

Jäntschi, L., Katona, G., Diudea, M.V.

## **Modeling Molecular Properties by Cluj Indices**

(2000) *Match*, 41, pp. 151-188.

<sup>a</sup> Fac. of Sci. and Eng. of Materials, Technical University

<sup>b</sup> Bios Research and Production Center

<sup>c</sup> Fac. of Chem. and Chem. Engineering, Babes-Bolyai University, 3400 Cluj, Romania

### **Abstract**

A new approach, leading to a fragmental property index family, FPIF, is presented. Indices are calculated as local descriptors of some molecular fragments and the global values are then obtained by summing the fragmental contributions. The modeling ability of FPIF is demonstrated by modeling some physico-chemical properties and biological activities on selected sets of organic compounds. The results are compared to those reported in some previous works.

### **References**

• Free, S.M., Wilson, J.W.

#### **A mathematical contribution to structure-activity studies**

(1964) *J. Med. Chem.*, 7, p. 395.

• Gao, C., Govind, R., Tabak, H.H.

#### **Application of the group contribution method for predicting the toxicity of organic chemicals**

(1992) *Environmental Toxicol. Chem.*, 11, pp. 631-636.

• Kalivas, J.H., Sutter, J.M., Roberts, N.

#### **Global optimization by simulated annealing with wavelength selection for ultraviolet-visible spectrophotometry**

(1989) *Anal. Chem.*, 61, pp. 2024-2030.

• Kalivas, J.H.

#### **Generalized simulated annealing for calibration sample selection from an existing set and orthogonization of undesigned experiments**

(1991) *J. Chemometrics*, 5, pp. 37-48.

- Sutters, J.M., Dixon, S.L., Jurs, P.C.  
**Automated descriptor selection for Quantitative Structure-Activity Relationships, using generalized simulated annealing**  
(1995) *J. Chem. Inf. Comput. Sci.*, 35, pp. 77-84.
- Leardi, R., Boggia, R., Terrile, M.  
**Genetic algorithms as a strategy for feature selection**  
(1992) *J. Chemom.*, 6, p. 267.
- Lucasius, C.B., Kateman, G.  
**Undersanding and using genetic algorithms. Part 1. Concepts, properties and context**  
(1993) *Chemom. Intell. Lab. Sys.*, 19, p. 1.
- Diudea, M.V.  
**Cluj matrix CJu : Source of various graph descriptors**  
(1997) *Commun. Math. Comput. Chem. (MATCH)*, 35, pp. 169-183.
- Diudea, M.V., Minailiuc, O., Katona, G., Gutman, I.  
**Szeged matrices and related numbers**  
(1997) *Commun. Math. Comput. Chem. (MATCH)*, 35, pp. 129-143.
- Diudea, M.V.  
**Cluj matrix invariants**  
(1997) *J. Chem. Inf. Comput. Sci.*, 37, pp. 300-305.
- Diudea, M.V., Pârv, B., Topan, M.I.  
**Derived Szeged and Cluj indices**  
(1997) *J. Serb. Chem. Soc.*, 62, pp. 267-276.
- Diudea, M.V., Gutman, I.  
**Wiener-type tpological indices**  
(1998) *Croat. Chem. Acta*, 71, pp. 21-51.
- Kiss, A.A., Katona, G., Diudea, M.V.  
**Szeged and Cluj matrices within the matrix operator W (M1,M2,M3)**  
(1997) *Coll. Sci. Papers Fac. Sci. Kragujevac*, 19, pp. 95-107.
- Gutman, I., Diudea, M.V.  
**Defining Cluj matrices and Cluj matrix invariants**  
(1998) *J. Serb. Chem. Soc.*, 63, pp. 497-504.
- Diudea, M.V., Parv, B., Gutman, I.  
**Detour-Cluj matrix and derived invariants**  
(1997) *J.Chem.Inf.Comput.Sci.*, 37, pp. 1101-1108.

- Diudea, M.V., Katona, G., Lukovits, I., Trinajstić, I.  
**Detour and Cluj-detour indices**  
(1998) *Croat. Chem. Acta*, 71, pp. 459-471.
- Minailiuc, O., Katona, G., Diudea, M.V., Strunje, M., Graovac, A., Gutman, I.  
**Szeged fragmental indices**  
(1998) *Croat. Chem. Acta*, 71, pp. 473-488.
- Horn, R.A., Johnson, C.R.  
(1985) *Matrix Analysis*,
- Diudea, M.V., Gutman, I., Jäntschi, L.  
*Molecular Topology*,
- Sears, F.W., Zemansky, M.W., Young, H.D.  
(1976) *University Physics, Fifth Edition*, Addison - Wesley Publishing Company,
- Golender, V., Vesterman, B., Vorpapel, E.  
**APEX-3D Expert system for drug design**  
(1996) *Network Science*,
- Rose, V.S., Wood, J.  
**Generalized cluster significance analysis and stepwise cluster significance analysis with conditional probabilities**  
(1998) *Quant. Struct.-Act. Relat.*, 17, pp. 348-356.
- Young, H.D.  
(1962) *Statistical Treatment of Experimental Data*,
- Reif, F.  
(1965) *Fundamentals of Statistical and Thermal Physics*,
- Diudea, M.V., Silaghi-Dumitrescu, I.  
**Valence group electronegativity as a vertex discriminator**  
(1989) *Rev. Roumaine Chim.*, 34, pp. 1175-1182.
- Diudea, M.V., Kacso, I.E., Topan, M.I.  
**Molecular topology. 18. A Qspr/Qsar study by using new valence group carbon-related electronegativities**  
(1996) *Rev. Roum. Chim.*, 41, pp. 141-157.
- Crawford Jr., F.S.  
(1968) *Waves*, 3.
- Ivanciuc, O.  
**3D QSAR models**

*QSPR/QSAR Studies by Molecular Descriptors,*

- Hopfield, J.J.  
**Neural networks and physical systems with emergent collective computational abilities**  
(1982) *Proc. Natl. Acad. Sci. U.S.A.*, 79, pp. 2554-2558.
- Zupan, J., Gasteiger, J.  
**Neural networks: A new method for solving chemical problems or just a passing phase?**  
(1991) *Anal. Chim. Acta*, 248, pp. 1-30.
- Gasteiger, J., Zupan, J.  
**Neural networks in chemistry**  
(1993) *Angew. Chem. Int. Ed. Engl.*, 32, pp. 503-527.
- Zupan, J., Gasteiger, J.  
(1993) *Neural Networks for Chemists*,
- Bulsari, A.B.  
(1995) *Neural Networks for Chemical Engineers*,
- Devillers, J.  
(1996) *Neural Networks in QSAR and Drug Design*, p. 279.
- Ivanciuc, O., Rabine, J.-P., Cabrol-Bass, D., Panaye, A., Doucet, J.P.  
**<sup>13</sup>C NMR chemical shift prediction of sp<sup>2</sup> carbon atoms in acyclic alkenes using neural networks**  
(1996) *J. Chem. Inf. Comput. Sci.*, 36, pp. 644-653.
- Ivanciuc, O.  
**Molecular graph descriptors used in neural network models**  
(1999) *Topological Indices and Related Descriptors in QSAR and QSPR*, pp. 697-777.
- Gakh, A.A., Gakh, E.G., Sumpter, B.G., Noid, D.W.  
**Neural network-graph theory approach to the prediction of the physical properties of organic compounds**  
(1994) *J. Chem. Inf. Comput. Sci.*, 34, pp. 832-839.
- Baskin, I.I., Palyulin, V.A., Zefirov, N.S.  
**A Neural device for searching direct correlations between structures and properties of chemical compounds**  
(1997) *J. Chem. Inf. Comput. Sci.*, 37, pp. 715-721.
- Kireev, D.B.  
**ChemNet: A Novel neural network based method for**

**graph/property mapping**

(1995) *J. Chem. Inf. Comput. Sci.*, 35, pp. 175-180.

- Topliss, J.G., Edwards, R.P.

**Chance factors in in studies of Quantitative Structure-Activity Relationships**

(1979) *J. Med. Chem.*, 22, p. 1238.

- Nikolić, S., Medić-Sarić, M., Matijević-Sosa, J.

**A QSAR study of 3-(Phtalimidoalkyl)-pyrazolin-5-ones**

(1993) *Croat. Chem. Acta*, 66, pp. 151-160.

- Nagy, P.I., Tokarski, J., Hopfinger, A.J.

**Molecular shape and QSAR analysis of a family of substituted dichlorodiphenyl aromatase inhibitors**

(1994) *J. Chem. Inf. Comput. Chem.*, 34, pp. 1190-1197.

- Wessel, M.D., Jurs, P.C.

**Prediction of normal boiling points for a diverse set of industrially important organic compounds from molecular structure**

(1995) *J. Chem. Inf. Comput. Sci.*, 35, pp. 841-850.

- Stanton, D.T., Jurs, P.C.

**Development and use of charged partial surface area structural descriptors for Quantitative Structure-Property Relationships studies**

(1990) *Anal. Chem.*, 62, p. 2323.

- Goll, E.S., Jurs, P.C.

**Prediction of the Normal Boiling Points of Organic Compounds from Molecular Structures with a Computational Neural Network Model**

(1999) *J. Chem. Inf. Comput. Chem.*, 39, pp. 974-983.

- Stuper, A.J., Brugger, W.E., Jurs, P.C.

(1979) *Computer-assisted Studies of Chemical Structure and Biological Function*,

- Trebst, A., Draber, W.

**Advan. Pest. Sci**

(1978) *Symp. Papers IV-th Int. Congress Pest. Chem.*, p. 223.

- Benigni, R., Gallo, G., Giorgi, F., Giuliani, A.

**On the equivalence between different descriptions of molecules: Value for computational approaches**

(1999) *J. Chem. Inf. Comput. Sci.*, 39, pp. 575-578.

### **Correspondence Address**

Diudea M.V.; Fac. of Chem. and Chem. Engineering; Babes-Bolyai University 3400 Cluj, Romania

**ISSN:** 03406253

**Language of Original Document:** English

**Abbreviated Source Title:** Match

## Modeling Molecular Properties by Cluj Indices

Lorentz Jäntschi,<sup>a</sup> Gabriel Katona<sup>b</sup> and Mircea V. Diudea<sup>c\*</sup>

<sup>a</sup>Faculty of Science and Engineering of Materials, Technical University

<sup>b</sup>Bios Research and Production Center

<sup>c</sup>Faculty of Chemistry and Chemical Engineering, “Babeş-Bolyai” University  
3400 Cluj, Romania

**Abstract.** A new approach, leading to a *fragmental property index family*, **FPIF**, is presented. Indices are calculated as local descriptors of some molecular fragments and the global values are then obtained by summing the fragmental contributions. The modeling ability of **FPIF** is demonstrated by modeling some physico-chemical properties and biological activities on selected sets of organic compounds. The results are compared to those reported in some previous works.

*Commun. Math. Comput. Chem. (MATCH)*, **2000**, *41*, 151-188

### Introduction

**QSPRs** (Quantitative Structure-Property Relationships) link quantitatively the physico-chemical properties of chemical compounds with the molecular structure. They provide mathematical models aimed to accurately predict a certain property from the structural attributes.

Some molecular properties (i.e. those of which numerical value vary with changes in the molecular structure) such as the normal boiling point, critical parameters, viscosity, solubility, retention chromatographic index, are often used for characterizing chemicals in databases. However, a certain property is not always available in tables or other reference sources. It is just the case of newly synthesized compounds. As a consequence, methods of estimation / prediction of physico-chemical properties from the structural features of organic

molecules become very important. The advent of combinatorial chemistry in the last decade required automated procedures for predicting various molecular properties.

Monitoring the environmental pollution needs the prediction of toxicity of chemicals in air, waste waters and soil. **QSARs** (Quantitative Structure-Property Relationships) can be used to predict the toxicity accurately, without using more expensive experimental methods. Drug research and production is also related to the **QSAR** techniques.

**QSPRs/QSARs** thus relate a molecular property, shown by a series of chemicals, to the structure encoded by a set of descriptors.

In the past works, investigators looked for easily obtainable descriptors.<sup>1,2</sup> As the computer technology developed, the attention of the descriptor designers turned toward more elaborated descriptors, with enhanced ability in modeling a certain molecular property/activity. Large pools of descriptors were thus created.

Nevertheless, a new problem arised: how to select, in real time, a subset of descriptors suitable for the optimal modeling of the chosen property? An algorithm doing such a task must be a simple one, rapid and convergent to a global optimum. Simulated annealing **SA** and genetic algorithms **GA** are already verified procedures.<sup>3-7</sup>

In this paper a new approach, leading to a *fragmental property index family*, **FPIF**, is presented. These indices are calculated as local descriptors of some fragments of the molecule and, a global index is then obtained by summing the fragmental contributions. The modeling ability of **FPIF** is demonstrated on selected sets of organic compounds.

### Fragmentation Criteria

The fragmentation criteria define the basic topological descriptors: **CJ**, **CF** and **Sz**. The fragments are just the entries in the Cluj and Szeged matrices,<sup>8-17</sup> respectively. Before defining these matrices, some graph-theoretical background is needed.

Let  $G = (V, E)$  be a connected graph, with  $V$  being the set of vertices and  $E \subset V \times V$  the set of edges.

