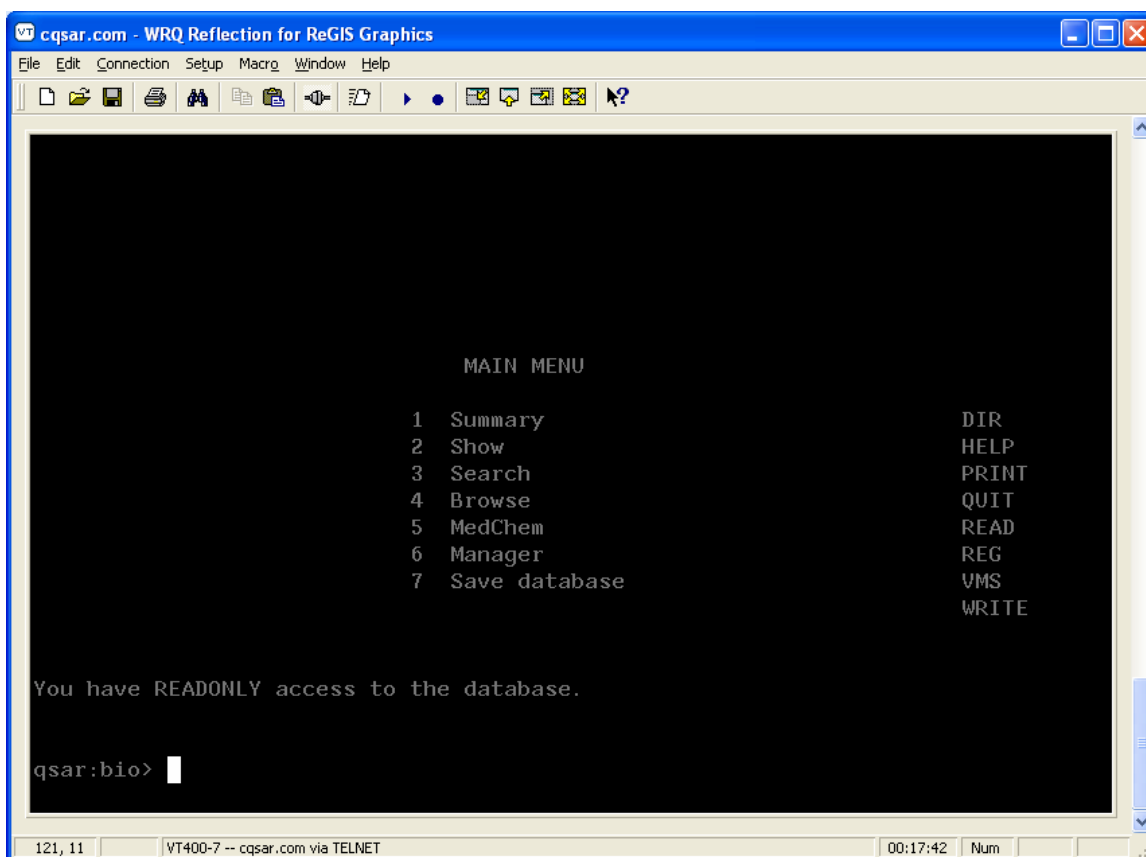
qsar>

qsar> **data bio** (or **data phys**) [press ENTER]

**QSAR Password:** [press ENTER]

**qsar: bio>**

The main menu window of C-QSAR program for BIO when the program is started (cqsar.com-WRQ Reflection for ReGIS Graphics) will appear (**Fig. 1**). [Similarly, one can obtain the main menu window of C-QSAR for PHYS by using (**qsar> data phys**) instead of (**qsar> data bio**)]



**Figure 1.** The main window of C-QSAR program when the program is started

3. Summary of BIO-Database: From the above main menu window (**Fig. 1**)

**qsar: bio> 1** [press ENTER] A summary of bio-database (**Table 3**) will appear.

Similarly, one can obtain the summary of physical-database from the main menu window of PHYS (**Table 3**).

### ? TROUBLESHOOTING

Commands are shown in **bold** for only identification purposes and not for actual use. In practice it should be used as normal typing, users can use it in the capital letter too.

## A. SEARCHING THE C-QSAR DATABASE

In this protocol, we will restrict the discussion about the search mode of this program.

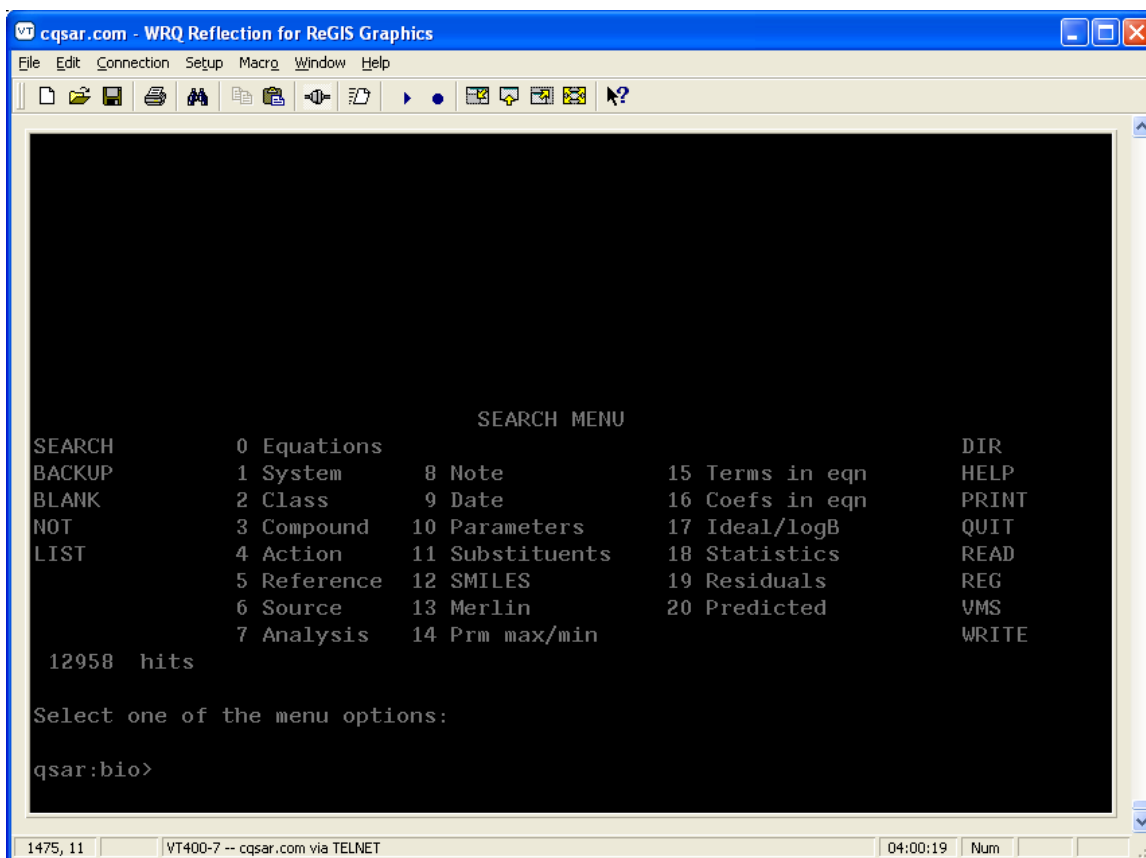
For those who become interested are referred to the earlier publications.<sup>2,3</sup> Briefly, the search mode can be approached in three ways:

- (i) String searching, based on words
- (ii) Search on parameters
- (iii) Chemical structure/molecule searching, using SMILES

### 1. Searching Bio-database

From the main menu window of C-QSAR program of Bio-database (**Fig. 1**)

**qsar: bio> 3** [press ENTER] A search menu window for Bio-database will appear (**Fig. 2**).



**Figure 2.** The search window of C-QSAR program for Bio-database

(i) **String searching:** String search is an important search mode, but one has to be careful. Some of the examples are as below:

(a) **Search using system (1):** users can search the database by using system (1) from the search menu window (**Fig. 2**)

**qsar:bio> 1 goldfish** [press ENTER] followed by

**qsar:bio>sea** [press ENTER] gave **17** hits.

Similarly, one can search on the following terms:

<u>Code</u>	<u>System</u>	<u>Hits</u> <sup>+</sup>
<u>1</u>	<u>Guppy</u>	<u>35</u>
<u>1</u>	<u>Mouse</u>	<u>461</u>
<u>1</u>	<u>Chloroplast</u>	<u>76</u>
<u>1</u>	<u>Liver</u>	<u>424</u>
<u>1</u>	<u>Human</u>	<u>1989</u> etc.

\*String search for the system shows the number of hits for the particular system as well as related to that system. Example: The search for the system 'Human' has 1989 hits; it means that the system must contained a word 'Human' that is, Human, Human Liver, Human Lymphocytes, Human Red Cell, Human Skin, Human Hemoglobin, Human Intestine, Human Platelets, Human Plasma, Human Thrombin etc.

**(b) Search using class (2):**

**qsar:bio> 2 B4C** [press ENTER] followed by **sea** command gave 2247 hits.

Similarly,

<u>Code</u>	<u>Class</u>	<u>Hits</u>
<u>2</u>	<u>B4V</u>	<u>424</u>
<u>2</u>	<u>B6F</u>	<u>208</u>
<u>2</u>	<u>B2A</u>	<u>1028</u> etc.

**(c) Search using compound (3):**

**qsar:bio> 3 phenol** [press ENTER] followed by

qsar:bio>sea [press ENTER] gave **281** hits, which will be narrow down as:

qsar:bio> **3 not miscellaneous** [press ENTER] followed by

qsar:bio>sea [press ENTER] gave **257** hits.

(d) Search using action (4):

<u>Code</u>	<u>Action</u>	<u>Hits</u>
<u>4</u>	<u>Pen Perm</u> (cell penetration)	<u>272</u>
<u>4</u>	<u>Hemolysis</u>	<u>52</u>
<u>4</u>	<u>Narcosis</u>	<u>86</u>
<u>4</u>	<u>I50</u>	<u>938</u>
<u>4</u>	<u>IC50</u>	<u>3790</u>
<u>4</u>	<u>Kill</u>	<u>199</u> etc.

(e) Search using reference (5):

**Example-1:** If users want to know about the number of QSARs in the bio-database from *Journal of Biochemistry* (stored as **J.BIOCHEM.**), the program will find all the sets from the journals where **BIOCHEM** occurs as leading or trailing or in between the names as shown as follows:

**BIOCHEM** as in **J.BIOCHEM.**

**BIOCHEM** as in **CAN.J.BIOCHEM.**

**BIOCHEM** as in **EUR.J.BIOCHEM.**

**BIOCHEM** as in **IND.J.BIOCHEM.BIOPHYS.**

**BIOCHEM** as in **J.BIOCHEM.MOL.BIOL.BIOPHYS.,** etc

The best way for any character search can be negated by prefacing it with **NOT**. This causes the result to be the reverse (logical complement) of what it would otherwise be. Thus, the search that picks only the word '**J.BIOCHEM.**' in the reference should be carried out from the search menu window of bio-database (Fig. 2).

