



QSPR modeling for enthalpies of formation of organometallic compounds by means of SMILES-based optimal descriptors

A.A. Toropov^{a,b,*}, A.P. Toropova^{a,b}, E. Benfenati^b

^a Institute of Geology and Geophysics, Khodzhibaev Street 49, 100041 Tashkent, Uzbekistan

^b Istituto di Ricerche Farmacologiche Mario Negri, 20156, Via La Masa 19, Milano, Italy

ARTICLE INFO

Article history:

Received 20 May 2008

In final form 8 July 2008

Available online 15 July 2008

ABSTRACT

A quantitative structure–property relationship (QSPR) model for a predicting gas-phase enthalpy of formation have been developed, using as chemical information descriptors based on the simplified molecular input line entry system (SMILES). The model is one-variable equation. The SMILES-based descriptors calculated with correlation weights of SMILES attributes which are obtained by the Monte Carlo method. The model addressed organometallic compounds. Statistical characteristics of the model are the following: $n = 104$, $R^2 = 0.9943$, $Q^2 = 0.9940$, $s = 19.9$ (kJ/mol), $F = 17701$ (training set); $n = 28$, $R^2 = 0.9908$, $Q^2 = 0.9892$, $s = 29.4$ (kJ/mol), $F = 2788$ (test set).

© 2008 Elsevier B.V. All rights reserved.

Organometallic compounds are an important class of chemicals, because many of them have vital biochemical features. There are studies on the quantitative structure–property/activity relationships (QSPR/QSAR) for organic substances [1–5], but only a few articles have deal with QSPR for organometallic compounds [6–9]. A major problem with organometallic compounds is that many of the software calculating chemical descriptors are not suitable for metals.

Simplified molecular input line entry system (SMILES) [9–13] has been used as an alternative for molecular graphs in the QSPR/QSAR analyses [14–17]. The aim of the present study is the estimation of the SMILES-based optimal (flexible) descriptors in QSPR modeling of enthalpies of formation from elements for organometallic compounds. Numerical data and the split into the training ($n = 104$) and test ($n = 28$) sets were taken from [6]. Canonical SMILES notations have been generated with ACD/ChemSketch software [9].

SMILES-based optimal descriptors were calculated as

$$DCW = \Pi CW(s_k) \quad (1)$$

where s_k are SMILES attributes (SA_k). The SMILES attribute (invariant) can be a symbol of the SMILES notation (for instance, 'c', 'C', 'N', ')', '=', etc.), or two symbols of the SMILES encoding a physico-chemical image (for instance, 'Cl', 'Si', 'Pb', etc.); $CW(x)$ is the correlation weight for a SMILES attribute x . The CWs are calculated by the Monte Carlo method optimization procedure [11–13] that provides CWs values which used in Eq. (1) give a maximum for the correlation coefficient between the descriptor and modelled

parameter. Symbols ')' and ']' have been replaced by '(' and '[', because these are indicators of the same molecular phenomenon (branching and ions, respectively [9–13]).

In fact, each SMILES attribute is a representation of molecular structure: i.e., chemical elements, covalent bonds (double, triple), *cis*- and *trans* isomerism, and others [10–13]. The additive scheme [18] is a well-known approach of modeling properties by the selection of special additive contributions for molecular fragments. The SMILES-based optimal descriptors that is calculated with Eq. (1) is a similar approach, but multiplicative one. In principle, additive optimal descriptors which are also calculated with SMILES can be utilized in QSPR/QSAR analysis [19,20]. The algorithm of the Monte Carlo optimization (used for both additive and multiplicative schemes) is described in Ref. [21]. The algorithm has been checked in many QSPR/QSAR studies [7,8,15–17].

The list of the SMILES attributes and their correlation weights for three probes of the Monte Carlo method optimization are represented in Table 1. The one-variable model for the gas-phase enthalpy of formation of organometallic compounds (probe 1) is the following:

$$\Delta_f H^0 \text{ (kJ/mol)} = -22091.00 (\pm 18.33) + 22007.12 (\pm 18.28) DCW \quad (2)$$

$n = 104$, $R^2 = 0.9943$, $Q^2 = 0.9940$, $s = 19.9$ (kJ/mol), $F = 17701$ (training set);

$n = 28$, $R^2 = 0.9908$, $Q^2 = 0.9892$, $s = 29.4$ (kJ/mol), $F = 2788$ (test set)

An example of the DCW calculation is given in Table 2. Calculations with Eq. (2) are demonstrated in Table 3. Structures of the organometallic compounds are represented in Supplementary material.