A *walk*  $w$  is an alternating string of vertices and edges,  $w_{1,n} = (v_1, e_1, v_2, e_2, \dots, v_{n-1}, e_{n-1}, v_n)$ ,  $v_i \in V(G)$ ,  $e_i \in E(G)$ ,  $m \geq n - 1$ , such that any subsequent pair of vertices  $(v_{i-1}, v_i) \in E(G)$ . Revisiting of vertices and edges is allowed. Then  $V(w_{1,n}) = \{v_1, v_2, \dots, v_{n-1}, v_n\}$  is the set of vertices of  $w_{1,n}$ . Similarly,  $E(w_{1,n}) = \{e_1, e_2, \dots, e_{n-1}, e_n\}$  is the set of edges of  $w_{1,n}$ .



The *length* of a walk,  $l(w_{1,n}) = |E(w_{1,n})| \geq |V(w_{1,n})| - 1$ , equals to the number of its traversed edges. The walk is *closed* if  $v_1 = v_n$  (i.e. its *endpoints* coincide) and is *open* otherwise. The *set of all walks* in  $G$  is denoted by  $W(G)$ .

A *path*  $p$  is a walk having all its vertices and edges distinct:  $v_i \neq v_j$ ,  $(v_{i-1}, v_i) \neq (v_{j-1}, v_j)$  for any  $1 \leq i < j \leq n$ . As a consequence, the revisiting of vertices and edges, as well as branching, is prohibited. The *length* of a path is  $l(p_{1,n}) = |E(p_{1,n})| = |V(p_{1,n})| - 1$ . A closed path is a *cycle* (i.e. *circuit*). The *set of all paths* in  $G$  is denoted by  $P(G)$ .

A *terminal path*  $tp_{1,n}$  is the path involving a walk  $w = v_1, e_1, v_2, \dots, v_n, e_n, v_k$ , that is *no more a path* in  $G$ , for any  $v_k \in V(G)$  such that  $(v_n, v_k) = e_n \in E$ .

A path is *Hamiltonian* if  $n = |V(G)|$ . In words, a Hamiltonian path visits once all the vertices in  $G$ . If such a path is a closed one, then it is a *Hamiltonian circuit*.

The *distance*,  $d_{ij}$ , between two vertices  $v_i$  and  $v_j$  is the length of a *shortest* path joining them, if *exists*:  $d_{ij} = \min l(p_{ij})$ ; otherwise  $d_{ij} = \infty$ . A shortest path is often called a *geodesic*. The *eccentricity* of a vertex  $i$ ,  $ecc_i$ , is the maximum distance between  $i$  and any vertex  $j$  of  $G$ :  $ecc_i = \max d_{ij}$ . The *radius* of a graph,  $r(G)$ , is the minimum eccentricity among all vertices  $i$  in  $G$ :  $r(G) = \min ecc_i = \min \max d_{ij}$ . Conversely, the *diameter*,  $d(G)$ , is the maximum eccentricity in  $G$ :  $d(G) = \max ecc_i = \max \max d_{ij}$ . The *set of all geodesics* (i.e. distances) in  $G$  is denoted by  $D(G)$ .

The *detour*,  $\delta_{ij}$ , between two vertices  $v_i$  and  $v_j$  is the length of a *longest* path joining these vertices, if *exists*:  $\delta_{ij} = \max l(p_{ij})$ ; otherwise  $\delta_{ij} = \infty$ . The *set of all detours* (i.e. longest paths) in  $G$  is denoted by  $\Delta(G)$ .

Resuming to the Cluj and Szeged matrix definition.

Let  $p \in D(G)$  or  $p \in \Delta(G)$ ; the *Cluj Fragments*,  $CJ$  and  $CF$  represent the sets of vertices obeying the relations

$$CJ_{i,j,p} = \{v \mid v \in V(G); d(G)_{v,i} < d(G)_{v,j}; \text{ and } \exists w \in W_{v,b} V(w) \cap V(p) = \{i\}\} \quad (1)$$

$$CF_{i,j,p} = \{v \mid v \in V(G); d(G_p)_{v,i} < d(G_p)_{v,j}; G_p = G - p \quad (2)$$

Here,  $G_p = G - p$  is the spanning subgraph, resulted by deleting the path  $p$  joining the vertices  $i$  and  $j$  (except its endpoints),  $d(G)$  and  $d(G_p)$  denote the topological distances measured in  $G$  and  $G_p$ , respectively, and  $D(G)$  and  $\Delta(G)$  have the above mentioned meaning.

The sets  $CJ_{i,j,p}$  and  $CF_{i,j,p}$  represent subgraphs (connected or not) in  $G$ , related to the endpoint  $i$  and referred to  $j$  and path  $p$ .

In *Cluj fragmentation criteria*, the path  $p$  plays the central role in selecting the fragments. In cycle-containing graphs, more than one path could join the pair  $(i,j)$  thus resulting more than one fragment related to  $i$ , so that we define the nondiagonal entries  $[UM]_{ij}$  in the Cluj matrices as the maximum cardinality of the sets defined by eqs 1 or 2

$$[UM]_{ij} = \max_p |V_{i,j,p}| \quad (3)$$

where  $M = CJD$  (Cluj-Distance,  $p \in D(G)$ );  $CJA$  (Cluj-Detour,  $p \in \Delta(G)$ ),  $CFD$  (Cluj-Fragmental-Distance,  $p \in D(G)$ ) and  $CFA$  (Cluj-Fragmental-Detour,  $p \in \Delta(G)$ ), and  $|V_{i,j,p}|$  is the cardinality of the set  $CJ_{i,j,p}$  or  $CF_{i,j,p}$ . The diagonal entries are zero. The above definitions hold for any connected graph.

The Cluj matrices are square arrays, of dimension  $N \times N$ , usually *unsymmetric* (excepting some symmetric regular graphs). They can be symmetrized, e.g., by the Hadamard product with their transposes

$$SM_p = UM \bullet (UM)^T \quad (4)$$

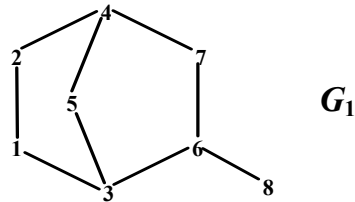
$$SM_e = SM_p \bullet A \quad (5)$$

The symbol  $\bullet$  indicates the Hadamard (pairwise) matrix product<sup>18</sup> ( $[M_a \bullet M_b]_{ij} = [M_a]_{ij} [M_b]_{ij}$ ). For the symmetric matrices, the letter  $S$  is usually missing. In eq 5, the Hadamard product between the path-defined matrix  $SM_p$  and the adjacency matrix  $A$  (i.e. the matrix having the non-diagonal entries unity for two adjacent vertices and zero otherwise) provides the corresponding edge-defined matrix,  $SM_e$ , which is a weighted adjacency matrix.

In trees,  $CJD$ ,  $CJA$ ,  $CFA$  and  $CFD$ , are identical, due to the uniqueness of the path joining a pair of vertices  $(i,j)$ . Some special properties of Cluj matrices were exposed elsewhere.<sup>10,12,15</sup>

As noted above, the sets  $CJ_{i,j,p}$  and  $CF_{i,j,p}$  represent subgraphs in  $G$ , either connected or not, related to the endpoint  $i$  and referred to  $j$  and path  $p$ . The connectivity of  $CF$  sets was

demonstrated elsewhere.<sup>19</sup> A case in which  $CJD_{i,j,p}$  is disconnected while  $CFD_{i,j,p}$  is connected is illustrated in Figure 1. Along with the Cluj matrices, the sets  $CJ_{i,j,p}$  and  $CF_{i,j,p}$  are presented as well.



$UCJD(G_1)$

	1	2	3	4	5	6	7	8
1	0	4	2	2	2	2	2	4
2	3	0	2	2	2	2	2	<b>3</b>
3	5	4	0	4	4	4	3	5
4	3	5	3	0	3	3	4	4
5	3	4	2	2	0	3	3	4
6	3	3	3	3	3	0	4	7
7	3	3	2	3	3	3	0	4
8	1	<b>1</b>	1	1	1	1	1	0

$UCFD(G_1)$

	1	2	3	4	5	6	7	8
1	0	4	2	2	2	2	3	5
2	3	0	2	2	2	3	2	<b>4</b>
3	5	4	0	4	4	4	3	6
4	3	5	3	0	3	3	4	5
5	5	5	2	2	0	4	4	5
6	3	4	3	3	3	0	4	7
7	3	3	2	3	3	3	0	6
8	1	<b>1</b>	1	1	1	1	1	0

	$CJD_{i,j,p}$	$CFD_{i,j,p}$
(2, 8) [2, 4, 7, 6, 8] {2,1,5} (disconnected)		
(2, 8) [2, 1, 3, 6, 8] {2, 4, 5}		
(8, 2) [8, 6, 3, 1, 2] {8}		
(8, 2) [8, 6, 7, 4, 2] {8}		

Figure 1. Unsymmetric Cluj matrices and fragmentation for the graph  $G_1$ .

UCJA( $G_1$ )									UCFA( $G_1$ )								
	1	2	3	4	5	6	7	8		1	2	3	4	5	6	7	8
1	0	1	1	1	1	2	1	2	1	0	1	1	1	1	2	1	2
2	1	0	1	1	1	1	2	1	2	1	0	1	1	1	1	2	1
3	2	2	0	3	<b>3</b>	2	2	2	2	2	2	0	3	<b>4</b>	2	2	2
4	2	2	2	0	2	2	2	3	2	2	2	0	4	4	2	2	3
<b>5</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>1</b>	<b>1</b>
6	3	2	2	2	2	0	2	7	3	2	2	2	2	2	0	2	7
7	1	3	1	1	1	1	0	1	1	3	1	1	1	1	1	0	1
8	1	1	1	1	1	1	1	0	1	1	1	1	1	1	1	1	0

$CJA_{i,j,p}$		$CFA_{i,j,p}$	
(3, 5) [3,1,2,4,5]		(3, 5) [3,1,2,4,5]	
{3,6,8}		(3, 5) [3,6,7,4,5]	
(3, 5) [3,6,7,4,5]		(3, 5) [3,6,7,4,5]	
{3,1}		{3,1,2}	
(5, 3) [5,4,2,1,3]		(5, 3) [5,4,2,1,3]	
{5}		{5}	
(5, 3) [5,4,7,6,3]		(5, 3) [5,4,7,6,3]	
{5}		{5}	

Figure 1. (continued)

An interesting property is shown by the detour-based matrices:  $CJA_p$  and  $CFA_p$ . Let consider the vertices 8 (of degree 1) and 5 (of degree 2) in  $G_1$ , Figure 1. The vertex 8 is an *external* vertex (with a terminal path ending in it) while the vertex 5 is an *internal* one (usually a terminal path not ending in it). An external vertex, like 8, shows all its entries in the Cluj matrices equal to 1 (see Figure 1). The same entries are shown by the internal vertex 5. This unusual property is called *the internal ending of all detours* joining a vertex  $i$  and the remaining vertices in  $G$ . Such a vertex is called an *internal endpoint*.<sup>15</sup> There exist graphs with all the vertices internal endpoints and their detours are *Hamiltonian paths* now. This kind of graph we call the *full Hamiltonian detour graph*,  $FHA$ .<sup>19</sup>

According to eq 3, sets of *maximal fragments*  $CJDS_{i,j}^M$ ,  $CJ\Delta S_{i,j}^M$ ,  $CFDS_{i,j}^M$  and  $CF\Delta S_{i,j}^M$  are expected in some graphs. This observation is important in the calculation of Cluj property indices (see the next section).

*Szeged Fragments*,  $SZD_{i,j}$  and  $SZ\Delta_{i,j}$  are defined by the equations:

$$SZD_{i,j} = \{v \mid v \in V(G); d(G)_{v,i} < d(G)_{v,j}\} \quad (6)$$

$$SZ\Delta_{i,j} = \{v \mid v \in V(G); \delta(G)_{v,i} < \delta(G)_{v,j}\} \quad (7)$$

These fragments represent the entries in the unsymmetric Szeged matrices, USZD and USZ\Delta.<sup>19</sup> In eq 5,  $\delta(G)_{v,i}$  represents the detour between the vertices  $i$  and  $v$ .

Note that in the definition of the Szeged fragments, the path joining the vertices  $i$  and  $j$  is irrelevant. Thus, in *Szeged fragmentation criteria*, **each pair  $(i,j)$  effects** one and only one fragment. It was demonstrated elsewhere<sup>19</sup> that the Szeged sets  $SZD_{i,j}$  represent connected subgraphs (i.e. fragments) while  $SZ\Delta_{i,j}$  are not necessary connected.

In any graph,  $CJD_e = CFD_e = SZD_e$ . In cyclic graphs,  $CJD_p \neq CFD_p \neq SZD_p$ ,  $CJ\Delta_p \neq CF\Delta_p \neq SZ\Delta_p$ .