**qsar:bio> 5 j.biochem.** [press ENTER] followed by

**qsar:bio> sea** [press ENTER] gave **46** hits.

**qsar:bio> 5 not can.j.biochem** [press ENTER] followed by

**qsar:bio> sea** [press ENTER] gave **44** hits.

**qsar:bio> 5 not eur.j.biochem** [press ENTER] followed by

**qsar:bio> sea** [press ENTER] gave **22** hits.

**qsar:bio> 5 not ind.j.biochem.biophys.** [press ENTER] followed by

**qsar:bio> sea** [press ENTER] gave **21** hits.

**qsar:bio> 5 not j.biochem.mol.biol.biophys.** [press ENTER] followed by

**qsar:bio> sea** [press ENTER] gave **20** hits.

**Example-2:** If user is interested to find out a particular paper [Selassie, C.D., Kapur, S., Verma, R.P. & Rosario, M. *J. Med. Chem.* **48**, 7234-7242 (2005)], which is available in the bio-database. It can be done from the search menu window of bio-database (Fig. 2) as follows:

**qsar:bio> 5 selassie** [press ENTER] followed by

**qsar:bio>sea** [press ENTER] gave **57** hits.

**qsar:bio> 5 2005** [press ENTER] followed by

**qsar:bio> sea** [press ENTER] gave 7 hits.

**qsar:bio> sh 5** [press ENTER] lists the same reference of Selassie et al. with R7518 and seven sets with their set numbers 12865, 12866, 12867, 12868, 12869, 12870 & 12871. The desired summary of these sets can be obtained by **show** command as follows:

**qsar:bio> sh 1 2 3 4 16 18** [press ENTER] will give the following details for all seven sets, but here we are giving the details for only one set:

Set No.	12865
System	L1210 CELLS
Class	B4C ; Cells in culture
Compound	X-PHENOLS
Action	IC50: CASPASE-MEDIATED APOPTOSIS OF X-PHENOLS
Equation	$\log 1/C = 1.06(\pm 0.12)B_{5_2} + 0.33(\pm 0.20)B_{5_3} - 0.18(\pm 0.09)\pi_{2,4} - 0.92(0.46)$ (1)
Stats	$n = 51$ $DF = 47$ $s = 0.349$ $r = 0.941$ $q^2 = 0.866$ $omit = 1$

Set No.	12865
System	L1210 CELLS
Class	B4C ; Cells in culture
Compound	X-PHENOLS
Action	IC50: CASPASE-MEDIATED APOPTOSIS OF X-PHENOLS
Equation	$\log 1/C = 1.06(\pm 0.12)B_{5_2} + 0.33(\pm 0.20)B_{5_3} - 0.18(\pm 0.09)\pi_{2,4} - 0.92(0.46)$ (1)
Stats	$n = 51$ $DF = 47$ $s = 0.349$ $r = 0.941$ $q^2 = 0.866$ $omit = 1$

Where,  $C$  represents the concentration of X-phenol that induces caspase-mediated apoptosis by 50%.  $B5_2$  is Verloop's sterimol descriptor and is a measure of the maximum width of the substituents in the ortho position, while  $B5_3$  represents the maximum width of the substituents in the meta position. The best way to open a data set from the main menu/regression mode that will be discussed later.

**(f) Search using source (6):** Person who entered the data sets.

**qsar:bio> 6 verma** [press ENTER] followed by

**qsar:bio> sea** [press ENTER] gave **1679** hits.

**(g) Search using analysis (7):** Person who analyzed (checked) the data sets.

**qsar:bio> 7 verma** [press ENTER] followed by

**qsar:bio> sea** [press ENTER] gave **63** hits.

**(h) Search using Prm max/min (14):** In this search, one can find all the sets in which every compound that has a  $\log 1/C$  of  $n$  or greater and also the possibility to find out the sets in which at least one compound has a  $\log 1/C$  of  $n$  or greater.

The search in the biological data base is as follows:

**qsar:bio> 14 log1/C>9** [press ENTER] followed by **sea** command gave **28** hits.

This indicates that the bio-database has a total number of **28** sets in which every compound that has a  $\log 1/C$  of **9** or greater. Similarly, the other command;

**qsar:bio> 14 log1/C@max>9** [press ENTER] followed by **sea** command gave **770** hits. This indicates that the bio-database has a total number of **770** sets in which at least one compound has a log 1/C of **9** or greater.

**(i) Search using Statistics (18):**

**qsar:bio> 18 2<terms<4** [press ENTER] followed by **sea** command gave **1909** hits – isolates all QSARs having 3 terms in the bio-database.

**qsar:bio> 18 n>75** [press ENTER] followed by **sea** command gave **78** hits – isolates all QSARs based on more than 75 data-points.

**qsar:bio> 18 r>.99** [press ENTER] followed by **sea** command gave **1075** hits – selects all QSARs with *r* greater than 0.99

- (ii) Search on parameters:** Six parameters (Clog *P*, Mlog *P*, CMR, NVE, MgVol, and MW) as well as forty-four parameters in **Table 4** can be automatically loaded for QSAR calculations. S stands for Hammett sigma  $\sigma$ ; -P and -M stand for para and meta values, respectively. In the broader sense para values are used for aromatic substituents conjugated with the reaction center and meta values for non-conjugated aromatic systems. These Hammett-type parameters [ $\sigma$ ,  $\sigma^+$ ,  $\sigma^-$ ,  $\sigma^*$  (S-star), and  $\sigma_I$  (S-inductive)] are obtained over half a century of study and testing on simple organic reaction mechanisms. Their use in the formulation of biological QSAR has already been discussed.<sup>4</sup> The resonance parameters (R) and field/inductive (F) have also been reviewed.<sup>5</sup> Molecular orbital parameters continue to be explored for the use in both biological and physical QSAR since

there are many instances where Hammett constants cannot be used.<sup>6-8</sup> Searching the biological database from their search menu window (**Fig. 2**):

**qsar:bio> 10 HOMO LUMO** [press ENTER] followed by

**qsar:bio>sea** [press ENTER] gave **166** hits. (HOMO or LUMO was tested)

**qsar:bio> 15 HOMO LUMO** [press ENTER] followed by

**qsar:bio>sea** [press ENTER] gave **78** hits. (HOMO or LUMO was used)

This figure **78** shows that in **88** of the examples, the molecular orbital parameters (HOMO or LUMO) were tested but found to be not as sound as Hammett constants. However, this statistic must be considered with caution since not all calculations were made with some of the more rigorous computational programs now available.

The crucial parameter for the initial success of the biological QSAR was the numerical value of hydrophobic interactions.<sup>9</sup> Despite the great complexity of studies of all types of chemicals reacting with various kinds of biological systems (from DNA to whole animals), the *n*-octanol/water partition coefficient used in log terms provides surprising insights. The hydrophobic parameter for the substituents ( $\pi$ ) can be of great importance in delineating local hydrophobic interactions at the receptor level.<sup>4</sup> Partition coefficients are rarely measured these days because it is costlier as well as time-consuming process. In the other hand, the use of data from the literature to formulate QSAR means that the compounds are not usually available for the measurement of their partition coefficients. C-QSAR program contains 12958 QSARs in bio-database, 6953 contain log *P* terms and 930 have  $\pi$  terms; hence, it is very important to have the best possible means

for their calculations. There are now a variety of methods for the calculation of  $\log P$ .<sup>10</sup> The most extensively used method is that of Leo.<sup>10,11</sup> The other important parameters are steric parameters i.e. MR-SUB, CMR, MgVol,  $E_s$ ,  $L$ ,  $B_1$ , and  $B_5$ .

In the biological QSAR  $\log 1/C$  is in molar terms except in a few cases marked by  $\log 1/C'$ . The following approaches are used in the search of QSARs of particular type, which contain the desired parameter and are illustrated as:

No.	Command (press ENTER followed by <u>sea</u> command)	Hits
1	qsar:bio> <u>15</u> " <u>log1/C</u> "	<u>8607</u>
2	qsar:bio> <u>15</u> <u>not</u> <u>**2</u> <u>bilin</u>	<u>6549</u>
3	qsar:bio> <u>15</u> " <u>logP</u> " " <u>ClogP</u> "	<u>3104</u>
4	qsar:bio> <u>15</u> <u>not</u> "S	<u>2638</u>
5	qsar:bio> <u>15</u> <u>not</u> <u>ES</u> <u>B1</u> <u>B5</u> <u>MR</u> <u>Pi</u> <u>PKA</u>	<u>2092</u>
6	qsar:bio> <u>16</u> <u>0.6</u> <LOGP<1< B>	<u>639</u>
7	qsar:bio> <u>16</u> <u>0</u> <CONST<0.5< B>	<u>55</u>

The first step ensures that  $1/C$  values are standard. The second step eliminates all QSARs with nonlinear terms, and the third step ensures that we have only  $n$ -octanol/water  $\log P$  values. Searches 4 and 5 eliminate parameters other than  $\log P$ . Step 6 selects only those QSARs where the coefficient with  $\log P$  is between 0.6 and 1.0, and 7 eliminates QSARs whose intercept is outside of 0 and 0.5.

Similarly, we can search the non-linear QSARs:

**(a) Optimal hydrophobicity:**

No.	Command (press ENTER followed by <u>sea</u> command)	Hits
1	<code>qsar:bio&gt; <u>15 logP</u></code>	<u>6953</u>
2	<code>qsar:bio&gt; <u>15 logP**2 bilin(logP) bilin(ClogP)</u></code>	<u>1786</u>
3	<code>qsar:bio&gt; <u>17 1.5&lt;logP&lt;2.5</u></code>	<u>170</u>

In this search,  $\log P^{**2}$  represents  $\log P^2$ . The third step narrows the catch to  $\log P_O$  (Optimal hydrophobicity) values between 1.5 and 2.5. Now, We can use the **show** command that is; `qsar:bio> sh 5` [press ENTER] list the results. For parabolic equations,  $\log P_O$  is displayed with its confidence limits, when it is possible to calculate them. One of the advantages of the parabolic model is that an estimate of  $\log P_O$  can be obtained without having data-points on the down side of the curve, which is necessary to derive the bilinear model.

**(b) Parabolic QSARs in terms of Clog P:**

No.	Command (press ENTER followed by <u>sea</u> command)	Hits
1	<code>qsar:bio&gt; <u>16 ClogP&gt;0</u></code>	<u>4377</u>
2	<code>qsar:bio&gt; <u>16 ClogP**2&lt;0</u></code>	<u>696</u>

Now, users can use the **show** command: `qsar:bio> sh 5` [press ENTER], which list the results.

**(c) Inverted Parabolic QSARs in terms of Clog P:**

No.	Command (press ENTER followed by <u>sea</u> command)	Hits
1	qsar:bio> <u>16 ClogP&lt;0</u>	<u>2795</u>
2	qsar:bio> <u>16 ClogP**2&gt;0</u>	<u>158</u>

Inverted parabolic QSAR may correspond to an allosteric reaction.

**(d) Bilinear QSARs in terms of Clog P:**

No.	Command (press ENTER followed by <u>sea</u> command)	Hits
1	qsar:bio> <u>16 ClogP&gt;0</u>	<u>4377</u>
2	qsar:bio> <u>16 bilin(ClogP)&lt;0</u>	<u>602</u>

**(e) Inverted Bilinear QSARs in terms of Clog P:**

No.	Command (press ENTER followed by <u>sea</u> command)	Hits
1	qsar:bio> <u>16 ClogP&lt;0</u>	<u>2795</u>
2	qsar:bio> <u>16 bilin(ClogP)&gt;0</u>	<u>57</u>

- (iii) **Chemical structure/molecule searching, using SMILES:** The SMILES search can be approached in two ways. One can identify every QSAR that contain a *specific* molecule, or one can use a MERLIN search that finds all *derivatives* of a given structures. For example: search for phenol (Oc1ccccc1) in bio-database **qsar:bio> 12 Oc1ccccc1** [press ENTER] followed by sea command gave **332** hits in the bio-database – finds **332** datasets that include un-substituted phenol.

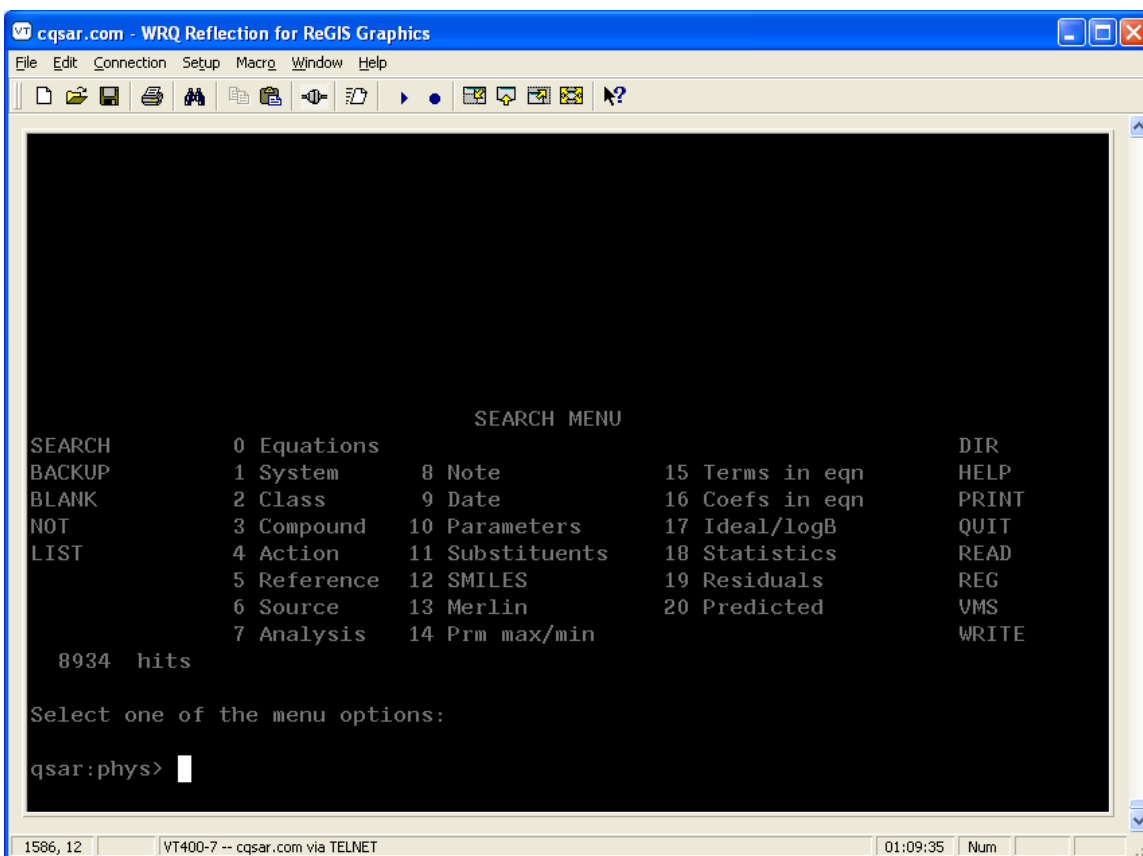
qsar:bio> 13 Oc1ccccc1 [press ENTER] followed by sea command gave **7805** hits – finds **7805** datasets that include at least one derivative of phenol.