* Corresponding author. Address: Institute of Geology and Geophysics, Khodzhibaev Street 49, 100041 Tashkent, Uzbekistan.

E-mail address: aatoropov@yahoo.com (A.A. Toropov).

Table 1

Correlation weights for SMILES attributes obtained in three probes of the Monte Carlo optimization

S_k	CW (S_k) in probe 1	CW (S_k) in probe 2	CW (S_k) in probe 3
#	1.0090556	1.0074964	1.0101665
(0.9998463	0.9998668	0.9998257
1	0.9996351	0.9999719	0.9998839
2	1.0009747	1.0008821	1.0010828
3	0.9996972	0.9998870	0.9996963
4	1.0001354	1.0004646	1.0004339
=	1.0042279	1.0039084	1.0050545
B	1.0017499	1.0026516	1.0025670
Al	1.0000149	1.0002363	0.9976983
C	0.9991618	0.9993168	0.9990561
As	1.0031238	1.0027897	1.0012775
Bi	1.0128181	1.0108004	1.0122084
Br	0.9940644	0.9951565	0.9932889
Cl	0.9925434	0.9938692	0.9915754
Ga	1.0012293	1.0012349	0.9990758
Ge	0.9990244	0.9993272	0.9964955
H	1.0022120	1.0018444	1.0024865
Hg	1.0059134	1.0050967	1.0044166
I	0.9961791	0.9968911	0.9956879
O	0.9963139	0.9964762	0.9954539
P	1.0029297	1.0038841	1.0041153
Pb	1.0090902	1.0077529	1.0080435
Sb	1.0055415	1.0048487	1.0040477
Se	1.0009714	1.0010277	0.9987955
Si	0.9947625	0.9958890	0.9917372
Sn	1.0038262	1.0033491	1.0019871
Te	1.0028037	1.0025096	1.0008769
[1.0020271	1.0021860	1.0037950
C	1.0008545	1.0006440	1.0009193

Table 2

Example of the DCW calculation: SMILES = 'CC(C)C[Al](CC(C)C)CC(C)C'; no. 1; DCW = 0.9927988

S_k	CW (S_k) in probe 1
C	0.9991618
C	0.9991618
(0.9998463
C	0.9991618
(0.9998463
C	0.9991618
[1.0020271
Al	1.0000149
[1.0020271
(0.9998463
C	0.9991618
C	0.9991618
(0.9998463
C	0.9991618
(0.9998463
C	0.9991618
(0.9998463
C	0.9991618
(0.9998463
C	0.9991618

The statistical characteristics of QSPR for the enthalpy of formation, with the same split into training and test sets, which have been reported in [6] are $R^2 = 0.988$, $s = 29.1$ (training set) and

Table 3Experimental and calculated with Eq. (2) gas-phase $\Delta_f H^0$ (kJ/mol) for the aliphatic organometallic compounds