## Fragmental Property Indices

### Model Parameters

It is well known that the physical laws govern the natural phenomena. Macroscopic interactions are interactions of field-type. The field is produced by a scalar function of potential. Let  $f(x, y, z)$  be such a scalar function. This function induces a field given in terms of the gradient of  $f$ :

$$\vec{\nabla} \cdot f = \left( \frac{\partial}{\partial x} \vec{i} + \frac{\partial}{\partial y} \vec{j} + \frac{\partial}{\partial z} \vec{k} \right) \cdot f(x, y, z) = \frac{\partial f}{\partial x} \vec{i} + \frac{\partial f}{\partial y} \vec{j} + \frac{\partial f}{\partial z} \vec{k} \quad (8)$$

For the potential of type

$$f(x, y, z) = pz \quad (9)$$

and applying eq 8 we obtain the associated field, in the form:

$$\begin{aligned}\bar{\nabla} \cdot \mathbf{f} &= \frac{\partial f}{\partial x} \bar{i} + \frac{\partial f}{\partial y} \bar{j} + \frac{\partial f}{\partial z} \bar{k} = \frac{\partial(\mathbf{p}z)}{\partial x} \bar{i} + \frac{\partial(\mathbf{p}z)}{\partial y} \bar{j} + \frac{\partial(\mathbf{p}z)}{\partial z} \bar{k} = \\ &= 0\bar{i} + 0\bar{j} + \mathbf{p}\bar{k} = \mathbf{p}\bar{k} = \bar{\mathbf{p}}\end{aligned}\quad (10)$$

This is the case of the well-known uniform gravitational field:

$$\bar{\mathbf{G}} = m\bar{\mathbf{g}} \quad (11)$$

the potential of which is given by

$$E_p = E_p(z) = mgz \quad (12)$$

where  $m$  is the mass of the probe and  $z$  is the reference coordinate.

Note that eq 10 is applicable not only to the Newtonian (gravitational) interactions but also to the Coulombian (electrostatic) interactions. In both cases the relation is valid if the mass  $m$  (or the charge  $q$ ) that generates the potential  $f$  and associated field  $\bar{\nabla} \cdot \mathbf{f}$  is far enough ( $r \gg z$ ) for the approximation  $(r+z)^2/r^2 = (r^2 + 2rz + z^2)/r^2 = 1 + 2z/r + (z/r)^2 \cong 1$  be applied in the equation of field produced by  $m$  or  $q$  (see below).

For the potential of type:

$$f(x, y, z) = p/z \quad (13)$$

eq 8 leads to the associated field:

$$\begin{aligned}\bar{\nabla} \cdot \mathbf{f} &= \frac{\partial f}{\partial x} \bar{i} + \frac{\partial f}{\partial y} \bar{j} + \frac{\partial f}{\partial z} \bar{k} = \frac{\partial(\mathbf{p}/z)}{\partial x} \bar{i} + \frac{\partial(\mathbf{p}/z)}{\partial y} \bar{j} + \frac{\partial(\mathbf{p}/z)}{\partial z} \bar{k} = \\ &= 0\bar{i} + 0\bar{j} + \frac{-\mathbf{p}}{z^2} \bar{k} = -\frac{\mathbf{p}}{z^2} \bar{k} = -\frac{\mathbf{p}}{z^3} \bar{z} = -\frac{\bar{\mathbf{p}}}{z^2}\end{aligned}\quad (14)$$

This is the case of well-known (non-uniform) gravitational field given by:

$$\bar{\mathbf{G}} = \bar{\mathbf{G}}(m, r) = -k \frac{m}{r^3} \bar{\mathbf{r}} \quad (15)$$

and the associated potential of the form:

$$U = U(m, r) = k \frac{m}{r} \quad (16)$$

where  $m$  is the mass of the probe and  $r$  is the position relative to the location of the point producing the field.

For the Coulombian field eq 14 becomes:

$$\bar{\mathbf{F}}_C = \bar{\mathbf{F}}_C(r) = -k \frac{q}{r^3} \bar{\mathbf{r}} \quad (17)$$

and the potential associated to the Coulombian field:

$$U = U(\mathbf{q}, \mathbf{r}) = k \frac{q}{r} \quad (18)$$

For *fragmental property indices* four models of interaction are implemented:<sup>19</sup> two of them are *topological* (dense topological and rare topological) and two others are *geometric* (dense geometric and rare geometric).

The *models* are related to two types of field interactions: one of *weak dependence on distance* for the potential of the type (9) generating a uniform field (10), and the second, of *strong dependence on distance* for the potential of the type (13) that generates a non-uniform field (14).

The variables in the models are *metrics of distance*  $\mathbf{d}$  (topological  $\mathbf{d}_T$  and geometric  $\mathbf{d}_E$ ), *properties*  $\Phi$  (mass  $M$ , electronegativity  $E$ , cardinality  $C$ , partial charge or any other atomic property  $P$ ), *property descriptors*  $\Omega(\mathbf{p}, \mathbf{d}, \mathbf{pd}, 1/\mathbf{p}, 1/\mathbf{d}, \mathbf{p}/\mathbf{d}, \mathbf{p}/\mathbf{d}^2, \mathbf{p}^2/\mathbf{d}^2)$  and method of *superposition*  $\Psi(\mathbf{S}, \mathbf{P}, \mathbf{A}, \mathbf{G}, \mathbf{H})$ .

Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be rational numbers; the (mathematical) *superposition* is

$$\Psi: \mathbf{S} = \sum_{i=1}^n \mathbf{x}_i; \mathbf{P} = \prod_{i=1}^n \mathbf{x}_i; \mathbf{A} = \mathbf{S}/n; \mathbf{G} = (\text{sgn}(\mathbf{P}))^n \cdot \sqrt[n]{\text{abs}(\mathbf{P})}; \mathbf{H} = \left( \sum_{i=1}^n \frac{1}{\mathbf{x}_i} \right)^{-1} \quad (19)$$

The expressions for the *property descriptors* are:

$$\Omega: \mathbf{p} = \mathbf{p}; \mathbf{d} = \mathbf{d}; \mathbf{pd} = \mathbf{p} \cdot \mathbf{d}; 1/\mathbf{p} = \frac{1}{\mathbf{p}}; 1/\mathbf{d} = \frac{1}{\mathbf{d}}; \mathbf{p}/\mathbf{d} = \frac{\mathbf{p}}{\mathbf{d}}; \mathbf{p}/\mathbf{d}^2 = \frac{\mathbf{p}}{\mathbf{d}^2}; \mathbf{p}^2/\mathbf{d}^2 = \frac{\mathbf{p}^2}{\mathbf{d}^2} \quad (20)$$

where  $\mathbf{p}$  is any property ( $\mathbf{p} \in \Phi$ ) and  $\mathbf{d}$  is any metric of distance.

These variables are most frequently used in our models by the following reasons:

- The expressions of the *property descriptor*  $\Omega$  simulate the most occurring physical interactions (e.g.  $\mathbf{p}, \mathbf{pd}, \mathbf{p}/\mathbf{d}, \mathbf{p}/\mathbf{d}^2, \mathbf{p}^2/\mathbf{d}^2$ )<sup>20</sup> and the most usual descriptors in topological and geometric models. The property descriptor is used either in the calculation of the *vertex descriptor* (when  $\mathbf{d}$  is the distance from the vertex  $\mathbf{v}$  to  $\mathbf{j}$  and  $\mathbf{p}$  is any atomic property) or in the evaluation of the *fragment descriptor* (when  $\mathbf{d}$  is the distance between the center of property of the fragment and  $\mathbf{j}$ , while  $\mathbf{p}$  is a calculated fragment property).

- The (mathematical) superposition is applied upon a string of vertex descriptors for giving a fragment descriptor. Note that  $\Psi$  means  $\mathbf{S} = \text{sum}$ ;  $\mathbf{P} = \text{product}$ ;  $\mathbf{A} = \text{arithmetic mean}$ ;  $\mathbf{G} = \text{geometric mean}$  and  $\mathbf{H} = \text{harmonic sum}$ . The summation is suitable in the case of any additive property (mass, volume, partial charges, electric capacities, etc.).<sup>21</sup> The multiplication

occurs in concurrent phenomena (probabilistically governed).<sup>22</sup> The arithmetic mean is useful in evaluating some mean contributions (corresponding to some uniform probabilistic distribution).<sup>23,24</sup> The geometric mean is used in calculating the group electronegativities.<sup>25,26</sup> Finally, the harmonic sum is important in connection with the elastic forces, electric fields and group mobility in viscous media.<sup>27</sup>

## Model Descriptions

Let  $(i,j)$  be a pair of vertices and  $Fr_{i,j}$  any fragment related to  $i$  and referred to  $j$ .

### *Dense Topological Model*

Let  $v$  be a vertex in the fragment  $Fr_{i,j}$ . The property descriptor applies to the vertex property  $p_v$  and topological distance  $d_{T v,j}$ . The fragmental *property descriptor PD*, resulting by the vertex descriptor superposition, gives the interaction of all the points belonging to the fragment  $Fr_{i,j}$  with the point  $j$ :

$$PD(Fr_{i,j}) = \Psi_{v \in Fr_{i,j}} (\Omega(d_{T v,j}, p_v)) \quad (21)$$

The  $j$  point can be conceived as an *internal probe atom* (see the *CoMFA* approach).<sup>28</sup> However, the chemical identity of  $j$  is not considered.

### *Rare Topological Model*

Within this model, the property descriptor applies to the fragmental property and topological distance  $d_{T i,j}$ . The fragmental property descriptor models the interaction of the whole fragment  $Fr_{i,j}$  with the point  $j$  and looks the global property being *concentrated* in the vertex  $i$ :

$$PD(Fr_{i,j}) = \Omega(d_{T i,j}, \Psi_{v \in Fr_{i,j}}(p_v)) \quad (22)$$

### *Dense Geometric Model*

The fragmental property descriptor is the vector sum of the vertex descriptor vectors. It applies the property descriptor to the vertex property  $p_v$  and the Euclidean distance  $d_{E,v,j}$  in providing a *point of equivalent (fragmental) property* located at the Euclidean distance



$d_{E,CP_i,j}$  (with  $d_{E,CP_i,j}$  being the distance between the center of fragmental property  $CP$ , of  $Fr_{i,j}$  and the vertex  $j$ ). The vector of the fragmental property has the orientation of this distance vector. The model simulates the interactions in non-uniform fields (gravitational, electrostatic, etc):

$$PD(Fr_{i,j}) = \left\| \sum_{v \in Fr_{i,j}} \bar{\Omega}(d_{E,v,j} p_v) \right\|; \quad \bar{\Omega} = \Omega \frac{\bar{d}_{E,v,j}}{d_{E,v,j}}; \quad P(Fr_{i,j}) = \Psi(p_v);$$

$$d_{E,CP_i,j} = \Omega_p^{-1}(DG(Fr_{i,j}), P(Fr_{i,j})), \quad (23)$$

where  $d_{E,CP_i,j}$  is the distance that satisfies:  $\Omega(d_{E,CP_i,j}, P(Fr_{i,j})) = PD(Fr_{i,j})$

#### Rare Geometric Model

The scalar fragmental descriptor applies the property descriptor to the center of fragmental property and Euclidean distance between this center and the vertex  $j$ .

The model simulates the interactions in uniform fields (uniform gravitational, electrostatic, etc.):

$$PD(Fr_{i,j}) = \Omega(d_{E,CP_i,j}, \Psi(p_v));$$

$$CP_i(x_{CP_i,j}, y_{CP_i,j}, z_{CP_i,j}); \quad x_{CP_i,j} = \frac{\sum_{v \in Fr_{i,j}} x_v \cdot p_v}{\sum_{v \in Fr_{i,j}} p_v} \quad (24)$$

$$y_{CP_i,j} = \frac{\sum_{v \in Fr_{i,j}} y_v \cdot p_v}{\sum_{v \in Fr_{i,j}} p_v}; \quad z_{CP_i,j} = \frac{\sum_{v \in Fr_{i,j}} z_v \cdot p_v}{\sum_{v \in Fr_{i,j}} p_v}$$

#### Some Particular Fragmental Property Models

Let  $i, j$  be two vertices in  $V(G)$  and  $Fr_{i,j}$  any fragment related to  $i$  with respect to  $j$ .

#### Fragmental Mass

In evaluating the fragmental mass, the chosen property is  $\Phi = M$ , descriptor  $\Omega = p$ , superposition  $\Psi = S$ , and the model is *rare topological*, **RT**. The fragmental mass descriptor takes the form:

$$PD(Fr_{i,j}) = \sum_{v \in Fr_{i,j}} M_v \quad (25)$$

It models the molecular mass of the fragment. The *name* of the associated property matrix is **RTcDdM\_\_p\_\_S**, with the known meaning for *c* and *Dd*.

If  $c = s$  and  $Dd = Di$  then **RTsDiM\_\_p\_\_S**, it models the molecular mass of the Szeged Distance Fragments (eq 6). If  $c = f$  and  $Dd = Di$  then the matrix **RTfDiM\_\_p\_\_S** collects mean values of mass of all the fragments belonging to *i* (with respect to *j*) according to the *CF* criterion (eqs 2 and 3).

### *Fragmental Electronegativity*

The well known equalizing principle of electronegativity *E*, is here considered: the fragment electronegativity is the geometric mean of electronegativities of the *s* atoms joined to form that fragment.