## 2. Searching Phys-database

From the main window of C-QSAR program of Phys-database

qsar: phys> 3 [press ENTER] A search menu window for Phys-database will appear (Fig. 3).

Now, users can search the physical database as similar to that of bio-database using the search menu window (Fig. 3).



**Figure 3.** The search window of C-QSAR program for Phys-database

## ? TROUBLESHOOTING

If one used the search command for one kind followed by sea command to check the number of the hits and not used the show command then only blank command will be needed to return into the normal search window and other type of search should be followed. On the other hand, if one used the show command then quit command will give the main menu, 3 [ENTER] will give the previous search window (not the normal search window), now the use of blank command will return into the normal search window.

## B. LOADING DATA SET FROM 'DATABASE SEARCH' INTO 'WORK SPACE'

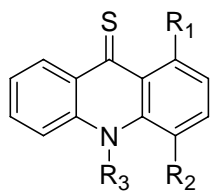
After a data set of interest has been found by means of the 'sea' and 'sh' commands, it can be further examined in details if transferred into regression mode via its set number. This is done by entering regression command from either the Bio or Phys database search menu window (Fig. 2 or 3):

```
qsar:bio> reg [press ENTER]
```

```
qsar>
```

(Note: This mode can also be obtained from either the main menu of Bio or Phys database using regression command)

**If users want to open a set no. 12675 from bio-database, which is for the binding of thioacridone derivatives (I) to DNA, then we need to follow the following steps in the regression mode:**



**I**

**qsar> load /d 12675** [press ENTER] will transfer the set to the workspace

**qsar> sum** [press ENTER] will list the key items of the data set as follows:

QSAR DB name: BIO  
 THOR DB name: QSAR\_MC\$BIGDATA:MASTER  
 Dataset name: BIO\_12675  
 Substituents: 15 Parameters: 7 SMILES: 15 Equations: 1  
 Active: 14 Starred: 1 Inactive: 0  
 System: DNA  
 Class: B1 ; Nonenzymatic Macromolecules:  
 Compound: THIOACRIDONE DERIVATIVES  
 Action: K: BINDING CONSTANT  
 Reference: DHEYONGERA, J.P., GELDENHUYS, W.J., DEKKER, T.G.,  
 VAN DER SCHYF, C.J., BIOORG. MED. CHEM., 13, 689-  
 698(2005) R7466  
 Source: RAJESHWAR PRASAD VERMA  
 Analysis: UNKNOWN  
 Date: 2006 October 26  
 Parameters: YPRED DEV LOGK CLOGP CMR NVE MGVOL

**qsar> seeeq** [press ENTER] will list the following equation:

$$\log K = -0.77(\pm 0.28)\text{Clog } P + 0.95(\pm 0.14)\text{CMR} - 4.67(\pm 1.76) \quad (2)$$

$$n = 14, \quad r = 0.981, \quad q^2 = 0.933, \quad \text{SS1} = 22.23, \quad \text{DEV}^+ = 6$$

$$\text{DF} = 11, \quad r^2 = 0.963, \quad s = 0.275, \quad \text{SS2} = 0.831, \quad \text{DEV}^- = 8$$

**qsar> pred** [press ENTER] will list the data in tabular form:

No.	R <sub>1</sub>	R <sub>2</sub>	R <sub>3</sub>	log K	Y <sub>pred.</sub>	Dev	Clog P	CMR
1	NH(CH <sub>2</sub> ) <sub>2</sub> N(CH <sub>3</sub> ) <sub>2</sub>	H	CH <sub>3</sub>	2.93	2.46	0.48	2.99	9.98
2	NH(CH <sub>2</sub> ) <sub>2</sub> N(CH <sub>3</sub> ) <sub>2</sub>	Cl	CH <sub>3</sub>	2.04	2.25	-0.21	3.86	10.47
3	NH(CH <sub>2</sub> ) <sub>2</sub> N(CH <sub>3</sub> ) <sub>2</sub>	CH <sub>3</sub>	CH <sub>3</sub>	2.36	2.51	-0.15	3.49	10.44
*4	NH(CH <sub>2</sub> ) <sub>2</sub> N(CH <sub>3</sub> ) <sub>2</sub>	H	H	4.41	2.56	1.85	2.28	9.52
5	NH(CH <sub>2</sub> ) <sub>2</sub> N(CH <sub>3</sub> ) <sub>2</sub>	Cl	H	2.08	2.21	-0.13	3.35	10.01
6	NH(CH <sub>2</sub> ) <sub>2</sub> N(CH <sub>3</sub> ) <sub>2</sub>	CH <sub>3</sub>	H	2.43	2.62	-0.19	2.78	9.98
7	N(CH <sub>2</sub> CH <sub>2</sub> Cl) <sub>2</sub>	H	H	2.04	1.88	0.16	3.93	10.13
8	N(CH <sub>2</sub> CH <sub>2</sub> Cl) <sub>2</sub>	Cl	H	1.95	1.62	0.33	4.86	10.62
9	N(CH <sub>2</sub> CH <sub>2</sub> Cl) <sub>2</sub>	CH <sub>3</sub>	H	1.78	1.93	-0.15	4.42	10.59
10	Cl	H	H	-0.1	-0.33	0.23	3.46	7.41
11	Cl	Cl	H	-0.22	-0.46	0.24	4.23	7.91
12	Cl	CH <sub>3</sub>	H	-0.15	-0.27	0.12	3.96	7.88
13	Cl	H	CH <sub>3</sub>	-0.52	-0.15	-0.37	3.8	7.88
14	Cl	Cl	CH <sub>3</sub>	-0.4	-0.25	-0.15	4.52	8.37
15	Cl	CH <sub>3</sub>	CH <sub>3</sub>	-0.3	-0.1	-0.2	4.3	8.34

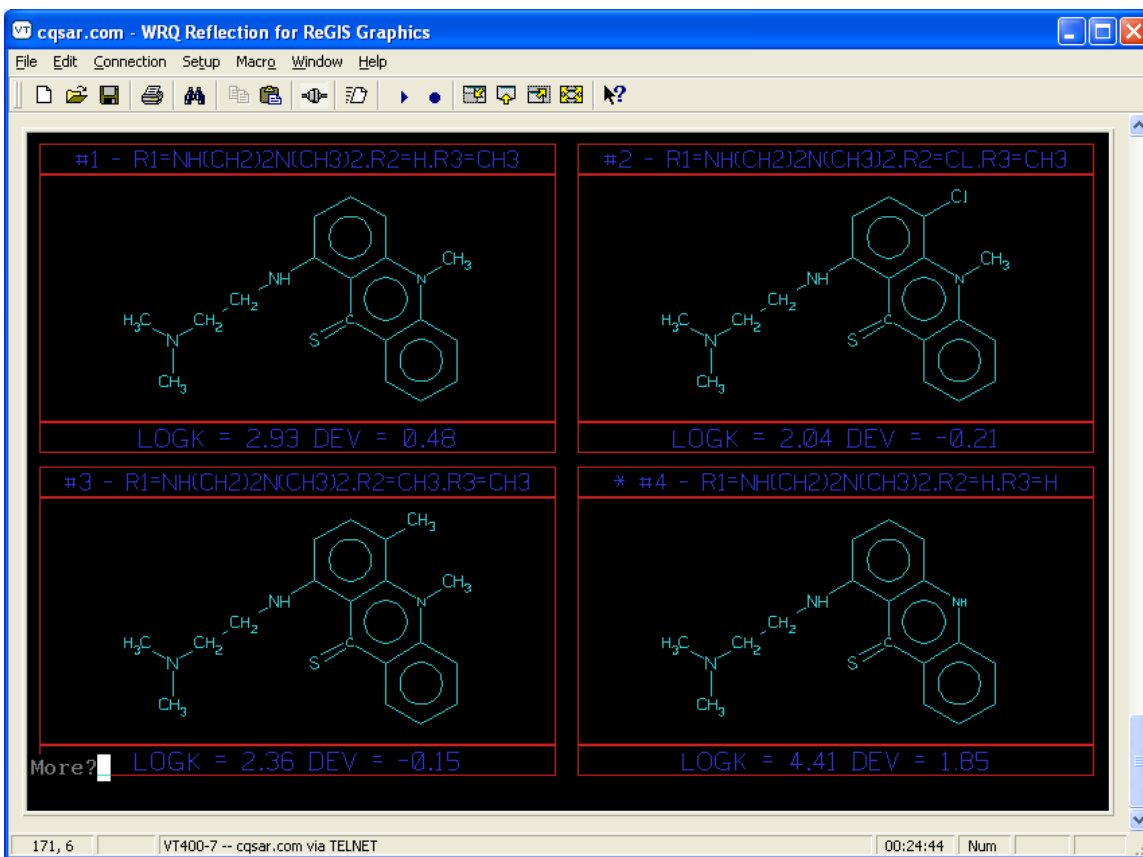
\*Not used in the development of QSAR that is outlier

In this set, 'Parameter' shows all parameters considered in this study. 'Ypred' is the predicted value from the stored equation. 'Dev' is the difference between this figure and the log of the observed value. In the QSAR equation,  $n$  is the number of data points,  $r$  is the correlation coefficient between observed values of the dependent and the values calculated from the equation,  $r^2$  is the square of the correlation coefficient represents the goodness of fit,  $q^2$  is the cross-validated  $r^2$  (a measure of the quality of the QSAR model and obtained by using leave-one-out procedure<sup>12</sup>), and  $s$  is the standard deviation. DF is the number of degree of freedom. SS1 is the sum of squares about the mean of the dependent variable and SS2 is the sum of squares from the deviations from the regression

line. DEV+ and DEV- are the number of positive and negative deviations respectively from the QSAR.

**To see the SMILES generated structures for the compounds of a data set:**

**qsar> depict ,** [press ENTER] will depict sequentially all of the compound structures as one presses **ENTER** after each panel of 4 structures (**Fig. 4**). The process can be stopped at any point by entering **q**. This can be very important in dealing with a large data set, say > 100 compounds. To view any particular structure enter **depict #** (compound number in the set). To check all structures following compound 5, enter **depict 5,** to view all structures up to 7, enter **depict ,7.** To see those between 4 and 7, enter **depict 4,7.**



**Figure 4.** Sequentially depiction of four structures of a data set (#12675)

### C. DERIVATION OF A QSAR MODEL

For several reason it is advisable to begin the regression program in the proper area, either **database bio** or **database phys**. Searching for the similar equations can then be carried out directly and if the developed equation is useful, it is easier to save in that area.

**qsar:bio> reg** [press ENTER]

**qsar> clear** [press ENTER] To be sure the workspace is empty

#### (i) Title Information (Set No. B12870)

**qsar> name Selassie-1** [press ENTER]

**qsar> t/sys CCRF (sensitive cell)** [press ENTER]

**qsar> t/comp 4-X-phenols** [press ENTER]

**qsar> t/act log1/C: Cytotoxicity** [press ENTER]

**qsar> ref Selassie,C.D., Kapur,S., Verma,R.P., Rosario,M., J. Med. Chem. 48,7234-7242 (2005)** [press ENTER]

**qsar> t/source Rajeshwar Prasad Verma** [press ENTER]

**qsar> t/class b4c** [press ENTER]

**qsar> sum** [press ENTER] To check the correctness of the above information

(ii) **Naming Parameters:** Automatic loading will be demonstrated in this example, and so only the dependent variable need to be entered.