No. ^a	SMILES	DCW	Exp.	Calcd.
<i>Training set</i>				
1	CC(C)C[Al](CC(C)C)CC(C)C	0.9927988	-231.500	-242.357
2	CCCC[Al](CCCC)CCCC	0.9937151	-216.000	-222.192
4	C[Al](C)C	1.0012427	-86.500	-56.531
5	CCC[Al](CCC)CCC	0.9962180	-153.300	-167.111
6	C[As](C)C	1.0043555	12.500	11.972
8	CC(C)CCB(CCC(C)C)CCC(C)C	0.9880130	-381.800	-347.680
9	CC(C)CB(CC(C)C)CC(C)C	0.9905015	-280.300	-292.915
10	CC(C)B(C(C)C)C(C)C	0.9929963	-251.400	-238.011
13	CCCCB(Br)CCCC	0.9888422	-300.400	-329.432
14	CCCCB(Cl)CCCC	0.9873291	-365.700	-362.729
15	CCCCB(I)CCCC	0.9909458	-225.500	-283.137
16	CCCCB(CCCC)CCCC	0.9914157	-287.000	-272.796
18	CCB(CC)CC	0.9964162	-152.700	-162.749
19	CCCCCB(CCCCC)CCCCC	0.9864403	-395.000	-382.291
20	CCCCCB(CCCCC)CCCCC	0.9839619	-454.800	-436.832
21	CB(C)C	0.9989259	-122.600	-107.518
22	CCCCCB(CCCCC)CCCCC	0.9814898	-514.600	-491.236
23	ClB(Cl)c1cccc1	0.9909081	-266.100	-283.967
25	CCCC[Bi](CCCC)CCCC	1.0064377	71.200	57.795
26	C[Bi](C)C	1.0140617	194.400	225.577
28	CCC[Bi](CCC)CCC	1.0089726	133.900	113.582
30	CCCC[Ga](CCCC)CCCC	0.9949219	-226.900	-195.635
31	CC[Ga](CC)CC	0.9999401	-62.300	-85.199
32	C[Ga](C)C	1.0024586	-35.900	-29.772
33	CCC[Ga](CCC)CCC	0.9974278	-125.500	-140.487
35	O=C(c1cccc1)[Ge](C(=O)c2cccc2)(C(=O)c3cccc3)C(=O)c4cccc4	1.0215911	388.200	391.278
36	CC[Ge](CC)CC	0.9977380	-111.900	-133.660
37	CC[Ge](CC)(CC)CC	0.9957599	-164.900	-177.192
38	Br[Ge](C)(C)C	0.9940083	-222.200	-215.740
41	c1ccc(cc1)[Ge]2(CCCC)c3cccc3	1.0100168	138.900	136.560
42	C#Cc1cccc1[Ge](c2cccc2)(c3cccc3)c4cccc4	1.0313527	578.900	606.104
43	C=C[Ge](c1cccc1)(c2cccc2)c3cccc3	1.0212093	362.600	382.877
44	c1cccc1[Ge](c2cccc2)(c3cccc3)c4cccc4	1.0241275	438.600	447.097
45	CCC[Ge](CCC)(CCC)CCC	0.9924257	-229.700	-250.570
46	[Hg](C#Cc1cccc1)C#Cc2cccc2	1.0364297	720.500	717.833
47	CC(C)CC[Hg]CCC(C)C	1.0009463	-82.700	-63.054

Table 3 (continued)