Let the property  $\Phi = E$  (electronegativity); descriptor  $\Omega = p$ ; superposition  $\Psi = G$ ; the model is *rare topological*, *RT*. The fragmental electronegativity descriptor is:

$$PD(Fr_{i,j}) = |Fr_{i,j}| \sqrt{\prod_{v \in Fr_{i,j}} E_v} \quad (26)$$

It models the electronegativity of the fragment  $Fr_{i,j}$ . The *name* of the property matrix associated with it is **RTcDdE\_\_p\_\_G**. Note that  $E_v$  is the group electronegativity for vertex *v* calculated with formula:

$$E_v = \sum_{j \in \Gamma_v} b(v,j) \sqrt{\prod_{j \in \Gamma_v} E_a^{b(v,j)}} \quad (27)$$

where  $b(v, j)$  is the *conventional bond order* between *v* and *j* (e.g. 1, 1.5, 2, 3 for single, aromatic, double and triple bonding, respectively),  $E_a$  is the atomic electronegativity (Sanderson) and  $j \in \Gamma_v$  is any atom (hydrogen atoms included) consisting the group  $\Gamma_v$ .

### *Fragmental Numbers*

The property  $\Phi = C$  (cardinality) was introduced for recovering some graph-theoretical quantities and/or graph theoretical analogue indices (see below).

For descriptor  $\Omega = p$ , superposition  $\Psi = \{P, A, G\}$ , and the model *rare topological*, *RT*, the cardinal numbering descriptor of  $Fr_{i,j}$  is:

$$PD(Fr_{i,j}) = \prod_{v \in Fr_{i,j}} 1 = \frac{\sum_{v \in Fr_{i,j}} 1}{|Fr_{i,j}|} = |Fr_{i,j}| \sqrt{\prod_{v \in Fr_{i,j}} 1} = 1 \quad (28)$$

The arithmetic mean  $A$ , geometric mean  $G$  and product  $P$  applied to  $\mathbf{1}$  (value for vertex property) leave it unchanged. The mean value for all fragments belonging to  $i$  vs.  $j$  ( $CJ$  and  $CF$  only) is also 1. All matrices  $\mathbf{RTcDdC\_p\_P}$ ,  $\mathbf{RTcDdC\_p\_A}$  and  $\mathbf{RTcDdC\_p\_G}$  have all their entries unity, except the main diagonal elements that are zero.

The corresponding path-calculated indices give the number of edges in the complete graph having the same number of vertices  $N$ , as the considered molecular graph:

$$\mathbf{RTcDdC\_p\_PP\_} = \mathbf{RTcDdC\_p\_AP\_} = \mathbf{RTcDdC\_p\_GP\_} = N(N-1)/2$$

Similarly, the edge-calculated indices,  $\mathbf{RTcDdC\_p\_PE\_}$ ,  $\mathbf{RTcDdC\_p\_AE\_}$ ,  $\mathbf{RTcDdC\_p\_GE\_}$  give the number of edges in the molecular structure.

Let now the property  $\Phi = C$ , descriptor  $\Omega = p$ , superposition  $\Psi = S$  and *rare topological* model. The value of cardinal numbering descriptor for  $Fr_{i,j}$  is:

$$PD(Fr_{i,j}) = \sum_{v \in Fr^{\sigma_t}_{i,j}} 1 = |Fr^{\sigma_t}_{i,j}| \quad (29)$$

It models the number of atoms in the fragment. The associated matrices are of the form  $\mathbf{RTcDdC\_p\_S}$ :

$$\mathbf{RTfDiC\_p\_S} = \mathbf{CFD}; \quad \mathbf{RTfDeC\_p\_S} = \mathbf{CFA}$$

$$\mathbf{RTjDiC\_p\_S} = \mathbf{CJD}; \quad \mathbf{RTjDeC\_p\_S} = \mathbf{CJA}$$

$$\mathbf{RTsDiC\_p\_S} = \mathbf{SZD}; \quad \mathbf{RTsDeC\_p\_S} = \mathbf{SZ\Delta}$$

and the corresponding indices are exactly the graph-theoretical descriptors corresponding to the Cluj and Szeged criteria (eqs 1-3, 6, 7).

#### *Uniform Field Gravity*

Let the property  $\Phi = M$ , descriptor  $\Omega = p/d^2$ , superposition  $\Psi = S$  and *rare geometrical* model.

The uniform gravity descriptor of  $Fr_{i,j}$  is calculated by:

$$PD(Fr_{i,j}) = \sum_{v \in Fr_{i,j}} \frac{M_v}{d_{v,j}^2} \quad (30)$$

It models the value of the gravitational field induced by the fragment  $Fr_{i,j}$  in the point  $j$ . Values given by eq 30 are collected in the matrix  $RGsDdM\_p/d2S$  while averaged values are considered in  $RGfDdM\_p/d2S$  and  $RGjDdM\_p/d2S$  matrices.

#### Non-Uniform Field Gravity

Let the property  $\Phi = M$ , descriptor  $\Omega = p/d2$ , superposition  $\Psi = S$  and *dense geometrical* model. The *distance* (vs.  $j$ ) of the *center of equivalent fragmental gravity* of  $Fr_{i,j}$  is:

$$d_{E,CP_{i,j}} = \sqrt{\left( \sum_{v \in Fr_{i,j}} M_v \right) / \left( \left( \sum_{v \in Fr_{i,j}} \frac{M_v \cdot \bar{d}_{v,j}}{d_{v,j}^2} \right) \cdot \left( \sum_{v \in Fr_{i,j}} \frac{M_v \cdot \bar{d}_{v,j}}{d_{v,j}^2} \right) \right)^{\frac{1}{2}}} \quad (31)$$

It models the distance at which a point mass equal to the fragment mass  $\sum_{v \in Fr_{i,j}} M_v$  should

be located vs.  $j$  such that the gravitational field induced by  $Fr_{i,j}$  in  $j$  be equal to the field induced by all atoms of the fragment. The associated matrix is of the form  $DGcDdM\_p/d2S$ .

#### Uniform Electrostatic field

Let the property  $\Phi = P$  ( $Q_P$  implicitly, in the Cluj Program), descriptor  $\Omega = p/d2$ , superposition  $\Psi = S$  and *rare geometrical* model. The uniform electrostatic field descriptor of  $Fr_{i,j}$  is:

$$PD(Fr_{i,j}) = \sum_{v \in Fr_{i,j}} \frac{Q_{Pv}}{d_{v,j}^2} \quad (32)$$

It models the value of electrostatic field induced by the fragment in  $j$ . The property matrix is of the form:  $RGcDdP\_p/d2S$ .

#### Non-Uniform Electrostatic Field

For the property  $\Phi = P$  ( $Q_P$  implicitly), descriptor  $\Omega = p/d2$ , superposition  $\Psi = S$  and *dense geometrical* model, the *distance* (vs.  $j$ ) of the *center of equivalent electrostatic field* of  $Fr_{i,j}$  is:

$$d_{E,CP_i,j} = \sqrt{\frac{\sum_{v \in Fr_{i,j}} Q_{P_v}}{\left( \left( \sum_{v \in Fr_{i,j}} \frac{Q_{P_v}}{d_{v,j}^2} \cdot \frac{\bar{d}_{v,j}}{d_{v,j}} \right) \cdot \left( \sum_{v \in Fr_{i,j}} \frac{Q_{P_v}}{d_{v,j}^2} \cdot \frac{\bar{d}_{v,j}}{d_{v,j}} \right) \right)^{\frac{1}{2}}} \quad (33)$$

It models the distance at which a point charge equal to the fragment charge  $\sum_{v \in Fr_{i,j}} Q_{P_v}$

be located vs.  $j$  such that the electrostatic field induced by it in  $j$  be equal to the field induced by the all atoms of the fragment. The associated matrix is of the form: **DGcDdP\_p/d2S**.

#### *Uniform Field Gravitational Potential*

It is obtained for the property  $\Phi = M$ , descriptor  $\Omega = p/d$ , superposition  $\Psi = S$  and *rare geometrical* model. The property descriptor of  $Fr_{i,j}$  is:

$$PD(Fr_{i,j}) = \sum_{v \in Fr_{i,j}} \frac{M_v}{d_{v,j}} \quad (34)$$

It models the value of the gravitational potential induced by the fragment in  $j$ . The property matrix is of the form: **RGcDdM\_p/d\_S**.

#### *Non-Uniform Field-Type Gravitational Potential*

For the property  $\Phi = M$ ; descriptor  $\Omega = p/d$ ; superposition  $\Psi = S$ ; *dense geometrical model*, the distance (vs.  $j$ ) of the *center of equivalent fragmental gravity* of  $Fr_{i,j}$  is:

$$d_{E,CP_i,j} = \sqrt{\frac{\left( \sum_{v \in Fr_{i,j}} M_v \right)}{\left( \left( \sum_{v \in Fr_{i,j}} \frac{M_v}{d_{v,j}} \cdot \frac{\bar{d}_{v,j}}{d_{v,j}} \right) \cdot \left( \sum_{v \in Fr_{i,j}} \frac{M_v}{d_{v,j}} \cdot \frac{\bar{d}_{v,j}}{d_{v,j}} \right) \right)^{\frac{1}{2}}} \quad (35)$$

It models the distance at which a point mass equal to the fragment mass ( $\sum_{v \in Fr_{i,j}} M_v$ ) should

be located vs.  $j$  such that the gravitational potential induced by it in  $j$  be equal to the potential induced by the all atoms of the fragment. The associated matrix is of the form **DGcDdM\_p/d\_S**.

#### *Uniform Field Coulombian Potential*

It is obtained for the property  $\Phi = P$  ( $Q_p$  implicitly), descriptor  $\Omega = p/d$ , superposition  $\Psi = S$  and *rare geometrical* model. The electrostatic potential descriptor of  $Fr_{i,j}$  is:

$$PD(Fr_{i,j}) = \sum_{v \in Fr_{i,j}} \frac{Q_{Pv}}{d_{v,j}} \quad (36)$$

It models the value of the electrostatic potential induced by the fragment in  $j$ . The property matrix is of the form: **RGcDdP\_p/d\_S**.

#### *Non-Uniform Field Electrostatic Potential*

For the property  $\Phi = P$  ( $Q_P$  implicitly); descriptor  $\Omega = p/d$ ; superposition  $\Psi = S$  and dense geometrical model, the distance (vs.  $j$ ) of the center of equivalent electrostatic potential of  $Fr_{i,j}$  is:

$$d_{E,CP_{i,j}} = \sqrt{\left( \sum_{v \in Fr_{i,j}} Q_{Pv} \right) / \left( \left( \sum_{v \in Fr_{i,j}} \frac{Q_{Pv}}{d_{v,j}} \cdot \frac{\vec{d}_{v,j}}{d_{v,j}} \right) \cdot \left( \sum_{v \in Fr_{i,j}} \frac{Q_{Pv}}{d_{v,j}} \cdot \frac{\vec{d}_{v,j}}{d_{v,j}} \right) \right)^{\frac{1}{2}}} \quad (37)$$

It models the distance at which a point charge equal to the fragment charge ( $\sum_{v \in Fr_{i,j}} Q_{Pv}$ )

should be located vs.  $j$  such that the electrostatic potential induced by it in  $j$  be equal to the potential induced by all the atoms of the fragment. The associated matrix is of the form **DGcDdP\_p/d\_S**.

In all the above models,  $j$  appears as a *virtual probe atom*. In the opposite to the CoMFA approach, whose descriptors are calculated as interactions of the molecule with external grid probe atoms, our approach makes use of internal probe atoms: the property of fragment  $Fr_{i,j}$  is viewed as the interaction of atoms forming the fragment  $Fr_{i,j}$  with the atom  $j$  (with no chemical identity, however). Other particular fragmental property models the reader can find in ref.<sup>19</sup>

### **Fragmental Property Matrices**

The fragmental property matrices are square matrices of the order  $N$  (i.e. the number of non-hydrogen atoms in the molecule). The non-diagonal entries in such matrices are fragmental properties corresponding to any pair of vertices ( $i,j$ ) by a chosen model.

In case of Cluj criteria, the fragmentation can supply more than one maximal fragment for the pair ( $i,j$ ). In such a case, the matrix entry is the arithmetic mean of the individual values.

Thus, if  $i, j \in V(G)$ ,  $i \neq j$  and  $P_{ij} = \{p_{i,j}^1, p_{i,j}^2, \dots, p_{i,j}^k\}$  paths joining  $i$  and  $j$ , then cf. *CJ* or *CF* definition (eqs 1-3), the fragments  $Fr_{i,j}^1, Fr_{i,j}^2, \dots, Fr_{i,j}^k$  are generated. Let  $m$  be the number of maximal fragments among all the  $k$  fragments,  $1 \leq m \leq k$ , and let  $\sigma_1, \dots, \sigma_m$  be the index for the maximal fragments. By applying any of the above models (eqs 21-29) for all the  $m$  maximal fragments we obtain  $m$  values (for example, by eq 29):

$$PD(Fr_{i,j}^{\sigma_1}), PD(Fr_{i,j}^{\sigma_2}), \dots, PD(Fr_{i,j}^{\sigma_m})$$

and consequently, the matrix entry associated to the pair  $(i, j)$  is the mean value:

$$PD_{i,j} = \frac{\sum_{t=1}^m PD(Fr_{i,j}^{\sigma_t})}{m} \quad (38)$$

The resulting matrices are in general *unsymmetric* but they can be symmetrized (see eqs 4, 5). The symbols for the fragmental property matrices will be detailed below.