**qsar> getp** [press ENTER]

**Label for parameter 3: log1/C** [press ENTER] (Since parameter 1 is reserved for predicted values and 2 is used for the deviation.)

Since, M.O. parameters, BDE, or pKa are not auto-loaded, they will be entered manually. Thus, in the present example the BDE parameter will be entered as parameter 4. This can also be entered in later by using **newp** command, Label for parameter #: '**BDE**' and finally the values of this parameter for each compound.

**Label for parameter 4: BDE** [press ENTER]

**Label for parameter 5: end** [press ENTER]

? **TROUBLESHOOTING:** In general the biological activity data is not published in logarithmic form with negative sign and in molar concentration. They can be entered as such and then converted into this form by using **gettran** command.

For example, if the biological activity is given in micro mole ( $\mu\text{M}$ ), then the data first entered as such under 'C' for the Label of parameter 3 and then converted into  $\log_1/C$  and in molar concentration by using **gettran** command as follows:

**Label for parameter 3: C** [press ENTER]

**Label for parameter 4: end** [press ENTER]

**qsar> newsub** [press ENTER]

**Label for substituent 1: Substituent 1 (or compound)** [press ENTER]

**C -- parameter value 3: activity value** [press ENTER]

**Label for substituent 2: Substituent 2 (or compound)** [press ENTER]

**C -- parameter value 3: activity value** [press ENTER]

Similarly, all the labels and parameter 3 values should be enter and then

**Label for substituent #: end** [press ENTER]

**qsar> seed** [press ENTER] To see the correctness of the data

**qsar> gett** [press ENTER] A box (Enter new label) will appeared. Type **log1/C** in the box and press ENTER. Now a second box (Enter transformation) will obtain. In this box, type the following: **6 - log C** [press ENTER] The transformed activity will be in log1/C (molar concentration) and present in the next **Label for parameter** (i.e. 4). Now, one wants to delete the parameter 'C' (Label for parameter 3). It can be done by the following command:

**qsar> del /para 3** [press ENTER]

**3 C**

**All these parameters will be deleted. OK? (yes/no): press y**

**The number of items deleted: 1**

**qsar> save** [press ENTER] To save the data

Similarly, the other transformation should be carried out by using **gettran** command.

<b><u>C</u></b> <b><u>(concentration of the activity)</u></b>	<b><u>Enter transformation</u></b> <b><u>(log1/C in molar concentration)</u></b>
<b>Nano molar (nM)</b>	<b><u>9 - log C</u></b>
<b>Mili molar (mM)</b>	<b><u>3 - log C</u></b>

If the biological activity concentration is in microgram/ml (mcg/ml), then the data first entered as such under 'X' for the Label of parameter **3** and then converted into 'C' in micro molar concentration by using **gettran** command as follows:

**qsar> gett** [press ENTER] A box (Enter new label) will appeared. Type **C** in the box and press ENTER. Now a second box (Enter transformation) will obtain. In this box, type the following: **X / MW** [press ENTER] The transformed activity

will be in C (micro molar concentration) and present in the next **Label for parameter** (i.e. 4). Now this will be converted into log1/C (molar concentration) as shown above.

**(iii) Naming and Entering Substituents:**

**qsar> newsub** [press ENTER]

**Label for substituent 1: NH2** [press ENTER]

**Log1/C -- parameter value 3: 4.61** [press ENTER]

**BDE -- parameter value 4: -9.25** [press ENTER]

**Label for substituent 2: OC6H13** [press ENTER]

**Log1/C -- parameter value 3: 5.34** [press ENTER]

**BDE -- parameter value 4: -6.30** [press ENTER]

Similarly, all the labels and parameters 3 and 4 values should be entered and then

**Label for substituent 11: end** [press ENTER]

? **TROUBLESHOOTING:** It is important to note that there must be no spaces within the label, i.e., 2,4-di-CL not 2 4 di-CL.

**qsar> seed** [press ENTER] Yields the following table:

No.	Substituent	log1/C	BDE
1	NH <sub>2</sub>	4.61	-9.25
2	OC <sub>6</sub> H <sub>13</sub>	5.34	-6.30
3	Me	3.82	-2.22
4	CMe <sub>3</sub>	3.86	-1.54
5	H	3.27	0.00
6	OMe	4.41	-6.01
7	C <sub>3</sub> H <sub>7</sub>	3.72	-2.01
8	C(Me) <sub>2</sub> C <sub>6</sub> H <sub>4</sub> -4-OH	4.18	-1.88
9	C <sub>2</sub> H <sub>5</sub>	3.79	-1.90
10	C <sub>8</sub> H <sub>17</sub>	4.66	-2.17

? **TROUBLESHOOTING:** If entry errors need to be corrected, users can enter **editsub** for editing substituents or **editdata** for editing data.

**Editing Substituent:** Example- If 2<sup>nd</sup> substituent OC<sub>6</sub>H<sub>13</sub> should be OC<sub>6</sub>H<sub>5</sub>, then  
**qsar> editsub 2** [press ENTER] A substituent label box will be obtained, Now one can simply correct the substituent and then press ENTER.

**qsar> save** [press ENTER]

**Editing Data:** Example- If log1/C value for cpd #1 should be 3.61 and not 4.61, then

**qsar> editdata 1** [press ENTER] A substituent box (having typed 1) will be obtained. Again pressing ENTER gave a blank parameter box, type **3** in the box (because log1/C is the parameter 3) and presses ENTER. This will give a blank box for the parameter value. Now user can simply type the correct value of log1/C and press ENTER.

**qsar> save** [press ENTER]

It is advisable to save the data entry frequently by entering **save** command.

#### (iv) Entering Structures via SMILES:

(a) If the data set is not based on a parent structure, then the SMILES for each compound should be entered one at a time and auto-loading of parameters is not possible. *(This step should be avoided if the data set is based on parent structure)*

**qsar> getsmi** [press ENTER] will provide a panel with a prompt for structure **1**.

When the SMILES of structure #1 is entered, 'pressing **ENTER**' will display the 2-D structure. If the structure needs editing, enter '**y**' if not, enter **n** and the

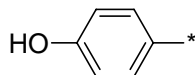
prompt for the second SMILES will appear. Since there are a large number of SMILES stored in the database together with name(s), the name can be entered at this point and the SMILES will be picked up from the database. When all the SMILES have been added, then enter **end** at the panel for next compound and [press ENTER], the prompt will return to qsar (**qsar>**).

**qsar> save** [press ENTER]

? **TROUBLESHOOTING:** Common name of the compound/drug should be used to enter their structure, i.e. acetic acid not ethanoic acid; p-chlorophenol not 4-chlorophenol. This can be a very time saving procedure for complex structures such as strychnine.

(b) If the data set is based on a **parent structure**, then the SMILES for the parent compound should be entered and auto-loading of parameters is possible. The present example is a set of 4-X-phenols, automatic loading is to be used, and the parent structure is entered via **getsmi /p**. A panel is displayed into which one enters the SMILES with an \* for each substituent position. For the present example:

**qsar> getsmi /p** [press ENTER] A panel is displayed into which one enters the SMILES with an \* for the substituent at 4-position, and a proper SMILES for the parent is: **Oc1ccc(\*)cc1**. Pressing ENTER should then display:



If the parent structure is not correct, enter **y** for editing. When the parent structure is correct, enter **n** and the prompt return to qsar (**qsar>**).

**qsar> getsmi** [press ENTER] A panel is displayed for the entry of first compound. **\*N** [press ENTER] to see 4-aminophenol. If editing is not needed, enter **n** and then the panel will ask to enter the second structure and on so on:

<u>*OCCCCC</u>	4- <i>n</i> -hexyloxyphenol
<u>*C</u>	4-methylphenol
<u>*C(C)(C)C</u>	4- <i>t</i> -butylphenol
<u>*H</u>	phenol
<u>*OC</u>	4-methoxyphenol
<u>*CCC</u>	4- <i>n</i> -propylphenol
<u>*C(C)(C)c2ccc(O)cc2</u>	4-CMe <sub>2</sub> C <sub>6</sub> H <sub>4</sub> (4-OH)-phenol
<u>*CC</u>	4-ethylphenol
<u>*CCCCCCC</u>	4- <i>n</i> -octylphenol

Now enter **end** at the panel for next compound and [press ENTER], the prompt will return to qsar (**qsar>**).

**qsar> depict ,** [press ENTER] To check that all SMILES correspond to the substituent name. If a particular SMILES is incorrect, then

**qsar> editsmi #** (the compound number) [press ENTER] and make the correction

**qsar> save** [press ENTER]

? **TROUBLESHOOTING:** The SMILES (**Simplified Molecular Input Line Entry System**) is a language for linear entry of complex structures of organic compounds into computer, which was invented by David Weininger.<sup>13</sup>

**(v) Auto-Loading of Parameters:**

**(a) Physicochemical parameters for the whole molecule (Clog *P*, Mlog *P*, CMR, NVE, MgVol, and MW):** Auto-loading of these six parameters is for the both types of data sets, which is either based on the parent structure or not.

**qsar> addcal** [press ENTER] This will auto-load four parameters (Clog *P*, CMR, NVE, and MgVol)

**qsar> add mlogp** [press ENTER] Mlog *P* will be auto-loaded

**qsar> add mw** [press ENTER] MW will be auto-loaded

**qsar> seed** [press ENTER] To see all the parameters in tabular form

**qsar> save** [press ENTER] To save the data.

**(b) Physicochemical parameters for the substituents (forty-four parameters):**

These physicochemical parameters for the substituents (**Table 4**) can be auto-loaded by using **f**etch command if the structures of the data set are entered by the use of PARENT SMILES. For the present example it can be demonstrated as:

**qsar> f** [press ENTER] This will show **Oc1ccc(\*)cc1**  
**1**

**Pick one or more positions to parameterize: 1** [press ENTER] (For this example, only one choice is possible) A list of parameters obtained (**Table 4**).

**Enter rangelist of parameters: 16** [press ENTER] (users can enter here more parameter numbers as required, but one space must be between two parameter numbers. Pressing ENTER will auto-load all the these parameters)

**10 values added for S-P+-1** (The values of  $\sigma^+$  for the substituents at 4-positions are auto-loaded) Now we need to edit **S-P+-1** to **S+**, which will be possible as follows:

**qsar> editp 11** (parameter level 11) [press ENTER] A box of parameter label will obtained, one should edit S-P+-1 to S+ by simply the use of BACKSPACE.

**qsar> save** [press ENTER] To save the data.

**(vi) Plotting Data:** Any two parameters can be plotted against each other by the command:

**qsar> x gr y** , i.e. **3 gr 4** [press ENTER] gives some idea about the nature of fit to the most important variables that is linear, parabolic, or bilinear. In the present instance this is of little help.

**(vii) Permuting:**

**qsar> 3 perm 4 5 6 7 8 10 11** [press ENTER] derives all possible equations for 1, 2, and 3 variables. S.D. is the standard deviation. As this value decreases, the quality of the fit increases. CONST is the value of the constant in each QSAR equation. Clearly, BDE is the most important parameter and then  $\sigma^+$  and Clog *P*.

### 1 TERM REGRESSIONS

	S.D.	BDE	Clog <i>P</i>	CMR	NVE	MgVol	MW	S+	CONST
1	.437	-.15							3.66
2	.457							-1.09	3.66
3	.525				.02				3.12
4	.531						.01		3.12
5	.538					.76			3.21
6	.550			.21					3.24
7	.594		.14						3.79

## 2 TERM REGRESSIONS

	S.D.	BDE	Clog <i>P</i>	CMR	NVE	MgVol	MW	S+	CONST
1	.142	-.20	.27						2.74
2	.176		.28					-1.52	2.69
3	.205	-.17				.89			2.50
4	.218	-.16			.02				2.51
5	.234	-.17		.25					2.51
6	.238	-.16					.01		2.51
7	.241					.90		-1.21	2.47
8	.249				.02			-1.16	2.48
9	.268						.01	-1.16	2.49
10	.269			.25				-1.22	2.49

## 3 TERM REGRESSIONS

	S.D.	BDE	Clog <i>P</i>	CMR	NVE	MgVol	MW	S+	CONST
1	.130	-.49	.25					2.16	2.84
2	.147	-.73				.84		4.23	2.65
3	.150	-.20	.24				.0		2.68
4	.150	-.20	.23		.0				2.68
5	.152	-.20	.24			.10			2.70
6	.152	-.20	.25	.02					2.70
7	.156	-.16		-.97		4.23			2.64
8	.170	-.16			.12		-.04		2.70
9	.171	-.71			.02			4.13	2.66
10	.171	-.79		.23				4.72	2.66

### (viii) Checking for Parameter Collinearity and derivation of QSAR:

We can check for overall collinearity problems as follows:

**qsar> corr 4 5 6 7 8 10 11** [press ENTER]

## CORRELATION MATRIX: $R^2$ and N

	BDE	Clog $P$	CMR	NVE	MgVol	MW	S+
BDE	.	.125	.015	.004	.012	.057	.996
Clog $P$	10	.	.756	.743	.901	.693	.379
CMR	10	10	.	.996	.991	.988	.131
NVE	10	10	10	.	.990	.995	.073
MgVol	10	10	10	10	.	.975	.123
MW	10	10	10	10	10	.	.005
S+	10	10	10	10	10	10	.