No. ^a	SMILES	DCW	Exp.	Calcd.
48	CC(C)[Hg]CC(C)C	1.0026264	-38.300	-26.081
49	CC(C)[Hg]C(C)C	1.0043092	37.000	10.954
50	CC(C)[Hg]Br	1.0011705	-53.800	-58.121
51	Br[Hg]CC	1.0023185	-30.300	-32.856
53	Br[Hg]CCC	1.0014784	-51.400	-51.345
56	CC(C)[Hg]Cl	0.9996386	-86.900	-91.834
57	CC[Hg]Cl	1.0007849	-67.700	-66.607
58	C[Hg]Cl	1.0016244	-55.100	-48.132
60	CC[Hg]CC	1.0066138	75.000	61.671
61	CC(C)[Hg]I	1.0033003	-5.800	-11.250
63	C[Hg]I	1.0052934	21.600	32.612
64	CC[C]HgI	1.0036089	-3.700	-4.459
65	C[Hg]C	1.0083034	92.400	98.853
66	c1cccc1[Hg]c2cccc2	1.0216447	394.200	392.458
67	CC[C]Hg]CCC	1.0049271	30.300	24.551
69	C[Pb](C)(C)C	1.0091720	136.100	117.970
70	Br[Pb](Br)(c1cccc1)c2cccc2	1.0121183	177.700	182.810
71	I[Pb](I)(c1cccc1)c2cccc2	1.0164292	288.200	277.679
72	Br[Pb](c1cccc1)(c2cccc2)c3cccc3	1.0227730	406.300	417.289
73	I[Pb](c1cccc1)(c2cccc2)c3cccc3	1.0249488	455.400	465.172
76	c1cccc1P(c2cccc2)c3cccc3	1.0187786	328.400	329.384
78	CC[Sb](CC)CC	1.0042467	48.700	9.578
80	c1cccc1[Sb](c2cccc2)c3cccc3	1.0255770	435.400	478.996
81	CC[C]Sb](CC)CCC	1.0017236	-37.700	-45.948
82	CC(C)[Se]C(C)C	0.9993752	-108.000	-97.631
84	CC[Se]CC	1.0016684	-57.300	-47.163
85	C[Se]C	1.0033497	17.800	-10.164
86	CCCC[Se]CCCC	0.9966415	-172.700	-157.790
87	c1cccc1[Se]c2cccc2	1.0166255	289.700	281.999
88	CC[Si](CC)(CC)CC	0.9915120	-265.700	-270.677
90	Br[Si](C)(C)C	0.9897678	-297.500	-309.061
91	C[Si](C)(C)Cl	0.9882534	-354.000	-342.389
92	C[SiH](Cl)Cl	0.9850069	-415.000	-413.835
93	Cl[Si](Cl)(c1cccc1)c2cccc2	0.9946967	-208.800	-200.590
94	CC(C)C[Sn](CC(C)C)(CC(C)C)CC(C)C	0.9926349	-272.200	-245.964
95	CC(C)[Sn](C(C)C)(C(C)C)C(C)C	0.9959699	-119.400	-172.571
96	CCCC[Sn](Br)(CCCC)CCCC	0.9912768	-270.600	-275.852
97	CCCC[Sn](CCCC)(CCCC)CCCC	0.9938566	-217.400	-219.078
98	CC[Sn](Cl)(CC)CC	0.9947522	-193.300	-199.369
99	CC[Sn](CC)CC	1.0025336	-54.700	-28.122
100	CC[Sn](CC)(CC)CC	1.0005460	-42.000	-71.864
101	CC[Sn](Cl)(Cl)Cl	0.9832651	-429.300	-452.168
102	CC[Sn](c1cccc1)(c2cccc2)c3cccc3	1.0217977	380.000	395.825
103	C[Sn](C)(Cl)Cl	0.9906519	-337.000	-289.604
104	C[Sn](C)(I)I	0.9979228	-150.000	-129.593
107	CC(C)C[Sn](C)(C)C	1.0007697	-67.000	-66.942
108	Br[Sn](C)(C)C	0.9987860	-138.100	-110.598
109	C[Sn](C)(C)C(=O)c1cccc1	1.0085495	90.400	104.270
110	C[Sn](C)(C)Cl	0.9972577	-174.900	-144.230
111	C[Sn](C)(C)CC	1.0030661	-26.300	-16.404
112	C[Sn](C)C	1.0050587	24.900	27.448
113	C[Sn](C)(C)I	1.0009107	-82.400	-63.838
115	C[Sn](C)(C)CCC	1.0022254	-46.800	-34.906
116	C[Sn](C)(C)C	1.0039076	-17.600	2.114
117	C[Sn](Cl)(Cl)Cl	0.9840899	-417.100	-434.016
118	C[Sn](c1cccc1)(c2cccc2)c3cccc3	1.0226549	406.000	414.689
119	c1ccc(cc1)[Sn]2(CCCC2)c3cccc3	1.0148714	301.800	243.396
121	C(#C)c1cccc1[Sn](c2cccc2)(c3cccc3)c4cccc4	1.0363099	736.900	715.197
122	C=[Sn](c1cccc1)(c2cccc2)c3cccc3	1.0261178	528.600	490.897
123	c1cccc1[Sn](c2cccc2)(c3cccc3)c4cccc4	1.0290499	575.400	555.426
124	CCC[Sn](CCC)(CCC)CCC	0.9971957	-142.900	-145.594
125	CC(C)C[Te]CCC(C)C	0.9978520	-148.000	-131.150
126	CC(C)[Te]C(C)C	1.0012045	-46.000	-57.371
127	CCCC[Te]CCCC	1.0001418	-56.000	-80.759
128	CC[Te]CC	1.0035020	-4.400	-6.811
129	C[Te]C	1.0051863	28.600	30.256
130	CCCC[Te]CCCC	0.9984659	-137.000	-117.641
132	CCC[Te]CCC	1.0018205	-42.500	-43.816
<i>Test set</i>				
3	CC[Al](CC)CC	0.9987272	-114.100	-111.890
7	c1cccc1[As](c2cccc2)c3cccc3	1.0231112	408.400	424.730
11	ClB(Cl)c1ccc(C)cc1	0.9897731	-294.400	-308.944
12	CC(CCCCC)B(C(C)CCCC)C(C)CCCC	0.9805848	-507.900	-511.153
17	C1CCCC1B(C2CCCC2)C3CCCC3	0.9870448	-397.900	-368.986
24	CCCB(CCC)CCC	0.9939128	-236.000	-217.842