### Operators for Calculating Fragmental Property Indices

Fragmental property indices are calculated at any fragmental property matrices above discussed, by applying four types of index operators:  $P_-$ ,  $P2$ ,  $E_-$ ,  $E2$  according to the relations:

$$\begin{aligned} P_-(M) &= \frac{1}{2} \sum \sum [M]_{ij} & ; & & P2(M) &= \frac{1}{2} \sum \sum [M]_{ij} [M]_{ji}; \\ E_-(M) &= \frac{1}{2} \sum \sum [M]_{ij} [A]_{ij} & ; & & E2(M) &= \frac{1}{2} \sum \sum [M]_{ij} [M]_{ji} [A]_{ij} \end{aligned} \quad (39)$$

where  $M$  is any property matrix, symmetric or unsymmetric.

### Name of the Fragmental Property Matrices and Indices

The name of *fragmental property matrices* is of the general form:

$$ABcDdEffffG \quad (40)$$

where:

$A \in \{D, R\}$ ;  $D$  = Dense;  $R$  = Rare;

$B \in \{T, G\}$ ;  $T$  = Topological;  $G$  = Geometric;

$c \in \{f, j, s\}$ ;  $f$  = *CF*-type;  $j$  = *CJ*-type;  $s$  = *Sz*-type;

$Dd \in \{Di, De\}$ ;  $Di$  = Distance;  $De$  = Detour;

$E \in \Phi$  (i.e.  $E \in \{M, E, C, P\}$  where  $M = mass$ ;  $E = electronegativity$ ;  $C = cardinality$ ;  $P =$  other atomic property - implicitly, *partial charge*; explicitly, a property given by manual input);

$ffff \in \Omega$  (i.e.  $ffff \in \{\_p\_, \_1/p\_, \_d\_, \_1/d\_, \_p.d\_, \_p/d\_, \_p/d2, p2/d2\}$ )

$G \in \Psi$  (i.e.  $G \in \{S, P, A, G, H\}$  with the known meaning (see above).

The name of *fragmental property indices* is of the general form:

$$ABcDdEffffGii \quad (41)$$

where:

$ii \in \{P\_ , P2, E\_ , E2\}$  with the known meaning (eq 39).

If an operator, such as  $f(x) = 1/x$  (inverse operator) or  $f(x) = \ln(x)$ , is applied the indices are labeled as follows:

$\ln ABcDdEffffGii := \ln(ABcDdEffffGii);$

$$1/ABcDdEffffGii := \frac{1}{ABcDdEffffGii} \quad (42)$$

For example, index  $\ln DGfDeM\_p\_SP\_$  is the logarithm of index  $DGfDeM\_p\_SP\_$  computed on the property matrix  $DGfDeM\_p\_S$ . The model used is dense, geometric, on fragment of type  $CF$ , with the cutting path being detour. The chosen property is the mass, the descriptor for property is even the property (mass) and the sum operator counts the vertex descriptors.

### *Model Degeneration and Computational Features*

The degeneration in the above models may occur in cases when the values of property are not diverse enough, like is case of cardinality (see *Fragmental Numbers*, this Section).

Another degeneration is in the case:  $RTfDiC\_p\_H = TfDiC\_1/p\_S$ .

The fragmental analysis was made by the aid of four original 16-bit windows computer programs. First program, **ClujTheor** calculates *topological descriptors* of *Cluj* and *Szeged* type and generates the fragments for the molecules. Second, **ClujProp** calculates the *fragmental properties*. The third one, **StatMon** makes *monovariate regressions* and sorts indices according to the correlation score. The forth program, **StatQ** performs *multi-linear regression* (2-variate, 4-variate, etc.) and saves on disk the best couples of indices. The total number of indices is given by:  $2560 \times 3$  (i.e.  $x, \ln(x), 1/x$ )  $\times 3$  (i.e. the criteria:  $CJ, CF, Sz$ )  $\times 2$  (i.e. the path criteria:  $Di, De$ ) = 46080. Note that in most cases, the degeneration induced by



property values and operators lead to a total number of distinct indices around 19,000. In bivariate regression, the first  $2^{14}-1=16383$  indices recording the best scores in monovariate regression are considered.

### Correlating Studies

The mathematical models of a certain property are performed by MLR (Multiple Linear Regression) and/or CNN (Computational Neural Networks).<sup>29-39</sup> In any case, the model is built by using a training set of structures that provides a calibration equation. Next, it is validated by a cross-validation procedure and also by using an external prediction set. In the following, the MLR procedure is presented.

**MLR**, for  $n$  observations and  $m$  independent variables is represented by

$$Y_i = b_0 + \sum_j^m b_{ij} X_{ij} \quad (43)$$

or, in matrix form as

$$\mathbf{Y} = \mathbf{bX} \quad (44)$$

where  $\mathbf{Y}$  is the  $n \times 1$  vector of responses,  $\mathbf{X}$  is an  $n \times (m + 1)$  matrix of independent variables and  $\mathbf{b}$  is the  $(m + 1) \times 1$  vector of regression coefficients. The regression coefficients can be determined by the least-squares solution of (44)

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad (45)$$

With  $\mathbf{b}$  calculated, eq 44 can be used for estimating the chosen property for other chemical structures.

To avoid the chance correlations, it is recommended that the number of descriptors submitted to regression be less than 60 % of the number of observations in the training set.<sup>40</sup>

#### Set 1. Substituted 3-(Phthalimidoalkyl)-pyrazolin-5-ones.

We tested the correlating ability of **FPIF** on a set of 17 molecular structures from the class of substituted 3-(Phthalimidoalkyl)-pyrazolin-5-ones.

The molecular structure of the selected chemicals is given in Figures 2.(a, b). It was performed by using the MM+ ( for 3D-geometries) and semiempirical AM1 (for partial charge calculation) procedures of the HyperChem Program (HyperCube Inc.).

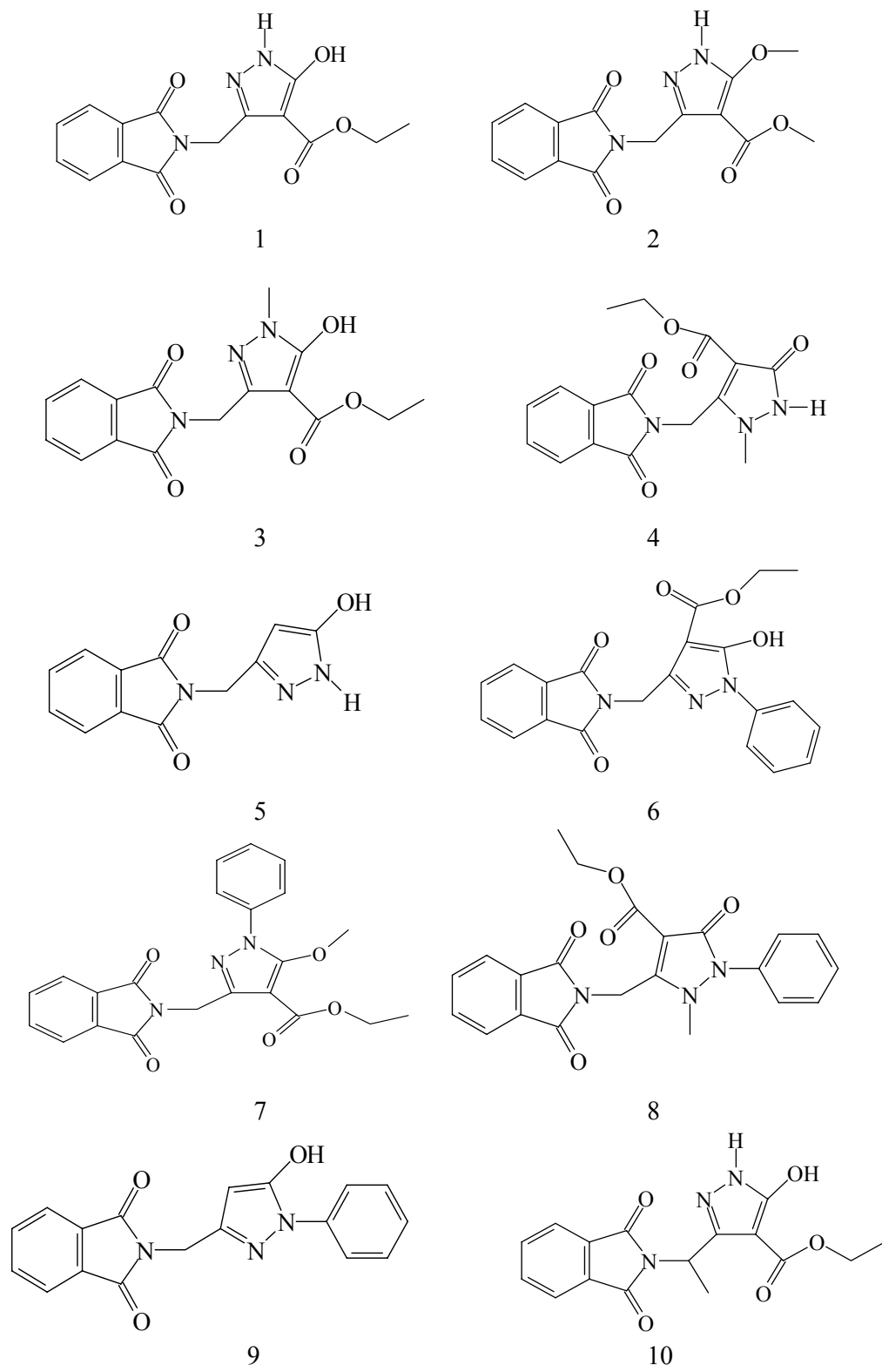


Figure 2.a. Structure of 17 substituted 3-(Phthalimidoalkyl)-pyrazolin-5-ones; molecules 1 to 10.

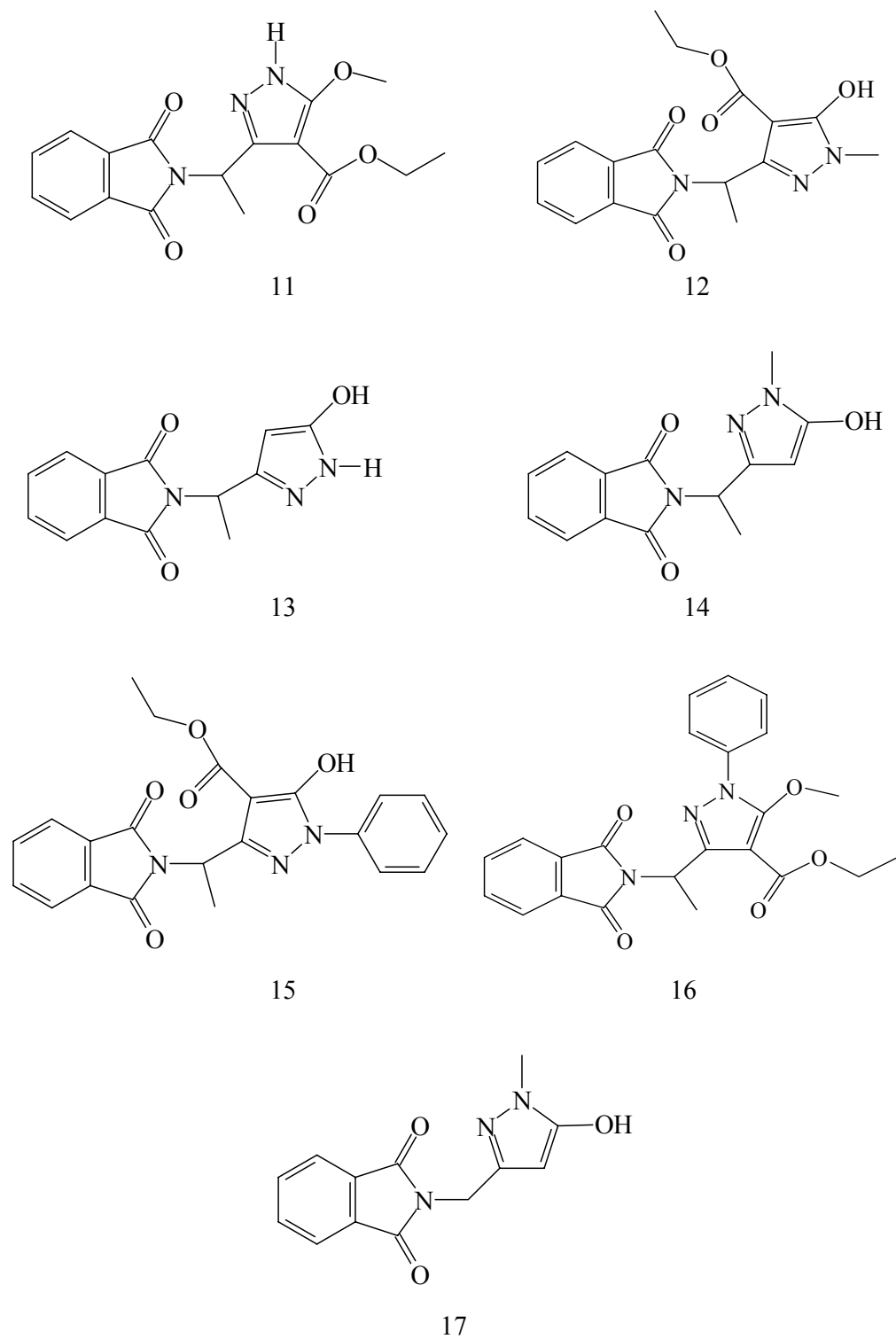


Figure 2.b. Structure of 17 substituted 3-(Phthalimidoalkyl)-pyrazolin-5-ones; molecules 11 to 17.