There is a high collinearity between BDE and  $\sigma^+$  and also the high collinearity among Clog  $P$ , CMR, NVE, MgVol, and MW. The upper part of the matrix gives correlation among the 10 different substituents. Next exploring the two best terms, the following QSAR can be derive:

**qsar> 3 reg 5 4** [press ENTER] will give the following equation:

$$\log 1/C = 0.27(\pm 0.08)\text{Clog } P - 0.20(\pm 0.04)\text{BDE} + 2.74(\pm 0.31) \quad (3)$$

$$n = 10, \quad r = 0.978, \quad q^2 = 0.894, \quad \text{SS1} = 3.236, \quad \text{DEV}^+ = 5$$

$$\text{DF} = 7, \quad r^2 = 0.956, \quad s = 0.142, \quad \text{SS2} = 0.141, \quad \text{DEV}^- = 5$$

Since, there is a high collinearity between BDE and  $\sigma^+$ , BDE can be replaced by  $\sigma^+$ , thus

**qsar> 3 reg 5 11** [press ENTER] will give the following equation:

$$\log 1/C = 0.28(\pm 0.10)\text{Clog } P - 1.52(\pm 0.39)\sigma^+ + 2.69(\pm 0.40) \quad (4)$$

$$n = 10, \quad r = 0.966, \quad q^2 = 0.807, \quad \text{SS1} = 3.236, \quad \text{DEV}^+ = 5$$

$$\text{DF} = 7, \quad r^2 = 0.933, \quad s = 0.176, \quad \text{SS2} = 0.217, \quad \text{DEV}^- = 5$$

**qsar> pred** [press ENTER] will list the data in tabular form:

No.	Substituent	log1/C	Ypred	Dev	Clog P	S+	BDE
1	NH <sub>2</sub>	4.61	4.74	-0.13	0.25	-1.30	-9.25
2	OC <sub>6</sub> H <sub>13</sub>	5.34	5.10	0.24	4.22	-0.81	-6.30
3	Me	3.82	3.71	0.11	1.97	-0.31	-2.22
4	CMe <sub>3</sub>	3.86	4.01	-0.15	3.30	-0.26	-1.54
5	H	3.27	3.10	0.17	1.48	0.00	0.00
6	OMe	4.41	4.32	0.09	1.57	-0.78	-6.01
7	C <sub>3</sub> H <sub>7</sub>	3.72	3.98	-0.26	3.03	-0.29	-2.01
8	C(Me) <sub>2</sub> C <sub>6</sub> H <sub>4</sub> -4-OH	4.18	4.15	0.03	3.67	-0.29	-1.88
9	C <sub>2</sub> H <sub>5</sub>	3.79	3.84	-0.05	2.50	-0.30	-1.90
10	C <sub>8</sub> H <sub>17</sub>	4.66	4.71	-0.05	5.68	-0.29	-2.17

? **TROUBLESHOOTING:** The following commands are used for the derivation of different types of QSAR in the MLR analysis:

(i) Linear Model

**qsar> x reg y** [press ENTER] (**x** and **y** are the parameter numbers) i.e.

**qsar> 3 reg 4** [press ENTER]

The QSAR may be obtained by the use of more than one parameter i.e.

**qsar> 3 reg 4 5 6 7 8** [press ENTER], etc

(ii) Parabolic / Inverted Parabolic Model

**qsar> x reg Py** [press ENTER] (**x** and **y** are the parameter numbers, and **P** is the command for parabolic derivation of the QSAR), i.e.,

**qsar> 3 reg P4** [press ENTER] The QSAR for more than one parameter can be derive as:

**qsar> x reg Py z** [press ENTER]

(iii) Bilinear / Inverted Bilinear Model

**qsar> x reg By** [press ENTER] (**x** and **y** are the parameter numbers, and **B** is the command for bilinear derivation of the QSAR), i.e.,

**qsar> 3 reg B4** [press ENTER] The QSAR for more than one parameter can be derive  
as:

**qsar> x reg By z** [press ENTER]

**(ix) Jackknifing:** This process is used for the detection and removal of the outliers. Since the present example (QSAR 3 or 4) is the good model and not have any outlier, jackknifing will not be used. To understand this process, we consider the data set B11019.

**qsar> load /d 11019** [press ENTER] will transfer the set to the workspace

**qsar> star /d** [press ENTER] Asterisks will be removed to restore all the data points for study.

**qsar> sum** [press ENTER] will list the key items of the data set as follows:

Dataset name: BIO\_11019

Substituents: 23 Parameters: 9 SMILES: 23 Equations: 1

Active: 23 Starred: 0 Inactive: 0

System: HT-1080 CELLS

Class: B4C ; Cells in culture

Compound: CAFFEIC ACID ESTERS

Action: EC50: ANTI-PROLIFERATIVE ACTIVITIES OF CAFFEIC ACID  
ESTERS

Reference: NAGAOKA,T., BANSKOTA,A.H., TEZUKA,Y., SAIKI,I.,  
KADOTA,S., BIOORG.MED. CHEM., 10,3351-3359(2002) R6979

Source: RAJESHWAR PRASAD VERMA

Analysis: UNKNOWN

Date: 2006 October 26

Parameters: YPRED DEV LOG1/C CLOGP CMR NVE MGVOL BILIN(CLOGP)  
CLOGP\*\*2

qsar> **3 reg B4** [press ENTER] will give the following equation:

$$\log 1/C = 0.14(\pm 0.06)\text{Clog } P - 0.33(\pm 0.11)\log (\beta \times 10^{\text{Clog } P} + 1) + 4.32(\pm 0.23)$$

(5)

$$n = 23, \quad r^2 = 0.719, \quad s = 0.132, \quad q^2 = 0.617$$

$$\text{Optimum Clog } P = 5.04 \quad \log \beta = -5.17$$

This is obviously not a good equation. We can now use jackknifing by the following command:

qsar> **3 j B4** [press ENTER] will derive all possible regression equations by dropping a different data point in each instance.

	Omitted	$r^2$	$s$
None		0.719	0.132
19	R=(CH <sub>2</sub> ) <sub>7</sub> CH <sub>3</sub>	0.782	0.119
22	R=(CH <sub>2</sub> ) <sub>13</sub> CH <sub>3</sub>	0.759	0.126
1	R=CH <sub>2</sub> Ph	0.739	0.129
7	R=(CH <sub>2</sub> ) <sub>8</sub> Ph	0.737	0.131
17	R=(CH <sub>2</sub> ) <sub>2</sub> CH <sub>3</sub>	0.73	0.133
3	R=(CH <sub>2</sub> ) <sub>3</sub> Ph	0.728	0.134
23	R=(CH <sub>2</sub> ) <sub>15</sub> CH <sub>3</sub>	0.723	0.13
9	R=CH <sub>2</sub> CHCHPh	0.722	0.135
21	R=(CH <sub>2</sub> ) <sub>11</sub> CH <sub>3</sub>	0.719	0.136
18	R=(CH <sub>2</sub> ) <sub>3</sub> CH <sub>3</sub>	0.719	0.136

Dropping compound #19 yields  $r^2 = 0.782$ . To delete this compound, the following command will be use:

**qsar> star /a 19** [press ENTER] This will place an asterisk on the data point #19 and it will not be used in deriving future equations.

**qsar> 3 reg B4** [press ENTER] will give the following equation:

$$\log 1/C = 0.16(\pm 0.06)\text{Clog } P - 0.36(\pm 0.10)\log (\beta \times 10^{\text{Clog } P} + 1) + 4.27(\pm 0.21)$$

(6)

$$n = 22, \quad r^2 = 0.782, \quad s = 0.119, \quad q^2 = 0.696$$

$$\text{Optimum Clog } P = 4.98 \quad \log \beta = -5.07$$

**qsar> 3 i B4** [press ENTER] will derive all possible regression equations by dropping a different data point in each instance.

	Omitted	$r^2$	$s$
None		0.782	0.119
22	R=(CH <sub>2</sub> ) <sub>13</sub> CH <sub>3</sub>	0.822	0.111
7	R=(CH <sub>2</sub> ) <sub>8</sub> Ph	0.804	0.117
1	R=CH <sub>2</sub> Ph	0.800	0.116
3	R=(CH <sub>2</sub> ) <sub>3</sub> Ph	0.797	0.119
23	R=(CH <sub>2</sub> ) <sub>15</sub> CH <sub>3</sub>	0.794	0.115
17	R=(CH <sub>2</sub> ) <sub>2</sub> CH <sub>3</sub>	0.793	0.120
9	R=CH <sub>2</sub> CHCHPh	0.789	0.121
8	R=(CH <sub>2</sub> ) <sub>12</sub> Ph	0.783	0.107
18	R=(CH <sub>2</sub> ) <sub>3</sub> CH <sub>3</sub>	0.783	0.123
20	R=(CH <sub>2</sub> ) <sub>9</sub> CH <sub>3</sub>	0.782	0.122

Dropping compound #22 yields  $r^2 = 0.822$ .

**qsar> star /a 22** [press ENTER] This will place an asterisk on the data point #22 and it will not be used in deriving future equations.

**qsar> 3 reg B4** [press ENTER] will give the following equation:

$$\log 1/C = 0.16(\pm 0.05)\text{Clog } P - 0.38(\pm 0.10)\log (\beta \times 10^{\text{Clog } P} + 1) + 4.27(\pm 0.20) \quad (7)$$

$$n = 21, \quad r^2 = 0.822, \quad s = 0.111, \quad q^2 = 0.729$$

$$\text{Optimum Clog } P = 4.97 \quad \log \beta = -5.09$$

**qsar> 3 i B4** [press ENTER] will derive all possible regression equations by dropping a different data point in each instance.

	Omitted	$r^2$	$s$
None		0.822	0.111
23	R=(CH <sub>2</sub> ) <sub>15</sub> CH <sub>3</sub>	0.865	0.096
1	R=CH <sub>2</sub> Ph	0.842	0.106
7	R=(CH <sub>2</sub> ) <sub>8</sub> Ph	0.841	0.108
3	R=(CH <sub>2</sub> ) <sub>3</sub> Ph	0.837	0.109
17	R=(CH <sub>2</sub> ) <sub>2</sub> CH <sub>3</sub>	0.834	0.111
9	R=CH <sub>2</sub> CHCHPh	0.829	0.112
21	R=(CH <sub>2</sub> ) <sub>11</sub> CH <sub>3</sub>	0.823	0.114
18	R=(CH <sub>2</sub> ) <sub>3</sub> CH <sub>3</sub>	0.823	0.114
4	R=(CH <sub>2</sub> ) <sub>4</sub> Ph	0.823	0.113
20	R=(CH <sub>2</sub> ) <sub>9</sub> CH <sub>3</sub>	0.822	0.114

Dropping compound #23 yields  $r^2 = 0.865$ .

**qsar> star /a 23** [press ENTER] This will place an asterisk on the data point #23 and it will not be used in deriving future equations.

**qsar> 3 reg B4** [press ENTER] will give the following equation:

$$\log 1/C = 0.16(\pm 0.05)\text{Clog } P - 0.43(\pm 0.09)\log (\beta \times 10^{\text{Clog } P} + 1) + 4.28(\pm 0.17)$$

(8)

$$n = 20, \quad r = 0.930, \quad q^2 = 0.819, \quad \text{SS1} = 1.089, \quad \text{DEV}^+ = 9$$

$$\text{DF} = 16, \quad r^2 = 0.865, \quad s = 0.096, \quad \text{SS2} = 0.147, \quad \text{DEV}^- = 11$$

$$\text{Optimum Clog } P = 5.10 \quad \log \beta = -5.33$$

Note that  $q^2$  is now approaching  $r^2$  in value. This implies that dropping another point would have little effect on the correlation. Thus, this equation (8) can be considered as the final QSAR model for this data set.