(continued on next page)

Table 3 (continued)

No. ^a	SMILES	DCW	Exp.	Calcd.
27	c1cccc1[Bi](c2cccc2)c3cccc3	1.0329986	600.600	642.325
29	CC(C)C[Ga](CC(C)C)CC(C)C	0.9940045	-239.500	-215.825
34	CCCC[Ge](CCCC)(CCCC)CCCC	0.9891025	-310.000	-323.702
39	C[Ge](C)(C)Cl	0.9924874	-266.100	-249.212
40	C[Ge](C)(C)C	0.9991054	-107.500	-103.568
52	Br[Hg]C	1.0031593	-18.600	-14.352
54	CCCC[Hg]CCCC	1.0032432	-32.400	-12.507
55	O=C([Hg]C(=O)c1cccc1)c2cccc2	1.0203788	278.000	364.599
59	CCC[Hg]Cl	0.9999460	-88.100	-85.068
62	CC[Hg]I	1.0044508	14.300	14.069
68	CC[Pb](CC)(CC)CC	1.0057928	109.600	43.604
74	c1cccc1[Pb](c2cccc2)(c3cccc3)c4cccc4	1.0344462	674.100	674.183
75	CP(C)C	1.0001024	-94.100	-81.627
77	CCCC[Sb](CCCC)CCCC	0.9992069	-100.500	-101.334
79	C[Sb](C)C	1.0067761	32.200	65.243
83	CCCC[Se]CCCC	0.9983144	-131.800	-120.976
89	C[Si](C)(Cl)Cl	0.9817072	-461.100	-486.451
105	C[Sn](C)(C)C=C	1.0073070	91.500	76.925
106	CC(C)[Sn](C)(C)C	1.0019172	-43.700	-41.687
114	C[Sn](C)(C)c1cccc1	1.0091750	113.600	118.034
120	c1cc(ccc1)[Sn]2(CCCCC2)c3cccc3	1.0140207	289.000	224.676
131	c1cccc1[Te]c2cccc2	1.0184865	306.800	322.954

^a The numbering is taken from Ref. [6].

$R^2 = 0.990$, $s = 32.0$ (test set). Molecular graphs are at the basis of descriptors used in [6]. Thus, the predictive abilities of the SMILES-based descriptors and regression analysis using descriptors calculated with molecular graphs [6] are very similar.

The criteria for validation of the QSPR/QSAR for the external test set suggested in Ref. [22] are the following:

$$k = \frac{\sum [Ye * Yp]^2}{\sum Ye^2} \quad k' = \frac{\sum [Ye * Yp]^2}{\sum Yp^2}$$

$$Y0e = k * Yp \quad Y0p = k' * Ye$$

$$Ro^2 = 1 - \frac{\sum [Yp - Y0e]^2}{\sum [Yp - YAp]^2} \quad Ro'^2 = 1 - \frac{\sum [Ye - Y0p]^2}{\sum [Ye - YAe]^2}$$

where Ye and Yp are experimental and predicted values for a test set; YAe and YAp are average values of the Ye and Yp .

According to Ref. [22], QSPR/QSAR models can be considered acceptable if they satisfy all the following conditions:

$$(R^2 - Ro^2)/R^2 < 0.1 \quad \text{and} \quad (R^2 - Ro'^2)/R^2 < 0.1;$$

$$0.85 < k < 1.15 \quad \text{and} \quad 0.85 < k' < 1.15$$

An indicator of the external predictability of QSPR/QSAR models, suggested in Ref. [23], is the following:

$$Rm^2 = R^2(1 - \text{abs}(R^2 - Ro^2)^{0.5}) > 0.5$$

These characteristics of the model calculated with Eq. (2) (for the test set) are the following:

$$k = 0.97921; k' = 1.01149;$$

$$(R^2 - Ro^2)/R^2 = -0.0092 < 0.1$$

$$(R^2 - Ro'^2)/R^2 = -0.0089 < 0.1$$

$$Rm^2 = 0.8962 > 0.5$$

Thus the Eq. (2) satisfies the criteria suggested in Refs. [22,23].