The modeled properties were the sum of one-electron energy calculated at the Extended-Huckel level and the inhibitory activity (in %) of a solution of 0.05 mg/ml pyrazolin-5-one on *Lepidium sativum* L. (Cresson). The data are listed in Table 1.

Table 1. The Sum of One-Electron Energy Calculated at Single Point Semi-Empirical Extended-Huckel and the Inhibitory Activity on *Lepidium sativum* L. (Cresson) for 17 Substituted 3-(Pthalimidoalkyl)-Pyrazolin-5-Ones\*

Molecule no.	Energy (kcal/mol)	Inhibition (%)
1	50978.19	28.4
8	64751.09	65.2
7	64752.65	49.4
6	62330.33	68.3
5	38604.68	14.3
4	53416.95	27.7
3	53441.43	30.4
2	51000.36	28
17	41057.46	15.1
16	67104.64	50.6
15	64701.39	71.7
14	43473.37	18.2
13	41020.54	12.2
12	55832.12	32.6
11	55729.99	28.9
10	53424.19	29.3
9	50012.42	46.9

\* Values of inhibition are taken from ref. <sup>41</sup>

#### *Monovariate Regression for Energy*

For the first six best indices in monovariate correlation, the equation is:

$$\text{Predicted energy} = b_0 + b_1 \cdot \ln \text{Index} \quad (46)$$

The indices are listed in Table 2 along with the Pearson correlation index R and the regression coefficients.

Table 2. The Bests Six Correlations of Energy in Monivariate and Divariate Regression

Index No.	Index Name	R	b <sub>0</sub>	b <sub>1</sub>
<b>1</b>	<b>lnDGjDeE_p/d2PE_</b>	<b>0.99973</b>	5370	3760
2	DTjDeEp2/d2SE_	0.99964	8802.3	138
3	lnDGjDeE_p/d2PE2	0.99964	1289.5	3671.1
4	DTjDeMp2/d2SE_	0.99939	9188.7	922.34
5	lnDGsDiM_p/d2PE_	0.99938	-59041	28397
6	lnRTjDeE_p*d_PP_	0.99927	3044	1938.6
1	lnDGjDeE_p/d2PE_	0.999688	8771.8	-21.648
2	DTjDeEp2/d2SE_			138.77
1	lnDGjDeE_p/d2PE_	0.999871	58810	3678
10175	lnDTsDeE_p/d_PP2			-6593.1
4	DTjDeMp2/d2SE_	0.999902	13056	1108.8
4315	DTfDiE_p/d_AP2			-95.598
34	RTsDiM_p/d2GP2	0.999969	-1193	1674.3
5947	DTjDeEp2/d2AP			-41.168
492	RTjDeM_p/d2SP_	0.999974	58267	46.095
1698	1/RTsDeM_p/d2AE2			-686800
<b>492</b>	<b>RTjDeM_p/d2SP_</b>	<b>0.999981</b>	56222	47.864
<b>1737</b>	<b>1/RTsDeM_p/d2AP2</b>			-711240

The best single variable QSPR (boldface in Table 2) was

$$\text{Predicted energy} = 5370 - 3760 * \ln DGjDeE\_p/d^2PE\_ \quad (47)$$

$$R = 0.99973; n = 17$$

This correlation could be satisfactory but usually a molecular property shows more than one dimension dependency. For this reason, we performed the bivariate regression.

#### *Bivariate Regression for Energy*

The first 16383 indices, labeled in decreasing order of their score in monivariate regression, are submitted for bivariate correlation. A procedure for finding subsets of optimal even number descriptors was developed. It is a simple, iterative technique that eludes the investigation of all possible descriptor combinations and reduces the time for drawing the best property model. More details will be presented in a future paper.

Here, the bivariate correlation for six pairs of indices is exemplified. The pairs are: (1, 2); (1, 10175); (4, 4315); (34, 5947); (492, 1698) and (492, 1737). The first two pairs are taken to show that the first scored index in monivariate regression does not provide the best

bivariate correlation. Selection of the pairs of indices for bivariate correlation must be done by traversing the whole pool (1...16383). For additional descriptors, our procedure for optimum descriptor selection avoid the mining of all possible index combinations.

The best bivariate score was provided by the pair (492, 1737) (Figure 3):

$$\begin{aligned} \text{Predicted energy} &= 56221.885 + 47.864 \cdot \text{RTjDeM\_p/d2SP\_} \\ &\quad - 711240.703 \cdot 1/\text{RTsDeM\_p/d2AP2} \\ R &= 0.99998; \quad s = 57.40; \quad n = 17 \end{aligned} \quad (48)$$

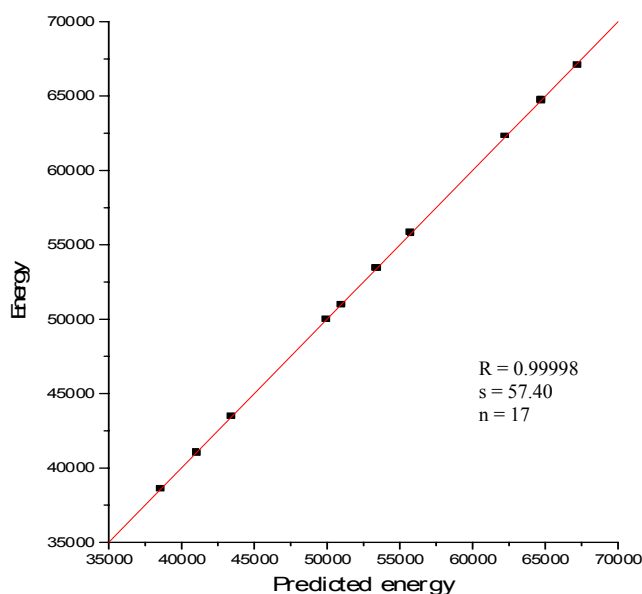


Figure 3. The plot: energy vs predicted energy of eq 48

An insight in Table 2 reveals that the best models (i.e. those showing  $R > 0.9999$ ) show a dependency of this energy by the molecular topology (topological models) and the nature of atoms (mass and electronegativity).

#### *Monovariate Regression for Inhibition*

For the first six best indices in monovariate regression, the equation of the model is:

$$\text{Predicted inhibition} = b_0 + b_1 \cdot \text{Index} \quad (49)$$

for which the indices and statistics are given in Table 3.

Table 3. The Best Six Correlations of Inhibition in Monivariate and Divariate Regression

Index No.	Index Name	R	b <sub>0</sub>	b <sub>1</sub>
1	lnDGsDeC_1/p_SE_	0.95389	-336.76	96.378
2	1/DGsDeC_1/p_SE_	-0.95231	137.01	-4754.8
3	lnDTjDeE_p*d_HE	0.95174	135.80	-493.02
4	1/DTjDeE_p*d_HE	-0.95169	175.92	-6826.3
5	1/DTjDeEp2/d2AE	-0.95144	186.07	-19546
6	DGsDeC_1/p_SE_	0.95126	-56.53	1.9019
1	lnDGsDeC_1/p_SE_	0.955154	-129.47	54.37
2	1/DGsDeC_1/p_SE_			-2121.1
1	lnDGsDeC_1/p_SE_	0.96935	-84.14	47.481
1369	1/RGsDeMp2/d2SE2			-3743200
2	1/DGsDeC_1/p_SE_	0.983967	-452.74	-4279.8
13227	lnRGsDeE_p/d_AE2			118.74
18	DTjDeEp2/d2AE_	0.988316	121.21	1.076
16842	RGjDeP_p/d_GP			-1.5194
37	DTjDeE_p/d_AE_	0.990564	-73.183	2.1644
11362	lnDGjDeP_p/d_PE2			-4.1769
4304	DTsDiM_p*d_HP_	0.99268	-26.846	1.5619
7649	DGjDeE_p/d2SE2			-1.7043

The best monivariate QSAR was

$$\text{Predicted inhibition} = -336.760 + 96.378 * \ln DGsDeC_{1/p\_SE\_} \quad (50)$$

$$R = 0.9539; n = 17$$

which is, of course, not satisfactory, despite in ref.<sup>41</sup> a value of  $R = 0.92$  was reported. Thus, we performed the bivariate regression.

#### *Bivariate Regression for Inhibition*

Six pairs of indices are considered here for bivariate correlation: (1, 2); (1, 1369); (2, 13227); (18, 16842); (37, 11362) and (4304, 7649).

As in the case of energy, the best scored index in monivariate correlation is not present in the pair of best bivariate correlation.

The best bivariate score was done by the pair (4304, 7649):

$$\text{Predicted inhibition} = -26.846 + 1.562 * DTsDiM_{p*d\_HP\_} - 1.704 * DGjDeE_{p/d2SE2}$$

$$R = 0.9927; s = 2.374; n = 17. \quad (51)$$

Figure 4. illustrates the plot of inhibition vs predicted inhibition of eq 51.

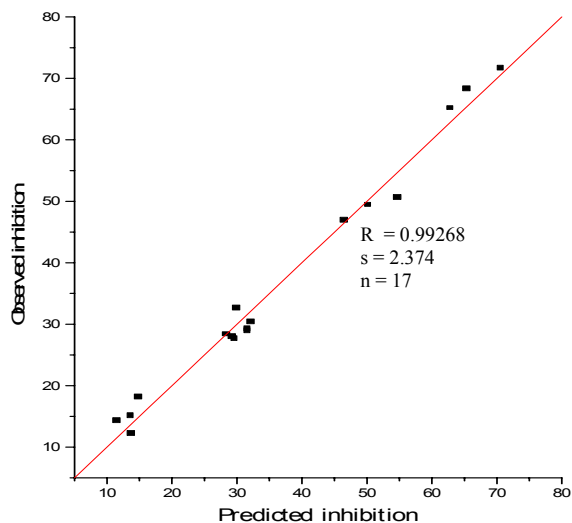


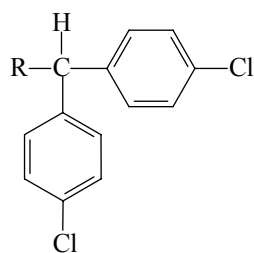
Figure 4. The plot: inhibition vs predicted inhibition of eq 51

The constant high correlation (see Table 3) between the best indices and the mitodepressive activity on *Lepidium Savitium L. (Cresson)* demonstrate ability of this family of indices to estimate the biological activity of the considered set of chemical structures. The models with  $R > 0.983$  suggest that the mitodepressive activity on *Lepidium Savitium L. (Cresson)* is dependent both on the geometric and topological features of molecules, the nature of atoms (mass and electronegativity) and the electrostatic field of atoms induced by their partial charges.

## Set 2. Aromatase Inhibitors.

A set of substituted dichlorodiphenyls (4, 4'-dichlorodiphenyl-methanes) inhibitors of aromatase<sup>42</sup> were considered. Enzymatic aromatization of androgens is involved in the biosynthesis of estrogens, and consequently in the estrogen-dependent diseases.



**Table 4.** Dichlorodiphenyl Methanes Aromatase Inhibitors.

No.	R	$-\log EC_{50}$ obs
1		7.43
2		8.03
3		8.06
4		5.70
5		5.71
6		5.30
7		5.30
8		6.80
9		5.30
10		7.26

For modeling the inhibition, the authors<sup>42</sup> used two dipole moment related descriptors. We modeled the inhibition in monivariate regression but no satisfactory correlation ( $R^2$  around 0.828) was found. In divariate regression, the correlation improved (Figure 5).

$$\text{Predicted inhibition} = 6.177 + 0.513 \cdot \ln \text{RTjDiP\_p\_HP2} - 0.071 \cdot 1/\text{DGjDeP\_1/p\_SP2}$$

$$R^2 = 0.9716; s = 0.205; n = 10 \quad (52)$$

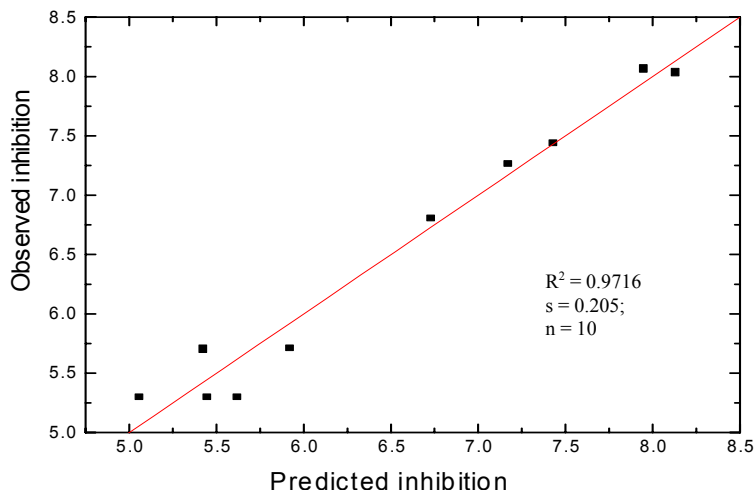


Figure 5. The plot of predicted vs observed inhibition of aromatase.