**qsar> pred** [press ENTER] will list the data in tabular form.

? **TROUBLESHOOTING:** Any number of data points can be omitted from this fashion:

**qsar> star /a 1 4 8 10,15** [press ENTER] would place asterisks on the data points # 1, 4, 8, and 10-15. An asterisk can be removed by the following command:

**qsar> star /d 15** [press ENTER] will remove asterisk and restore the data point #15

**qsar> star /d 10,15** [press ENTER] will remove asterisks and restore the data point #10-15

**qsar> star /d** [press ENTER] will remove all asterisks from the data set and restore all the data points.

(x) **Editing:** This is the important tool to correct the errors in the data set.

**qsar> editset** [press ENTER] This displays all the data and puts one into a general edit mode allowing: (1) movement of the cursor by the arrow key, (2) the use of the delete key to erase the character to the left of the cursor, and (3) the insertion of new characters at that spot. A new variable can also be added by first assigning it a symbol and then entering the values. When finished, one must exit from this mode, using **control Z**. For the minor change a quick mode is also available utilizing the following command:

**qsar> editdata** [press ENTER] will provide the prompt for substituent number. Entering that number and pressing ENTER will give the prompt for parameter number and the current value will display. Entering the parameter number and pressing ENTER will give a box for parameter value. The correct value of that parameter should be entered and then presses ENTER. After editing is complete, the command **seedata** enables one to check the results.

(xi) **Deleting:** It is convenient to use the delete command to clean up the set. However, parameters cannot be deleted if an equation has been saved. To check for saved equations:

**qsar> eq run** [press ENTER] will list all the saved equations for the present data set.

**qsar> del /eq** [press ENTER] will delete all the saved equations from the present data set.

**qsar> del /eq 2** [press ENTER] will delete only the 2<sup>nd</sup> equation from the present data set.

**qsar> seep** [press ENTER] will display the parameter numbers

**qsar> del /para 4** [press ENTER] will delete the parameter 4.

**qsar> del /para 4 8 10,16** [press ENTER] will delete the parameters 4, 8, and 10-16.

To delete a data point and the information associated with it including the SMILES, the following command should be used.

**qsar> del /sub 4** [press ENTER] will delete the data point #4 (compound #4). Similarly, more data points can be deleted.

## **D. VALIDATION OF QSAR MODEL (EXAMPLE: QSAR 8; SET NO. B11019)**

### **1. Statistical Diagnostics:**

(i) Number of descriptors: The ideal ratio = data points (compounds)/descriptor  $\geq 4$ .

$$\text{Number of data points/number of descriptors} = 20/3 = 6.67$$

(ii) Squared correlation coefficient ( $r^2$ ): Closer the value of  $r^2$  to unity, the better the QSAR model. According to the literature, the predictive QSAR model must have  $r^2 > 0.60$ .<sup>14</sup>

$$r^2 = 0.865$$

- (iii) Standard Deviation ( $s$ ): It is believed that the smaller the value of  $s$  ( $s \leq 0.3$ ), the better the QSAR model.<sup>15</sup>

$$s = 0.096$$

## 2. Internal Validation:

- (i) Cross-validated  $r^2$  ( $q^2$ ): It has been suggested that the value of  $q^2$  must be greater than 0.50 for a predictive QSAR model.<sup>14</sup>

$$q^2 = 0.819$$

- (ii) Quality factor ( $Q$ ): The  $Q$  is the quality factor (quality ratio), where  $Q = r/s$ .

$$Q = r/s = 0.930/0.096 = 9.688$$

- (iii) Fischer statistics ( $F$ ): The  $F$ -value (Fischer ratio) is the ratio between explained and unexplained variance for a given number of degree of freedom. It indicates a true relationship, or the significance level for the MLR models.

$F_{m, f, 0.05} = fr^2/[(1-r^2)m]$ , where  $f$  is the number of degree of freedom,  $f = n-(m+1)$ ,  $n$  = number of data points, and  $m$  = number of variables.

$$F_{3, 16, 0.05} = fr^2/[(1-r^2)m] = (16 \times 0.865)/[(1-0.865) \times 3] = 13.840/0.405 = 34.173$$

Alternatively, it can also be calculated as:

$$F_{x, DF, 0.05} = [(SS1-SS2) \times DF]/(SS2 \times x)$$

$$F_{3, 16, 0.05} = [(1.089-0.147) \times 16]/(0.147 \times 3) = 15.072/0.441 = 34.177$$

Where, DF is the number of degree of freedom and 'x' represents the number of variables. SS1 is the sum of squares about the mean of the dependent variable and SS2 is the sum of squares from the deviations from the regression line. 0.05 represents the 95% of the confidence limit.

(iv) Y-randomization Test: At present, we not have the facility to do this test automatically from this program, but it can be done manually. In this test, the dependent-variable vector (Y-vector) is randomly shuffled and a new QSAR model is developed using the original independent variable matrix. This process is repeated several times. It is expected that the resulting QSAR models should have low  $r^2$  and low  $q^2$  values. Alternatively, this test can be done automatically from our Bio-Loom program (<http://www.bio-loom.com> or <http://www.biobyte.com>) for any data set saved in C-QSAR program. For the present example the Y-randomization test results obtained from the Bio-Loom program are as follows:

Y Scramble Results	
s.d.: 0.243	$r^2$ : 0.130
s.d.: 0.228	$r^2$ : 0.232
s.d.: 0.255	$r^2$ : 0.044
s.d.: 0.222	$r^2$ : 0.275
s.d.: 0.258	$r^2$ : 0.022



### 3. External Validation:

It is a better method that removing a percentage of the training set into a test set. The QSAR model is derived using the reduced training set, and the properties of the test set predicted using this model. Following are the methods for the selection of training and test sets:

- (i) Random selection
- (ii) Selection based on biological activity data
- (iii) Various systematic clustering techniques
- (iv) Self-organizing map (SOM)
- (v) Kennard Stone method

- (vi) Formal statistical experimental design (factorial and D-Optimal)
- (vii) Sphere-exclusion algorithms

## REFERENCES:

1. C-QSAR Program, BioByte Corp., 201W. 4<sup>th</sup> st., Suit 204, Claremont, CA 91711, USA. [www.biobyte.com](http://www.biobyte.com)
2. Hansch, C., Hoekman, D., Leo, A., Weininger, D. & Selassie, C.D. Chem-bioinformatics: Comparative QSAR at the interface between chemistry and biology. *Chem. Rev.* **102**, 783-812 (2002).
3. Hansch, C. Hoekman, D. & Gao, H. Comparative QSAR: Toward a deeper understanding of chemicobiological interactions. *Chem. Rev* **96**, 1045-1075 (1996).
4. Hansch, C. & Leo, A. Exploring QSAR: Fundamentals and Applications in Chemistry and Biology, American Chemical Society, Washington, D.C. (1995).
5. Hansch, C. Leo, A. & Taft, R.W. A survey of Hammett substituent constants and resonance and field parameters. *Chem. Rev.* **91**,165-195 (1991).
6. Schultz, T.W. & Cronin, M.T.D. Response-Surface Analyses for Toxicity to *Tetrahymena pyriformis*: Reactive Carbonyl-Containing Aliphatic Chemicals. *J. Chem. Inf. Comput. Sci.* **39**, 304-309 (1999).
7. Zhang, L., Gao, H., Hansch, C. & Selassie, C.D. Molecular orbital parameters and comparative QSAR in the analysis of phenol toxicity to leukemia cells. *J. Chem. Soc. Perkin Trans 2*, 2553-2556 (1998).
8. Hu, J. Eriksson, L., Bergman, A., Jakobsson, E., Kolehmainen, E., Knuutinen, J., Suontamo, R. & Wei, X. Molecular orbital studies on brominated diphenyl ethers. Part II-reactivity and quantitative structure-activity (property) relationships. *Chemosphere* **59**,1043-1057(2005).
9. Hansch, C., Maloney, P.P., Fujita, T. & Muir, R.M. Correlation of biological activity of phenoxyacetic acids with Hammett substituent constants and partition coefficients. *Nature* **194**, 178-180 (1962).
10. Leo, A.J. Calculating log Poct from structures. *Chem. Rev.* **93**, 1281-306 (1993).
11. Leo, A.J. & Hansch, C. Role of hydrophobic effects in mechanistic QSAR. *Perspect. Drug Discovery Des.* **17**, 1-25 (1999)
12. Cramer III, R.D., Bunce, J.D., Patterson, D.E. & Frank, I.E. Cross validation, Bootstrapping and partial least squares compared with multiple regression in conventional QSAR studies. *Quant. Struct.-Act. Relat.* **7**, 18-25 (1988).
13. Weininger, D. SMILES, a chemical language and information system.1. Introduction to methodology and encoding rules. *J. Chem. Inf. Comput. Sci.* **28**, 31-36 (1988).
14. Golbraikh, A. & Tropsha, A. Beware of q<sup>2</sup>! *J. Mol. Graph. Modl.* **20**, 269-276 (2002).
15. Wold, S. Validation of QSAR's. *Quant. Struct.-Act. Relat.* **10**, 191-193 (1991).

**Table 1:** Available QSAR database as well as QSAR and molecular descriptors programs

No.	Software Name	Brief Description	URL/Reference
1	C-QSAR	Development of QSAR model and a database of over 21,000 QSARs of which 12,950+ pertain to biological systems and 8,900+ are from mechanistic organic chemistry	<b>Corwin Hansch</b> , Professor Emeritus, Department of Chemistry, Pomona College, Claremont, CA 91711 <a href="http://www.biobyte.com/bb/prod/cqsarad.html">http://www.biobyte.com/bb/prod/cqsarad.html</a>
2	Bio-Loom	To access BioByte's entire Masterfile database, which includes over 60,000 measured log <i>P</i> and log <i>D</i> values (in many solvent systems), as well as 14,000 pKas, including associated references. To access entire C-QSAR database and easy searching online. Still calculating hydrophobic and molecular refractivity parameter via Clog <i>P</i> and CMR calculations.	<a href="http://www.biobyte.com/bb/prod/bioloom.html">http://www.biobyte.com/bb/prod/bioloom.html</a>
3	Pharma Algorithm's QSAR Builder	QSAR/QSPR modeling and calculation of some descriptors	<b>Hugo Kubinyi</b> , Professor of Pharmaceutical Chemistry at the University of Heidelberg, Germany <a href="http://www.ap-algorithms.com">http://www.ap-algorithms.com</a>
4	Distill, HQSAR, Almond, Molconn-Z	Distill: Structure-activity relationship determination HQSAR: Automated QSAR analysis Almond: alignment of independent molecular descriptors Molconn-Z: QSAR/QSPR model descriptors	<a href="http://www.tripos.com/index.php?family=modules,SimplePage,sybyl_ligand_based_design">http://www.tripos.com/index.php?family=modules,SimplePage,sybyl_ligand_based_design</a>
5	ChemSAR	ChemSAR is a Chem3D Windows add-in for MS Excel with descriptive statistics and plots for structure-activity relationships.	<a href="http://www.mpassociates.gr/software/distrib/science/cambsoft/ChemOfficePro2005.html">http://www.mpassociates.gr/software/distrib/science/cambsoft/ChemOfficePro2005.html</a>