An important advantage of the suggested approach is the possibility to carry out QSPR/QSAR analyses directly from databases on physicochemical properties and/or biological activity available via the internet. Since each SMILES symbol has transparent physicochemical interpretation, this approach is able to hint ideas related to the mechanisms of physicochemical and biological phenomena which have influence on the property/activity of interest.

Acknowledgement

The authors thank the Marie Curie Fellowship for financial support (the contract ID 39036, CHEMPREDICT).

Appendix A. Supplementary material

Supplementary data associated with this article can be found, in the online version, at doi:10.1016/j.cplett.2008.07.027.

References

- [1] K.V. Balakin, N.P. Savchuk, I.V. Tetko, *Curr. Med. Chem.* 13 (2006) 223.
- [2] I.V. Tetko, V.Yu. Tanchuk, T.N. Kasheva, A.E.P. Villa, *J. Chem. Inf. Comput. Sci.* 41 (2001) 246.
- [3] P.R. Duchowicz, A. Talevi, C. Bellera, L.E. Bruno-Blanch, E.A. Castro, *Bioorg. Med. Chem.* 15 (2007) 3711.
- [4] F. Torrens, J. Sanchez-Marrn, I. Nebot-Gil, *J. Chromatogr. A* 827 (1998) 345.
- [5] J.T. Leonard, K. Roy, *Eur. J. Med. Chem.* 43 (2008) 81.
- [6] J. Jover, R. Bosque, J.A.M. Simoes, J. Sales, *J. Organomet. Chem.* 693 (2008) 1261.
- [7] A.A. Toropov, A.P. Toropova, *Russ. J. Coord. Chem.* 27 (2001) 574.
- [8] A.A. Toropov, A.P. Toropova, *Russ. J. Coord. Chem.* 28 (2002) 877.
- [9] ACD/ChemSketch Freeware, version 11.00, Advanced Chemistry Development, Inc., Toronto, ON, Canada, (www.acdlabs.com), 2007.
- [10] D. Weininger, *J. Chem. Inf. Comput. Sci.* 28 (1988) 31.
- [11] D. Weininger, A. Weininger, J.L. Weininger, *J. Chem. Inf. Comput. Sci.* 29 (1989) 97.
- [12] D. Weininger, *J. Chem. Inf. Comput. Sci.* 30 (1990) 237.
- [13] (<http://www.daylight.com>).
- [14] D. Vidal, M. Thormann, M. Pons, M. J. Chem. Inf. Model. 45 (2005) 386.
- [15] A.A. Toropov, A.P. Toropova, D.V. Mukhamedzhanova, I. Gutman, *Indian J. Chem. Sec. A* 44 (2005) 1545.
- [16] A.A. Toropov, E. Benfenati, *Comput. Biol. Chem.* 31 (2007) 57.
- [17] A.A. Toropov, E. Benfenati, *Eur. J. Med. Chem.* 42 (2007) 606.

- [18] I.G. Zenkevich, M. Moeder, G. Koeller, S.J. Schrader, *Chromatogr. A* 1025 (2004) 227.
- [19] A.A. Toropov, B.F. Rasulev, D. Leszczynska, J. Leszczynski, *Chem. Phys. Lett.* 444 (2007) 209.
- [20] A.A. Toropov, E. Benfenati, *Bioorg. Med. Chem.* 16 (2008) 4801.
- [21] A.A. Toropov, B.F. Rasulev, D. Leszczynska, J. Leszczynski, *Chem. Phys. Lett.* 457 (2008) 332.
- [22] A. Golbraikh, A. Tropsha, *J. Mol. Graph. Model.* 20 (2002) 269.
- [23] P.P. Roy, K. Roy, *QSAR Comb. Sci.* 27 (2008) 302.