The best reported<sup>42</sup> correlation for this subset was:  $R^2 = 0,89$ ;  $s = 0.44$ . In our model, both the topology and geometry (see the indices in eq 52) are important in modeling the aromatase inhibition by dichlorodiphenyl methanes.

### Set 3. N-containing compounds

A set of 90 N-containing compounds (Table 5) of industrial importance was taken from the paper.<sup>43</sup> The tested property was the normal boiling point, B.P. The authors modeled this property by using four categories of molecular descriptors: topological, geometric, electronic and charged-partial surface area descriptors (CPSA).<sup>44,45</sup> The nitrogen-containing compounds were problematic in modeling a diverse set of organic chemicals, so that the authors excluded such compounds from their initial model.

The best found MLR model involved ten descriptors (1. dipole moment; 2. partial negative surface area; 3. relative negative charge; 4. relative negative charged surface area; 5. number of aromatic bonds; 6. path 2 molecular connectivity index; 7. cluster 3 valence connectivity index; 8. sum of all path weights from heteroatoms; 9. surface area of donatable hydrogens and 10. charge of donatable hydrogens) and showed the following statistics:  $n = 90$

**Table 5.** Nitrogen-Containing Compounds and Their Boiling Points.

No.	Compound	BP	No.	Compound	BP
1.	2-ethylpyridine	422.2	46.	n-tetradecylamine	564.5
2.	2-ethylpiperidine	416.2	47.	acridine	619.2
3.	1-ethylpiperidine	404.2	48.	tri-n-butylamine	487.2
4.	2,2-dimethyl-1,3-diaminomethane	426.2	49.	n-dodecylamine	532.4
5.	N,N-dimethyl-1,3-diaminomethane	418.2	50.	diamylamine	476.1
6.	3,3-dimethylpiperidine	410.2	51.	tripropylamine	429.7
7.	p-fluorobenzylamine	456.2	52.	n-nonylamine	475.4
8.	cianogene	252	53.	quinoline	510.8
9.	m-bromoaniline	524.2	54.	acetonitrile	354.8
10.	o-bromoaniline	502.2	55.	isoquinoline	516.4
11.	N-ethylbutylamine	381.2	56.	n-octylamine	452.8
12.	triethylamine	362	57.	indole	526.1
13.	N,N-diethylamină	337.2	58.	n-heptylamine	430.1
14.	o-nitrotoluene	498.2	59.	p-nitrotoluene	511.7
15.	nitrocyclopentane	453.2	60.	benzonitrile	464.1
16.	N-allylaniline	492.2	61.	3-nitrobenzotrifluoride	475.9
17.	ethylamine	289.7	62.	di-n-propilamine	382
18.	p-nitrophenole	552.2	63.	nitrohexane	436.8
19.	cyclopentylamine	380.2	64.	phenilhidrazine	516.7
20.	2-methylbutylamine	368.7	65.	methylamine	266.8
21.	N-methylbutylamine	364.2	66.	3-methylpyridine	417.3
22.	benzylamine	457.7	67.	aniline	457.2
23.	p-methoxyaniline	514.7	68.	p-chloroaniline	503.7
24.	m-methoxyaniline	524.2	69.	m-chloroaniline	501.7
25.	o-methoxyaniline	498.2	70.	n-pentylamine	377.6
26.	t-pentylamine	350.2	71.	isobutylamine	340.9
27.	dimethylamine	280	72.	diethylamine	328.6
28.	1-(2-aminoethyl)-piperidine	459.2	73.	tert-butylamine	317.5
29.	1-(2-aminoethyl)-piperidine	493.2	74.	n-butylamine	350.6
30.	9-methyl carbazole	616.8	75.	pirolidine	359.7
31.	carbazole	627.8	76.	nitromethane	374.4
32.	4-methylaniline	473.6	77.	isobutyronitrile	376.8
33.	3-methylaniline	476.5	78.	n-butyronitrile	390.8
34.	2-methylaniline	473.5	79.	cis-crotonitrile	380.6
35.	2-propylamine	304.9	80.	trimethylamine	276
36.	1-naphtylamine	573.8	81.	2-nitropropane	393.4
37.	nitroethane	387	82.	1-nitropropane	404.3
38.	piperidine	376.4	83.	propionitrile	370.5
39.	4-methylpyridine	418.5	84.	acrylonitrile	350.5
40.	2-methylpyridine	402.5	85.	N-methylhexylamine	414.2
41.	pyridine	388.4	86.	n-heptylamine	428.2
42.	pyrole	402.9	87.	N-tert-butylisopropylamine	371.2
43.	2-butylamine	335.9	88.	2-aminoheptane	416.2
44.	triethylamine	516.2	89.	malononitrile	491.5
45.	ethylenimine	329	90.	hydrogen cyanide	298.8

compounds;  $R = 0.990$ ;  $s = 10.7$  K. The largest pairwise  $R$  value of descriptors was 0.83. The modeling was performed by the ADAPT system.<sup>46</sup>

Our aim was to verify the quality of our property descriptors exactly in the same conditions as given in ref.<sup>43</sup> Thus, we extracted from the initial set of 104 N-containing compounds the same subset of 90 structures.

Molecular geometries and partial charges were calculated by the semiempirical AM1 method. The set of 19350 descriptors were reduced to 16383 after the monivariate regression.<sup>19</sup> Our procedure for finding the optimal subset of descriptors led to a subset of 72 descriptors. The best scores in ten variate regression for the set of 90 compounds of Table 5 are listed in Table 6.

Table 6. The Best Multivariate Regressions for the 90 Structures of Table 5.

No	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>	X <sub>7</sub>	X <sub>8</sub>	X <sub>9</sub>	X <sub>10</sub>	R
1	7	4959									0.944600
2	4	4797									0.945140
3	5	4959									0.948670
4	5	4959	14250	14607							0.962734
5	5	4959	10990	9671							0.962809
6	5	4959	10990	7206							0.967113
7	5	4959	10990	9671	16	3320					0.974787
8	5	4959	10990	9671	3320	6422					0.976078
9	5	4959	10990	9671	3528	6422					0.976574
10	5	4959	10990	9671	3528	6422	7256	16148			0.979844
11	5	4959	10990	9671	3528	6422	16148	6895			0.979862
12	5	4959	10990	9671	3528	6422	6895	15789			0.980012
13	5	4959	10990	9671	3528	6422	6895	15789	16158	16225	0.983581
14	5	4959	10990	9671	3528	6422	16158	16225	15060	15789	0.984335
12	6	6895	4	16275	963	841	163	13920	1	4727	0.985431

The best model was:

$$\begin{aligned}
 BP_{\text{calc}} = & 225.441 - 59.627 \cdot \ln DTsDiP\_p/d2SE2 + 316.627 \cdot RTsDiPp2/d2AE\_ + & (53) \\
 & 1.124 \cdot DGfDePp2/d2PP\_ - 1729.562 \cdot 1/DTsDiE\_p*d\_HE2 - \\
 & 0.010 \cdot 1/DTsDePp2/d2SP2 - 49.623 \cdot 1/DGsDeP\_p*d\_HE\_ + \\
 & 8.846 \cdot \ln DGjDiPp2/d2GP\_ - 4.698 \cdot 1/RGjDeP\_p*d\_GP\_ - \\
 & 12.188 \cdot \ln DGjDeP\_p/d\_HP\_ + 33.597 \cdot DGjDeE\_p\_SE2
 \end{aligned}$$

$R = 0.98543$ ;  $s = 13.149$ ;  $n = 90$

The plot corresponding to eq (53) is given in Figure 6.

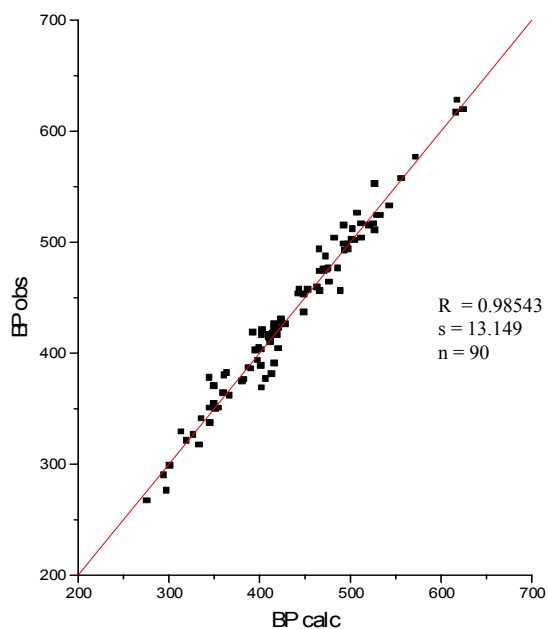


Figure 6. The plot of calculated vs observed normal boiling points  
(the set of Table 5)

Our result is slightly lower ( $s = 13.149$  K) than that reported in ref.<sup>43</sup> ( $s = 10.7$  K). It is possible to further improve the model by mining the whole descriptor pool not only within the limits of a heuristic procedure. Another possibility is to use different training subset selection and outlier elimination. Such procedures will be reported in a future paper.

#### Set 4. Nitrophenols.

A set of 25 nitrophenols<sup>47</sup> showing herbicidal activity (Table 7) was considered for correlation with the Cluj Property indices. Nitrophenols are known to inhibit the electronic flux of photosynthesis.

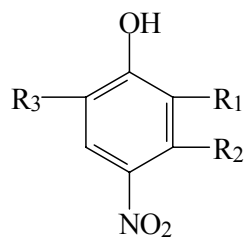


Table 7. Nitrophenols and Their Herbicidal Activity

No	R <sub>1</sub>	R <sub>2</sub>	R <sub>3</sub>	pI <sub>50</sub>
1	H	methyl	methyl	3.3
2	H	methyl	isopropyl	4.1
3	H	H	t-butyl	5.7
4	H	H	phenyl	4.35
5	H	H	cyclohexyl	4.85
6	Cl	methyl	methyl	4.89
7	Cl	methyl	isopropyl	6.07
8	Cl	H	t-butyl	6.88
9	Cl	H	phenyl	6.45
10	Cl	H	cyclohexyl	6.52
11	Br	methyl	methyl	5.25
12	Br	methyl	isopropyl	6.70
13	Br	H	t-butyl	6.15
14	Br	H	phenyl	6.52
15	Br	H	cyclohexyl	6.75
16	I	methyl	methyl	6.24
17	I	methyl	isopropyl	6.70
18	I	H	t-butyl	7.03
19	I	H	phenyl	6.86
20	I	H	cyclohexyl	6.65
21	NO <sub>2</sub>	H	H	3.00
22	NO <sub>2</sub>	H	methyl	3.70
23	NO <sub>2</sub>	H	s-butyl	5.10
24	NO <sub>2</sub>	H	t-butyl	5.79
25	NO <sub>2</sub>	H	cyclohexyl	6.05

Table 8 list the best scores of correlation in decreasing order. From this table it can be seen that the monivariate and divariate regression are not satisfactory. Additional variables are needed for good statistics (entries 6-13).

Table 8. Mono- and Multivariate Regression for Nitrophenols

No	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>	R
1	1						0.94991
2	6	155					0.97089
3	7	174					0.96974
4	11	10620					0.96654
5	5	260					0.96617
6	13028	15806	91	15636			0.99078
7	13028	15806	12	15398			0.99013
8	12	15806	13891	15749			0.99067
9	7	13028	382	14214			0.98932
10	13028	15806	12	15398	15228	15865	0.99850
11	12	15806	13891	15749	15648	16378	0.99674
12	7	13028	382	14214	13186	16282	0.99650
13	13028	15806	91	15636	14064	14943	0.99562

The best model is given in eq 54 (see also entry 10, Table 8):

$$\begin{aligned}
 \text{Predicted activity} = & 8.062 - 0.003 \cdot \text{RGsDeM\_p/d2PE2} + 0.395 \cdot 1/\text{DTsDiP\_p*d\_HE\_} \\
 & - 0.000008 \cdot 1/\text{RTsDiPp2/d2HE2} - 229.564 \cdot 1/\text{DGjDeMp2/d2PE2} \\
 & + 0.003 \cdot \text{RGjDiPp2/d2HP\_} + 0.004 \cdot \text{DTsDeP\_p*d\_HP2} \quad (54) \\
 & R = 0.9985; s = 0.067; n = 25
 \end{aligned}$$

The plot of the predicted vs observed herbicidal activity, cf. eq 54, is shown in Figure

7.