6	Partek QSAR Solution	Data mining and QSAR	<a href="http://www.partek.com/html/products/products.html">http://www.partek.com/html/products/products.html</a>
7	MCASE/MC4PC	QSAR development	<a href="http://www.multicase.com/products/prod01.htm">http://www.multicase.com/products/prod01.htm</a>
8	FieldAlign	Align molecules in field space for SAR and QSAR	<a href="http://www.cresset-bmd.com/productindex.html">http://www.cresset-bmd.com/productindex.html</a>
9	MDL® QSAR	A comprehensive QSAR modeling system	<a href="http://www.mdl.com/products/predictive/qsar/index.jsp">http://www.mdl.com/products/predictive/qsar/index.jsp</a>
10	Strike	Strike (Statistical tool for revealing insight and knowledge): Software for statistical modeling and QSAR	<a href="http://www.schrodinger.com/ProductDescription.php?mID=6&amp;sID=17">http://www.schrodinger.com/ProductDescription.php?mID=6&amp;sID=17</a>
11	TerraQSAR	QSAR and Terra Tox database	<a href="http://www.terrabase-inc.com">http://www.terrabase-inc.com</a>
12	Cerius <sup>2</sup>	QSAR modeling, includes ADME, docking and scoring	<a href="http://www.accelrys.com">http://www.accelrys.com</a>
13	Bioreason ClassPharmer	QSAR/QPSR modeling	<a href="http://www.bioreason.com">http://www.bioreason.com</a>
14	Cheminformatics (Chemical computing Group Inc.)	QSAR/QSPR, HTS-QSAR, and Binary QSAR modeling, ADME assessments, and calculation of over 300 molecular descriptors.	<a href="http://www.chemcomp.com/software-chem.htm">http://www.chemcomp.com/software-chem.htm</a>
15	CODESSA	QSAR program	<a href="http://www.semichem.com/codessa/default.php">http://www.semichem.com/codessa/default.php</a>
16	AMBIT	QSAR database search and QSAR model development	<a href="http://ambit.acad.bg">http://ambit.acad.bg</a>
17	HASL	3D QSAR	<a href="http://www.bio.com/store/product.jhtml?id=prod300024">http://www.bio.com/store/product.jhtml?id=prod300024</a>
18	GOLPE	3D QSAR	<a href="http://www.miasrl.com/golpe.htm">http://www.miasrl.com/golpe.htm</a>
19	Biograf <sup>3R</sup>	Multi-dimentional QSAR tools Quasar	<a href="http://www.biograf.ch/index.php">http://www.biograf.ch/index.php</a>
20	DRAGON	Calculation of 1664 molecular descriptors	<a href="http://www.taletе.mi.it/main_net.htm">http://www.taletе.mi.it/main_net.htm</a>
21	Molinspiration Toolkit	Java based software and free online calculations of fragments and basic properties/descriptors	<a href="http://www.molinspiration.com">http://www.molinspiration.com</a>

22	PETRA	A free online service for calculating various physicochemical properties e.g. heats of formation, bond dissociation energies, sigma charge distribution, pi charge distribution, inductive effect, resonance effect, delocalization energies, and polarizability etc.	<a href="http://www2.ccc.uni-erlangen.de/software/petra">http://www2.ccc.uni-erlangen.de/software/petra</a>
23	Molecule Evuator and Polar Surface Area	Molecule Evuator: Automatic calculation of physicochemical properties of molecules Polar Surface Area: A free software for the calculation of polar surface area (PSA) of the molecules	<a href="http://www.cidrux.com">http://www.cidrux.com</a>
24	Spartan'06	Calculation of HOMO and LUMO Energies, Polar Surface Area, Electronegativity, Hardness, Dipole, etc	<a href="http://www.wavefunction.com">http://www.wavefunction.com</a>
25	Topological Polar Surface Area (TPSA)	A free online calculation of TPSA for the organic molecules	<a href="http://www.daylight.com/meetings/emug00/Ertl/tpsa.html">http://www.daylight.com/meetings/emug00/Ertl/tpsa.html</a>

**Table 2:** The information associated with each data set, which is either in biological or physical database

Field	Title	Description
Input Data		
1	SYSTEM	Biological/physical system
2	CLASS	Classification of the system (Tables 5 and 6)
3	COMPOUND	Parent compound (if any)
4	ACTION	Measured action/activity
5	REFERENCE	Journal reference/other source of data set
6	SOURCE	The person who entered the data set
7	ANALYSIS	The person who analyzed the data set
8	NOTE	Additional information for the data set
9	DATE	The date on which set was saved into database
10	PARAMETERS	List of the parameters <sup>a</sup>
11	SUBSTITUENTS	Labels of the substituents
12	SMILES	The topological description of the compounds
13	DATA	Table of the parameter values
14	PRM MAX/MIN	The maximum and minimum of each parameter
Output Data (Equation)		
15	TERMS IN EQN	Parameters in the regression equation
16	COEFS IN EQN	The regression coefficients for each of the parameter
17	IDEAL/LOG $\beta$	An ideal (or optimal) log $P$ , and their confidence limits
18	STATISTICS	$n$ , $r$ , $r^2$ , $q^2$ , $s$ , $DF$ , $SS1$ , $SS2$ , $DEV+$ , $DEV-$
19	RESIDUALS	Deviations between $y$ -predicted and the log of the observed value
20	PREDICTED	Predicted values of the dependent parameter

<sup>a</sup>Examined, even if not used in the final equation

**Table 3.** Summary of the Biological and Physical Database in C-QSAR program

No.	Summary	Biological Database	Physical Database	TOTAL
1	Sets (non-empty)	12958	8934	21892
2	Last dataset	12958	8934	21892
3	Equations	12958	8934	21892
4	Total compounds	203503	89652	293155
5	Total SMILES	203090	88899	291989
6	Titles: system	7019	1605	8624
7	Titles: classification	213	112	325
8	Titles: compound	6201	4475	10676
9	Titles: action	5745	4995	10740
10	Titles: reference	6303	5165	11468

**Table 4.** Physicochemical parameters of the substituents from C-QSAR program

No.	Parameter	Description	No.	Parameter	Description
0	CPI	Calculated $\pi$ (relative to the parent)	22	S-O-	Sigma ortho minus
1	PI	$\pi$ (measured hydrophobic parameter for the substituents)	23	S-INDUC	Sigma inductive
2	MR-SUB	Substituent molar refractivity	24	S-AN-RS	Sigma resonance, anilines
3	F	Field effect (from S-L)	25	S-RES+	Sigma resonance plus
4	R	Resonance effect (from S-L)	26	S-'	Sigma prime
5	R+	Resonance plus	27	S-PARNO	Sigma para normalized
6	R-	Resonance minus	28	S-ORTH+	Sigma ortho plus
7	ES	$E_S$ from Taft	29	S-PHOSP	Sigma phosphoric acid
8	ES-HYBO	$E_S$ from Hydroboration	30	S-L	Sigma localized (Charton)
9	ES-V	$E_S$ from Charton	31	S-ZTWST	Sigma orthogonal twist
10	ES-A	$E_S$ from Austel	32	S-STAR	Sigma star from Taft
11	L	Length sterimol	33	S-IND.P	Sigma inductive (Phosphorus)
12	B1	Width sterimol (minimum width of the substituents)	34	S-RES.P	Sigma resonance (Phosphorus)
13	B5	Width sterimol (maximum width of the substituents)	35	ER-P	Electronic radical, para
14	O-STER	Ortho quats with MeI	36	ER-M	Electronic radical, meta
15	S-P	Sigma para	37	S.DOT-P	Sigma dot, para
16	S-P+	Sigma para plus	38	S.DOT-M	Sigma dot, meta
17	S-P-	Sigma para minus	39	S.-DOT-P	Sigma dot, para (JJ)
18	S-M	Sigma meta	40	S.-DOT-M	Sigma dot, meta (JJ)
19	S-M+	Sigma meta plus	41	S.P-C	Sigma para (C)
20	S-M-	Sigma meta minus	42	S.M-C	Sigma meta (C)
21	S-O	Sigma ortho	43	CMR-SUB	Calculated MR for the substituents (relative to the parent)

## SUPPLEMENTARY INFORMATION

### **Classification**

The classification of C-QSAR program into its biological and physical databases has been shown in **Table S1** and **Table S2** respectively. In biological database, there are six classes that are further divided into 34 sub-classes (**Table S1**). On the other hand, the physical database contains 24 classes and 15 sub-classes (**Table S2**). Most important parameters for biological QSAR are hydrophobic, electronic and steric whereas the same for physical QSAR are electronic and steric.

### **Physicochemical Parameters:**

There are six physicochemical parameters (molecular descriptors) for the whole molecule (Clog *P*, Mlog *P*, CMR, NVE, MgVol, and MW) and forty-four for the substituents (**Table 4**) that can be auto-loaded in the system for deriving QSAR models. The other physicochemical parameters calculated from different software can also be used by this program in the development of QSAR models. A brief definition for the physicochemical parameters available in C-QSAR program is as follows:

(a) *Physicochemical parameters for the whole molecule*: Clog *P* is the calculated partition coefficient in *n*-octanol/water and is a measure of hydrophobicity of the whole molecule while Mlog *P* is the measured partition coefficient.<sup>4,16</sup> There are a number of methods for the calculation of log *P*.<sup>10</sup> The most extensively used method is that of Leo.<sup>10,11</sup> The quality of his method is illustrated by eq. S1.<sup>17</sup>

$$\text{Mlog } P = 0.96(\pm 0.003) \text{ Clog } P + 0.07(\pm 0.008) \quad (\text{S1})$$

$$n = 12510, \quad r^2 = 0.973, \quad s = 0.300, \quad q^2 = 0.973$$

CMR is the calculated molar refractivity for the whole molecule. MR is calculated from the Lorentz-Lorenz equation and is described as follows:  $[(n^2-1)/(n^2+2)](MW/\delta)$ , where  $n$  is the refractive index, MW is the molecular weight, and  $\delta$  is the density of the substance. MR is dependent on volume and polarizability. It can be used for a substituent or for the whole molecule. A new polarizability parameter, NVE, was developed, which is shown to be effective at delineating various chemico-biological interactions.<sup>18-21</sup> NVE represents the total number of valence electrons and is calculated by simply summing up the valence electrons in a molecule, that is, H = 1, C = 4, Si = 4, N = 5, P = 5, O = 6, S = 6 and halogens = 7. It may also be represented as:  $NVE = n_\sigma + n_\pi + n_n$ , where  $n_\sigma$  is the number of electrons in  $\sigma$ -orbital,  $n_\pi$  is the number of electrons in  $\pi$ -orbitals, and  $n_n$  is the number of lone pair electrons. MgVol is the molar volume for the whole molecule,<sup>22</sup> and MW is their molecular weight.

(b) *Physicochemical parameters for the substituents*: At present, the C-QSAR program contains forty-four physicochemical parameters for the substituents (**Table 4**) that can be auto-loaded in the system for the derivation of QSAR models. The more molecular descriptors should be added very soon. A brief definition for these physicochemical parameters is as follows:

0. CPI. Calculated  $\pi$  (relative to parent)
1. PI ( $\pi$ ). Hydrophobic parameter for the substituents defined by partitioning of X-C<sub>6</sub>H<sub>5</sub> between *n*-octanol and water.<sup>4</sup>  

$$\pi_X = \log P_{X-C_6H_5} - \log P_{C_6H_6}$$
2. MR-SUB. Molar refractivity of a substituent defined analogously to  $\pi$ .  

$$MR = [(n^2-1)/(n^2+2)](MW/\delta)$$

- Where  $n$  = refractive index, MW = molecular weight, and  $\delta$  = density. MR values are scaled by 0.1 and it is highly collinear to substituent volume.<sup>4</sup>
3. F. Swain-Lupton inductive/field effect parameter for aromatic systems.<sup>5,23</sup>
  4. R. Corresponding Swain-Lupton resonance parameter.<sup>5,23</sup>
  5. R+. Taft resonance parameter for substituent delocalization of a +ve charge.<sup>5</sup>
  6. R-. Taft resonance parameter for substituent delocalization of a -ve charge.<sup>5</sup>
  7. ES. Classic steric parameter for substituents ( $E_S$ ) defined by Taft.<sup>24</sup>
  8. ES-HYBO. An  $E_S$  type parameter obtained from the hydroboration of substituted ethylenes.
  9. ES-V. Charton's steric parameter.<sup>25</sup>
  10. ES-A. Austel version of steric parameter.<sup>26</sup>
  11. *L*. Verloop sterimol parameter for substituent length.<sup>4,27</sup>
  12. *B1*. Sterimol parameter for the minimum width of the substituent.<sup>4,27</sup>
  13. *B5*. An estimate of the maximum width of the substituent.<sup>4,27</sup>
  14. O-STER. There are 30 substituents having this parameter for the effect of adjacent substituents inhibiting the reaction of pyridines with  $\text{CH}_3\text{I}$ .<sup>28</sup>
  15. S-P ( $\sigma_p$ ). Normal Hammett constant for para substituents. It is based on the ionization constants of benzoic acids.<sup>4</sup>
  16. S-P+ ( $\sigma_p^+$ ). Brown parameter where substituents delocalize a +ve charge or radical via resonance.<sup>4,29</sup>
  17. S-P- ( $\sigma_p^-$ ). Hammett constant where substituents delocalize a -ve charge via resonance. It is derived from the ionization constants of phenol.<sup>4</sup>
  18. S-M ( $\sigma_m$ ). Hammett constant for meta substituents (non conjugate substituents).<sup>4</sup>

19. S-M+ ( $\sigma_m^+$ ). Brown parameter for meta substituents. There is little difference between  $\sigma_m$  and  $\sigma_m^+$ .<sup>4</sup>
20. S-M- ( $\sigma_m^-$ ). Hammett constant for meta substituents (non conjugate substituents).<sup>4</sup>
21. S-O ( $\sigma_o$ ). Hammett constant for ortho substituents. It correlates poorly with  $\sigma_p$ , for 51 substituents  $r^2 = 0.303$ .<sup>30</sup>
22. S-O- ( $\sigma_o^-$ ). This is the parameter for ortho substituents.
23. S-INDUC ( $\sigma_I$ ). This parameter is for the field/inductive effect. Originally defined from 4-X-bicyclo[2.2.2]octane-1-carboxylic acids.<sup>4</sup>
24. S-AN-RS. Resonance parameter ( $\sigma^-$ ) obtained from anilines.
25. S-RES+. Resonance parameter for delocalization of +ve charge.<sup>5</sup>
26. S-'. Field/inductive effect parameters from bicycle[2.2.2]oct-ene-1-carboxylic acid, 4-X-dibenzobicyclo[2.2.2]octa-2,4-diene-1-carboxylic acids and cubanedicarboxylic acids.<sup>4</sup>
27. S-PARNO. A set of normalized  $\sigma_p$  values.<sup>31</sup>
28. S-ORTH+.  $\sigma^+$  for ortho substituents.<sup>31</sup>
29. S-PHOSP.  $\sigma$  for substituents attached to phosphorus.<sup>32</sup>
30. S-L.  $\sigma_L$  for field/inductive effect.<sup>4</sup>
31. S-ZTWST. Effect on  $\sigma$  resonance by twisting substituent  $90^\circ$  out of plane.
32. S-STAR. Classic  $\sigma^*$  defined by Taft.<sup>4</sup>
33. S-IND.P. Field/inductive parameter for substituents attached to phosphorus.<sup>32</sup>
34. S-RES-P. Resonance parameter for substituents attached to phosphorus.<sup>32</sup>
35. ER-P. Radical parameters at para position defined by Yamamoto and Otsu.<sup>33</sup>
36. ER-M. Radical parameters at meta position defined by Yamamoto and Otsu.<sup>33</sup>

37. S.DOT-P. Radical parameters ( $\sigma^\bullet$ ) at para position defined by Dust and Arnold.<sup>34</sup>
38. S.DOT-M. Radical parameters ( $\sigma^\bullet$ ) at meta position defined by Dust and Arnold.<sup>34</sup>
39. S.-DOT-P. Radical parameters ( $\sigma_{JJ}^\bullet$ ) at para position defined by Jiang and Ji.<sup>35</sup>
40. S.-DOT-M. Radical parameters ( $\sigma_{JJ}^\bullet$ ) at meta position defined by Jiang and Ji.<sup>35</sup>
41. S.P-C. Radical parameters ( $\sigma^\bullet$ ) at para position defined by Creary.<sup>36</sup>
42. S.M-C. Radical parameters ( $\sigma^\bullet$ ) at meta position defined by Creary.<sup>36</sup>
43. CMR-SUB. Calculated MR for substituents.

These parameters were collected over the past 40 years. Many of these are absolute and have been kept for the historical reasons. The most important and frequently used parameters are: Clog  $P$ , Mlog  $P$ , CMR, NVE, MgVol,  $C\pi$ ,  $\pi$ , MR-SUB, CMR-SUB,  $F$ ,  $E_s$ ,  $L$ ,  $B1$ ,  $B5$ , S-P ( $\sigma_p$ , classic  $\sigma$  constant for aromatic substituents), S-P+ ( $\sigma_p^+$ , when there is a direct resonance between substituent and the reaction center involving delocalization of a +ve charge), S-P- ( $\sigma_p^-$ , through resonance involving delocalization of a -ve charge), S-M ( $\sigma_m$ , Hammett constant for meta substituents), S-INDUC ( $\sigma_I$ , for aliphatic reactions), S-STAR ( $\sigma^*$ , for aliphatic reactions) etc.

The above parameters can be divided into three basic categories that elucidate the important features of the chemical entity. These three categories are hydrophobic, electronic and steric.

Hydrophobic:  $C\pi$ ,  $\pi$ , Clog  $P$ , Mlog  $P$

Electronic:  $\sigma$ ,  $\sigma^+$ ,  $\sigma^-$ ,  $\sigma_I$ ,  $\sigma^*$ ,  $F$ , NVE

Steric: MR-SUB, CMR-SUB, CMR, MgVol,  $E_s$ ,  $L$ ,  $B1$ ,  $B5$

## REFERENCES:

16. Hansch, C., Leo, A. & Hoekman, D. Exploring QSAR: Hydrophobic, Electronic, and Steric Constants, American Chemical Society, Washington, D.C. (1995).
17. Leo, A.J. Unpublished result
18. Hansch, C., Steinmetz, W.E., Leo, A.J., Mekapati, S.B., Kurup, A. & Hoekman, D. On the role of polarizability in chemical-biological interactions. *J. Chem. Inf. Comput. Sci.* **43**, 120-125 (2003).
19. Hansch, C. & Kurup, A. QSAR of chemical polarizability and nerve toxicity. 2. *J. Chem. Inf. Comput. Sci.* **43**, 1647-1651 (2003).
20. Verma, R.P., Kurup, A. & Hansch, C. On the role of polarizability in QSAR. *Bioorg. Med. Chem.* **13**, 237-255 (2005).
21. Verma, R.P. & Hansch, C. A comparison between two polarizability parameters in chemical-biological interactions. *Bioorg. Med. Chem.* **13**, 2355-2372 (2005).
22. Abraham, M.H. & McGowan, J.C. The use of characteristic volumes to measure cavity terms in reversed phase liquid chromatography. *Chromatographia* **23**, 243-246 (1987).
23. Swain, C.G. & Lupton, E.C., Jr. Field and resonance components of substituent effects. *J. Am. Chem. Soc.* **90**, 4328-4337 (1968).
24. Unger, S.H. & Hansch, C. Quantitative models of steric effects. *Prog. Phys. Org. Chem.* **12**, 91-118 (1976).
25. Charton, M. Quantitative treatment of the ortho effect. *Prog. Phys. Org. Chem.* **8**, 235-317 (1971).
26. Austel, V., Kutter, E. & Kalbfleisch, W. A new easily accessible steric parameter for structure-activity relationships. *Arzneim. Forsch.* **29**, 585-587 (1979).
27. Verloop, A. The Sterimol Approach to Drug Design, Marcel Dekker: New York (1987).
28. Taft, R.W. & Grob, C.A. Separation of polar and resonance effects in the ionization of 4-substituted pyridinium ions. *J. Am. Chem. Soc.* **96**, 1236-1238 (1974).
29. Okamoto, Y. & Brown, H.C. Quantitative treatment of electrophilic reactions of aromatic derivatives. *J. Org. Chem.* **22**, 485-494 (1957).
30. Charton, M. Electrical effect substituent constants for correlation analysis. *Prog. Phys. Org. Chem.* **13**, 119-251 (1981).
31. Le Guen, M.M.J. & Taylor, R. Electrophilic aromatic substitution. Part XVII. Protiodetritiation of some cycloalkyl- and secondary alkyl-benzenes. A linear free energy relation for ortho-aromatic substitution. *J. Chem. Soc. Perkin Trans 2*, 559-565 (1976).
32. Mastryukova, T.A. & Kabachnik, M.I. Correlation constants in the chemistry of organophosphorus compounds. *J. Org. Chem.* **36**, 1201-1205 (1971).
33. Yamamoto, T. & Otsu, T. Effects of substituents in radical reactions: Extension of the Hammett equation. *Chem. Ind.* 787-789 (1967).
34. Dust, J.M. & Arnold, D.R. Substituent effects on benzyl radical ESR hyperfine coupling constants. The  $\sigma_a$  scale based upon spin delocalization. *J. Am. Chem. Soc.* **105**, 1221-1227 (1983).

35. Jiang, X.-K. & Ji, G.-Z. A self-consistent and cross-checked scale of spin-delocalization substituent constant, the  $\sigma_{\bullet J}$  scale. *J. Org. Chem.* **57**, 6051-6056 (1992).
36. Creary, X., Mehrsheikh-Mohammadi, M.E. & McDonald, S. Methylenecyclopropane rearrangement as a probe for free radical substituent effects.  $\sigma_{\bullet}$  values for commonly encountered conjugating and organometallic groups. *J. Org. Chem.* **52**, 3254-3263 (1987).

**Table S1.** Classification of the C-QSAR Biological Database: October, 2006 (Number of Sets in Parentheses)<sup>a</sup>

Code	Class	Code	Class
<b>B0</b>	<b>unknown</b>	<b>B4</b>	<b>Single-Celled Organisms</b>
		B4A	Algae (51)
<b>B1</b>	Nonenzymatic <b>Macromolecules</b> (DNA, fibrin, hemoglobin, soil, albumin, etc.) (395)	B4B	Bacteria (982)
		B4C	Cells in culture (2247)
<b>B2</b>	<b>Enzymes</b>	B4E	Erythrocytes (79)
B2A	Oxidoreductases (1028)	B4F	Fungi, Molds (315)
B2B	Transferases (448)	B4P	Protozoa (179)
B2C	Hydrolases (1536)	B4V	Viruses (424)
B2D	Lyases (43)	B4Y	Yeasts (123)
B2E	Isomerases (56)		
B2F	Ligases (21)	<b>B5</b>	<b>Organs/Tissues</b>
B2G	Receptors (2672)	B5C	Cancer (531)
		B5G	Gastrointestinal tract (100)
<b>B3</b>	<b>Organelles</b>	B5H	Heart (97)
B3A	Mitochondria (94)	B5I	Internal/soft organs (69)
B3B	Microsomes (115)	B5N	Nerves, Brain, Muscles (402)
B3C	Chloroplasts (86)	B5S	Skin (62)
B3M	Membranes (162)	B5L	Liver (35)
B3R	Ribosomes (0)		
B3S	Synaptosomes (29)	<b>B6</b>	<b>Multicellular Organisms</b>
		B6A	Animal (vertebrates) (702)
		B6B	Insects (254)
		B6F	Fish (208)
		B6H	Human (56)
		B6I	Invertebrates (non-insect) (109)
		B6P	Plants (126)

<sup>a</sup>These numbers are constantly changing as new data are added daily

**Table S2.** Classification of the C-QSAR Physical Database: October, 2006 (Number of Sets in Parentheses)<sup>a</sup>

Code	Class	Code	Class
<b>PT</b>	<b>Theoretical</b> (65)	<b>P7</b>	<b>Addition</b>
		P7D	Dimerization (12)
<b>PO</b>	<b>Unknown</b> (0)	P7E	Electrophilic addition (149)
		P7N	Nucleophilic addition (289)
<b>P1</b>	<b>Ionization</b> (1618)	P7P	Polymerization (10)
P1P	Ionization potential (38)		
P1X	Proton exchange (74)	<b>P8</b>	<b>Elimination</b> (169)
		<b>P9</b>	<b>Rearrangement</b> (219)
<b>P2</b>	<b>Hydrolysis</b> (1032)	<b>P10</b>	<b>Oxidation</b> (569)
		<b>P12</b>	<b>Radical Reactions</b> (606)
<b>P3</b>	<b>Solvolysis</b> (644)	<b>P13</b>	<b>Complex Formation</b> (104)
<b>P4</b>	<b>Spectra</b>	<b>P14</b>	<b>Partitioning</b> (130)
P4I	Ionization spectra (60)	P14C	Chromatography (21)
P4E	ESR spectra (6)		
P4M	Mass spectra (12)	<b>P15</b>	<b>Pyrolysis</b> (88)
P4N	NMR spectra (194)	<b>P16</b>	<b>H-Bonding</b> (34)
P4R	IR spectra (9)	<b>P17</b>	<b>Electrochemical</b> (300)
P4U	UV spectra (23)	<b>P18</b>	<b>Brønsted</b> (119)
		<b>P19</b>	<b>Esterification</b> (237)
<b>P5</b>	<b>Miscellaneous Reactions</b> (522)	<b>P20</b>	<b>Photochemical</b> (48)
		<b>P21</b>	<b>Hydrogenation</b> (15)
<b>P6</b>	<b>Substitution</b>	<b>P22</b>	<b>Isokinetic</b> (2)
P6E	Electrophilic substitution (263)	<b>P23</b>	<b>Reduction</b> (99)
P6N	Nucleophilic substitution (1192)		

<sup>a</sup>These numbers are constantly changing as new data are added daily