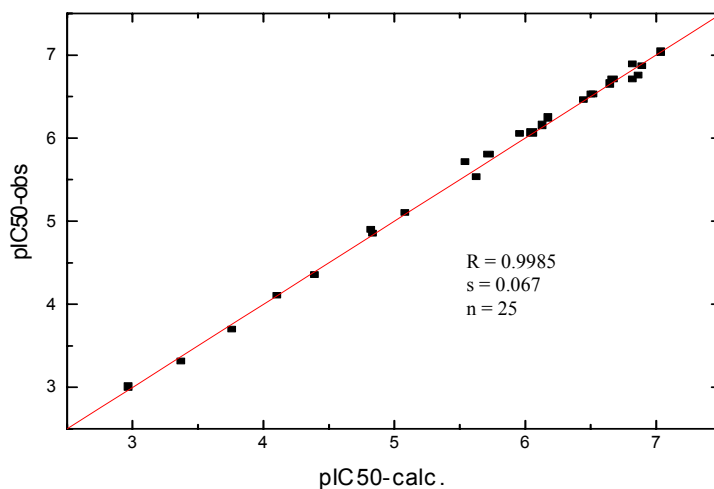


Figure 7. The plot of predicted vs. observed herbicidal activity

The descriptors involved in eq 54 show a rather low inter-correlation (Table 9). The average absolute value of the pairwise correlation coefficients was 0.2200.

Table 9. Intercorrelation of the indices in entry 10, Table 8.

	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>
X <sub>1</sub>	0.216489	0.697018	0.067367	0.121666	0.171049
X <sub>2</sub>		0.137108	0.015049	0.060902	0.264608
X <sub>3</sub>			0.201832	0.039716	0.269802
X <sub>4</sub>				0.682383	0.22193
X <sub>5</sub>					0.133726

## Discussion

The fragmental property indices take into account the chemical nature of atoms (mass and electronegativity), various kinds of interactions between the fragments of molecules and the 3D geometry of molecular structures.

There exist an analogy between *CoMFA* and *FPIF*: both of them calculate the interaction of a chemical structure (or substructure) with a *probe atom* in the 3D space. The property of fragment  $Fr_{i,j}$  is viewed as the interaction of atoms forming the fragment  $Fr_{i,j}$  with



the atom  $j$ . The major difference is that *CoMFA* uses external probe atoms (with defined chemical identity) whereas *FPIF* considers internal probe atoms with no chemical identity. Only the fragments (i.e. substructures) are chemically well defined.

Bivariate correlation with indices belonging to *FPIF* offer good quality models for quite diverse molecular properties such as the inhibition of mitodepressive activity on *Lepidium Savitium L.* ( $R > 0.99$  - see set 1) and the aromatase inhibition (set 2) as well. The same is true for the sum of one-electron energy calculated at the Extended-Huckel level ( $R > 0.9999$ ).

Multivariate regression provided good models for the boiling points of a very diverse set of N-containing organic molecules (set 3) or for the herbicidal activity (set 4).

Note that there is no causal relationship between descriptors and a chosen property, despite each descriptor encodes some aspects of intra- and/or intermolecular interactions. In case of a highly noncongeneric database, (the case of the set 3) more than one model, with equivalent accuracy are expected. Our results provided different vectorial descriptions that are, in some extent, isomorphic.

The occurrence of a certain descriptor type can be however, indicative, for a possible causal relation between the structure and the investigated property. In the set 4, the best model (eq 54) is described by the mass (M) for  $X_1$  and partial charge (P) for the remainder descriptors. It is just the expected case: the herbicidal activity of nitrophenols (a congeneric set) is controlled by the acidity of the phenol group, which increases with decreased negative partial charge on the oxygen atom.

It appears that, in large ensembles of molecules, the correlation is not strongly dependent of the type of chemical description, which is not the case in the more compact sampling of congeneric sets. At the detailed scale, the observable properties strictly depend on the particular data set and require more specified description. It justifies the conclusion<sup>48</sup> that there exists different mapping behavior of the chemical space at different scales and which one is more suitable in a given problem remains at the latitude of the researcher.

The above presented results demonstrate the good correlating ability of *FPIF*. It represents a promise for further **QSPR/QSAR** studies.

**Acknowledgement.** One author (M. D.) gratefully thanks to Professor A. Kerber for fruitful discussions in the period he visited the University of Bayreuth, Germany.

## References

1. Free, S. M.; Wilson, J. W. A mathematical contribution to structure-activity studies, *J. Med. Chem.* **1964**, *7*, 395.
2. Gao, C.; Govind, R.; Tabak, H. H. Application of the group contribution method for predicting the toxicity of organic chemicals. *Environmental Toxicol. Chem.* **1992**, *11*, 631-636.
3. Kalivas, J. H.; Sutter, J. M.; Roberts, N. Global optimization by simulated annealing with wavelength selection for ultraviolet-visible spectrophotometry, *Anal. Chem.* **1989**, *61*, 2024-2030.
4. Kalivas, J. H. generalized simulated annealing for calibration sample selection from an existing set and orthogonalization of undesigned experiments. *J. Chemometrics*, **1991**, *5*, 37-48.
5. Sutters, J. M.; Dixon, S. L.; Jurs, P. C. Automated descriptor selection for Quantitative Structure-Activity Relationships, using generalized simulated annealing, *J. Chem. Inf. Comput. Sci.* **1995**, *35*, 77-84.
6. Leardi, R.; Boggia, R.; Terrile, M. Genetic algorithms as a strategy for feature selection. *J. Chemom.* **1992**, *6*, 267.
7. Lucasius, C. B.; Kateman, G. Understanding and using genetic algorithms. Part 1. Concepts, properties and context. *Chemom. Intell. Lab. Sys.* **1993**, *19*, 1.
8. Diudea, M.V. Cluj matrix  $CJ_u$  : source of various graph descriptors, *Commun. Math. Comput. Chem. (MATCH)*, **1997**, *35*, 169-183.
9. Diudea, M.V.; Minailiuc, O.; Katona, G.; Gutman, I. Szeged matrices and related numbers, *Commun. Math. Comput. Chem. (MATCH)*, **1997**, *35*, 129-143.
10. Diudea, M.V. Cluj matrix invariants, *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 300-305.
11. Diudea, M.V.; Pârv, B.; Topan, M.I. Derived Szeged and Cluj indices, *J. Serb. Chem. Soc.* **1997**, *62*, 267-276.
12. Diudea, M.V.; Gutman, I. Wiener-type topological indices, *Croat. Chem. Acta* **1998**, *71*, 21-51.
13. Kiss, A.A.; Katona, G.; Diudea, M.V. Szeged and Cluj matrices within the matrix operator  $W_{(M1,M2,M3)}$  *Coll. Sci. Papers Fac. Sci. Kragujevac* **1997**, *19*, 95-107.
14. Gutman, I.; Diudea, M.V. Defining Cluj matrices and Cluj matrix invariants, *J. Serb. Chem. Soc.* **1998**, *63*, 497-504.
15. Diudea, M.V.; Parv, B.; Gutman, I. Detour-Cluj matrix and derived invariants, *J.Chem.Inf.Comput.Sci.* **1997**, *37*, 1101-1108.
16. Diudea, M.V.; Katona, G.; Lukovits, I.; Trinajstić, I. Detour and Cluj-detour indices, *Croat. Chem. Acta* **1998**, *71*, 459-471.

17. Minailiuc, O.; Katona, G.; Diudea, M.V.; Strunje, M., Graovac, A.; Gutman, I. Szeged fragmental indices, *Croat. Chem. Acta* **1998**, *71*, 473-488.
18. Horn, R.A.; Johnson, C.R. *Matrix Analysis*; Cambridge Univ. Press, Cambridge, **1985**.
19. Diudea, M. V.; Gutman, I.; Jäntschi, L. *Molecular Topology*, Nova Science, Huntington, New York, (in press).
20. Sears, F.W.; Zemansky, M.W.; Young, H.D. *University Physics*, fifth edition, Addison – Wesley Publishing Company, USA-Canada, **1976**.
21. Golender, V.; Vesterman, B.; Vorpapel, E. APEX-3D Expert system for drug design. *Network Science*, **1996**, *Jan*, <http://www.netsci.org/Science/Compchem/feature09.html>.
22. Rose, V.S.; Wood, J. Generalized cluster significance analysis and stepwise cluster significance analysis with conditional probabilities. *Quant. Struct.-Act. Relat.* **1998**, *17*, 348-356.
23. Young, H.D. *Statistical Treatment of Experimental Data*. McGraw-Hill, New York, **1962**.
24. Reif, F. *Fundamentals of Statistical and Thermal Physics*, McGraw-Hill, Chap.1 New York, **1965**.
25. Diudea, M.V.; Silaghi-Dumitrescu, I.; Valence group electronegativity as a vertex discriminator. *Rev. Roumaine Chim.* **1989**, *34*, 1175-1182.
26. Diudea M.V.; Kacso I.E.; Topan M.I. Molecular topology. 18. A Qspr/Qsar study by using new valence group carbon-related electronegativities. *Rev. Roum. Chim.* **1996**, *41*, 141-157.
27. Crawford, F.S.Jr.; *Waves*, Berkeley Physics Course, Newton, Massachusetts, vol. 3, **1968**.
28. Ivanciuc, O. *3D QSAR models*, in *QSPR/QSAR studies by molecular descriptors*, Ed. Diudea, M.V. Nova Science, Huntington, New York, (in press).
29. Hopfield, J.J. Neural networks and physical systems with emergent collective computational abilities. *Proc. Natl. Acad. Sci. U.S.A.* **1982**, *79*, 2554-2558.
30. Zupan, J.; Gasteiger, J. Neural networks: a new method for solving chemical problems or just a passing phase? *Anal. Chim. Acta* **1991**, *248*, 1-30.
31. Gasteiger, J.; Zupan, J. Neural networks in chemistry. *Angew. Chem. Int. Ed. Engl.* **1993**, *32*, 503-527.
32. Zupan, J.; Gasteiger, J. *Neural Networks for Chemists*; VCH: Weinheim, **1993**.
33. Bulsari, A.B. *Neural Networks for Chemical Engineers*; Elsevier: Amsterdam, **1995**.
34. Devillers, J. *Neural Networks in QSAR and Drug Design*. Academic Press, London, **1996**, p. 279.
35. Ivanciuc, O.; Rabine, J.-P.; Cabrol-Bass, D.; Panaye, A.; Doucet, J.P. <sup>13</sup>C NMR chemical shift prediction of sp<sup>2</sup> carbon atoms in acyclic alkenes using neural networks. *J. Chem. Inf. Comput. Sci.* **1996**, *36*, 644-653.
36. Ivanciuc, O. Molecular graph descriptors used in neural network models. In:

- Topological Indices and Related Descriptors in QSAR and QSPR*; Devillers, J., Balaban, A.T., Eds.; Gordon and Breach Science Publishers: The Netherlands, **1999**, pp. 697-777.
37. Gakh, A.A.; Gakh, E.G.; Sumpter, B.G.; Noid, D.W. Neural network-graph theory approach to the prediction of the physical properties of organic compounds. *J. Chem. Inf. Comput. Sci.* **1994**, *34*, 832-839.
38. Baskin, I.I.; Palyulin V.A.; Zefirov, N.S. A Neural device for searching direct correlations between structures and properties of chemical compounds. *J. Chem. Inf. Comput. Sci.* **1997**, *37*, 715-721.
39. Kireev, D.B. ChemNet: A Novel neural network based method for graph/property mapping. *J. Chem. Inf. Comput. Sci.* **1995**, *35*, 175-180.
40. Topliss, J. G.; Edwards, R. P. Chance factors in studies of Quantitative Structure-Activity Relationships. *J. Med. Chem.* **1979**, *22*, 1238.
41. Nikolić, S.; Medić-Sarić, M.; Matijević-Sosa, J. A QSAR study of 3-(Phtalimidoalkyl)-pyrazolin-5-ones, *Croat. Chem. Acta*, **1993**, *66*, 151-160.
42. Nagy, P.I.; Tokarski, J.; Hopfinger, A. J. Molecular shape and QSAR analysis of a family of substituted dichlorodiphenyl aromatase inhibitors, *J. Chem. Inf. Comput. Chem.* **1994**, *34*, 1190-1197.
43. Wessel, M. D.; Jurs, P. C. Prediction of normal boiling points for a diverse set of industrially important organic compounds from molecular structure, *J. Chem. Inf. Comput. Sci.* **1995**, *35*, 841-850.
44. Stanton, D. T.; Jurs, P. C. Development and use of charged partial surface area structural descriptors for Quantitative Structure-Property Relationships studies. *Anal. Chem.* **1990**, *62*, 2323.
46. Goll, E.S.; Jurs, P.C. Prediction of the Normal Boiling Points of Organic Compounds from Molecular Structures with a Computational Neural Network Model, *J. Chem. Inf. Comput. Chem.* **1999**, *39*, 974-983.
47. Stuper, A. J.; Brugger, W. E.; Jurs, P. C. *Computer-assisted studies of chemical structure and biological function*, Wiley-Interscience: New York, 1979.
47. Trebst, A.; Draber, W. *Advan. Pest. Sci., Symp. Papers IV-th Int. Congress Pest. Chem.*, Zurich, Swiss, 1978, pp. 223.
48. Benigni, R.; Gallo, G.; Giorgi, F.; Giuliani, A. On the equivalence between different descriptions of molecules: value for computational approaches. *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 575-578.